

Artificial intelligence in pharmacovigilance

CIOMS Working Group report
Draft, 1 May 2025

This report was posted for comment on 1 May 2025 at:

https://cioms.ch/working_groups/working-group-xiv-artificial-intelligence-in-pharmacovigilance/.

Please note that the layout will be improved in the final version, and best efforts will be made to correct remaining typographical and/or grammatical errors, as well those pertaining to references.

Permissions are being sought to reproduce some of the illustrative materials included in this report. We welcome responses from organisations that own any of these materials and have not yet been contacted in this regard.

Please submit your comments using the form posted on the CIOMS website at

https://cioms.ch/working_groups/working-group-xiv-artificial-intelligence-in-pharmacovigilance.

The timeline for submission of comments is 6 June 2025.

Thank you.

Suggested citation: Artificial intelligence in pharmacovigilance. CIOMS Working Group report. Geneva, Switzerland: Council for International Organizations of Medical Sciences (CIOMS), 2026.

Copyright © 2026 by the Council for International Organizations of Medical Sciences (CIOMS)

DOI *number tbc*

All rights reserved. CIOMS publications may be obtained directly from CIOMS through its publications e-module at <https://cioms.ch/publications/>. Further information can be obtained from CIOMS, P.O. Box 2100, CH-1211 Geneva 2, Switzerland, www.cioms.ch, e-mail: info@cioms.ch.

This publication is freely available on the CIOMS website at:

<https://cioms.ch/publications/product-category/recently-published/>

Disclaimer: The authors alone are responsible for the views expressed in this publication, and those views do not necessarily represent the decisions, policies or views of their respective institutions or companies.

Acknowledgements

(to follow)

Table of contents

Table of contents	iv
List of figures and tables	vii
Abbreviations	ix
Preface	xii
Executive summary	13
Chapter 1: Introduction	16
Scope	21
Chapter 2: Landscape analysis	23
Use of artificial intelligence in pharmacovigilance to date	23
Regulatory considerations	26
Principles for integrating and implementing artificial intelligence within pharmacovigilance	37
Chapter 3: Risk-based approach	39
Introduction	39
Risk assessment	41
Issue detection and risk mitigation	43
Documentation	44
Chapter 4: Human oversight	46
Introduction	46
Considerations on human involvement and oversight	47
Transformation of traditional roles	48
Chapter 5: Validity & Robustness	51
Introduction	51
Specification and design	52
Performance evaluation	54
Reproducibility	56
Assessing artificial intelligence solutions with human-in-the-loop	56
Continuous integration and deployment	57
Chapter 6: Transparency	59
Introduction	59
Disclosing use of artificial intelligence	59
Transparency regarding the artificial intelligence model	60
Explainability	61
Transparency regarding performance	63

Chapter 7: Data privacy	66
Introduction	66
Ethical considerations	66
Data privacy regulations	68
Practical considerations to support data privacy	72
Conclusions	75
Chapter 8: Fairness & Equity	77
Introduction: general concepts and considerations	77
Fairness and equity considerations and pharmacovigilance	78
Sources of potential threat to fairness and equity	78
Risk, impact, and mitigation measures	81
Key mitigation strategies	81
Chapter 9: Governance & Accountability	83
Introduction	83
Governance framework	84
Traceability and version control	88
Chapter 10: Future considerations for development and deployment of artificial intelligence in pharmacovigilance	91
The evolution and future of artificial intelligence in pharmacovigilance	91
Transformative role of pharmacovigilance long-term and beyond: from detection to prevention	91
Future development and deployment of AI and the guiding principles	92
Conclusions to the future vision of artificial intelligence	96
APPENDIX 1: Glossary	98
APPENDIX 2: Comparison table of guiding principles	112
APPENDIX 3: Use cases	121
Use Case A: Large Language Models data extraction for case processing	121
Use Case B: Case deduplication	124
Use Case C: Artificial intelligence translation assistant	127
Use case D: Large language models for context-aware Structured Query Language	129
Use Case E: Causality assessment of adverse drug reactions	131
Use Case F: Process efficiencies supporting signal detection	134
Use Case G: Generative artificial intelligence for enhanced and intelligently structured outputs from large pharmacovigilance document libraries	137
Use Case H: Artificial intelligence to support diagnosis and prediction of (hydroxy)chloroquine retinopathy	140
APPENDIX 4: Content related to explainability and to Fairness & Equity	143
Content related to explainability	143

Content related to Fairness and Equity	150
APPENDIX 5: CIOMS Working Group membership and meetings	152
APPENDIX 6: Public consultation commentators	155

List of figures and tables

List of figures

[Figure 1]: Traditional signal management process.....	18
Figure 2: Growth over time of VigiBase, the World Health Organization global database of adverse event reports for medicines and vaccines	19
Figure 3: Growth over time of the FDA Adverse Event Reporting System (FAERS) database	20
Figure 4: Multi- modal input models	25
Figure 5: Model risk matrix	43
Figure 6: Relationship of timing of large language model introduction, parameter size and attentiveness to data privacy.....	75
Figure 8: Flowchart summarizing the implementation of the automatic selection of controls and the dismissal of false positive signals when using a conditional inference tree	134
Figure 9: An outline of our initial artificial intelligence architecture	137
Figure 10: FDA's Information Visualization Platform user interface illustrates the system capabilities, focusing on the features that positively contribute to classification for assessability...	144

List of tables

Table 1: Examples of deployed artificial intelligence solutions in pharmacovigilance described in the public domain	25
Table 2: Comparison of CIOMS Working Group XIV guiding principles for artificial intelligence across regional and country government institutions, and international organizations	27
Table 3: Key aspects of an artificial intelligence model to disclose to stakeholders	60
Table 4: Relevant aspects to disclose to ensure transparency regarding the estimated performance of an artificial intelligence solution.....	64
Table 5: Data privacy regulations for using secondary data in Germany	70
Table 6: Data privacy regulations for using secondary data in China	71
Table 7: Data privacy regulations for using secondary data in Japan.....	71
Table 8: Governance framework grid	84
Table 9: Comparison of CIOMS Working Group XIV guiding principles for artificial intelligence across regional and country government institutions, and international organizations – Extracted description of principles	112
Table 10: Use case A: Alignment with the governance framework (detail)	122
Table 12: Use case C: Alignment with the governance framework (detail)	128
Table 14: Use case E: Modeling Data	131
Table 15: Use case E: Alignment with the governance framework (detail).....	132
Table 16: Use case F: Alignment with the governance framework (detail).....	135
Table 17: Use case G: Alignment with the governance framework (detail)	138

Table 18: Use case H: Alignment with the governance framework (detail) 141

Abbreviations

2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

ADR	Adverse Drug Reaction
AE	Adverse Event
AI	Artificial Intelligence
AIA	Algorithmic Impact Assessment
AIDA	Artificial Intelligence and Data Act (Canada)
ALTAI	Assessment List for Trustworthy Artificial Intelligence
ATC	Anatomical Therapeutic Chemical
BLEU	Bilingual Evaluation Understudy
CDC	Centers for Disease Control and Prevention
CIOMS	Council for International Organizations of Medical Sciences
CNN	Convolutional Neural Networks
CRISPR	Clustered Regularly Interspaced Short Palindromic Repeats
CSV	Computerized System Validation
EEA	European Economic Area
EHR	Electronic Health Records
EMA	European Medicines Agency
ETHER	Event-based Text-mining of Health Electronic Records
EU AI Act	European Artificial Intelligence Act
EU	European Union
FAERS	Food and Drug Administration Adverse Event Reporting System (USA)
GAMP 5	Good Automated Manufacturing Practice 5
GDPR	General Data Protection Regulation
GenAI	Generative Artificial Intelligence
GPT	Generative Pre-trained Transformer
GSD	Global Safety Database
GVP	Good Pharmacovigilance Practices
GxP	Good [x] Practices
HIC	Human-in-Command
HIPAA	Health Insurance Portability and Accountability Act
HITL	Human-in-the-Loop
HOTL	Human-on-the-Loop
ICH	The International Council for Harmonisation of Technical Requirements for Pharmaceuticals for Human Use

37	ICSR	Individual Case Safety Report
38	InfoViP	Information Visualization Platform
39	IRB	Institutional Review Board
40	ISIC	International Skin Imaging Collaboration
41	ISPE	International Society for Pharmaceutical Engineering
42	IT	Information Technology
43	KPI	Key Performance Indicator
44	LIME	Local Interpretable Model-Agnostic Explanations
45	LLM	Large Language Models
46	MAH	Marketing Authorisation Holder
47	MedDRA	Medical Dictionary for Regulatory Activities
48	MHRA	Medicines and Healthcare products Regulatory Agency (UK)
49	ML	Machine Learning
50	ML-DSF	Machine Learning-Enabled Device Software Functions
51	NIST	National Institute of Standards and Technology
52	NLP	Natural Language Processing
53	OECD	Organisation for Economic Co-operation and Development
54	OTC	Over the Counter
55	PASS	Post-Approval Safety Studies
56	PD	Pharmacodynamic
57	PDMP	Prescription Drug Monitoring Program
58	PHI	Protected Health Information
59	PK	Pharmacokinetic
60	PPV	Positive Predictive Value
61	PSMF	Pharmacovigilance System Master File
62	PSUR	Periodic Safety Update Reports
63	PV	Pharmacovigilance
64	QA	Quality Assurance
65	RCT	Randomized controlled trial
66	RMF	Risk Management Framework
67	RPA	Robotics Process Automation
68	RWD	Real-World Data
69	RWE	Real-World Evidence
70	SARAH	Smart Artificial Intelligence Resource Assistant for Health
71	SHAP	Shapley Additive exPlanations
72	SME	Subject Matter Expert

CIOMS Working Group XIV: Abbreviations

73	SmPC	Summary of product characteristics
74	SOP	Standard Operating Procedures
75	US FDA	U.S. Food and Drug Administration
76	UTC	United Therapeutics Corporation
77	WHO	World Health Organization
78	xAI	Explainable Artificial Intelligence

79

80

Preface

81 The Council for International Organizations of Medical Sciences (CIOMS) has played a pivotal role in
82 the advancement of modern pharmacovigilance (PV) by developing guidelines that address ethical
83 and scientific aspects of drug development and safety. Notably, CIOMS has published guidance
84 documents that have supported a structured approach for the collection and reporting of adverse
85 drug reactions (ADRs) in addition to guidance on practical aspects of signal detection in PV, fostering
86 international collaboration and standardization in drug safety monitoring.

87 The thalidomide tragedy of the early 1960s exposed severe deficiencies in global drug safety
88 practices, highlighting the need for comprehensive data collection and international harmonization.
89 In response, the World Health Organization (WHO) established the Program for International Drug
90 Monitoring in 1968, initiating efforts to share individual case reports between countries and
91 harmonize data practices. Building on these foundational efforts, the late 1980s and the 1990s saw
92 pivotal CIOMS reports like the *Monitoring and Assessment of Adverse Drug Effects* (1985) and the
93 *International Reporting of Adverse Drug Reactions* (1987), both by the CIOMS Working Group I, and
94 the *Current Challenges in Pharmacovigilance: Pragmatic Approaches* (1999) by the CIOMS Working
95 Group V. Subsequent CIOMS Working Group reports and the establishment of the International
96 Council for Harmonisation of Technical Requirements for Pharmaceuticals for Human Use (ICH)
97 aimed to address the fragmented approaches to drug safety identified decades earlier, providing a
98 framework for standardized adverse event collection and reporting in addition to signal detection
99 processes.

100 Advancements in technology have transformed PV, with artificial intelligence (AI) playing an
101 increasing role in healthcare. At this juncture, regulatory harmonization and transparency remain
102 critical to leveraging AI effectively, ensuring that data is reliable and that AI systems deliver accurate
103 and trustworthy results, while ensuring safe and responsible AI use.

104 As AI continues to evolve and impact biomedical research, its increased integration in and impact on
105 PV practice is inevitable. AI's potential for transformative disruption, compels us to engage in critical
106 discourse: how do we wish to see AI developed, validated, and deployed within this domain?

107 Since the CIOMS expert Working Group XIV on AI in PV was established in early 2022, there has been
108 significant progress in the field, marked by the rapid development and widespread availability of
109 generative AI (GenAI). While there is growing interest in exploring GenAI for PV applications, we
110 recognize the need to focus on its appropriate use, which brings specific challenges in highly
111 regulated domains such as PV, and we look to distinguish where possible and beneficial from general
112 issues of AI use. In light of this swiftly evolving context, this report aims to offer a general framework
113 of principles and good practices for developing and using AI in PV. Rather than offering technical
114 guidance, the aim is to ensure continued relevance as AI capabilities advance. The report focuses on
115 applications that are specific to PV or issues of particular importance to PV rather than general use
116 of AI, highlighting issues of high importance or priority to PV, even if they may not have widespread
117 attention or relevance in other contexts.

118 This report aims to provide guidance to those working in PV, in addition to organisations and
119 vendors developing AI solution for the PV domain.

120

121

122

123

124

Executive summary

125 The CIOMS report on artificial intelligence (AI) in pharmacovigilance (PV) tackles a rapidly emerging
126 cross-disciplinary field that is at the intersection of PV, computer science, regulation, law, medicine,
127 human rights, psychology and social science. Consequently, just as with medicinal products, it is
128 important to establish the approved indications, posology, side effects, and warnings and
129 precautions for use of AI in PV. This must be clearly defined and understood by many people from
130 diverse backgrounds to propel research and practical implementation in an effective, safe and
131 responsible manner. The diverse pool includes PV professionals, researchers, and decision makers
132 working in PV in the industry, government, academia and software vendors, who are developing AI
133 solutions in the PV space, including signal management and all aspects of Individual Case Safety
134 Report (ICSR) processing. This report provides the requisite terminology and conceptual
135 understanding to actively engage in this space, either by participating in the applied scientific
136 research and public discourse, or by performing evaluations and making decisions at their respective
137 organizations.

138 Perhaps more than other Council for International Organizations of Medical Sciences (CIOMS) report
139 topics, the potential adverse effects of AI in PV and related points are major elements of our key
140 results because rapidly evolving, advanced and often opaque technologies may generate a rush of
141 excited promotion and initial over-estimation of utility, observed in so called technology “hype
142 cycles”, that does not correspond with the practical realities. There is a corresponding safety net of
143 core guiding principles for human protection elaborated by multiple organizations, through which AI
144 PV must grow. This report provides a set of guiding principles and corresponding organizations that
145 have elaborated each one. These principles form the bulk of the report: a risk-based approach,
146 human oversight, validity and robustness, transparency, data privacy, and governance and
147 accountability. Key points to consider for these guiding principles are elaborated throughout the
148 report and summarized concisely below.

149 Similar to prior CIOMS reports, this one benefits from a consensus position from multiple
150 stakeholders, including those based in regulatory agencies, academia, and industry. The Working
151 Group (WG) recognized that the field of AI is progressing so rapidly, that a prescriptive document
152 would likely be quickly outdated. Instead, the WG decided to focus upon a set of common principles
153 that were expected to be useful for years to come for PV professionals. PV is but one of a myriad of
154 AI applications that are now transforming many aspects of modern life. As such, this reports benefits
155 as well from the increasing interest in AI by national governments, several of which have issued
156 legislation and guidances not only on AI in drug development but also more broadly on the general
157 use of AI.

158 **Risk-based approach.** Integrating AI into PV processes needs to take into account the potential
159 inaccuracies and variability of AI algorithms, and corresponding impacts on the safety and well-being
160 of individuals and society. The level of risk, and corresponding intensity of required oversight, will be
161 a function of how high stakes the decision made by the AI is and whether the machine is intended to
162 be used in an unchecked stand-alone mode or with human-machine interaction. A sound risk-based
163 approach, in which the human oversight in the development and deployment of AI is commensurate
164 with these risks, enables organisations to make the most of AI capabilities while ensuring that
165 neither patient safety nor PV stakeholders are adversely affected. The risk-based approach applies to
166 the human oversight modalities, the validity and robustness strategy, the level of transparency, and
167 the efforts to uphold fairness and equity, and data privacy. The risk assessment should consider the
168 AI system itself, the context of use, and the potential impact and likelihood of risks materialising. A
169 risk-based approach should be reviewed at regular intervals and adapted if needed.

170 **Human oversight.** Human oversight supports performance optimization of AI in PV and increases
171 trustworthiness and accountability. The extent and nature of human oversight for an AI solution

172 should be risk-based, incorporating quality assurance principles. The human oversight might be
 173 “human-in-the-loop” meaning that the decision is the end results of a human-machine interaction,
 174 while in “human-on-the-loop”, the machine autonomously makes a decision or otherwise returns a
 175 result that is checked by a human. Human oversight is necessary to define fit-for-purpose levels of
 176 performance for the intended task (i.e. validity). It involves predefining acceptable performance
 177 levels, selecting appropriate data for model development in a realistic setting, an ongoing quality
 178 assessment process and retraining of the model as needed. Increased use of automation and AI in
 179 PV will transform traditional roles and competencies, requiring appropriate change management
 180 and training strategies.

181 **Validity & Robustness.** PV stakeholders must learn to critically appraise proposed AI solutions.
 182 Performance evaluation must demonstrate acceptable and robust results for intended use under
 183 realistic conditions. Such an evaluation should be both qualitative and quantitative, and a cross-
 184 disciplinary exercise and span a diverse range of relevant examples. Evaluations should use a
 185 sufficient representation of relevant data types to detect biases, promote adequate and
 186 generalizable performance across the intended deployment domain, assess usability, and identify
 187 circumstances associated with underperformance. Enrichment strategies to obtain representative
 188 test sets with high enough prevalence of the outcome may be required. Special care should be taken
 189 to ensure that performance evaluation results generalize to real-world settings.

190 **Transparency.** Declaring when and how AI solutions are used is critical for building trust among
 191 stakeholders. The nature of AI solutions deployed for core PV tasks should be communicated,
 192 including model architectures, expected inputs and outputs, and the nature of human-computer
 193 interaction. To fully characterize an AI solution’s effectiveness and limitations, performance
 194 evaluation results should describe the scope and nature of the test set(s) used including reference
 195 standards and sampling strategies. Performance metrics should be relevant for the intended tasks,
 196 compared with relevant benchmarks, and complemented by qualitative review of representative
 197 examples of correct and incorrect output. Explainability is an important concept relevant to those
 198 models whose internal decision pathways are so intricate and non-linear that they remain
 199 inscrutable even to technically literate persons – so called black boxes of the first kind. Explainable AI
 200 are a set of techniques that “look under the hood” and return plausible hypothesis about these
 201 pathways – roughly how the black box arrived at its outputs. To be able to do this can be
 202 advantageous to model building/trouble shooting, building trust, establishing auditability and
 203 accountability, including providing a basis for a human to challenge an AI result that may be
 204 adversely impacting them, and regulatory compliance and scientific hypothesis generation.
 205 However, explainable AI methods have limitations, and they only provide plausible hypothesis, but
 206 are no guarantee that the AI in fact used the hypothesized decision pathways.

207 **Data Privacy.** The ethical framework to evaluate the use of AI in PV is embedded within the standard
 208 principles for research activities involving human subjects. A crucial principle for the use of AI in
 209 routine PV is the sanctity of data privacy. With the increasing power of both the hardware and
 210 software that power AI, there is a vast potential to build large, linked databases, and the potential
 211 inherent in LLMs for patient reidentification. These may pose an ongoing challenge to the traditional
 212 safeguards that protect data privacy. In this context, there are multiple opportunities to reveal
 213 highly sensitive personal and health information to a broad, cross-disciplinary range of stakeholders
 214 throughout the AI development and deployment workflow. Consequently, countries have been
 215 enacting legislation and guidances intended to protect these data. PV professionals should recognize
 216 that existing procedures used to assure regulatory compliance may need to be reevaluated due to
 217 the heightened risks of GenAI to compromise data privacy.

218 **Fairness & Equity.** Supporting fairness and equity, avoiding propagating or amplifying harmful
 219 explicit biases underserving certain subpopulations, discrimination and inaccurate results during
 220 model development and deployment are regulatory and ethical imperatives. Equity may be
 221 advanced by taking measures to assure that AI in PV returns outputs that are relevant to populations

222 anticipated to have exposure to the specific medicinal product being evaluated. Screening,
223 identifying and excising explicit or potential bias when possible is key to mitigating risk, determining
224 AI applicability and limitations, and defining acceptable performance. Scrutinize training and
225 performance evaluation of reference data sets for adequate representation and evaluate
226 performance in relevant subgroups when possible. Inadequate reference data is often the cause of
227 inadequate fairness and equity.

228 **Governance & Accountability.** Robust governance and clear accountability are crucial for the
229 success of AI initiatives. These principles help ensure that AI systems are used safely, responsibly and
230 ethically, and in compliance of all applicable legal and regulatory mandates while fostering trust and
231 transparency among stakeholders. Clearly defined roles and responsibilities are crucial to enable all
232 stakeholders to understand their accountability and obligations in order to effectively oversee AI
233 systems.

234 As AI technology evolves, governance and accountability frameworks will need to be adapted. New
235 risks and challenges will emerge, requiring updated principles and practices. Continuous review and
236 adaptation are essential for staying ahead of these changes. This includes the adaptation refinement
237 of the proposed grid for practical use.

238 **Future considerations for development and deployment of artificial intelligence in**
239 **pharmacovigilance.** Increasing deployment of AI in PV is expected to prioritize and accommodate
240 rapid data collection, assessment and reporting for signal detection in real or quasi real time. This
241 may also be accompanied by a relative shift from warm-start to cold-start prediction scenarios (i.e.
242 post-approval to early-stage drug development). This could fundamentally change the way we work
243 to take advantage of these technological advances, for example, streamlining processes and causing
244 changes in the wider healthcare environment and beyond, including patient privacy. We also expect
245 to see increasing deployment of AI in PV in the clinic, where it will support primary, secondary and
246 tertiary prevention of adverse drug reactions. The extent to which humans remain in or on the loop
247 will be determined by the nature of the task (e.g. routinized tasks versus those requiring expert
248 clinical and scientific judgement), consistent with the elaborated risk-based approach, but it is
249 possible that some AI-based expert systems will eventually develop refined medical and scientific
250 judgement.

251 It is critical that the guiding principles outlined in this report remain as core considerations, but they
252 will need to evolve and adapt with advancements and application of AI in PV and medicine in
253 general. This is to ensure AI use in PV remains unbiased, transparent, and secure to prevent misuse
254 or accidental harm. The appropriate human oversight, including regulatory and ethical safeguards,
255 will be as crucial as the technological advancements being applied.

256

257 Chapter 1: Introduction

258 An artificial intelligence (AI) system is a machine-based system that, for explicit or implicit
259 objectives, infers, from the input it receives, how to generate outputs such as predictions,
260 content, recommendations, or decisions that can influence physical or virtual environments.¹
261 Different AI systems vary in their levels of autonomy and adaptiveness after deployment. In the
262 context of pharmacovigilance, the use of AI systems and activities is aimed at enhancing drug safety
263 monitoring, patient safety and regulatory compliance. PV is practiced not only at pharmaceutical
264 companies, health authorities, drug monitoring centres and academia, but also in the clinic, and AI is
265 finding applications to PV in all these settings.²

266 An AI solution is designed to address specific objectives within PV. The overall AI solution could be
267 developed with one or many AI systems. An AI system encompasses not only the model itself but
268 also includes the components necessary for utilization including user interfaces and data processing
269 pipelines. At the core of these systems are AI models. These models utilize parameters to learn
270 relationships within data, enabling the systems to adapt and improve model performance over time.

271 Simpler AI systems, such as statistical methods for signal detection, have been widely utilized in PV
272 for decades.³ However, the past decade(s) have seen drastic improvements in AI capabilities,
273 particularly in image analysis and natural language processing. These advancements have resulted in
274 a significant increase in their use. In addition, tremendous and continual advances in computing
275 power and model architectures have enabled the development and aggregation of large electronic
276 databases with potential for linkage. These have enabled the field of AI to be applied to an increasing
277 number of disciplines, including the life sciences.⁴ Within the life sciences, AI is being applied to a
278 growing number of areas, such as drug discovery and development, medical imaging and diagnostics,
279 genomics, precision medicine, public health, and healthcare delivery.⁵

280 Partly due to advances in AI, the pharmaceutical field is poised for rapid transformation across
281 clinical, regulatory and PV practices, aiming to streamline end-to-end processes to accelerate
282 product development and market delivery. Similarly, there is a growing emphasis focusing on
283 enhancing clinical and post marketing safety and risk management activities to enable proactive
284 identification (or even prediction) of safety signals and benefit-risk evaluation. In the clinic, AI is
285 being tested or deployed for early diagnosis (and thus secondary and tertiary prevention) of various
286 adverse drug reactions. Examples include early detection of hydroxychloroquine retinopathy,⁶
287 digoxin toxicity,⁷ and drug-induced movement disorders in Parkinsons patients.⁸

288 These advancements leverage massive integrated datasets and inductive logic, enabling AI models to
289 make informed decisions by utilizing accumulated data, rather than relying solely on explicit rules or
290 human intervention. This approach facilitates the development of AI tools that provide new,
291 improved, or complementary solutions. A critical enabler for AI success within PV will be the ability to
292 link and analyze large volumes of heterogeneous data of varying quality from diverse data sources,
293 such as electronic health records (EHRs), claims databases, registries, Internet of Things (IoT), and
294 connected devices. The ability to leverage health data can lead to potentially faster development of
295 new treatments, improved patient outcomes, and reduced healthcare costs, including the potential
296 for unlocking novel, useful, and actionable insights that might not have been identified otherwise.
297 Hence, there is an acute need to effectively communicate the key importance of data access to
298 support patient safety outcomes.

299 Incorporating AI into PV necessitates a thorough assessment of its potential benefits and risks,
300 helping stakeholders understand its implications for existing practices. Given the rapid pace of

301 change, this document does not prescribe specific uses for AI in PV but rather establishes and
302 promotes guiding principles for utilizing AI including ML.

303 The start of systematic safety surveillance predated the advent of the internet and widespread
304 electronic reporting capabilities. As such, it was a largely manual process that relied upon computing
305 for purposes such as summarizing data.

306 ICSRs are a key component of PV and remain a cornerstone of post-market safety surveillance as
307 they provide crucial safety information for an approved pharmaceutical product, which is important
308 to mitigate patient harm when assessed within a broader signal management system.

309 The processing of ICSRs involves several steps: collection, triage, data entry, quality review, medical
310 assessment, a with further transmission to other safety databases (e.g. regulatory authorities). As the
311 number of product approvals and the patient exposure grow, so do the number of reported adverse
312 events. The increased volume of ICSRs, coupled with stringent safety regulations, create significant
313 challenges in ICSR processing and compliance.

314 Once a signal is detected as a result of individual or aggregate analysis of adverse event reports, it
315 needs to be systematically investigated through sequential steps, which include signal triage,
316 validation, and, based on scientific assessment, formal evaluation using independent data sets, such
317 as hypothesis-testing research studies.⁹ Such investigation must be conducted in an integrated,
318 holistic fashion with all available scientific evidence and logic, offering wider opportunities for use of
319 AI for data insights (see Figure 1).

320 Traditional PV methods for analysis of adverse event reports include:^{10,11}

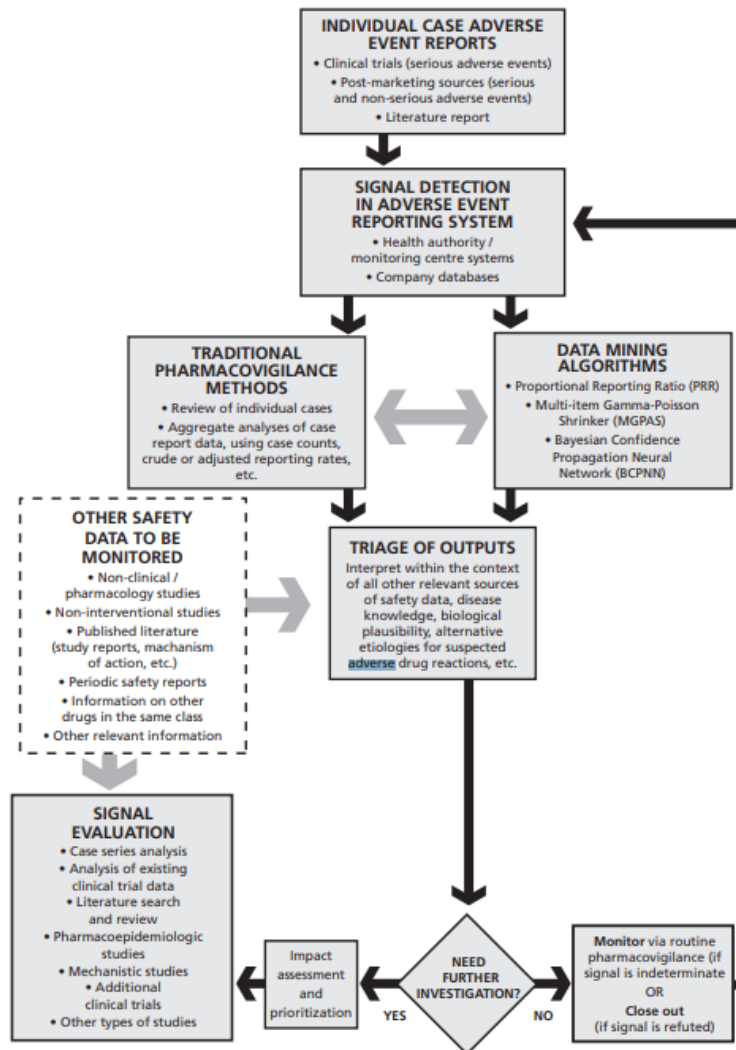
- 321 • Review of individual cases safety reports (ICSRs) or case series in a PV database or in
322 published medical or scientific literature; and
- 323 • Aggregate analyses of case reports using absolute case counts, simple reporting rates,
324 proportions or estimated exposure-adjusted reporting rates.

325 While ICSRs are fundamental to PV, other data streams are also considered throughout the PV
326 lifecycle. These streams may be directly linked or conceptually related and include pharmacokinetic/
327 pharmacodynamic (PK-PD) data, other real-world data (RWD), literature, and information from
328 clinical trials.

329 Once safety concerns (important identified risks or important potential risks) are identified, it is
330 essential to communicate them appropriately to a wide range of stakeholders. This is achieved
331 through documents such as aggregate reports, risk management plans, labelling information and
332 Dear Healthcare Professional communications (DHCPs).

333

334 **[Figure 1]: Traditional signal management process**
 335 Source¹²

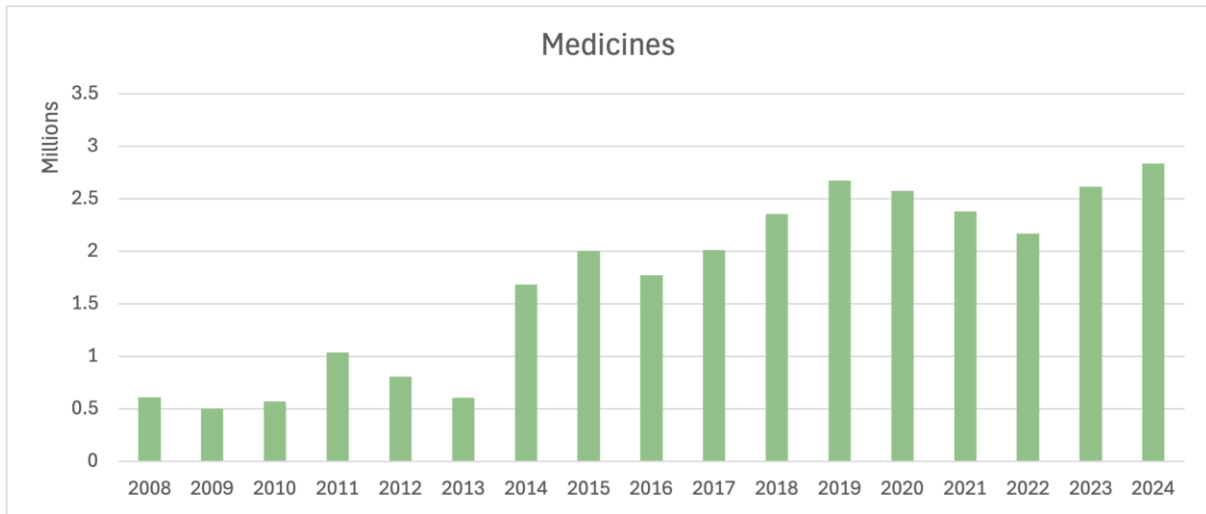


336
 337
 338
 339
 340
 341
 342
 343

The COVID-19 pandemic has further emphasized the need for advanced methods in PV, as it has led to a significant rise in safety reports (see Figures 2 and 3).^{13,14} As public awareness and expectations regarding drug safety continue to rise, there is a greater demand for robust PV systems that can effectively identify and mitigate potential risks associated with medications.

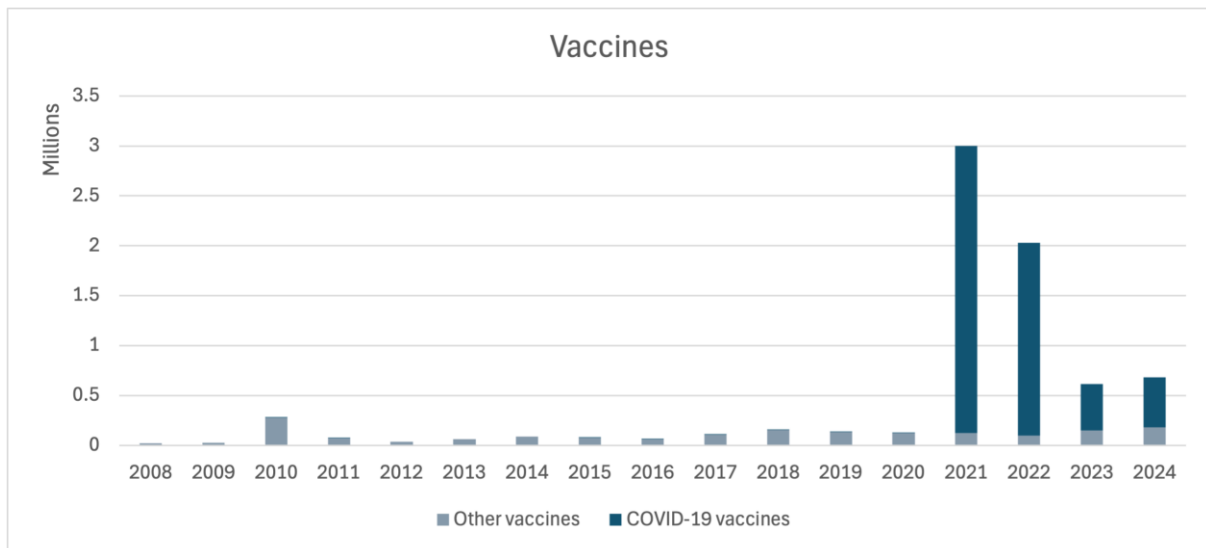
344 **Figure 2: Growth over time of VigiBase, the World Health Organization global database of adverse**
345 **event reports for medicines and vaccines**

346 Source: VigiBase accessed April 2025.
347
348



349

350



351

352

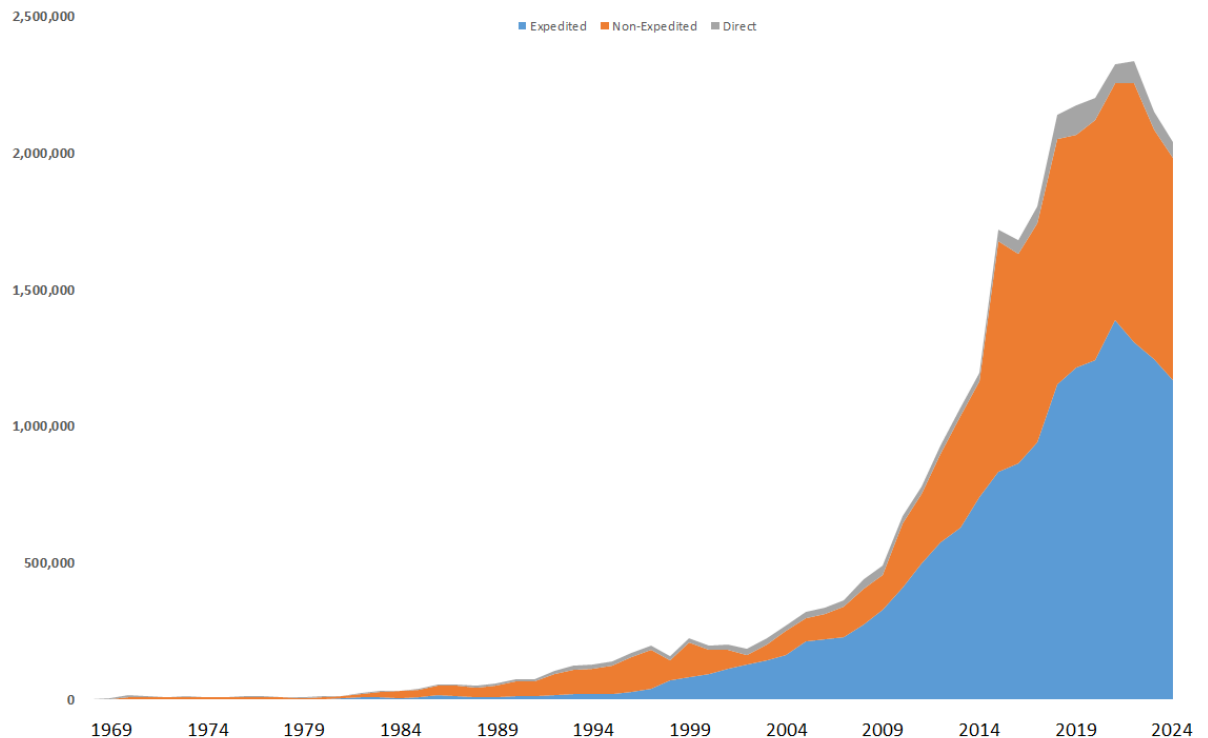
353

354 **Figure 3: Growth over time of the FDA Adverse Event Reporting System (FAERS) database**

355 Source: Constructed using FDA FAERS database.¹⁵

356

357



358

359

360 Nevertheless, the challenges of establishing and maintaining progressively more complex PV systems
 361 in a globally diverse and evolving regulatory environment are increasing. There is a need to rethink
 362 traditional PV strategies based on existing pressures on the one hand (e.g. managing increasing
 363 volumes and increasing regulatory complexity), and increasing and data sources on the other.

364 Technology solutions are already vital for the evolution of PV. While this notion of technology as a
 365 transformative enabler spans across all areas of product development, it is evident that applying
 366 innovative automation tools and processes to PV is no longer an option but an essential capability.

367 Rapid evolution of artificial intelligence

368 Traditional AI methods (e.g. K-means clustering, decision trees, support vector machines etc.) have
 369 traditionally been tailored for specific tasks, primarily utilizing supervised learning techniques. In
 370 contrast, deep neural networks such as BERT¹ have played a significant role in natural language
 371 processing, where they are pre-trained on large datasets and subsequently fine-tuned for specific
 372 applications delivering predictable outputs.

373 However, the landscape is evolving beyond this framework thanks to emerging technologies like
 374 Generative AI (GenAI). GenAI models are trained on expansive and varied text corpuses, often
 375 incorporating phases of human reinforcement learning. These models can perform specific tasks
 376 using sophisticated prompts, adopting zero-shot or few-shot learning techniques.

¹ BERT: Bidirectional Encoder Representations from Transformers

377 **Scope**

378 This document aims to guide those working in PV in addition to organisations developing AI solutions
379 for the PV domain, such as medicinal product regulators, medicinal products industry professionals,
380 software vendors, international and national PV organizations, researchers, and health care
381 professionals.

382 This report proposes a broad framework of principles and best practices for integrating and
383 implementing AI within PV, not technical guidance. Recognizing the rapid evolution and application
384 of AI technology, the CIOMS Working Group XIV developed this document to guide the development
385 and integration of AI tools into PV activities.

386 Our scope focuses on all aspects, direct and indirect, of the optimal collection, organization, analysis,
387 and communication of ICSRs from any source, including RWD, medical literature, randomized
388 controlled trials, and social media). Additionally, it includes productivity enhancers closely linked to
389 PV, such as tools that improve querying of safety databases¹⁶ or capabilities that enable faster, more
390 effective, or consistent data entry into a safety database which indirectly contributes to better safety
391 surveillance.¹⁷

392 The scope deliberately excludes broader healthcare data applications outside the direct purview of
393 safety, such as pharmacoepidemiology and other real world evidence study designs and conduct that
394 fall outside the realm of ICSRs. Similarly, the general use of AI as a productivity enhancer, if not
395 directly connected to PV activities (e.g. for email support), is excluded, as considerations may differ.

396 The scope has been intentionally limited to provide a practical guidance organised as principles and
397 their applications of AI in PV, rather than detailed guidance to ensure longevity. As AI is progressing
398 extremely rapidly, future opportunities and considerations are described in a later chapter.

399

400

References

¹ OECD (2024), "Explanatory memorandum on the updated OECD definition of an AI system", OECD Artificial Intelligence Papers, No. 8, OECD Publishing, Paris, <https://doi.org/10.1787/623da898-en>.

² Hauben M. Artificial Intelligence in pharmacovigilance: Do we need explainability? *Pharmacopidemiol Drug Saf.* 2022;Dec;31(12):1311-1316. ([Journal full text](#))
<https://doi.org/10.1002/pds.5501>

³ Kompa B, Hakim JB, Palepu A, Kompa KG, Smith M, Bain PA, Woloszynek S, Painter JL, Bate A, Beam AL. Artificial intelligence based on machine learning in pharmacovigilance: A scoping review. *Drug Safety.* 2023;Apr46(4):433. ([Journal full text](#)) <https://doi.org/10.1007/s40264-023-01273-9>

⁴ Webb S. Deep learning for biology. *Nature.* 2018;Feb22;554:555-557 ([Journal abstract](#)) <https://doi.org/10.1038/d41586-018-02174-z>

⁵ Noorbakhsh-Sabet N, Zand R, Zhang Y, Abedi V. Artificial intelligence transforms the future of health care. *The American Journal of Medicine.* 2019;Jul1;132(7):795-801. ([Journal full text](#)) <https://doi.org/10.1016/j.amjmed.2019.01.017>

⁶ Wright T, Yan P, Easterbrook M. Machine learning to identify multifocal ERG deficits in patients taking hydroxychloroquine. *Investigative Ophthalmology & Visual Science.* 2019;Jul22;60(9):5959. ([Journal full text](#))

⁷ Chang DW, Lin CS, Tsao TP, et al. Detecting digoxin toxicity by artificial intelligence-assisted electrocardiography. *Int J Environ Res Public Health.* 2021;Apr6;18(7). ([Journal full text](#)) <https://doi.org/10.3390/ijerph18073839>

⁸ Li MH, Mestre TA, Fox SH, Taati B. Vision-based assessment of parkinsonism and levodopa-induced dyskinesia with pose estimation. *J Neuro Engineering Rehabil.* 2018;Dec15:1-3. ([Journal full text](#)) [HYPERLINK "https://doi.org/10.1186/s12984-018-0446-z"](https://doi.org/10.1186/s12984-018-0446-z)<https://doi.org/10.1186/s12984-018-0446-z>

⁹ Practical Aspects of Signal Detection in Pharmacovigilance. CIOMS Working Group VIII report. Council for International Organizations of Medical Sciences (CIOMS), 2010.

¹⁰ U.S. Food and Drug Administration. Good Pharmacovigilance Practices and Pharmacoepidemiologic Assessment. Guidance for Industry. March 2005. ([Full text](#) accessed 21 March 2025).

¹¹ *Practical Aspects of Signal Detection in Pharmacovigilance: Report of CIOMS Working Group VIII 2010*

¹² *Practical Aspects of Signal Detection in Pharmacovigilance. CIOMS Working Group VIII report. Council for International Organizations of Medical Sciences (CIOMS), 2010.*

¹³ V. Kěpuska and G. Bohouta, "Next-generation of virtual personal assistants (Microsoft Cortana, Apple Siri, Amazon Alexa and Google Home)," 2018 IEEE 8th Annual Computing and Communication Workshop and Conference (CCWC), Las Vegas, NV, USA, 2018, pp. 99-103, doi: 10.1109/CCWC.2018.8301638.

¹⁴ Alvager T, Smith TJ, Vijai F. Neural-network applications for analysis of adverse drug reactions. *Biomed Instrum Technol.* 1993 Sep-Oct;27(5):408-11. PMID: 8220635.

¹⁵ U.S. Food and Drug Administration. FDA Adverse Event Reporting System (FAERS) Public Dashboard. ([Webpage](#) accessed 21 March 2025)

¹⁶ Painter JL, Chalamalasetti VR, Kassekert R, Bate A. Automating pharmacovigilance evidence generation: using large language models to produce context-aware structured query language. *JAMIA open.* 2025;Feb;8(1):ooaf003. ([Journal full text](#)) <https://doi.org/10.1093/jamiaopen/ooaf003>

¹⁷ Painter JL, Mahaux O, Vanini M, Kara V, Roshan C, Karwowski M, Chalamalasetti VR, Bate A. Enhancing drug safety documentation search capabilities with large language models: a user-centric approach. In 2023 International Conference on Computational Science and Computational Intelligence (CSCI). 2023;Dec13:49-56. ([Journal abstract](#)) <https://doi:10.1109/CSCI62032.2023.00015>

401

402

403 Chapter 2: Landscape analysis

404 Use of artificial intelligence in pharmacovigilance to date

405 AI may directly or indirectly impact all aspects of PV (see [Figure 1: Traditional signal management](#)
406 [process](#)). In this chapter, we discuss solutions that incorporate elements of AI and have been
407 developed or deployed for a variety of tasks across PV, focusing on those that have been
408 implemented specifically for PV or have accounted for attributes or features especially prominent in
409 PV applications. For example, AI solutions for general translation tasks are out of scope, but PV
410 specific translations, e.g. of adverse event reports, are in scope. Further, the landscape analysis
411 reflects the overall scope of the document, which focuses on collection, processing, and analysis of
412 adverse event reports. For this reason, research on AI methods to identify covariates for inclusion in
413 propensity score models for epidemiological studies are out of scope. Rather than seeking to provide
414 an exhaustive enumeration, the aim here is to illustrate the range and variety of current applications.
415 Additional examples can be found in recent review articles.¹⁸ The reader is also referred to the many
416 perspectives and commentaries that discuss the use of AI in PV^{19,20,21,22} and the cautionary notes that
417 have been provided.²³

418 Adverse event reporting and capture

419 AI solutions have been proposed for a variety of tasks related to natural language processing of social
420 media content to identify references to (personal experiences of) medicine use and adverse events.
421 These tasks include identifying relevant posts,^{24,25} identifying relevant parts of such posts,²⁶
422 normalizing descriptions of adverse events or medicinal products within such posts to standardized
423 terminologies like MedDRA or ATC,²⁷ and classifying the relationship between adverse events and
424 drugs mentioned in the same posts.²⁸ Similarly, screening the scientific literature for adverse events
425 is an investigated AI application.²⁹

426 Individual Case Safety Report Processing

427 An area of ICSR processing where AI solutions have been in routine use by some organizations since
428 at least the 2010's is duplicate detection, which relates to the identification of multiple unlinked
429 records describing the same adverse event in a particular patient.³⁰ Duplicate detection methods
430 based on ML and probabilistic record linkage have been implemented for VigiBase,³¹ FAERS,³² and
431 EudraVigilance.³³ The use of natural language processing to improve duplicate detection by
432 extracting and incorporating information from free text has also been explored.^{34,35} Rule-based
433 methods are more widely used and would be easier to implement but do not perform as well.³⁶

434 Another area where AI has been used to support ICSR processing is in the encoding of information on
435 adverse events^{37,38} or medicinal products³⁹ in standard terminologies based on verbatim fields and /
436 or free-flowing case narratives. Natural language processing has also been applied to extract relevant
437 information from case narratives and map it to structured fields^{40,41,42,43,44,45} and for ICSR
438 translation.⁴⁶ A major challenge is the lack of the data homogeneity⁴⁷ that could be improved by
439 adhering to existing common data models and standards that have already demonstrated value in
440 the PV domain, such as the Observational Medical Outcomes Partnership (OMOP) Common Data
441 Model for longitudinal observational health data⁴⁸ and Fast Healthcare Interoperability Resources
442 (FHIR), a standard for health care data exchange, published by HL7®, especially if they are extended
443 to support a broader spectrum of PV cases.^{49,50}

444 Several organisations who process large numbers of case reports have also automated repetitive,
445 labour-intensive tasks using so-called robotic process automation (RPA) technologies.⁵¹ These
446 operate on the user interface of other computer systems like humans would.⁵²

447 Other applications of AI solutions during ICSR case processing include methods have been developed
448 to help support triage incoming reports for human review,^{53,54} individual case causality assessment,⁵⁵
449 and automated redaction of person names in case narratives.⁵⁶

450 **Signal detection and analysis**

451 The earliest examples of real-world use of (simple) AI solutions in PV are from the late 1990s. At this
452 point, disproportionality analysis, first conceptualized in the 1970s,⁵⁷ began to be implemented as
453 part of triage algorithms to help direct the attention of PV specialists in their analysis of large
454 national and international collections of individual case reports.^{58,59,60,61} Since then, various
455 incremental improvements have been introduced and evaluated including automated adjustment for
456 confounding through e.g. regression,^{62,63} or propensity scores,⁶⁴ extensions to drug-drug
457 interactions,^{65,66,67,68,69,70} and other possible risk factors for adverse reactions.^{71,72} Methods to detect
458 adverse events associated with the production process or with substandard or counterfeit medicines
459 have also be explored.^{73,74,75} In addition, there have been efforts to develop predictive models for
460 statistical signal detection that account for other aspects of a case series, such as its geographic
461 spread and the quality and content of individual reports,⁷⁶ the time-to-onset of the reported
462 reactions,⁷⁷ or a combination of e.g. Naranjo scores and the proportions of reports on a drug-adverse
463 event combination coming from healthcare professionals and marketing authorization holders,
464 respectively.⁷⁸

465 Natural language processing has been applied to mine regulatory information,⁷⁹ scientific literature,
466 and clinical notes^{80,81,82,83} for information on already known/unknown and potentially serious adverse
467 effects. This may support and streamline decision making, especially during early signal assessment
468 and prioritization.

469 Some published AI-based signal detection exercises provide tantalizing glimpses of how elegant AI
470 solutions may uncover truly novel adverse events.⁸⁴ At the same caution is warranted in that highly
471 technical and elegant methods may be associated with overly-optimistic interpretations of, and
472 corresponding messaging about, the results, which may disseminate widely.⁸⁵

473 Several organizations have developed predictive models for ICSR prioritization to assess causality
474 associations between drugs and adverse events⁸⁶ and/or inform a regulatory action.⁸⁷ These can be
475 used to prioritize reports for human review during signal assessment and/or case processing.
476 Semantic search has been developed for case narratives to support signal detection and
477 assessment^{88,89} and there have been efforts to provide machine learning-based decision support for
478 signal validation⁹⁰ and to automatically visualize relevant information on case reports to facilitate
479 human review during signal assessment.⁹¹ Machine learning has been used to help estimate the
480 proportion of patients with a genotype associated with drug toxicity based on the phenotypical
481 manifestations reported in ICSRs.⁹²

482 Applications of unsupervised learning have been developed to support signal detection and analysis,
483 especially seeking to bring together reports describing similar or related adverse events. These
484 include network analyses of adverse events (and to a lesser extent drugs)^{93,94,95,96,97} cluster analysis of
485 adverse event reports,⁹⁸ and data-driven derivations of semantic representations of adverse events
486 and drugs.^{99,100}

487 Datasets with information about drug side effects and indications such as DrugBank¹⁰¹ and SIDER,¹⁰²
488 as well as those with information on pharmacology and chemical structures such as Bio2RDF,¹⁰³ have
489 been leveraged to enhance PV signal detection and analysis,^{104,105} or derive knowledge graphs that
490 can serve as downstream inputs for AI-based predictive signal detection.¹⁰⁶

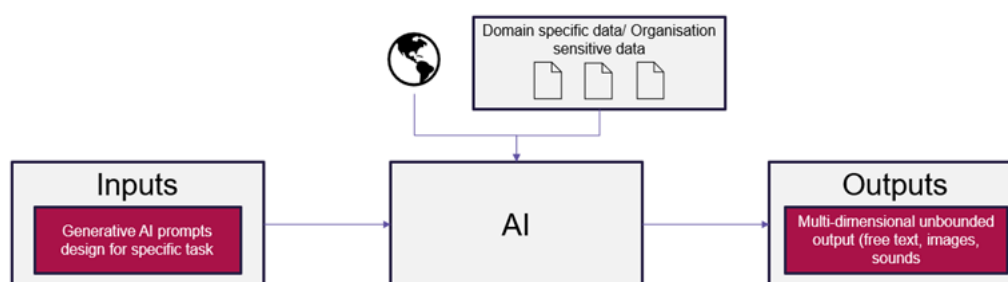
491 **Early applications of generative AI in pharmacovigilance**

492 Early applications of generative AI in PV have started to be explored. Researchers have applied large
 493 language models to a variety of PV tasks and reported their experiences. Examples to date include
 494 use of LLMs to simplify the patient communication from a regulatory authority,¹⁰⁷ summarization for
 495 drug labelling documents,¹⁰⁸ named entity recognition in scientific literature and social media,¹⁰⁹
 496 search of drug safety documentation,¹¹⁰ Q&A for drug labelling,¹¹¹ PV context-aware generation of
 497 SQL code¹¹², and drafting follow-up letters to reporters.¹¹³

498 However, the landscape is evolving beyond this framework thanks to emerging technologies like
 499 GenAI. These models, developed from deep neural networks with up to hundreds of billions of
 500 network parameters provided with vast and opaque large text corpuses, give rise to boundless
 501 number of multi-modal inputs and outputs (see Figure 4), therefore deciphering useful from
 502 misleading outputs for the end user can become a challenge. The non-deterministic and ‘black box’
 503 nature of such algorithms as well as the lack of ability to fully understand training can make
 504 developing and maintaining trust potentially harder as well as the ease of communicating such to
 505 external parties. As noise and biases may infiltrate the results, leading to potentially unreliable and
 506 misleading conclusions, necessitating guidance on how these technologies should be used
 507 appropriately in PV is needed where societal expectations are greater.

508 **Figure 4: Multi- modal input models**

509 Source: Variation of submitted article to Therapeutic Advances in Drug Safety



510

511 **Examples of deployed AI solutions**

512 Much of the research and development of AI solutions for PV to date has been experimental, with
 513 either no real-world deployment yet or only limited experimental use, for example in the form of
 514 pilot studies. However, **Error! Reference source not found.** presents examples of AI solutions that
 515 have been adopted for routine use in PV by various PV organizations and are described in the public
 516 domain. The deployment of AI solutions by pharmaceutical companies may on the other hand to a
 517 large extent be based on software vendor implementations, which are not described in the public
 518 domain.

519 **Table 1: Examples of deployed artificial intelligence solutions in pharmacovigilance**
 520 **described in the public domain**

521 Source: CIOMS XIV working group

522

AI solution	Pharmacovigilance context / database
Automated coding of medicinal products	VigiBase ¹¹⁴
Duplicate detection	FAERS, ¹¹⁵ VigiBase ¹¹⁶
Automated triages of individual case reports	Swedish Medical Products Agency ¹¹⁷ , pharmaceutical companies ¹¹⁸
Automated triages for quantitative signal detection	Databases of various regulatory authorities, international organizations, and pharmaceutical companies

Predictive models for quantitative signal detection	VigiBase, ^{119,120} Netherlands pharmacovigilance centre Lareb ¹²¹
Adverse event cluster analysis for signal detection and assessment	VigiBase ^{122,123}
Literature surveillance for safety data	EudraVigilance Netherlands pharmacovigilance centre Lareb ¹²⁴

523 Regulatory considerations

524 Introduction

525 Since 2017, countries around the world have been developing national AI strategies in order to adapt
 526 to technological advancements and their impact on society and the economy (OECD).¹²⁵ Countries
 527 have developed different regulatory frameworks and guiding principles to ensure the ethical use and
 528 trustworthiness of AI systems, and legislation of AI are being implemented (i.e. EU AI Act, AIDA, US
 529 Algorithmic Accountability Act and Executive Order: Promoting the Use of Trustworthy Artificial
 530 Intelligence in the Federal Government) (OECD).¹²⁶ In addition, there have been published reflection
 531 and discussion papers on the use of AI in medicinal products by the EMA and FDA, as well as a draft
 532 guidance on AI use to support regulatory decision making for drug products by the FDA.

533

534 *Guiding Principles for AI in Pharmacovigilance*

535 There are numerous published guiding principles for safe and responsible use of AI by governments,
 536 regulatory bodies and international organisations such as the WHO and OECD, that have been
 537 reviewed by the CIOMS Working Group XIV. These publications all define guiding principles and
 538 recommend best practices for safe and responsible AI use in regulated fields; however, the majority
 539 of these publications were not developed specifically for PV. Furthermore, it should be
 540 acknowledged that some discretion was used to establish the guiding principles by the various
 541 organizations, as some of the principles were described in conjunction with other principles.
 542 Nonetheless, these principles can be applied to the field of PV. The table below provides an overall
 543 comparison of the guiding principles, and a non-exhaustive description of the principles is presented
 544 in [Appendix 2](#):

545

546 **Table 2: Comparison of CIOMS Working Group XIV guiding principles for artificial**
 547 **intelligence across regional and country government institutions, and international**
 548 **organizations**

549 Source: CIOMS Working Group XIV

550

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{127,128}	Australia ¹²⁹	Canada ¹³⁰	Singapore ¹³¹	UK ¹³²	US ¹³³	PAHO ¹³⁴	WHO ¹³⁵	OECD ¹³⁶
Human Oversight	✓	✓	✓		✓	✓	✓	✓	✓
Validity & Robustness	✓	✓	✓		✓		✓	✓	✓
Data Privacy	✓	✓			✓	✓	✓		
Transparency	✓	✓	✓	✓	✓		✓	✓	✓
Accountability	✓	✓	✓	✓	✓	✓	✓	✓	✓
Societal well-being	✓	✓		✓			✓	✓	✓
Environmental Well-being	✓	✓						✓	✓
Fairness & Equity	✓	✓	✓	✓	✓	✓	✓	✓	✓
Explainability	✓	✓		✓	✓	✓		✓	✓
Safety	✓	✓	✓		✓	✓		✓	✓
Governance	✓				✓			✓	

551

552 *EMA Reflection Paper on the Use of Artificial Intelligence (AI) in the Medicinal Product Lifecycle*

553 On September 9, 2024, the European Medicines Agency finalized its Reflection paper on the use of AI
 554 in the medicinal product lifecycle.¹³⁷ The reflection paper addresses the use of AI/ML in the safe and
 555 effective development, manufacturing and use of medicines.

556 EMA advocates a risk-based approach for the development, deployment and monitoring of AI and
 557 ML tools throughout the system lifecycle. The paper uses the terms 'high patient risk' for systems
 558 affecting patient safety and 'high regulatory impact' for cases with a substantial impact on regulatory
 559 decision making. It is expected that applicants/marketing authorisation holders (MAHs) and
 560 developers of AI and ML systems will perform a regulatory impact and risk analysis. The level of
 561 scrutiny of the AI and ML systems will be dependent on the assessment of risk level and regulatory
 562 impact.

563 The paper provides technical and regulatory considerations on the use of AI and ML throughout the
 564 lifecycle of medicinal products, from drug discovery and development to post-authorisation settings.
 565 Specifically for PV, the paper foresees that AI/ML tools can effectively support activities such as
 566 adverse event report management and signal detection, in line with applicable GVP requirements.
 567 Applications within PV may allow a more flexible approach to AI/ML modelling and deployment than
 568 other domains, for example, to improve severity scoring of adverse event reports and signal
 569 detection. It is, however, the responsibility of the MAH to validate, monitor and document model
 570 performance and include AI/ML operations in the PV system, to mitigate risks related to all
 571 algorithms and models used.

572 Generally, the applicant or MAH is responsible for ensuring that all elements of the AI and ML
573 applications (i.e. algorithms, models, datasets, and data processing pipelines) are fit for purpose and
574 comply with GxP standards and current EMA scientific guidelines. Member State data protection
575 authorities are responsible for the supervision and monitoring of data protection compliance of
576 AI systems. Applicants or MAHs and developers are recommended to engage with EMA on
577 experimental technology, especially for AI and ML models that may have a high impact on the
578 regulatory decision making.¹³⁸

579 The EMA is planning to develop further guidance on the use of AI in the medicines lifecycle,
580 including in PV.¹³⁹

581 *FDA Discussion Paper on Using artificial Intelligence & Machine Learning in the Development of Drug*
582 *& Biological Products*

583 In May 2023, the US FDA published a discussion paper on “Using artificial Intelligence & Machine
584 Learning in the Development of Drug & Biological Products”.¹⁴⁰ The FDA acknowledges the increased
585 use of AI/ML in the lifecycle of drug development with novel approaches in data mining, analyzing
586 large multi-omics, PK/PD modeling, real world data, data collection from wearable devices and other
587 datasets (e.g. *in vitro* and *in vivo* studies, mechanistic studies, and multi-organ chip systems. In the
588 post-marketing safety surveillance, the FDA sees the potential to: i) automate the processing and
589 prioritization of individual case safety report (ICSR) using AI/ML, due to the increasing volume of
590 reports and complexity of data sources; ii) classifying ICSRs on the likelihood of causal relationship
591 between the drug and adverse event; iii) determine the seriousness of the outcome of ICSRs; and iv)
592 automate aggregate reports for multiple adverse events for a particular product.

593 *FDA Draft Guidance on Considerations for the Use of Artificial Intelligence to Support Regulatory*
594 *Decision-Making for Drug and Biological Products*

595 The FDA published a draft guidance titled “Considerations for the Use of Artificial Intelligence to
596 Support Regulatory Decision-Making for Drug and Biological Products Guidance for Industry and
597 Other Interested Parties” in January 2025. It elaborates a risk-based credibility assessment
598 framework for AI. The scope of the document focuses on the support of regulatory decision making
599 pertaining to the safety, effectiveness, or quality for drugs. Out of scope are drug discovery or
600 scenarios in which AI is deployed for operational efficiencies. The draft guidance is not a regulatory
601 mandate, rather the FDA’s current thinking and recommendations on AI use. It considers AI broadly,
602 i.e. not limited to specific subsets of AI such as machine learning. There are three major segments of
603 the draft guidance:

- 604 1. Establishing a risk-based credibility assessment framework (see also Chapter on [Risk-](#)
605 [based approach](#));
- 606 2. Lifecycle credibility maintenance;
- 607 3. Options for sponsors for engaging with the agency to discuss AI model development.

608 The risk-based credibility assessment framework has seven steps as below.

- 609 1. Define the question to be addressed by an AI model.
- 610 2. Define the “context of use (COU)” defined as “...the specific role and scope of the AI
611 model used to address a question of interest.” Importantly this includes whether the
612 questions being answered, and any ensuing classifications or decisions, are based solely
613 on the AI outputs versus the AI being used in conjunction with other information (i.e.
614 “model influence”). This is important because it helps define the associated risk in the
615 subsequent step.
- 616 3. Define model risk. This is determined by model influence as defined in COU and decision
617 consequence - i.e. the consequences of an incorrect decision. The risk is highest when

- 618 the AI model is operating in a stand-alone capacity and incorrect decisions present a
 619 major hazard. The required level of oversight throughout the development and
 620 production cycle is positively correlated with the risk.
- 621 4. Develop a plan to establish AI model credibility within the COU.
 - 622 5. Execution of the plan.
 - 623 6. Document the results of the credibility assessment plan and discuss deviations from the
 624 plan.
 - 625 7. Determine the Adequacy of the AI Model for the COU.¹⁴¹

626 *FDA Emerging Drug Safety Technology Program (EDSTP)*

627 The FDA has established the Emerging Drug Safety Technology Program (EDSTP) in June 2024 to
 628 engage with industry stakeholders on AI and other emerging novel technologies used in PV and the
 629 lifecycle of the drug product. The three goals of the EDSTP include discussion between industry and
 630 FDA, knowledge dissemination of emerging AI/ML models or other emerging novel technologies, and
 631 to inform potential regulatory or policy development within the context of PV.¹⁴²

632 *Guidance on use of Large Language Models*

633 Since the release of ChatGPT on November 30, 2022, there has been significant work in exploring
 634 how generative AI could be adapted to a variety of tasks (such as text and image generation, coding,
 635 brainstorming, and research) for productivity gains. Given the potential use and broad applicability of
 636 GenAI, regulatory agencies and organizations have developed high level guides and best practices on
 637 the safe and responsible use of GenAI by their own staff and broader stakeholder groups,
 638 respectively, which aligns with established guiding principles for AI:

- 639 • Guiding principles on the use of large language models in regulatory science and for
 640 medicines regulatory activities (EMA¹⁴³);
- 641 • Guide on the use of generative artificial intelligence (Canada¹⁴⁴);
- 642 • Initial policy considerations for generative artificial intelligence (OECD¹⁴⁵);
- 643 • WHO Ethics and governance of artificial intelligence for health: Guidance on large
 644 multi-modal models. (WHO¹⁴⁶).

645 *Guidelines for safe AI*

646 Other related regulatory and international organization (e.g. WHO and OECD) published guidelines
 647 for safe AI include:

- 648 • Regulatory considerations on artificial intelligence for health, WHO 2023;¹⁴⁷
- 649 • Ethics guidelines for trustworthy AI, European Commission 2019;¹⁴⁸
- 650 • Recommendation of the Council on Artificial Intelligence, OECD 2019, amended
 651 2023;¹⁴⁹
- 652 • Good machine learning practice for medical device development: Guiding Principles,
 653 U.S. Food & Drug Administration FDA, Health Canada, Medicines & Healthcare
 654 products Regulatory Agency MHRA 2021.¹⁵⁰

655

References

¹⁸ Salas M, Petracek J, Yalamanchili P, Aimer O, Kasthuril D, Dhingra S, Junaid T, Bostic T. The Use of Artificial Intelligence in Pharmacovigilance: A Systematic Review of the Literature. *Pharm Med* 36, 2022;Oct;36(5):295–306. [\[Journal abstract\] https://doi.org/10.1007/s40290-022-00441-z](https://doi.org/10.1007/s40290-022-00441-z)

- ¹⁹ Bate A, Stegmann JU. Artificial intelligence and pharmacovigilance: What is happening, what could happen and what should happen?. *Health Policy and Technology*. 2023;Jun1;12(2):100743. ([Journal full text](#)) <https://doi.org/10.1016/j.hlpt.2023.100743>
- ²⁰ Basile AO, Yahi A, Tatonetti NP. Artificial intelligence for drug toxicity and safety. *Trends in Pharmacological Sciences*. 2019;Sep1;40(9):624-635. ([Journal full text](#)) <https://doi.org/10.1016/j.tips.2019.07.005>
- ²¹ Painter JL, Kassekert R, Bate A. Corrigendum: An industry perspective on the use of machine learning in drug and vaccine safety. *Frontiers in Drug Safety and Regulation*. 2023;Dec4;3:1244115. ([Journal full text](#)) <https://doi.org/10.3389/fdsfr.2023.1110498>
- ²² Kassekert R, Grabowski N, Lorenz D, Schaffer C, Kempf D, Roy P, Kjoersvik O, Saldana G, ElShal S. Industry perspective on artificial intelligence/machine learning in pharmacovigilance. *Drug Safety*. 2022;May;45(5):439-448. ([Journal full text](#)) <https://doi.org/10.1007/s40264-022-01164-5>
- ²³ Hauben M, Reynolds R, Caubel P. Deconstructing the Pharmacovigilance Hype Cycle. *Clin Ther*. 2018;Dec;40(12):1981-1990.e3. ([Journal full text](#)) doi:10.1016/j.clinthera.2018.10.021.
- ²⁴ Freifeld CC, Brownstein JS, Menone CM, Bao W, Filice R, Kass-Hout T, Dasgupta N. Digital drug safety surveillance: monitoring pharmaceutical products in twitter. *Drug Safety*. 2014;May;37:343-50. ([Journal full text](#)) <https://doi.org/10.1007/s40264-014-0155-x>
- ²⁵ Alvaro N, Conway M, Doan S, Lofi C, Overington J, Collier N. Crowdsourcing Twitter annotations to identify first-hand experiences of prescription drug use. *Journal of Biomedical Informatics*. 2015;Dec1;58:280-287. ([Journal full text](#)) <https://doi.org/10.1016/j.jbi.2015.11.004>
- ²⁶ Nikfarjam A, Sarker A, O'connor K, Ginn R, Gonzalez G. Pharmacovigilance from social media: mining adverse drug reaction mentions using sequence labeling with word embedding cluster features. *Journal of the American Medical Informatics Association*. 2015;May1;22(3):671-681. ([Journal full text](#)) <https://doi.org/10.1093/jamia/ocu041>
- ²⁷ Karimi S, Metke-Jimenez A, Nguyen A. CADEminer: a system for mining consumer reports on adverse drug side effects. In *Proceedings of the eighth workshop on exploiting semantic annotations in information retrieval*. 2015;Oct22;47-50. ([Journal abstract](#)) <https://doi.org/10.1145/2810133.2810143>
- ²⁸ White RW, Harpaz R, Shah NH, DuMouchel W, Horvitz E. Toward enhanced pharmacovigilance using patient-generated data on the internet. *Clinical Pharmacology & Therapeutics*. 2014;Aug;96(2):239-246. ([Journal full text](#)) <https://ascpt.onlinelibrary.wiley.com/doi/full/10.1038/clpt.2014.77>
- ²⁹ Park J, Djelassi M, Chima D, Hernandez R, Poroshin V, Iliescu AM, Domalik D, Southall N. Validation of a natural language machine learning model for safety literature surveillance. *Drug Safety*. 2024;Jan;47(1):71-80. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-023-01367-4>
- ³⁰ Tregunno PM, Fink DB, Fernandez-Fernandez C, Lázaro-Bengoa E, Norén GN. Performance of probabilistic method to detect duplicate individual case safety reports. *Drug Safety*. 2014;Apr;37:249-258. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-014-0146-y>
- ³¹ Norén GN, Orre R, Bate A, Edwards IR. Duplicate detection in adverse drug reaction surveillance. *Data Mining and Knowledge Discovery*. 2007;Jun;14:305-328. ([Journal abstract](#)) <https://doi.org/10.1007/s10618-006-0052-8>
- ³² Kreimeyer K, Menschik D, Winięcki S, Paul W, Barash F, Woo EJ, Alimchandani M, Arya D, Zinderman C, Forshee R, Botsis T. Using probabilistic record linkage of structured and unstructured data to identify duplicate cases in spontaneous adverse event reporting systems. *Drug Safety*. 2017;Jul;40:571-582. ([Journal abstract](#))
- ³³ Personal communication from Tom Paternoster-Howe (March 2025).
- ³⁴ Kreimeyer K, Menschik D, Winięcki S, Paul W, Barash F, Woo EJ, Alimchandani M, Arya D, Zinderman C, Forshee R, Botsis T. Using probabilistic record linkage of structured and unstructured data to identify duplicate cases in spontaneous adverse event reporting systems. *Drug Safety*. 2017;Jul;40:571-582. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-017-0523-4>
- ³⁵ Kreimeyer K, Dang O, Spiker J, Gish P, Weintraub J, Wu E, Ball R, Botsis T. Increased confidence in deduplication of drug safety reports with natural language processing of narratives at the us food and drug administration. *Frontiers in Drug Safety and Regulation*. 2022;Jun15;2:918897. ([Journal full text](#)) <https://doi.org/10.3389/fdsfr.2022.918897>
- ³⁶ Tregunno PM, Fink DB, Fernandez-Fernandez C, Lázaro-Bengoa E, Norén GN. Performance of probabilistic method to detect duplicate individual case safety reports. *Drug Safety*. 2014;Apr;37:249-258. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-014-0146-y>
- ³⁷ Combi C, Corzi M, Pozzani G, Moretti U, Arzenton E. From narrative descriptions to MedDRA: automagically encoding adverse drug reactions. *Journal of Biomedical Informatics*. 2018;Aug1;84:184-199. ([Journal full text](#)) <https://doi.org/10.1016/j.jbi.2018.07.001>
- ³⁸ Gurulingappa H, Mateen-Rajpu A, Toldo L. Extraction of potential adverse drug events from medical case reports. *Journal of Biomedical Semantics*. 2012;Dec;3:1-10. ([Journal full text](#)) <https://doi.org/10.1186/2041-1480-3-15>

- ³⁹ Meldau EL, Bista S, Rofors E, Gattepaille LM. Automated Drug Coding Using Artificial Intelligence: An Evaluation of WHODrug Koda on Adverse Event Reports. *Drug Safety*. 2022;May;45(5):549-561. ([Journal full text](#)) <https://doi.org/10.1007/s40264-022-01162-7>
- ⁴⁰ Abatemarco D, Perera S, Bao SH, Desai S, Assuncao B, Tetarenko N. Training augmented intelligent capabilities for pharmacovigilance: applying deep-learning approaches to individual case safety report processing. *Pharmaceut Med*. 2018;32 (6):391–401. ([Journal full text](#)) <https://doi.org/10.1007/s40290-018-0251-9>
- ⁴¹ Wang W, Kreimeyer K, Woo EJ, Ball R, Foster M, Pandey A, Scott J, Botsis T. A new algorithmic approach for the extraction of temporal associations from clinical narratives with an application to medical product safety surveillance reports. *Journal of Biomedical Informatics*. 2016;Aug1;62:78-89. ([Journal full text](#)) <https://doi.org/10.1016/j.jbi.2016.06.006>
- ⁴² Kreimeyer K, Foster M, Pandey A, Arya N, Halford G, Jones SF, Forshee R, Walderhaug M, Botsis T. Natural language processing systems for capturing and standardizing unstructured clinical information: a systematic review. *Journal of Biomedical Informatics*. 2017;Sep1;73:14-29. ([Journal full text](#)) <https://doi.org/10.1016/j.jbi.2017.07.012>
- ⁴³ Pham P, Cheng C, Wu E, Kim I, Zhang R, Ma Y, Kortepeter CM, Muñoz MA. Leveraging case narratives to enhance patient age ascertainment from adverse event reports. *Pharmaceutical Medicine*. 2021;Sep;35(5):307-316. <https://doi.org/10.1007/s40290-021-00398-5>
- ⁴⁴ Dang V, Wu E, Kortepeter CM, Phan M, Zhang R, Ma Y, Muñoz MA. Evaluation of a natural language processing tool for extracting gender, weight, ethnicity, and race in the US food and drug administration adverse event reporting system. *Frontiers in Drug Safety and Regulation*. 2022;Nov14;2:1020943. ([Journal full text](#)) <https://doi.org/10.3389/fdsfr.2022.1020943>
- ⁴⁵ Abatemarco, D., Perera, S., Bao, S.H. et al. Training Augmented Intelligent Capabilities for Pharmacovigilance: Applying Deep-learning Approaches to Individual Case Safety Report Processing. *Pharm Med* 32, 391–401 (2018). <https://doi.org/10.1007/s40290-018-0251-9>
- ⁴⁶ Römning H-J, Pushparajan R. AI Translation Assistant for Pharmacovigilance. Poster presented at DIA Europe 2021. (accessed 21 March 2025)
- ⁴⁷ Young IJ, Luz S, Lone N. A systematic review of natural language processing for classification tasks in the field of incident reporting and adverse event analysis. *International Journal of Medical Informatics*. 2019;Dec1;132:103971. ([Journal abstract](#)) <https://doi.org/10.1016/j.ijmedinf.2019.103971>
- ⁴⁸ OMOP Common Data Model for longitudinal observational health data <https://www.ohdsi.org/data-standardization/>
- ⁴⁹ Hauben M, Aronson JK, Ferner RE. Evidence of misclassification of drug–event associations classified as gold standard ‘negative controls’ by the observational medical outcomes partnership (OMOP). *Drug Safety*. 2016;May;39:421-432. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-016-0392-2>
- ⁵⁰ Wang X, Lehmann H, Botsis T. Can FHIR support standardization in post-market safety surveillance?. In *Public Health and Informatics*. 2021;(pp. 33-37). IOS Press. ([Journal full text](#)) <https://doi.org/10.3233/SHTI210115>
- ⁵¹ Ghosh, R., Kempf, D., Pufko, A. et al. Automation Opportunities in Pharmacovigilance: An Industry Survey. *Pharm Med* 34, 7–18 (2020). <https://doi.org/10.1007/s40290-019-00320-0>
- ⁵² van der Aalst, W.M.P., Bichler, M. & Heinzl, A. Robotic Process Automation. *Bus Inf Syst Eng* 60, 269–272 (2018). <https://doi.org/10.1007/s12599-018-0542-4>
- ⁵³ Gosselt HR, Bazelmans EA, Lieber T, van Hunsel FP, Härmark L. Development of a multivariate prediction model to identify individual case safety reports which require clinical review. *Pharmacoepidemiology and Drug Safety*. 2022;Dec;31(12):1300-1307. ([Journal full text](#)) <https://doi.org/10.1002/pds.5553>
- ⁵⁴ Bergman E, Dürlich L, Arthurson V, Sundström A, Larsson M, Bhuiyan S, Jakobsson A, Westman G. BERT based natural language processing for triage of adverse drug reaction reports shows close to human-level performance. *PLOS Digital Health*. 2023;Dec 6;2(12):e0000409. ([Journal full text](#)) <https://doi.org/10.1371/journal.pdig.0000409>
- ⁵⁵ Cherkas Y, Ide J, van Stekelenborg J. Leveraging machine learning to facilitate individual case causality assessment of adverse drug reactions. *Drug Safety*. 2022;May;45(5):571-582. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-022-01163-6>
- ⁵⁶ Meldau EL, Bista S, Melgarejo-González C, Norén GN. Automated redaction of names in adverse event reports using transformer-based neural networks. *BMC Medical Informatics and Decision Making*. 2024;Dec23;24(1):401. ([Journal full text](#)) <https://doi.org/10.1186/s12911-024-02785-9>
- ⁵⁷ Finney DJ. Systematic Signalling of Adverse Reactions to Drugs. *Methods of Information in Medicine*. 1974;13(01):01-10. ([Journal full text](#)) <https://doi.org/10.1055/s-0038-1636131>
- ⁵⁸ Practical Aspects of Signal Detection in Pharmacovigilance. CIOMS Working Group VIII report. Council for International Organizations of Medical Sciences (CIOMS), 2010.
- ⁵⁹ Bate A, Lindquist M, Edwards IR, Olsson S, Orre R, Lansner A, De Freitas RM. A Bayesian neural network method for adverse drug reaction signal generation. *European Journal of Clinical Pharmacology*. 1998;Jul;54:315-21. ([Journal abstract](#)) <https://doi.org/10.1007/s002280050466>

- ⁶⁰ DuMouchel W. Bayesian data mining in large frequency tables, with an application to the FDA spontaneous reporting system. *The American Statistician*. 1999;1;53(3):177-190. ([Journal full text](#)) <https://doi.org/10.1080/00031305.1999.10474456>
- ⁶¹ Evans SJ, Waller PC, Davis S. Use of proportional reporting ratios (PRRs) for signal generation from spontaneous adverse drug reaction reports. *Pharmacoepidemiology and Drug Safety*. 2001;Oct;10(6):483-486. ([Journal full text](#)) <https://doi.org/10.1002/pds.677>
- ⁶² Caster O, Norén GN, Madigan D, Bate A. Large-scale regression-based pattern discovery: the example of screening the WHO global drug safety database. *Statistical Analysis and Data Mining: The ASA Data Science Journal*. 2010;Aug;3(4):197-208. ([Journal abstract](#)) <https://doi.org/10.1002/sam.10078>
- ⁶³ Harpaz R, Perez H, Chase HS, Rabadan R, Hripcsak G, Friedman C. Biclustering of adverse drug events in the FDA's spontaneous reporting system. *Clinical Pharmacology & Therapeutics*. 2011;Feb;89(2):243-250.
- ⁶⁴ Tatonetti NP, Ye PP, Daneshjou R, Altman RB. Data-driven prediction of drug effects and interactions. *Science Translational Medicine*. 2012;Mar14;4(125):125ra31-. ([Journal abstract](#)) <https://doi:10.1126/scitranslmed.3003377>
- ⁶⁵ Thakrar BT, Grundschober SB, Doesseger L. Detecting signals of drug–drug interactions in a spontaneous reports database. *British Journal of Clinical Pharmacology*. 2007;Oct;64(4):489-495. ([Journal full text](#)) <https://doi.org/10.1111/j.1365-2125.2007.02900.x>
- ⁶⁶ Norén GN, Sundberg R, Bate A, Edwards IR. A statistical methodology for drug–drug interaction surveillance. *Statistics in Medicine*. 2008;Jul 20;27(16):3057-3070. ([Journal full text](#)) <https://doi.org/10.1002/sim.3247>
- ⁶⁷ Tatonetti NP, Ye PP, Daneshjou R, Altman RB. Data-driven prediction of drug effects and interactions. *Science Translational Medicine*. 2012;Mar14;4(125):125ra31-. ([Journal abstract](#)) <https://doi:10.1126/scitranslmed.3003377>
- ⁶⁸ Hauben M, Hung E, Chen Y. Potential signals of COVID-19 as an effect modifier of adverse drug reactions. *Clinical Therapeutics*. 2024;Jan1;46(1):20-29. ([Journal full text](#)) <https://doi.org/10.1016/j.clinthera.2023.10.002>
- ⁶⁹ Hauben M. Artificial intelligence and data mining for the pharmacovigilance of drug–drug interactions. *Clinical Therapeutics*. 2023;Feb1;45(2):117-133. ([Journal full text](#)) <https://doi.org/10.1016/j.clinthera.2023.01.002>
- ⁷⁰ van Puijenbroek, E., Egberts, A., Heerdink, E. et al. Detecting drug–drug interactions using a database for spontaneous adverse drug reactions: an example with diuretics and non-steroidal anti-inflammatory drugs. *Eur J Clin Pharmacol* 56, 733–738 (2000). <https://doi.org/10.1007/s002280000215>
- ⁷¹ Sandberg L, Taavola H, Aoki Y, Chandler R, Norén GN. Risk factor considerations in statistical signal detection: using subgroup disproportionality to uncover risk groups for adverse drug reactions in VigiBase. *Drug Safety*. 2020;Oct;43(10):999-1009. ([Journal full text](#)) <https://doi.org/10.1007/s40264-020-00957-w>
- ⁷² Hauben M, Hung E, Chen Y. Potential signals of COVID-19 as an effect modifier of adverse drug reactions. *Clinical Therapeutics*. 2024;Jan1;46(1):20-29. ([Journal full text](#)) <https://doi.org/10.1016/j.clinthera.2023.10.002>
- ⁷³ Juhlin, K., Karimi, G., Andér, M. et al. Using VigiBase to Identify Substandard Medicines: Detection Capacity and Key Prerequisites. *Drug Saf* 38, 373–382 (2015). <https://doi.org/10.1007/s40264-015-0271-2>
- ⁷⁴ Trippe, Z.A., Brendani, B., Meier, C. et al. Identification of Substandard Medicines via Disproportionality Analysis of Individual Case Safety Reports. *Drug Saf* 40, 293–303 (2017). <https://doi.org/10.1007/s40264-016-0499-5>
- ⁷⁵ Olivia Mahaux, Vincent Bauchau, Ziad Zeinoun, Lionel Van Holle, Tree-based scan statistic – Application in manufacturing-related safety signal detection, *Vaccine*, Volume 37, Issue 1, 2019, Pages 49-55, ISSN 0264-410X, <https://doi.org/10.1016/j.vaccine.2018.11.044>.
- ⁷⁶ Caster O, Juhlin K, Watson S, Norén GN. Improved statistical signal detection in pharmacovigilance by combining multiple strength-of-evidence aspects in *vigiRank*: retrospective evaluation against emerging safety signals. *Drug Safety*. 2014;Aug37:617-628. ([Journal full text](#)) <https://doi.org/10.1007/s40264-014-0204-5>
- ⁷⁷ Van Holle L, Bauchau V. Use of logistic regression to combine two causality criteria for signal detection in vaccine spontaneous report data. *Drug Safety*. 2014;Dec;37:1047-1057. ([Journal full text](#)) <https://doi.org/10.1007/s40264-014-0237-9>
- ⁷⁸ Scholl JH, van Hunsel FP, Hak E, van Puijenbroek EP. A prediction model-based algorithm for computer-assisted database screening of adverse drug reactions in the Netherlands. *Pharmacoepidemiology and drug safety*. 2018;Feb;27(2):199-205. ([Journal full text](#)) <https://doi.org/10.1002/pds.4364>
- ⁷⁹ Ly T, Pamer C, Dang O, Brajovic S, Haider S, Botsis T, Milward D, Winter A, Lu S, Ball R. Evaluation of Natural Language Processing (NLP) systems to annotate drug product labeling with MedDRA terminology. *Journal of Biomedical Informatics*. 2018;Jul1;83:73-86. ([Journal full text](#)) <https://doi.org/10.1016/j.jbi.2018.05.019>
- ⁸⁰ Avillach P, Dufour JC, Diallo G, Salvo F, Joubert M, Thiessard F, Mouglin F, Trifirò G, Fourrier-Réglat A, Pariente A, Fieschi M. Design and validation of an automated method to detect known adverse drug reactions in MEDLINE: a contribution from the EU–ADR project. *Journal of the American Medical Informatics Association*. 2013;May1;20(3):446-452. ([Journal abstract](#)) <https://doi.org/10.1136/amiainl-2012-001083>

- ⁸¹ Shetty KD, Dalal SR. Using information mining of the medical literature to improve drug safety. *Journal of the American Medical Informatics Association*. 2011;Sep1;18(5):668-674. ([Journal abstract](#)) <https://doi.org/10.1136/amiainl-2011-000096>
- ⁸² Silverman AL, Sushil M, Bhasuran B, Ludwig D, Buchanan J, Racz R, Parakala M, El-Kamary S, Ahima O, Belov A, Choi L. Algorithmic Identification of Treatment-Emergent Adverse Events From Clinical Notes Using Large Language Models: A Pilot Study in Inflammatory Bowel Disease. *Clinical Pharmacology & Therapeutics*. 2024;Jun;115(6):1391-1399. ([Journal full text](#)) <https://doi.org/10.1002/cpt.3226>
- ⁸³ Wu J, Ruan X, McNeer E, Rossow KM, Choi L. Developing a natural language processing system using transformer-based models for adverse drug event detection in electronic health records. *medRxiv*. 2024;Jul10:2024-07. ([Journal full text](#)) <https://doi.org/10.1101/2024.07.09.24310100>
- ⁸⁴ Tatonetti NP, Denny JC, Murphy SN, Fernald GH, Krishnan G, Castro V, Yue P, Tsau PS, Kohane I, Roden DM, Altman RB. Detecting drug interactions from adverse-event reports: interaction between paroxetine and pravastatin increases blood glucose levels. *Clinical Pharmacology & Therapeutics*. 2011 Jul;90(1):133-142. Erratum in: *Clin Pharmacol Ther*. 2011;Sep;90(3):480. Tsau, P S (corrected to Tsao, P S). ([Journal full text](#))<https://doi.org/10.1038/clpt.2011.83>
- ⁸⁵ Hauben M, Reynolds R, Caubel P. Deconstructing the pharmacovigilance hype cycle. *Clinical Therapeutics*. 2018;Dec1;40(12):1981-1990. ([Journal full text](#)) <https://doi.org/10.1016/j.clinthera.2018.10.021>
- ⁸⁶ Kreimeyer K, Dang O, Spiker J, Muñoz MA, Rosner G, Ball R, Botsis T. Feature engineering and machine learning for causality assessment in pharmacovigilance: Lessons learned from application to the FDA Adverse Event Reporting System. *Computers in Biology and Medicine*. 2021;Aug1;135:104517. ([Journal abstract](#)) <https://doi.org/10.1016/j.compbiomed.2021.104517>
- ⁸⁷ Muñoz MA, Dal Pan GJ, Wei YJ, Delcher C, Xiao H, Kortepeter CM, Winterstein AG. Towards automating adverse event review: a prediction model for case report utility. *Drug Safety*. 2020;Apr;43:329-338. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-019-00897-0>
- ⁸⁸ Singh LL, Sudarsan SD, Jetley RP, Fitzgerald B, Milanova M. Semantic search tool for adverse event reports of medical devices. *InSWWS 2011:proceedings of the 2011 international conference on semantic web & web services*. 2011;July18-21:129-135. ([Journal full text](#)) <https://doi.org/10.13140/RG.2.2.22712.29443>
- ⁸⁹ Zekarias A, Meldau EL, Bista S, Félix China J, Sandberg L. Narrative Search Engine for Case Series Assessment Supported by Artificial Intelligence Query Suggestions. *Drug Safety*. 2025;Mar15:1-3. ([Journal full text](#)) <https://doi.org/10.1007/s40264-025-01529-6>
- ⁹⁰ Imran M, Bhatti A, King DM, Lerch M, Dietrich J, Doron G, Manlik K. Supervised machine learning-based decision support for signal validation classification. *Drug Safety*. 2022;May;45(5):583-596. ([Journal full text](#)) <https://doi.org/10.1007/s40264-022-01159-2>
- ⁹¹ Spiker J, Kreimeyer K, Dang O, Boxwell D, Chan V, Cheng C, Gish P, Lardieri A, Wu E, De S, Naidoo J. Information visualization platform for postmarket surveillance decision support. *Drug Safety*. 2020;Sep;43:905-915. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-020-00945-0>
- ⁹² Pinheiro LC, Durand J, Dogné JM. An Application of Machine Learning in Pharmacovigilance: Estimating Likely Patient Genotype from Phenotypical Manifestations of Fluoropyrimidine Toxicity. *Clin Pharmacol Ther*.2020;107(4):944-947. ([Journal full text](#)) <https://doi.org/10.1002/cpt.1789>
- ⁹³ Orre R, Bate A, Norén GN, Swahn E, Arnborg S, Edwards IR. A Bayesian recurrent neural network for unsupervised pattern recognition in large incomplete data sets. *International Journal of Neural Systems*. 2005;Jun;15(03):207-222. ([Journal abstract](#)) <https://doi.org/10.1142/S0129065705000219>
- ⁹⁴ Ball R, Botsis T. Can network analysis improve pattern recognition among adverse events following immunization reported to VAERS?. *Clinical Pharmacology & Therapeutics*. 2011;Aug;90(2):271-278. ([Journal abstract](#)) <https://doi.org/10.1038/clpt.2011.119>
- ⁹⁵ Harpaz R, Perez H, Chase HS, Rabadan R, Hripcsak G, Friedman C. Biclustering of adverse drug events in the FDA's spontaneous reporting system. *Clinical Pharmacology & Therapeutics*. 2011;Feb;89(2):243-250. ([Journal abstract](#)) <https://doi.org/10.1038/clpt.2010.285>
- ⁹⁶ Botsis T, Scott J, Goud R, Toman P, Sutherland A, Ball R. Novel algorithms for improved pattern recognition using the US FDA Adverse Event Network Analyzer. *Ine-Health—For Continuity of Care*. 2014;1178-1182. ([Journal full text](#)) <https://doi.org/10.3233/978-1-61499-432-9-1178>
- ⁹⁷ Fusaroli M, Polizzi S, Menestrina L, Giunchi V, Pellegrini L, Raschi E, Weintraub D, Recanatini M, Castellani G, De Ponti F, Poluzzi E. Unveiling the burden of drug-induced impulsivity: a network analysis of the FDA adverse event reporting system. *Drug Safety*. 2024;Aug15:1-8. ([Journal full text](#)) <https://doi.org/10.1007/s40264-024-01471-z>
- ⁹⁸ Norén GN, Meldau EL, Chandler RE. Consensus clustering for case series identification and adverse event profiles in pharmacovigilance. *Artificial Intelligence in Medicine*. 2021;Dec1;122:102199. ([Journal full text](#))[HYPERLINK https://doi.org/10.1016/j.artmed.2021.102199](https://doi.org/10.1016/j.artmed.2021.102199)"<https://doi.org/10.1016/j.artmed.2021.102199>

- ⁹⁹ Simms AM, Kanakia A, Sipra M, Dutta B, Southall N. A patient safety knowledge graph supporting vaccine product development. *BMC medical informatics and decision making*. 2024;Jan4;24(1):10. ([Journal full text](#)) <https://doi.org/10.1186/s12911-023-02409-8>
- ¹⁰⁰ Bean DM, Wu H, Iqbal E, Dzahini O, Ibrahim ZM, Broadbent M, Stewart R, Dobson RJ. Knowledge graph prediction of unknown adverse drug reactions and validation in electronic health records. *Scientific Reports*. 2017;Nov27;7(1):16416. ([Journal full text](#)) <https://doi.org/10.1038/s41598-017-16674-x>
- ¹⁰¹ Knox C, Wilson M, Klinger CM, Franklin M, Oler E, Wilson A, Pon A, Cox J, Chin NE, Strawbridge SA, Garcia-Patino M. DrugBank 6.0: the DrugBank knowledgebase for 2024. *Nucleic Acids Research*. 2024;Jan5;52(D1):D1265-1275. ([Journal full text](#)) <https://doi.org/10.1093/nar/gkad976>
- ¹⁰² Kuhn M, Letunic I, Jensen LJ, Bork P. The SIDER database of drugs and side effects. *Nucleic Acids Research*. 2016;Jan4;44(D1):D1075-2079. ([Journal full text](#)) <https://doi.org/10.1093/nar/gkv1075>
- ¹⁰³ Belleau F, Nolin MA, Tourigny N, Rigault P, Morissette J. Bio2RDF: towards a mashup to build bioinformatics knowledge systems. *Journal of Biomedical Informatics*. 2008;Oct1;41(5):706-716. ([Journal full text](#)) <https://doi.org/10.1016/j.jbi.2008.03.004>
- ¹⁰⁴ Zhang W, Liu F, Luo L, Zhang J. Predicting drug side effects by multi-label learning and ensemble learning. *BMC Bioinformatics*. 2015;Dec;16:1-1. ([Journal full text](#)) <https://doi.org/10.1186/s12859-015-0774-y>
- ¹⁰⁵ Low YS, Caster O, Bergvall T, Fourches D, Zang X, Norén GN, Rusyn I, Edwards R, Tropsha A. Cheminformatics-aided pharmacovigilance: application to Stevens-Johnson Syndrome. *Journal of the American Medical Informatics Association*. 2016;Sep1;23(5):968-78. ([Journal full text](#)) <https://doi.org/10.1093/jamia/ocv127>
- ¹⁰⁶ Hauben M, Rafi M, Abdelaziz I, Hassanzadeh O. Knowledge graphs in pharmacovigilance: a scoping review. *Clinical therapeutics*. 2024;Jul9. ([Journal full text](#)) <https://doi.org/10.1016/j.clinthera.2024.06.003>
- ¹⁰⁷ Sridharan K, Sivaramakrishnan G. Enhancing readability of USFDA patient communications through large language models: a proof-of-concept study. *Expert Review of Clinical Pharmacology*. 2024;Aug2;17(8):731-741. ([Journal abstract](#)) <https://doi.org/10.1080/17512433.2024.2363847>
- ¹⁰⁸ Ying L, Liu Z, Fang H, Kusko R, Wu L, Harris S, Tong W. Text summarization with ChatGPT for drug labeling documents. *Drug Discovery Today*. 2024;May7:104018. ([Journal full text](#)) <https://doi.org/10.1016/j.drudis.2024.104018>
- ¹⁰⁹ Dietrich J, Hollstein A. Performance and Reproducibility of Large Language Models in Named Entity Recognition: Considerations for the Use in Controlled Environments. *Drug Safety*. 2025;Mar;48(3):287-303. ([Journal full text](#)) <https://doi.org/10.1007/s40264-024-01499-1>
- ¹¹⁰ Painter JL, Mahaux O, Vanini M, Kara V, Roshan C, Karwowski M, Chalamalasetti VR, Bate A. Enhancing drug safety documentation search capabilities with large language models: a user-centric approach. *International Conference on Computational Science and Computational Intelligence (CSCI)*. 2023;Dec13:49-56. ([Journal abstract](#)) <http://doi.org/10.1109/CSCI62032.2023.00015>
- ¹¹¹ Wu L, Xu J, Thakkar S, Gray M, Qu Y, Li D, Tong W. A framework enabling LLMs into regulatory environment for transparency and trustworthiness and its application to drug labeling document. *Regulatory Toxicology and Pharmacology*. 2024;May1;149:105613. ([Journal full text](#)) <https://doi.org/10.1016/j.yrtph.2024.105613>
- ¹¹² Painter JL, Chalamalasetti VR, Kassekert R, Bate A. Automating pharmacovigilance evidence generation: using large language models to produce context-aware structured query language. *JAMIA open*. 2025;Feb;8(1):ooaf003. ([Journal full text](#)) <https://doi.org/10.1093/jamiaopen/ooaf003>
- ¹¹³ Benaïche A, Billaut-Laden I, Randriamihaja H, Bertocchio JP. Assessment of the Efficiency of a ChatGPT-Based Tool, MyGenAssist, in an Industry Pharmacovigilance Department for Case Documentation: Cross-Over Study. *J Med Internet Res* 2025;27:e65651. doi: 10.2196/65651
- ¹¹⁴ Meldau EL, Bista S, Rofors E, Gattepaille LM. Automated Drug Coding Using Artificial Intelligence: An Evaluation of WHODrug Koda on Adverse Event Reports. *Drug Safety*. 2022;May;45(5):549-561. ([Journal full text](#)) <https://doi.org/10.1007/s40264-022-01162-7>
- ¹¹⁵ Kreimeyer K, Dang O, Spiker J, Gish P, Weintraub J, Wu E, Ball R, Botsis T. Increased confidence in deduplication of drug safety reports with natural language processing of narratives at the us food and drug administration. *Frontiers in Drug Safety and Regulation*. 2022;Jun15;2:918897. ([Journal full text](#)) <https://doi.org/10.3389/fdsfr.2022.918897>
- ¹¹⁶ Norén GN, Orre R, Bate A, Edwards IR. Duplicate detection in adverse drug reaction surveillance. *Data Mining and Knowledge Discovery*. 2007;Jun;14:305-28. ([Journal abstract](#)) <https://doi.org/10.1007/s10618-006-0052-8>
- ¹¹⁷ Bergman E, Dürlich L, Arthurson V, Sundström A, Larsson M, Bhuiyan S, Jakobsson A, Westman G. BERT based natural language processing for triage of adverse drug reaction reports shows close to human-level performance. *PLOS Digital Health*. 2023;Dec6;2(12):e0000409. ([Journal full text](#)) <https://doi.org/10.1371/journal.pdig.0000409>
- ¹¹⁸ Painter J, Haguinet F, Cranfield C, Bate A. MSR20 NLP and Machine Learning to Automate Identification of Suspected Medication Errors from Real World Unstructured Narratives. *Value in Health, Volume 26, Issue 6, S281*. <https://doi.org/10.1016/j.jval.2023.03.1556>

- ¹¹⁹ Caster O, Juhlin K, Watson S, Norén GN. Improved statistical signal detection in pharmacovigilance by combining multiple strength-of-evidence aspects in *vigiRank*: retrospective evaluation against emerging safety signals. *Drug Safety*. 2014;Aug;37:617-628. ([Journal full text](#)) <https://doi.org/10.1007/s40264-014-0204-5>
- ¹²⁰ Caster O, Sandberg L, Bergvall T, Watson S, Norén GN. *vigiRank* for statistical signal detection in pharmacovigilance: first results from prospective real-world use. *Pharmacoepidemiology and Drug Safety*. 2017;Aug;26(8):1006-1010. ([Journal full text](#)) <https://doi.org/10.1002/pds.4247>
- ¹²¹ Scholl JH, van Hunsel FP, Hak E, van Puijenbroek EP. A prediction model-based algorithm for computer-assisted database screening of adverse drug reactions in the Netherlands. *Pharmacoepidemiology and Drug Safety*. 2018;Feb;27(2):199-205. ([Journal full text](#)) <https://doi.org/10.1002/pds.4364>
- ¹²² Norén GN, Meldau EL, Chandler RE. Consensus clustering for case series identification and adverse event profiles in pharmacovigilance. *Artificial Intelligence in Medicine*. 2021;Dec1;122:102199. ([Journal full text](#)) <https://doi.org/10.1016/j.artmed.2021.102199>
- ¹²³ Rudolph A, Mitchell J, Barrett J, Sköld H, Taavola H, Erlanson N, Melgarejo-González C, Yue QY. Global safety monitoring of COVID-19 vaccines: How pharmacovigilance rose to the challenge. *Therapeutic Advances in Drug Safety*. 2022;Aug;13:20420986221118972. ([Journal full text](#)) <https://doi.org/10.1177/20420986221118972>
- ¹²⁴ Habets PC, van IJzendoorn DG, Vinkers CH, Härmark L, de Vries LC, Otte WM. Development and validation of a machine-learning algorithm to predict the relevance of scientific articles within the field of teratology. *Reproductive Toxicology*. 2022;Oct1;113:150-154. ([Journal full text](#)) <https://doi.org/10.1016/j.reprotox.2022.09.001>
- ¹²⁵ OECD (2023), “The state of implementation of the OECD AI Principles four years on”, OECD Artificial Intelligence Papers, No. 3, OECD Publishing, Paris, <https://doi.org/10.1787/835641c9-en>.
- ¹²⁶ OECD (2023), “The state of implementation of the OECD AI Principles four years on”, OECD Artificial Intelligence Papers, No. 3, OECD Publishing, Paris, <https://doi.org/10.1787/835641c9-en>.
- ¹²⁷ European Medicines Agency. Reflection Paper on the Use of Artificial Intelligence (AI) in the Medicinal Product Lifecycle. 2024 (Full text accessed 21 March 2025).
- ¹²⁸ European Commission: Directorate-General for Communications Networks, Content and Technology, The Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self assessment, Publications Office, 2020. ([Full text](#) accessed 21 March 2025) <https://data.europa.eu/doi/10.2759/002360>
- ¹²⁹ Department of Industry, Science, Energy and Resources Australia. Australia’s Artificial Intelligence Ethics Principles. ([Webpage](#) accessed 21 March 2025).
- ¹³⁰ Government of Canada. Artificial Intelligence and Data Act (AIDA) Companion Document. ([Full text](#) accessed 21 March 2025).
- ¹³¹ Ministry of Health Singapore. 1.0 Artificial Intelligence in Healthcare Guidelines (AIHGLE). 2021. ([Full text](#) accessed 21 March 2025).
- ¹³² Government of the United Kingdom. Department of Science, Innovation & Technology. Implementing the UK AI Regulatory Principles: Guidance for Regulators. 2024. ([Full text](#) accessed 21 March 2025).
- ¹³³ White House Office of Science and Technology Policy (OSTP). Blueprint for an AI Bill of Rights. ([Webpage](#) accessed 21 March 2025).
- ¹³⁴ Pan American Health Organization. Artificial Intelligence of Public Health. Digital Transformation Toolkit, Knowledge Tools. ([Full text](#) accessed 21 March 2025).
- ¹³⁵ World Health Organization. Ethics and governance of artificial intelligence for health. 2021. ([Full text](#) accessed 21 March 2025).
- ¹³⁶ Organisation for Economic Co-operation and Development. OECD AI Principles Overview. 2019, updated in 2024. ([Full text](#) accessed 21 March 2025).
- ¹³⁷ European Medicines Agency. Reflection Paper on the Use of Artificial Intelligence (AI) in the Medicinal Product Lifecycle. 2024. ([Full text](#) accessed 21 March 2025).
- ¹³⁸ European Medicines Agency. Reflection Paper on the Use of Artificial Intelligence (AI) in the Medicinal Product Lifecycle. 2024. ([Full text](#) accessed 21 March 2025).
- ¹³⁹ Heads of Medicines Agencies - European Medicines Agency Joint Big Data Steering Group. Multi-annual artificial intelligence Workplan 2023-2028. 2023. ([Full text](#) accessed 21 March 2025).
- ¹⁴⁰ U.S. Food and Drug Administration. Using Artificial Intelligence & Machine Learning in the Development of Drug and Biological Products. 2023, revised 2025. ([Full text](#) accessed 21 March 2025)
- ¹⁴¹ U.S. Food and Drug Administration. Center for Drug Evaluation and Research. Artificial Intelligence in Drug Manufacturing. 2023. ([Full text](#) accessed 21 March 2025)
- ¹⁴² U.S. Food and Drug Administration. CDER Emerging Drug Safety Technology Program (EDSTP). 2024. ([Webpage](#) accessed 21 March 2025)

¹⁴³ European Medicines Agency. *Harnessing AI in medicines regulation: use of large language models (LLMs)*. 2024. ([Full text](#) accessed 21 March 2025)

¹⁴⁴ Government of Canada. *Guide to the use of generative AI*. 2023. ([Full text](#) accessed 21 March 2025)

¹⁴⁵ Organisation for Economic Co-operation and Development. *Initial policy considerations for generative artificial intelligence*. 2023. ([Full text](#) accessed 21 March 2025)

¹⁴⁶ World Health Organization. *Ethics and governance of artificial intelligence for health: Guidance on large multi-modal models*. ([Full text](#) accessed 21 March 2025)

¹⁴⁷ World Health Organization (WHO). *Regulatory considerations on artificial intelligence for health*. Geneva: World Health Organization; 2023. ([Full text](#) accessed 21 March 2025)

¹⁴⁸ European Commission. *Ethics guidelines for trustworthy AI*. 2019. (accessed 21 March 2025).

¹⁴⁹ OECD Legal Instruments, *Recommendation of the Council on Artificial Intelligence*, OECD/LEGAL/0449 2019, amended 2024. ([Full text](#) accessed 21 March 2025)

¹⁵⁰ U.S. Food and Drug Administration, Health Canada, Medicines & Healthcare products Regulatory Agency. *Good machine learning practice for medical device development: Guiding Principles*. 2021. ([Full text](#) accessed 21 March 2025)

Principles for integrating and implementing artificial intelligence within pharmacovigilance

656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698

This is a summary of the broad framework of principles for integrating and implementing AI within PV that will be addressed in the following chapters in greater detail.

Risk-based approach

A risk-based approach acknowledges the potential hazards that AI systems can pose and recognises that different use cases present varying types and levels of risk. This necessitates a risk assessment that identifies, prioritises, and manages risks that could negatively affect a PV system's behaviour and results, taking into consideration existing process controls. A risk is characterised by both the anticipated impact and the likelihood of negative outcomes.

This approach also supports procedures to identify and reduce errors and biases in a way that is proportionate to their risk. It influences the implementation strategies of AI systems (including documentation, compliance, and record-keeping), which should generally be commensurate with the identified risk.

Human oversight

Human oversight refers to the expected role of humans in the design, implementation, monitoring, and analysis of AI systems in PV. It requires a framework to manage performance and to detect and mitigate potential issues related to the AI system.

Validity

Validity means that a system achieves its intended purpose within acceptable parameters. It requires predefining acceptable performance levels, selecting appropriate data for model training and/or testing, assessing model performance in a realistic setting, and integrating the system into an ongoing quality assessment process.

Robustness

Robustness means that a system reliably achieves its intended objectives (while accounting for variations in data).

Transparency

Transparency regarding AI involves disclosing information between organizations or individuals. This includes sharing relevant documentation of the AI system lifecycle (i.e. design, development, evaluation, deployment, operation, re-training, maintenance and decommission) to facilitate traceability and providing stakeholders with enough information to have a general understanding of the AI system, its use, risks, limitations, and impact on their rights.

Data privacy

Data privacy refers to the fundamental right of an individual to control how their personal information is collected, stored, shared, and used. It is an aspect of the principle of "respect for persons" that is foundational to the conduct of biomedical research. Regulations, legislation and guidance documents provide measures intended to preserve the confidentiality, anonymity, autonomy and control of sensitive and potentially personally identifiable health data in the setting of PV.

Fairness & Equity

Fairness and equity require awareness of and adherence to impartiality, equality, non-discrimination, diversity, justice, and lawfulness. The benefits of AI in PV should be equitable across all relevant

699 populations and groups. Throughout the AI lifecycle, it is important to avoid and mitigate unfair bias,
700 and any discriminatory practices and unjust social wellbeing and environmental impacts.

701 *Governance*

702 Governance refers to the human management and oversight used to control and direct the use of AI
703 in the PV system. An AI governance framework requires implementation of risk management
704 practices and policies to ensure adherence to the AI guiding principles.

705 *Accountability*

706 Accountability applies to clearly defined roles, responsibilities and liability for organisations and/or
707 individuals deploying, operating and managing AI systems. It requires the adoption of appropriate
708 governance measures by relevant stakeholders, including but not limited to regulators, vendors,
709 users, developers, data providers or pharmaceutical companies involved in setting policy, developing,
710 deploying and managing AI systems. This ensures operations remain within expected parameters
711 throughout the AI lifecycle while addressing any unforeseen consequences.

712

713 Chapter 3: Risk-based approach

714 *Principle*

715 A risk-based approach acknowledges the potential hazards that AI systems can pose and recognises
716 that different use cases present varying types and levels of risk. This necessitates a risk assessment
717 that identifies, prioritises, and manages risks that could negatively affect a PV system's behaviour
718 and results, taking into consideration existing process controls. A risk is characterised by both the
719 anticipated impact and the likelihood of negative outcomes.¹⁵¹

720 This approach also supports procedures to identify and reduce errors and biases in a way that is
721 proportionate to their risk. It influences the implementation strategies of AI systems (including
722 documentation, compliance, and record-keeping), which should generally be commensurate with the
723 identified risk.

724 *Key messages*

- 725 • Integrating AI into PV processes needs to take into account that the performance of both AI
726 algorithms, and humans, is imperfect.
- 727 • The risks potentially associated with the use of AI in PV may affect patient safety, the trust
728 and engagement of PV users, the efficiency of PV processes as well as compliance with
729 regulatory standards and ethical principles.
- 730 • By focusing efforts and resources where they most matter, a sound risk-based approach
731 enables organisations to make the most of AI capabilities while ensuring that neither patient
732 safety nor PV stakeholders are adversely affected.
- 733 • The risk-based approach applies to the human oversight modalities, the validity and
734 robustness strategy, the level of transparency, and the efforts to uphold fairness and equity,
735 and data privacy.
- 736 • The risk assessment should consider the AI system itself, the context of use, and the
737 potential impact and likelihood of risks materialising.
- 738 • Comparative performance to current best practice should also be considered.
- 739 • Assessment should be end-to-end with an emphasis on end objective.
- 740 • A risk-based approach should be reviewed at regular intervals and adapted if needed.

741 Introduction

742 *Regulatory considerations*

743 Regardless of the integration of AI elements, PV systems are expected to comply with existing
744 regulations and good pharmacovigilance practices (GVP).^{152,153} In accordance with GVP, a wide range
745 of PV processes are considered critical for business continuity purposes, including collection and
746 handling of ICSRs, signal management, and PSURs.¹⁵⁴

747 Regulatory frameworks generally recommend a risk-based approach in the development,
748 deployment, monitoring, documentation and regulatory oversight of AI systems, to ensure that
749 relevant risks are anticipated, identified and mitigated throughout the system lifecycle.^{155,156,157} The
750 European Artificial Intelligence Act (EU AI Act)¹⁵⁸ introduces four risk categories for AI systems: low or
751 minimal risk, limited risk (transparency obligations), high risk, and unacceptable risk (prohibited AI
752 practices). High-risk AI systems, which include e.g. AI-based medical software/devices or AI systems
753 used for staff recruitment, are associated with strict requirements and obligations on providers and
754 deployers, including risk-mitigation systems, high quality data sets for training, validation and testing,
755 logging of activity, detailed documentation, clear user information, human oversight, and a high level
756 of robustness, accuracy, and cybersecurity. While the guiding principles advocated throughout this

757 report overlap with the EU AI Act's requirements for high-risk AI systems, determining the applicable
758 EU AI Act's risk category of an AI system considered for integration into an organisation's PV process
759 will likely require a careful case-by-case assessment, with legal advice as appropriate. Within the
760 medicines' lifecycle, EMA foresees AI systems with 'high patient risk' in use cases where patient
761 safety is affected and AI systems with 'high regulatory impact' in use cases where impact on the
762 regulatory decision making is substantial.¹⁵⁹ The Artificial Intelligence and Data Act (AIDA) was
763 developed to ensure the development of responsible AI in Canada, with a risk-based approach
764 aligned with international norms, including the EU AI Act, the OECD AI Principles, and the US National
765 Institute of Standards and Technology (NIST) Risk Management Framework (RMF).¹⁶⁰

766 During development and other stages of an AI solution's lifecycle as relevant, and depending on the
767 level of risk to individual patients, public health or the regulatory decision making, applicants and
768 developers should consider engaging actively with regulatory authorities and seek suitable scientific
769 advice. Where necessary, technical qualification of the AI technology through appropriate channels
770 should be sought based on legislative or regulatory requirements applicable to medicinal products,
771 medical devices and/or software development.^{161,162} Due to its fast-moving nature, the use of AI
772 technology in the medicines' lifecycle including in PV will pose challenges to both regulators,
773 required to adapt and keep abreast of this evolving field,¹⁶³ and industry PV stakeholders, required to
774 maintain regulatory compliance (see Chapter on [Future Vision](#)).

775 **Motivation and interplay with other guiding principles**

776 A sound, risk-based approach will allow organisations to focus their efforts and resources where they
777 matter most to maximise their AI capabilities while ensuring that guiding principles are upheld, as
778 described earlier.

779 Rather than a self-standing principle, a risk-based approach is applicable to the other guiding
780 principles presented in this report. Notably, a risk-based approach will inform where, when, how and
781 how much human oversight should be implemented within PV processes involving AI in addition to
782 other risk mitigation activities and conversely, an AI solution may be 'risk-assessed' assessed taking
783 into account the degree and nature of existing human oversight (see Chapter on [Human oversight](#)). A
784 risk-based approach should be applied to the testing and verification of AI systems (see Chapter on
785 [Validity & Robustness](#)) and the level of documentation and record-keeping (see Chapter on
786 [Transparency](#)). A risk-based approach is also relevant to data privacy and fairness and equity. For
787 example, AI systems should be assessed for any risks that specific groups may be under-served or
788 biased against, and those risks should be appropriately mitigated (see Chapter on [Fairness & Equity](#)).

789 **Types of risks**

790 This section briefly outlines some of the risks potentially associated with the use of AI solutions in PV.

791 *Risks to patient safety and public health*

792 Inadequate use of AI solutions in PV, or their poor performance, may impede the fulfilment of PV
793 objectives: detection, assessment, understanding and prevention of adverse effects of drugs or
794 vaccines, which may come at the cost of patient safety or public health. Unreliable outputs produced
795 by an AI system, including but not limited to false negatives or false positives, or unfair bias, could
796 negatively impact PV activities with e.g. relevant adverse events not captured, events misclassified
797 during case processing, or signals missed. This could result in safety issues not being identified or
798 being identified with delay, potentially putting patients at risk. In rare scenarios, the late detection of
799 new, unexpected safety signals could have a major public health impact ('Black swan' events).¹⁶⁴ An
800 initially robust AI tool could also start underperforming over time due to e.g. model drift, or become

801 inoperative due to an IT incident or system failure, which would impede the PV activity that the AI
802 tool is intended to support.

803 *Risks to user trust and engagement*

804 The lack of transparency and interpretability of certain AI algorithms may hinder trust and
805 acceptance by users, including PV professionals (see Chapter on [Transparency](#)). Lack of trust from
806 users may also result from poor previous experience with AI systems of insufficient validity and
807 robustness, leading to mistrust of AI solutions in general. In clinic-based PV settings, a more subtle
808 potential source of mistrust is 'uniqueness neglect', in which patients prefer a human clinician over a
809 more accurate computer due to a belief that machines do not fully accommodate their personal
810 human uniqueness.¹⁶⁵ Other possible sources of mistrust include poor performance for certain
811 subpopulations or failure to protect confidentiality of personal data during the development or
812 operation of an AI system. Conversely, some users may put excessive trust in AI systems, leading to
813 automation bias (especially if those have shown robust performance upon validation) and the
814 resulting unconscious bias to accept erroneous outputs. Additionally, integrating AI solutions into
815 existing workflows and systems may pose technical, organisational, and cultural challenges, with a
816 risk of degraded job motivation or satisfaction in the absence of adequate training and change
817 management strategies (see Chapter on [Human oversight](#)).

818 *Risks to efficiency*

819 Although the integration of AI in PV processes is generally aimed at increasing efficiency,
820 substandard AI solutions may cause more manual work than they save, if for instance, significant
821 time is required to understand and verify the AI outputs or bring them up to acceptable standards.
822 Uncertainties, such as false positives, in interpretations and actions based on AI outputs might add to
823 inefficiencies or suboptimal use of limited resources.

824 *Legal, ethical and other risks*

825 Other risks may be related to data privacy, cybersecurity threats, intellectual property, liability, or
826 economic and reputational aspects.

827 This chapter mainly focusses on the impact that the use AI tools in PV processes could have on
828 patient safety. Potential concerns related to user acceptance and other challenges are further
829 discussed in the chapters on [Data Privacy](#), [Fairness & Equity](#), [Transparency](#), [Human oversight](#) and
830 [Governance & Accountability](#).

831 **Risk assessment**

832 **General considerations**

833 Organisations planning to deploy AI-based tools to support PV processes are expected to perform a
834 thorough risk analysis. This assessment should be performed for each AI system and should form the
835 basis for a risk-proportionate approach applied throughout the AI system's lifecycle from
836 development to routine use.

837 When determining the level of risk related to the implementation of an AI tool within a PV system,
838 key considerations include the AI technology itself, the context of use, the likelihood of risks
839 materialising and their potential impact.

840 *Artificial intelligence technology*

841 The level of risk may depend on the type of tool used (e.g. static vs dynamic model), the underlying
842 data quality, the novelty of the technology or the maturity of the system (i.e. lifecycle stage).

843 Particular caution should be exercised with the integration of GenAI models within PV processes.
844 Compared to simpler or more explainable AI approaches, the non-deterministic nature of GenAI and

845 similar AI models, the opacity of training data and the potential for hallucinations, may make the
846 detection and mitigation of issues more challenging and require consideration of further guardrails.

847 As the AI landscape continues to evolve, so will AI-related risk areas. New risks may emerge while
848 others may become less prominent; for instance, the current challenges associated with GenAI/LLMs
849 may be solved, making their integration into sensitive PV processes safer.

850 *Context of use and degree of influence*

851 These broadly refer to the place and importance of the AI solution within the overall PV system,
852 including:

- 853 • Whether or not the AI solution is used in a critical PV process or high-risk context (e.g.
854 emergency Public Health use, novel substance, clinical trial cases);
- 855 • At which stage within a particular process the AI tool intervenes (e.g. automated triage of
856 relevant cases as a preliminary step to signal review) and whether the tool is assistive or
857 directly supports a PV process;
- 858 • The extent of human involvement and oversight in the process (see Chapter on [Human
859 oversight](#)).

860 *Impact and likelihood*

861 Not all occurrences of system malfunction or suboptimal model performance are as likely, nor will
862 they have the same impact. For instance, a duplicate detection tool applied to a very large database
863 is not expected to detect 100% of duplicates but missed duplicates will have no or limited
864 consequences in terms of patient safety, whereas the late detection of a very serious signal in a
865 context of mass patient exposure happens very rarely but may have dramatic public health
866 consequences (i.e. black swan event).

867 **Examples of structured approaches**

868 Risk-based assessment frameworks have been proposed in various domains and may provide
869 inspiration to organisations wishing to deploy AI solutions within PV systems. Selected examples are
870 briefly described hereafter.

871 *Credibility assessment framework*

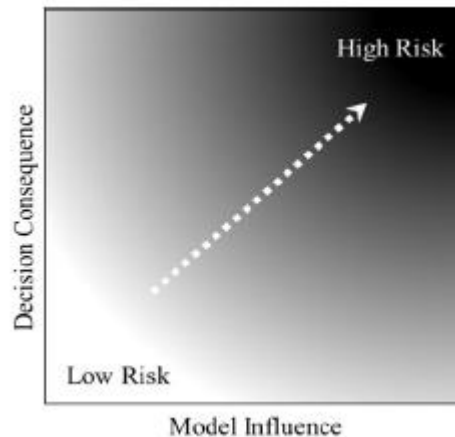
872 The FDA proposes a stepwise approach to demonstrate the credibility of AI models to produce
873 information or data intended to support regulatory decision making regarding the safety,
874 effectiveness, or quality of drugs (see also Chapter [Landscape analysis](#)).¹⁶⁶ Similar frameworks have
875 been proposed for the use of computational models in medical device submissions¹⁶⁷ or drug
876 development.¹⁶⁸ The preliminary steps of the credibility assessment, as outlined below, help assess
877 the model risk.

- 878 1. *Define the question of interest*: This describes the specific question, decision, or concern
879 to be addressed by the AI model.
- 880 2. *Define the context of use*: this is a description of how the model will be used to address
881 the question of interest, i.e. the specific role and scope of the AI model.
- 882 3. *Assess the AI model risk*: this is defined by (i) the contribution of the evidence derived
883 from the AI model relative to other contributing evidence used to inform the question of
884 interest, i.e. model influence; and (ii) the significance of an adverse outcome resulting
885 from an incorrect decision concerning the question of interest, i.e. decision
886 consequence. The ratings for decision consequence and model influence are
887 independently determined, but are shaped by the context of use, thus enabling model
888 risk to be case specific. The AI model risk is assessed by subject matter expertise and
889 judgement on the possibility of model output leading to an adverse outcome, rather than

890 the intrinsic risk of the model itself. As illustrated in Figure 4, the model risk moves from
 891 low to high as decision consequence or model influence increases.

892 **Figure 5: Model risk matrix**

893 Source: U.S. Food and Drug Administration.¹⁶⁹



894

895 *Algorithmic impact assessment*

896 In Canada, the mandatory algorithmic impact assessment (AIA) tool¹⁷⁰ is designed to help
 897 departments and agencies better understand and manage the risks associated with automated
 898 decision systems. It is composed of questions in various formats that consider many factors (e.g.
 899 system's design, algorithm, decision type, impact, data) within risk and mitigation areas and
 900 contribute to a scoring system. The value of each question is weighted based on the level of risk it
 901 introduces or mitigates in the automation project. The resulting impact levels (from I: little impact, to
 902 IV: very high impact) determine the mitigations required under the Directive on Automated Decision-
 903 Making.¹⁷¹

904 **Issue detection and risk mitigation**

905 Defining when to mitigate requires knowing how to detect issues based on a pre-defined risk-
 906 proportionate testing and verification plan which is laid out during the development of the AI system.
 907 Testing and verification are essential steps of Computerized System Validation (CSV), which considers
 908 different levels based on AI system maturity. The latest version of the Good Automated
 909 Manufacturing Practice 5 (GAMP 5) of the International Society for Pharmaceutical Engineering
 910 (ISPE), a framework widely adopted by pharmaceutical companies and health authorities, contains an
 911 appendix focusing on AI and ML.¹⁷² Testing should be based on pre-defined key performance
 912 indicators and acceptance criteria, considering the human performance, and account for the
 913 identified risk areas, e.g. low quality data. Issues can be detected using mechanisms such as 'golden
 914 questions', i.e. known truths which are checked each time an AI system undergoes a change.

915 After the AI system has been proven fit for purpose and deployed, an ongoing process should be in
 916 place to monitor its performance and trigger mitigation measures when issues are detected.

917 A risk-based approach may be very conservative in the initial stages of deployment with additional
 918 pre-determined mitigation measures in place, for example, high percentage of human-in-the-loop. As
 919 confidence in the routine performance increases over time, based on pre-defined indicators and
 920 examination of sample outputs by human experts, a gradual reduction in the frequency, amount (e.g.
 921 number of samples) or depth of human controls may be considered.

922 When issues or performance deviations are detected, risk-based mitigation measures may include:

- 923 • *Human-in-the-loop*: Increased or full human review/quality control, indefinitely or until
 924 performance levels are back within acceptance criteria, e.g. if a seriousness detection

- 925 algorithm fails to detect seriousness criteria in some cases, e.g. false negatives, mitigation
 926 could involve reviewing all cases classified as non-serious until the issue is understood and
 927 addressed;
- 928 • *AI improvement* strategies such as model re-training or hallucination mitigation strategies;¹⁷³
 - 929 • Articulation of the *level of uncertainty* in AI outputs;
 - 930 • Approaches to *combat automation bias or complacency*,¹⁷⁴ e.g. mock data simulations or
 931 injection of simulated false positive outputs for verification/assessment;
 - 932 • *Decommissioning* of the tool when mitigation options appear inefficient or costly, in which
 933 case alternative approaches should be considered.
- 934 Finally, AI components, especially those deployed in critical PV processes, should be included in the
 935 organisation's business continuity plan.
- 936 The above aspects are further developed in the Chapters on [Validity & Robustness](#), [Human](#)
 937 [Oversight](#), and [Governance & Accountability](#).
- 938 The risk-based approach should be reviewed at regular pre-determined intervals to adapt oversight
 939 measures based on performance data but also as new technical options for risk mitigation emerge.

940 Documentation

- 941 The key components of the AI-related risk management strategy should be documented, including:
- 942 • AI system risk assessment;
 - 943 • Testing plan with key performance indicators and acceptance criteria including any
 944 comparative assessments;
 - 945 • Planned mitigation measures including human in-the-loop strategy and criteria for more
 946 stringent or reduced quality control, and continual monitoring after deployment;
 - 947 • Plans for periodic re-assessment and update of the risk management strategy;
 - 948 • Business continuity plan.

References

- ¹⁵¹ Council for International Organizations of Medical Sciences. *Practical Approaches to Risk Minimisation for Medicinal Products*. Geneva, Switzerland: Council for International Organizations of Medical Sciences; 2014. ([Full text](#) accessed 21 March 2025)
- ¹⁵² European Medicines Agency. *Guideline on good pharmacovigilance practices (GVP): Module I – Pharmacovigilance systems and their quality systems*. (Full text accessed 21 March 2025)
- ¹⁵³ U.S. Food and Drug Administration. *Good Pharmacovigilance Practices and Pharmacoeconomic Assessment. Guidance for Industry*. March 2005. ([Full text](#) accessed 21 March 2025).
- ¹⁵⁴ European Medicines Agency. *Guideline on good pharmacovigilance practices (GVP): Module I – Pharmacovigilance systems and their quality systems*. (Full text accessed 21 March 2025)
- ¹⁵⁵ World Health Organization. *Regulatory considerations on artificial intelligence for health*. Geneva: World Health Organization; 2023. Licence: CC BY-NC-SA 3.0 IGO. (accessed 21 March 2025)
- ¹⁵⁶ European Medicines Agency. *Reflection paper on the use of Artificial Intelligence (AI) in the medicinal product lifecycle*. ([Webpage](#) accessed 21 March 2025)
- ¹⁵⁷ U.S. Food and Drug Administration. *Considerations for the Use of Artificial Intelligence To Support Regulatory Decision-Making for Drug and Biological Products. Draft Guidance for Industry and Other Interested Parties*. January 2025. ([Full text](#) accessed 21 March 2025)
- ¹⁵⁸ European Parliament and Council of the European Union. *Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828*. ([Webpage](#) accessed 21 March 2025)
- ¹⁵⁹ European Medicines Agency. *Reflection paper on the use of Artificial Intelligence (AI) in the medicinal product lifecycle*. ([Webpage](#) accessed 21 March 2025)

- ¹⁶⁰ Government of Canada. *The Artificial Intelligence and Data Act (AIDA) – Companion document.* ([Webpage](#) accessed 21 March 2025)
- ¹⁶¹ European Medicines Agency. *Reflection paper on the use of Artificial Intelligence (AI) in the medicinal product lifecycle.* ([Webpage](#) accessed 21 March 2025)
- ¹⁶² U.S. Food and Drug Administration. *Considerations for the Use of Artificial Intelligence To Support Regulatory Decision-Making for Drug and Biological Products. Draft Guidance for Industry and Other Interested Parties. January 2025.* ([Full text](#) accessed 21 March 2025)
- ¹⁶³ Hines PA, Herold R, Pinheiro L, Frias Z, Arlett P. *Artificial intelligence in European medicines regulation.* *Nat Rev Drug Discov.* 2023;Feb;22(2):81-82. ([Journal full text](#)) <https://doi.org/10.1038/d41573-022-00190-3>.
- ¹⁶⁴ Kjoersvik O, Bate A. *Black Swan Events and Intelligent Automation for Routine Safety Surveillance.* *Drug Saf.* 2022;May;45(5):419-427. ([Journal full text](#)) <https://doi.org/10.1007/s40264-022-01169-0>.
- ¹⁶⁵ Hauben M. *Artificial intelligence in pharmacovigilance: Do we need explainability?* *Pharmacoepidemiology Drug Saf.* 2022;Dec31(12):1311-1316. ([Journal full text](#)) <https://doi.org/10.1002/pds.5501>.
- ¹⁶⁶ U.S. Food and Drug Administration. *Considerations for the Use of Artificial Intelligence To Support Regulatory Decision-Making for Drug and Biological Products. Draft Guidance for Industry and Other Interested Parties. January 2025.* ([Full text](#) accessed 21 March 2025)
- ¹⁶⁷ U.S. Food and Drug Administration. *Assessing the Credibility of Computational Modeling and Simulation in Medical Device Submissions.* 2023. ([Full text](#) accessed 21 March 2025)
- ¹⁶⁸ International Conference on Harmonisation (ICH). *General Principles for Model-Informed Drug Development. M15. Draft version endorsed on 06 November 2024.* ([Full text](#) accessed 21 March 2025)
- ¹⁶⁹ U.S. Food and Drug Administration. *Considerations for the Use of Artificial Intelligence To Support Regulatory Decision-Making for Drug and Biological Products. Draft Guidance for Industry and Other Interested Parties. January 2025.* ([Full text](#) accessed 21 March 2025)
- ¹⁷⁰ Government of Canada. *Algorithmic Impact Assessment tool.* ([Webpage](#) accessed 21 March 2025)
- ¹⁷¹ Government of Canada. *Directive on Automated Decision-Making.* ([Webpage](#) accessed 21 March 2025)
- ¹⁷² The International Society for Pharmaceutical Engineering (ISPE). *GAMP 5 Guide 2nd Edition.* 2022. ([Abstract](#) accessed 21 March 2025).
- ¹⁷³ Hakim JB, Painter JL, Ramcharran D, Kara V, Powell G, Sobczak P, Sato C, Bate A, Beam A. *The Need for Guardrails with Large Language Models in Medical Safety-Critical Settings: An Artificial Intelligence Application in the Pharmacovigilance Ecosystem.* *arXiv preprint arXiv:2407.18322.* 2024;Jul1. ([Journal full text](#)) <https://doi.org/10.48550/arXiv.2407.18322>
- ¹⁷⁴ Adler-Milstein J, Redelmeier DA, Wachter RM. *The Limits of Clinician Vigilance as an AI Safety Bulwark.* *JAMA.* 2024;331(14):1173-1174. ([Journal full text](#)) <https://doi.org/10.1001/jama.2024.3620>.

949

950 Chapter 4: Human oversight

951 *Principle*

952 Human oversight refers to the expected role of humans in the design, implementation, monitoring,
953 and analysis of AI systems in PV. It requires a framework to manage performance and to detect and
954 mitigate potential issues related to the AI system.

955 *Key messages*

- 956 • Human oversight supports the optimisation of the performance of AI systems deployed in PV
957 and increases trustworthiness and accountability.
- 958 • The extent and nature of human oversight for an AI solution should follow a risk-based
959 approach.
- 960 • Quality assurance principles should apply to the conduct of the human oversight of AI
961 systems in PV.
- 962 • The increased use of automation and AI to support PV processes will require redefining
963 skillsets to integrate AI with human expertise, ensuring robustness and reliability in decision-
964 making processes. This will lead to a transformation of traditional roles and competencies
965 that requires appropriate change management and training strategies.

966 Introduction

967 *Motivation*

968 Human oversight is required to minimise the risk that an AI system undermines human autonomy or
969 causes other negative or unintended effects.¹⁷⁵ Human agency and oversight are a key requirement
970 of trustworthy AI according to several regulatory frameworks, including the Assessment List for
971 Trustworthy Artificial Intelligence (ALTAI), the EU AI Act, and the Canadian Artificial Intelligence and
972 Data Act (AIDA), for high-risk systems^{176,177,178} (see also Chapter on [Landscape analysis](#)). Although
973 human review by itself does not guarantee full accuracy of outputs, human oversight is essential to
974 monitor the performance of AI systems and make corrections if needed, thereby increasing
975 trustworthiness and accountability for the AI system, especially in some high-risk applications.

976 AI systems are often intended to help eliminate manual, labour-intensive or complicated work
977 performed by humans, or to enhance human performance when used as intelligence augmentation
978 tools. However, due to the complexity and sensitivity of certain PV tasks, and the complex and
979 variable nature of PV data, AI components will exhibit increasingly good but imperfect performance.
980 This may require more extensive human intervention during the development, evaluation and
981 deployment of some AI solutions in PV to monitor and mitigate risks.

982 A key challenge and important starting point for defining an AI quality assurance (QA) approach is to
983 strike a balance between the efficiency boost that an AI system is intended to provide and the level
984 of human intervention that may be required to ensure a high-quality output. In plain words, ideally a
985 human expert should not do work that a machine can do well, and a machine should not do poorly
986 the work that a human expert can do well.¹⁷⁹

987 Human oversight is fundamental to a sound risk-based approach (see Chapter on [Risk-based
988 approach](#)). The level of monitoring of the performance of AI systems by humans should be
989 proportional to the potential impact on patient safety of an undetected mistake or spurious output
990 by the AI system.

991 **Considerations on human involvement and oversight**

992 **Multidisciplinary expertise**

993 The successful integration of AI solutions into PV systems requires that multidisciplinary human
994 expertise is mobilised as appropriate throughout the lifecycle of the tools, from development to
995 routine use. This multidisciplinary expertise is usually obtained through a close collaboration
996 between PV professionals, QA staff, data scientists, statisticians, AI/ML engineers, data engineers,
997 prompt engineers, IT specialists, cybersecurity experts, platform analysts, software engineers, ethics
998 specialists, legal experts, data protection officers, project managers, senior management, etc. (see
999 also Chapter on [Governance & Accountability](#)).

1000 PV professionals, i.e. staff performing core tasks in ICSR management, signal detection and analytics
1001 or risk management, hold robust 'domain' or 'subject matter' expertise, which is instrumental to
1002 effective integration of AI capabilities into PV processes. As such, PV professionals should be engaged
1003 in the design, development, pre-deployment and testing/piloting of AI solutions to ensure that they
1004 are fit for purpose and widely accepted by the end-users that they - PV professionals - will ultimately
1005 be.

1006 **Mechanisms of human oversight**

1007 Human oversight may serve different objectives and be achieved through governance mechanisms at
1008 different stages.¹⁸⁰ There are various possible approaches based on the activity monitored and how
1009 much autonomy is granted to an AI system. Depending on the scope, extent and intensity of human
1010 intervention, the European Commission's Ethics Guidelines for trustworthy AI describe three main
1011 governance mechanisms: human-in-the-loop (HITL), human-on-the-loop (HOTL) and human-in-
1012 command (HIC). HITL refers to the capability for human intervention in every decision cycle of the AI
1013 system. HOTL, which foresees a higher autonomy of the AI system, refers to the capability for human
1014 intervention during the design of an AI system and monitoring of its operation. HIC refers to the
1015 capability to oversee the overall activity of an AI system, including its broader economic, societal,
1016 legal and ethical impact, and the ability to decide when and how to use an AI system. This may
1017 include the decision not to use an AI system in a particular situation, to establish levels of human
1018 discretion during its use, or to ensure the ability to override a decision made by the system.¹⁸¹ The
1019 delineations of these three terms may vary according to sources¹⁸² and their practical
1020 implementation may differ according to individual organisations and use cases. For example, during
1021 the lifecycle of a given AI system in PV, human oversight may be exercised at early stages to help
1022 define the system's context of use or support the identification or development of reference datasets
1023 (HOTL) and, when deployed, to perform quality controls of the system (HOTL) or as part of its
1024 execution in case of a semi-automated system (HITL).

1025 As a rule, some level of human oversight is always required and the absence of a human-in/on-the-
1026 loop in any major or supporting PV process should be substantiated by a risk assessment, with risk
1027 mitigation measures in place.

1028 **Monitoring and interacting with deployed artificial intelligence systems**

1029 The level, frequency, means and modalities of human intervention required to monitor and interact
1030 with AI systems depend on the complexity of the task, the risks associated with suboptimal outputs,
1031 the type of AI system, and the performance as assessed during validation (see Chapters on [Risk-
1032 based approach](#) and [Validity & Robustness](#)). As experience with AI evolves, further clarity, guidance,
1033 and consistency in assessing these factors are likely to develop. As suggested above, the respective
1034 roles of the human and AI components in a particular process could be seen as a continuum, from an

1035 AI tool merely performing preparatory work to support assessment and decision making by a human,
1036 to a near-fully automated system merely monitored by a human who performs quality controls.
1037 Intermediate approaches may also be envisaged where, for instance, an AI solution flags cases it
1038 struggles with to a human specialist.

1039 The metrics and KPIs used to monitor the performance of deployed AI systems should be pre-defined
1040 as part of the testing and verification plan (see Chapters on [Validity & Robustness](#) and [Risk-based
1041 approach](#)).

1042 In situations where the standalone performance of an AI system is suboptimal (e.g. if it cannot match
1043 the established human performance, or when the associated risks are unacceptably high), one or
1044 more manual process steps must be considered, with a human fully in control of the final output.
1045 Even when a static AI-based system exceeds human performance upon validation, monitoring after
1046 development is still recommended to ensure that the performance does not fall below acceptable
1047 levels over time (see Chapter on [Validity & Robustness](#)). Each time an AI system undergoes
1048 modifications, human oversight should be directed at the change i.e. change-specific samples should
1049 be prioritised.

1050 There are different ways the performance of an AI system can be monitored once deployed. In a
1051 static AI system, one could perform a one-off retrospective analysis by checking a sample or the
1052 totality of generated outputs against expected outputs (see Chapter on [Validity & Robustness](#)). This
1053 may be followed by post hoc corrections and re-training or re-validation of the model. A more
1054 dynamic real-time, in-process interaction can also be envisaged where independent human
1055 assessment is applied to confirm or correct the AI output in a decision-support setting. In such a
1056 dynamic AI application, the interaction provides an opportunity for immediate feedback to the
1057 algorithm to continuously learn and adjust if needed. Running an independent model in parallel to
1058 the main AI system may also be an option in a one-off or continuous manner (AI-assisted human
1059 oversight).

1060 Caution is required in the monitoring of GenAI/LLM-based systems. Humans-in/on-the-loop should
1061 be aware of the inherent variability of outputs, limited explainability and risk of hallucinations, and
1062 not overly rely on the AI system's results. Processes must be robust, demonstrated to be effective,
1063 and maintain their dependability even in the event of erroneous outputs. Hallucinations, specifically,
1064 may lead to seemingly coherent and convincing outputs that may be deceiving for humans-in/on-
1065 the-loop. Regardless of the underlying AI technology, PV professionals should be empowered to
1066 challenge the system's outputs based on their experience and avoid falling for automation bias. On
1067 the other hand, they should be aware of the possibility of confirmation bias and remain open to the
1068 possibility that an AI output, albeit unexpected, is correct. AI solutions with high performance may
1069 also warrant specific monitoring strategies as humans are more prone to miss very rare¹⁸³

1070 **Transformation of traditional roles**

1071 As the PV landscape continues to embrace AI capabilities, a reduced dependency on large workforces
1072 with PV expertise is expected due to the replacement of some of the activities traditionally
1073 performed by PV professionals. Indeed, the increased use of automation and AI within PV processes
1074 will unburden PV professionals from certain repetitive, time-consuming, manual activities. This may
1075 render certain roles obsolete and thereby reduce the size, diversity and experience of the PV
1076 workforce, not unlike the impact on staff observed when organisations offshore activities. This may
1077 create legitimate concerns and anxiety about job displacement and employment prospects in the PV
1078 space, but also around work culture, motivation and fulfilment. Perceived unfairness may also ensue
1079 from the fact that some AI models are trained using historical datasets and documented decisions
1080 based on the work originally performed by PV professionals.

1081 On a brighter side, the introduction of AI in PV brings opportunities for growth for PV professionals.
1082 With fewer menial time-consuming tasks, PV experts will be able to focus on more scientifically

1083 complex and intellectually stimulating PV activities. In addition, the business needs associated with AI
1084 solutions will bring new roles in governance and human oversight. As mentioned earlier, PV
1085 professionals will be increasingly involved in the testing, evaluation, implementation, oversight and
1086 use of AI models. They will often be best placed to identify those activities in need of automation and
1087 suggest AI use cases accordingly. They may have to participate in design and development activities
1088 including model training and validation, participate in user acceptance testing, manage the
1089 challenges of automating and modifying existing processes, perform monitoring and quality control
1090 activities, identify and resolve issues related to inconsistent assessments, and interact with
1091 automation experts and vendors.

1092 Contributing to the development, use, and maintenance of AI systems will allow PV professionals to
1093 evolve with the changing PV landscape, but this will require that they extend their skillsets beyond
1094 core PV competencies.¹⁸⁴ These new skills include specific competencies around the use of the new
1095 systems and the critical evaluation of their outputs, as well as more general literacy around data
1096 science and AI, including a good understanding of AI capabilities, risks and limitations. Regulatory
1097 frameworks such as the EU AI Act impose an obligation on organisations to ensure a sufficient level
1098 of AI literacy of staff operating or using deployed AI systems.¹⁸⁵

1099 Beyond PV professionals, staff working in QA also need to develop an understanding of the
1100 organisation's human oversight strategy in addition to some AI literacy, to ensure that human
1101 oversight activities are adequate. Likewise, AI experts involved in the design and development of AI
1102 in PV solutions will need to develop an understanding of PV processes and the implications of
1103 operating in a regulated environment.

1104 Change management and readiness strategies are a key responsibility of organisations, which should
1105 put PV staff at the centre of role redefinition and upskilling opportunities. Adequate change
1106 management and training plans are a pre-requisite to a seamless, safe and successful integration of
1107 AI systems into PV processes, with a wide engagement and adoption by staff and smooth
1108 interactions between various roles (see also Chapter on [Governance](#)).

1109 Training programs should be carefully crafted, documented and evaluated so that their content and
1110 format (including materials and methods) meet the learning needs of the target audience (e.g. PV
1111 end-users, QA staff, AI experts). Human training in a decision-support context is an approach that
1112 may be drawn on to train staff monitoring and interacting with AI systems. It generally refers to
1113 programs designed to educate staff to use specific tools and make informed decisions effectively.
1114 This involves not only showing staff how to use the software front-end but also explaining the back-
1115 end functionalities and helping them build the skillset for critically evaluating the automated output.
1116 Training modalities (e.g. classroom-based vs online, live vs asynchronous) should be adapted to the
1117 system's complexity, the supported use case or task, and the specific needs of the organisation or the
1118 individual.^{186,187,188,189}

References

¹⁷⁵ European Commission. *Ethics guidelines for trustworthy AI*. 2019. Available from: <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>. [Accessed 16 August 2024] European Commission. *Ethics guidelines for trustworthy AI*. 2019. ([Webpage](#) accessed 21 March 2025)

¹⁷⁶ European Commission. *Ethics guidelines for trustworthy AI*. 2019. Available from: <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>. [Accessed 16 August 2024] European Commission. *Ethics guidelines for trustworthy AI*. 2019. ([Webpage](#) accessed 21 March 2025)

¹⁷⁷ European Parliament and Council of the European Union. *Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828*. ([Webpage](#) accessed 21 March 2025)

¹⁷⁸ Government of Canada. *The Artificial Intelligence and Data Act (AIDA) – Companion document*. ([Webpage](#) accessed 21 March 2025)

- ¹⁷⁹ Ball R, Dal Pan G. "Artificial Intelligence" for Pharmacovigilance: Ready for Prime Time? *Drug Saf.* 2022;May;45(5):429-438. ([Journal full text](#)) <https://doi.org/10.1007/s40264-022-01157-4>.
- ¹⁸⁰ Enqvist L. 'Human oversight' in the EU artificial intelligence act: what, when and by whom? *Law, Innovation and Technology.* 2023;Jul 3;15(2):508-535. ([Journal full text](#)) <https://doi.org/10.1080/17579961.2023.2245683>
- ¹⁸¹ European Commission. *Ethics guidelines for trustworthy AI.* 2019. Available from: <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>. [Accessed 16 August 2024] European Commission. *Ethics guidelines for trustworthy AI.* 2019. ([Webpage](#) accessed 21 March 2025)
- ¹⁸² Enqvist L. 'Human oversight' in the EU artificial intelligence act: what, when and by whom? *Law, Innovation and Technology.* 2023;Jul 3;15(2):508-535.
- ¹⁸³ Rich AN, Kunar MA, Van Wert MJ, Hidalgo-Sotelo B, Horowitz TS, Wolfe JM. Why do we miss rare targets? Exploring the boundaries of the low prevalence effect. *Journal of Vision.* 2008;Nov1;8(15):15. ([Journal full text](#))[HYPERLINK "https://doi.org/10.1167/8.15.15"](https://doi.org/10.1167/8.15.15)<https://doi.org/10.1167/8.15.15>
- ¹⁸⁴ Danysz K, Cicirello S, Mingle E, Assuncao B, Tetarenko N, Mockute R, et al. Artificial Intelligence and the Future of the Drug Safety Professional. *Drug Saf.* 2019;Apr;42(4):491-497. ([Journal full text](#))<https://doi.org/10.1007/s40264-018-0746-z>.
- ¹⁸⁵ European Parliament and Council of the European Union. Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828. ([Webpage](#) accessed 21 March 2025)
- ¹⁸⁶ Power DJ. A brief history of decision support systems. *DSSResources.com.* 2007;Mar10;3. ([Webpage](#) accessed 21 March 2025)
- ¹⁸⁷ Wang RY. A product perspective on total data quality management. *Communications of the ACM.* 1998;Feb1;41(2):58-65. ([Journal full text](#)) <https://doi.org/10.1145/269012.269022>
- ¹⁸⁸ Turban E. *Decision support and business intelligence systems.* Pearson Education India; 2011.
- ¹⁸⁹ Keen PG. *Decision support systems: a research perspective.* *Decision Support Systems: Issues and Challenges,* International Institute for Applied Systems Analysis (IIASA) Proceedings Series. 1980 Jun 23;11:23-7. ([E-book abstract](#))

1119

1120 **Chapter 5: Validity & Robustness**1121 *Principle*

- 1122 • Validity means that a system achieves its intended purpose within acceptable parameters. It
- 1123 requires predefining acceptable performance levels, selecting appropriate data for model
- 1124 training and/or testing, assessing model performance in a realistic setting, and integrating
- 1125 the system into an ongoing quality assessment process.
- 1126 • Robustness means that a system reliably achieves its intended objectives (while accounting
- 1127 for variations in data).

1128 *Key messages*

- 1129 • PV professionals and decision makers must learn to critically appraise proposed AI solutions
- 1130 whether they acquire them or participate in their development.
- 1131 • A performance evaluation able to demonstrate acceptable and robust results for the
- 1132 intended use under realistic conditions is crucial. Such an evaluation should cover a wide
- 1133 enough range of relevant examples to interrogate the model's objective and is often based
- 1134 on statistical metrics.
- 1135 • There should be a focus on looking to ensure sufficient representation of relevant types of
- 1136 data in the test set(s) to detect biases, promote adequate and generalizable performance
- 1137 across the intended deployment domain, assess usability, and identify circumstances where
- 1138 the model may underperform.
- 1139 • Many PV applications focus on very rare events or patterns (e.g. emerging safety signals,
- 1140 reports of a certain kind – such as related to pregnancy, duplicated reports etc.) and may
- 1141 require enrichment strategies to obtain representative test sets with high enough prevalence
- 1142 of the event of interest. If so, special care should be taken to ensure that performance
- 1143 evaluation results generalize to real-world settings.

1144 **Introduction**

1145 Ensuring the validity and robustness of AI solutions is central to building trust and achieving the best
 1146 possible value for end-users. To invest resources optimally, PV professionals and decision makers
 1147 must learn to critically appraise and evaluate proposed AI solutions regardless of whether they
 1148 develop them in-house or acquire them from other organizations. This requires familiarity with basic
 1149 principles for performance evaluation and some of the common pitfalls that may mislead
 1150 expectations on real-world performance in prospective use.

1151 AI solutions will often be embedded in broader computer systems supporting the PV use case. These
 1152 should be subjected to general computer system-validation according to standard practices for the
 1153 organization. In general, this will be considered orthogonal to ensuring the validity and robustness of
 1154 the core AI solution (and is out of scope for this document). However, some special considerations
 1155 regarding validation of systems that include dynamic AI models that continually learn from and adapt
 1156 to incoming data are presented in the Section on [Continuous integration and deployment](#). Our focus
 1157 will be on key considerations related to establishing the validity and robustness of AI models
 1158 themselves, including their dependency on underlying data for training and deployment and the
 1159 need for probabilistic / statistical performance evaluation.

1160 The nature of PV data may in some instances impact our ability and approach to leveraging AI
 1161 solutions. AI models depend heavily on the quality of the data they are trained on and the data they
 1162 use for ongoing predictions. PV data suffer from inconsistencies, incomplete entries, and
 1163 inaccuracies, and may vary substantially depending on the source. For example, the contribution of
 1164 individual case reports is for the most part voluntary, and reporting practices vary over time,

1165 between organizations and types of reporters. This may impact the types of adverse events that are
1166 reported, which information is captured, and how it is encoded. Inconsistencies and inaccuracies can
1167 lead to models that are less accurate, and systematic variability can reduce the generalizability of AI
1168 solutions to adjacent domains and make them more sensitive to data drift. They may also make it
1169 more difficult to ensure consistent performance across regions and organizations (see also Chapter
1170 on [Fairness & Equity](#)).

1171 Generally, the variable quality and consistency of individual case reports and the complex nature of
1172 the studied drug-event relationships may require more extensive human involvement than in other
1173 domains to ensure the validity and robustness of AI solutions for PV (see also the Chapter on [Human
1174 oversight](#)). The practice of PV is subject to medicines regulation, and regulatory expectations
1175 regarding validity and robustness may differ from those of the business itself. For example, even if
1176 from a business perspective an organization were prepared to adopt an AI solution whose output
1177 was not fully reproducible in repeated execution using the same input, the organisation would not
1178 consider the acceptability to regulators, as sufficient controls may exist to de-risk the process to
1179 which it is applied.

1180 Performance evaluation and testing are crucial considerations in ensuring the validity and robustness
1181 of AI solutions. While it will usually be more effective to account for these considerations also during
1182 development and training of AI solutions, to do so might not be feasible or required. For example,
1183 LLMs can be capable of zero-shot learning, with solid performance on language tasks for which they
1184 have not been specifically trained. What is important even then is to demonstrate adequate
1185 performance on the relevant tasks in independent testing with conditions reflecting the intended
1186 use.

1187 **Specification and design**

1188 **Use case and deployment domain**

1189 The intended use case and deployment domain for AI solutions in PV should be clearly defined, and
1190 the performance evaluation targeted to these, as far as possible. For example, in evaluating methods
1191 for PV signal detection, historical safety signals would typically be a more relevant basis for
1192 performance evaluation than well-known, already labelled adverse drug reactions since their
1193 reporting patterns differ in important ways¹⁹⁰ Similarly, if an AI solution for recognizing adverse
1194 events in free text is intended for broad use, its evaluation should include reports related to various
1195 medicinal products and adverse events, from both patients and health professionals, in relevant
1196 languages, etc. Ideally, there should be sufficient representation of relevant types of data in the test
1197 set to detect biases, promote adequate and generalizable performance across the intended
1198 deployment domain, assess usability, and identify circumstances where the model may
1199 underperform.

1200 **Multi-disciplinary collaboration**

1201 Ensuring the validity and robustness of AI solutions often requires collaboration across disciplines,
1202 including not only PV decision makers and practitioners, but also for example, data scientists and AI
1203 experts, and individuals with experience in computer systems validation. Diverse perspectives and
1204 expertise, in-depth understanding of a model's intended integration into the PV system and defined
1205 desired benefits and associated risks can help ensure that deployed AI solutions are effective, and
1206 address identified needs over their lifecycle.

1207 Applications of AI solutions pertaining to the complex relationships between drugs and adverse
1208 events often require a human-in-the-loop, especially in view of the variable quality and provenance
1209 of the underlying PV data. AI outputs in such applications need to be interpreted considering the
1210 broader clinical context, known pharmacological mechanisms, and possible alternative explanations

1211 which may not be captured in the data at hand and that the AI might not fully account for. Human
1212 intervention ensures that the final output is clinically meaningful and scientifically sound. On the
1213 other hand, more basic tasks such as redaction of personal data or drug and adverse event encoding
1214 may lend themselves to automation with minimal human intervention.

1215 **Definition of reference standards**

1216 Test sets must be aligned with intended deployment domain and able to demonstrate performance
1217 under realistic conditions. Reference standards relevant to the intended use need to be clearly
1218 defined and kept up to date. In many PV applications, these may be based on human execution of
1219 the task in question. Approaches to mitigate inconsistencies are often required, for example by
1220 having multiple human assessors annotate (parts of) the same data. When legacy human annotations
1221 are used as the reference standard, efforts should be made to clarify the definitions of relevant
1222 categories in the reference standard retrospectively, and to ensure that all included historical
1223 annotations adhere to these standards and are relevant for the intended future use. This may require
1224 the omission of available annotations that were developed following outdated principles or were
1225 based on different types of data. If reference standards are to be developed de novo, an explicit
1226 annotation guideline is recommended. This in turn may require a strengthening and clarification of
1227 existing processes and guidelines for human decision making on the PV task of interest, sometimes
1228 bringing value by harmonizing and making explicit decision processes that may otherwise remain
1229 implicit and variable within an organization. In smaller projects with a limited number of participants,
1230 special care may be required to ensure that annotations of the test set used for performance
1231 evaluation are independent of the development of the AI solution, for example annotations may be
1232 performed preferably by individuals with limited insights and vested interest in a specific AI solution
1233 to avoid conflicts of interest and confirmation bias. Similarly, if testing human-AI teams, the
1234 qualifications of human team member(s) should match those of the intended use case and
1235 deployment domain.

1236 Sometimes, boundaries between reference standard categories are not clear, which yields additional
1237 sources of possible ambiguity. For example, different organizations may have different requirements
1238 on how strong the conviction should be that two individual case reports refer to the same event for
1239 them to be classified as suspected duplicates. This may vary even within an organization depending
1240 on the intended use case, for example one may cast a wider net in highlighting suspected duplicates
1241 if each highlighted pair will be reviewed by a human before action and be more conservative if
1242 suspected duplicates will be automatically removed prior to statistical signal detection. Similar
1243 ambiguities exist in natural language processing tasks seeking to map free text to standard
1244 terminologies such as MedDRA where there may be multiple acceptable terms/codes for a specific
1245 verbatim, and it may be inappropriate to treat terms adjacent to the reference standard annotation
1246 as false positives.

1247 For unsupervised learning like cluster analysis and representation learning, and for applications of
1248 zero-shot learning like text summarization, formal test sets often cannot be obtained and other
1249 approaches to performance evaluation must be considered. In some cases, one may rely on human
1250 subjective review and assessment of an AI solution's output, but then potential biases must be
1251 considered and mitigated. For example, one may present the results of several different AI solutions
1252 to a blinded, domain expert and ask which they prefer. There are also performance evaluation
1253 approaches specifically designed for unsupervised learning like intruder detection analysis.¹⁹¹

1254 A general challenge in PV has been ensuring sustainable and reusable access to reference sets. As AI
1255 technologies rapidly advance, the necessity for consistent and frequent testing becomes increasingly
1256 important to safeguard against unintended consequences of AI usage. The potential for widespread
1257 impact of AI solutions in PV underscores the importance of maintaining up-to-date, accessible
1258 reference standards with clarity on how they were developed and related assumptions.

1259 Performance evaluation

1260 Performance evaluation is necessary for critical appraisal of AI solutions. The ability to carry out or
1261 assess performance evaluations are crucial skills for those who develop AI solutions and for those to
1262 whom AI solutions are proposed.

1263 Many of the metrics relevant to performance evaluation for AI solutions in PV come from
1264 information retrieval and apply primarily to use cases that can be viewed as binary classification
1265 tasks. In binary classification, we may refer to those instances that we want a method to retrieve as
1266 *positive controls* and those that we do not want it to retrieve as *negative controls*. We use this
1267 terminology throughout the description below (sometimes replacing positive controls by *target*
1268 *events*), acknowledging that other use cases may require different frameworks of evaluation, for
1269 example considering an ordering.

1270 AI solutions, like humans, will not always achieve perfect performance on complex tasks. Therefore,
1271 performance is typically assessed statistically for a sample of cases referred to as the test set,
1272 considering measures such as precision (= positive predictive value) and recall (= sensitivity) relative
1273 to the selected reference standard. The balance between precision and recall (and correspondingly
1274 between sensitivity and specificity) should be determined based on the relative costs of different
1275 types of errors (and utilities associated with correct decisions). Composite metrics like the F1 score
1276 (the harmonic mean of precision and recall) provide single-dimensional measures of predictive
1277 accuracy accounting for both precision and recall under some assumptions (for the F1 score that
1278 precision and recall are of equal importance and false positives as costly as false negatives). Test sets
1279 need to be large, diverse, and representative enough to reflect a sufficient portion of the intended
1280 deployment domain and to provide statistically robust estimates of performance. They should
1281 include different populations and consider possible scenarios in line with the intended use.²

1282 Since the primary interest is the expected performance of an AI solution in prospective use (as part
1283 of an overall system), performance evaluation should be independent of any data directly used
1284 during its development (this is in addition to any cross-validation or other separation of data for
1285 training and validation during development). Various potential sources of dependence between
1286 development and evaluation must be considered and eliminated, the most obvious being the risk
1287 that the same individual data points are considered in both phases. More subtle forms of
1288 dependence, can occur and lead to optimistic performance estimates, for example there may be a
1289 disproportional overlap in scope between the training and test sets compared with the deployment
1290 domain e.g. if training and test sets cover the same subset of drugs and adverse events, which can be
1291 referred to as a specific form of data leakage.¹⁹²

1292 Selection of machine learning models may also account for indirect performance characteristics such
1293 as an AI model's susceptibility to overfitting, computational cost and robustness to outliers,
1294 especially if test sets are not large and diverse enough to be reliably capture their impact during
1295 performance evaluation. When comparing different types of methods, any user-driven design
1296 decisions should be fixed and finalized before developers first access test sets. This is especially
1297 important for more complex methods with numerous analytical choices regarding model
1298 architecture, hyper-parameters, and model initialization.¹⁹³ On a related note, complex methods
1299 highly dependent on skilful design and deployment by human experts may not readily transfer to
1300 adjacent application areas without access to the same expertise. In routine deployment one is less
1301 concerned about whether one method is theoretically better than another but rather with which one
1302 is likely to perform best for a given purpose, irrespective of what design/analytical choices one made.

² For a continually updated inventory, see for example <https://oecd.ai/en/catalogue/metrics>

1303 **Benchmarking**

1304 Ideally, performance should be compared against relevant benchmark methods, if available. For
1305 example, AI-based signal detection methods may be compared against standard disproportionality
1306 measures. In the case of more complex benchmark methods, including those based on AI models,
1307 special care must be taken to ensure that the benchmark methods have been appropriately
1308 instantiated and fine-tuned to the task at hand to serve as a relevant comparator.

1309 When public benchmark test sets exist, performance may be evaluated against these, ideally as a
1310 complement to performance evaluation targeted to the deployment domain of interest. At present,
1311 public benchmarks exist only for some specific applications in PV. They include sets of emerging
1312 safety signals,^{194,195} sets of established adverse drug reactions,^{196,197,198,199,200} and clinically relevant
1313 drug-drug interactions.²⁰¹ However, continual access to benchmark reference sets over time can be a
1314 challenge and the degree to which they are maintained and kept up to date varies.

1315 To complement overall performance estimates, subgroup analyses can provide useful information on
1316 the strengths and weaknesses of the AI solution for different parts of the deployment domains (See
1317 also chapter on [Fairness & Equity](#)). Along the same lines, sensitivity analyses can help assess the
1318 robustness of the AI solution and its evaluation to variations in specification and design.

1319 **Special considerations for low-prevalence settings**

1320 Many PV applications focus on rare patterns and events. For example, in a case retrieval task most
1321 reports will typically not be relevant for a given topic, such as pregnancy, medication errors, positive
1322 rechallenge interventions, or drug-induced liver injury. Similarly, for PV signal detection, most drug-
1323 event combinations are not true adverse drug reactions, let alone recently detected safety signals.
1324 Managing and analyzing these rare events effectively requires reliable reference datasets, however,
1325 existing resources, such as SIDER, are often limited by outdated and static information, underscoring
1326 the need for alternative solutions.²⁰² As an even more extreme example, pairs of duplicate reports
1327 are vanishingly rare among all possible pairs of reports in large collections of individual case reports –
1328 if 10% of the reports in a database of 1 million reports have a (single) duplicate, the chance that a
1329 randomly selected pair would be duplicates is only 1 in 10 million.³

1330 This low prevalence of positive controls (i.e. class imbalance) limits our ability to achieve accurate
1331 performance evaluation and requires special care and consideration. For example, a balance may
1332 need to be struck between the quality of each annotation and the resulting size of the test sets (or
1333 the cost/time to develop them), for example related to whether double annotations by multiple
1334 assessors are feasible to increase quality or evaluate consistency. The heterogeneity in skills and
1335 preferences among different human reviewers poses a significant challenge to achieving consistent
1336 quality in annotations. The use of intelligent automation to support decision making can mitigate
1337 inconsistency and reduce subjective bias in the evaluation process. Moreover, straight random
1338 samples of test cases often contain too few positive controls whereas test sets enriched with positive
1339 controls can lead to misleading estimates of precision and recall.

1340 *Recall* measures how many of the target events are correctly identified (*recalled*) by the AI solution.
1341 *Sensitivity* is a synonym. If heuristics are used to increase the proportion of target events in the test
1342 set, then recall may be over-estimated since target events which are harder to identify for the AI
1343 solution, may less likely be included. This does not mean that rebalancing approaches should
1344 necessarily be avoided but if they are used, this should be acknowledged and critically assessed.
1345 *Precision* is the proportion of target events among all events highlighted by the AI solution. *Positive*
1346 *predictive value* (PPV) is a synonym. It is highly dependent on the prevalence of target events in the
1347 test set, and if test sets have been enriched with target events, naive test set precision estimates will
1348 be optimistic. For reliable precision estimates, the prevalence of positive controls in the test sets

³ $0.10 \cdot 1/10^6$

1349 should match as far as possible that of the intended deployment. For a specific AI solution, this is
1350 straightforward to obtain by applying the AI solution to a random sample and annotating all
1351 highlighted instances. However, such test sets are tied to the AI solution in question and will need to
1352 be developed again, or at least extended, if the solution is modified.

1353 Estimates of precision and recall depend on the selected decision threshold, and performance
1354 evaluation should be targeted at decision thresholds relevant to the intended deployment domain,
1355 i.e. with a relevant balance between false positives and false negatives.

1356 **Beyond summary statistics**

1357 Summary statistics as captured by the metrics described in the previous section go only so far in
1358 enabling us to assess and understand the performance of an AI solution. Access to and ability to
1359 inspect representative, concrete examples of an AI solution's classification of individual instances in a
1360 test set is also valuable. Examining false positives and false negatives in an error analysis step can
1361 each give useful insights regarding the strengths and limitations of the AI solution and its evaluation.
1362 For example, if a false negative in de-identification corresponds to a full name preceded by 'Mr', this
1363 may undermine end-users' trust in the solution, even if overall recall is excellent. On the other hand,
1364 if the false negative is 'AF' and it is hard to know from the surrounding text if these are initials or an
1365 abbreviation for *atrial fibrillation*, then the overall precision metric may be viewed as potentially
1366 conservative. Review of correctly classified instances may in turn give insights regarding an AI
1367 solution's capacity to solve challenging tasks. Does it correctly classify more difficult cases or just the
1368 trivial ones? This may be especially important when there is no baseline comparator, and we may not
1369 understand from overall performance metrics the difficulty of the task at hand. When there is a
1370 baseline comparator method, one may focus on instances that are differentially classified by the two
1371 methods, to better understand the nature of any improved performance of the proposed solution
1372 over the comparator.

1373 **Reproducibility**

1374 Reproducibility for an AI solution requires that it will generate the same output for a given input.
1375 Predictive models like support vector machines and decision trees fulfil this requirement. So do
1376 certain LLMs and other deep neural networks, once their weights have been fixed at the end of
1377 training / fine-tuning, even though their model fitting may include stochastic components so that
1378 new weights may result if a model is re-trained on the same data.

1379 GenAI solutions on the other hand include stochastic components also in their execution and will
1380 typically generate different outputs for the same prompt, without changes to the underlying models.
1381 The same is true for other methods such as mixture model-based cluster analysis and semantic
1382 vector representations of adverse events. For such solutions, stability is a key additional performance
1383 metric, reflecting how similar the results of fully replicated analyses are. While replicability of results
1384 can sometimes be artificially ensured through seeding the pseudo random number generator, this
1385 can be non-trivial to do for proprietary models and does not improve the inherent (in)stability of the
1386 AI solution, which should be evaluated.

1387 Reproducibility, as described here, pertains to an organization that has full access to a specific AI
1388 model and the relevant reference sets. Full reproducibility by, say, the broader scientific community,
1389 in addition requires transparency.

1390 **Assessing artificial intelligence solutions with human-in-the-loop**

1391 Many AI solutions aim for intelligence augmentation, i.e. to support and enhance human decision
1392 making. In this context, the relevant focus of performance evaluation would be of the human-AI
1393 team. To date, we have limited experience of such studies in PV applications, but at a minimum, they

1394 would need to account for the variability in skills and preferences between different human
1395 members of the team. Defining a relevant test set may also present new challenges: for example, for
1396 signal detection applications, human domain experts could not be blinded to historical safety signals;
1397 and it may be difficult to obtain a reference standard if the aim is for the human-AI team to *exceed*
1398 the quality of classification by unassisted human domain experts.

1399 What constitutes acceptable performance may need to account for how the AI solution is integrated
1400 with the PV system and whether there is a human in the loop.²⁰³ For example, performance
1401 evaluation for an NLP-based system to identify and extract adverse events from source documents
1402 might in addition to the overall performance evaluation consider whether errors can be readily
1403 spotted in the results and whether the end-to-end hybrid process performs better than a fully
1404 manual approach (for an example see Park et al 2023²⁰⁴).

1405 **Continuous integration and deployment**

1406 Deployed models should be monitored in real-world use with a focus on maintained or improved
1407 performance. There may be reasons to revise and update performance criteria in production as the
1408 business understanding of the task is refined or the conditions for the task itself change due to
1409 external factors.

1410 For deployed AI solutions that incorporate ML components, there should be appropriate processes
1411 and quality controls for periodical re-training to manage risks of performance degradation or
1412 negative impact from dataset drift. In some instances, the retraining may consist of incremental fine-
1413 tuning within existing model architectures whereas more substantial changes to the deployment
1414 domain may require changes to the architecture of the AI solution. The latter could result from a
1415 change in scope from medicines to vaccines, revisions of the underlying medical terminologies or
1416 data structures, updated regulation or conventions and more.

1417 Continual performance evaluation can be relevant regardless of whether an AI solution incorporates
1418 ML components or not. Its frequency should follow the risk-based approach and depending on the
1419 application, may include data-driven safeguards to identify, for example, substantial data drift or
1420 performance degradation triggering remedial actions that could include additional evaluation, and
1421 possible retraining, stopping use of the algorithm and/or introducing quality control measures to
1422 maintain confidence in its results. Documentation of activities and acceptance criteria for re-
1423 introducing AI solutions under such circumstances may also be required. As an example using
1424 mechanisms such as 'ground truths', i.e. known input/output pairs which are checked each time an AI
1425 system undergoes a change, or mechanisms to guard against automation bias.

1426 One of the potential benefits of ML is the ability to improve performance through iterative
1427 modifications, including by learning from real-world data. To support this approach, the US FDA,
1428 Health Canada, and MHRA described a “Predetermined Change Control Plan” for ML-enabled device
1429 software functions (ML-DSF). Their general principles might conceivably be applied to AI solutions in
1430 PV. A Predetermined Change Control Plan generally includes: 1) a detailed description of the specific,
1431 planned modifications; 2) the associated methodology to develop, validate, and implement those
1432 modifications in a manner that ensures the continued acceptable performance of the algorithm; and
1433 3) an Impact Assessment of the benefits and risks of the planned modifications and risk mitigations.
1434 The detailed description of the planned modification should include changes to the characteristics
1435 and performance of the algorithm resulting from the implementation of the modifications. An
1436 example of a modification might include retraining a ML model. A protocol providing the details of
1437 the data and methods used to develop, evaluate, and implement such a modification should be
1438 created and adhered to. An Impact Assessment of the modification should be carried out and risk
1439 mitigation measures developed to ensure that any identified risks will be controlled. This approach
1440 should be further incorporated into the quality management system governing the PV process being
1441 modified.

1442 There should be straightforward means to report issues or anomalies encountered, and these should
1443 be addressed promptly. Ideally, the response would include acknowledging receipt of feedback,
1444 providing updates on investigations, and implementing necessary changes to the AI system.

References

- ¹⁹⁰ Norén GN, Caster O, Juhlin K, Lindquist M. Zoo or savannah? Choice of training ground for evidence-based pharmacovigilance. *Drug Safety*. 2014;Sep;37:655-659. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-014-0198-z>
- ¹⁹¹ Jonathan Chang, Jordan Boyd-Graber, Sean Gerrish, Chong Wang, et al. (2009) Reading Tea Leaves: How Humans Interpret Topic Models. [Unpublished manuscript] ([PDF](#) accessed 27 April 2025)
- ¹⁹² Rudin C, Radin J. Why Are We Using Black Box Models in AI When We Don't Need To? A Lesson from an Explainable AI Competition. *Harv Data Sci Rev*. 2019;1;1-9. ([Journal full text](#)) <https://doi.org/10.1162/99608f92.5a8a3a3d>
- ¹⁹³ Duin RP. A note on comparing classifiers. *Pattern Recognition Letters*. 1996;May1;17(5):529-536. ([Journal abstract](#)) [https://doi.org/10.1016/0167-8655\(95\)00113-1](https://doi.org/10.1016/0167-8655(95)00113-1)
- ¹⁹⁴ Harpaz R, Odgers D, Gaskin G, DuMouchel W, Winnenburgh R, Bodenreider O, Ripple A, Szarfman A, Sorbello A, Horvitz E, White RW. A time-indexed reference standard of adverse drug reactions. *Scientific Data*. 2014;Nov 11;1(1):1-0. ([Journal full text](#)) <https://doi.org/10.1038/sdata.2014.43>
- ¹⁹⁵ Sartori D, Aronson JK, Norén GN, Onakpoya IJ. Signals of adverse drug reactions communicated by pharmacovigilance stakeholders: a scoping review of the global literature. *Drug Safety*. 2023;Feb;46(2):109-120.) <https://doi.org/10.1007/s40264-022-01258-0>
- ¹⁹⁶ Coloma PM, Avillach P, Salvo F, Schuemie MJ, Ferrajolo C, Pariente A, et al. A reference standard for evaluation of methods for drug safety signal detection using electronic healthcare record databases. *Drug Saf*. 2013;Jan;36(1):13-23. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-012-0002-x>
- ¹⁹⁷ Ryan PB, Schuemie MJ, Welebob E, Duke J, Valentine S, Hartzema AG. Defining a reference set to support methodological research in drug safety. *Drug Saf* 2013;Oct;36(Suppl 1):S33-S47. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-013-0097-8>
- ¹⁹⁸ Kuhn M, Letunic I, Jensen LJ, Bork P. The *SIDER* database of drugs and side effects. *Nucleic Acids Research*. 2016;Jan4;44(D1):D1075-9. ([Journal full text](#)) <https://doi.org/10.1093/nar/gkv1075>
- ¹⁹⁹ Demner-Fushman D, Shooshan SE, Rodriguez L, Aronson AR, Lang F, Rogers W, Roberts K, Tonnig J. A dataset of 200 structured product labels annotated for adverse drug reactions. *Scientific Data*. 2018;Jan30;5(1):1-8. ([Journal full text](#)) <https://doi.org/10.1038/sdata.2018.1>
- ²⁰⁰ Bayer S, Clark C, Dang O, Aberdeen J, Brajovic S, Swank K, Hirschman L, Ball R. ADE eval: an evaluation of text processing systems for adverse event extraction from drug labels for pharmacovigilance. *Drug Safety*. 2021;Jan;44:83-94. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-020-00996-3>
- ²⁰¹ Kotsioli E, Maskell S, Dutta B, et al. A reference set of clinically relevant adverse drug-drug interactions. *Sci Data* 9, 72 (2022). <https://doi.org/10.1038/s41597-022-01159-y>
- ²⁰² Painter J, Powell G, Bate A. PVLens: Enhancing Pharmacovigilance Through Automated Label Extraction. *arXiv:2503.20639v2 [cs.CL]* <https://doi.org/10.48550/arXiv.2503.20639>
- ²⁰³ The US Food and Drug Administration. *Good machine learning practice for medical device development: guiding principles*. The US Food and Drug Administration: Silver Spring, MD, USA, 2021. ([Webpage](#) accessed 21 March 2025)
- ²⁰⁴ Park J, Djelassi M, Chima D, Hernandez R, Poroshin V, Iliescu AM, Domalik D, Southall N. Validation of a natural language machine learning model for safety literature surveillance. *Drug Safety*. 2024;Jan;47(1):71-80. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-023-01367-4>

1445 **Chapter 6: Transparency**

1446 *Principle*

1447 Transparency regarding AI involves disclosing information between organizations or individuals. This
1448 includes sharing relevant documentation of the AI system lifecycle (i.e. design, development,
1449 evaluation, deployment, operation, re-training, maintenance and decommission) to facilitate
1450 traceability and providing stakeholders with enough information to have a general understanding of
1451 the AI system, its use, risks, limitations, and impact on their rights.

1452 *Key messages*

- 1453 • Declaring when and how AI solutions are used for core PV tasks is critical for building trust
1454 among domain experts, decision makers, regulatory authorities, and the public.
- 1455 • The nature of AI solutions deployed for core PV tasks should be described including their
1456 model architectures, expected inputs and outputs, and the level and type of human-
1457 computer interaction.
- 1458 • To give a comprehensive picture of an AI solution’s effectiveness and limitations, the
1459 presentation of performance evaluation results should describe the scope and nature of the
1460 test set(s) used including definitions of their reference standards and sampling strategies.
1461 Presented performance metrics should be relevant for the intended deployment domain,
1462 compared with relevant benchmarks, and complemented by qualitative review of
1463 representative examples of correct and incorrect output.
- 1464 • If possible, a description of the general principles and logic by which an AI model functions
1465 and arrives at its outcomes / predictions should be shared, or the lack of such explainability
1466 should be acknowledged and its implications discussed.

1467 **Introduction**

1468 Transparency provides stakeholders with relevant information regarding the nature and use of an AI
1469 solution. It reflects what information is shared with key stakeholders by those who develop or deploy
1470 it. The main purposes of transparency are to build trust, to enable individuals and organizations not
1471 involved in their development to inspect and scrutinize the design and performance of AI solutions,
1472 and to ensure regulatory compliance.

1473 As further elaborated on in the chapter on [Governance & Accountability](#), the primary direction of
1474 transparency and disclosure of information varies during the phases of the AI solution lifecycle. For
1475 example, during the design phase the organization should be transparent toward developers
1476 regarding the specification and requirements for an AI solution, whereas the main direction of
1477 transparency is the opposite in the pre-deployment phase. During routine use, the most important
1478 form of transparency may be from the organization toward end users (and in some cases regulatory
1479 authorities).

1480 **Disclosing use of artificial intelligence**

1481 It is essential to disclose why, when and how AI is being used in different PV tasks. This is to maintain
1482 trust and accountability among stakeholders, including developers, PV professionals and decision
1483 makers, regulatory authorities, healthcare professionals, and patients.

1484 Regulatory bodies require disclosure of AI use to assure compliance with applicable laws and
1485 regulations. To this end, software vendors and internal development groups need to be transparent
1486 toward PV organizations, who in turn need to be transparent toward regulatory authorities. At the
1487 same time, those individuals who utilize AI solutions to process or analyze PV data must be informed
1488 about the AI’s role in their workflows to help them integrate AI into their processes in an informed

1489 manner to support its effective application and ensure that they can identify any issues arising from
1490 AI use.

1491 PV professionals should also communicate the provenance of data elements and whether AI
1492 solutions contributed to their capture or development. Human interpretation of PV data may depend
1493 on how it was ascertained. For example, signal assessors may lend different weight to a case
1494 narrative that was auto generated from structured elements compared with one that documents the
1495 patients or health professionals' verbatim description of the adverse event. There is also a risk of a
1496 vicious circle where AI generated information is used as part of a reference standard in subsequent
1497 AI model development, if its provenance is not properly disclosed.

1498 **Transparency regarding the artificial intelligence model**

1499 Ensuring transparency of the AI models used in PV (to the extent possible), is critical to fostering
1500 trust, facilitating informed decision making, and ensuring that these models are applied
1501 appropriately. Ideally, transparency should be extended to also capture decisions made by PV
1502 professionals resulting from the AI model. Model transparency is not only a technical requirement
1503 but also an ethical imperative, ensuring that all parties understand the tools they are working with
1504 and can make informed decisions based on their outputs. Below are key aspects of an AI model that
1505 should be disclosed to stakeholders. The rationale behind the design choices should also be
1506 explained, to help ensure that the model is aligned with its intended use and stakeholder needs.

1507 **Table 3: Key aspects of an artificial intelligence model to disclose to stakeholders**

1508 Source: CIOMS Working Group XIV

1509

Intended Use	The intended use of each AI model should be clearly defined and communicated. This includes specifying the PV tasks the model is designed to assist with or perform, such as adverse event recognition in free text, signal detection, or case triage.
Human-Computer Interaction	The level and type of interaction between humans and the AI models should be communicated. This includes specifying whether the AI model is executed autonomously, has a human in-the-loop (and what their required competence would be), or aims to provide decision support to down-stream human specialists.
Model architecture	The type of AI model and its general architecture should be disclosed, such as whether it is rule based, uses linear models, or specific types of neural networks, or combines different ML models in an ensemble, etc. Additionally, relevant details about the model's structure, such as the type and depth of a neural network architecture, should be shared.
Model parameters	At a minimum, key predictors or features that drive the decisions of an AI model should be disclosed, if they are known. If feasible, the full set of model weights and parameters can be shared, to enable external replication and external performance evaluation. See further discussion regarding this, at the end of this section. For AI solutions based on GenAI, any predefined prompts should be specified along with any pre- or post-processing steps.
Explainability	If possible, a description of the general principles and logic by which an AI model functions and arrives at its outcomes / predictions should be shared, or the lack of explainability should be acknowledged and its implications discussed. (See also Section on Explainability).
Training set	Details about the training set(s) based on which ML components have been developed should be disclosed. This would include their size, scope, and creation date, along with reflections on how well they align with the intended deployment domain.

Standard AI Components	If the AI model incorporates public standard components, such as pre-trained ML models, libraries, or frameworks, or datasets, this should be disclosed, including the specific versions used, date of access, and any custom parameter settings.
Acceptable Inputs	The types of inputs that the AI model expects should be specified. This provides insights regarding the basis for the AI model's outputs and ensures that it is only fed data it is designed to handle, thereby maintaining the accuracy and reliability of its outputs.
Type(s) of Output	The types of output generated by the AI model should be described. Examples may be risk scores, classifications, alerts, or free text.
Known Limitations	Any known limitations regarding the nature of the AI model should be communicated, including e.g. features or types of interactions, which it is unable to account for.

1510 To allow other developers and researchers to fully replicate an AI model and possibly even modify it
1511 for further use, an organization might choose to publish its full set of parameters and weights or
1512 even share the source code. This level of openness supports peer review and validation by external
1513 experts, which can enhance trust in the model's reliability and foster innovation. However, it will not
1514 always be feasible due to considerations regarding intellectual property, competitive advantage, or
1515 the sheer complexity of large models. Moreover, for many stakeholders, access to raw code and
1516 parameters of a complex AI model may not enhance their understanding and will need to be
1517 complemented by the other measures for model transparency described above. Understanding the
1518 rationale, assumptions, and subjective decisions made in the implementation can be more important
1519 for gaining meaningful insights into the model's function and effectiveness. For full scientific
1520 reproducibility, developers may also need to share the relevant reference sets, at least those used
1521 for performance evaluation. However, depending on the use case and stakeholders involved this may
1522 conflict with the data privacy principle.

1523 Explainability

1524 A specific form of transparency relates to disclosure of the general principles and logic by which an AI
1525 solution operates and has arrived at a specific output. This may help nurture trust, allow affected
1526 individuals to understand and influence outcomes, support down-stream human decision making
1527 and facilitate human oversight and regulatory compliance. In this context, *explainability* and
1528 *interpretability* are important concepts, which partly overlap.

1529 The set of *Guidelines on the testing of AI-based systems* in the ISO standard for Software testing in
1530 Software and systems engineering, characterizes explainability as a "*level of understanding how the*
1531 *AI-based system ... came up with a given result*" and interpretability to a "*level of understanding how*
1532 *the underlying (AI) technology works*".²⁰⁵

1533 Similarly, the AI Risk Management Framework of the US National Institutes of Standards and
1534 Technology includes the following statement: "Explainability refers to a representation of the
1535 mechanisms underlying AI systems' operation, whereas interpretability refers to the meaning of AI
1536 systems' output in the context of their designed functional purposes".²⁰⁶

1537 The Organisation for Economic Co-operation and Development (OECD) Transparency and
1538 Explainability Principle 1.3 states:²⁰⁷

1539 "Explainability means enabling people affected by the outcome of an AI system to understand how it
1540 was arrived at. This entails providing easy-to-understand information to people affected by an AI
1541 system's outcome that can enable those adversely affected to challenge the outcome, notably – to the
1542 extent practicable – the factors and logic that led to an outcome."

1543 For the context of this report, we adopt a similar perspective and use *explainability* in a broader
 1544 sense to reflect the degree to which humans can understand the factors and logic that have led to a
 1545 specific outcome or that play a role in the general operation of an AI solution.

1546 Concrete examples which illustrate the role of explainability in different PV use cases are provided in
 1547 the [Appendix 4](#).

1548 **Benefits of explainability**

1549 Explainability can be beneficial because it may:

- 1550 • Nurture trust in an AI solution, by enabling stakeholders to make sense of and contextualize
 1551 an AI solution's output;
- 1552 • Allow individuals affected by an AI solution's output to challenge and influence the outcome;
- 1553 • support and speed up human decision-making which builds on or integrates an AI solution
 1554 output;^{208,209}
- 1555 • Propose scientific hypotheses for consideration by end users - individual or combinations of
 1556 features such as drugs, diseases, and demographics that are included in the proposed
 1557 explanation of the findings may provide signals of adverse drug reactions, and adverse drug-
 1558 disease interactions worthy of evaluation, as well as potential biological mechanisms of
 1559 adverse drug reactions;²¹⁰
- 1560 • Enable more complete documentation, audit, and human oversight of AI solutions;
- 1561 • Contribute to regulatory compliance especially when it is possible to retain and examine the
 1562 human decision together with the AI output and the explanation upon which the decision
 1563 was based;
- 1564 • Facilitate troubleshooting by revealing issues such as possible biases or likely spurious
 1565 correlations;^{211,212}
- 1566 • Contribute towards model assessment and selection by uncovering what is causing different
 1567 models trained on the same data to perform differently.

1568
 1569 Referring to the definition above, the individuals who could challenge the output of the PV AI system
 1570 and require explainability are more likely to be stakeholders who are directly involved in the PV
 1571 process rather than members of the public.²¹³ They may range from the PV and quality assurance
 1572 staff who are directly interacting with the AI, the developers who are building or maintaining an AI
 1573 system to the regulators who are inspecting it. Examples on how different stakeholders in the PV
 1574 process can benefit from explainability are provided in [Appendix 4](#).

1575 **Inherent vs post hoc explainability**

1576 AI models of limited complexity may be inherently explainable, allowing the basis for their output to
 1577 be deduced from direct inspection of their model architectures and parameters.²¹⁴ This is also
 1578 referred to as *ante-hoc* explainability. Examples may include simple decision trees, rule-based
 1579 classifiers, and regression models.

1580 In contrast, a growing field of research seeks to obtain *post-hoc* explainability (referred to by the
 1581 acronym xAI) for more opaque AI solutions, including deep neural networks with complex
 1582 architectures and more parameters than a human can survey or comprehend. With such approaches,
 1583 a separate layer of methods and techniques are applied top of the AI solution²¹⁵ to trace and explain
 1584 the basis for a specific, already generated output. Some xAI approaches seek to explain the output of
 1585 complex AI models by estimating relative feature importance and others do so by determining the
 1586 minimal change in one or more features required to change a given output. There are also xAI
 1587 methods that provide explainability for a specific output by fitting simpler, inherently interpretable

1588 models to the local context of a specific output. For examples of specific xAI methods in use at the
1589 time of writing this report, please see [Appendix 4](#).

1590 When an xAI method is used to gain explainability, it must itself comply to the applicable regulatory
1591 requirements for computerised systems. In other words, the xAI method must be validated and
1592 proven to be fit for purpose, and processes must be in place to ensure that it continues to be so.
1593 Explanations provided by xAI are an ‘approximate understanding of the relationship between the
1594 data and the predictions’ according to Pinheiro et al (2022).²¹⁶ The explanations can be imperfect or
1595 incomplete and/or provide only a partial explanation.

1596 **Challenges related to explainability**

1597 Stakeholders are advised to critically consider what type of explainability is required for the intended
1598 use case, for whom, and for what purpose, and whether the system they are considering can provide
1599 it.²¹⁷

1600 The level of explainability of an AI solution’s output should not be the sole determining factor for
1601 model selection. Some applications (e.g. machine translation) may not require inherent explainability
1602 and may depend on the capabilities and improved performance offered by more complex AI
1603 solutions. In such cases, the negative impact of limited explainability may be mitigated by ensuring
1604 high transparency regarding other aspects of the AI solution and its performance, coupled with extra
1605 care to achieve validity and robustness and human oversight.²¹⁸ At the same time, it should not be
1606 assumed that explainability necessarily leads to lower performance and that a trade-off between the
1607 two needs to be made.²¹⁹

1608 Similarly, while explainability of an AI solution’s output can sometimes help identify issues with
1609 validity & robustness or fairness & equity, explainability alone does not prove that the system is fit
1610 for purpose, nor does it vouch for the trustworthiness of the system.²²⁰ It attempts to clarify what
1611 factors led to a specific output but is not indicative of an AI solution’s general performance or of its
1612 fairness and equity. For example, even if an inherently explainable AI solution does not include age
1613 as one of its explicit features, it could bias against an age group, if this bias is mediated by other
1614 features. In fact, explanations may make stakeholders more susceptible to overreliance on model
1615 outputs, so called automation bias.²²¹ Also, explainability is no guarantee of transparency – an
1616 organization may, for example, choose not to disclose the key features and inner logic of an
1617 inherently explainable model such as a decision tree.

1618 Explainability is not the same for all stakeholders. What is understood by model developers could be
1619 incomprehensible for other stakeholders²²² and cognitive capabilities must be considered to ensure
1620 that explanations are comprehensible to humans.²²³ Since humans will need to process and
1621 contextualize any explanations provided, they should also be informed about and aware of *their own*
1622 possible biases and blind spots which may influence their ability to leverage the explanations.
1623 Related to this, it should be noted that, in the worst case, a plausible explanation for an incorrect AI
1624 output may increase the likelihood that it is accepted without the appropriate critical review by some
1625 end users.

1626 **Transparency regarding performance**

1627 Transparency regarding an AI solution’s assessed performance not only communicates how well an
1628 AI model operates in practice but also provides stakeholders with insights into the design,
1629 implementation, and decision-making processes behind the model. As such, it provides a bridge
1630 between theoretical capability and practical utility. Without a clear view of how an AI solution
1631 behaves under realistic conditions, stakeholders cannot fully assess its suitability for use or be
1632 confident in its robustness and validity. Performance transparency ensures that all stakeholders,
1633 from end-users to regulatory authorities, have a clear understanding of an AI solution’s strengths,

1634 limitations, and expected behaviour in the contexts where it will be deployed. This is particularly
 1635 important in PV, where AI systems support information processing and decision making, with the aim
 1636 of safeguarding patient safety and public health.

1637 By recording and being able to share detailed performance evaluations with relevant stakeholders,
 1638 organisations offer clarity on the strengths and limitations of the AI system, as well as the key
 1639 assumptions made during their design, including quantitative metrics, qualitative examples, and
 1640 comparisons to benchmarks, and thereby organizations provide the necessary context to build trust
 1641 and appropriate reliance on AI systems. This transparency allows for informed decision making,
 1642 ensures that AI systems are used within their intended scope, and helps identify areas where
 1643 adaptations or special measures may be required. Additionally, it supports continuous improvement
 1644 by highlighting areas where the model may need further refinement or retraining.

1645 In support of this, there should exist a clear documentation of the data used for performance
 1646 evaluation, including data acquisition, cleaning and transformation, and processes for managing
 1647 missing or erroneous data.

1648 Table 4 outlines relevant aspects to disclose to ensure transparency regarding the estimated
 1649 performance of an AI solution. For further elaboration, see the Chapter on [Validity & Robustness](#).

1650 ***Table 4: Relevant aspects to disclose to ensure transparency regarding the estimated***
 1651 ***performance of an artificial intelligence solution***

1652 Source: CIOMS Working Group XIV
 1653

Scope of evaluation	Describe the nature of the reference sets used for performance evaluation, acknowledging any known deviations from the intended deployment domain (e.g. over- or under-representation of certain drugs, adverse events, patient populations etc.). Relevant information would include the types of data and from where they have been derived.
Sampling	Describe the prevalence of positive and negative controls in the reference set and how this relates to the intended use. If they are different, describe how performance evaluation was adjusted to account for this. Describe any use of data augmentation for performance evaluation.
Reference standard	Disclose the definitions of different categories of classification used in performance evaluation (for example, positive and negative controls in a binary classification task). Share any annotation guidelines used to improve quality and consistency of human annotations in developing the reference standard.
Human input	Describe the qualifications of human assessors contributing to test set development and any use of parallel annotations and evaluations of concordance during this phase. If the AI solution includes a human-in-the-loop during operation, then state the qualifications of those individuals who participated during performance evaluation.
Summary metrics	Present standard performance evaluation metrics when suitable or motivate the use of customized metrics. Place emphasis in this presentation on levels of the decision threshold relevant to the intended use and deployment domain (e.g. with a realistic balance between false positives and false negatives). Complement composite performance metrics with their components (e.g. precision & recall for an F-score).
Benchmarks	Present comparisons against relevant benchmark methods (including human-level performance) and/or standard benchmark reference sets, when available.
Subsets & sensitivity analyses	Present the results of any subset or sensitivity analyses during performance evaluation or acknowledge the lack thereof.
Qualitative review	Provide representative examples of correct classifications and representative examples of incorrect classifications (false positives and false negatives).

References

- ²⁰⁵ International Organization for Standardization (ISO). ISO/IEC TR 29119-11:2020. Software and systems engineering — Software testing. Part 11: Guidelines on the testing of AI-based systems. ([Abstract accessed 21 March 2025](#))
- ²⁰⁶ National Institute of Standards and Technology. AI RMF Trustworthy & Responsible AI Resource Center. Gaithersburg, MD: National Institute of Standards and Technology. ([Webpage accessed 21 March 2025](#))
- ²⁰⁷ Organisation for Economic Co-operation and Development. Transparency and Explainability (Principle 1.3). OECD AI Policy Observatory. ([Webpage accessed 21 March 2025](#))
- ²⁰⁸ Spiker J, Kreimeyer K, Dang O, Boxwell D, Chan V, Cheng C, Gish P, Lardieri A, Wu E, De S, Naidoo J. Information visualization platform for postmarket surveillance decision support. *Drug Safety*. 2020;Sep;43:905-915. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-020-00945-0>
- ²⁰⁹ Botsis T, Dang O, Kreimeyer K, Spiker J, De S, Ball R. A Decision-Support Platform Powered by AI and Humans-in-the-Loop Boosts Efficiency and Assures Quality in FDA's Pharmacovigilance. *International Society of Pharmacovigilance, 23rd Annual Meeting 2024, Montreal, Canada*. ([Webpage accessed 21 March 2025](#))
- ²¹⁰ Hauben M. Artificial intelligence in pharmacovigilance: Do we need explainability? *Pharmacoepidemiol Drug Saf*. 2022;Dec31(12):1311-1316. ([Journal full text](#)) <https://doi.org/10.1002/pds.5501>.
- ²¹¹ Albahri AS, Duham AM, Fadhel MA, Alnoor A, Baqer NS, Alzubaidi L, Albahri OS, Alamoodi AH, Bai J, Salhi A, Santamaría J. A systematic review of trustworthy and explainable artificial intelligence in healthcare: Assessment of quality, bias risk, and data fusion. *Information Fusion*. 2023;Aug1;96:156-191. ([Journal full text](#)) <https://doi.org/10.1016/j.inffus.2023.03.008>.
- ²¹² Ribeiro MT, Singh S, Guestrin C. "Why Should I Trust You?": Explaining the Predictions of Any Classifier [Internet]. [arXiv:1602.04938](https://arxiv.org/abs/1602.04938) 2016. ([Journal full text](#)) <https://doi.org/10.48550/arXiv.1602.04938>
- ²¹³ Hauben M. Artificial intelligence in pharmacovigilance: Do we need explainability?. *Pharmacoepidemiology and Drug Safety*. 2022;Dec;31(12):1311-1316. ([Journal full text](#)) <https://doi.org/10.1002/pds.5501>
- ²¹⁴ Ghassemi M, Oakden-Rayner L, Beam AL. The false hope of current approaches to explainable artificial intelligence in health care. *The Lancet Digital Health*. 2021;Nov1;3(11):e745-50. ([Journal full text](#)) [https://doi.org/10.1016/S2589-7500\(21\)00208-9](https://doi.org/10.1016/S2589-7500(21)00208-9)
- ²¹⁵ Reddy S. Explainability and artificial intelligence in medicine. *Lancet Digit Health*. 2022;Apr;4(4):e214–215. Cited by: Cutillo CM, Sharma KR, Foschini L, Kundu S, Mackintosh M, Mandl KD. Machine intelligence in healthcare—perspectives on trustworthiness, explainability, usability, and transparency. *npj Digit Med*. 2020;Mar;26;3(1):1-5. ([Journal full text](#)) <https://doi.org/10.1038/s41746-020-0254-2>
- ²¹⁶ Pinheiro LC, Kurz X. Artificial intelligence in pharmacovigilance: a regulatory perspective on explainability. *Pharmacoepidemiol Drug Saf*. 2022;Dec1;31(12):1308-1310. ([Journal full text](#)) <https://doi.org/10.1002/pds.5524>
- ²¹⁷ Royal Society. Explainable AI. [Royalsociety.org](https://royalsocietypublishing.org/journal/rsos). 2024. ([Website accessed 11 August 2024](#))
- ²¹⁸ Ghassemi M, Oakden-Rayner L, Beam AL. The false hope of current approaches to explainable artificial intelligence in health care. *The Lancet Digital Health*. 2021;Nov1;3(11):e745-50. ([Journal full text](#)) [https://doi.org/10.1016/S2589-7500\(21\)00208-9](https://doi.org/10.1016/S2589-7500(21)00208-9)
- ²¹⁹ Rudin C, Radin J. Why are we using black box models in AI when we don't need to? A lesson from an explainable AI competition. *Harvard Data Science Review*. 2019;Nov22;1(2):1-9. ([Journal full text](#)) <https://doi.org/10.1162/99608f92.5a8a3a3d>
- ²²⁰ Royal Society. Explainable AI. [Royalsociety.org](https://royalsocietypublishing.org/journal/rsos). 2024. ([Website accessed 11 August 2024](#))
- ²²¹ Goddard K, Roudsari A, Wyatt JC. Automation bias: a systematic review of frequency, effect mediators, and mitigators. *Journal of the American Medical Informatics Association*. 2012;Jan1;19(1):121-127. ([Journal full text](#)) <https://doi.org/10.1136/amiainl-2011-000089>
- ²²² Hauben M. Artificial intelligence in pharmacovigilance: Do we need explainability?. *Pharmacoepidemiology and Drug Safety*. 2022;Dec;31(12):1311-1316. ([Journal full text](#)) <https://doi.org/10.1002/pds.5501>
- ²²³ Ribeiro MT, Singh S, Guestrin C. "Why Should I Trust You?": Explaining the Predictions of Any Classifier [Internet]. [arXiv:1602.04938](https://arxiv.org/abs/1602.04938) 2016. ([Journal full text](#)) <https://doi.org/10.48550/arXiv.1602.04938>

1655 Chapter 7: Data privacy

1656 *Principle*

1657 Data privacy refers to the fundamental right of an individual to control how their personal
1658 information is collected, stored, shared, and used. It is an aspect of the principle of “respect for
1659 persons” that is foundational to the conduct of biomedical research. Regulations, legislation and
1660 guidance documents provide measures intended to preserve the confidentiality, anonymity,
1661 autonomy and control of sensitive and potentially personally identifiable health data in the setting of
1662 PV.

1663 *Key messages*

- 1664 • The ethical framework to evaluate the use of data privacy protected applications of AI in PV
1665 is embedded within the standard principles for research activities involving human subjects.
- 1666 • The use of certain AI applications in PV requires additional attention to assure data privacy.
- 1667 • The applications of ethical principles most relevant for the use of AI in routine PV are data
1668 privacy, fairness, and equity.
- 1669 • PV professionals should recognize that existing procedures used to assure regulatory
1670 compliance may need to be re-evaluated due to the heightened risks of GenAI to
1671 compromise data privacy and for ML to amplify biases.

1672 Introduction

1673 Although data privacy has been recognized as an implicit legal right for well over a century,²²⁴ it was
1674 not until the 1970’s that this topic began to receive formal international attention. Advances in
1675 computer technology began to facilitate the large-scale collection, organization, and evaluation of
1676 amounts of data that had previously relied upon paperwork. In the absence of any laws regulating
1677 how public bodies could collect, store, or share personal data, the first data privacy law was passed
1678 in 1980.²²⁵ In the US, public concerns about the potential misuse of collected data led to the US
1679 Privacy Act (1974),²²⁶ which provided boundaries for the collection, integrity, and use of personal
1680 data. These same issues raised concerns about transfer of large amounts of personal data across
1681 borders, which led to the first international guidelines to protect data privacy in the context of
1682 international trade.²²⁷ Similar to this CIOMS Working Group report, the OECD guidelines laid out a set
1683 of core principles; however, its intent was to assist governments, business and consumer
1684 representatives with the objective of supporting data transfer to facilitate commerce while
1685 protecting personal data privacy. Each of these documents were framed as legal responses to
1686 technological threats, and did not explicitly highlight the ethical foundations of data privacy. Over
1687 subsequent decades, the guidelines have influenced most subsequent data protection
1688 regulations/laws, such as the Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule
1689 1996 and General Data Protection Regulation (GDPR) 2016, both of which are discussed later in this
1690 chapter. As noted in Appendix 2, considerations to protect data privacy are specifically identified in a
1691 survey of recent major national and international reports on the use of AI, generally and in
1692 pharmaceutical development.

1693 Ethical considerations

1694 While data privacy concerns are widely recognized in the use of AI, at the time of this publication, the
1695 authors are unaware of a standard reference focusing specifically upon the ethical considerations in
1696 the use of AI applied to PV. Many publications that refer to ethics and AI, such as the WHO Guidance,
1697 Ethics and Governance of Artificial Intelligence for Health,²²⁸ emphasize several basic principles that
1698 were first elaborated in the Belmont Report (1979).²²⁹ That report was federally commissioned (USA)
1699 to determine basic ethical principles that are foundational to biomedical and behavioural research

1700 involving human participants. These principles underlying modern human research protection bear
1701 upon clinical research as well as certain post-marketing PV activities globally, recognizing that some
1702 PV activities are not technically considered research.

1703 The Belmont Report identified three basic principles that are foundational to interventional and
1704 behavioural research involving human participants: respect for persons; beneficence; and justice.

1705 Respect for Persons refers to the obligation for the research subject to enter research knowingly and
1706 of their own free will. It also considers that some individuals, such as those who are cognitively
1707 impaired, may not be capable of making this decision independently and based on self-
1708 determination. From a practical perspective, respect for persons is immediately recognized by the
1709 use of informed consent documents, which require that research subjects be advised that they are
1710 being invited to participate in a research study, that they are made aware of the purpose of the
1711 study, including its potential risks and benefits, and that the information be transmitted so that it is
1712 comprehensible to the participant. Participants must be able to freely choose whether to participate,
1713 including the option to exit the research project. Participants cannot be unduly influenced or coerced
1714 into participation. Informed consent is waived for routine PV activities through statutes recognizing
1715 the pre-eminence of societal/public health interest.

1716 The principle of Beneficence describes the need for research to be designed to maximize its potential
1717 benefits while minimizing potential harms. Research in which the harms of participation are known
1718 to outweigh potential benefits would not be justifiable. This concept also captures the core tenet of
1719 the Hippocratic Oath to “do no harm”. For those involved in clinical trials, equipoise is a familiar
1720 concept, based on the premise that there is genuine uncertainty about which treatment arm in a
1721 clinical trial is best with respect to safety and/or efficacy. It captures the concept the research is not
1722 frivolous, including that potential risks have been minimized. PV incorporates beneficence into
1723 ongoing benefit/risk assessments. It should be acknowledged that there are applications of AI in PV
1724 that do not contain personal information, and these activities are outside the scope of data privacy
1725 discussed in this chapter.

1726 Finally, Justice refers to the ethical obligation to see that research is conducted among those who
1727 might benefit and that involvement in research is sensitive to ensuring that certain populations, e.g.
1728 prisoners, are not preferentially selected for research out of convenience, and that the burdens and
1729 benefits of research should be born equally among those who might benefit.

1730 The intent of the Belmont Report was to address considerations in medical research, intended to
1731 advance generalizable knowledge, and draw distinctions with the distinct responsibilities of medical
1732 practice, intended to support individual patient care. The report did not discuss public health
1733 activities, but its key principles are used in public health activities ranging from disease surveillance
1734 to PV. Certain public health activities are mandated by statute, including PV. Drug safety includes
1735 monitoring safety during clinical trials, post-marketing surveillance, and post-approval safety studies
1736 (PASS). Clinical trials are defined as research and are a focus of the Belmont Report. PASS is a form of
1737 real-world evidence, often required as a condition of product licensure (RWE is the focus of CIOMS
1738 Working Group XIII).

1739 The ethical implications for the use of AI in PV for approved products that are most pertinent include
1740 data privacy (derived from the principle of Respect for Persons) and fairness and equity (derived
1741 primarily from that of Justice). Respect for persons indicates that every individual has the right to
1742 control information about themselves. In the context of PV, Justice captures the concept of fairness,
1743 meaning that the benefits of PV knowledge should be equitably distributed among those who may
1744 use specific medicinal products, i.e. lack of discrimination. Fairness as applied to PV means that the
1745 activity is conducted in a non-discriminatory manner, ideally so that the methods support the
1746 representativeness of the population being evaluated and equitable to provide insight into the
1747 population that may be exposed to the product. In the context of PV, the principle of Justice is

1748 applied through efforts to assure Fairness and Equity, which is the subject of a separate chapter in
1749 this report.

1750 **Data privacy regulations**

1751 Many countries have established laws to protect the data privacy rights of the individual. These laws
1752 share the common principle that personal data requires protection, and that this should be
1753 accomplished through mechanisms that mitigate risk to the individual while requiring accountability
1754 of the entity using the data. Two of the most frequently cited are the HIPAA, 1996, used in the United
1755 States, and the GDPR, 2016, which is employed in the European Union. These examples will be used
1756 to illustrate commonalities and differences between data privacy regulations, and their implications
1757 for the application of AI to PV.

1758 *Example: Health Insurance Portability and Accountability Act*

1759 As the name suggests, HIPAA (1996) originally focused on health insurance data²³⁰ and was
1760 developed to ensure data privacy as medical information moved from analog to digital. At the same
1761 time, the legislation was intended to advance common administrative standards for health care data.
1762 In short order, the rapid adoption of digital technologies in health care (e.g. electronic health
1763 records) and the interest in using electronic data for research and other purposes led to follow on
1764 legislation to support the use of electronic health records according to standards that would ensure
1765 administrative efficiency while protecting patient privacy and security (HIPAA Privacy Rule, 2000;
1766 Security Rule, 2003). The HIPAA Privacy rule defined individually identifiable health information
1767 (Protected Health Information – “PHI”) and defined safeguards to protect the privacy of this
1768 information. HIPAA also delineates Business Associate Agreements, requiring covered entities that
1769 handle PHI to comply with its privacy and security rules.

1770 HIPAA emphasizes the confidentiality, integrity and availability of health data, and includes
1771 provisions focused on the “minimum data necessary standard”, i.e. limiting data to those necessary
1772 for the specific purpose. It specifies patients' rights to access and amend medical records. To protect
1773 patient confidentiality, HIPAA recognizes types of data that could be used to identify individuals and
1774 specifies 18 unique PHI identifiers. The list underscores the range of common data types that are
1775 largely unrelated to health care and which contain identifiable information that could compromise
1776 patient identity: name(s), geographic subdivisions smaller than a state, dates (except year, e.g. date
1777 of birth), telephone numbers, fax numbers, email addresses, social security numbers, medical record
1778 numbers, health plan beneficiary number, account numbers, certificate/license numbers, vehicle
1779 identifiers, device identifiers, web URLs, internet protocol (IP) addresses, biometric identifiers (e.g.
1780 fingerprints); full face photographs, as well as any other unique identifier that could be used to trace
1781 the identify of an individual. Once these identifiers are stripped from a source record, the record can
1782 be used without restrictions imposed by HIPAA as the record no longer contains PHI.
1783 Public health often balances societal interest with personal rights. Based on overriding societal needs
1784 for the safety, effectiveness, and quality of medicinal products licensed for use in the US, routine PV
1785 activities conducted by license holders are typically exempt from some HIPAA requirements for
1786 patient authorization to disclose and use PHI. Medicinal products are governed in the US by Food and
1787 Drug Administration (FDA) regulations that require monitoring the quality and safety of FDA-
1788 regulated products, which is conducted in part through adverse event reporting, product tracking,
1789 recalls, and post-marketing surveillance. While some PV activities are exempt from HIPAA, data
1790 privacy protections remain, including: use of the minimum necessary data standard (collecting only
1791 data essential to fulfil the PV responsibility), de-identification and/or anonymization of data
1792 (employed where possible); use of technical, administrative, and physical safeguards to prevent
1793 unauthorized access, use, and disclosure); and safeguarding by Business Associate Agreements
1794 (where vendors or partners are engaged). Within the US, the FDA and Centers for Disease Control
1795 and Prevention (CDC, Atlanta) have complementary responsibilities to support and advance health.
1796 Among its responsibilities, the FDA is responsible for “protecting the public health by assuring the

1797 safety, efficacy, and security of human (and veterinary drugs), biological products”, and medical
1798 devices. The CDC responsibilities include protecting “America from health, safety and security
1799 threats, both foreign and in the U.S.”, including those associated with domestic and international
1800 diseases and chronic or acute diseases. This is accomplished in part by conducting “critical science
1801 and providing health information that protects (the US) against expensive and dangerous health
1802 threats”. To fulfil their responsibilities, public health authorities have a somewhat broader remit than
1803 the private sector and are able to conduct their responsibilities and use PHI without patient
1804 authorization; however, they implement data privacy safeguards both within their organizations and
1805 when collaborating with others.

1806 In contrast to public health authorities and the private sector, academic involvement in PV is
1807 generally conducted as a research activity (e.g. PASS), and are subject to different oversight,
1808 including use of institutional review boards (IRBs, aka ethical review boards) to assure that studies
1809 meet appropriate ethical standards, and are conducted with mitigation for data use and privacy, data
1810 use agreements, where applicable with other organizations the use of de-identified and limited data
1811 sets, and compliance with both HIPAA and the Common Rule.

1812 *Example: General Data Protection Regulation*

1813 The GDPR (Regulation [EU] 2016/679) is a comprehensive set of rules overseeing personal data
1814 protection in the European Union and succeeds the earlier Data Protection Directive (Directive
1815 95/46/EC), which was issued contemporaneously with HIPAA, at the dawn of the internet age. The
1816 scope of the GDPR is much broader than HIPAA as it pertains to the use of personal data affecting all
1817 manner of human interaction, including processing by automated means as well, and stems from the
1818 1950 European Convention on Human Rights: “Everyone has the right to respect for his private and
1819 family life, his home and his correspondence”.²³¹

1820 The GDPR incorporates principles such as lawfulness, fairness, transparency, accuracy and integrity,
1821 purpose limitation, data minimization, confidentiality, and storage limitation. Compliance is a major
1822 feature of the GDPR with organizations such as pharmaceutical companies required to have a Data
1823 Protection Officer responsible for overseeing compliance. Penalties for non-compliance are
1824 significantly greater than those under HIPAA, with fines up to 4% of global turnover. To offer
1825 practical examples, based on the requirements and obligations under the GDPR, four risk-based
1826 categories may be differentiated:

- 1827 1. Prohibited use (e.g. genetic data that might be used for the purpose of drug
1828 development would require explicit consent from the patient/participant); it is therefore
1829 prohibited without such consent;
- 1830 2. Restricted use (e.g. the use of biometric data to identify employees for security purposes
1831 only, in line with necessity, proportionality and other requirements;
- 1832 3. Permitted use with safeguards (e.g. use of purchasing history to create targeted
1833 advertisements would be typically be permissible - provided that efforts were made to
1834 protect customer identity, customers would need to consent and have the opportunity
1835 to opt-out, and efforts were made to identify and mitigate and risks to customer
1836 privacy); and
- 1837 4. General permitted use (e.g. processing customer information for online billing purposes,
1838 provided that the data are limited to those needed to fulfil the transaction, and that
1839 security measures are taking to prevent unauthorized access).

1840 Additionally, several safeguard measures may be used, such as data encryption (preventing access
1841 without a decryption key), pseudonymisation (replacing identifiable information with pseudonyms to
1842 mask identity), and use of Data Protection Impact Assessments to identify and mitigate risks in data

1843 processing to protect the individual. In contrast to HIPAA, GDPR incorporates a “right to be
1844 forgotten”, permitting individuals to request deletion of their personal data. In the case of special
1845 categories of personal data, such as health data, explicit consent may be required for data processing
1846 under GDPR and, where collected, such consent may be revocable.

1847 Similar to HIPAA, the rules of the GDPR allow for pharmaceutical companies to meet their legal
1848 obligations to conduct PV activities, monitor and report adverse events without consent in order to
1849 ensure oversight of the safety and effectiveness of medicinal products – provided that certain
1850 safeguards are in place. These responsibilities may limit data protection rights normally in place
1851 under the GDPR, e.g. the “right to be forgotten”. Other safeguards include requirements for data
1852 minimization as well as administrative, technical and organizational measures to protect personal
1853 data.

1854 In fulfilling its responsibilities to assure the safety, effectiveness, and quality of medicinal products
1855 authorized for use in the European Union, the European Medicines Agency is empowered to assure
1856 PV oversight in a manner that acknowledges that certain data protection rights, such as the right to
1857 be forgotten, may be limited for specific PV activities. The EMA emphasizes the principles of data
1858 minimization, purpose limitation, lawfulness, fairness and transparency in its data use. In contrast to
1859 HIPAA, the GDPR has special provisions for international data transfers, imposing restrictions in
1860 exporting data collected for EU citizens (regardless of domicile) outside the European Economic Area
1861 and applies safeguards to provide an appropriate level of data protection.

1862 Data privacy expectations for PV research conducted by academia in the EU are analogous to those
1863 for the US, with IRB oversight and an emphasis upon adherence to principles of data minimization
1864 and purpose limitation. Additionally, international collaborations that involve data transfers outside
1865 of the EEA require safeguards that typically include contractual language to assure compliance with
1866 GDPR rules.

1867 Other data privacy regulations

1868 Although the FDA and EMA data privacy regulations are currently the most widely followed, it should
1869 be noted that there are an increasing number of country-specific differences, which pose particular
1870 challenges for the use of multinational AI model development using secondary data. Comparison of
1871 regulations in place in Germany, China, and Japan illustrate this point.

1872 **Table 5: Data privacy regulations for using secondary data in Germany**

1873 Source: CIOMS Working Group XIV
1874

Aspect	Germany
Governing Law	General Data Protection Regulation (GDPR) (EU-wide), Federal Data Protection Act (BDSG) (Germany)
Health Data Classification	“Special category data” (Art. 9 GDPR)
Consent Requirements (Research & PV)	Usually required; exceptions for public interest (e.g. PV, RWE)
Secondary Use of Data (e.g. RWE)	Allowed if legal basis exists (public health, scientific research, etc.) with safeguards
De-identification Standards	Pseudonymization encouraged; full anonymization for broader reuse
Cross-border Data Transfer	Allowed to countries with adequacy or with SCCs/BCRs
Pharmacovigilance Exemptions	Explicitly exempt from consent under public health/legal obligation

Regulator Guidance on Biomedical Use	Extensive EMA and national ethics bodies guidance
Oversight Body	German DPAs + European Data Protection Board (EDPB)
Key References	GDPR (Regulation EU 2016/679); EDPB Guidelines 03/2020; EMA Module VI (GVP); BDSG (Germany)

Table 6: Data privacy regulations for using secondary data in China

Source: CIOMS Working Group XIV

Aspect	China
Governing Law	Personal Information Protection Law (PIPL)
Health Data Classification	“Sensitive personal information”
Consent Requirements (Research & PV)	Explicit consent generally required; strict interpretation
Secondary Use of Data (e.g. RWE)	Permitted with new consent or proper anonymization
De-identification Standards	Anonymization required to avoid consent; “irreversible” standard
Cross-border Data Transfer	Strict rules: security assessments, contracts, individual consent; limited adequacy
Pharmacovigilance Exemptions	AE reporting permitted but must minimize identifiable data
Regulator Guidance on Biomedical Use	PIPL + draft health data governance rules; evolving
Oversight Body	Cyberspace Administration of China (CAC)
Key References	PIPL (2021); CAC draft regulations on health data; State Council health data measures (2022)

Table 7: Data privacy regulations for using secondary data in Japan

Source: CIOMS Working Group XIV

Aspect	Japan
Governing Law	Act on the Protection of Personal Information (APPI)
Health Data Classification	“Special care-required personal information”
Consent Requirements (Research & PV)	Consent generally required, but pseudonymized data may be used for public interest or research
Secondary Use of Data (e.g. RWE)	Allowed with pseudonymization/anonymization and research purpose declaration
De-identification Standards	Recognizes both anonymized and pseudonymized data; latter still regulated
Cross-border Data Transfer	Permitted to “adequate” countries (EU, UK); otherwise, consent or contracts needed

Pharmacovigilance Exemptions	AE reporting allowed without consent under regulatory mandate
Regulator Guidance on Biomedical Use	MHLW guidance on clinical research and PV under APPI
Oversight Body	Personal Information Protection Commission (PPC)
Key References	APPI (2020 amendment); PPC Guidelines; MHLW guidance on GPSP and human research ethics

1875 **Practical considerations to support data privacy**

1876 As these examples indicate, regulations have been developed to assure appropriate data privacy
 1877 within the framework required to conduct routine PV activities. The list of 18 unique identifiers
 1878 enumerated by HIPAA highlights the breadth of the types of data that can be used to identify
 1879 individuals. In the years following the introduction of HIPAA and the GDPR, there has been
 1880 recognition that additional measures may be required to anonymize data.

1881 As a regulated industry, pharmaceutical companies must comply with the data privacy and reporting
 1882 requirements of all countries in which their products are licensed. As an example, the EMA requires
 1883 adherence to Good Pharmacovigilance Practices (GVP) and to data protection principles from the
 1884 GDPR. Ensuring compliance requires attention to evolving country-specific regulations as well as the
 1885 oversight of vendors that support companies (in some cases conducting certain PV activities for
 1886 individual companies) as well as business partnerships, e.g. where a combination therapy is co-
 1887 developed by more than one company. Regulatory authorities may have different requirements for
 1888 reporting patient information, necessitating some customization and additional oversight to assure
 1889 adherence to local requirements. For example, Australia requires reporting of ethnicity (to support
 1890 fairness/equity), while it is prohibited in France out of concerns of discrimination.

1891 To support compliance with global data privacy requirements, contractual arrangements with third
 1892 parties (e.g. vendors, partners) include privacy-specific provisions and language. In the US,
 1893 contractual arrangements with vendors/partners require Business Associate Agreements. Under
 1894 GDPR, binding corporate rules (BCRs) may be implemented to enable multinational companies to
 1895 move personal data within their companies across borders; BCRs are legally binding and require
 1896 approval from EU authorities. In the EU, an additional layer of oversight is imposed through the
 1897 required use of in-house data privacy officers for certain businesses such as pharmaceutical
 1898 companies. Globally, there are a range of potential consequences for data breaches, from
 1899 requirements for notification to data protection authorities up to and including significant fines and
 1900 penalties

1901 In the EU, an additional layer of oversight is imposed through the required use of in-house data
 1902 protection officers for certain businesses such as pharmaceutical companies. In many countries,
 1903 including in the European Economic Area, there are penalties for data breaches.

1904 **Risks to maintaining data privacy as artificial intelligence is employed in pharmacovigilance**

1905 One of the promises of AI is that it will permit more efficient processing of large amounts of routine
 1906 PV data, e.g. individual case reports. Additionally, large language models (LLMs), whether open or
 1907 closed, permit nearly instantaneous planned (or unplanned) linking of data sources that would
 1908 otherwise not have occurred, or have been difficult to accomplish. As discussed below, these risks
 1909 are substantively greater for open vs closed LLMs. GenAI models are also useful for extrapolation –
 1910 finding patterns that might otherwise not have been recognized. These attributes raise the question
 1911 of whether current data privacy tools are sufficient to prevent re-identification of deidentified data.

1912 *Adequacy of de-identification measures*

1913 In 1990 (six years prior to HIPAA) a US researcher used census data to identify 87% of the US
 1914 population based on three readily available data elements: five-digit mailing (zip) code, gender and
 1915 date of birth, illustrating that few data points were needed to uniquely identify individuals.²³² Though
 1916 mailing codes were subsequently classified as PHI under HIPAA, the point remains that just a few
 1917 generally accessible data points may be needed to compromise data privacy.

1918 A study using data from a children’s hospital in Ontario, Canada, demonstrated that the risk of re-
 1919 identification of individuals based upon de-identified pharmacy data could be minimized, or even
 1920 eliminated, by reducing the precision of values in selected data elements, such as replacing the
 1921 admission and discharge dates with the quarter and year of admission. However, the maximum
 1922 amount of acceptable generalization in the data element values must be determined by formally
 1923 examining not only the risk of re-identification and breach of patient privacy but also the intended
 1924 analysis, which may not be conducted without the appropriate level of precision.²³³

1925 The EMA and Health Canada now require public sharing of clinical trial reports as part of the drug
 1926 approval process. Standards for data anonymization have been issued. Applying these standards,
 1927 researchers evaluated the risk of re-identification associated with a clinical study report for a
 1928 nonsteroidal anti-inflammatory drug, grading suspected cases based on the likelihood of accurate
 1929 matching.²³⁴ The authors found six suspected matches out of 500 reviewed cases and observed that
 1930 identifying the matches was time-consuming (24.2 hours per case). Matching was best informed by
 1931 social media and death records. Based on the 0.09 probability risk threshold of re-identification
 1932 established by EMA²³⁵ and accepted by Health Canada, the authors concluded that existing
 1933 anonymization guidance was sufficient to provide an adequate level of data protection and advised
 1934 review the mechanisms by which re-identification had occurred. With rapid advances in AI, the time
 1935 required to replicate the reidentification exercise reported in that study (published in 2020) will likely
 1936 have decreased by the date of this report, and will continue to do so.

1937 *Risks of data breaches*

1938 The potential consequences of re-identification of de-identified data are amplified by numerous
 1939 examples of data breaches that have occurred throughout the world. A few examples illustrate the
 1940 breadth of some recent data breaches, along with their potential consequences:

- 1941 • A purposeful attack on a US financial firm leading to access of more than 100 million
 1942 customer accounts and credit card applications;²³⁶
- 1943 • An apparently politically motivated international attack by a foreign government on a credit
 1944 reporting agency in the US resulting in the release of names, birth dates, and social security
 1945 numbers of nearly half the US population, purportedly with the intent of using AI to
 1946 compromise US government officials;²³⁷
- 1947 • A purposeful domestic data breach intended to embarrass a political opponent, involving a
 1948 cyber-attack on an Asian health care plan result in 1.5 million patient records;²³⁸
- 1949 • An unintentional release of Indian government biometric, and other personal data, in a
 1950 database containing records of 1.2 billion individuals.²³⁹ In this instance, a criminal group
 1951 exploited the data breach and offered individual patient records for sale. Approximately
 1952 100,000 persons are known to have had their data accessed.

1953 Each example occurred before the widespread use of generative AI, a technology that has been
 1954 advancing rapidly, and which has the potential to efficiently link publicly available data sources with
 1955 those obtained maliciously, leading to enhanced risk for re-identification and compromising data
 1956 privacy.

1957 *Individual responsibility to protect personal data*

1958 In addition to processes to ensure data privacy to conform to data privacy regulations, individuals
 1959 play a role in protecting their own data. This responsibility grows in importance with the ever-

1960 increasing number of digital tools (e.g. Smartphones, wearables), apps, and software (e.g. AI-assisted
1961 translation tools, GenAI) that provide opportunities for individuals to disclose personal data that may
1962 not be sufficiently protected. In many instances, there are legal requirements to support an
1963 individual's data privacy (e.g. through the GDPR); however, there remain opportunities for lapses in
1964 data privacy, and these are particularly worrisome in the use of GenAI. Common mechanisms to
1965 advise persons of data use policies may include terms of use (e.g. End User Licensing Agreements –
1966 aka EULA), data privacy notices, and, in some instances, the requirement of explicit consent for use
1967 of personal data. Individuals may share personal data (e.g. names, phone numbers), when submitting
1968 queries without understanding the consequences of disclosing such information. In many instances,
1969 notably those using open GenAI tools, these data may no longer be private. Individuals may also
1970 share context-specific information, such as a recent illness (or as in a noted example, a motor vehicle
1971 accident) that might be used to identify them.²⁴⁰ Users may also unintentionally provide personal
1972 identifying information by sharing (e.g. in social media) context-specific data outputted by GenAI.
1973 These data may subsequently be leveraged to identify the individual even if personal data was not
1974 directly entered depending upon the data privacy policies of the respective platform.

1975 *Potential risks to data privacy using large language models in pharmacovigilance, and approaches for*
1976 *mitigation*

1977 In principle, existing data privacy regulations, (or legislation, such as the GDPR), should provide the
1978 basis for protection of data used in AI applications to PV. Model development may require the use of
1979 PV data (such as ICSRs) to data scientists and machine learning engineers for training as well as
1980 execution. All involved parties, which may include both pharmaceutical companies and vendors, may
1981 have access to data that is protected, creating the risk for exposure to larger groups. All parties
1982 should be aware of data privacy requirements. The risk is potentially greater with LLMs, as the
1983 underlying mechanism of these models provides the potential for some re-identification that would
1984 otherwise be unlikely. Those organizations that maintain closed LLMs exercise control of prompts as
1985 well as the data contained in the models; in contrast, open LLM models do not have this safeguard
1986 and run a greater risk of re-identification, due to the LLM capability of drawing from data sources
1987 that may be opaque and not intrinsically within the purview of data privacy regulations (for example,
1988 containing the sort of data described above under Individual responsibility to protect personal data).
1989 Leaks may occur through prompts bypass data privacy considerations or through models that are
1990 trained on personal data. Additionally, as noted above, re-identification can occur even with
1991 presumed de-identified data. PV requires review of potentially identifiable and sensitive information
1992 that includes basic demographics (including birth date) associated with sensitive data elements
1993 including medical or health information (including medicine and vaccine exposure), ethnicity, race,
1994 sexual orientation, genetic information, biometric data, physical characteristics, lifestyle information,
1995 etc., requiring heightened safeguarding measures.

1996 Among the types of challenges posed by GenAI (as well as in some cases ML) for PV are the following.

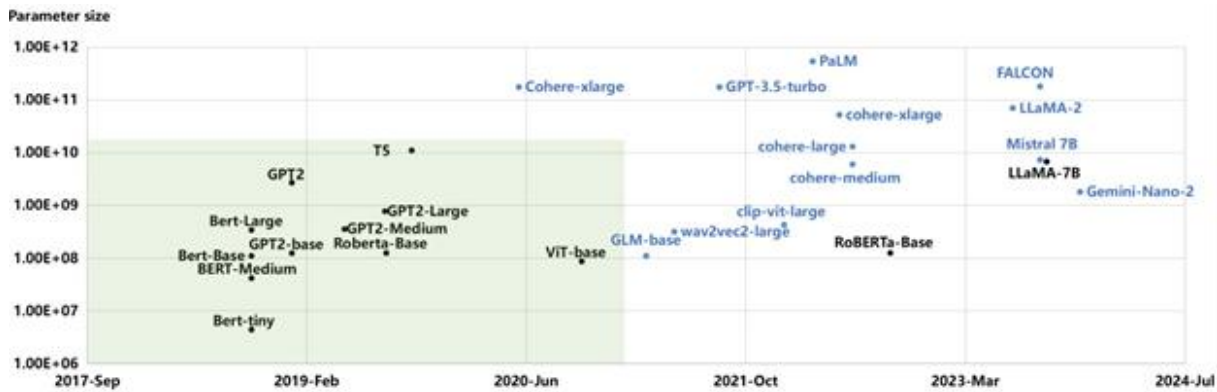
- 1997 • Algorithms may be developed within open LLMs, without attentiveness to applicable data
1998 privacy requirements, thereby posing a potential privacy risk. If these LLMs are then adopted
1999 for use within closed LLMs, there is the potential risk for disclosure of protected information.
- 2000 • Within a closed LLM, attention should be paid to different sources that may be added to the
2001 LLM for unrelated purposes. If genetic data has been collected (with participant consent) for
2002 a study and is added to a LLM for a specific analysis, measures would need to be taken to
2003 ensure that it is not used for a different purpose outside of the original consent. The
2004 accepted practice is to seek consent for additional uses of those data (as the data would now
2005 be part of the LLM).
- 2006 • Generative AI programs can integrate otherwise discrete data sources such as such as census
2007 and vital statistics and public health data, which may be linked to a deidentified health
2008 record data set (e.g. in the setting of an active PV activity (e.g. a post-authorization safety
2009 study) leading to the possibility of reidentification.

2010 These types of risks are amplified in settings where data privacy regulations are lax or poorly
 2011 enforced. LLMs that are smaller and introduced earlier have tended to have more scrutiny for data
 2012 privacy than larger and more recent LLMs (see Figure 6). Reasons may include: 1) lack of public
 2013 availability of newer, larger LLMs; and 2) privacy technologies have struggled to keep up with these
 2014 newer, larger LLMs.

2015 **Figure 6: Relationship of timing of large language model introduction, parameter size and**
 2016 **attentiveness to data privacy**

2017 Source: ²⁴¹

2018



2019

2020 The horizontal axis represents the time of LLMs release, while the vertical axis represents the size of
 2021 model parameters. Blue dots signify LLM instances not addressed in the literature pertaining to
 2022 privacy protection, whereas black dots indicate those that have been examined in such literature.
 2023 The green backdrop delineates the central cluster zone of LLMs with the potential to facilitate
 2024 privacy protection.

2025 **Conclusions**

2026 The right to data privacy resides within the well-established framework of basic ethical principles for
 2027 human research protection articulated in the Belmont Report. National regulatory authorities have
 2028 provided requirements intended to protect data privacy, indicating the types of data that can be
 2029 made available, along with safeguards (such as data minimization, anonymization, de-identification
 2030 and data encryption) along with potential penalties for non-compliance.

2031 Despite data privacy laws and the increasing sophistication of technical measures employed by
 2032 companies entrusted with personal information, attempted and successful data breaches have been
 2033 occurring with increasing frequency and often at enormous scale (affecting in some cases >100
 2034 million individuals), suggesting both failures in oversight along with technical advances to outwit
 2035 cybersecurity measures and break into secure data sources.

2036 Ever-increasing computational power, larger linked databases, and the introduction of GenAI, are
 2037 occurring in parallel with an increasingly globalized PV landscape involving more numerous and
 2038 complex interdependencies (e.g. business partnerships, international vendors conducting PV
 2039 activities). The ongoing challenge for PV professionals and regulatory agencies and industry, as well
 2040 as colleagues and academia will be to assure that within this rapidly evolving data science landscape,
 2041 data privacy measures are monitored and regularly updated to properly protect personal data.

2042 A potential risk in applying GenAI for PV is patient reidentification, suggesting a need to reconsider
 2043 the specificity of de-identified data, along with risks associated with open LLMs in which some data
 2044 sources may be outside the control of the user.

2045 In addition to following existing data privacy regulations, efforts to mitigate risks to data privacy
 2046 when applying AI to PV may include the following.

- 2047
- 2048
- 2049
- 2050
- 2051
- 2052
- 2053
- 2054
- 2055
- 2056
- 2057
- 2058
- 2059
- 2060
- Recognition that the technology is advancing rapidly, requiring ongoing monitoring, e.g. to assure that data de-identification measures are adequate.
 - Understanding that open and closed LLMs pose somewhat different challenges to data privacy. Operating closed LLMs in safeguarded environments within institutional firewalls and carefully examining the risks of sharing these models with third parties should be helpful in risk mitigation.
 - Audits to evaluate whether only the minimum required personal information is included in reports, that any re-use of data for secondary purposes is consistent with the purposes for which that data was collected and that adequate measures are in place to support compliance with data protection requirements by all entities (e.g. vendors) contributing to PV. Insofar as PV activities may be conducted by a network of collaborating organizations, the organization with the weakest oversight of data privacy may present a risk.
 - Oversight regarding access to LLMs for PV practices to assure that use by trained PV professionals is fit for purpose.

References

- ²²⁴ Warren and Brandeis. "The Right to Privacy". *Harvard Law Review*. Vol. IV. December 15, 1890, No. 5. ([Webpage accessed 27 April 2025](#))
- ²²⁵ Gesetz zum Schutz personenbezogener Daten (Hessisches Datenschutzgesetz – HDSG), Gesetz- und Verordnungsblatt für das Land Hessen I, 1970, S. 61.
- ²²⁶ Privacy Act of 1974, Pub. L. No. 93-579, § 552a, 88 Stat. 1896 (1974) (codified as amended at 5 U.S.C. § 552a).
- ²²⁷ OECD, *Guidelines on the Protection of Privacy and Transborder Flows of Personal Data*, OECD Doc. C(80)58/FINAL (Sept. 23, 1980.)
- ²²⁸ World Health Organization. *Ethics and governance of artificial intelligence for health: WHO guidance*. Geneva: World Health Organization; 2021. (Full text accessed 25 March 2025)
- ²²⁹ U.S. Department of Health & Human Services. *The Belmont Report: Ethical Principles and Guidelines for the Protection of Human Subjects of Research*. Washington, D.C.: U.S. Department of Health & Human Services. ([Full text accessed 21 March 2025](#))
- ²³⁰ U.S. Department of Health & Human Services. *Health Information Privacy. HIPAA for Professionals*. ([Webpage accessed 21 March 2025](#))
- ²³¹ GDPR.eu. *What is GDPR, the EU's new data protection law?* ([Webpage accessed 21 March 2025](#))
- ²³² Sweeney L. Simple demographics often identify people uniquely. *Carnegie Mellon University, Data Privacy Working Paper*. Health (San Francisco). 2000;Jan1;671:1-34. ([Full text accessed 21 March 2025](#))
- ²³³ El Emam K, Dankar FK, Vaillancourt R, Roffey T, Lysyk M. Evaluating the risk of re-identification of patients from hospital prescription records. *The Canadian Journal of Hospital Pharmacy*. 2009;Jul;62(4):307. ([Journal full text](#)) <https://doi.org/10.4212/cjhp.v62i4.812>
- ²³⁴ Branson J, Good N, Chen JW, Monge W, Probst C, El Emam K. Evaluating the re-identification risk of a clinical study report anonymized under EMA Policy 0070 and Health Canada Regulations. *Trials*. 2020;Dec;21:1-9. ([Journal full text](#)) <https://doi.org/10.1186/s13063-020-4120-y>
- ²³⁵ European Medicines Agency, Guidotti T. *External guidance on the implementation of the European Medicines Agency policy on the publication of clinical data for medicinal products for human use*. 2018. ([Webpage accessed 21 March 2025](#))
- ²³⁶ Starks T. U.S. charges Chinese military hackers with massive Equifax breach. *CNN*. ([Article accessed 21 March 2025](#))
- ²³⁷ Geller E. U.S. charges Chinese military hackers with massive Equifax breach. *New York Times*. 2020. ([Article accessed 21 March 2025](#))
- ²³⁸ *The Straits Times*. Personal info of 1.5m SingHealth patients, including PM Lee, stolen in Singapore's most serious cyber attack. *The Straits Times*; 2020. ([Article accessed 21 March 2025](#))
- ²³⁹ Jain M. *The Aadhaar Card: Cybersecurity Issues with India's Biometric Experiment*. Seattle, WA: Henry M. Jackson School of International Studies, University of Washington; 2019. ([Article accessed 21 March 2025](#))
- ²⁴⁰ Janmey V, Elkin PL. Re-Identification Risk in HIPAA De-Identified Datasets: The MVA Attack. *AMIA Annu Symp Proc*. 2018;Dec5;2018:1329-1337. ([Journal full text](#))
- ²⁴¹ Yan B, Li K, Xu M, Dong Y, Zhang Y, Ren Z, Cheng X. On protecting the data privacy of Large Language Models (LLMs) and LLM agents: A literature review. *High-Confidence Computing*. 2025;Feb28:100300. ([Journal full text](#)) <https://doi.org/10.1016/j.hcc.2025.100300>

2061 Chapter 8: Fairness & Equity

2062 *Principle*

2063 Fairness and equity require awareness of and adherence to impartiality, equality, non-discrimination,
2064 diversity, justice, and lawfulness. The benefits of AI in PV should be equitable across all relevant
2065 populations and groups. Throughout the AI lifecycle, it is important to avoid and mitigate unfair bias,
2066 and any discriminatory practices and unjust social wellbeing and environmental impacts.

2067 *Key Messages*

- 2068 • Consider the development and application of AI impacting fairness and equity, whose lack or
2069 imbalance may result in discriminatory harm to subpopulations underserved by an AI
2070 solution, explicit biases resulting in negative impact, or impact performance by providing
2071 inaccurate results.
- 2072 • Plan and implement mitigation strategies where possible for areas that bias may be
2073 introduced reducing potential underperformance; avoid discriminatory harm to underserved
2074 populations.
- 2075 • Equity may be advanced by taking measures (e.g. assess for representative data sets) to
2076 assure that AI applications to PV result in outputs (e.g. assessments, aggregated data outputs
2077 used for product safety assessments, etc.) that are relevant to populations anticipated to
2078 have exposure to the specific medicinal product being evaluated.
- 2079 • Screening and identifying explicit or potential bias when possible is key to implementing
2080 mitigation measures to reduce risk, determining AI applicability and limitations, and
2081 establishing expected performance acceptance criteria.
- 2082 • Scrutinize training and performance evaluation reference data sets for adequate
2083 representation and evaluate performance in relevant subgroups when possible. Inadequate
2084 reference data is often the cause of inadequate fairness and equity.
- 2085 • There is limited fairness and equity in ICSRs, with some countries reporting significantly more
2086 than others and providing more contextual data available for analysis, such as real-world
2087 data (RWD). Consequently, our understanding of routine usage is often limited among
2088 underserved populations.

2089 **Introduction: general concepts and considerations**

2090 In the context of PV, adherence to established laws and regulations such as privacy laws and PV
2091 regulations must remain intact with the introduction of AI. What has changed is the increasing
2092 awareness of the need for consideration, governance, and mitigation of potential factors that may
2093 influence or impact fairness and equity with the use of AI.

2094 Not all fairness and equity concepts, considerations or negative consequences associated with the
2095 use of AI will be uniquely specific to PV. That does not negate the need to address these
2096 considerations. Biases within AI solutions is a general problem which may impact performance, and
2097 not all forms of statistical bias will result in a negative impact on fairness and equity. Within PV, the
2098 focus will be regarding unfair bias introduced through data collection, selection, model development
2099 and human involvement in the design, development and use of AI that could potentially result in
2100 unfairness, discrimination, or inequality.

2101 This chapter will not address the impact of development and use of AI on justice and lawfulness, on
2102 individuals' access to essential services, lack of public resources for financing and implementing AI
2103 tools and the required ecosystem, and impact on social well-being, because while these are
2104 important issues, they are not unique to PV. While not unique to PV, access to AI and the required
2105 ecosystem can be a significant barrier for low- and middle-income countries that can result in
2106 inequality and underserved populations. Potential workforce implications with introduction of AI in

2107 PV will not be addressed here as it is discussed in the Chapter on Human Oversight under the Section
2108 on [Transformation of traditional roles](#). In addition, the rapid acceleration in the use of GenAI is
2109 associated with significantly increased energy demand and environmental consequences, both of
2110 which are acknowledged as having a broad societal impact, but which are beyond the scope of this
2111 report.²⁴²

2112 Fairness and equity considerations are challenging and can be influenced by cultural differences,
2113 historical inequalities, perceptions, and risk factors may appear subjective. To reduce bias,
2114 intentional actions are required throughout the AI lifecycle, from design through implementation,
2115 and while in production, to reduce discriminatory risk. There are numerous factors that may
2116 introduce bias.

2117 **Fairness and equity considerations and pharmacovigilance**

2118 Fairness and equity principles are fundamental in identifying and addressing discriminatory biases
2119 arising from the use of AI systems. In PV, proactive measures are essential to detect, assess,
2120 understand and prevent adverse effects to ensure safe and effective use of medicines. When
2121 utilizing AI systems in PV, it is essential to implement proactive strategies to mitigate potential
2122 harm caused by high-impact AI systems.

2123 Thorough evaluation and ongoing monitoring are required throughout the AI system lifecycle to
2124 identify and quantify potential areas of risk and mechanisms through which bias may be introduced
2125 as a first step to define strategies to mitigate discrimination biases arising from the use of AI systems
2126 in PV. Monitoring for biases is required from conception, development, testing, and following
2127 solution deployment. The frequency of monitoring for bias and appropriate modification needs to be
2128 defined based on risk assessment, solution results, and potential external factors that may impact
2129 model bias and performance.

2130 It is crucial to acknowledge the possibility of bias that may lead to unfair practices or unequal
2131 treatment of patients when using AI in activities related to detecting, collecting, assessing,
2132 monitoring, understanding, and preventing adverse effects or any issues related to medicinal
2133 products.

2134 The PV professional is responsible for ensuring that the AI solution meets the defined business
2135 requirements, supports patient safety activities, and does not introduce bias that may inadvertently
2136 place patients at risk, a disadvantaged position or at potential for discrimination, e.g. denied the
2137 potential benefits of a medicinal product, through exclusion based on race, gender, age, or socio-
2138 economic factors. However, we acknowledge that the current system is not perfect; it is essential to
2139 ensure that AI does not exacerbate existing issues, but instead, contributes to improving fairness and
2140 equity.

2141 **Sources of potential threat to fairness and equity**

2142 Inherently, humans are biased and can introduce that bias throughout the AI system lifecycle (e.g.
2143 requirements gathering, model training, monitoring may not detect poor performance, incorrect
2144 results, or missed scenarios). AI experts and developers can have unconscious bias, and potentially if
2145 not identified and addressed, the output can have limitations, be discriminatory and may not be
2146 recognized as biased. Conversely, the output could be accurate, fair, and equitable; however, results
2147 may be rejected by the human with a bias as being poor performance.

2148 **Inadequate training and/or testing data set(s)**

2149 Bias is primarily introduced in the data used to develop and test AI solutions, which can perpetuate
2150 bias and discrimination resulting in harm. Incorrect conclusions can occur when there are data

2151 limitations such as when it is not a complete dataset or does not represent the population where the
2152 AI is being applied. The inappropriate or unintended application of AI to populations not represented
2153 can occur if data limitations are not transparent or recognized.

2154 **Explicit negative bias**

2155 Some of the main risk to fairness and equity is introduced from explicit biases (e.g. a case
2156 prioritization algorithm down-prioritising reports for men because of patterns represented in the
2157 training set).

2158 The lack of robustness and availability of data, e.g. health records not digitally available globally, can
2159 lead to underserved populations or underperforming models. When data representation is
2160 inadequate, the available data does not correspond to the population and consideration is required
2161 to remediate under-represented groups or lack of available organized data, e.g. regions with less
2162 developed PV reporting systems, otherwise scenarios will be biased toward groups represented by
2163 the training data, and since the data does not represent all groups, e.g. ethnicities, results in bias
2164 against these groups, e.g. minorities.

2165 Inadequate data - whether as a result of data not being available or organized in a usable format or
2166 structure, lack of data robustness, or inadequate representation of all variables - may result in an
2167 under-performing model, or worse, incorrect conclusions as a result of model limitations not being
2168 recognized, and this may negatively impact patients' health outcomes. Imbalance of data
2169 representation can potentially skew data, amplify imbalances, and it may be difficult to identify and
2170 assess bias when reviewing an AI solution's output.

2171 Historically, there have been examples of bias influencing PV activities because of data limitations
2172 such as known under-reporting or stimulated reporting of adverse events, with inadequate data or
2173 imbalance of data providing artifact that have had negative impacts. Litigation such as class action
2174 lawsuits that are pursuing product liability claims can result in stimulating high volume of reported
2175 adverse events during the process of legal firms identifying potential plaintiffs that could overshadow
2176 unsolicited reports and the imbalance of data could be a threat to fairness and equity considerations
2177 if the data imbalance results in incorrect conclusions with groups that are under-represented as a
2178 result of skewed data. Local prescribing practices' impact on adverse event reporting and data
2179 availability should be considered. These data biases if introduced into an AI solution will potentially
2180 magnify the negative impact and remain undetected with difficult identification of underlying bias.

2181 **Underserved groups**

2182 Under-representation can directly result in underserved population segment(s) and potentially not
2183 recognize nuances of subpopulations. Population-specific segmentation can be done by
2184 demographics, disease processes, genetic variability, health practices variability, and cultural
2185 differences for medical regimens and patient expectations. Such differences can introduce bias with
2186 a negative impact if data is exclusive to a specific group, if data is exclusionary, or if nuances of a
2187 subgroup are not understood, such as a case prioritization algorithm that underperforms in reports
2188 from certain countries in Asia, because they were under-represented in a training data set and
2189 differed in important ways from other countries represented.

2190 During clinical trials, such potential harm may be overlooked if subgroups are under-represented in
2191 the study population or receive a lesser level of care, e.g. have limited access to medical
2192 professionals or facilities. There may be more focus on preventing false negatives to not miss
2193 significant information, e.g. the failure of the PV process to detect potential harm restricted to or
2194 over-represented in certain subgroups. In the post-marketing period, deployment of an AI solution
2195 working less well in certain patient subpopulations could lead to an inability to detect adverse events
2196 from these populations. Conversely, false positives may be of greater concern in duplicate detection

2197 where a higher rate of reports falsely flagged as suspected duplicates in a specific country could lead
2198 to missed or delayed safety signals there.

2199 Special populations frequently not represented, such as age related (paediatric, geriatric), pregnant
2200 women, and infrequent or under-reported events such as rare diseases, and events with social
2201 stigmas need to be considered when assessing bias. In the example of an AI solution implemented to
2202 support signal detection activities, with limited data from special populations (e.g. pregnancy), the
2203 negative impact would be magnified with misinterpreted or missed signals.

2204 Reliance of decision making on data not representative of respective populations (e.g. post-approval
2205 risk minimization activities based on data with limited representation of served population) could
2206 result in minimization measures not properly addressing safety of patients in the population. If
2207 unable to mitigate lack of representation in AI solution, it may require reliance on historical PV
2208 approaches and safety measures (e.g. robust monitoring measures for special populations).

2209 Detailed identification of “groups” that could be disfavoured or low volume events proportionate to
2210 data set and comprehensive strategies to address data inadequacies can reduce potential bias,
2211 discrimination, and underserved populations.

2212 **Artificial intelligence solution design**

2213 Algorithms should not perpetuate existing bias or discrimination, and the algorithmic design can lead
2214 to unintended consequences. When AI was used to develop a model to predict what patients would
2215 benefit from proactive intervention in the care of their chronic illness, its results directed more
2216 resources to white patients than black patients, because the data set used for training was based on
2217 utilization, not need.²⁴³ Given a healthcare system and a universe of healthcare data that is likely to
2218 carry country-specific biases, any naïve use of AI will reproduce these biases in its predictions. The
2219 likelihood of adverse consequences is more likely because of the apparent opacity of AI, hype about
2220 its capabilities, limited understanding of how it works, and unclear pathways to question its
2221 conclusions.

2222 Parameters defined by a programmer and how a model processes data could introduce bias or
2223 produce inaccurate results by skewing the result. If a developer selects or designs features for an AI
2224 solution based on their own conscious or unconscious bias, the resulting output could be suboptimal
2225 or even incorrect. In the case of GenAI prompt engineering development, the potential to introduce
2226 bias based on the prompt design, lack of specificity, context, or omission of a required prompt could
2227 result in an output with a negative bias. Individual preferences influence decisions and subsequently
2228 influence data selection and model development. This could occur due to the model developer
2229 having an affinity to subgroups like their own profile (e.g. developer is a young person and may select
2230 data that does not account for paediatric or geriatric populations).

2231 The development strategy should have a conscious systematic approach to avoid bias and achieve
2232 complete and accurate data representation accounting for diverse groups. Documenting how distinct
2233 groups are represented in the training and test data may provide insight to limitations, bias, and
2234 potential impact supporting implementing mitigation measures. When considering the population of
2235 respective groups, confirmation that the data are representative of the global population is needed
2236 to ensure balance, demographic parity, and appropriate distribution and allocation.

2237 AI is increasingly employed in the field of medicine to identify patterns and anomalies, such as
2238 consistencies, inconsistencies, and outliers in the identification of safety issues and communications.
2239 For example, examining sentiment consistency can help flag and mitigate human-induced
2240 discrepancies. This proactive approach reduces the risk of unfairness and bias, enhancing the
2241 reliability and objectivity.²⁴⁴

2242 Risk, impact, and mitigation measures

2243 The consequences of AI on fairness and equity are dependent upon the application of AI within PV,
2244 the usability, performance, and the risk of where the AI is being used within the process. When there
2245 is discrimination and bias embedded in the AI model through data limitations and/or algorithm
2246 development, the negative impact of the resulting biased model is magnified in its application. The
2247 model may amplify or skew outcomes resulting in incorrect conclusions, incorrect introduction of an
2248 advantage or disadvantage, inequalities, or discrimination of groups or populations.

2249 Evaluating an AI solution pre and post deployment for explicit or potential bias allows for mitigation
2250 measures to reduce risk. AI solution explainability may highlight explicit bias and understanding the
2251 profile of training data provides a degree of insight into potential areas where bias may be
2252 introduced into the solution, determine appropriate use, solution limitations, degree of human
2253 oversight required, and expected performance. When evaluating for bias, consideration should be
2254 given to post deployment data annotation processes for future retraining activities and mitigate
2255 when possible.

2256 Within PV signalling activities, omitted results could cause misrepresentation of a product
2257 benefit/risk profile and have a detrimental impact, leading to incorrect human conclusions or
2258 decisions impacting patient safety.

2259 Sensitivity analysis of performance across different subgroups can be important to highlight groups
2260 or populations underserved by an AI solution. A risk-based approach when selecting subgroups to
2261 evaluate performance may be necessary when an exhaustive sensitivity analysis is not feasible and
2262 may be dependent upon data limitations for training and test data for subgroups or populations.

2263 Key mitigation strategies

- 2264 • Evaluate each AI solution for fairness and equity, outlining the assessment method, results,
2265 and any measures taken to mitigate.
- 2266 • Review training and test data sets thoroughly for completeness and adequate group
2267 representation.
- 2268 • Perform sensitivity analysis when possible, evaluating AI model results for equality by
2269 changing subgroups/populations to confirm expected results (e.g. modify gender input and
2270 evaluate impact to the output) and highlight potentially underserved populations. This is
2271 especially important when an AI solution has a lower level of explainability.
- 2272 • Review AI solution design, parameters, and feature selection for bias when an AI solution is
2273 explainable, and the results are not as expected.
- 2274 • Ensure training data description is transparent highlighting explicit bias and allow clarity on
2275 model limitations to reduce inappropriate application or incorrect conclusions.
- 2276 • Determine level of human involvement required in development and monitoring activities
2277 providing required input to ensure accurate performance and fair results.

2278 Identification of potential risk areas is challenging but key to preventing bias, discrimination, and
2279 suboptimal model performance. Avoidance of data limitations is not always possible and providing
2280 visibility of data characteristics allows appropriate application and opportunity to mitigate risk. It is
2281 important to understand the model limitations and communicate to the user community and group
2282 monitoring AI performance of limitations and potential bias.

2283

References

²⁴² Government of Canada. *Guide on the use of generative artificial intelligence*. Ottawa: Treasury Board of Canada Secretariat; 2023. ([Webpage](#) accessed 21 March 2025)

²⁴³ Ziad Obermeyer et al. *Dissecting racial bias in an algorithm used to manage the health of populations*. *Science* 366,447-453(2019). DOI:10.1126/science.aax2342

²⁴⁴ Bergman E, Sherwood K, Forslund M, Arlett P, Westman G (2022) *A natural language processing approach towards harmonisation of European medicinal product information*. *PLoS ONE* 17(10): e0275386. <https://doi.org/10.1371/journal.pone.0275386>

2284

2285 **Chapter 9: Governance & Accountability**

2286 *Governance - Principle*

2287 Governance refers to the human management and oversight used to control and direct the use of AI
2288 in the PV system. An AI governance framework requires implementation of risk management
2289 practices and policies to ensure adherence to the AI guiding principles.

2290 *Accountability - Principle*

2291 Accountability applies to clearly defined roles, responsibilities and liability for organisations and/or
2292 individuals deploying, operating and managing AI systems. It requires the adoption of appropriate
2293 governance measures by relevant stakeholders, including but not limited to regulators, vendors,
2294 users, developers, data providers or pharmaceutical companies involved in setting policy, developing,
2295 deploying and managing AI systems. This ensures operations remain within expected parameters
2296 throughout the AI lifecycle while addressing any unforeseen consequences.

2297 *Key messages*

- 2298 • Governance requires a comprehensive approach across all lifecycle stages of an AI solution as
2299 well as the processes it impacts and should therefore be established as early as possible.
- 2300 • Accountability requires clearly defined roles and responsibilities for stakeholders involved in
2301 AI systems for PV; AI systems themselves cannot be held accountable.
- 2302 • Systems and processes, along with service providers and software vendors, need to be
2303 qualified.
- 2304 • Regular reviews of AI systems and how they adhere to the AI principles are necessary to
2305 ensure ongoing regulatory compliance and performance.
- 2306 • A governance framework grid for an AI system in PV can serve as a structured guide to help
2307 relevant parties to document e.g. assessments, actions and references processes, such as
2308 Standard Operating Procedures (SOPs) etc. throughout the lifecycle of the AI system.
- 2309 • Governance and accountability should be independent of the business' utilization and value
2310 proposition of the AI system to facilitate unbiased decision making.

2311 **Introduction**

2312 Previous chapters have discussed in detail the importance of taking a risk-based approach, providing
2313 adequate human oversight, demonstrating validity and robustness, and addressing transparency,
2314 data privacy, fairness and equity when integrating and implementing AI solutions into the overall PV
2315 system. This chapter outlines the guiding principles of governance and accountability in AI-enhanced
2316 PV. We will discuss the importance of these two principles, the stages of the AI lifecycle that require
2317 specific governance actions, the roles and responsibilities of various stakeholders, regulatory
2318 oversight, and the need for ongoing training in the rapidly evolving field of AI technology.

2319 Robust governance and clear accountability are crucial for the success of AI initiatives. These
2320 principles help ensure that AI systems are used responsibly and ethically, are compliant with
2321 regulations, while fostering trust and transparency among stakeholders. Clearly defined roles and
2322 responsibilities enable all stakeholders to understand their obligations and effectively oversee AI
2323 systems.

2324 As AI technology evolves, governance and accountability frameworks will need to be adapted. New
2325 risks and challenges will emerge, requiring updated principles and practices. Continuous review and
2326 adaptation are essential for staying ahead of these changes. This includes the refinement of the
2327 proposed governance framework grid for practical use.

2328 Governance framework

2329 A governance framework grid (referred to as “grid”) for AI solutions in PV (see Table 2) is a structured
2330 guide designed to identify key considerations to address each of the principles throughout the
2331 lifecycle of the AI system, including concept, development, deployment, and monitoring phases of
2332 the AI model developed for PV use.

2333 In addition to serving as a structured guide for planning and overseeing AI solutions, the grid can also
2334 aid in self-assessment. By detailing where each action or process is recommended, the grid helps
2335 ensure that the principles such as transparency, accountability, and a risk-based approach are
2336 consistently adhered to, facilitating the integration of AI into PV systems. Regular reviews of KPIs by a
2337 governance body, aimed at ensuring adherence to the AI guiding principles, can facilitate
2338 identification of gaps and drive improvements in the AI solution. While a governance body with
2339 expertise and focus on AI’s use in PV is needed in early phases, integration of the governance process
2340 into the overall PV system oversight mechanisms should be considered when the AI solutions enter
2341 routine use phase. If a risk emerges that warrants significant modification to the AI solution, the AI
2342 focused governance body may need to be re-engaged.

2343 Consultation with the grid can occur in multiple ways. A unit within a PV organization may have an
2344 idea for AI-based automation and specify governance requirements upfront when commissioning a
2345 vendor or internal development team. Alternatively, a vendor might present a ready-made AI
2346 solution to a PV organization, which then can be evaluated against the AI guiding principles for
2347 example by applying this grid. Early consideration of governance principles is crucial for the
2348 successful implementation of an AI system. These principles should guide the development or
2349 selection of a vendor solution, deployment, and ongoing management. Early planning should be
2350 focused on identifying potential risks and determining mitigation strategies. Furthermore, it can
2351 stimulate focus on alignment with ethical and regulatory standards of the AI system from the outset,
2352 setting the foundation for a robust and compliant AI system.

2353 The grid is composed of five lifecycle phases of the AI solution: an initial requirement specification
2354 phase where business units typically provide input, followed by development, pre-deployment, post-
2355 deployment, and routine use. These phases are valid for both initial qualification and iterative
2356 changes of the AI solution. In each phase, the AI guiding principles should be considered, and in the
2357 grid, each principle constitutes a cell for relevant documentation hereof. When the grid is used for a
2358 specific AI solution, each cell is intended to provide information about actions, considerations, or
2359 references to where these actions are documented, such as SOPs, working instructions, or
2360 repositories containing log files, reviewed performance metrics, or names of accountable
2361 persons/review bodies. Illustrations of how each guiding principle is applied throughout the lifecycle
2362 phases can be found below, and examples of how to put this grid into practice can be found in
2363 [Appendix 3: Use cases](#).

2364 **Table 8: Governance framework grid**

2365 Source: CIOMS Working Group XIV
2366

	<i>Collection of specifications, requirements</i>	<i>Development & change management</i>	<i>Pre-deployment & post-change sign-off</i>	<i>Post-deployment & post-change hyper-care</i>	<i>Routine use</i>	<i>General considerations</i>
<i>Risk-based approach</i>	Risk assessment (theoretical) - AI model - Context of use - Impact & likelihood of risks	Risk mitigation plan	Risk assessment (empirical) Adjustments to risk mitigation plan based on performance evaluation	Intensive or targeted monitoring for risk assessment (empirical), target high risk areas Mitigation if needed	Routine monitoring (e.g. risk for model drift) Mitigation if needed	Review and refine risk-based approach at regular intervals

Human oversight	Multidisciplinary expertise HOTL (human-on-the-loop; design)	Define HITL (human-in-the-loop; strategy) Change management, including staff training plan	Fine-tune HITL strategy Staff training roll-out	Implement HITL strategy Intensive or targeted intervention	Routine HITL activities Adjust and fine-tune(?) HITL strategy as needed	HIC (human-in-command; holistic oversight) HOTL (general monitoring)
Validity & robustness	Specification of use case & deployment domain Specification of reference standard(s) Specification of benchmarks Requirements on reproducibility	Training & validation Development or acquisition of reference standards	Performance evaluation Benchmark comparisons	Performance monitoring	Continuous integration and deployment Periodic performance monitoring	Special considerations for low-prevalence settings Reproducibility Assessing AI solutions with human-in-the-loop
Transparency	[From organization to developer] Model requirements - Intended use - Human-computer interaction - Explainability - Expected outputs Performance evaluation requirements - Scope - Reference standard	[From developer to organization] Model - Architecture - Parameters - Acceptable inputs - Expected outputs - Standard AI components - Training & validation - Known limitations Explainability in support of model development, debugging, and documentation	[From developer to organization] Performance evaluation - Scope - Sampling - Reference standard - Human input - Summary metrics - Benchmarks - Subsets & sensitivity analyses - Qualitative review Explainability in support of assessing Validity & robustness and Fairness & equity	[From developer to organization and from organization to end user] Performance evaluation - Deviations Explainability in support of assessing Validity & robustness and Fairness & equity	[From organization to end user (and regulatory authorities)] Disclosing use of AI Explainability in support of building trust with end users	
Data privacy	Specification of use -Specification of data sources -Identification of data elements that contain identifiers -Jurisdictions / provenance of data Data privacy by design (data minimization)	Training data set selection, algorithm design	Test set (if publishing is intended, e.g. as a public benchmark, need to assure data privacy consistent with local legislation/regulations/guidance)	Adherence to data privacy considerations in running the models with full / accruing data sets Greater attention to deviations in hypercare	Ongoing processes to identify and rectify data privacy issues in routine use	Data privacy legislation/regulations vary by country/region, potentially leading to inconsistencies (emerging risks of new prompts)

	-Data protection impact assessment					
Fairness & equity	Context of use -acceptable application	Training data set selection, algorithm design and cognitive bias Avoid -Explicit or potential unfair bias -inadequate data inclusion	Pre-deployment performance evaluation -Reference data sets inadequate (e.g. unavailable, inadequate representation) -Algorithm design - Human/cognitive bias		Routine Monitoring -poor performance related to model shift, inadequate training data (underrepresented populations, special populations)	Ensure model and training data description is transparent, and limitations highlighted to reduce inappropriate application or incorrect conclusions
Governance & accountability	Consideration on how AI solution fits into existing PV system Key roles (non-exhaustive examples) - PV experts - AI experts	Agreement on KPIs to support implementation Key roles (non-exhaustive examples) - PV experts - Data scientists - AI experts - IT specialists - Ethics specialists	Refinement of KPIs to support implementation, based on performance evaluation Key roles (non-exhaustive examples) - PV experts - Data scientists - AI experts - Senior management	Approval of risk mitigation strategies Key roles (non-exhaustive examples) - PV experts - AI experts - IT specialists - Data protection officers - Cybersecurity experts	Integration into overall PV quality management system including oversight of KPIs Key roles (non-exhaustive examples) - PV experts - AI experts	

2367 Description of each lifecycle phase included in the grid is presented below:

2368 **Collection of specifications, requirements:** This is the initial phase where the stakeholders are
 2369 identified and engaged, and the project’s objectives, scope and features are defined. The
 2370 multidisciplinary team of PV professionals, data scientists, AI/ML engineers, software engineers, IT
 2371 specialists, and other domain experts (also refer to the chapter on [Human Oversight](#)), is typically
 2372 managed by system developers, software vendors, or an internal IT development team. This phase
 2373 provides a roadmap for developers and end-users, and lays the foundation for the entire
 2374 development process. Like traditional software, as an AI solution evolves, the requirement
 2375 specifications may also require iterations, and consequently, the grid may need to be reconsidered
 2376 accordingly.

2377 **Development & Change Management:** In this phase, the multidisciplinary team focuses on acquiring,
 2378 creating or modifying AI systems, ensuring they are built with the necessary functionality and
 2379 adherence to governance principles. Whether developing an AI system or selecting a vendor solution,
 2380 these principles will apply throughout.

2381 **Pre-Deployment & Post Change “Sign-Off”:** At this phase, the AI system transitions from the
 2382 development stage to deployment into the PV process. Before implementation, a thorough
 2383 validation and approval process is required to ensure the AI system, or any changes hereof, is ready
 2384 for deployment. Typically, a PV expert becomes accountable for the results produced by the AI
 2385 solution and for adapting the processes in which the solution will be used. Documentation of this
 2386 phase may include risk assessments, review of sufficient adherence to principles, sign-off forms,
 2387 validation reports, and many references to SOPs detailing the sign-off procedure, etc.

2388 **Post-Deployment & Post Change "Hypercare"**: Following deployment, this phase is critical for the
2389 immediate monitoring of the AI model's performance or the latest changes' impact. It is a period of
2390 intensive observation to promptly identify and resolve any unanticipated issues, as real-life
2391 application of the AI system in the PV process might surface issues due to various reasons such as
2392 incorrect assumptions, design flaws, unintended bias, in earlier stages. This phase, dominated by
2393 hypercare, differs from traditional software hypercare; for AI solutions, immediate fixes may not be
2394 feasible and other measures such as human intervention or increase in human oversight might be
2395 needed. Documentation is expected and may include incident logs and performance analysis reports
2396 specific to the most recent change while under observation.

2397 **Routine**: This phase signifies the full integration of the AI solutions into the PV process. It involves
2398 ongoing monitoring, maintenance, and documentation to ensure full oversight and allows for the
2399 identification of trends through the monitoring of pre-defined KPIs. This phase may reference routine
2400 reports, logs of ongoing actions, and which SOP or working instruction manages this review process,
2401 reflecting the model's full operational status.

2402 Of note, discoveries during post-deployment or routine use phases may necessitate the AI solution
2403 being sent back to pre-deployment for enhancements.

2404 The following, non-exhaustive examples illustrate aspects to consider for each guiding principle in
2405 relation to the lifecycle phases in the grid:

2406 Transparency: In the Development phase, there is a focus on creating comprehensive documentation
2407 of the development activities including reason for changes and data used in model training. In Pre-
2408 Deployment, transparency is further enhanced by adding model performance evaluation, and
2409 empirical evidence for fairness and equity. Also, the documentation created should ensure consistent
2410 understanding of the intended use among different stakeholders. In routine use, the most important
2411 transparency is toward the end-users and those responsible for the continual performance
2412 evaluation and monitoring.

2413 Accountability: Throughout all phases, there is a consistent need to assign and document
2414 responsibility, whether it is to IT, vendors, or to PV experts. This ensures clarity about who is
2415 accountable for the AI model's development, change management, deployment, and performance at
2416 any time.

2417 Risk-based approach and human oversight: This begins with identifying the level of risks associated
2418 with development of the AI solution. When relevant, it may involve the development of clear
2419 annotation guidelines for human domain experts to ensure solid method development and
2420 performance evaluation. The next step is to propose appropriate mitigation strategies such as
2421 defining "human-in-the-loop" within an AI solution and other oversight measures up to eventually
2422 creating risk mitigation requirements in the user interface. It continues with redefining human
2423 oversight in Pre-Deployment, and further refining these concepts in the Routine phase based on real-
2424 life observations. This sequential approach highlights the need for evolving risk management as the
2425 AI model advances through its lifecycle. A risk-based approach in general is recommended for all
2426 measures taken to adhere to AI guiding principles.

2427 Any changes to the AI system must undergo the same rigorous governance considerations as the
2428 initial deployment. This ensures that modifications do not compromise the system's integrity or
2429 performance. Documentation and validation are essential to maintain transparency and
2430 accountability. Change management processes should be in place to handle updates and
2431 modifications effectively and account for a post-deployment phase, that based on "hypercare", will
2432 confirm performance and quality beyond routine monitoring.

2433 As computer system validation requirements need to be met at the same time, it is advisable to de-
2434 couple AI model version control from the rest of the software versioning.

2435 Governance body and accountability assignments

2436 To effectively manage the review and agree on actions and risk assessments towards the different
2437 principles, it is advisable to nominate a governance body. This group ideally should be a diverse,
2438 cross-functional team that has sufficient awareness of the end-to-end process and the extent of
2439 automation within it. It should include representatives from all relevant stakeholders and
2440 representation from the software vendor may also be considered. This diversity ensures a broad and
2441 balanced review of the AI solution. The governance body oversees the development, deployment,
2442 and ongoing management of the AI system to ensure that all actions align with guiding principles and
2443 regulatory standards. The governance body also determines accountable persons for the respective
2444 lifecycle phases, which includes sign-off of the documentation prior to deployment of the AI system
2445 into the PV process. Because business cases are often drivers of AI initiatives, the governance body
2446 should also include the respective project managers or sponsors to ensure adequate resourcing of
2447 governance measures during each phase of the lifecycle.

2448 Unlike traditional software, the governance body of an AI system should review the adherence to the
2449 AI use guiding principles in defined intervals, and ad-hoc if needed, to ensure the assessments are
2450 still valid. This is due to the rapid evolution of the field and the inherent risks of AI models that
2451 changing inputs, rules or other unforeseen issues may disrupt the solution at varying degrees, some
2452 significantly. The appropriate frequency and scope of reassessment of a deployed AI solution should
2453 be assessed. There should be measures ready to intervene or even disable the AI solution if
2454 necessary. Once the AI solution reaches the routine use phase, governance can be handed off to
2455 process owner to be integrated in the overall PV system monitoring process. Nevertheless, if a risk
2456 emerges that warrants significant modification to the AI solution, the AI focused governance body
2457 may need to be re-engaged. The introduction of version control for the governance framework grid
2458 should also be considered.

2459 Traceability and version control

2460 Traceability and version control are crucial aspects of managing AI solutions, particularly in a
2461 regulated field like PV where errors could impact patient safety or public health. They can enable
2462 evaluation and reproducibility of earlier versions of an AI solution and are often required for audit
2463 purposes (however, see also the discussion of AI solutions with stochastic components in the Chapter
2464 on [Validity & Robustness](#)). General best practices from existing version control frameworks can offer
2465 orientation for the version control of AI models, which should be documented alongside other
2466 relevant systems involved in the end-to-end process. They should include clear change control
2467 processes within both a user acceptance testing environment and the production environment.

2468 Documentation of an AI model should comprise its entire lifecycle, and may cover the justification,
2469 initial scoping and conception, development, deployment, validation, and post-deployment. It should
2470 allow for the retrieval and reproducibility of essential steps and decisions, including justifications and
2471 reasoning for deviating from pre-specified plans. As in traditional computer system validation,
2472 experiments conducted in, or before, the development environment are not required to be
2473 documented step by step. However, when the outcome of such an experiment or analysis impacts
2474 how an AI solution is evaluated or deployed, the justification for such decisions should be
2475 documented. If a decision is based on certain results or insights from the development stage, this
2476 should be documented.

2477 During the development phase, AI models undergo continual experimentation and iterative
2478 improvement. Transparency between the development team and the PV organization is crucial to
2479 ensure efficiency and that the solution is fit-for-purpose. Developers may create multiple versions of
2480 a model, test various features, and experiment with different training sets. In this context, focus
2481 should be on maintaining clear records of significant milestones – such as major changes in model

2482 architecture, the introduction of new datasets, or significant shifts in performance metrics. This
2483 allows developers to track the evolution of the model and understand the implications of key
2484 changes without being overwhelmed by the sheer volume of minor tweaks and experiments.

2485 Once an AI model moves from development to routine use in a production environment, the need
2486 for rigorous traceability and version control increases substantially. Deployed versions of the model
2487 should be documented in detail. In addition to the source code for each version, its underlying model
2488 architecture, training and test sets, and performance evaluation results should also be documented.
2489 From a regulatory perspective, the appropriate place to declare this would be in a document such as
2490 the PV System Master File (PSMF), in the EU.

2491 The continual improvement and adaptation of AI models post-deployment should also be
2492 documented. It may be triggered by human domain experts or built into the deployment of the AI
2493 solution itself including pre-specified monitoring of deterioration of performance or model drift.
2494 Some challenges related to this for Software as Medical Device have been described by FDA [65].

2495 When integrating external AI components (such as pre-trained models or libraries), it is important to
2496 document the versions of these components, particularly if they play a critical role in the model's
2497 performance. However, it may be sufficient to document these components at the time of significant
2498 milestones rather than during every iteration. As an example, for AI-based static systems, previous
2499 work proposes a specific documentation approach with proposed considerations for documentation
2500 within the different stages of the AI solution lifecycle.

2501 **Roles and responsibilities in artificial intelligence-enhanced pharmacovigilance systems**

2502 Organizations are accountable for the quality processes associated with their PV system, including
2503 the oversight of the AI components by the system owner. Oversight activities may be executed by a
2504 third party under appropriate supervision. Regulations, e.g. EU AI Act, may require organizations to
2505 establish specific roles, such as those to promote AI literacy, and facilitate fairness and equity.
2506 AI systems themselves cannot be held accountable. Human oversight is essential for ensuring the
2507 safe and responsible use of AI. Clear roles and responsibilities must be defined for all stakeholders
2508 involved in AI initiatives.

2509 The roles of PV experts are evolving with the introduction of AI. Already now, AI introduces new
2510 tasks, such as overseeing AI systems and interpreting their outputs. PV experts must adapt to these
2511 changes and develop new skills and competencies (see chapter on [Human Oversight](#)) to fulfil their
2512 obligations. This is especially relevant for members of the governance body and persons nominated
2513 as accountable for a lifecycle phase. The governance framework grid allows stakeholders to assess
2514 whether certain new activities will become relevant at specific steps, highlighting training needs
2515 early.

2516 Just like with traditional software providers, the collaboration between vendors of AI solutions and
2517 PV experts is crucial. This collaboration can facilitate that AI systems meet PV requirements and
2518 governance principles. Regular audits of vendors and AI systems are essential for maintaining
2519 compliance and ensuring development standards. Effective collaboration and audits foster
2520 transparency and accountability. This can ensure AI systems that are reliable, meet regulatory
2521 standards and are inspection ready.

2522 Regulatory authorities also play a role in monitoring AI in PV. They oversee that AI systems comply
2523 with regulatory standards and governance principles through inspections. Regulatory authorities are
2524 also developing guidance on the use of AI in the drug lifecycle, including PV (see Chapter on
2525 [Landscape analysis](#)). Integration of AI tools into the PV system must include appropriate regulatory
2526 documentation, such as in the Pharmacovigilance System Master File (PSMF) (see Chapter on
2527 [Transparency](#)).

2528 PV inspections are likely to increasingly focus on AI systems, with inspectors reviewing AI-related
2529 documentation, performance metrics, and governance practices. Inspectors will need adequate

2530 competencies to evaluate these systems effectively. This includes technical knowledge of AI and data
2531 science. As a result, continuous development and training are needed for inspectors to fulfil their
2532 role in new and fast-evolving areas.

2533 Balancing innovation with regulatory compliance and adherence to guiding principles is important for
2534 the success of AI initiatives. This involves fostering a culture of responsible innovation. These goals
2535 can be achieved by establishing effective governance processes that include regular reviews of AI
2536 solution KPIs.

2537

2538 **Chapter 10: Future considerations for development and** 2539 **deployment of artificial intelligence in pharmacovigilance**

2540 **The evolution and future of artificial intelligence in pharmacovigilance**

2541 The chapter explores the continuing transformative impact of AI on PV from the current application
2542 to a vision of how AI might impact PV in the future. The CIOMS Working Group XIV's discussion in the
2543 earlier chapters of this report is grounded in common principles. Use cases (presented in [Appendix 3](#))
2544 detail various AI systems under evaluation, various stages of deployment and assessment of their
2545 effectiveness within the discipline. To try to predict into the future, it is essential to recognize that
2546 the trajectory of AI is dynamic and highly unpredictable. Indeed, the only truly predictable elements
2547 are that AI will be ubiquitously deployed and is set to revolutionize many aspects, arguably all, of
2548 drug development and medical practice, from bench to bedside – as well as PV. For this reason, this
2549 chapter is grounded on the further developments of AI in PV described in earlier chapters of this
2550 report and anticipates how applications based on the principles might need to evolve as AI use in PV
2551 becomes more prevalent and sophisticated.

2552 The chapter provides considerations for PV stakeholders, including regulators and healthcare
2553 professionals and other industry stakeholders to ensure AI's safe and equitable deployment in PV.
2554 The skillsets needed by PV professionals today will likely differ from those required in the future
2555 necessitating involvement in the design, development, deployment, and routine use of AI in PV. The
2556 examples illustrate the direction and immense potential of AI adoption in PV; however, these
2557 examples are speculative to a certain extent and are not meant to be exhaustive. AI is set not only to
2558 potentially revolutionize PV, dissolving traditional boundaries of PV, but also expand its footprint far
2559 more broadly across medical sciences.

2560 The current decade represents a nascent phase for AI adoption in PV, and it is worthwhile
2561 acknowledging that the broad field of AI, particularly GenAI, is currently advancing rapidly. Further
2562 and more extensive deployment of AI may necessitate changes in how we think or approach PV
2563 strategies in the years ahead, driving the discipline of PV beyond its traditional frameworks and
2564 transforming it into self-detecting, real-time monitoring of safety data that aligns with the evolution
2565 of AI-driven medical science; for example, with the ability to rapidly analyse and extract vast
2566 quantities of safety data for case reporting and signal detection purposes. By leveraging this
2567 capability of AI, PV will evolve from a reactive focus on reporting and assessment to a forward-
2568 looking approach centred on proactive prediction and prevention and real / near real time learning
2569 systems.

2570 The initial phase of evolution has started to impact PV's core activities including case management
2571 and safety surveillance, as it continues to move from reporting and assessment towards prevention,
2572 enabled by advancements in AI-enhanced healthcare and into radically new areas of medicine. These
2573 technologies have the potential to reduce manual workload and the burden on PV professionals. For
2574 example, by accelerating response times for priority events, increased capabilities to sift through
2575 large and varied sources of safety data including literature, automatically creating case reports, and
2576 performing signal detection.^{245,246,247,248,249}

2577 **Transformative role of pharmacovigilance long-term and beyond: from** 2578 **detection to prevention**

2579 Advanced AI systems are poised to take PV beyond current boundaries and into automated or
2580 augmented decision making.²⁵⁰ AI, with capabilities for approximate reasoning, could handle
2581 ambiguity and partial truth values, for instance, in assessing safety data from social media entries or

2582 from fragmented safety-relevant data across different systems. Such AI systems may enable PV
2583 professionals to make nuanced decisions in case classification (e.g. assigning causality) or other PV
2584 situations requiring medical decision making. This would be particularly useful for cases with
2585 incomplete or conflicting data, where the “gray area” requires sophisticated, context-aware analysis
2586 and/or medical judgment.^{251,252}

2587 In the future, an expert AI system, designed specifically for PV, may emulate the judgement and
2588 decision-making processes of seasoned professionals or organizations with deep expertise in the
2589 field. These systems may not supplant human experts but augment their capabilities, enabling more
2590 nuanced and efficient decision making. An expert PV AI system would ideally be tailored to
2591 incorporate advanced analytical preciseness specific to therapeutic areas such as oncology,
2592 immunology, vaccines and medical devices as well as very different therapeutic options such as
2593 digital therapies, ensuring they adapt to the complexities of specific therapies, diseases, and patient
2594 populations, within those therapeutic areas, while performing or supporting PV work. While the
2595 development of such systems requires significant investment, their potential to drive the next
2596 generation of targeted PV solutions positions them as a critical innovation in advancing patient
2597 safety.

2598 By mid-century, it is possible that traditional PV will have transitioned from primarily detecting and
2599 processing adverse effects to a frontline technology-driven discipline that is engineering technologies
2600 that can detect, evaluate and share the information with “self” (human or organ: heart, kidney, liver,
2601 lungs etc.).^{253,254,255,256}

2602 This may then allow HCPs and patients to take a more active role in vigilance and prevention by
2603 taking corrective actions before adverse symptoms arise. As a discipline, PV leveraging AI is likely to
2604 evolve into a function that develops technologies enabled by AI to perform a proactive assessment of
2605 anomalies, self-report and self-learn on how to prevent the presence of such anomalies in the future
2606 and continue to promote patient wellness and safety. This will include “true” AI-enabled proactive
2607 self-regulated vigilance and risk mitigation.

2608 **Future development and deployment of AI and the guiding principles**

2609 The CIOMS Working Group XIV members made the careful decision to structure the report around
2610 common principles for the use of AI in PV, based in part upon the recognition that this
2611 transformative technology is in a period of exponential growth. A report that was prescriptive and
2612 overly reliant upon current examples would quickly become outdated, especially if AI technologies
2613 from other healthcare domains are leveraged. The authors expect that common principles for the
2614 use of AI in PV will be durable for the foreseeable future. What is less certain is how the guiding
2615 principles may be applied. Although the principles are robust and are expected to endure, it is likely
2616 that they will evolve in parallel with the technical AI advances and their use in detection, prevention
2617 and decision making by the individual human or subject going under medical treatment. The
2618 potential implications are discussed for each of the guiding principles below.

2619 2620 *Risk-based approach*

2621 Chapter 3 discusses risk-based approaches including risk mitigation, and also considers the regulatory
2622 framework required.

2623 The proliferation and advancement of AI may lead to continuous self-learning and potentially
2624 autonomous AI systems, with potentially great advancement in PV and benefit to patients and HCPs.

2625 Nevertheless, such systems come with potential concerns and risks. For example, a significant
2626 concern is the potential for AI to distort our understanding of a medicine's benefit-risk profile in real-
2627 world settings. Traditionally, these profiles are evaluated through carefully designed frameworks
2628 involving spontaneous reporting systems and planned surveillance studies. However, AI-driven tools

2629 may inadvertently restrict prescribing practices; for instance, by limiting access to AI-enhanced PV
2630 systems for high-risk patients or preventing off-label use.

2631 Further complicating matters, the adoption and availability of such tools may vary across healthcare
2632 systems and regions, introducing inconsistencies in data patterns that are challenging to interpret.
2633 This fragmented landscape can obscure the true influence that AI systems exert on prescribing
2634 decisions, making it difficult to assess their actual impact on patient outcomes. In addition, incorrect
2635 interpretation and poor utilization of AI is likely to significantly hamper patient safety. The principles
2636 of human factors and ergonomics (HFE) can assist in simplifying AI design and consecutively optimize
2637 human performance ensuring better understanding of AI outcome.²⁵⁷

2638 The oversight and risk mitigation of such advanced AI systems demand a dynamic risk assessment
2639 framework; one that integrates near-real-time monitoring and adaptive evaluation processes.
2640 Ensuring effective communication of these evolving risks to all stakeholders, including patients, will
2641 be crucial. As part of risk mitigation, healthcare leaders must embrace flexible governance models
2642 that account for AI's evolving nature, ensuring that transparency, accountability, and equitable
2643 access remain at the forefront.

2644

2645 *Human oversight*

2646 Chapter 4 covers human oversight including the changing and transformation of traditional roles in
2647 PV as AI use becomes increasingly embedded and ubiquitous.

2648 As AI systems become increasingly pervasive and autonomous, the role of human oversight will
2649 inevitably shift. While maintaining a "human-in-the-loop" approach will likely remain essential, this
2650 may prove insufficient for highly complex or higher-risk applications — including aspects of PV.
2651 Conversely, in some scenarios, human oversight may substantially change and become less relevant,
2652 as AI systems surpass human capabilities in reviewing data and regulating their own
2653 processes.^{258,259,260}

2654 This evolving landscape will require PV professionals to develop new skillsets and undergo
2655 specialized training to effectively oversee AI-driven systems. The focus must extend beyond
2656 traditional oversight methods to include competencies in understanding, interpreting, and guiding AI
2657 behaviours. By cultivating these skills, PV professionals can ensure that human oversight remains
2658 meaningful and effective in safeguarding patient safety and public health.

2659

2660 *Validity & Robustness*

2661 Chapter 5 discusses validity and oversight and considers multi-disciplinary collaborations required as
2662 well as reference standards and performance evaluation that might be needed to ensure robust and
2663 valid AI systems.

2664 As AI becomes more embedded and sophisticated, the challenge is to develop appropriate methods
2665 and systems that validate and ensure data integrity in tandem with the developments. For example,
2666 with the potential for generating vast amounts of data in real time or near real time, there is a need
2667 for more appropriate validation methods to avoid the risk of false signals. This may require PV
2668 individuals to develop new skill sets or even new specific scientific disciplines. AI use with some
2669 advanced technologies would need the creation of new standards and validation methods for the
2670 outputs and real-time / near-real-time safety data generated, e.g. neurotechnology such as
2671 implantable chips, smart organs, nanotechnology and smart organs.

2672

2673 *Transparency*

2674 Chapter 6 covers transparency and explainability of AI systems and related challenges.

2675 As AI becomes increasingly pervasive, our ability to track its deployment and understand its decision-
2676 making processes may diminish, posing significant challenges to explainability and transparency. AI
2677 systems may mirror complex statistical processes and advance programming or AI-coded programs.
2678 Consequently, the necessity, and even practicality, of full transparency may face new challenges.
2679 Expectations of transparency may need to evolve as trust in AI systems strengthens and meets
2680 predefined confidence thresholds.

2681 Much like AI's role in data analysis, statistics, and signal detection today, tracing AI's precise
2682 influence on downstream decisions may become increasingly difficult. Just as the complexities of
2683 prior distributions in Empirical Bayes Geometric Mean (EBGM) disproportionality models are widely
2684 accepted yet rarely scrutinized, established trust in AI-generated outputs may drive a shift in focus,
2685 with the expectation that errors or miscalculations will still prompt corrective actions to ensure
2686 sound decision making.

2687 In parallel, as trust in AI solidifies, the emphasis on explainability may similarly evolve. While
2688 transparency will remain important, its most critical value may emerge during incidents or errors.
2689 Much like the role of flight data recorders in aviation, explainability may become vital for
2690 understanding failures and enhancing system improvements rather than serving as a constant
2691 requirement.

2692 This shift may significantly influence PV decision making, emphasizing timely interventions and near-
2693 real-time root cause analysis. Looking ahead, organizations may need to balance the benefits of
2694 enhanced-AI performance against the degree of transparency required, carefully weighing improved
2695 efficiency with the need for interpretability in high-stakes decisions.

2696

2697 *Data privacy*

2698 The right to control one's personal data is durable and has been widely adopted internationally.
2699 What is likely to occur in the coming years is that preserving data privacy will become more
2700 challenging. As noted in Chapter 7, leaks of personal data have been increasing in frequency, with
2701 some at enormous scale.²⁶¹ The increasing use of online platforms for communications and services
2702 has been accompanied (in some countries) by a common lack of understanding into how collected
2703 data are used along with an acquiescence to the risk of data breaches. Breaches have occurred for
2704 reasons ranging from neglect to criminal intent. In the case of health care data, the release of
2705 personal data contrary to individual approval carries risks for emotional well-being, stigmatization,
2706 and discriminatory treatment.

2707 The pressures to amass and link large health care data sources are compelling, both on account of
2708 operational efficiencies (assuring consistencies in clinical care as well as medical care costs) and the
2709 advancement of scientific knowledge. At this time, the use of GenAI is in its infancy, and the only
2710 certainty is that it will both improve in quality and accelerate in use, as it is applied to many areas of
2711 biomedical research and clinical practice, and indeed in our daily lives. The use of open LLMs carries
2712 particular risks for the unintended disclosure of personal data, a topic that is likely to receive
2713 attention in coming years as the risk becomes clearer.

2714 Societies will need to balance the pressures for the commoditization of data to maximize learning
2715 and therefore better outcomes for patients with AI, with protections against unintended disclosure.
2716 One possibility is that data sharing will be automated, but that systems have built-in checks and an
2717 obligation to maximise the demonstrable value of the data for the patients and/or patients' carers.
2718 Security measures to support anonymization might incorporate blockchain or similar technology to
2719 make complete anonymization possible without a patient key to allow all care-relevant data to be
2720 safely shared with complete confidence and assurance that the Individual's data are anonymised.
2721 Without the appropriate regulatory checks and balances, it is also easy to see that these data could
2722 easily be misappropriated or abused.

2723 AI's evolution may usher in an era where access to underlying safety data becomes instantaneous,
2724 enhancing real-time insights and facilitating seamless data sharing. These advancements could
2725 significantly improve the timeliness and accuracy of safety assessments. However, an opposing
2726 scenario is equally plausible, one in which data sharing becomes increasingly restricted due to
2727 proprietary concerns, legal complexities, or public mistrust. As awareness grows regarding data's
2728 value as a commercial asset, particularly in insurance and other industries, heightened caution may
2729 further constrain data flow.

2730 Balancing these dynamics will be critical. Establishing transparent frameworks that foster trust,
2731 ensure data integrity, and promote responsible data sharing will be essential to fully realize AI's
2732 potential while safeguarding public confidence.

2733

2734 *Fairness & Equity*

2735 The fairness and equity Chapter 8 considers how and what type of discriminatory biases might be
2736 identified, addressed and/or prevented arising from the use of AI systems.

2737 Fairness and equity should mean that patients and health care professionals should have equal
2738 access to all the new and advanced AI technologies.

2739 It is important, as described in the data privacy section above, to ensure that PV with ubiquitous AI
2740 use is deployed equitably, and that data sharing does not put individuals at risk for e.g. higher costs
2741 associated with more advance monitoring, genetic profiling and/ or personalized risk / remediation,
2742 e.g. avoiding the risk of discrimination for insurance or treatment purposes.

2743 AI should help to ensure equal understanding of safety data and its relevance to all patients,
2744 irrespective of social circumstances and background, and the understanding of benefits and risks to
2745 specific individuals or subgroups of the population.

2746 This presents a unique challenge for vigilance, particularly in identifying rare, unexpected anomalies,
2747 the so-called Black Swan incidents.²⁶² While current PV systems are well-equipped to anticipate,
2748 assess, and manage common safety risks, they must also adapt in detecting these outlier events,
2749 particularly where advanced AI systems are deployed.

2750

2751 *Governance & Accountability*

2752 Chapter 9 of this report covers Governance & Accountability including a governance framework grid
2753 for the lifecycle phases of AI solutions in PV.

2754 The accelerated integration of AI underscores the need for dynamic, risk-based governance
2755 frameworks capable of near-real-time interventions.

2756 This is especially true as AI systems become more autonomous and self-determining, for example,
2757 with automated patient or HCP alerts, which will self-monitor their function and output and take
2758 preventative measures based on self-detected alerts. Such advancements raise critical questions:
2759 how will governance, accountability, and human oversight of PV of these new technologies evolve in
2760 tandem with these capabilities?

2761 Ideally, regulatory authorities and industry leaders in PV will establish robust oversight mechanisms
2762 to ensure that AI systems in PV are developed and deployed responsibly. Safeguards must be in place
2763 to protect against data misuse, uphold privacy standards, and ensure these technologies ultimately
2764 enhance outcomes for patients.

2765 The growing autonomy of AI in PV further emphasizes the need for adaptable regulatory
2766 frameworks. Continuous surveillance, proactive auditing, and rigorous inspection protocols will be
2767 essential to mitigate risks, uphold patient safety, and protect public health. Achieving this will require

2768 a shift toward governance models that are as agile and responsive as the technologies they seek to
2769 manage.

2770 **Conclusions to the future considerations for development and deployment of** 2771 **artificial intelligence in pharmacovigilance**

2772 Proliferation and deployment of AI and its integration into PV is set to give a paradigm shift in this
2773 discipline, which is likely to be focused on rapid or real-time data collection, assessment and
2774 reporting; for example, with the ability to analyse and extract vast quantities of safety data for case
2775 reporting and signal detection purposes at a rapid pace. This could fundamentally change the way we
2776 work to take advantage of these technological advances, for example, streamlining processes and
2777 causing changes in the wider healthcare environment and beyond, including patient privacy.

2778 Along with the enormous potential for AI in PV, there are many challenges which warrant future
2779 consideration, particularly around oversight of autonomous AI systems, and how AI may impact data
2780 privacy and ethical frameworks. It is critical that the guiding principles outlined in this report remain
2781 as core considerations, but with the understanding that they will need to evolve and adapt with
2782 advancements and application of AI in PV and medicine in general. This is to ensure AI use in PV
2783 remains unbiased, transparent, and secure to prevent misuse or accidental harm. The appropriate
2784 human oversight, including regulatory and ethical safeguards, will be as crucial as the technological
2785 advancements being applied.

References

- 245 Kjoersvik O, Bate A. *Black Swan Events and Intelligent Automation for Routine Safety Surveillance*. *Drug Saf*. 2022 May;45(5):419-427. doi: 10.1007/s40264-022-01169-0. ([Article](#) accessed 28 April 2025)
- 246 Ventola CL. *The nanomedicine revolution: part 1: emerging concepts*. *P T*. 2012 Sep;37(9):512-25. ([Abstract](#) accessed 28 April 2025)
- 247 Ventola CL. *The nanomedicine revolution: part 2: current and future clinical applications*. *P T*. 2012 Oct;37(10):582-91. PMID: 23115468; PMCID: PMC3474440. ([Abstract](#) accessed 28 April 2025)
- 248 Ventola CL. *The nanomedicine revolution: part 3: regulatory and safety challenges*. *P T*. 2012 Nov;37(11):631-9. ([Abstract](#) accessed 28 April 2025)
- 249 Ventola CL. *Big Data and Pharmacovigilance: Data Mining for Adverse Drug Events and Interactions*. *P T*. 2018 Jun;43(6):340-351. ([Abstract](#) accessed 28 April 2025)
- 250 Shneiderman, Ben, *Human-Centered Artificial Intelligence: Reliable, Safe & Trustworthy*, arXiv.org (February 23, 2020). <https://arxiv.org/abs/2002.04087v1> (Extract from forthcoming book by the same title) Copyright Ben Shneiderman 2020 ([Webpage](#) accessed 28 April 2025)
- 251 *Eliminating Bias in AI-Assisted Decisions: Chouldechova, A., & Roth, A. (2020). The Frontiers of Fairness in Machine Learning, Communications of the ACM, discusses the ways AI can counteract cognitive biases in human decision-making.* ([Website](#) accessed 28 April 2025)
- 252 Hauben, Manfred. *Artificial Intelligence and Data Mining for the Pharmacovigilance of Drug–Drug Interactions, Clinical Therapeutics, Volume 45, Issue 2, 117 - 133. DOI: 10.1016/j.clinthera.2023.01.002* ([Article](#) accessed 28 April 2025)
- 253 Xiaolu Zhu, Zheng Wang, Fang Teng, *A review of regulated self-organizing approaches for tissue regeneration, Progress in Biophysics and Molecular Biology, Volume 167, 2021, Pages 63-78, ISSN 0079-6107, <https://doi.org/10.1016/j.pbiomolbio.2021.07.006>.* ([Article](#) accessed 28 April 2025)
- 254 Wang C, He T, Zhou H, Zhang Z, Lee C. *Artificial intelligence enhanced sensors - enabling technologies to next-generation healthcare and biomedical platform*. *Bioelectron Med*. 2023 Aug 2;9(1):17. doi: 10.1186/s42234-023-00118-1. ([Article](#) accessed 28 April 2025)
- 255 Mir, M., Permyakova, A., Das, R., & Ünal, A. B. (2018). *Smart biomaterials for tissue engineering: Functional and antibacterial nanostructures*. *Biomaterials Science*. Volume: 6, Issue: 9, Pages: 2382–2395.
- 256 *Discussion on advanced sensors for healthcare 2022-2024: “The Emergence of Smart Organs: Science and Fiction?” at National Institutes of Health, U.S. Department of Health and Human Services. This provided an overview of developing “smart” organ and tissue technology.*
- 257 Choudhury, A and Asan, O, *Human Factors: Bridging Artificial Intelligence and Patient Safety (October 5, 2020). Proceedings of the International Symposium on Human Factors and Ergonomics in Health Care. 2020;9(1):211-215. doi:10.1177/2327857920091007* ([Article](#) accessed 28 April 2025))

²⁵⁸ Kurzweil, Ray (2005). *The Singularity is Near: When Humans Transcend Biology*. Viking Press.

²⁵⁹ Bonabeau, E, Marco D, and Guy T, *Swarm Intelligence: From Natural to Artificial Systems* (New York, 1999; online edn, Oxford Academic, 12 Nov. 2020), <https://doi.org/10.1093/oso/9780195131581.001.0001> ([Abstract](#) accessed 28 April 2025)

²⁶⁰ *Imagine you are Living in the Age of Singularity*: Karakas, F (2021), *Creative Adventures*; 75 ([Webpage](#) accessed 28 April 2025)

²⁶¹ *Cybersecurity Stats: Facts And Figures You Should Know*. *Forbes*. ([Webpage](#) accessed 28 April 2025)

²⁶² Kjoersvik O, Bate A. *Black Swan Events and Intelligent Automation for Routine Safety Surveillance*. *Drug Saf*. 2022 May;45(5):419-427. doi: 10.1007/s40264-022-01169-0. ([Article](#) accessed 28 April 2025)

2786

2787

2788 APPENDIX 1: Glossary

2789 This glossary provides definitions specific to terms within the context of Artificial Intelligence use in
2790 pharmacovigilance. Refer to the International Council for Harmonization of Technical Requirements
2791 for Pharmaceuticals for Human Use (ICH) compiled by CIOMS in the [Glossary of ICH Terms and](#)
2792 [Definitions](#) and all other relevant glossaries available for any additional terms not described within
2793 this glossary.

2794

2795 **Accountability**

2796 Accountability applies to clearly defined roles, responsibilities and liability for organizations
2797 and/or individuals deploying, operating and managing artificial intelligence systems. It requires
2798 the adoption of appropriate governance measures by relevant stakeholders (including but not
2799 limited to Regulators, Vendors, Users, Developers, Data Providers or Pharmaceutical Company)
2800 involved in setting policy, developing, deploying, maintaining and managing artificial
2801 intelligence systems. This ensures operation within expected parameters throughout the
2802 artificial intelligence lifecycle as well as managing any unforeseen consequences.

2803 Proposed by CIOMS Working Group XIV.

2804

2805 **Adverse event**

2806 Any untoward medical occurrence in a patient or clinical investigation subject administered a
2807 pharmaceutical product and which does not necessarily have a causal relationship with this
2808 treatment. An adverse event (AE) can therefore be any unfavourable and unintended sign
2809 (including an abnormal laboratory finding), symptom, or disease temporally associated with
2810 the use of a medicinal (investigational) product, whether or not related to the medicinal
2811 (investigational) product.

2812 Adopted from: Council for International Organizations of Medical Sciences (CIOMS) Glossary of ICH terms and
2813 definitions. Geneva: Council for International Organizations of Medical Sciences. 2024. ([Full text accessed 4 April](#)
2814 [2025](#))

2815

2816 **Adverse reaction**

2817 A response to a medicinal product that is noxious and unintended, meaning a causal
2818 relationship between the product and the event is at least a reasonable possibility.

2819 Adopted from: Council for International Organizations of Medical Sciences (CIOMS) Glossary of ICH terms and
2820 definitions. Geneva: Council for International Organizations of Medical Sciences. 2024. ([Full text accessed 4 April](#)
2821 [2025](#))

2822

2823 **Artificial intelligence**

2824 An artificial intelligence (AI) system is a machine-based system that, for explicit or implicit
2825 objectives, infers, from the input it receives, how to generate outputs such as predictions,
2826 content, recommendations, or decisions that can influence physical or virtual environments.

2827 Adopted from: OECD (2024), "Explanatory memorandum on the updated OECD definition of an AI system", OECD
2828 Artificial Intelligence Papers, No. 8, OECD Publishing, Paris, <https://doi.org/10.1787/623da898-en>.

2829

2830 Note: In the context of pharmacovigilance, the use of AI systems and activities is aimed at enhancing drug
2831 safety monitoring, patient safety and regulatory compliance.

2832

2833 Artificial intelligence literacy

2834 Having the essential abilities needed to understand, learn and work in a digital world through
2835 AI-driven technologies.

2836 Adopted from: Davy Tsz Kit Ng, Jac Ka Lok Leung, Samuel Kai Wah Chu, Maggie Shen Qiao, Conceptualizing AI
2837 literacy: An exploratory review, *Computers and Education: Artificial Intelligence*, Volume 2, 2021, 100041, ISSN
2838 2666-920X, <https://doi.org/10.1016/j.caeai.2021.100041>.

2839

2840 Augmented intelligence / Intelligence augmentation

2841 Augmented intelligence is a conceptualization of artificial intelligence that focuses on artificial
2842 intelligence's assistive role. It emphasizes the use of artificial intelligence for enhancing, i.e.
2843 augmenting or amplifying human intelligence, rather than replacing it. Inherent in this view is
2844 the recognition that artificial intelligence and humans work together in a human-centered
2845 partnership, where each one can perform certain tasks better than either could alone.

2846 Combined from:

2847 - Madni AM. Augmented intelligence: A human productivity and performance amplifier in systems engineering and
2848 engineered human-machine systems. *Systems engineering for the digital age: practitioner perspectives*.
2849 2023;Oct8:375-391. ([Chapter abstract](#)) <https://doi/10.1002/9781394203314.ch17>

2850 - World Medical Association (WMA). *WMA Statement on augmented intelligence in medical care*.
2851 2019. ([Webpage](#) accessed 3 April 2025)

2852

2853 Automation bias or automation complacency

2854 Automation bias and automation complacency are overlapping manifestations of automation-
2855 induced phenomena, where human attention plays a central role. Both refer to the human
2856 tendency to favour suggestions from automated decision-making systems over non-automated
2857 contradictory information even when it is correct. They can involve attentional bias directed
2858 toward the automated output, or insufficient attention and monitoring of the automated
2859 output, especially in context of multi-tasking where manual tasks compete with the human
2860 expert's attention.

2861 Combined from:

2862 - Parasuraman R, Manzey DH. Complacency and bias in human use of automation: An attentional
2863 integration. *Human factors*. 2010Jun;52(3):381-410. ([Journal full](#)
2864 [text](https://doi.org/10.1177/0018720810376055)) <https://doi.org/10.1177/0018720810376055>

2865 - Cummings ML. Automation bias in intelligent time critical decision support systems. In *Decision making in aviation*.
2866 2017;Jul5;289-294. Routledge. ([Chapter abstract accessed 4 April 2025](#))

2867

2868 Bias

2869 The tendency of a measurement process to over- or under-estimate the value of a population
2870 parameter.

2871 Adopted from: Council for International Organizations of Medical Sciences (CIOMS). *Glossary of ICH terms and*
2872 *definitions*. Geneva: Council for International Organizations of Medical Sciences. 2024. ([Full text accessed 4 April](#)
2873 [2025](#))

2874 In AI, bias may be systematic difference in treatment of certain objects, people, or groups in
2875 comparison to others [18]. Bias can be introduced into study design, conduct or analysis.
2876 Sources of bias include selection bias (of study sample), operational bias, and analyses that do
2877 not account for missing data.

2878 Adopted from: International Medical Device Regulators Forum (IMDRF). *Machine Learning-enabled Medical*
2879 *Devices—A subset of Artificial Intelligence-enabled Medical Devices: Key Terms and Definitions*. 2021. ([Full](#)
2880 [text](#) accessed 3 April 2025)

2881 In the context of artificial intelligence, bias can occur when the artificial intelligence data or
2882 algorithms reflect or perpetuate existing social inequalities, leading to discriminatory or unfair
2883 artificial intelligence outputs.

2884 Adopted from: University of Saskatchewan. Generative Artificial Intelligence: Glossary of AI Related Terms.
2885 ([Webpage](#) accessed 4 April 2025)

2886

2887 **Black-Box model**

2888 An analytics model that provides results based on received data but the logic used to provide
2889 those results cannot be determined or inferred on how it achieved those results.

2890 Proposed by CIOMS Working Group XIV.

2891

2892 **Black Swan event**

2893 Event of extreme impact that, although outside the realm of regular expectations (i.e.
2894 prospectively unpredictable), prompts humans to concoct explanations for its occurrence after
2895 the fact, making it seemingly explainable and predictable (i.e. retrospectively distorted).

2896 Combined from:

2897 - Kjoersvik O, Bate A. Black swan events and intelligent automation for routine safety surveillance. *Drug*
2898 *Safety*.2022;May;45(5):419-427. ([Journal full text](#))

2899 - Taleb NN. Black swans and the domains of statistics. *The American statistician*. 2007;Aug 1;61(3):198-200. (Journal
2900 abstract) <https://doi.org/10.1198/000313007X219996>

2901

2902 **Business continuity plan**

2903 Set of provisions and systems for the prevention of / recovery from events that could severely
2904 impact on an organisation's staff and infrastructure in general or on the structures and
2905 processes for pharmacovigilance in particular, including the urgent exchange of information
2906 within an organisation, amongst organisations sharing pharmacovigilance tasks as well as
2907 between marketing authorisation holders and competent authorities.

2908 Adopted from: Heads of Medicines Agencies (HMA). European Medicines Agency (EMA). Guideline on good
2909 pharmacovigilance practices (GVP): Module I – Pharmacovigilance systems and their quality systems. 2012. ([Full](#)
2910 [text](#) accessed 3 April 2025)

2911

2912 **Change management**

2913 Change Management describes processes, methods and techniques designed and used to
2914 plan, implement and control changes to organizational structures and/or business processes.
2915 Methodologies span around people, process and culture.

2916 Typically Change Management includes following components: Leadership alignment,
2917 Stakeholder engagement, Communication, Training, Impact Assessment, Continuous
2918 improvement.

2919 Adopted from: International Organization for Standardization (ISO). *What is change management: a Quick guide*.
2920 ([Webpage](#) accessed 3 April 2025)

2921

2922 **Class imbalance**

2923 Imbalance between categories in classification tasks. This affects model performance metrics,
2924 e.g. by the fact that a model always predicting the same outcome will be 99% accurate if 99%
2925 of test cases belong to the corresponding class.

2926 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
2927 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

2928

2929 Cluster analysis

2930 An unsupervised machine learning method that groups data elements based on similarities to
2931 identify patterns that are not immediately evident.

2932 Proposed by CIOMS Working Group XIV.

2933

2934 Computerized system validation

2935 Process of establishing and documenting that the specified requirements of a computerized
2936 system are fulfilled consistently from design until decommissioning of the system and/or
2937 transition to a new system. The approach to validation should focus on a risk assessment that
2938 takes into consideration the intended use of the system and the potential of the system to
2939 affect human subject protection and reliability of trial results.

2940 Adopted from: International Council for Harmonisation of Technical Requirements for Pharmaceuticals for Human
2941 Use (ICH). *Integrated Addendum to ICH E6(R1): Guideline for good clinical practice E6(R2)*. 2016. ([Full text](#) accessed
2942 3 April 2025)

2943

2944 Confirmation bias

2945 Confirmation bias is the tendency to give greater weight to data that support preliminary
2946 assumptive results, while failing to seek or dismissing contradictory evidence.

2947 Adopted from: Elston DM. Confirmation bias in medical decision-making. *Journal of the American Academy of*
2948 *Dermatology*. 2020;Mar1;82(3):572. ([Journal full text](#)) <https://doi:10.1016/j.jaad.2019.06.1286>

2949

2950 Cross-validation

2951 Resampling method used to assess the generalisation ability of a machine learning model and
2952 prevent overfitting.

2953 Adopted from: D. Cross-Validation. *Preprint submitted to Encyclopedia of Bioinformatics and Computational Biology,*
2954 *2nd edition (Elsevier)*. 2019;542-545. ([Full text](#) accessed 3 April 2025).

2955

2956 Note: This is an alternative to maintaining separate training and validation data sets to provide a more efficient
2957 use of data during development.

2958

2959 Data anonymization

2960 Anonymisation of personal data is the process whereby both direct and indirect personal
2961 identifiers are removed, and technical safeguards are used to ensure zero risk of re-
2962 identification.

2963 Adopted from: World Health Organization (WHO). *Ethics and governance of artificial intelligence for*
2964 *health: Guidance on large multi-modal models*. Geneva: World Health Organization. 2024. ([Webpage](#) accessed 3
2965 April 2025)

2966

2967 Data drift

2968 Change in the input data distribution a deployed model receives over time, which can cause
2969 the model's performance to degrade. This occurs when the properties of the underlying data
2970 change. Data drift can affect the accuracy and reliability of predictive models.

2971 Adopted from: U.S. Food and Drug Administration. *FDA Digital Health and Artificial Intelligence Glossary–*
2972 *Educational Resource*. 2024. ([Webpage](#) accessed 3 April 2025)

2973

- 2974 **Data privacy**
- 2975 Data privacy refers to measures taken to protect the fundamental right of individuals to the
- 2976 protection of their personal information. In the setting of PV, these measures emphasize the
- 2977 protection of sensitive and personal data (including health data).
- 2978 Proposed by CIOMS Working Group XIV.
- 2979
- 2980 **Decision tree**
- 2981 A model that uses a tree-like structure where data is progressed through various pre-defined
- 2982 tests or attributes to reach a final decision or prediction.
- 2983 Proposed by CIOMS Working Group XIV.
- 2984
- 2985 **Deep learning**
- 2986 A variant of machine learning involving neural networks with multiple layers of processing
- 2987 units known as artificial neurons, or ‘perceptrons’ (nodes), which together facilitate extraction
- 2988 of higher features of unstructured input data (for example, images, video and text).
- 2989 Adopted from: Thirunavukkarasu AJ, Ting DS, Elangovan K, Gutierrez L, Tan TF, Ting DS. Large language models in
- 2990 medicine. *Nature medicine*. 2023;Aug;29(8):1930-1940. (Journal full text) [https://doi.org/10.1038/s41591-023-](https://doi.org/10.1038/s41591-023-02448-8)
- 2991 [02448-8](https://doi.org/10.1038/s41591-023-02448-8)
- 2992 Approach to creating rich hierarchical representations through the training of neural networks
- 2993 with many hidden layers.
- 2994 Adopted from: European Medicines Agency (EMA). Reflection paper on the use of Artificial Intelligence (AI) in the
- 2995 medicinal product lifecycle. 2024. (Full text accessed 3 April 2025)
- 2996
- 2997 **Explainability**
- 2998 The degree to which humans can understand the factors and logic that have led to a specific
- 2999 outcome or that play a role in the general operation of an AI system.
- 3000 Proposed by CIOMS Working Group XIV.
- 3001
- 3002 **Fairness and equity**
- 3003 Fairness and Equity requires awareness and adherence to the ideas of impartiality, equality,
- 3004 non-discrimination, diversity, justice and lawfulness. Avoidance and mitigation of unfair bias,
- 3005 discriminatory or unjust social wellbeing and environmental impacts and/or outcomes should
- 3006 be considered throughout the whole artificial intelligence lifecycle.
- 3007 Proposed by CIOMS Working Group XIV.
- 3008
- 3009 **False negative**
- 3010 The determination of a data point not belonging to a class of interest when the reference or
- 3011 test set states that does belong.
- 3012 Proposed by CIOMS Working Group XIV.
- 3013
- 3014 **False positive**
- 3015 The determination of a data point belonging to a class of interest when the reference or test
- 3016 set states that does not belong.
- 3017 Proposed by CIOMS Working Group XIV.
- 3018

- 3019 **Feature**
- 3020 A measurable property or characteristic of the data or engineered through data processing or
- 3021 transformation of the data that is used to train a model.
- 3022 Proposed by CIOMS Working Group XIV.
- 3023 **Error! Bookmark not defined.**
- 3024 **Few-shot learning**
- 3025 AI developed to complete tasks with exposure to only a few initial examples of the task, with
- 3026 accurate generalization to unseen examples.
- 3027 Adopted from: Thirunavukkarasu AJ, Ting DS, Elangovan K, Gutierrez L, Tan TF, Ting DS. Large language models in
- 3028 medicine. *Nature medicine*. 2023;Aug;29(8):1930-1940. ([Journal full text](https://doi.org/10.1038/s41591-023-02448-8)) [https://doi.org/10.1038/s41591-023-](https://doi.org/10.1038/s41591-023-02448-8)
- 3029 [02448-8](https://doi.org/10.1038/s41591-023-02448-8)
- 3030
- 3031 **Fuzzy logic**
- 3032 An approach to variable processing that allows for multiple possible truth values to be
- 3033 processed through the same variable. Fuzzy logic attempts to solve problems with an open,
- 3034 imprecise spectrum of data and heuristics that makes it possible to obtain an array of accurate
- 3035 conclusions.
- 3036 Proposed by CIOMS Working Group XIV.
- 3037
- 3038 **Generative artificial intelligence**
- 3039 Category of artificial intelligence techniques in which algorithms are trained on data sets that
- 3040 can be used to generate new content, such as text, images or video.
- 3041 Proposed by CIOMS Working Group XIV.
- 3042
- 3043 **Governance**
- 3044 A governance framework requires implementation of risk management practices and policies,
- 3045 responsible use, security, openness, fairness, and ethical practices to ensure adherence to the
- 3046 Artificial Intelligence Guiding Principles in the report of the CIOMS Working Group XIV.
- 3047 Proposed by CIOMS Working Group XIV.
- 3048
- 3049 **Hallucination**
- 3050 In natural language generation tasks, hallucinations are generated content that is either
- 3051 nonsensical or unfaithful to the provided source content.
- 3052 Adopted from: Huang L, Yu W, Ma W, Zhong W, Feng Z, Wang H, Chen Q, Peng W, Feng X, Qin B, Liu T. A survey on
- 3053 hallucination in large language models: Principles, taxonomy, challenges, and open questions. arXiv preprint
- 3054 arXiv:2311.05232. 2023 Nov 9.
- 3055
- 3056 **Human agency**
- 3057 Human agency is the capacity for human beings to make choices out of their own volition and
- 3058 to follow those choices to action.
- 3059 Proposed by CIOMS Working Group XIV.
- 3060
- 3061 **Human-in-command**
- 3062 The capability of a human to oversee the overall activity of an artificial intelligence system,
- 3063 including its broader economic, societal, legal and ethical impact, and the ability to decide if,
- 3064 when, and how to use an artificial intelligence system.

3065 Adopted from: European Commission. *Ethics guidelines for trustworthy AI*. 2019. ([Webpage](#) accessed 3
3066 April 2025)

3067

3068 **Human-in-the-loop**

3069 The capability for human intervention in every decision cycle of the artificial intelligence
3070 system.

3071 Adopted from: European Commission. *Ethics guidelines for trustworthy AI*. 2019. ([Webpage](#) accessed 3
3072 April 2025)

3073

3074 **Human-on-the-loop**

3075 The capability for human intervention during the design of an artificial intelligence system and
3076 monitoring of its operation.

3077 Adopted from: European Commission. *Ethics guidelines for trustworthy AI*. 2019. ([Webpage](#) accessed 3 April 2025)

3078

3079 **Human oversight**

3080 Human oversight refers to the expected role of humans in the design, implementation,
3081 monitoring, and analysis of AI in PV. It requires a framework to manage performance and to
3082 detect and mitigate potential issues related to the AI system.

3083 Proposed by CIOMS Working Group XIV.

3084

3085 **Individual Case Safety Report**

3086 A report containing information about a suspected adverse drug reaction related to the
3087 administration of one or more medicinal products to a specific patient at a specific time.

3088 Adopted from: Council for International Organizations of Medical Sciences (CIOMS) Glossary of ICH terms and
3089 definitions. Geneva: Council for International Organizations of Medical Sciences. 2024. ([Full text accessed 4 April](#)
3090 [2025](#))

3091

3092 **Interpretability (*See also Explainability*)**

3093 A description of the general principles and logic by which an AI model functions and arrives at
3094 its outcomes / predictions should be shared, or the lack of explainability should be
3095 acknowledged and its implications discussed.

3096 Proposed by CIOMS Working Group XIV.

3097 **Error! Bookmark not defined.**

3098 **Knowledge graph**

3099 A heterogeneous knowledge base consisting of triples (facts) each comprised of object pairs
3100 and connecting relationships modelled through graphs and ontologies (a standardized machine
3101 readable semantic framework for representing all objects, and their properties and
3102 relationships in a domain of knowledge), which extract new insights from existing data sets via
3103 their integration.

3104 Adopted from: Manfred Hauben, Mazin Rafi, Knowledge Graphs in Pharmacovigilance: A Step-By-Step Guide,
3105 Clinical Therapeutics, Volume 46, Issue 7, 2024, Pages 538-543.

3106

3107 **Large language model**

3108 A type of artificial intelligence model using deep neural networks to learn the relationships
3109 between words in natural language, using large datasets of text to train, these include those
3110 with or without decoders.

3111 LLMs can be classified according to their availability and accessibility into two main categories:
3112 open-source models (freely accessible to the public) and closed-source models (developed as
3113 commercial products and often necessitating licenses or subscriptions for use).

3114 Adopted from: Heads of Medicines Agencies (HMA). European Medicines Agency (EMA). *Guiding principles on the*
3115 *use of large language models in regulatory science and for medicines regulatory activities*. 2024. ([Full text](#) accessed
3116 4 April 2025)

3117

3118 **Local Interpretable Model-Agnostic Explanations (LIME)**

3119 A technique that approximates a black box machine learning model with a local, interpretable
3120 model to explain each individual prediction.

3121 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
3122 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

3123

3124 **Machine learning**

3125 Computational process of optimising the parameters of a model from data, which is a
3126 mathematical construct generating an output based on input data. Machine learning
3127 approaches include, for instance, supervised, unsupervised and reinforcement learning, using a
3128 variety of methods including deep learning with neural networks.

3129 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
3130 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

3131

3132 **(AI) Model**

3133 Mathematical or computational method with parameters (weights) arranged in an architecture
3134 that allows learning of patterns (features) from training data.

3135 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
3136 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

3137

3138 **Model drift**

3139 A process where the model performance changes overtime either in a positive or negative
3140 performance outcome.

3141 Adopted from: S. Wang, S. Schlobach, M. Klein, Concept drift and how to identify it J Web Semant: Sci Serv Agents
3142 World Wide Web, 9 (2011), [10.1016/j.websem.2011.05.003](https://doi.org/10.1016/j.websem.2011.05.003)

3143

3144 **Natural language processing**

3145 Field of artificial intelligence focusing on the interaction between computers and human
3146 language.

3147 Adopted from: Thirunavukkarasu AJ, Ting DS, Elangovan K, Gutierrez L, Tan TF, Ting DS. Large language models in
3148 medicine. *Nature medicine*. 2023;Aug;29(8):1930-1940. ([Journal full text](#)) [https://doi.org/10.1038/s41591-023-](https://doi.org/10.1038/s41591-023-02448-8)
3149 [02448-8](https://doi.org/10.1038/s41591-023-02448-8)

3150

3151 **Negative controls**

3152 A real-world data point sampled as not belonging to the class of interest or deliberately
3153 created to not trigger a positive response from an artificial intelligence model.

3154 Proposed by CIOMS Working Group XIV.

3155

3156 Neural network

3157 Computing system inspired by biological neural networks, comprising (nodes), edges/weights,
3158 activation functions usually arranged in layers, communicating with one another and
3159 performing transformations upon input data.

3160 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
3161 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

3162

3163 Overfitting

3164 Learning details from training data that cannot be generalised to new data.

3165 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
3166 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

3167

3168 Parameter, hyper- parameter

3169 Variable within a machine learning model that is updated — usually automatically — during
3170 training to maximize performance. In deep learning, parameters are the ‘weights’ or data
3171 transforming functions comprising neural network nodes.

3172 Adopted from: Thirunavukkarasu AJ, Ting DS, Elangovan K, Gutierrez L, Tan TF, Ting DS. Large language models in
3173 medicine. *Nature medicine*. 2023;Aug;29(8):1930-1940. ([Journal full text](#)) [https://doi.org/10.1038/s41591-023-](https://doi.org/10.1038/s41591-023-02448-8)
3174 [02448-8](https://doi.org/10.1038/s41591-023-02448-8)

3175 Hyper-parameters are parameters that are used to configure a model. Unlike model
3176 parameters, they cannot be directly estimated from data learning and must be set before
3177 training a machine learning model. Hyper-parameter tuning is a step often required to build
3178 effective ML models.

3179 Adopted from: Yang L, Shami A. On hyperparameter optimization of machine learning algorithms: Theory and
3180 practice. *Neurocomputing*. 2020 Nov 20;415:295-316.

3181

3182 Performance degradation

3183 When results from an artificial intelligence system either fail or diminish in their ability to
3184 achieve the expected or required results as achieved earlier.

3185 Proposed by CIOMS Working Group XIV.

3186

3187 Personal data

3188 ‘Personal data’ means any information relating to an identified or identifiable natural person
3189 (‘data subject’). Information such as a name, an identification number, location data, an online
3190 identifier or to one or more factors specific to the physical, physiological, genetic, mental,
3191 economic, cultural or social identity of that natural person are examples of personal data.
3192 Sensitive (personal) data refers to special categories of personal data.

3193 Adopted from: European Parliament, Council of the European Union. Regulation (EU) 2016/679 of the European
3194 Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of
3195 personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection
3196 Regulation), Art. 4(1). *Official Journal of the European Union*. 2016; L 119. ([Webpage](#) accessed 4 April 2025)

3197

3198 Pharmacovigilance

3199 The science and activities relating to the detection, assessment, understanding and prevention
3200 of adverse effects or any other drug related problem.

3201 Adopted from: CIOMS Cumulative Glossary, with a focus on Pharmacovigilance (Version 2.1). Geneva: Council for
3202 International Organizations of Medical Sciences. 2024. ([Full text](#)) <https://doi.org/10.56759/ocf1297>

3203

3204 Pharmacovigilance system

3205 System used by an organisation to fulfil its legal tasks and responsibilities in relation to
3206 pharmacovigilance and designed to monitor the safety of authorised medicinal products and
3207 detect any change to their risk-benefit balance.

3208 Adopted from: Heads of Medicines Agencies (HMA). European Medicines Agency (EMA). *Guideline on good*
3209 *pharmacovigilance practices (GVP) Module I – Pharmacovigilance systems and their quality systems*. 2012. ([Full](#)
3210 [text](#) accessed 3 April 2025)

3211

3212 Precision

3213 Proportion of retrieved samples which are annotated as positive controls in the reference set,
3214 calculated as the ratio between correctly classified positive controls and all samples assigned
3215 to that class. Precision is also known as positive predictive value (PPV).

3216 Adopted from: Hicks SA, Strümke I, Thambawita V, Hammou M, Riegler MA, Halvorsen P, Parasa S. On evaluation
3217 metrics for medical applications of artificial intelligence. *Scientific reports*. 2022;Apr8;12(1):5979. ([Journal full](#)
3218 [text](#))<https://doi.org/10.1038/s41598-022-09954-8>

3219

3220 Positive controls

3221 A real-world data point sampled as belonging to the class of interest or deliberately created to
3222 trigger a positive response from an artificial intelligence model.

3223 Proposed by CIOMS Working Group XIV.

3224

3225 Predictive models

3226 A machine learning algorithm that analyzes data to identify patterns and trends, allowing it to
3227 make predictions about future outcomes or events based on input data.

3228 Adopted from: De Hond AA, Leeuwenberg AM, Hooft L, Kant IM, Nijman SW, van Os HJ, Aardoom JJ, Debray TP,
3229 Schuit E, van Smeden M, Reitsma JB. Guidelines and quality criteria for artificial intelligence-based prediction
3230 models in healthcare: a scoping review. *NPJ digital medicine*. 2022;Jan10;5(1):2. ([Journal full text](#))
3231 <https://doi.org/10.1038/s41746-021-00549-7>

3232

3233 Quality management system

3234 Part of the pharmacovigilance system and consists of its own structures and processes. It shall
3235 cover organisational structure, responsibilities, procedures, processes and resources of the
3236 pharmacovigilance system as well as appropriate resource management, compliance
3237 management and record management.

3238 Adopted from: D. Cross-Validation. *Preprint submitted to Encyclopedia of Bioinformatics and Computational Biology*,
3239 *2nd edition (Elsevier)*. 2019;542-545. ([Full text](#) accessed 3 April 2025).

3240

3241 Real-world data

3242 Data relating to patient health status and/or the delivery of health care routinely collected
3243 from a variety of sources. Examples of RWD include data derived from electronic health
3244 records (EHRs); medical claims and billing data; data from product and disease registries;
3245 patient-generated data, including from mobile devices and wearables; and data gathered from
3246 other sources that can inform on health status (e.g. genetic and other biomolecular
3247 phenotyping data collected in specific health systems).

3248 Adopted from: Council for International Organizations of Medical Sciences (CIOMS) Glossary of ICH terms and
3249 definitions. Geneva: Council for International Organizations of Medical Sciences. 2024. ([Full text accessed 4 April](#)
3250 [2025](#))

3251

3252 Recall

3253 Proportion of positive controls correctly classified as such, calculated as the ratio between
3254 correctly classified positive controls and all positive controls. Also known as sensitivity or true
3255 positive rate (TPR).

3256 Adopted from: Hicks SA, Strümke I, Thambawita V, Hammou M, Riegler MA, Halvorsen P, Parasa S. On evaluation
3257 metrics for medical applications of artificial intelligence. *Scientific reports*. 2022;Apr8;12(1):5979. ([Journal full](#)
3258 [text](https://doi.org/10.1038/s41598-022-09954-8))<https://doi.org/10.1038/s41598-022-09954-8>

3259

3260 Reproducibility

3261 The ability to achieve consistent results when analysis is repeated under the same conditions.
3262 Data and computer codes are used to regenerate the results.

3263 Adopted from: National Academies of Sciences, Engineering, and Medicine; Policy and Global Affairs; Committee on
3264 Science, Engineering, Medicine, and Public Policy; Board on Research Data and Information; Division on Engineering
3265 and Physical Sciences; Committee on Applied and Theoretical Statistics; Board on Mathematical Sciences and
3266 Analytics; Division on Earth and Life Studies; Nuclear and Radiation Studies Board; Division of Behavioral and Social
3267 Sciences and Education; Committee on National Statistics; Board on Behavioral, Cognitive, and Sensory Sciences;
3268 Committee on Reproducibility and Replicability in Science. Reproducibility and Replicability in Science. *Washington*
3269 *(DC): National Academies Press (US)*; 2019; May7. Chapter 3, Understanding Reproducibility and Replicability.
3270 ([Chapter full text](#) accessed 4 April 2025)

3271

3272 Risk-based approach

3273 A risk-based approach acknowledges the potential hazards that artificial intelligence systems
3274 can pose and recognises that different use cases present varying types and levels of risk. This
3275 necessitates a risk assessment that identifies, prioritises, and manages potential risks that
3276 could negatively impact a pharmacovigilance system's behaviour and results, taking into
3277 consideration existing process controls. A risk is characterised by both the anticipated impact
3278 and the likelihood of negative outcomes.

3279 This approach also supports procedures to identify and reduce errors and biases in a way that
3280 is proportionate to their risk. It influences the implementation strategies of AI systems
3281 (including documentation, compliance, and record-keeping), which should generally be
3282 commensurate with the identified risk.

3283 Proposed by CIOMS Working Group XIV.

3284

3285 Robustness

3286 A system reliably performs to its intended objectives or requirements accounting for known
3287 variances in data.

3288 Proposed by CIOMS Working Group XIV.

3289

3290 Semantic vector

3291 A mathematical representation of a word, phrase, or document as an identifier, where the
3292 identifier's position in the high-dimensional space captures the meaning or relationship of that
3293 word/phrase, allowing artificial intelligence systems to understand the context and similarity
3294 between different pieces of text based on their meaning.

3295 Adopted from: Cohen T, Widdows D. Empirical distributional semantics: methods and biomedical
3296 applications. *Journal of biomedical informatics*. 2009;Apr1;42(2):390-405. ([Journal full](#)
3297 [text](https://doi.org/10.1016/j.jbi.2009.02.002)) <https://doi.org/10.1016/j.jbi.2009.02.002>

3298

3299 Sensitivity analysis

3300 An assessment technique used to evaluate how changes in input data or model parameters
3301 affect the output of an artificial intelligence model.

3302 Proposed by CIOMS Working Group XIV.

3303

3304 **SHapley Additive exPlanations (SHAP)**

3305 Explainable artificial intelligence framework that can provide model-agnostic local
3306 explainability for tabular, image, and text datasets.

3307 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
3308 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

3309

3310 Note: It is derived from cooperative game theory of payouts to groups of cooperating individuals.

3311

3312 **Signal**

3313 Information that arises from one or multiple sources (including observations and experiments),
3314 that suggests a new potentially causal association, or a new aspect of a known association,
3315 between an intervention and an event or set of related events, either adverse or beneficial,
3316 that is judged to be of sufficient likelihood to justify further action to verify.

3317 For the purposes of pharmacovigilance a signal is Information on a new or known side effect
3318 that may be caused by a medicine and is typically generated from more than a single report of
3319 a suspected side effect. It is important to note that a signal does not indicate a direct causal
3320 relationship between a side effect and a medicine, but is essentially only a hypothesis that,
3321 together with data and arguments, justifies the need for further assessment.

3322 Adopted from: Council for International Organizations of Medical Sciences (CIOMS) Glossary of ICH terms and
3323 definitions. Geneva: Council for International Organizations of Medical Sciences. 2024. ([Full text accessed 4 April](#)
3324 [2025](#))

3325

3326 **(Bioartificial) Smart organ technology**

3327 A series of enabling techniques that can be used to produce human organs based on bionic
3328 principles.

3329 Proposed by CIOMS Working Group XIV.

3330

3331 **Supervised learning**

3332 Machine learning that makes use of labelled data during training. (ISO/IEC DIS 22989).

3333 Adopted from: International Medical Device Regulators Forum (IMDRF). *Machine Learning-enabled Medical*
3334 *Devices—A subset of Artificial Intelligence-enabled Medical Devices: Key Terms and Definitions*. 2021. ([Full](#)
3335 [text](#) accessed 3 April 2025)

3336

3337 **Target levels**

3338 A numerical value that serves as a goal or benchmark for artificial intelligence systems to
3339 achieve or surpass during their performance evaluation.

3340 Proposed by CIOMS Working Group XIV.

3341

3342 **Test dataset**

3343 A subset of the data that is never shown to the machine learning model during training, used
3344 to verify what the model has learned. (Modified from ISO/IEC DIS 22989).

3345 Adopted from: International Medical Device Regulators Forum (IMDRF). *Machine Learning-enabled Medical*
3346 *Devices—A subset of Artificial Intelligence-enabled Medical Devices: Key Terms and Definitions*. 2021. (Full
3347 text accessed 3 April 2025)

3348 Data used to evaluate the performance of the AI system, before its deployment. It is expected
3349 to be similar to production data, and proper evaluation needs test data to be disjointed from
3350 any data used during development.

3351 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
3352 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

3353

3354 **Traceability (AI)**

3355 The ability to track and document the data and processes used to create an artificial
3356 intelligence model.

3357 Proposed by CIOMS Working Group XIV.

3358

3359 **Training**

3360 Process intended to establish or to improve the parameters of a machine learning model,
3361 based on a machine learning algorithm, by using training data. (Modified from ISO/IEC DIS
3362 22989).

3363 Adopted from: International Medical Device Regulators Forum (IMDRF). *Machine Learning-enabled Medical*
3364 *Devices—A subset of Artificial Intelligence-enabled Medical Devices: Key Terms and Definitions*. 2021. ([Full](#)
3365 [text](#) accessed 3 April 2025)

3366

3367 **Training dataset**

3368 Data used specifically in the context of machine learning: it serves as the raw material from
3369 which the machine learning algorithm extracts its model to address the given task.

3370 Adopted from: European Medicines Agency (EMA). *Reflection paper on the use of Artificial Intelligence (AI) in the*
3371 *medicinal product lifecycle*. 2024. ([Full text](#) accessed 3 April 2025)

3372

3373 **Transparency**

3374 Transparency regarding AI involves disclosing information between organizations or
3375 individuals. This includes sharing relevant documentation of the AI system lifecycle (i.e. design,
3376 development, evaluation, deployment, operation, re-training, maintenance and
3377 decommission) to facilitate traceability and providing stakeholders with enough information to
3378 have a general understanding of the AI system, its use, risks, limitations, and impact on their
3379 rights.

3380 Proposed by CIOMS Working Group XIV.

3381

3382 **Unsupervised learning**

3383 Machine learning that makes use of unlabelled data during training. (ISO/IEC DIS 22989)

3384 Adopted from: International Medical Device Regulators Forum (IMDRF). *Machine Learning-enabled Medical*
3385 *Devices—A subset of Artificial Intelligence-enabled Medical Devices: Key Terms and Definitions*. 2021. ([Full](#)
3386 [text](#) accessed 3 April 2025)

3387

3388 **Validity**

3389 Validity means that a system achieves its intended purpose within acceptable parameters. It
3390 requires predefining acceptable performance levels, selecting appropriate data for model
3391 training and/or testing, assessing model performance in a realistic setting and integrating the
3392 system into an ongoing quality assessment process.

3393 Proposed by CIOMS Working Group XIV.

3394

3395 **Validation dataset**

3396 Data used to tune hyperparameters or to validate some algorithmic choices (rule design, etc.).

3397 Adopted from: International Organization for Standardization (ISO). *ISO/IEC DIS 22989. Information technology —*
3398 *Artificial intelligence — Artificial intelligence concepts and terminology*. 2022. ([Webpage](#) accessed 4 April 2025)

3399

3400 **Zero-shot learning**

3401 Artificial intelligence developed to complete tasks without exposure to any previous examples
3402 of the task.

3403 Adopted from: Thirunavukkarasu AJ, Ting DS, Elangovan K, Gutierrez L, Tan TF, Ting DS. Large language models in
3404 medicine. *Nature medicine*. 2023;Aug;29(8):1930-1940. ([Journal full text](https://doi.org/10.1038/s41591-023-02448-8)) [https://doi.org/10.1038/s41591-023-](https://doi.org/10.1038/s41591-023-02448-8)
3405 [02448-8](https://doi.org/10.1038/s41591-023-02448-8)

3406

3407 **APPENDIX 2: Comparison table of guiding principles**

3408

3409 **Table 9: Comparison of CIOMS Working Group XIV guiding principles for artificial intelligence across regional and country government**
 3410 **institutions, and international organizations – Extracted description of principles**

3411 Source: CIOMS Working Group XIV
 3412

	Examples of regional - and country government institutions’, and international organisations’ principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
Human Oversight	AI systems should support human agency and human decision making, as prescribed by the principle of respect for human autonomy.	When an AI system significantly impacts a person, community, group or environment, there should be a timely process to allow people to challenge the use or outcomes of the AI system.	Human Oversight means that high-impact AI systems must be designed and developed in such a way as to enable people managing the operations of the system to exercise meaningful oversight. This includes a level of interpretability appropriate to the context.		AI systems should function in a robust, secure and safe way throughout the AI life cycle, and risks should be continually identified, assessed and managed. Where appropriate, users, impacted third parties and actors in the AI lifecycle should be able to contest an AI decision or outcome that is harmful or creates.	Automated systems...”in design and development, pre-development and on-going disparity testing and mitigation, and clear organizational oversight”	Formal processes for human control and review of automated decisions are mandatory.	The principle of autonomy requires that any extension of machine autonomy not undermine human autonomy. In the context of health care, this means that humans should remain in full control of health-care systems and medical decisions. Human oversight may depend on the risks associated with an AI system but should always be meaningful	Mechanisms should be in place, as appropriate, to ensure that if AI systems risk causing undue harm or exhibit undesired behaviour, they can be overridden, repaired, and/or decommissioned safely as needed.

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
								and should thus include effective, transparent monitoring of human values and moral considerations.	
Validity & Robustness	Technical robustness requires that AI systems are developed with a preventative approach to risks and that they behave reliably and as intended while minimising unintentional and unexpected harm as well as preventing it where possible. This should also apply in the event of potential changes in their operating environment or the presence of other agents (human or artificial) that may interact	AI systems should reliably operate in accordance with their intended purpose.	Validity means a high-impact AI system performs consistently with intended objectives. Robustness means a high-impact AI system is stable and resilient in a variety of circumstances.		Consider how the associated actors on the AI supply chain can regularly test or carry out due diligence on the functioning, resilience and security of a system. Provide tools and guidance for undertaking AI-related safety risk assessments and implementing appropriate mitigations.		AI interventions should follow scientific best practice including being reliable, reproducible, fair, honest, and accountable.	All algorithms should be tested rigorously in the settings in which the technology will be used in order to ensure that it meets standards of safety and efficacy. The examination and validation should include the assumptions, operational protocols, data properties and output decisions of the AI technology. There should be robust, independent	AI systems must function in a robust, secure and safe way throughout their lifetimes, and potential risks should be continually assessed and managed.

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
	with the AI system in an adversarial manner.							oversight of such tests and evaluation to ensure that they are conducted safely and effectively.	
Data Privacy	Principle of prevention of harm is privacy, a fundamental right particularly affected by AI systems. Prevention of harm to privacy also necessitates adequate data governance that covers the quality and integrity of the data used, its relevance in light of the domain in which the AI systems will be deployed, its access protocols and the capability to process data in a manner that	AI systems should respect and uphold privacy rights and data protection, and ensure the security of data.			Encourage AI developers and deployers (within their remit) to mitigate and build resilience to cybersecurity related risks throughout the AI life cycle. Encourage AI developers and deployers to consider and mitigate where possible potential malicious or criminal use of AI products and services.	“Data privacy...protections are included by default, including ensuring that data collection conforms to reasonable expectations and that only data strictly necessary for the specific context is collected”.	Privacy, confidentiality, and security of data use must be foundational to every AI development.	Data protection laws are “rights-based approaches” that provide standards for regulating data processing that both protect the rights of individuals and establish obligations for data controllers and processors.	

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
	protects privacy.								
Transparency	...transparency which encompasses three elements: 1) traceability, 2) explainability and 3) open communication about the limitations of the AI system.	There should be transparency and responsible disclosure so people can understand when they are being significantly impacted by AI, and can find out when an AI system is engaging with them.	Transparency means providing the public with appropriate information about how high-impact AI systems are being used. The information provided should be sufficient to allow the public to understand the capabilities, limitations, and potential impacts of the systems.	End-users of AI-MD (e.g. medical practitioners, patients) should be informed that they are interacting with an AI-MD.	Encourage AI developers and deployers (within their remit) to implement appropriate transparency and explainability measures.		Transparent approaches must always be used and communicated when developing AI algorithms. Everything must be as open and sharable as possible. Tools and underlying concept of Openness must be a feature and a critical success factor of any AI development.	AI should be intelligible or understandable to developers, users and regulators. Two broad approaches to ensuring intelligibility are improving the transparency and explainability of AI technology.	This principle is about transparency and responsible disclosure around AI systems to ensure that people understand when they are engaging with them and can challenge outcomes.
Accountability	The principle of accountability necessitates that mechanisms be put in place to ensure responsibility for the development, deployment	People responsible for the different phases of the AI system lifecycle should be identifiable and accountable for the	Accountability means that organizations must put in place governance mechanisms needed to ensure compliance with all legal obligations of high-impact AI systems in the context in which they will be used.	While developers should be responsible for the proper design of algorithms used in the AIMD, organisations using AI-MD to deliver care will be responsible for the decision to implement the AI-	See Governance.	"...should have access to timely human consideration and remedy by a fallback and escalation process if an automated system fails, it produces an error, or you would like to	See Validity & Robustness.	Although AI technologies perform specific tasks, it is the responsibility of human stakeholders to ensure that they can perform those tasks and that	Organisations and individuals developing, deploying or operating AI systems should be held accountable for their proper functioning.

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
	and/or use of AI systems.	outcomes of the AI systems, and human oversight of AI systems should be enabled.	This includes the proactive documentation of policies, processes, and measures implemented.	MD and the clinical outcomes arising from the use of AI-MD in ensuring that safe care is delivered. Similar to the implementation of any other MD, the use of AI-MD does not change the liability of the implementing institution or the individual medical professional in their provision of appropriate and safe care.		appeal or contest its impacts on you.”		they are used under appropriate conditions. Institutions have not only legal liability but also a duty to assume responsibility for decisions made by the algorithms they use, even if it is not feasible to explain in detail how the algorithms produce their results.	
Societal well-being	...the broader society, other sentient beings and the environment should be considered as stakeholders throughout the AI system's life cycle. Ubiquitous exposure to social AI systems in all	AI systems should benefit individuals, society and the environment.		Safeguards in the design, development, and implementation of AI-MD should be put in place to ensure that patients' interests, including their safety and well-being, are protected.			Actions and solutions must be people centred and not be used solely by itself. As one of many technologies to aid public health AI should respect the rights of the individual.	AI technologies should not harm people. They should satisfy regulatory requirements for safety, accuracy and efficacy before deployment, and measures should be in place to ensure quality control	This Principle highlights the potential for trustworthy AI to contribute to overall growth and prosperity for all – individuals, society, and planet – and advance global development objectives.

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
	areas of our lives (be it in education, work, care or entertainment) may alter our conception of social agency, or negatively impact our social relationships and attachment.							and quality improvement. Thus, funders, developers and users have a continuous duty to measure and monitor the performance of AI algorithms to ensure that AI technologies work as designed and to assess whether they have any detrimental impact on individual patients or groups.	
Environmental well-being	See Societal well-being.	See Societal well-being.						AI systems should be designed to minimize their ecological footprints and increase energy efficiency, so that use of AI is consistent with society's efforts to reduce the impact of human beings on the earth's	See Societal well-being.

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
								environment, ecosystems and climate.	
Fairness & Equity	...enable inclusion and diversity throughout the entire AI system's life cycle. AI systems (both for training and operation) may suffer from the inclusion of inadvertent historic bias, incompleteness, and bad governance models.	AI systems should be inclusive and accessible, and should not involve or result in unfair discrimination against individuals, communities or groups.	Fairness and Equity means building high-impact AI systems with an awareness of the potential for discriminatory outcomes. Appropriate actions must be taken to mitigate discriminatory outcomes for individuals and groups.	The development and implementation of AI-MD should not result in discriminatory or unjust clinical impact on patients across different demographic lines (e.g. race, gender, etc.).	AI systems should not undermine the legal rights of individuals or organisations, discriminate unfairly against individuals or create unfair market outcomes. Actors involved in all stages of the AI life cycle should consider descriptions of fairness that are appropriate to a system's use, outcomes and the application of relevant law.	"Algorithmic discrimination protections...should include proactive equity assessments as part of the system design, use of representative data and protection against proxies for demographic features, ensuring accessibility for people with disabilities...".	Fairness, equality and inclusiveness in impact and design should always form the foundation of any AI initiative for Public Health. Discussions, developments, and implementation must be grounded in the globally-agreed ethical principles of human dignity, beneficence, nonmaleficence and justice.	Inclusiveness requires that AI used in health care is designed to encourage the widest possible appropriate, equitable use and access, irrespective of age, gender, income, ability or other characteristics. AI developers should be aware of the possible biases in their design, implementation and use and the potential harm that biases can cause to individuals and society.	AI systems should be designed in a way that respects the rule of law, human rights, democratic values and diversity, and should include appropriate safeguards to ensure a fair and just society.
Explainability	...the ability to explain both the technical processes of the AI system and	See Transparency		The decisions or recommendations from an AI-MD should endeavour to be explainable	See Transparency	"Automated systems should provide explanations that are technically		See Transparency	See Transparency

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
	the reasoning behind the decisions or predictions that the AI system makes.			and reproducible. The level of explainability is dependent on the varying expectations of the end user and the risks of the AI-MD. End-users should be consulted during the development or adoption of the AI-MD to ensure the explainability meets their expectations.		valid, meaningful and useful to you and to any operators or others who need to understand the system, and calibrated to the level of risk based on the context... in plain language and assessments of the clarity and quality of the notice and explanations should be made public whenever possible."			
Safety	See Validity & Robustness.	See Validity & Robustness.	Safety means that high-impact AI systems must be proactively assessed to identify harms that could result from use of the system, including through reasonably foreseeable misuse. Measures must be taken to mitigate the risk of harm.		Enable AI deployers (within their remit) and end users to make informed decisions about the safety of AI products and services. Communicate the level of safety related risk in their remit by appropriately identifying, monitoring, communicating	"Automated systems...should be designed to proactively protect you from harms stemming from unintended, yet foreseeable, uses or impacts...".		Preventing harm requires that use of AI technologies does not result in any mental or physical harm.	See Validity & Robustness.

	Examples of regional - and country government institutions', and international organisations' principles								
Principle	EU ^{263,264}	Australia ²⁶⁵	Canada ²⁶⁶	Singapore ²⁶⁷	UK ²⁶⁸	US ²⁶⁹	PAHO ²⁷⁰	WHO ²⁷¹	OECD ²⁷²
					and acting upon risks.				
Governance	See Data Privacy.				Governance measures could be put in place to ensure effective oversight of the supply and use of AI systems, with clear lines of accountability established across the AI life cycle.			Human rights standards, data protection laws and ethical principles are all necessary to guide, regulate and manage the use of AI for health by developers, governments, providers and patients.	

3413

References

²⁶³ Reflection paper on the use of Artificial Intelligence (AI) in the medicinal product lifecycle. European Medicines Agency, Committee for Medicinal Products for Human Use (CHMP), 9 September 2024. ([PDF](#) accessed 27 April 2025)

²⁶⁴ European Commission: Directorate-General for Communications Networks, Content and Technology, The Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self assessment, Publications Office, 2020, <https://data.europa.eu/doi/10.2759/002360> ([PDF](#) accessed 27 April 2025)

References

²⁶⁵ Australian Government, Department of Industry, Science and Resources, Australia's Artificial Intelligence Ethics Principles ([Webpage](#) accessed 27 April 2025)

²⁶⁶ Government of Canada, The Artificial Intelligence and Data Act (AIDA) – Companion document ([Webpage](#) accessed 27 April 2025)

²⁶⁷ Artificial intelligence in healthcare guidelines. Ministry of Health, Singapore Government Agency. ([PDF](#) accessed 27 April 2025)

²⁶⁸ Implementing the UK's AI Regulatory Principles Initial Guidance for Regulators. Department for Science, Innovation & Technology. UK Government. February 2024. ([PDF](#) accessed 27 April 2025)

²⁶⁹ Blueprint for an AI Bill of Rights, Making automated systems work for the American people. The White House. ([Website](#) accessed 27 April 2025)

²⁷⁰ Artificial Intelligence in Public Health, Digital Transformation Toolkit. Pan American Health Organization, 2021. PAHO/EIH/IS/21-011. ([PDF](#) accessed 27 April 2025)

²⁷¹ Ethics and governance of artificial intelligence for health: WHO guidance. Geneva: World Health Organization; 2021. Licence: CC BY-NC-SA 3.0 IGO. ([Website](#) accessed 27 April 2025)

²⁷² OECD AI Principles overview © 2025 OECD. ([Website](#) accessed 27 April 2025)

3414 APPENDIX 3: Use cases

3415 Use Case A: Large Language Models data extraction for case processing

3416 Source:²⁷³

3417 Area of PV: ICSR Processing

3418 A1. Business rational and challenges

3419 As per Good Pharmacovigilance Practices (GVP), pharmaceutical companies must act on potential
3420 adverse reactions to drugs. With significant increases in the number of case reports in recent years,
3421 case intake/processing operations face complex challenges beyond the number of cases, such as
3422 handling very diverse data sources including unstructured texts and scanned documents or managing
3423 sudden peak inflows with a finite workforce. With the complexity of the relevant data points ranging
3424 from simple demographics to more complex lab values, simpler technology approaches like Named
3425 Entity Recognition were unsuccessful in consistently improving case intake/processing operations
3426 under real-world circumstances. The use of Large Language Models (LLMs) in case intake/processing
3427 provides potential to advance processes without compromising quality.

3428 A2. Solution

3429 A pharmaceutical company executed a proof-of-concept study to assess the feasibility as well as the
3430 quantitative and qualitative business impact of utilizing LLMs for case intake purposes. Specifically,
3431 LLMs were applied for data extraction from source documents for case intake and processing while
3432 covering regulatory and compliance aspects.

3433 To process the selected source documents and extract pre-defined pieces of information, a 3-step
3434 semi-automatic processing pipeline was set up. The pipeline consisted of (1) pre-processing steps to
3435 unify the input for the LLM (OpenAI's GPT-4), (2) a JSON-formatted extraction template that guided
3436 the LLM in structuring the information as well as providing hints regarding the location of the
3437 information in the source data, and (3) post-processing steps to match the model output with fields
3438 where predefined values were applicable. Original source documents were augmented by references
3439 and highlighting of extracted key terms.

3440 For the assessment of the business impact of using LLMs for case intake, a selection representative
3441 cases was identified. A graphical user interface (GUI) was designed for the purpose of comparing the
3442 processing performance of (a) the fully manual process vs (b) the manual process augmented by
3443 fields pre-filled by the results of the LLM extraction pipeline. Four experienced professionals were
3444 randomly assigned to either process version (a) or (b). The processing times were tracked for each
3445 source document to derive the overall processing time regarding extraction of the representative set
3446 of fields.

3447 A3. Results

3448 In this study, two key results were derived from the implementation of LLMs in the case intake and
3449 processing operations:

3450 The first result focused on the performance of the LLM model, measured through the match scores
3451 of all extracted fields and averaged across cases of a category for the full number of source
3452 documents in scope of this study. The statistical evaluation revealed that the model achieved match
3453 scores, ranging from 85% to 100% for clinical studies, and 60% to 100% for patient support programs
3454 (PSP) cases. For literature cases, while the sample size precludes a robust statistical evaluation,
3455 model performance ranges from 67 to 100%, suggesting qualitative results that align with the other
3456 types.

3457 The high match scores achieved by the model demonstrate its capability to extract accurate and
 3458 relevant information from unstructured sources. This can be translated into tangible efficiency gains
 3459 for business operations.

3460 The second result highlighted the efficiency gains identified in the business impact assessment. The
 3461 implementation of LLM in case intake led to an estimated efficiency gain of 39%, translating to time
 3462 savings of approximately 20 minutes per case. Specifically, the study found that the average number
 3463 of data points extracted per case was 69.4, with only 2.4 data points requiring manual correction.

3464 Implementing LLMs is not just a technical enhancement; it represents a strategic move towards
 3465 improving operational efficiency and ensuring high-quality outcomes in PV practices.

3466 *A4. Compliance with the governance framework*

3467 **Table 10: Use case A: Alignment with the governance framework (detail)**

3468 Source: CIOMS Working Group XIV

3469

Principle	Activities	SPEC	DEV	PreD	PstD	RU
Risk-based approach	A risk-based approach has been followed thoroughly and a 100% QC by human interaction has been applied.	A	A	N/A	N/A	N/A
Human oversight	Implementation of dedicated features to support human oversight, including user-friendly interfaces and references to the source data. The 100% human QC ensures robustness of all extraction outputs.	A	A	N/A	N/A	N/A
Validity and robustness	No continuous learning is applied; rather, the model is used in a locked state. Releases of new versions are quality assured on a sufficiently broad test set to derive.	A	A	N/A	N/A	N/A
Transparency	Model performance has been measured with match score. The correction of the failures can be used as feedback in regular intervals to improve the prompting strategy.	A	A	N/A	N/A	N/A
Data Privacy	The service is established on a private cloud. Access is provided only to project team members. Personally identifiable information is redacted prior to the actual data extraction step.	A	A	N/A	N/A	N/A
Fairness and Equity	Not applicable. The application is not providing any data consolidation or decision support. The 1:1 match of the data extraction is verified by the human QC.	N/A	N/A	N/A	N/A	N/A
Governance & Accountability	The LLM model is using tailored prompting strategy maintained on vendor domain to test the data extraction. The model is provided by Open AI, and is powered by a selection of large language models ("LLMs"). The case intake and processing team takes over the accountability and perform the 100% human QC process. The ultimate accountability remains with the pharma company.	A	A	N/A	N/A	N/A

3470

3471 **Abbreviations**
 SPEC: Collection of specifications, requirements

3472 DEV: Development and change management

3473 PreD: Pre-deployment & post-change sign-off

3474 PstD: Post-deployment & post-change hyper-care

3475 RU: Routine Use

3476 A: Applicable

3477 NA: Not Applicable

References

²⁷³ Römning, H.-J. *et. How LLMs can advance safety case intake – points to consider and insights from a Proof of Concept. Submitted to “Therapeutic Advances in Drug Safety” as editorial article in Mar-2025.*

3478

3479 Use Case B: Case deduplication

3480 Source:²⁷⁴

3481 Area of PV: ICSR Processing

3482 B1. Business rational and challenges

3483 Adverse event reporting systems (AERS) are essential in PV as they support the identification and
3484 evaluation of safety signals related to the use of medical products. Expert review in safety monitoring
3485 involves several steps, such as data mining and case series analysis, which are significantly affected
3486 by the AERS data quality. A representative example of quality issues is duplication, where more than
3487 one report describes the same patient case and the same adverse event experience for the same
3488 product. Duplicate reports may result in false or missed safety signals and increase the workload for
3489 safety evaluators by misinterpreting the actual number of true adverse events and making a product-
3490 event relationship look weaker or stronger.

3491 B2. Solution

3492 A regulatory agency that maintains an AERS for drugs and biologics with >28 million historical reports
3493 and an average of 8,000 new submissions daily sought an efficient solution to deduplicate all
3494 historical and incoming adverse event reports. The regulatory agency collaborated with an academic
3495 partner to address this issue by developing a **deduplication pipeline** relying on modern technologies
3496 (mainly, natural language processing, network analysis, and cloud computing) and utilizing structured
3497 data and free-text narratives. The pipeline executes an initial pass to filter down the pairs of reports
3498 by placing minimum requirements on similarity based on demographic data and other features.
3499 Subsequently, a pairwise streamlined worker implementing a **duplicate detection algorithm**
3500 performs a probabilistic comparison of all qualifying report pairs and calculates two scores, a
3501 probabilistic weight score and a second component score value, that together rate how similar the
3502 two reports are. In the third step, the pairs exceeding a preselected validated threshold that was
3503 specified in a dedicated analysis are merged into networks (a.k.a. groups) of potentially duplicate
3504 reports and split into tightly linked communities (a.k.a groups) of actual duplicates. Finally, a
3505 reference case selection component identifies the most representative report in each duplicate
3506 group based on several parameters and the remaining reports in the group are flagged as duplicates
3507 and they are excluded from subsequent data mining calculations. An existing decision-support tool
3508 developed to support the case series analysis allows for evaluating the groups of duplicate reports
3509 and verifying the reference case, keeping medical reviewers in the loop.

3510 B3. Results

3511 In an early research study, the **duplicate detection algorithm** was applied to two datasets of post-
3512 market reports, one including vaccine product reports and one containing reports for biologics,
3513 identifying 77% of and 13% of known duplicate pairs, respectively, with (nearly) perfect precision in
3514 both cases (95% and 100%, respectively) (<https://pubmed.ncbi.nlm.nih.gov/28293864/>). This
3515 algorithm was refined in subsequent steps to reach acceptable levels of performance that, in some
3516 cases and based on new evaluations using drug adverse event reports, supported the detection of
3517 duplicate pairs with an F-measure >0.9. The medical reviewers who participated in this new
3518 evaluation round felt confident about the algorithm and expressed their interest in using it, as
3519 discussed in the corresponding publication
3520 (<https://www.frontiersin.org/articles/10.3389/fdsfr.2022.918897/full>). Subsequently, the medical
3521 reviewers generated a gold standard of 2300 reports with labelled duplicates in a systematic process
3522 to support the validation of the recently built **deduplication pipeline**, which was then compared with
3523 existing deduplication approaches used at the regulatory agency. The deduplication pipeline
3524 outperformed these approaches (results unpublished) and was approved for processing all historical
3525 reports and incoming live data in an ETL process. As of December 20, 2024, the pipeline, installed on

3526 the AWS environment and tightly integrated with the agency's AERS, has screened >30 million
 3527 historical reports and continues deduplicating an average of 8,000 new submissions daily.

3528 *B4. Compliance with the governance framework*

3529 **Table 11: Use case B: Alignment with the governance framework (detail)**

3530 Source: CIOMS Working Group XIV

3531

Principle	Activities	SPEC	DEV	PreD	PstD	RU
Risk-based approach	A risk-based approach has been discussed extensively, especially regarding missed or false positive duplicate reports. It has been determined that implementing the pipeline in the decision-support system, with humans-in-command, eliminates any risks for the case series analyses. What remains to be done is acknowledging any risks for data mining calculations and potential noise in signal detection; this part has not yet been fully developed and mostly affects the routine use of deduplication for data mining calculations and not its use in case series analyses that is currently fully implemented.	A	A	A	A	A
Human oversight	Human experts actively provided feedback to the software engineers during the development stage and evaluated the deduplication output to refine and validate the pipeline. Human experts can confirm or modify the reference case selection using an existing decision-support tool while conducting their case series analyses in the routine use setting. On the other hand, data mining calculations incorporate deduplication output without humans being involved.	A	A	A	A	A
Validity & Robustness	The deduplication pipeline has been evaluated and validated to ensure it meets expectations and serves its intended purpose. The effect of deduplicated data on data mining calculations and the discovery of potential safety signals, which is one of the major uses of deduplication output, has not yet been investigated.	A	A	A	A	A
Transparency	Several publications, technical reports, and other documentation describe the pipeline and results of all evaluations conducted with safety reviewers' assistance.	A	A	A	A	A
Data Privacy	Fully complying with the principle as all processing occurs in a secure cloud environment.	A	A	A	A	A
Fairness & Equity	The deduplication pipeline has been evaluated and validated in several rounds and is closely monitored in the post-deployment phase. The pipeline is fully migrated to the production environment to be routinely used at the time of writing this report; it is therefore marked as partially aligned since this process has not been completed yet.	A	A	A	A	A
Governance & Accountability	System administrators have full control and continuously monitor the deduplication pipeline as well as the use of its output in the decision-support tool. A plan has also been developed to incorporate the deduplication output in the data mining calculations.	A	A	A	A	A

	<p>Clearly defined roles were specified in the development, pre-deployment, and post-deployment stages, where the Contractor led the pipeline's construction and incorporation into the decision-support tool and the existing environment at the regulator's site, assisted by the end users and other stakeholders. Roles have not yet been fully assigned in the routine use setting.</p>					
--	--	--	--	--	--	--

3532

Abbreviations

3533

SPEC: Collection of specifications, requirements

3534

DEV: Development and change management

3535

PreD: Pre-deployment & post-change sign-off

3536

PstD: Post-deployment & post-change hyper-care

3537

RU: Routine Use

3538

A: Applicable

3539

NA: Not Applicable

References

²⁷⁴ Kreimeyer K, Spiker J, Dang O, De S, et al. Deduplicating the FDA adverse event reporting system with a novel application of network-based grouping. *J Biomed Inform.* 2025 May;165:104824. doi: 10.1016/j.jbi.2025.104824. ([PubMed](#) accessed 29 April 2025)

3540

3541 Use Case C: Artificial intelligence translation assistant**3542 Source:**²⁷⁵

3543 Area of PV: ICSR reporting

3544 *C1. Business rational and challenges*

3545 The pharma company had engaged with a vendor to consolidate and streamline the global case
3546 intake and translation process. The vendor had established 2 hubs in Europe and Asia to cover 16
3547 languages across 32 countries replacing a distributed network of multiple local country organizations
3548 and local vendors. To further increase productivity, the vendor had been requested to automate the
3549 translation process.

3550 *C2. Solution*

3551 While processing foreign language adverse event reports, about half of the effort was required for
3552 accurate translation of source documents from local languages to English, enabling centralized case
3553 management in English and subsequent submission to authorities. The pharma company and the
3554 vendor formed a common project team consisting of experts on ML and PV associates to pilot an AI-
3555 powered translation assistant based on commercially available technology. The team had set up a
3556 private cloud environment to store learning data (source texts and human-edited translations) and
3557 developed a user interface to input original text and retrieve and (if necessary) edit the result. The
3558 system automatically stores and analyses any modifications done by the users to enable further
3559 learning iteration and improvement of the first-time quality of the AI translation assistant. A 100%
3560 QC (Quality Control) by a human translator of all the translations was established to always verify the
3561 accuracy of the translation. The solution facilitates continuous learning through the automated
3562 integration of the manual edits into the translation model in defined regular intervals. With each
3563 model update the relevant quality measures (BLEU scores, see below) are re-calculated.

3564 *C3. Results*

3565 The translation's quality was assessed by BLEU scores. BLEU (Bilingual Evaluation Understudy) is a
3566 metric for evaluating machine-translated text. The BLEU score is a number between zero and one
3567 that measures the similarity of the machine-translated text compared to a set of high-quality
3568 reference translations. Within 6 months the AI translation assistant mimicked the quality of a human
3569 translator (i.e. BLEU equal or greater than 0.6).²⁷⁶

3570 The results of the AI Translation Assistant pilot for the first language (Portuguese) were leading to a
3571 reduction of translation efforts by ca. 30%. Hence the solution was extended to 5 further languages
3572 (Chinese, Spanish, French, German, and Dutch). Pharma company and vendor teams are jointly and
3573 continuously evaluating the BLEU score to monitor the quality of the solution.

3574 Improving the AI model is a function of case volume as every revised sample translation provided by
3575 the QC team helps to improve the model. More samples make better models, and better models
3576 finally reduce the effort for the team, allowing them to work through more cases, faster, and with
3577 greater consistency.

3578 Since the initial setup of the AI Assistant for Translation in 2021 the technology has further matured.
3579 The translations obtained by the assistant are generally of high quality and less than 10% corrections
3580 are still necessary by the human translators. An analysis is ongoing to move from 100% human QC to
3581 a sampling approach as a continuous monitoring measure. Depending on the detailed results per
3582 language, the sample size may be adjusted within reasonable borders (up to 100%) to ensure high
3583 quality translations continuously.
3584

3585 *C4. Compliance with the governance framework*3586 **Table 12: Use case C: Alignment with the governance framework (detail)**

3587 Source: CIOMS Working Group XIV

3588

Principle	Activities	SPEC	DEV	PreD	PstD	RU
Risk-based approach	A risk-based approach has been followed thoroughly and a 100% QC by human translators has been applied. With further maturing of the system and under close monitoring of the overall quality a reduction of human QC for individual translations should be possible.	A	A	A	A	A
Human oversight	To ensure human oversight a 100% human QC of the translated text by the vendor translators was established from the beginning. The BLEU scores are regularly measured for each language to identify changes in the overall performance.	A	A	A	A	A
Validity and robustness	The system has been implemented following the vendors standard validation approach. The 100% human QC ensures validation of all translation outputs. Any failure of the translation assistant would be immediately detected and corrected.	A	A	A	A	A
Transparency	Transparency of the translation performance is obtained as all translations are tracked by the system as well as any edits by the human translator. These edits are used in regular intervals to improve the model.	A	A	A	A	A
Data Privacy	The service is established on a private cloud. Access is provided only to project team members. Personally identifiable information is redacted prior to the actual translation process. The original source document remains available only for the local team who received the initial information and who may have to follow-up with the initial reporter.	A	A	A	A	A
Fairness and Equity	Not applicable. The application is not providing any data consolidation or decision support. The 1:1 match of the translation is verified by the human QC.	NA	NA	NA	NA	NA
Governance & Accountability	The translation assistant is a standalone tool owned by the vendor. Hence, the regular life-cycle governance is executed by the vendor and available on request to the pharma company. It concerns, e.g. the update of the model based on learning progress. While the responsibility for the execution of the translation lies with the vendor the ultimate accountability remains with the pharma company. Hence, in addition to the 100% human QC process by the vendor, the pharma company is doing a defined sample QC of the overall case intake results, including the translation.	A	A	A	A	A

3589 **Abbreviations**

3590 SPEC: Collection of specifications, requirements

3591 DEV: Development and change management

3592 PreD: Pre-deployment & post-change sign-off

3593 PstD: Post-deployment & post-change hyper-care

3594 RU: Routine Use

3595 A: Applicable

3596 NA: Not Applicable

References

²⁷⁵ Römning H-J, Pushparajan R. AI Translation Assistant for Pharmacovigilance. Poster presented at DIA Europe 2021. ([Full text](#) accessed 21 March 2025)

²⁷⁶ Google Cloud. Evaluate AutoML Translation models. Google Cloud. ([Webpage](#) accessed 21 March 2025)

3598 Use case D: Large language models for context-aware Structured Query 3599 Language

3600 **Source Article:**²⁷⁷

3601 **Area of PV:** Safety analysis

3602 *D1. Business rational and challenges*

3603 Safety scientists are often reliant on technical teams for safety query formulation and extraction of
3604 data from safety databases using SQL, which can introduce delays in assessment. The aim therefore
3605 was to enhance the accuracy of information retrieval from PV databases by employing LLMs to
3606 convert natural language queries (NLQs) into Structured Query Language (SQL) queries, leveraging a
3607 business context document.

3608 *D2. Solution*

3609 A sandboxed version of OpenAI's GPT-4 model was utilized within a retrieval-augmented generation
3610 (RAG) framework, enriched with a business context document, to transform NLQs into executable
3611 SQL queries. The study was conducted in three phases, varying query complexity, and assessing the
3612 LLM's performance both with and without the business context document.

3613 *D3. Results*

3614 The integration of a business context document markedly improved the LLM's ability to generate
3615 accurate SQL queries (i.e. both executable and returning semantically appropriate results), increasing
3616 from 8.3% with the database schema alone to 78.3% with the business context document. This
3617 enhancement was consistent across low, medium, and high complexity queries.

3618 The method is an assistive method to enable non-technical users to perform complex data queries,
3619 potentially enhancing timeliness of PV data analysis and reporting.

3620 *D4. Compliance with the governance framework*

3621 **Table 13: Use case D: Alignment with the governance framework (detail)**

3622 Source: CIOMS Working Group XIV

3623

Principle	Activities	SPEC	DEV	PreD	PstD	RU
Risk-based approach	Within this study, the intent is to demonstrate use of natural language to generate SQL queries to retrieve data from a safety database. The risk based approach should consider the feasibility of implementation and controls and processes needed to ensure its appropriate and trusted use.	A	A	N/A	N/A	N/A
Human oversight	Within the PoC human experts have reviewed the relevance of the outputs against a reference standard, within a product setting consideration will need to be given to how human oversight will be provided to ensure robustness of the outputs, including monitoring performance over time and at defined intervals (e.g. change in model version). Human oversight could be applied using a risk-based approach e.g. taking into consideration the intended use of the output.	A	A	N/A	N/A	N/A
Validity & Robustness	The tool has been evaluated against a curated reference standard, beyond the PoC consideration should be given to generalisability in production use.	A	A	N/A	N/A	N/A

Transparency	Whilst there is transparency of the GPT model, the use of RAG and context specific documentation provides transparency of the pipeline and how data is processed to achieve the output.	A	A	N/A	N/A	N/A
Data Privacy	This is an assistive tool not using individual patient data to generate SQL outputs.	N/A	N/A	N/A	N/A	N/A
Fairness & Equity	This is an assistive tool not using individual patient data to generate SQL outputs.	N/A	N/A	N/A	N/A	N/A
Governance & Accountability	During the proof of concept (PoC), the accountability of the methodology remains with the developer. However, if the methodology is integrated into a production setting, accountability would transition to the human subject matter expert. Governance within a PoC ensures that scientific integrity principles are adhered to, while future product use governance should cover how the tool fits into the overall PV system and quality management system (QMS).	A	A	N/A	N/A	N/A

3624 **Abbreviations**
3625 SPEC: Collection of specifications, requirements
3626 DEV: Development and change management
3627 PreD: Pre-deployment & post-change sign-off
3628 PstD: Post-deployment & post-change hyper-care
3629 RU: Routine Use
3630 A: Applicable
3631 NA: Not Applicable

References

²⁷⁷ Painter JL, Chalamalasetti VR, Kassekert R, Bate A. Automating pharmacovigilance evidence generation: using large language models to produce context-aware structured query language. *JAMIA open*. 2025;Feb;8(1):ooaf003. ([Journal full text](#)) <https://doi.org/10.1093/jamiaopen/ooaf003>

3632

3633 Use Case E: Causality assessment of adverse drug reactions

3634 **Source:**²⁷⁸

3635 **Area of PV:** Causality Assessment

3636 *E1. Business rational and challenges*

3637 Assessing the causal relationship between an adverse event and the patient’s exposure to a drug is a
3638 critical part of the PV process. Causality assessment is a time-consuming process requiring manual
3639 review by medical experts who evaluate data in the case with data from external sources (e.g. drug
3640 labels, scientific publications, drug mechanism of action, and disease symptoms). As the volume of
3641 adverse events to be reviewed increases an opportunity exists to create solutions that leverage ML
3642 to support the medical experts by predicting causality assessments.

3643 *E2. Solution*

3644 The authors of this paper created a modelling feature set comprising of various data attributes from
3645 solicited cases from the pharmaceutical company’s safety database relevant to causality assessment
3646 of drug-event combinations. This was supplemented by engineered data features comprising
3647 external data and data from other internal sources. The resulting training data schema (shown
3648 below) was selected as it provides a comprehensive set of features relevant to the causality
3649 assessment process.

3650 **Table 14: Use case E: Modeling Data**

3651 Source:²⁷⁹
3652

Modeling Data		
Case Level Data		External Sourced Data
Causality Label	Medical History Exclusions	Disproportionality
Rechallenge	Drug Exclusions	Anatomical Therapeutic Class & System Organ Class
Labeledness		Temporal Relationship
Reporter Causality		Temporal Compatible

3653 In parallel, a separate decision support tool (CASCADE) was developed and validated through
3654 consultation with experienced drug safety physicians. A decision tree structure was adopted due to
3655 its increased transparency and interpretability when compared to other causality assessment
3656 algorithms. This increased transparency and interpretability allow a clear statement of the rationale
3657 for the assessment to be written (e.g. “The case is deemed causally related as it is (a) Labelled for the
3658 event (b) The event has a plausible temporal relationship, etc.”).

3659 The work on the decision tree provided a basis for the subsequent predictive model, informing
3660 contributing factors and the topology of the resulting Bayesian Network model. The authors rationale
3661 for selecting this type of model include: the ability to combine multiple sources of information with
3662 expert knowledge, transparency and interpretability, and their capability to model complex
3663 frameworks with causal dependencies where a lot of uncertainty exists. Model training utilized an
3664 annotated dataset of 50k cases, with a separate test dataset of 20k cases. Both the training and test
3665 dataset represented a broad range of drug classes and event categories. All cases had been
3666 previously assessed by medical experts and were taken from a period where the causality
3667 assessment practices were consistent.

3668 *E3. Results*

3669 The model demonstrated high performance (sensitivity was 0.900, with Positive Predictive Value of
3670 0.778) in predicting the causality assessment of drug–event pairs compared with clinical judgment
3671 using global introspection. The authors also explored a learned topology Bayesian Network model

3672 with the same training data. The learned topology model was found to have inferior performance
 3673 compared to their CASCADE-based model.

3674 *E4. Compliance with the governance framework*

3675 **Table 15: Use case E: Alignment with the governance framework (detail)**

3676 Source: CIOMS Working Group XIV

3677

Principle	Activities	SPEC	DEV	PreD	PstD	RU
Risk-based Approach	A risk-based approach was used to limit the scope of the model to solicited, post-marketing cases as automating the causality of assessment of these cases was determined to have a lower risk/ impact to the PV system.	A	A	N/A	N/A	N/A
Human Oversight	Drug safety physicians and SMEs were involved in the data review and model development activities, ensuring the applicability of the model to its intended purpose. As this was a POC/study, there was no discussion about the creation of a quality management framework to support human oversight for future production use.	A	A	N/A	N/A	N/A
Validity & Robustness	The use case and deployment domain are described in the paper. The data selection, model training and testing activities used in model development are discussed in detail, as is the approach used for performance assessment. The authors consider areas for investigation that may be used to further demonstrate the model's validity and improve robustness.	A	A	N/A	N/A	N/A
Transparency	The paper provides information about intended use of the model and its design, including the decision tree (CASCADE). Data, results, areas for further investigation, and how the model could be applied in a PV system are discussed.	A	A	N/A	N/A	N/A
Data Privacy	In alignment with the principles in this book, the data used for the development, training, and validation of the model is from the company's internal post-marketing safety database suggesting it was obtained with the patient's/reporter's consent and in compliance with relevant privacy laws and regulations.	A	A	N/A	N/A	N/A
Fairness & Equity	Based on the article, it is not possible to comment on whether model development aligns with this guiding principle	N/A	N/A	N/A	N/A	N/A
Governance & Accountability	There is no discussion of governance and accountability activities, as defined in this guidance, in the paper. The authors do acknowledge the need for models to remain compliant with regulatory frameworks and guidelines. Further, the CASCADE decision tree created is referenced as a causality assessment support tool implying accountability for the final causality assessment decision remains with the drug safety SME.	N/A	N/A	N/A	N/A	N/A

3678 **Abbreviations**

3679 SPEC: Collection of specifications, requirements

3680 DEV: Development and change management

3681 PreD: Pre-deployment & post-change sign-off
3682 PstD: Post-deployment & post-change hyper-care
3683 RU: Routine Use
3684 A: Applicable
3685 NA: Not Applicable

References

²⁷⁸ Cherkas, Y., Ide, J. & van Stekelenborg, J. *Leveraging Machine Learning to Facilitate Individual Case Causality Assessment of Adverse Drug Reactions*. *Drug Saf* 2022;45;571-582. ([Journal abstract](https://doi.org/10.1007/s40264-022-01163-6)) <https://doi.org/10.1007/s40264-022-01163-6>

²⁷⁹ Modified from Cherkas, Y., Ide, J. & van Stekelenborg, J. *Leveraging Machine Learning to Facilitate Individual Case Causality Assessment of Adverse Drug Reactions*. *Drug Saf* 45, 571–582 (2022). <https://doi.org/10.1007/s40264-022-01163-6>

3686

3687

3688 Use Case F: Process efficiencies supporting signal detection

3689 **Source Article:**²⁸⁰

3690 **Area of PV:** Signal Detection

3691 *F1. Business rational and challenges*

3692 One of the most time- and resource-demanding procedures for dismissing safety signals is the
3693 identification of alternative causes for the reported adverse events (AEs) in ICSRs after signals of
3694 disproportionate reporting have been identified. This includes the screening of co-reported drugs to
3695 identify alternative potential causes for the newly identified drug–event pair.

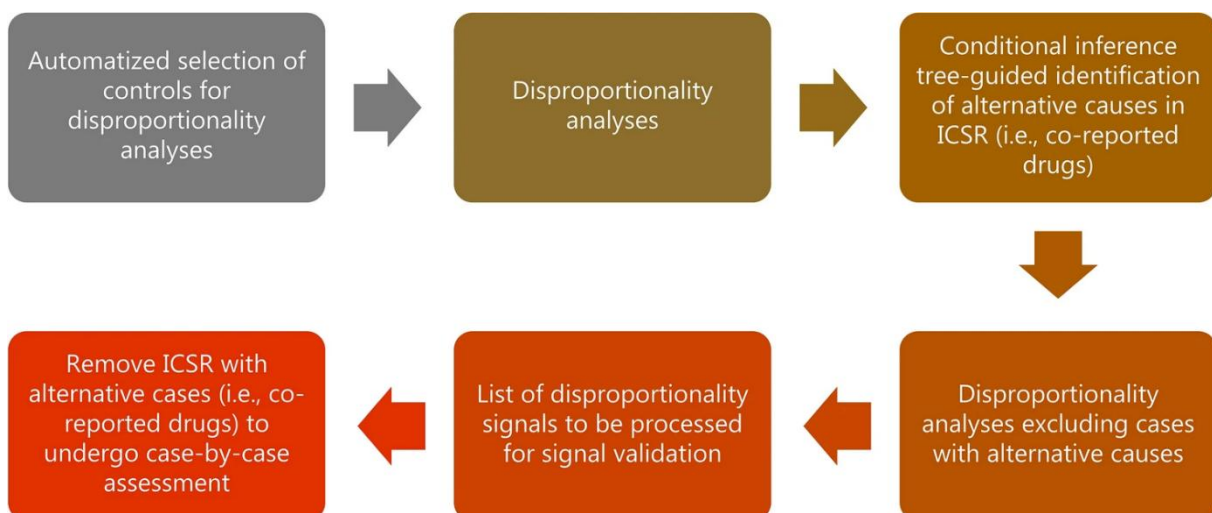
3696 *F2. Solution*

3697 This study aimed to develop an AI-based framework to automate (1) the selection of control groups
3698 in disproportionality analyses and (2) the identification of co-reported drugs serving as alternative
3699 causes, to look to dismiss false-positive disproportionality signals.

3700 The implementation of automatic selection of controls and dismissal of false positive signals using a
3701 conditional inference tree is summarized in the flowchart below.

3702 **Figure 8: Flowchart summarizing the implementation of the automatic selection of controls and the
3703 dismissal of false positive signals when using a conditional inference tree**

3704 Source:²⁸¹



3705
3706

3707 A dual approach combining the Anatomical Therapeutic Chemical (ATC) classification system code
3708 and the approved therapeutic indication in the US Summary of Product Characteristics (SmPC) of
3709 galcanezumab was used for automatizing the selection of controls for disproportionality analysis
3710 when using FAERs. All active ingredients with the same therapeutic target (i.e. CGRP antagonists) as
3711 galcanezumab were identified using the 4th level of the ATC code, or rather the chemical subgroup.
3712 DrugBank was used to identify controls with the same approved therapeutic indication but with
3713 active ingredient outside the chemical subgroup of galcanezumab, aiming to avoid masking due to
3714 drug class effect and confounding by indication.

3715 Disproportionality signals were further analyzed by using conditional inference trees to identify
3716 alternative cause co-reported drugs. The SmPC of disproportionately co-reported drugs was screened
3717 to identify those drugs that listed in the SmPC the AE in disproportionality signal mimicking
3718 procedures performed during signal validation. The disproportionality analysis was conducted again
3719 by removing cases with co-reported drugs for which the AE under investigation was listed in the
3720 SmPC as these cases had alternative causes for the AE.

3721 *F3. Results*

3722 AI could significantly ease some of the most time-consuming and labour-intensive steps of signal
 3723 detection and validation. The AI-based approach showed promising results, however, future work is
 3724 needed to validate the framework.

3725 *F4. Compliance with the governance framework*3726 **Table 16: Use case F: Alignment with the governance framework (detail)**

3727 Source: CIOMS Working Group XIV

3728

Principle	Activities	SPEC	DEV	PreD	PstD	RU
Risk-based approach	When developing tools to support signal detection developers should consider the overall impact the tool may have to the PV system within a QMS and consider the need for mitigations that maybe required to support broader deployment.	A	A	N/A	N/A	N/A
Human oversight	Human oversight within the proof-of-concept setting would ensure that the outputs of the conditional inference tree are robust for use in augmenting the quantitative signal detection processes. In a production setting it would include monitoring mechanisms and sampling approaches to ensure ongoing validity.	A	A	N/A	N/A	N/A
Validity & Robustness	The outputs of the conditional inference trees can be subject to existing validation procedures to compare inputs to known outputs.	A	A	N/A	N/A	N/A
Transparency	Transparency would involve clear communication and documentation of the AI system's use in the process of signal detection and validation from ICSRs. Within this proof of concept this includes detailed information about the business case, methodology and training datasets, should the tool move into a production setting further documentation and information relating to how the tool is integrated into the wider PV system and it's QMS would be required.	A	A	N/A	N/A	N/A
Data Privacy	The method uses publicly available information on labelling	N/A	N/A	N/A	N/A	N/A
Fairness & Equity	The method is not using individual patient data.	N/A	N/A	N/A	N/A	N/A
Governance & Accountability	During the proof of concept (PoC), the accountability of the methodology remains with the developer. However, if the methodology is integrated into a production setting, accountability would transition to the human subject matter expert. Governance within a PoC ensures that scientific integrity principles are adhered to, while future product use governance should cover how the tool fits into the overall PV system and quality management system (QMS).	A	A	N/A	N/A	N/A

3729 **Abbreviations**

3730 SPEC: Collection of specifications, requirements

3731 DEV: Development and change management

3732 PreD: Pre-deployment & post-change sign-off

3733 PstD: Post-deployment & post-change hyper-care

3734 RU: Routine Use

3735 A: Applicable

3736

NA: Not Applicable

References

²⁸⁰ Al-Azzawi F, Mahmoud I, Haguinet F, Bate A, Sessa M. Developing an Artificial Intelligence-Guided Signal Detection in the Food and Drug Administration Adverse Event Reporting System (FAERS): A Proof-of-Concept Study Using Galcanezumab and Simulated Data. *Drug Saf.* 2023;Aug;46(8):743-751 ([Journal full text](#)) <https://doi.org/10.1007/s40264-023-01317-0>

²⁸¹ Al-Azzawi F, Mahmoud I, Haguinet F, Bate A, Sessa M. Developing an Artificial Intelligence-Guided Signal Detection in the Food and Drug Administration Adverse Event Reporting System (FAERS): A Proof-of-Concept Study Using Galcanezumab and Simulated Data. *Drug Saf.* 2023;Aug;46(8):743-751 ([Journal full text](#)) <https://doi.org/10.1007/s40264-023-01317-0>

3737 **Use Case G: Generative artificial intelligence for enhanced and intelligently**
 3738 **structured outputs from large pharmacovigilance document libraries**

3739 **Source:** Internal to CIOMS Working Group member organization

3740 **Area of PV:** PV document retrieval

3741 *G1. Business rational and challenges*

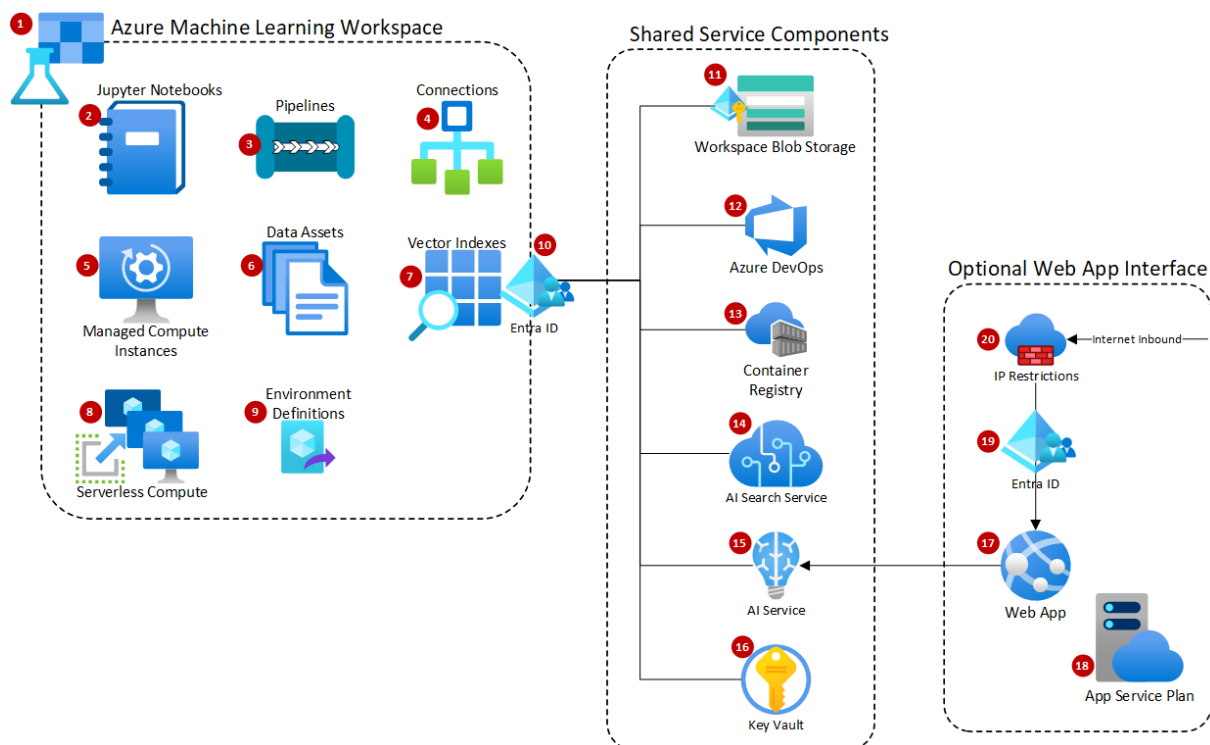
3742 Use of AI represents a burgeoning field that has gained significant traction across various industries
 3743 including the pharmaceutical industry. However, due to its relative novelty, concrete business use
 3744 cases remain somewhat scarce, leaving many organizations searching for innovative ways to leverage
 3745 AI technology effectively. LLMs, such as the Generative Pre-trained Transformer (GPT) models, offer
 3746 a potentially valuable resource for businesses seeking historical references and streamlined
 3747 document management solutions, with capabilities to summarize & synthesize large sources of data.

3748 *G2. Solution*

3749 Harnesses the power of LLMs to optimize our document filing structure. With hundreds of thousands
 3750 of documents pertaining to patient safety and PV activities, in this use case an AI tool was employed
 3751 to function as a powerful system for searching, extracting and generating information in an effective
 3752 and structured manner specific to the parameters and requirements of the (human) user. By
 3753 automating the process of retrieving relevant information, subject matter experts (SMEs) can
 3754 redirect their time towards value-added endeavours rather than manual data sifting.

3755 **Figure 9: An outline of our initial artificial intelligence architecture**

3756 Source: Internal to CIOMS Working Group member organization



3757
 3758 The implementation of LLMs not only aids day-to-day tasks but can also enhance efficiency in many
 3759 different vigilance activities across all different verticals of vigilance. Currently this tool has the ability
 3760 to help support audit and inspection requests, performative AI researcher for PV projects that
 3761 require detailed and structured retrieval from safety document libraries, e.g. searching for regulatory
 3762 / safety correspondence.

3763 Development and deployment of a meticulously crafted vector index of relevant company file
 3764 structures, allows us to leverage an LLM to easily and efficiently navigate documents, files and data
 3765 available within those files. By transitioning to newer models as they release, the organization has
 3766 observed improved accuracy and utility in responses, validated through a manual verification of
 3767 processes.

3768 Regular feedback sessions are conducted to refine the AI tool's performance and uncover additional
 3769 use cases. This iterative approach ensures continual improvement and minimizes errors. Guidelines
 3770 were also developed for framing questions to the model effectively, further enhancing the tool's
 3771 usability.

3772 Looking ahead, there are plans to expand this application of GenAI, focusing on areas such as
 3773 summarization of safety literature for example signal detection purposes and summary outputs for
 3774 case reporting. By training additional models on relevant datasets, including Individual Case Safety
 3775 Reports (ICSR), the organization aims to automate safety summary generation and aggregate
 3776 analysis, trending, any other ad hoc safety questions that can be answered from the data from or
 3777 within the global safety database (GSD) thereby streamlining decision-making processes.

3778 *G3. Results*

3779 The generative AI tool in this use case has demonstrated potential as an AI 'research assistant'
 3780 enabling PV workers to quickly and efficiently search hundreds and thousands of safety documents
 3781 to provide structured and intelligent outputs. The adoption of AI technology like LLMs could
 3782 empower organizations using this technology to optimize operational efficiency and prioritize tasks
 3783 that impact product benefit-risk assessments. By leveraging AI-driven document management
 3784 solutions, subject matter experts can devote more time to critical medical and scientific evaluations,
 3785 ultimately enhancing product safety and efficacy. In addition, the planned expansion of use of LLMs
 3786 in GenAI capabilities to evaluate transactional data and provide scientific analysis and summary has
 3787 the potential to be a game changer in evaluating product or patient level safety risk in real time,
 3788 simply by "asking a question" with well referenced and documented summaries/conclusions.

3789 *G4. Compliance with the governance framework*

3790 During the development and pre-deployment phases, the GenAI project is carefully developed and
 3791 managed with a limited scope, ensuring full alignment with most guiding principles. As the project
 3792 transitions into production and its exposure and scope expand, careful consideration will be given to
 3793 maintaining close alignment with these principles. Therefore, compliance is indicated as closely
 3794 aligned, laying a foundation of trust in the solution's ability to perform vigilance tasks with adaptive
 3795 and growth capabilities.

3796 **Table 17: Use case G: Alignment with the governance framework (detail)**

3797 Source: CIOMS Working Group XIV

Principle	Activities	SPEC	DEV	PreD	PstD	RU
Risk-based approach	GenAI use is a closed environment used for training and testing during development and pre-development. However, there is communication of potential inaccuracies and pitfalls during these phases. Currently there is no anticipated (patient risk) for post-deployment or routine use. As GenAI achieves more general and expanded risks will be regularly reassessed. Therefore, all phases are considered partially aligned with this guiding principle.	A	A	A	A	A
Human oversight	Fully aligned in development phase. Although there is still human oversight from the user,	A	A	A	A	A

	there is probably partial alignment during post deployment and routine use as GenAI would not be routinely scrutinized by human oversight and management					
Validity & Robustness	Validation and testing is extensive during development and pre-deployment and therefore in full alignment with the guiding principle (based on smaller sample set during these stages. Once in post deployment and routine use the data sets are very large and it is not possible to 100% fully test and validate all use cases (therefore in partial alignment). However, if users are noting inaccuracies in information retrieval, feedback will be provided to the human developers of the GenAI to refine and update the system.	A	A	A	A	A
Transparency	In full alignment during development and pre-development. As GenAI system expands in scope and complexity during post-deployment and routine use, further realignment is anticipated. Transparency in relation to the public is not applicable as this is a closed system	A	A	A	A	A
Data Privacy	Fully alignment during all phases. All of data remains internal; therefore, no transfer of information and external vendors/individuals.	A	A	A	A	A
Fairness & Equity	Fully aligned with this guiding principle at all phases. No relevant groups or individuals are being excluded or disadvantaged.	A	A	A	A	A
Governance & Accountability	Accountability from system usage and implementation during development and pre-deployment, e.g. if system is clearly not useful then it will be discontinued / upgraded. Ultimately regulatory accountability resides with subject matter expert / user as they are responsible to review and verify content. Therefore, partial alignment is anticipated from post-deployment onwards	A	A	A	A	A

3799
3800
3801
3802
3803
3804
3805
3806

Abbreviations

SPEC: Collection of specifications, requirements
DEV: Development and change management
PreD: Pre-deployment & post-change sign-off
PstD: Post-deployment & post-change hyper-care
RU: Routine Use
A: Applicable
NA: Not Applicable

3807

3808 Use Case H: Artificial intelligence to support diagnosis and prediction of 3809 (hydroxy)chloroquine retinopathy

3810 **Source:**^{282,283}

3811 **Area of PV:** Pharmacovigilance in The Clinic

3812 *H1. Business rational and challenges*

3813 PV in the clinic is concerned with the prevention and treatment of adverse drug reactions in
3814 individuals. Prevention may be primary, which can be achieved through identifying potential complex
3815 or non-obvious combinations of patient characteristics that are predictive of adverse drug reactions
3816 to guide optimum medication selection. (i.e. precision medicine) It also encompasses secondary and
3817 tertiary prevention (i.e. early diagnosis of adverse drug reactions, and ensuing interventions) to
3818 mitigate the impacts of ADRs. Examples follow.

3819 Chloroquine and hydroxychloroquine are important drugs in rheumatology. Although relatively well
3820 tolerated compared to some other therapeutic options, retinal toxicity, is a risk which can result in
3821 serious visual impairment if not detected early so that the drug may be discontinued in a timely
3822 manner. Even so, by the time of retinopathy diagnosis there may be irreversible retinal damage.
3823 Conversely, if predictive AI can provide sufficient leading indicators or progression, therapy duration
3824 and attendant therapeutic benefits might be maximized. Historically the gold standard for screening
3825 and detection has been fundus photography and automated perimetry. More recently, multifocal
3826 electroretinography (mfERG) and optical coherence tomography (OCT) have been added to the
3827 diagnostic armamentarium. Each of these are routinely assessed by human readers, ideally retinal
3828 specialists, but subtle changes, including temporal patterns, can be missed, and not all locales have
3829 the necessary instrumentation or available retinal specialists. It would be ideal to augment human
3830 visual assessors to identify early functional changes indicative of retinopathy prior to onset of
3831 irreversibility or better predict progression. AI has shown potential in detecting or predicting various
3832 ocular diseases based on retinal images/fundus photography, such AMD, DR. More AI has been
3833 retrospectively developed and tested to diagnose or predict (hydroxy) chloroquine retinopathy.

3834 *H2. Solution*

3835 AI has been applied to colour fundus photographs, OCT and multifocal electroretinographic tracings
3836 for diagnosing hydroxychloroquine retinopathy, Fan et al studied hyperspectral imaging (HIS) of 176
3837 fundus photographs from retinopathy positive (25) versus retinopathy negative (66) patients at a
3838 referral clinic using four deep learning models for the detection of retinopathy Kulyabin et al
3839 compared deep learning-based classification of raw mfERGs versus models based on conventional
3840 readout parameters of the mfERG for classification, and for prediction (regression) of visual field
3841 sensitivities from 53 predominantly female patients (35 retinopathy negative, 9 minimal retinopathy,
3842 and 9 manifest retinopathy) monitored with mfERGs and perimetry for period of 0.7-20.9 years. Kalra
3843 et al used random forests for automated diagnosis and prediction of disease progression using
3844 clinical features and features based on spectral domain OCT (SD-OCT) obtained from 388 eyes / 368
3845 patients, a majority being female. Habib et al trained support vector machines (SVM) to identify
3846 hydroxychloroquine retinopathy in 1463 eligible eyes (748 predominantly female patients), of which
3847 95 eyes (48 patients) were eligible for inclusion as controls.

3848 *H3. Results*

3849 The best performing deep learning models in the study of Fan et al achieved accuracy, precision,
3850 recall, specificity, and F1-scores of ≥ 0.95 , with superior performance using hyperspectral images
3851 versus the original retinal images. Habib et al's SVM returned a specificity of 84.0% with sensitivity of
3852 90.9%. Performance could be calibrated to place a premium on sensitivity for screening or specificity
3853 for diagnosis. Kalra reported a mean AUC of 0.97, a sensitivity 95% and specificity of 91% for
3854 detection, and mean AUC=0.89, recall of 90% and specificity of 80% for progression prediction.

3855 Kulyabin reported that AI-based models using full mfERG traces had a balanced accuracy of up to
3856 0.795, precision of up to 0.844, recall of up to 0.866, and F1-score of up to 0.771.

3857 *H4. Compliance with the governance framework*

3858 In considering the alignment of the reviewed studies with the governance framework we note
3859 several points up front. The studies were retrospective and feasibility/pilot studies, without reported
3860 advancement to routine use in the clinic. Because the drugs under study are for autoimmune
3861 disorders, in addition the study populations often being small, the subjects were predominantly
3862 female, consistent with the treatment indication. The use of eyes as the unit of observation raises
3863 the question of pseudo-replication and its potential impacts of performance estimates, though
3864 confidence intervals were not typically presented. At least one study noted under-representation of
3865 Asian in the study sample and the need for further assessment in larger and more diverse
3866 populations. Finally as is often the case in AI applications involving retinal pathology, other ocular
3867 pathologies were excluded, which limits generalizability to more diverse patient populations that
3868 have multiple retinal pathologies (e.g. diabetic retinopathy and drug-induced retinopathy).

3869 **Table 18: Use case H: Alignment with the governance framework (detail)**

3870 Source: CIOMS Working Group XIV

Principle	Activities	SPEC	DEV	PreD	PstD	RU
Risk-based approach	Not aligned. Risk assessment and risk mitigation plans not provided in these pilot studies. But placement within a human-in-the-loop framework was explicitly considered in one or more studies	N/A	N/A	N/A	N/A	N/A
Human oversight	Partial alignment. One or more of the publications, which report feasibility/pilot studies in clinical settings, discuss the proper deployment with respect to human oversight, such as HITL. However, change management and staff training plans not discussed. Discussed is the fact that the available human oversight in some locations may be provided by generalists with less experience and expertise than retinal specialists, affording more opportunity for incremental benefits in underserved settings.	A	A	N/A	N/A	N/A
Validity & Robustness	Partial alignment. Reference standards defined. One or more studies note the limitation of the imbalanced data sets used that impair generalizability. Also, in one/more studies patients with other ocular pathology excluded so the two classes were HCQ retinopathy present versus normal retina, which limits generalizability to screening in patients with other coexistent ocular disorders that may affect the retina. Source population (deployment domain) not clearly defined in all studies No discussion of integrating data pre-processing (e.g. cropping retinal images) into routine use). In some studies unit of observation was "eyes" raising questions about pseudo-replication.	A	A	A	N/A	N/A
Transparency	One/more papers report adherence to tenets of the Declaration of Helsinki and obtained Institutional Review Board approval. One or more papers described explanations of results such as heatmaps of feature distributions.	A	A	A	N/A	N/A

Data Privacy	One or more of the referenced studies declared adherence to tenets of the Declaration of Helsinki and obtained Institutional Review Board Approval.	A	A	A	N/A	N/A
Fairness & Equity	One/more of the referenced studies report adherence to tenets of the Declaration of Helsinki and obtaining Institutional Review Board approval. One or more of papers acknowledge that data under-represents specific groups of persons such as Asians, who may display different findings and recommends further assessment with larger data sets with more diverse representations. Further discussion involved scenarios in which retinal specialists may not be available, such as under-resourced or under-represented locales, as also discussed in human oversight above.	A	A	A	N/A	N/A
Governance & Accountability	These studies which occurred in clinical settings were conducted according to the guidelines of the Declaration of Helsinki and approved by the respective Institutional Review Board o	A	N/A	N/A	NA	N/A

3872
3873
3874
3875
3876
3877
3878
3879

Abbreviations

SPEC: Collection of specifications, requirements
DEV: Development and change management
PreD: Pre-deployment & post-change sign-off
PstD: Post-deployment & post-change hyper-care
RU: Routine Use
A: Applicable
NA: Not Applicable

References

- ²⁸² Kulyabin M, Kremers J, Holbach V, Maier A, Huchzermeyer C. Artificial intelligence for detection of retinal toxicity in chloroquine and hydroxychloroquine therapy using multifocal electroretinogram waveforms. *Scientific Reports*. 2024;Oct22;14(1):24853. ([Journal full text](https://doi.org/10.1038/s41598-024-76943-4)) <https://doi.org/10.1038/s41598-024-76943-4>
- ²⁸³ Kalra G, Talcott KE, Kaiser S, Ugwuegbu O, Hu M, Srivastava SK, Ehlers JP. Machine learning–based automated detection of hydroxychloroquine toxicity and prediction of future toxicity using higher-order OCT biomarkers. *Ophthalmology Retina*. 2022;Dec1;6(12):1241-1252. ([Journal abstract](https://doi.org/10.1016/j.oret.2022.05.031)) <https://doi.org/10.1016/j.oret.2022.05.031>

3880

3881 **APPENDIX 4: Content related to explainability and to** 3882 **Fairness & Equity**

3883

3884 **Content related to explainability**

3885 *Illustrative examples*

3886 As stated above, diverse stakeholders in PV may require and benefit from explainability. However,
3887 the explainability that is required and how it is used, differs depending on the setting, i.e. who is
3888 asking and for what, in which task or process and in which system lifecycle stage.²⁸⁴ In the following
3889 sections, examples are provided below to illustrate the different scenarios. Finally, the benefits of
3890 explainability described in the example are summarised at the end of this section.

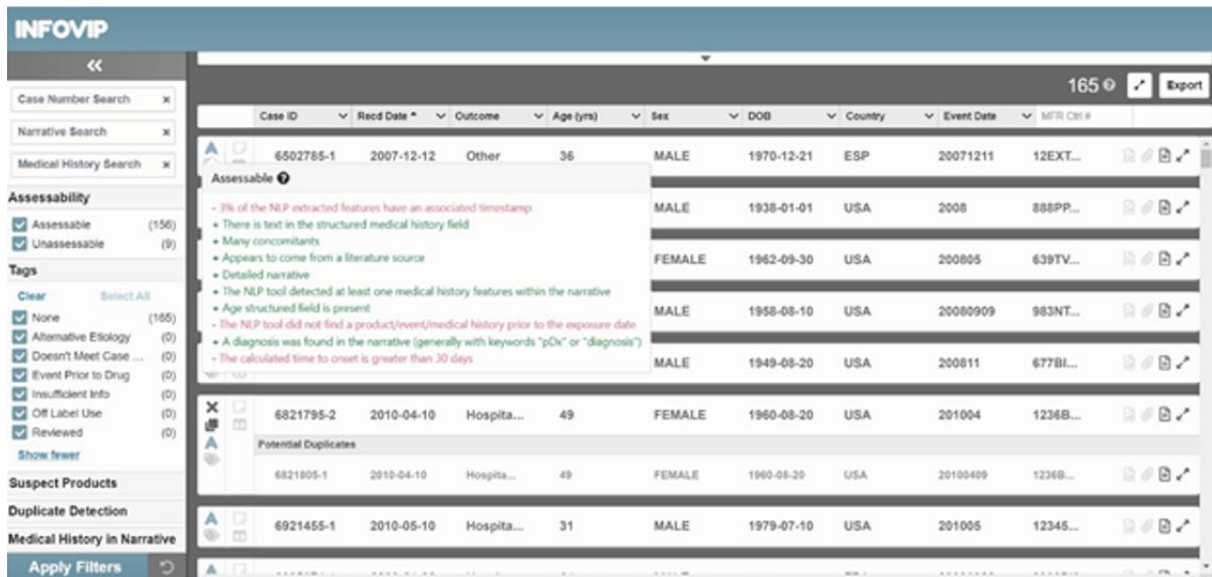
3891 *Examples of explainability in artificial intelligence-supported pharmacovigilance tasks*

3892 Take for example a setting in which a PV officer is reviewing a case that has been selected by the AI
3893 as a potential signal but does not immediately see why the case was flagged. In such cases, the
3894 reviewer could benefit from being able to see which text in the case data led the AI to this
3895 conclusion. An actual example of this is described below.

3896 The Information Visualization Platform (InfoViP) developed for the FDA's Center for Drug Evaluation
3897 and Research is an example of explainability supporting the human expert engaged in signal
3898 detection and assessment.²⁸⁵ InfoViP uses Natural Language Processing (NLP) and several other
3899 components to process post-market data and visualize information, i.e. explanations, to support
3900 medical reviewers who detect and evaluate potential signals from the millions of adverse event
3901 reports submitted to the FDA's FAERS database. The NLP component, the Event-based Text-mining of
3902 Health Electronic Records (ETHER), coupled with modern frontend techniques, provide visual
3903 information by colour-coded highlighting of relevant text in the case narrative to help reviewers
3904 focus on signal-related information. An informed model further identifies cases containing enough
3905 information to assist reviewers assess the report quality and provides concrete explanations of these
3906 selections. All these functionalities combined with case deduplication and several filtering options
3907 facilitate speedy review by the medical reviewers, an otherwise humanly impossible task across
3908 millions of reports.²⁸⁶

3909

3910 **Figure 10: FDA's Information Visualization Platform user interface illustrates the system**
 3911 **capabilities, focusing on the features that positively contribute to classification for assessability.**
 3912 **Source:**²⁸⁷
 3913



3914
 3915 The example above illustrates a core benefit of explainability described by Albahri et al (2023).²⁸⁸
 3916 Explainability can facilitate human experts in making 'sound and reliable' decisions. And ultimately,
 3917 when the human decision and the accompanying explanation are retained, this information would
 3918 nurture trust of the system owner and quality assurance staff who are tasked with ensuring
 3919 compliance as well as the trust of regulators who may wish to inspect why certain cases are selected
 3920 or rejected as signals.

3921 Also, it is conceivable that explanations could lead a user to notice a bias or spurious correlation that
 3922 is leading to incorrect predictions. Reporting this back to the development team can contribute
 3923 towards future improvement. In this way, explainability is useful for ongoing vigilance against bias
 3924 risk and performance issues that may appear post-deployment and for continually ensuring the
 3925 trustworthiness of the decisions made. As a result, XAI explanations have resulted in increased trust
 3926 and the perception of fairness in AI-supported decision making.²⁸⁹

3927 *Examples of pharmacovigilance stakeholders benefitting from explainability*

3928 While the likelihood of an individual from the general public requiring explainability in a PV setting is
 3929 considered to be small, the possibility of this happening cannot be excluded entirely as the use of AI
 3930 becomes more commonplace. Some conceivable scenarios are described below:

- 3931 • When a reporter (healthcare professional or patient) directly reports a serious adverse event
 3932 and the report is subsequently processed as a non-serious case by an AI triage system, upon
 3933 being informed of this, the reporter undergoes a negative experience and may request an
 3934 explanation from the MAH. When this scenario takes place in a non-AI setting or an AI-
 3935 assisted triage setting, the reporter could receive an explanation from the PV officer who has
 3936 made the final triage decision. However, when this takes place in a fully-automated AI triage
 3937 process, unless some form of explainability is provided, the lack of explainability may impact
 3938 trust and acceptance of the result by the reporter.
- 3939 • When a consumer or a healthcare professional interacts with a medical information chatbot
 3940 used by the MAH to provide drug product information and collect safety data and quality
 3941 complaints, the individual may question and challenge any information that doesn't make
 3942 sense.
- 3943 • GenAI could be used in assisting a pharmacist in medication therapy management to prevent
 3944 drug interactions. Both the pharmacist and the patient are directly exposed to the AI's

3945 recommendations in this case.²⁹⁰ Here, questions concerning the AI recommendations could
3946 be raised by both parties.

3947 *Examples of explainability in system development*

3948 A developer who is training the AI system benefits from explainability when it reveals which features
3949 are used by the AI to reach a specific prediction or when it reveals a bias in the training data.
3950 Especially in complex systems which lack inherent interpretability and explainability, an XAI tool
3951 could provide explanations that facilitate troubleshooting by revealing what to change or exclude in
3952 order to ‘flip’ the outcome.²⁹¹ Using the XAI explanation, one could discover that needle in the deep
3953 neural network haystack. However, in most cases, tweaking the system architecture of a deep neural
3954 network or specific features will be quite challenging even when XAI points the way. The XAI output
3955 is more likely to identify hidden biases in the training data which can be corrected as illustrated in
3956 the example below.

3957 An example presented by Ribeiro et al (2016)²⁹² demonstrates how the XAI explanation provided by
3958 using Local Interpretable Model-Agnostic Explanations LIME could be understood by humans and
3959 reveal the likely cause of incorrect predictions. In this experiment, the model that was trained to
3960 distinguish images of dogs and wolves was first intentionally trained to associate wolves with
3961 snowscapes. In other words, the training data was deliberately biased by excluding images of wolves
3962 in other seasons. This resulted in predictions that included a wolf against a green background
3963 identified as a dog and a husky in a snowscape identified as a wolf. LIME was used to show subjects
3964 which areas of the image were used as features by the AI in its predictions to see if the subjects could
3965 identify the cause of the misidentification. The subjects successfully identified background snow as
3966 the potential feature that led the AI to make the incorrect predictions. Thus, XAI can be used to
3967 explain a prediction made by an inscrutable deep neural network and uncover the underlying issue in
3968 the training data and the resulting spurious correlation that led to the incorrect output.

3969 In the context of PV, similar techniques could be used to highlight words in the text which are picked
3970 up by the AI as relevant features. In a real-life but unpublished example in which an AI triage system
3971 was misidentifying some serious cases, developers benefited from seeing which terms in the case
3972 were considered by the AI in its seriousness predictions. In this case, the XAI explanation revealed a
3973 focus on the drug name. Combined with the fact that the missed serious cases concerned Over The
3974 Counter (OTC) drugs, the developers discovered that the AI was basing decisions on the drug name
3975 and a learned spurious assumption that OTC drugs are not likely to cause serious events. Using the
3976 insight gained from explainability, the developers could reject the model in favour of another one,
3977 examine the training data for bias such as the lack of serious OTC cases or when there is no bias,
3978 solve the issue through feature engineering by instructing the AI not to consider the drug name in its
3979 decisions.

3980 Explainability, therefore, can help developers make informed decisions when assessing AI models by
3981 uncovering hidden biases as well as features and spurious correlations that are resulting in incorrect
3982 predictions. Explainability may also reveal the underlying factors that result in performance
3983 differences between models that are trained on the same training data and aid the developer in
3984 model selection. In turn, transparent documentation of this process will go a long way towards
3985 nurturing trust in the system, not only for the developers but also for the system owners, users and
3986 the regulators.

3987 *Examples of explainability in artificial intelligence-systems interacting with health care professionals
3988 and patients*

3989 Another hypothetical example can be the case of a healthcare professional (HCP) who is requesting
3990 product-specific information via a chatbot provided by a marketing authorisation holder (MAH). Such
3991 a chatbot could have multiple objectives ranging from the provision of drug product information to
3992 the collection of adverse event and quality defect reports. When the HCP notices that the chatbot
3993 response is inadequate, i.e. not considering key medical terms or adverse events, or providing

3994 questionable information, the HCP may contact the MAH for an explanation. Explainability, if it is
3995 available, may help the MAH troubleshoot the system and/or provide information back to the HCP.

3996 While the scenario above is a fictive example, one example of a chatbot that is currently available is
3997 the Smart AI Resource Assistant for Health (SARAH) on the World Health Organization website. This is
3998 a prototype chatbot that is intended to provide tips on health topics and not medical advice as
3999 clearly stated on the landing page of SARAH.²⁹³ On one hand, SARAH exemplifies how such an
4000 application could be of service to the public as it is available 24/7 and in eight languages. On the
4001 other hand, incidents of the chatbot providing inaccurate or incorrect information or being unable to
4002 answer some queries have unfortunately been reported in the media and taken up in the OECD AI
4003 incidents monitoring database.²⁹⁴ This illustrates how, when a chatbot is deployed, the interacting
4004 patient or healthcare provider or the media may challenge the information that is provided. It is
4005 therefore conceivable that in PV, when a MAH deploys an AI system that interacts directly with the
4006 public, explainability for the public will also be required.

4007 Finally, any system that interacts directly with the public in a medical setting warrants extra attention
4008 in that a HCP is likely to notice medically incorrect information, but most consumers and patients
4009 may not be able to do this. Individuals without a medical background will be at risk of accepting and
4010 acting on medically incorrect information. To illustrate this point, in a study of trust and medical
4011 advice provided by ChatGPT, persons without a medical background have been found to trust the
4012 chatbots for lower-risk health topics.²⁹⁵ Without the medical background, a layperson is at increased
4013 risk of harm by not being able to recognize incorrect information. Thus, aside from being able to
4014 provide an explanation to an individual from the public who is challenging the AI output, system
4015 owners must thoroughly consider and mitigate the risks of an AI system that interfaces with the
4016 general public. This also touches on the subject of accountability since it is not the chatbot that is
4017 held accountable for any harm that befalls the individual.

4018 *Example where explainability is not available and not required*

4019 To illustrate a situation in which explainability is not necessary nor possible, consider first how the
4020 use of publicly available machine translation tools is now commonplace and how the public generally
4021 do not require explanations into how the AI translated the text. Consider also how translations in a
4022 GxP-regulated environment require a quality check regardless of whether the translation was carried
4023 out by a human or a machine. Furthermore, the quality check is normally carried out by an individual
4024 who is proficient in both the source and destination languages. Therefore, not only will the human be
4025 able to spot errors, but also in some cases, understand the underlying reasons for machine
4026 translation errors even without any explanation. An example of a translation issue with a self-evident
4027 root cause is the case of biased gender assignment that occurs when translating a genderless
4028 language such as Finnish to English.²⁹⁶ The bilingual human reviewing the translations would easily
4029 notice the gender bias, understand why this has been introduced by the AI and can correct this
4030 accordingly. See also chapters x on human oversight, data privacy and fairness and equity.

4031 *Example where explainability could be available but is not required*

4032 In PV, another example of AI use which would not require explainability would be automatic de-
4033 identification of case narratives presented by Meldau et al 2024.²⁹⁷ Automatic redaction of case
4034 narratives could ensure data privacy while allowing case narratives sharing which are crucial for
4035 providing more complete information to signal detection assessors. In this case, a system using a
4036 neural network combined with hand-engineered rules was trained to automatically detect names in
4037 case narratives for the purpose of redaction. Using a test set of over 5000 Yellow Card narratives
4038 from the MHRA and over 500 open source annotated and unannotated deidentified patient
4039 discharge summaries from i2b2.org, the system was able to identify 96% of names longer than three
4040 characters and 88% of all names for redaction. The false positive rate was 0.2%. It is not conceivable
4041 that stakeholders would require explanations during routine use of such a system. Even if
4042 stakeholders notice certain names are not being identified, it is conceivable that the reasons could be
4043 self-evident. Taking this system as an example, since it is trained to detect names that consist of

4044 three letters or more, it may miss two-letter names which are common in Asia. When such a lapse is
4045 noticed, the developers and stakeholders would not require an explanation in order to take the
4046 necessary corrective measures such as re-training using a new dataset or adding new rules.

4047 *Example of recommended explainability combined with downstream human oversight*

4048 Even when there is downstream human quality control and the final decision rests with the human,
4049 explainability remains useful for supporting human decision and nurturing trust in the machine and
4050 in the decision that is made.

4051 While the machine translation example above describes one end of the spectrum where
4052 explainability is not needed when risks are adequately covered human review, for AI in PV and other
4053 GxP environments, it is not conceivable to deploy a system that expects the human to make
4054 decisions based on AI predictions without any explanation. After all, the human needs to understand
4055 why the suggestion presented by the AI is trustworthy. Explainability is therefore essential for making
4056 a well-informed decision. The FDA's InfoVIP system (see above) demonstrates how a medical officer
4057 can make an informed decision supported by XAI visualisation of relevant features.

4058 Another example from the WHO Collaborating Centre for International Drug Monitoring illustrates
4059 how explanations can be used in reviewing AI case deduplication. In this case, a pair of Norwegian
4060 cases were identified as probable duplicates, while on first glance for the reviewers the duplication
4061 was not obvious. Due to missing outcomes, onset dates and ages that were close but not matching,
4062 and no matches between the registered adverse drug reaction terms, these cases would not have
4063 been identified as duplicates if they were not flagged by the AI. The AI explanation revealed that the
4064 match score was based on six different drug substances which were identical between the cases in
4065 addition to the fact that these six drug substances are commonly not co-reported. The researchers
4066 verified that the cases were indeed duplicates by contacting the Norwegian national centre. The
4067 cases concerned the same incident but were reported by two different physicians from the same
4068 hospital, thus accounting for the differences.²⁹⁸ This example illustrates how the explanation could be
4069 essential for the human in order to understand, further investigate the veracity of the prediction and
4070 finally reject or in this case accept the prediction.

4071 *Examples of xAI methods*

4072 Some xAI methods in use at the time of writing this report include:

- 4073 • Local Interpretable Machine-Agnostic Explainability LIME
4074 See the example of Ribeiro et al (2016)²⁹⁹ described earlier in this chapter.
- 4075 • Shapley Additive exPlanations SHAP

4076 An example of SHAP explainability in a supervised ML model used to support signal validation is
4077 presented by Imran et al (2024).³⁰⁰

- 4078 • Trust scores that indicate the model's uncertainty for the output.³⁰¹
- 4079 • Confidence scores are a metric that is usually available and can be used to flag output that is
4080 uncertain for human review.³⁰²
- 4081 • Visualisation through highlighting of text that was considered by the AI in its prediction and
4082 saliency maps using a heat map overlay to indicate areas of the input image that are relevant
4083 for the model's prediction.

4084 Although assessing and processing images is not a mainstream activity in PV, saliency maps are
4085 mentioned as another example of XAI to complete the view of the current landscape. See examples
4086 in Plass et al (2023).³⁰³

4087 In the case of PV where the data is predominantly text based, visual explanations are likely to take
4088 the form of highlighting relevant text within the case data. See also the FDA InfoViP described above
4089 as an example of explainability benefits.^{304,305}

4090
4091

- Counterfactuals explanations or examples that show what characteristics in the input data when changed, would result in different output.³⁰⁶

References

- ²⁸⁴ Royal Society. *Explainable AI*. [Royalsociety.org](https://royalsocietypublishing.org/journal/rsos). 2024. ([Website](#) accessed 11 August 2024)
- ²⁸⁵ Spiker J, Kreimeyer K, Dang O, Boxwell D, Chan V, Cheng C, Gish P, Lardieri A, Wu E, De S, Naidoo J. Information visualization platform for postmarket surveillance decision support. *Drug Safety*. 2020;Sep;43:905-915. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-020-00945-0>
- ²⁸⁶ Botsis T, Dang O, Kreimeyer K, Spiker J, De S, Ball R. A Decision-Support Platform Powered by AI and Humans-in-the-Loop Boosts Efficiency and Assures Quality in FDA's Pharmacovigilance. *International Society of Pharmacovigilance, 23rd Annual Meeting 2024, Montreal, Canada*. ([Webpage](#) accessed 21 March 2025)
- ²⁸⁷ Botsis T, Dang O, Kreimeyer K, Spiker J, De S, Ball R. A Decision-Support Platform Powered by AI and Humans-in-the-Loop Boosts Efficiency and Assures Quality in FDA's Pharmacovigilance. *International Society of Pharmacovigilance, 23rd Annual Meeting 2024, Montreal, Canada*.
- ²⁸⁸ Albahri AS, Duham AM, Fadhel MA, Alnoor A, Baqer NS, Alzubaidi L, Albahri OS, Alamoodi AH, Bai J, Salhi A, Santamaría J. A systematic review of trustworthy and explainable artificial intelligence in healthcare: Assessment of quality, bias risk, and data fusion. *Information Fusion*. 2023;Aug1;96:156-191. ([Journal full text](#)) <https://doi.org/10.1016/j.inffus.2023.03.008>.
- ²⁸⁹ Albahri AS, Duham AM, Fadhel MA, Alnoor A, Baqer NS, Alzubaidi L, Albahri OS, Alamoodi AH, Bai J, Salhi A, Santamaría J. A systematic review of trustworthy and explainable artificial intelligence in healthcare: Assessment of quality, bias risk, and data fusion. *Information Fusion*. 2023;Aug1;96:156-191. ([Journal full text](#)) <https://doi.org/10.1016/j.inffus.2023.03.008>.
- ²⁹⁰ Roosan D, Padua P, Khan R, Khan H, Verzosa C, Wu Y. Effectiveness of ChatGPT in clinical pharmacy and the Role of Artificial Intelligence in medication therapy management. *J Am Pharm* 2023;Dec 2;64(2);422-428. ([Journal full text](#)) <https://doi.org/10.1016/j.japh.2023.11.023>
- ²⁹¹ Hauben M. Artificial intelligence in pharmacovigilance: Do we need explainability?. *Pharmacoepidemiology and Drug Safety*. 2022;Dec;31(12):1311-1316. ([Journal full text](#)) <https://doi.org/10.1002/pds.5501>
- ²⁹² Ribeiro MT, Singh S, Guestrin C. "Why Should I Trust You?": Explaining the Predictions of Any Classifier [Internet]. [arXiv:1602.04938](https://arxiv.org/abs/1602.04938) 2016. ([Journal full text](#)) <https://doi.org/10.48550/arXiv.1602.04938>
- ²⁹³ World Health Organization. *Using AI to lead a healthier lifestyle*. ([Website](#) accessed 21 March 2025).
- ²⁹⁴ Organisation for Economic Co-operation and Development. *OECD Legal Instruments. Recommendation of the Council on Artificial Intelligence*. [Oecd.org](https://www.oecd.org/). 2019. ([Webpage](#) accessed 21 March 2025)
- ²⁹⁵ Nov O, Singh N, Mann D. Putting ChatGPT's medical advice to the (Turing) test: survey study. *JMIR Medical Education*. 2023;Jul10;9:e46939. ([Journal full text](#)) <https://doi.org/10.2196/46939>
- ²⁹⁶ Savoldi B, Gaido M, Bentivogli L, Negri M, Turchi M. Gender bias in machine translation. *Transactions of the Association for Computational Linguistics*. 2021;Aug;18;9:845-874. ([Journal full text](#)) <https://doi.org/10.1162/tacl-a-00401>
- ²⁹⁷ Meldau EL, Bista S, Melgarejo-González C, Niklas Norén G. WHO Uppsala Monitoring Centre. Automated De-identification of Case Narratives Using Deep Neural Networks for UK Yellow Card [Internet]. ([Full text](#) accessed 21 March 2025)
- ²⁹⁸ Norén GN, Orre R, Bate A, Edwards IR. Duplicate detection in adverse drug reaction surveillance. *Data Mining and Knowledge Discovery*. 2007;Jun;14:305-328. ([Journal full text](#)) <https://doi.org/10.1007/s10618-006-0052-8>
- ²⁹⁹ Ribeiro MT, Singh S, Guestrin C. "Why Should I Trust You?": Explaining the Predictions of Any Classifier [Internet]. [arXiv:1602.04938](https://arxiv.org/abs/1602.04938) 2016. ([Journal full text](#)) <https://doi.org/10.48550/arXiv.1602.04938>
- ³⁰⁰ Imran M, Bhatti A, King DM, Lerch M, Dietrich J, Doron G, Manlik K. Supervised machine learning-based decision support for signal validation classification. *Drug Safety*. 2022 May;45(5):583-596. ([Journal full text](#)) <https://doi.org/10.1007/s40264-022-01159-2>
- ³⁰¹ Plass M, Kargl M, Kiehl TR, Regitnig P, Geißler C, Evans T, Zerbe N, Carvalho R, Holzinger A, Müller H. Explainability and causability in digital pathology. *The Journal of Pathology: Clinical Research*. 2023;Jul;9(4):251-260. ([Journal full text](#)) <https://doi.org/10.1002/cjp2.322>
- ³⁰² Royal Society. *Explainable AI*. [Royalsociety.org](https://royalsocietypublishing.org/journal/rsos). 2024. ([Website](#) accessed 11 August 2024)
- ³⁰³ Plass M, Kargl M, Kiehl TR, Regitnig P, Geißler C, Evans T, Zerbe N, Carvalho R, Holzinger A, Müller H. Explainability and causability in digital pathology. *The Journal of Pathology: Clinical Research*. 2023;Jul;9(4):251-260. ([Journal full text](#)) <https://doi.org/10.1002/cjp2.322>
- ³⁰⁴ Spiker J, Kreimeyer K, Dang O, Boxwell D, Chan V, Cheng C, Gish P, Lardieri A, Wu E, De S, Naidoo J. Information visualization platform for postmarket surveillance decision support. *Drug Safety*. 2020;Sep;43:905-915. ([Journal abstract](#)) <https://doi.org/10.1007/s40264-020-00945-0>

³⁰⁵ Botsis T, Dang O, Kreimeyer K, Spiker J, De S, Ball R. A Decision-Support Platform Powered by AI and Humans-in-the-Loop Boosts Efficiency and Assures Quality in FDA's Pharmacovigilance. *International Society of Pharmacovigilance, 23rd Annual Meeting 2024, Montreal, Canada.* ([Webpage](#) accessed 21 March 2025)

³⁰⁶ Royal Society. *Explainable AI.* [Royalsociety.org](https://royalsocietypublishing.org/journal/rsos). 2024. ([Website](#) accessed 11 August 2024)

4092 **Content related to Fairness and Equity**

4093

4094 *While not all of the examples provided below are specific to PV, they illustrate the potential*

4095

4096 *Impact of inadequate data, bias from underrepresented populations and explicit bias potentially*
4097 *leading to unfair treatment of specific populations, underserved populations, and potential treatment*
4098 *inequality.*

4099

4100 *Example of inadequate training of AI solutions and/or inadequate data sets that introduced unfair*
4101 *bias and resulted in inequity.*

4102 In the US, prescription opioids are tracked through electronic databases, Prescription Drug
4103 Monitoring Programs (PDMPs). While not a PV specific example, Bamboo Health's NarxCare[®] is an
4104 example of an AI-powered tool that leverages PDMPs to calculate an opioid risk metric to predict the
4105 likelihood of a potential overdose, and although it is intended to support medical decisions, there
4106 have been observations that patients who are high health care utilizers with complex medical
4107 conditions may be discriminated against and underserved for pain management because of a "high
4108 risk score".³⁰⁷ The score is calculated based on limited data available in the PDMP and does not
4109 consider any other factors when calculating the risk score. One factor that influences the score is the
4110 number of prescribers. Patients treated at teaching hospitals with multiple healthcare prescribers
4111 may have "too many prescribing physicians" and they may be interpreted as seeking treatment from
4112 multiple physicians to obtain multiple prescriptions. An April 2021 study in Drug and Alcohol
4113 dependence found that "common data driven algorithms" misclassified 20% of patients with cancer
4114 who often see multiple specialists as patients seeking multiple physicians in an effort to obtain
4115 multiple opioid prescriptions. As noted by the authors, the PDMP data lacks diagnostic information
4116 and other critical patient context limiting ability to distinguish misuse from appropriate clinical use.
4117 An October 2021 study published in Drug and Alcohol Dependence conducted an independent
4118 validation study found that the NarxCare tool had a 17.2% false positive and 13.4% false negative.³⁰⁸

4119 Bias introduced because of data limited to PDMP, interpretation of data, e.g. patients with complex
4120 medical conditions, multiple prescribers are perceived as a high risk for abuse, lack of context for
4121 patient population, e.g. cancer patients requiring prolonged opioid use, can influence high scores.
4122 The threat to fairness and equity for patients within subgroups who have a high score assigned
4123 because of bias, potentially may not receive adequate pain management when the high score is
4124 considered in isolation.

4125 Within PV, the risk to fairness and equity are primarily from explicit biases that may result in negative
4126 impact or may result in discriminatory harm to subpopulations underserved by an AI solution. The
4127 NarxCare example, while not PV related, demonstrates both explicit bias from inadequate data, lack
4128 of context and implicit bias because negative stereotypes associated with "high health care utilizers"
4129 were applied.³⁰⁹

4130 *Example of bias applied because of under-represented populations*

4131 In Brazil, the assertiveness outcomes of the skin's lesions classification using artificial neural network
4132 in Caucasian patients and Brazilian patients were compared. The skin lesions were classified using
4133 basic architecture of convolutional neural networks (CNN). The ISIC database was used to train the
4134 neural network. This database was applied in about 25 thousand images of skin lesions. These images
4135 have included melanoma, melanocytic nevus, basal cell carcinoma, actinic keratosis, benign
4136 keratosis, dermatofibroma, vascular lesion, and squamous cell carcinoma lesions. The tests

4137 performed with ISIC patients had accuracy rates close to 90%. However, the accuracy rate was less
4138 than 40% when the tests were carried out with Brazilian patients. Thus, there was a great
4139 discrepancy in the assertiveness of the Artificial Neural Network when applied to Caucasian patients
4140 and Brazilian patients (Ref. K. Mundim et al).

4141 *Example of Explicit negative bias*

4142 In Appendix 3, “Examples of explainability in system development” included an example describing
4143 an AI triage system that incorrectly identified serious cases. The AI solution incorrectly learned to
4144 predict any adverse events associated with an OTC drug of interest as being non-serious because
4145 serious events were under-represented in the training data. This can also be considered an example
4146 of explicit negative bias with the inadequate data set resulting in incorrect assessments not
4147 recognized because of an explicit bias that it was not likely the OTC products in question would have
4148 serious adverse events associated with the use of the products. Since populations that may not have
4149 the same means to seek treatment at a medical facility or access to a healthcare professional may be
4150 reliant on OTC products, and these groups have a high likelihood of being from minority groups, a
4151 systematic misclassification of serious reports for OTC products as being non-serious could be seen
4152 as a threat to fairness and equity.

References

³⁰⁷ Sammer MBK, Akbari YS, Barth RA, Blumer SL, Dillman JR, Farmakis SG. Use of Artificial Intelligence in Radiology: Impact on Pediatric Patients, a White Paper from the ACR Pediatric AI Workgroup. *J Am Coll Radiol.* 2023;Aug;20(8):731-742. [\[Journal full text\] https://doi.org/10.1016/j.jacr.2023.06.003](https://doi.org/10.1016/j.jacr.2023.06.003)

³⁰⁸ Cochran G, Brown J, Yu Z, Frede S, Bryan MA, Ferguson A, et al. Validation and threshold identification of a prescription drug monitoring program clinical opioid risk metric with the WHO alcohol, smoking, and substance involvement screening test. *Drug Alcohol Depend.* 2021;228:109067. [\[Journal full text\] https://doi.org/10.1016/j.drugalcdep.2021.109067](https://doi.org/10.1016/j.drugalcdep.2021.109067)

³⁰⁹ Siegel Z. In a world of stigma and bias, can a computer algorithm really predict overdose risk? *Ann Emerg Med.* 2022;Jun;79(6):A16-A19. [\[Journal full text\] https://doi.org/10.1016/j.annemergmed.2022.04.006](https://doi.org/10.1016/j.annemergmed.2022.04.006)

4153

4154 **APPENDIX 5: CIOMS Working Group membership and**
 4155 **meetings**

4156

4157 The CIOMS Working Group XIV on *Artificial intelligence in* pharmacovigilance included the following
 4158 groups of stakeholders: academics, pharmaceutical companies, regulatory authorities, as well as
 4159 national and international organisations.

4160

Academia		
Name	Company/Organisation	Country
Altman, Russ	Stanford University	USA
Botsis, Taxiarchis	Johns Hopkins University School of Medicine	USA
Dogné, Jean-Michel	University of Namur	Belgium
Pharmaceutical companies		
Name	Company/Organisation	Country
Amelio, Justyna	AbbVie	UK
Barrios, Mariane	Merck Sharp & Dohme	Colombia
Bate, Andrew	GSK	UK
Bellur, Arvind	CSL Behring	USA
Berridge, Adrian	Takeda Development Center Americas, Inc	USA
Carroll, Hua	Biogen	USA
Cherkas, Yauheniya	Johnson & Johnson	USA
Comfort, Shaun	Genentech	USA
Cooper, Selin	AbbVie	UK
Diniz, Mariane	Bayer	Brazil
Domalik, Douglas	AstraZeneca	UK
Franco, Piero Francesco	Pfizer	Italy
Girod, Julie	Sanofi	USA
Grabowski, Neal	Sanofi	USA
Hauben, Manfred	Merck KGaA, Darmstadt, Germany	USA
Henn, Thomas	United Therapeutics	USA
Kara, Vijay	GSK	UK
Kempf, Dieter	Genentech	USA
Kidos, Kostadinos	Formerly Takeda Development Center Americas, Inc	USA
Lorenz, Denny	formerly Bayer AG	Germany

APPENDIX 5: CIOMS Working Group membership and meetings

MacEntee Pileggi, Beth	Johnson & Johnson	USA
Patel, Ravi	United Therapeutics	USA
Reinhard Pietzsch, John	Bayer	Germany
Römning, Hans-Jörg	Merck KGaA, Darmstadt, Germany	Germany
Savage, Elizabeth	Johnson & Johnson	USA
Straus, Walter	Moderna	USA
Whitehead, James	AstraZeneca	UK
Regulatory authorities		
Name	Company/Organisation	Country
Ball, Robert	United States Food and Drug Administration (US FDA)	USA
Buch, Brian	Medicines and Healthcare products Regulatory Agency (MHRA)	UK
Durand, Julie	European Medicines Agency (EMA)	
Egebjerg Juul, Kirsten	Danish Medicines Agency (DKMA)	Denmark
Harrison, Kendal	Medicines and Healthcare products Regulatory Agency (MHRA)	UK
Hirokawa-Voorburg, Satoko	Health and Youth Care Inspectorate (HYCI)	The Netherlands
Horst, Alexander	Swissmedic	Switzerland
Jensen, Morten	Danish Medicines Agency (DKMA)	Denmark
Kikuchi, Yuki	Pharmaceuticals and Medical Devices Agency (PMDA)	Japan
Kjær, Jesper	Danish Medicines Agency (DKMA)	Denmark
Ling, Benny	Health Canada	Canada
Da Luz Carvalho Soares, Monica	Brazilian Health Regulatory Agency (ANVISA)	Brazil
Maniwa, Harumi	Pharmaceuticals and Medical Devices Agency (PMDA)	Japan
Matsunuga, Yusuke	Pharmaceuticals and Medical Devices Agency (PMDA)	Japan
McAteer, Richard	Health Canada	Canada
Mentzer, Dirk	Paul-Ehrlich-Institut (PEI)	Germany
Messelhäußer, Manuela	Formerly Paul-Ehrlich-Institut (PEI)	Germany
Moreira Cruz, Flávia	Brazilian Health Regulatory Agency (ANVISA)	Brazil
Perez, Nicolas	Swissmedic	Switzerland
Scholz, Irene	Swissmedic	Switzerland
Stammschulte, Thomas	Swissmedic	Switzerland

Tregunno, Phil	Medicines and Healthcare products Regulatory Agency (MHRA)	UK
National and international organisations		
Name	Company/Organisation	Country
Mathur, Roli	Indian Council of Medical Research	India
Meldau, Eva-Lisa	Uppsala Monitoring Centre/World Health Organization	Sweden
Norén, Niklas	Uppsala Monitoring Centre/World Health Organization	Sweden
Rosenfeld, Stephen	North Star Review Board	USA
Yau, Brian	World Health Organization	Switzerland
CIOMS		
Name	Company/Organisation	Country
Heaton, Stephen	Individual expert	Germany
Hill, Sanna	CIOMS	Switzerland
Le Louët, Hervé	CIOMS	Switzerland
Rägo, Lembit	CIOMS	Switzerland
Rannula, Kateriina	CIOMS	Estonia
Tsintis, Panos	CIOMS	UK

4161

4162 The Working Group XIV will have met ten times from 2022 to 2025 at the time of writing and most of
4163 the meetings at this stage will have been hybrid in nature.

4164

- | | | |
|------|-------------------------|--|
| 4165 | 1. Geneva, Switzerland | 18-19 May 2022 |
| 4166 | 2. Geneva, Switzerland | 10-11 October 2022 |
| 4167 | 3. Virtual meeting | 19 January 2023 |
| 4168 | 4. Virtual meeting | 12 April 2023 |
| 4169 | 5. Zurich, Switzerland | 6-7 June 2023 |
| 4170 | 6. Virtual meeting | 8 November 2023 |
| 4171 | 7. Virtual meeting | 11 January 2023 |
| 4172 | 8. Geneva, Switzerland | 7-8 March 2024 |
| 4173 | 9. Darmstadt, Germany | 24-25 September 2024 |
| 4174 | 10. Geneva, Switzerland | 25-26 June 2025 (planned at time of writing) |

4175

4176

4177 **APPENDIX 6: Public consultation commentators**

4178

4179 (to follow)

4180