



1  
2  
3

**Draft**

IMDRF/AIML WG/N93 DRAFT: 202X

# **Technical Framework for Artificial Intelligence Life Cycle Management**

**AUTHORING GROUP**

**Artificial Intelligence/Machine Learning-enabled Working  
Group**

4  
5

7 April 2026

# 6 Preface

7 © Copyright 202X by the International Medical Device Regulators Forum.

8 This work is copyright. Subject to these Terms and Conditions, you may download, display, print,  
9 translate, modify and reproduce the whole or part of this work for your own personal use, for research,  
10 for educational purposes or, if you are part of an organisation, for internal use within your organisation,  
11 but only if you or your organisation do not use the reproduction for any commercial purpose and retain  
12 all disclaimer notices as part of that reproduction. If you use any part of this work, you must include the  
13 following acknowledgement (delete inapplicable):

14 “[Translated or adapted] from [insert name of publication], [year of publication], International Medical  
15 Device Regulators Forum, used with the permission of the International Medical Device Regulators  
16 Forum. The International Medical Device Regulators Forum is not responsible for the content or  
17 accuracy of this [adaption/translation].”

18 All other rights are reserved and you are not allowed to reproduce the whole or any part of this work in  
19 any way (electronic or otherwise) without first being given specific written permission from IMDRF to do  
20 so. Requests and inquiries concerning reproduction and rights are to be sent to the IMDRF Secretariat.

21 Incorporation of this document, in part or in whole, into another document, or its translation into  
22 languages other than English, does not convey or represent an endorsement of any kind by the IMDRF.

23

24 **[Name], IMDRF Chair**

25

# Contents

27	<b>1. Introduction</b>	<b>4</b>
28	<b>2. Purpose and Scope</b>	<b>5</b>
29	2.1. Purpose of the document	5
30	2.2. Scope of the document	5
31	<b>3. References</b>	<b>7</b>
32	<b>4. Universal Concepts that apply across the AI Life Cycle</b>	<b>9</b>
33	4.1. Quality Management System (QMS)	9
34	4.2. Risk Management	9
35	4.3. Human Oversight	11
36	4.4. Cybersecurity	12
37	<b>5. AI-enabled Medical Device Life Cycle Steps</b>	<b>14</b>
38	5.1. Planning and Design	15
39	5.2. Data Collection and Management	16
40	5.3. Model Building and Tuning	19
41	5.4. Verification and Validation, including Clinical Evaluation	22
42	5.5. Deployment	25
43	5.6. Operations and Monitoring	27
44	5.7. Real-World Performance Evaluation	29
45	5.8. Sunsetting	30
46	<b>6. Transparency and Labelling</b>	<b>31</b>
47	<b>7. Conclusion</b>	<b>32</b>
48	<b>Appendix A: Traceability between GMLP and Life Cycle Steps/Document Sections</b>	<b>33</b>
49		
50	<b>Appendix B: Examples of Common Evaluation Metrics</b>	<b>35</b>
51	<b>Appendix C: Labelling Elements</b>	<b>36</b>
52		

# 1. Introduction

54 The International Medical Device Regulators Forum (IMDRF) recognizes the growing role of artificial  
55 intelligence (AI) technologies in medical devices and has continued its efforts to establish a  
56 harmonized approach for the oversight, evaluation, and use of AI-enabled medical devices. Building  
57 on its previous work, including the publication IMDRF/AIML WG (Artificial Intelligence/Machine  
58 Learning-enabled Working Group)/N88 FINAL:2025: *Good Machine Learning Practice (GMLP) for  
59 Medical Device Development: Guiding Principles*, IMDRF is introducing this document focused on  
60 considerations for AI-enabled medical devices across their total product life cycle.

61 The GMLP principles describe foundational best practices for the development of AI-enabled medical  
62 devices, emphasizing areas such as data quality, model transparency, performance evaluation, and  
63 the role of multidisciplinary expertise. These principles underpin each step of the AI life cycle, and this  
64 document provides relevant GMLP references to help provide a foundational understanding of  
65 applicable principles.

66 Together with the broader IMDRF mission, this work contributes to the establishment of globally  
67 harmonized considerations that, when applied across a device's total product life cycle, can foster  
68 innovation while protecting public health.

69

70

## 71 2. Purpose and Scope

### 72 2.1. Purpose of the document

73 This document provides an internationally harmonized technical framework to help promote  
74 responsible innovation and patient-centricity by facilitating the secure, safe, ethical, and effective  
75 design, development, deployment, maintenance, use, and decommission of AI-enabled medical  
76 devices. .

77 The purpose of this document is to:

- 78 • Provide foundational information on AI-enabled medical device life cycle management;
- 79 • Highlight universal concepts applicable to all life cycle steps;
- 80 • Provide an overview of concepts and considerations for each step of the AI-enabled medical  
81 device life cycle; and
- 82 • Provide references to applicable, internationally-recognized standards and resources.

83 The document is not meant to:

- 84 • Replace or conflict with existing IMDRF publications such as those published by the Software  
85 as a Medical Device or Cybersecurity Working Groups (WG.H). However, it may in part relate  
86 to or overlap with those publications and is intended to be complementary in those  
87 circumstances. Where relevant, this document references other IMDRF publications on  
88 related topics (e.g. Risk Management);
- 89 • Serve as regulation nor guidance. It is not intended to be an interpretation of any jurisdiction's  
90 laws and regulations, does not provide recommendations for market submission within a  
91 specific jurisdiction, and does not imply a convergence of regulations across jurisdictions.  
92 Instead, this document aims to describe harmonized concepts and general considerations for  
93 the AI-enabled medical device life cycle. Individual jurisdictions may apply and align some or  
94 all of these concepts to their particular regulatory framework; and
- 95 • Provide specific information on how to prepare a regulatory submission. However, in some of  
96 its sections, documentation is identified that may be helpful when communicating with  
97 regulatory authorities.

### 98 2.2. Scope of the document

99 This document is intended for manufacturers that are developing AI-enabled medical devices to  
100 identify considerations as they make critical technical and governance choices. It is intended to  
101 provide considerations specific to where AI-enabled medical devices may necessitate new or different  
102 approaches than those more broadly applicable to medical devices and medical device software,  
103 building upon existing work and harmonization efforts IMDRF has published. Rather than exhaustively  
104 cover all steps of the AI-enabled medical device life cycle, the document highlights key considerations  
105 as well as internationally-recognized standards and resources for manufacturers to additionally  
106 consider. Given the rapid pace of advancement in the field of AI, this document attempts to reflect the  
107 most current considerations for manufacturers of AI-enabled medical devices.

108 This document applies to AI-enabled medical devices, which include Machine Learning-enabled  
109 Medical Devices (MLMD). MLMD are medical devices that use machine learning, in part or in whole,  
110 to achieve their intended purpose<sup>1</sup>. Machine Learning (ML) models are developed by ML training  
111 algorithms through analysis of data, without models being explicitly programmed<sup>1</sup>. While this  
112 document is intended to apply generally to MLMD, certain types of MLMD that are enabled by or that  
113 incorporate generative AI, autonomous or adaptive models may warrant additional considerations.  
114 While this document makes some recommendations for these subsets of MLMD, some aspects of this  
115 AI life cycle may differ for such technologies and this document is not intended to be comprehensive.

116 The concepts of a “model” and “AI-enabled medical device” appear throughout this document but are  
117 not used interchangeably. A model is a “mathematical construct that generates an inference or  
118 prediction based on new input data, and is the result of an ML training algorithm learning from data”.<sup>2</sup>  
119 As stated above, AI-enabled medical devices, that include MLMD, are devices that use AI or ML, in  
120 part or in whole, to achieve their intended medical purpose, and incorporate one or more “models”  
121 into their design to achieve this purpose. Certain steps of the life cycle described in this document  
122 focus on the model (e.g., Data Collection and Management and Model Building and Tuning), while  
123 other steps describe considerations for AI-enabled medical devices that incorporates the model(s)).  
124 This distinction will be highlighted in the sections throughout this document.

125 Explainability and Interpretability are also discussed throughout this document. Explainability “refers  
126 to a representation of mechanisms underlying AI systems’ operation”<sup>3</sup>. It includes not only  
127 understanding how decisions are made, but also why they are appropriate and trustworthy in the  
128 clinical context. Interpretability “refers to the meaning of AI systems’ output in the context of their  
129 designed functional purposes”<sup>3</sup>. In other words, the ability for clinicians and healthcare professionals,  
130 patients, and other users to comprehend how the device arrives at its outputs or recommendations.  
131 Understanding the difference between these concepts is important when applying considerations  
132 outlined throughout this document.

133 This document acknowledges that AI-enabled medical devices have a variety of users and can be  
134 deployed in a variety of environments, including at medical institutions, hospitals, or in other  
135 healthcare settings (collectively referred to as “sites”) as well as on a patient’s own consumer  
136 electronics or in a patient’s home. While the considerations in this document are intended to apply  
137 across the spectrum of intended use environments, some considerations are specific to the deployed  
138 environment and are indicated as such throughout the document.

139

140

---

<sup>1</sup> IMDRF/AIMD WG/N67 (Edition 1):2022 *Machine Learning-enabled Medical Devices: Key Terms and Definitions*

<sup>2</sup> IMDRF/AIMD WG/N67 (Edition 1):2022 *Machine Learning-enabled Medical Devices: Key Terms and Definitions*

<sup>3</sup> National Institute of Standards and Technology. (2023). *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*. <https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf>.

## 3. References

- 142 IMDRF/AIML WG/N88 FINAL:2025: *Good Machine Learning Practice (GMLP) for Medical Device*  
 143 *Development: Guiding Principles*
- 144 IMDRF/AIMD WG/N67 (Edition 1):2022 *Machine Learning-enabled Medical Devices: Key Terms and*  
 145 *Definitions*
- 146 National Institute of Standards and Technology. (2023). *Artificial Intelligence Risk Management*  
 147 *Framework (AI RMF 1.0)*. <https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf>.
- 148 ISO 13485 Third Edition 2016-03-01 *Medical Devices—Quality Management Systems—*  
 149 *Requirements for Regulatory Purposes*
- 150 IMDRF/SaMD WG/N23 FINAL:2015 *Software as a Medical Device (SaMD): Application of Quality*  
 151 *Management System*
- 152 ANSI/AAMI/ISO 14971:2019 *Medical devices—Application of risk management to medical devices*
- 153 ISO/TR24971:2020 *Medical devices — Guidance on the application of ISO 14971*
- 154 AAMI TIR 34971:2023 *Application of ISO 14971 to machine learning in artificial intelligence— Guide*
- 155 ENISA. (2021, December 14). *Securing Machine Learning Algorithms*, Enisa Europa.  
 156 <https://www.enisa.europa.eu/publications/securing-machine-learning-algorithms>
- 157 Apostol Vassilev (NIST), Alina Oprea (Northeastern University), Alie Fordyce (Robust  
 158 Intelligence), Hyrum Anderson (Robust Intelligence) (2024, January) NIST AI 100-2 E2023  
 159 *Adversarial Machine Learning: A Taxonomy and Terminology of Attacks and Mitigations*. NIST.  
 160 <https://csrc.nist.gov/pubs/ai/100/2/e2023/final>
- 161 IEC 62304:2006/A1:2016 *Medical device software – Software life cycle processes*
- 162 IMDRF/MC/N79 DRAFT:2023 *Guiding Principles to Support Medical Device Health Equity*
- 163 ISO/IEC 5259-4:2024 *Artificial intelligence — Data quality for analytics and machine learning (ML) —*  
 164 *Part 4: Data quality process framework*
- 165 IMDRF/CYBER WG/N60 FINAL:2020 *Principles and Practices for Medical Device Cybersecurity*  
 166
- 167 IMDRF/CYBER WG/N70 FINAL:2023 (Edition 1) *Principles and Practices for the Cybersecurity of*  
 168 *Legacy Medical Devices*
- 169
- 170 IMDRF/CYBER WG/N73 FINAL:2023 (Edition 1) *Principles and Practices for Software Bill of*  
 171 *Materials for Medical Device Cybersecurity*
- 172
- 173 IMDRF/SaMD WG/N41 FINALFINALN41 FINAL:2017 *Software as a Medical Device (SaMD): Clinical*  
 174 *Evaluation*
- 175 ISO/IEC 23053:2022 *Framework for Artificial Intelligence (AI) Systems Using Machine Learning (ML)*
- 176 ISO/IEC TS 25058:2024 *Systems and software engineering — Systems and software Quality*  
 177 *Requirements and Evaluation (SQuaRE) — Guidance for quality evaluation of artificial intelligence*  
 178 *(AI) systems*
- 179 ISO/IEC 25059:2023 *Software engineering — Systems and software Quality Requirements and*  
 180 *Evaluation (SQuaRE) — Quality model for AI systems*
- 181 IEEE 2941-2021 IEEE Standard for Artificial Intelligence (AI) Model Representation, Compression,  
 182 Distribution, and Management
- 183 IMDRF/UDI WG/N48 FINAL:2019 *Unique Device Identification system (UDI system) Application*  
 184 *Guide*

- 185   IMDRF/GRRP WG/N52 FINAL:2024 (Edition 2) *Principles of Labelling for Medical Devices and IVD*  
186   *Medical Devices*
- 187   IMDRF/SaMD WG/N81 DRAFT:2024 *Medical Device Software: Considerations for Device and Risk*  
188   *Characterization*
- 189
- 190   ISO 14155:2020 *Clinical investigation of medical devices for human subjects - Good clinical practice*
- 191   ISO/IEC 62366-1:2015 *Medical devices Part 1: Application of usability engineering to medical devices*
- 192

DRAFT

# 4. Universal Concepts that apply across the AI Life Cycle

This section highlights universal concepts applicable to all steps of an AI-enabled medical device’s life cycle. Many of these universal concepts and their considerations are not new or unique for AI-enabled medical devices, but can be adapted for the unique characteristics of these devices.

## 4.1. Quality Management System (QMS)

A QMS is a comprehensive framework of policies, processes, and procedures that medical device manufacturers implement to ensure consistent quality throughout the device life cycle. It is important for manufacturers to implement a QMS that complies with applicable requirements and established standards as defined by regulatory authorities in their regulations. This may include, for example, the international standard ISO 13485 Third Edition 2016-03-01 *Medical Devices—Quality Management Systems— Requirements for Regulatory Purposes* (hereafter referred to as ISO 13485:2016). Furthermore, *IMDRF/SaMD WG/N23 FINAL:2015 Software as a Medical Device (SaMD): Application of Quality Management System* (hereafter referred to as *IMDRF/SaMD WG/N23 FINAL:2015*) provides guidance to manufacturers and regulators on QMS practices specific to SaMD, which can generally be applied to AI-enabled medical devices.

AI-enabled medical devices benefit from implementation of scalable life cycle support processes that emphasize safety-focused risk management throughout all life cycle steps. For example, QMS requirements management captures functional specifications as well as clinical environment considerations, such as how AI outputs will be interpreted by healthcare providers or patients. Additionally, configuration management and control processes maintain traceability of training data and AI model versions, which is especially important for models that adapt to new data in their intended use environment. Finally, comprehensive post-market monitoring and surveillance mechanisms can help actively monitor performance of the AI-enabled medical device in real-world environments and enable rapid response to emerging safety concerns or performance degradation over time.

## 4.2. Risk Management

Risk management principles and processes for AI-enabled medical devices follow the same fundamental framework established for other medical devices, including SaMD. As with all medical devices, AI-enabled medical devices carry risks that are necessary for manufacturers to systematically address. Manufacturers are responsible for evaluating these risks through comprehensive risk analysis and assessment processes, implementing appropriate controls to reduce identified risks to acceptable levels, ensuring medical devices do not pose unacceptable risk to patients, users, or others, and demonstrating that clinical benefits outweigh residual risks through appropriate evidence and documentation.

This section highlights key risks that are unique to AI-enabled medical devices and that warrant careful consideration and mitigation throughout the life cycle. The examples provided below serve as starting points for AI-related risk management activities and do not constitute an exhaustive list of all possible AI-related risks that a medical device may encounter.

#### 234 **4.2.1. Risks Related to Information**

235 The “black box” nature of some AI models makes it challenging to understand how and why certain  
236 outputs are produced or why certain decisions are made by the model (when a model is embedded  
237 and controls hardware or autonomously makes decisions), creating unique transparency,  
238 explainability, interpretability concerns. Risks related to information can include:

- 239 • Inaccurate or misleading outputs
- 240 • Incomplete information presentation, including lack of explanation on probability of error

#### 241 **4.2.2. Risks Related to the Human-AI Interaction**

242 The integration of AI-enabled medical devices into the clinical workflow can significantly alter  
243 interaction dynamics in healthcare settings. One such manifestation of this is automation bias, which  
244 refers to a tendency of users to rely on automated systems over their own knowledge or skills when  
245 making decisions. AI-enabled medical devices may also be perceived as having human-like cognitive  
246 or reasoning abilities, potentially leading to over-reliance that compromises independent reasoning  
247 and decision-making. These types of risk may develop over time and may not be detectable in testing  
248 or early use of a device. Risks related to the human-AI interaction can include:

- 249 • Over-reliance or automation bias
- 250 • Under-reliance or dismissal
- 251 • Workflow disruption
- 252 • Verification fatigue
- 253 • Information overload (e.g., alert fatigue)
- 254 • De-learning of clinical knowledge (due to over-reliance over time)

#### 255 **4.2.3. Risks Related to Model Training and Data Quality**

256 AI-enabled medical devices may be uniquely vulnerable to performance degradation from, for  
257 example, data drift when they are exposed to real-world data distributions that significantly differ from  
258 the training data on which their model(s) were developed. Often, the substantial volume and  
259 complexity of data required for training introduces scalability challenges in data curation, annotation,  
260 and quality assurance, which limits generalizability. Risks related to model training and data quality  
261 can include:

- 262 • Training data bias (i.e., underperformance due to lack of training data for certain demographic  
263 groups)
- 264 • Data drift, including changes in the context of use, such as patient demographics
- 265 • Incomplete or missing data
- 266 • Labelling or annotation errors
- 267 • Out of distribution (OOD) inputs
- 268 • Knowledge corpora fragmentation or misalignment (e.g., discrepancies between the model’s  
269 training or reference data and current medical standards or the specific clinical environment  
270 where it will be deployed)

271 **4.2.4. Risks Related to Deployment and Post-Market Monitoring and**  
272 **Performance**

273 AI-enabled medical devices may rely on complex algorithms, substantial computational resources,  
274 and interconnected system dependencies that create novel failure modes where infrastructure  
275 inadequacies, integration issues, or deployment errors can compromise medical device performance  
276 and potentially lead to patient harm if not properly identified and controlled. Risks related to  
277 deployment and post-market monitoring and performance can include:

- 278 • Integration, interoperability, or compatibility issues
- 279 • Performance degradation due to changes in third-party, general-purpose models (see Section  
280 5.3 for more information on these types of models) incorporated into a device
- 281 • Software governance issues including lack of version control when the AI model or device is  
282 modified
- 283 • Infrastructure inadequacy and computational or scalability constraints
- 284 • Misaligned performance and calibration if context-of-use differs from validation

285 **Additional Standards and References to be Considered**

- 286 • ANSI/AAMI/ISO 14971:2019 *Medical devices—Application of risk management to medical*  
287 *devices*
- 288 • ISO/TR24971:2020 *Medical devices — Guidance on the application of ISO 14971*
- 289 • AAMI TIR 34971:2023 *Application of ISO 14971 to machine learning in artificial intelligence—*  
290 *Guide*

291 **4.3. Human Oversight**

292 Human oversight, including that of clinicians, health care providers, patients, and lay users, is  
293 essential throughout the entire life cycle of AI-enabled medical devices to ensure that human and  
294 clinical expertise informs model development, validates real-world performance, and maintains  
295 appropriate human-AI collaboration that prioritizes patient safety and effective clinical decision-  
296 making.

297 For example, human involvement in identifying user needs for an AI-enabled medical device can help  
298 to ensure that the device design accounts for the user’s ability to interpret AI outputs, override or  
299 reverse automated recommendations, and intervene or interrupt automation. As with all medical  
300 devices, usability and human factors testing with representative users can help identify potential use-  
301 related hazards or validate the design of controls for use-related hazards.

302 Clinician expertise and input can be leveraged in activities such as feature selection, validating the  
303 clinical relevance of models, data labelling and annotation and identifying potential biases. Clinical  
304 involvement also ensures that clinical decisions remain accurate, safe, and contextually appropriate  
305 and that complex or atypical cases are adequately understood and addressed.

306 Post-market monitoring and surveillance benefits from human oversight in monitoring to detect  
307 performance degradation, identify unexpected failure modes, and assess real-world effectiveness  
308 across patient populations and clinical settings. This oversight can also help manufacturers to remain  
309 aware of the possible tendency of automatically relying or over-relying on the output produced by the  
310 AI-enabled medical device (automation bias) and make adjustments for such biases.

311 **Additional Standards and References to be Considered**

- 312 • ISO/IEC 62366-1:2015 *Medical devices Part 1: Application of usability engineering to medical*  
313 *devices*

## 314 4.4. Cybersecurity

315 Similar to other medical devices, cybersecurity is an important consideration throughout the AI-  
316 enabled medical device life cycle. The large quantity of data needed for the development and  
317 validation of AI-enabled medical devices and the sensitive patient data processed during use can  
318 make these devices an attractive target for data theft. Additionally, maliciously poisoned data could  
319 negatively impact the AI-enabled medical device's performance. This section focuses on the  
320 cybersecurity related to the potential harm to an affected person, though it acknowledges that data  
321 extraction and exploitation are also possible issues.

322 Mechanisms to help protect access to the data used for the development and validation during the AI-  
323 enabled medical device life cycle include, but are not limited to, data anonymization, data separation,  
324 system separation, data validation, data encryption, access control, logging, auditing, model  
325 simplification, anomaly detection techniques, and periodic audits. For AI-enabled medical devices  
326 specifically, it is essential to follow best practices for mapping data sources and controlling data  
327 suppliers and labellers as part of QMS processes. This can help control unwanted bias and data  
328 poisoning stemming from outside sources in both development and post-market performance  
329 optimization (when carried out with appropriate regulatory approval).

330 In some instances, due to the complex nature of AI-enabled medical device design, it might be difficult  
331 to recognize when data poisoning occurs because the cause-effect relationship between input data  
332 and AI-enabled medical device output is not always transparent and explainable, such that malicious  
333 degradation may be thought to be within performance specifications of the model or mistaken for  
334 natural performance drift. In addition to data derived vulnerabilities and complexities with  
335 transparency and explainability, the AI model itself can be under threat potentially leading to model  
336 inversion, extraction and evasion occurring. As outlined below, robust security controls and monitoring  
337 activities can help mitigate these vulnerabilities.

338 Using a secure product development framework to manage cybersecurity risks can help identify and  
339 reduce the number and severity of vulnerabilities in devices. Using device design processes to  
340 support secure product development and maintenance may include Threat Modeling, Cybersecurity  
341 Risk Assessments, interoperability considerations, third party software components, Cybersecurity  
342 unresolved anomalies and Risk Management.

343 Deployed AI-enabled medical devices are vulnerable to cybersecurity threats due to their potential  
344 reliance on continuous data flows, cloud connectivity for model updates, and complex software  
345 architectures that may create multiple entry points for malicious actors. Correct deployment and  
346 continuous post-market monitoring and surveillance, which is described in more detail in the  
347 Deployment (Section 5.5) and Operations and Monitoring (Section 5.6) life cycle steps below, is an  
348 important consideration to minimize cybersecurity threats in the post-market setting. Unwanted  
349 access may be mitigated via mechanisms such as; limiting the frequency of updates to batches,  
350 hosting the AI-enabled medical device locally, if possible and ensuring security controls are  
351 centralized to manage monitoring, auditing and validation more reliably. It is important to note that  
352 hosting of AI-enabled medical devices on their manufacturers' server may introduce additional  
353 reliability considerations that are not unique to AI-enabled devices.

354 Furthermore, when deploying AI-enabled medical devices that contain adaptive models that continue  
355 to learn after their initial release, it is important for manufacturers to consider how data inputs to the  
356 device in the deployed environment are controlled and validated, and who is handling and accessing  
357 the data prior to its input into the AI-enabled medical device. The primary responsibility for validation  
358 of input data resides with the legal manufacturer; in practice, additional parties such as end users and  
359 sites at which the device is being deployed may also be involved. If the model is autonomously  
360 adapting and improving, it may warrant further scrutiny of interconnected systems that deliver this  
361 data and their attack surfaces and it remains critical that updates are verified for safety prior to  
362 implementation. Again, performing post-market validation checks on algorithm performance after  
363 updates with adequate human oversight as described in Section 4.3, may help identify attempts at  
364 data poisoning or other disruptive attacks related to the AI-enabled medical device adapting to new  
365 data.

### 366 **Additional Standards and References to be Considered**

367  
368  
  
369  
370  
371  
372  
373  
374  
375  
  
376  
  
377

- ENISA. (2021, December 14). *Securing Machine Learning Algorithms*, Enisa Europa. <https://www.enisa.europa.eu/publications/securing-machine-learning-algorithms>
- Apostol Vassilev (NIST), Alina Oprea (Northeastern University), Alie Fordyce (Robust Intelligence), Hyrum Anderson (Robust Intelligence) (2024, January) NIST AI 100-2 E2023 *Adversarial Machine Learning: A Taxonomy and Terminology of Attacks and Mitigations*. NIST. <https://csrc.nist.gov/pubs/ai/100/2/e2023/final>
- IMDRF/CYBER WG/N60 FINAL:2020 *Principles and Practices for Medical Device Cybersecurity*

DRAFT

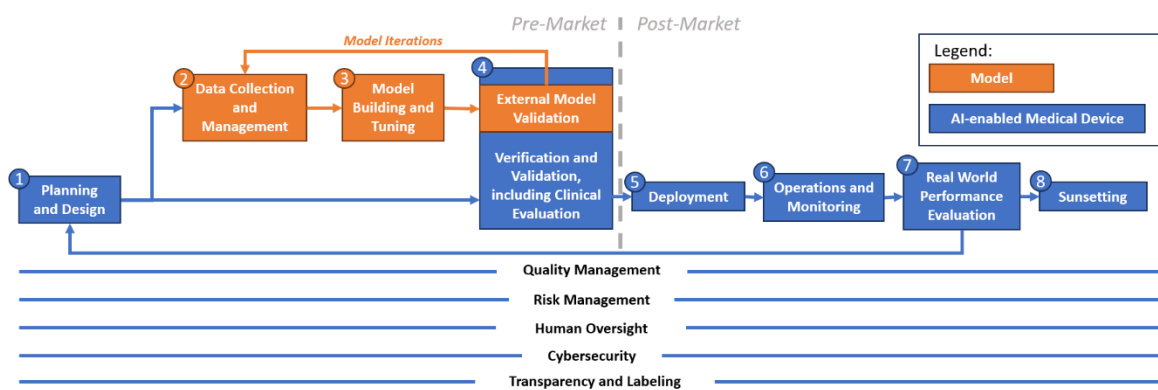
# 5. AI-enabled Medical Device Life Cycle Steps

378

379

380 The AI-enabled medical device life cycle presented in this document builds upon established  
 381 frameworks, such as ISO 13485:2016, IMDRF/SaMD WG/N23 FINAL:2015, and the International  
 382 Electrotechnical Commission (IEC) 62304 A1:152015 *Medical device software – Software life cycle  
 383 processes*<sup>4</sup> standard with additional considerations for the uniqueness of incorporating AI into a  
 384 medical device, including AI’s data-driven nature and complexity of the human-AI interaction. The  
 385 core principles of medical device and SaMD life cycles of systematic planning, requirements  
 386 management, traceability, risk management and validation remain essential. The figure below depicts  
 387 the steps of the life cycle described in this document, as well as the universal concepts described in  
 388 Section 4 above.

389 **Figure 1: AI-Enabled Medical Device Life Cycle**



390

391 *This figure represents the general life cycle outlined in this document. However, manufacturers may follow different*  
 392 *iterations of this cycle depending on what might be needed for the model and/or device. One example of this flexibility*  
 393 *is the pathway from Step 7 directly to Step 1 and then to Step 4. This sequence illustrates scenarios where real-world*  
 394 *performance evaluation identifies incidents or performance drift following deployment. In such cases, a model update*  
 395 *may not be warranted but could instead involve labelling modifications or a comprehensive re-evaluation of the model*  
 396 *and/or device to reestablish baseline performance metrics.*

397

398

<sup>4</sup> IEC 62304:2006 *Medical device software – Software life cycle processes* is one commonly used software development life cycle standard that can be used to develop an AI-enabled medical device software life cycle process

## 399 5.1. Planning and Design

400 Similar to of the planning for other medical devices, and, per IMDRF/SaMD WG/N23 FINAL:2015,  
401 "the objective of planning is to provide a roadmap to be followed during the product development life  
402 cycle" following "a methodical and rigorous plan for managing projects such as a plan-do-check-act  
403 approach". For AI-enabled medical devices, GMLP Guiding Principle 1<sup>5</sup> emphasizes that the intended  
404 use / intended purpose of the device is well understood, and multi-disciplinary expertise will be  
405 leveraged throughout the total product life cycle. This can be achieved with a comprehensive product  
406 definition, which includes understanding the clinical objectives, existing standards of care, use  
407 environment, potential risks and controls, relevant clinical workflows and user needs and constraints  
408 as well as team resource planning.

409 While the Planning and Design life cycle step encompasses many of the preliminary plans made for  
410 subsequent life cycle steps, this section generally highlights key considerations for this step, with  
411 additional details on each particular step in the sections below. Furthermore, this section points out  
412 considerations for both the AI-enabled medical device and its model(s), which both have their own  
413 unique planning and design considerations.

414 *Model Selection and Anticipating Risks:* As discussed in Section 4.2 above, the risks and  
415 implementation challenges of incorporating AI may exceed its clinical benefits for certain medical  
416 uses. It is important for manufacturers to consider whether an AI model (versus, for example, a rules-  
417 based approach) is the right tool to achieve the intended purpose of the device, and if so, what kind of  
418 model is best suited for that purpose. If the decision is made to move forward with incorporating an AI  
419 model into a medical device, selection of the appropriate model up front may streamline the risk  
420 control process; by enabling faster and more accurate identification of potential risks, manufacturers  
421 can implement preventative risk control measures rather than relying solely on corrective actions after  
422 development. This strategy is often more effective than applying risk control measures retrospectively  
423 to compensate for suboptimal model choices. Manufacturers may also want to consider opting for  
424 simplicity if explainability and interpretability is critical to the intended use environment.

425 *Data Availability and Suitability:* AI-enabled medical devices and their models rely heavily on fit-for-  
426 purpose, representative<sup>6</sup> data, and it is important for manufacturers to begin thinking about data  
427 availability and suitability, and generally the feasibility of data collection, during the Planning and  
428 Design step. This can help manufacturers preliminarily understand the quantity, quality and  
429 distribution of available data that would be needed for the design, development, and validation of their  
430 AI model(s). Furthermore, it is important for manufacturers to consider quantity of data needed such  
431 that training datasets can be independent from test data sets.<sup>7</sup>

432 *Model, Infrastructure, and Use Requirements:* In this step, it is important for manufacturers to  
433 contemplate the spectrum of technical considerations, from the architecture and model design to be  
434 employed and their requirements, such as interpretability, stability, and performance requirements, to  
435 the infrastructure needed for deployment, monitoring and maintenance.<sup>8</sup> For AI-enabled medical  
436 devices specifically, deployment may include site-specific localization or customization (with  
437 appropriate regulatory approval), and it is important that the infrastructure account for version control  
438 and resets or rollbacks in the event that recalls or issues occur. It is beneficial for manufacturers to  
439 not only define requirements from a technological perspective, but also for how the AI-enabled  
440 medical device and its model(s) will be used in practice, which may include additional design and  
441 labelling considerations<sup>9</sup> (as further described in Section 6 of this document), additional human  
442 oversight requirements, possible workflow challenges, and the security of the model and its data once  
443 deployed.

---

<sup>5</sup> IMDRF/AIML WG/N88 FINAL:2025: *Good Machine Learning Practice (GMLP) for Medical Device Development: Guiding Principles*

<sup>6</sup> See GMLP Guiding Principle 3.

<sup>7</sup> See GMLP Guiding Principle 4.

<sup>8</sup> See GMLP Guiding Principle 2.

<sup>9</sup> See GMLP Guiding Principle 9.

444 *Validation and Clinical Evaluation Needs*: The term validation has been used to represent different  
445 concepts within the fields of medical device development and AI model development<sup>10</sup>. Validation, in  
446 this context and aligned with IMDRF/AIMD WG/N67 (Edition 1):2022, means “confirmation by  
447 examination and provision of objective evidence that the particular requirements for a specific  
448 intended use can be consistently fulfilled”. In this step, it is essential to establish preliminary validation  
449 strategies and appropriate evaluation metrics (see Appendix B), including specifying acceptance  
450 criteria and thresholds for the chosen metrics as well as statistical analysis methods, which can  
451 impact how data is collected and used.

452 *Post-Market Monitoring and Surveillance*: Requirements planning for continuous monitoring of the  
453 ongoing performance of the AI-enabled medical device in the clinical workflow is crucial to ensure  
454 patient safety and device effectiveness. It is important for manufacturers to begin making preliminary  
455 plans for post-market monitoring and surveillance for both the AI-enabled medical device and its  
456 model(s) during this step. These plans can address the unique challenges of AI-enabled medical  
457 devices, such as performance degradation due to drift, and enable real-world performance evaluation.  
458 Regulatory requirements in the particular jurisdiction(s) in which the AI-enabled medical device will be  
459 deployed may help to scope what might be needed for this activity.

#### 460 **Additional Standards and References to be Considered**

- 461 • ISO/IEC 23053:2022 *Framework for Artificial Intelligence (AI) Systems Using Machine*  
462 *Learning (ML)*

## 463 **5.2. Data Collection and Management**

464 Fit-for-purpose, representative<sup>11</sup> data is crucial for the appropriate design and development of AI  
465 models in medical devices, while also important for minimizing risks such as unwanted bias. While  
466 there is no shortage of recognized industry best practices for data collection and management, this  
467 step covers considerations for appropriate data collection and management practices necessary for  
468 training, tuning, and validation of AI models that can help support the model’s safety and  
469 performance.

470 It is important for manufacturers to review the legal and regulatory requirements in the jurisdiction(s)  
471 in which data is created, collected, analyzed, secured, and stored (such as data privacy requirements  
472 and consent as well as requirements on data collection and data management) and follow good  
473 clinical practices<sup>12</sup>, which are not extensively covered in this document. Furthermore, as it is common  
474 to outsource activities associated with data acquisition, tools, and services, it is important for  
475 manufacturers to consider recommendations in section 7.6 *Managing Outsourced Processes,*  
476 *Activities and Products* in IMDRF/SaMD WG/N23 FINAL:2015 for these activities.

477 The rest of this section highlights key considerations relevant to data collection and management for  
478 AI models and AI-enabled medical devices. The recommendations included in this section can be  
479 applied across the AI-enabled medical device life cycle wherever data is being collected, used, and  
480 stored. This includes the processes that occur during the Verification and Validation life cycle step, for  
481 example, as well as for data collected in the post-market setting to maintain, monitor, or further train  
482 and improve model performance.

---

<sup>10</sup> Per IMDRF/AIMD WG/N67 (Edition 1):2022 *Machine Learning-enabled Medical Devices: Key Terms and Definitions*, “MLMD manufacturers, regulators, and users should be aware of the conflicting interpretations of the term validation and ensure that communication regarding the development phases and the associated datasets is clear to avoid confusion between data validation, ML model tuning, and medical device validation. Alternatively, the use of the term validation that refers to the training and tuning process should be avoided in the context of medical device development. It is recommended that the use of the term ‘validation’ be accompanied by the context when referring to ML model tuning, data curation, and the associated datasets”.

<sup>11</sup> See GMLP Guiding Principle 3.

<sup>12</sup> See, for example, ISO 14155:2020 *Clinical investigation of medical devices for human subjects - Good clinical practice*

483 *Data Suitability and Collection/Generation:* While this may be generally addressed in the Planning and  
484 Design step, it is important for manufacturers to consider data suitability when further investigating  
485 sources of data that are readily available or need to be generated. It is important for data to align with  
486 the model's particular requirements, including input, intended patient population, disease state, and  
487 reference standard<sup>13</sup> and account for what may be needed to ensure adequate transparency to the  
488 end user (for more information on transparency, see Section 6 Transparency and Labelling). For  
489 example, for some use cases, age and sex data may not be needed to train the model but can help  
490 users understand performance across these demographics. When collecting and/or generating data,  
491 clearly defining inclusion and exclusion criteria is a critical activity to complete prior to collection or  
492 analysis (i.e., in the case of a retrospective data analysis). If multiple data sources will be used, for  
493 example a combination of a data repository and data from a prospective clinical study, it is important  
494 that there is consistency across these datasets (including, for example, data pre-processing), where  
495 possible, to ensure discrepancies will not cause issues in training or tuning.

496 *Data Augmentation and Use of Synthetic/Simulated Data:* As AI models become more complex and  
497 challenges arise with accessing data, manufacturers may use various methods, including data  
498 augmentation and synthetic/simulated data, to incorporate new and modified data into their datasets.  
499 It is important to note that the acceptability of these approaches depends on whether the jurisdiction  
500 in which the model and device are developed, validated, and deployed allows for such uses. Data  
501 augmentation techniques may be used in this content to modify existing datasets by applying  
502 transformations that preserve the underlying characteristics while increasing dataset size and  
503 diversity. For example, techniques such as applying controlled rotations, brightness adjustments, and  
504 geometric transformations to X-ray or MRI images can be used to increase dataset diversity while  
505 preserving clinically relevant features. Manufacturers may also seek to use synthetic or simulated  
506 data to supplement actual, human data for model development and validation. Use of synthetic or  
507 simulated data can help address ethical and privacy concerns and accelerate data generation,  
508 particularly for underrepresented populations such as patients with rare diseases.

509 Because data that is augmented and synthetic or simulated data may not accurately reflect the  
510 nuances and complexity of data that may be encountered in the real-world use of the model and  
511 device, using such data in AI model development and validation may introduce uncertainty, especially  
512 about whether synthetic data can serve as an adequate surrogate for human data in validation  
513 studies. In light of these concerns, it is important for manufacturers to carefully document how data  
514 was augmented or synthetically generated (considering both data provenance and transformation  
515 methods), how it is used, and document justification that the data and approach to its use is fit-for-  
516 purpose. Such documentation and justification can be particularly helpful to provide to regulatory  
517 authorities to help provide clarity on the appropriateness and risk-based approach taken to the use of  
518 such data. In addition, manufacturers can consider tagging such data as such at the point of  
519 generation to aid in traceability between, for example, synthetic and human-derived data sources.

520 *Data Representativeness and Bias Mitigation:* It is important that data collected and used for training,  
521 tuning and validation represents the full spectrum of data that the model is intended to encounter in  
522 clinical use in order to ensure that the data is fit-for-purpose. Furthermore, it is important that this data  
523 represents not only clinical subtypes, such as demographics of the intended use population and  
524 disease conditions (both cases and controls) but also nonclinical subtypes, such as data acquisition  
525 equipment and their parameters and collection sites relevant to the device's intended purpose.

---

<sup>13</sup> See GMLP Guiding Principle 5.

526 In addition to clinical and non-clinical subtypes, it is important for manufacturers to consider how  
527 under-represented and priority populations within the intended use population (e.g., specific ethnic  
528 groups, Indigenous populations, or rural and remote communities) are reflected in the data used for  
529 training, tuning, and validation. Generally, where adequate representation cannot reasonably be  
530 achieved, it is important to transparently document and relay these limitations to users, including  
531 potential impacts on performance and generalizability, and to implement appropriate mitigation  
532 measures such as additional local validation, cautious labelling, or risk-proportionate bridging  
533 evidence. Collecting data from multiple sources to represent the intended use population(s),  
534 environment(s), and context(s)-of-use may be a solution to potentially reduce under-representation as  
535 well as improve representativeness across other jurisdictions whilst reducing the need for subsequent  
536 bridging studies in this regard. However, it is important for manufacturers to evaluate different data  
537 sources to address specific potential risks of unwanted biases and confounding factors. For example,  
538 data from electronic health records (EHRs) recording routine health delivery may represent clinical  
539 decisions that take into consideration various factors related to a particular patient encounter such as  
540 cost and affordability issues or patient preferences. Hence, this data may not generalize to all health  
541 care settings.

542 Manufacturers can consider sampling strategies to ensure all relevant subgroups are appropriately  
543 represented in the dataset as well as to address potential algorithmic biases that may arise from use  
544 of specific datasets (e.g., sampling bias, selection bias, measurement bias, algorithm bias) through  
545 careful sampling, stratification, and bias mitigation techniques. When developing adaptive AI models  
546 that continue to learn after their initial deployment, it is important for manufacturers to eliminate or  
547 reduce as far as possible the risk of biased outputs influencing input for future operations (feedback  
548 loops), and to ensure that any such feedback loops are addressed with appropriate mitigation  
549 measures.

550 *Data Cleaning/Quality Assurance:* Data Cleaning, or the process of correcting and/or deleting  
551 incomplete, incorrect, or irrelevant records from a database or table,<sup>14</sup> may be needed to ensure data  
552 collected is fit-for-use. Trained professionals (e.g., domain experts, data scientists) and/or those with  
553 appropriate expertise are best suited to perform data cleaning to minimize introducing unwanted bias  
554 into the dataset. Furthermore, it is important that cleaned data be representative of the data the model  
555 will encounter in its intended use environment and, where applicable, that any anomalies or  
556 discrepancies identified during cleaning tasks inform risk control strategies as similar anomalies may  
557 be encountered with model inputs after deployment. Other key activities include prospectively defining  
558 data cleaning procedures with contextually relevant stop thresholds, setting limits on correcting errors  
559 (to be done to the extent possible without infringing on privacy), preventing duplication of data when  
560 aggregating from multiple sources, and maintaining documentation of all data cleaning and quality  
561 assurance steps to ensure reproducibility and traceability. Manufacturers may choose to implement  
562 automated processes for data extraction, transformation, and loading (ETL) to maintain data integrity  
563 while supporting scalable operations. This can include, for example, real-time data quality monitoring  
564 with automated anomaly detection to enable rapid identification of issues that could compromise  
565 model performance and patient safety. It is also important for manufacturers to validate data  
566 processing pipelines to ensure consistent and reproducible results.

567 *Data Lineage and Provenance Documentation:* Building upon the documentation recommendations  
568 mentioned for data cleaning, it is important for manufacturers to establish comprehensive data lineage  
569 and provenance tracking throughout the AI-enabled medical device life cycle to support traceability  
570 and remediation of data quality issues. Data lineage documents the data origin, transformations, and  
571 movement over time, providing visibility and providing manufacturers the ability to trace errors in  
572 device performance or other issues that occur in the post-market setting back to potential root causes.  
573 Data provenance documents the inputs, entities, systems, and processes involved in creating the  
574 data, providing a historical record of data origins.

---

<sup>14</sup> Artificial Intelligence-Data quality for analytics and machine learning (ML) – Part 4: Data quality process framework (ISO/IEC 5259-4)

575 *Data Retention, Retirement and Life Cycle Management:* Effective data life cycle management spans  
576 from initial collection through active use, archival storage, and final disposal, including well-defined  
577 criteria and processes for data retirement. Data use tracking mechanisms help monitor access  
578 patterns and identify potential overuse, which can warrant earlier retirement. Depending on the type of  
579 data, there may be distinct life cycle paths based on sensitivity, ownership, and intended use. While it  
580 is important for data to be retained in order to allow manufacturers to investigate adverse events,  
581 analyze device performance patterns, trace the source of malfunctions or safety concerns, and  
582 provide evidence for regulatory inquiries or corrective actions, retention periods typically reflect a  
583 balance of regulatory and legal requirements, clinical or research utility, and privacy obligations, as  
584 applicable to a particular jurisdiction. Robust governance frameworks that include documented  
585 policies, designated data stewards, and automated monitoring ensure consistent application of  
586 retention and disposal rules.

#### 587 **Additional Standards and References to be Considered**

- 588 • ISO/IEC 5259-4:2024 *Artificial intelligence — Data quality for analytics and machine learning*  
589 *(ML) — Part 4: Data quality process framework*
- 590 • ISO/IEC 23053:2022 *Framework for Artificial Intelligence (AI) Systems Using Machine*  
591 *Learning (ML)*

### 592 **5.3. Model Building and Tuning**

593 In this next step of the life cycle, AI model(s) are developed and refined using the collected data,  
594 including selection of appropriate architecture, feature engineering, and performance optimization.

595 GMLP Guiding Principle 6 highlights that model choice and design are tailored to both the available  
596 data and the device's intended use, with design decisions evaluated to actively mitigate known risks  
597 like overfitting and performance degradation while supporting clinically meaningful performance goals.  
598 The rest of this section provides additional considerations relevant to model building and tuning for  
599 manufacturers to consider during this stage of development.

600 *Model Design and Architecture Selection:* Manufacturers are best positioned to make decisions when  
601 developing and tuning models, particularly for AI-enabled medical devices. As introduced in Section  
602 5.1, design choices including feature selection, model type, and evaluation metrics, can inadvertently  
603 introduce or reinforce risks, such as inaccurate or biased outputs as described in Section 4.2.  
604 Avoiding overly complex models, e.g., limiting the number of parameters, layers, or features, may be  
605 optimal for addressing such risks, ensuring the model is fit-for-use and potentially reducing the  
606 amount of data needed during this step. Justification for model choices and rationale for decisions  
607 made with respect to the intended use of the device can not only support regulatory submissions but  
608 also help ensure that the device is and continues to be well-aligned in the clinical context and  
609 workflow across the life cycle. For instance, a developer working on a medical diagnostic AI-enabled  
610 medical device might explain their choice of a convolutional neural network (CNN) over a random  
611 forest model due to the CNN's superior performance in image recognition tasks, which is crucial for  
612 accurately identifying anomalies in medical scans.

613 *Model Explainability and Interpretability:* Building upon considerations for model selection described in  
614 Section 5.1, when choosing to employ a certain model, it is important for manufacturers to consider  
615 the explainability and interpretability of the model, including as it relates to the functionality of the AI-  
616 enabled medical device in which it is incorporated. Models that are explainable can enable users to  
617 understand their underlying reasoning and limitations and can also support the ability of the  
618 manufacturer to create an appropriate evaluation plan. However, model explainability may have  
619 potential trade-offs with complexity and performance. As such, manufacturers can consider the trade-  
620 offs between the potential risks and benefits of unexplainable models in the context of the clinical use,  
621 including the environment, the characteristics of the users, and the system into which the model will  
622 be deployed. There may also be a higher expected level of clinical evaluation for less explainable  
623 models to gain clinician trust for real-world adoption. By considering these topics during this step,  
624 manufacturers can avoid identifying challenges and risks later in development when it is difficult to  
625 address the underlying structure that might be leading to risks or reduced trust from the clinical  
626 community. A model's interpretability can also be considered to ensure that users will be able to  
627 comprehend how the model arrives at its outputs or recommendations.

628 *Data and Feature Preprocessing:* Data preprocessing is a critical step for optimizing data prior to  
629 model training. It includes activities such as data cleaning (described in Section 5.22 above),  
630 normalization, and transformation to ensure the data is consistent, reliable, and suitable for modelling.  
631 Regarding feature preprocessing, while not exhaustive, both feature engineering and selection as well  
632 as dimensionality reduction<sup>15</sup> are some of the many approaches that can be employed. When feature  
633 engineering and selection is used, it is beneficial for features within the raw data to be chosen with  
634 considerations for biological and scientific feasibility and evidence because these will inform clinical  
635 evaluation and explainability in terms that clinicians and patients can understand, providing a "valid  
636 clinical association"<sup>16</sup> as described in IMDRF/SaMD WG/ N41 Final:2017 *Software as a Medical  
637 Device (SaMD): Clinical Evaluation*. In the absence of full scientific evidence, it is possible that other  
638 steps of the life cycle may compensate for the lack of clarity in features selected, such as additional  
639 validation on a wide variety of datasets, which could be used to ensure the model performs as  
640 expected across the intended use population. Additionally, choosing features that are generalizable  
641 can help manufacturers avoid introducing unwanted bias for particular demographic groups and/or  
642 subgroups, including both clinical and non-clinical subtypes. For example, when choosing features for  
643 an AI-enabled medical device for predicting the risk of pre-term birth, manufacturers may want to  
644 select features that are biologically relevant and generalizable across populations, including factors  
645 such as previous pre-term births, maternal age, body mass index (BMI), and presence of certain  
646 infections. Finally, regardless of the approach to feature preprocessing used, it is important that  
647 manufacturers ensure the decision-making process of choosing the approach is transparent and  
648 justified, which can inform evidence for future regulatory inquiries and transparency for users.

649 *Selection of Evaluation Metrics:* It is important to select appropriate metrics to ensure reliable model  
650 performance. A loss function is the difference in error between the model's predictions and target (or  
651 reference) values and can be used during model training and building to guide adjustments to the  
652 model's parameters, such as weights and hyperparameters. An evaluation metric is used to evaluate  
653 a model's performance during tuning and can be used throughout the life cycle (when feasible) to  
654 ensure that the model is continuously meeting performance goals. Appropriate loss function and  
655 evaluation metrics chosen for model training, tuning and validation are relevant to the intended use  
656 and are generalizable to the target patient population. Selection of generalizable evaluation metrics  
657 may be particularly important when features could not be selected based on biological or scientific  
658 evidence. Examples of common evaluation metrics can be found in Appendix B.

---

<sup>15</sup> Dimensionality reduction helps to improve computational efficiency, reduce noise in the data, mitigate overfitting, and enhance interpretability of the model and is another option when performing data preprocessing.

<sup>16</sup> Per IMDRF/SaMD WG/ N41Final:2017 *Software as a Medical Device (SaMD): Clinical Evaluation*, "valid clinical association, also known as scientific validity, is used to refer to the extent to which the SaMD's output (concept, conclusion, measurements) is clinically accepted or well-founded (based on an established scientific framework or body of evidence), and corresponds accurately in the real world to the healthcare situation and condition identified in the SaMD definition statement.

659 *Deployment Considerations:* Specific requirements of the final deployment environment may result in  
660 considerations at the model tuning stage to address risks such as infrastructure inadequacy and  
661 computational and scalability constraints, as described in Section 4.2.4. For example, some clinical  
662 applications may pose hardware resource constraints or require fast inference times that limit the  
663 model design options. Such constraints may require, for example, compression/distillation of the  
664 initially designed model, creating smaller, more cost-efficient models by transferring knowledge from  
665 the initial, complex models. It is important that decisions regarding trade-offs inherent to addressing  
666 such design considerations be documented as part of the risk/benefit analysis.

667 *Leveraging general-purpose and off-the-shelf models:* If manufacturers choose to incorporate or  
668 leverage third-party, general-purpose models, including large language models (LLMs), other  
669 architectures or foundation models, into their AI-enabled medical devices, there are several factors to  
670 consider. For example, as with all medical devices, it is important to implement robust supplier  
671 management practices and follow industry best-practices for software of unknown provenance  
672 (SOUP)<sup>17</sup> and third-party oversight. Other key considerations include assessing the third-party  
673 supplier credibility, model provenance, model's longevity, and life cycle management processes,  
674 including the model's version history, update cadence, and support commitments. It is important for  
675 manufacturers to also assess and document the model's known limitations, potential biases, and  
676 uncertainties. When applicable, additional considerations include understanding the generative  
677 behaviour of the model, such as variability in outputs, susceptibility to hallucination, and dependence  
678 on input prompts or context.

679 When thinking about incorporating these types of models, it is important for manufacturers to evaluate  
680 feasibility of managing model and documentation updates, evaluate available information about the  
681 model while understanding how information gaps may impact safety, and implement appropriate risk  
682 controls, such as human oversight and guardrails, as well as those that may be needed for  
683 appropriate deployment of such models. Even when information is available on the training data set,  
684 manufacturers may find that general-purpose models are typically not trained on data that is  
685 representative of the intended use population or clinical use case. Manufacturers can consider using  
686 methods such as transfer learning and fine tuning as part of a comprehensive risk strategy to address  
687 these representativeness gaps. For generative AI-enabled models, it is also important to monitor for  
688 emergent behaviours or unintended content generation over time as the model or its environment  
689 changes.

690 Providing transparent information about the use of general-purpose and off-the-shelf models is  
691 important for both regulatory authorities and users to understand the foundation, reliability, and  
692 limitations of an AI-enabled device that incorporates these types of models. See Section 6 for more  
693 information on transparency and labelling for AI-enabled medical devices.

#### 694 **Additional Standards and References to be Considered**

- 695 • Section 8.3 Development in IMDRF/SaMD WG/N23 FINAL: 2015
- 696 • ISO/IEC 23053:2022 *Framework for Artificial Intelligence (AI) Systems Using Machine*  
697 *Learning (ML)*
- 698 • ISO/IEC TS 25058:2024 Systems and software engineering — Systems and software Quality  
699 Requirements and Evaluation (SQuaRE) — Guidance for quality evaluation of artificial  
700 intelligence (AI) systems
- 701 • ISO/IEC 25059:2023 Software engineering — Systems and software Quality Requirements  
702 and Evaluation (SQuaRE) — Quality model for AI systems

---

<sup>17</sup> IEC 62304:2015 defined the term as follows: “A software item that is already developed and generally available and that has not been developed for the purpose of being incorporated into the medical device (also known as “off- the-shelf software”) or software previously developed for which adequate records of the development processes are not available. NOTE A MEDICAL DEVICE SOFTWARE SYSTEM in itself cannot be claimed to be SOUP”

- 703 • IEEE 2941-2021 IEEE Standard for Artificial Intelligence (AI) Model Representation,  
704 Compression, Distribution, and Management

## 705 **5.4. Verification and Validation, including Clinical Evaluation**

706 As with all medical devices, verification and validation (also referred to as “V&V”) in the AI-enabled  
707 medical device life cycle builds trust in the performance of the device and establishes the baseline  
708 assurance that the device has met its design requirements and can be used safely and effectively for  
709 its intended purpose.

710 At this point in the life cycle, the model(s) have been built and tuned. Evaluation metrics (described in  
711 Section 5.3 with examples included in Appendix B) used as part of tuning can continue to establish  
712 the model and AI-enabled medical device’s performance. It is important to note that same underlying  
713 model can be used to produce multiple types of outputs, which may be used, understood and  
714 therefore, evaluated in different ways. An AI model’s outputs may include, for example, continuous  
715 numerical outputs (such as a thermometer that produces an estimate of core body temperature in  
716 degrees Celsius), categorical outputs (such as a thermometer that produces a result of “subnormal,”  
717 “normal,” or “fever”) or binary outputs (such as a thermometer that produces a result of “positive” or  
718 “negative” for fever). The model’s outputs will have a significant influence on the approach used in  
719 V&V processes to adequately evaluate the safety and effectiveness of an AI-enabled medical device.  
720 Furthermore, typically, models themselves are not deployed but are integrated into a device (i.e., an  
721 AI-enabled medical device) or system (e.g., hosted in cloud server) with additional functions that may  
722 impact the context or utility of the model outputs and will also impact what V&V strategies will be most  
723 appropriate.

724 This section will address considerations for V&V activities related to the model (Section 5.4.1),  
725 including external validation of the model, as well as considerations for V&V activities related to the  
726 AI-enabled medical device, including clinical evaluation<sup>18</sup> (Section 5.4.2). It is important to note that  
727 clinical evaluation distinguishes itself from external validation of the model, where clinical evaluation  
728 covers the entire medical device, including the intended user and workflow and not just the  
729 generalizability of the AI model(s).

730 It is also important for manufacturers to conduct verification and validation activities for an AI-enabled  
731 device in accordance with applicable design control processes for medical devices, (e.g., ISO  
732 13485:2016). As part of these verification and validation activities, manufacturers may want to  
733 consider the recommendations in section 8.0 (SaMD Realization and Use Processes) of Software as  
734 a Medical Device (SaMD): Application of Quality Management System (IMDRF/SaMD WG/N23  
735 FINAL:2015), as well as the definition of validation provided in Machine Learning-enabled Medical  
736 Devices: Key Terms and Definitions (IMDRF/AIMD WG/N67 (Edition 1):2022). Furthermore, IMDRF’s  
737 GMLP principles provide some common principles that can be considered when validating medical  
738 devices that incorporate AI, such as using validation data that is independent from the data that was  
739 used to develop the model and is representative of the intended use population<sup>19</sup>.

740 Documentation throughout the V&V processes is essential to track and monitor the model’s strengths,  
741 limitations, and impact. This recordkeeping can help develop comprehensive and transparent  
742 documentation for the AI-enabled medical device and its model(s) so that clinicians and patients can  
743 understand the V&V outcomes and make informed choices regarding AI-enabled medical device  
744 deployment and use. See Section 6 for more information on transparency and labelling for AI-enabled  
745 medical devices.

---

<sup>18</sup> Per IMDRF/SaMD WG/N41FINAL:2017 *Software as a Medical Device (SaMD): Clinical Evaluation*, “clinical evaluation is a systematic and planned process to continuously generate, collect, analyze, and assess the clinical data pertaining to a SaMD in order to generate clinical evidence verifying the clinical association and the performance metrics of a SaMD when used as intended by the manufacturer” while “clinical validation measures the ability of a SaMD to yield a clinically meaningful output associated to the target use of SaMD output in the target health care situation or condition identified in the SaMD definition statement”.

<sup>19</sup> See GMLP Guiding Principle’s 3 and 4.

#### 746 **5.4.1. V&V Activities related to the Model**

747 *External Validation of the Model:* Prior to widespread clinical evaluation and deployment, it is  
748 important to evaluate whether the model's performance generalizes beyond the training environment  
749 and remains performant in real-world clinical applications. This involves, for example, validating the  
750 model on a dataset that is separate from training data, and assessing the performance, accuracy, and  
751 reliability of the finalized AI model. It is also beneficial for external validation to address the socio-  
752 technical environment where the AI model will be deployed, as emphasized in the SaMD quality  
753 management framework, which requires "thorough understanding of the socio-technical environment  
754 (clinical perspective), and the technology and system environment (software perspective)" to prevent  
755 "incorrect, inaccurate, and/or delayed diagnoses and treatments" (IMDRF/SaMD WG/N23  
756 FINAL:2015).

757 *Model Robustness:* Model robustness assesses the model's ability to continue to perform in  
758 conditions outside of the ideal use condition. Techniques such as stress testing and sensitivity  
759 analysis can be used to evaluate the model's resilience against variations in input data, noise,  
760 outliers, and adversarial attacks. Identifying such potential vulnerabilities early may allow for proactive  
761 adjustments that strengthen the model's stability and ensures that identified operational bounds and  
762 limitations can be conveyed to users to control for risks associated with improper or unsupported use.  
763 Additionally, red teaming, a process where a team intentionally probes the model for weaknesses and  
764 exploits potential vulnerabilities, is a valuable method for identifying issues and validating  
765 implemented risk controls and failure/error handling mechanisms. Demonstrating model robustness  
766 may allow for a better understanding of the risks related to device performance in edge cases and  
767 whether adequate risk controls are in place.

#### 768 **5.4.2. V&V Activities related to the AI-enabled Medical Device**

769 *Verification of the AI-enabled Medical Device:* Once a manufacturer externally validates a model, it  
770 may be incorporated into an AI-enabled medical device and verified to ensure that requirements are  
771 met. In general, this can follow typical practices for medical device software. Additional verification  
772 activities may need to be performed according to medical device requirements.

773 *Clinical Evaluation Study Design:* It is important for manufacturers to ensure that the approach taken  
774 to clinically evaluate an AI-enabled medical device is tailored to its intended use and unique  
775 characteristics, including the intended population. As introduced above, clinical evaluation activities  
776 can depend on the outputs of the model and whether or not the model serves as the medical device's  
777 primary intended use or purpose or facilitates other device functionality. For example, in diagnostic  
778 medical devices like a radiology computer-aided detection or diagnosis medical device software, the  
779 model's output may be a prediction of whether and where evidence of cancer appears in an image.  
780 Diagnostic medical device evaluations may therefore focus specifically on determining whether the  
781 output of the model is accurate and clinically useful based on metrics like positive and negative  
782 percent agreement or predictive value. In contrast, for a therapeutic medical device such as an  
783 automated insulin delivery system, a model may use sensor and behaviour data (e.g., carbohydrates  
784 eaten, exercise planned) to estimate and predict a user's current and future blood glucose values,  
785 which may be used by another device function to determine how much insulin should be delivered to  
786 the user. Therefore, therapeutic device evaluations may focus more broadly on the therapeutic  
787 effectiveness of the overall device, including determining whether outcome-related measures were  
788 improved. Furthermore, the usability and evidence supporting the output, such as transparency on  
789 why a given output was provided, may also affect clinical outcomes and may be included in clinical  
790 evaluations where relevant. This is an especially important consideration for non-autonomous AI-  
791 enabled medical devices that function as decision support and rely on the user to make the final  
792 therapeutic, diagnostic or other clinical decision. When patient populations, disease characteristics  
793 and medical practices have similar documented characteristics across jurisdictions, it may in some  
794 cases be appropriate to extrapolate the clinical evaluation information to inform decisions in another  
795 appropriate patient population. High quality clinical evaluations that are consistent with the GMLP  
796 principles can reduce duplicative efforts, thereby further promoting and enabling regulatory reliance.

797 Finally, it is critical for the clinical evaluation study to consider the intended use and risk-level of the  
798 AI-enabled medical device with respect to the associated range of potential, appropriate strategies for  
799 validating its use, including from robust retrospective studies with suitable reference standards, silent  
800 mode studies, or simulation-based evaluations to prospective studies in the intended use  
801 environment. Where the device is intended to support or influence human clinical decision-making, it  
802 is important for the manufacturer to consider study designs that explicitly assess the performance and  
803 safety of the human-AI team<sup>20</sup>.

804 *Usability and Human Factors:* To ensure that the AI-enabled medical device is effective and usable in  
805 its intended environment, it is important for manufacturers to consider what is needed for human  
806 factors and usability validation. This involves assessing whether the device aligns with the skills,  
807 expertise, and needs of its users, ensuring that they can correctly interpret and use its outputs and  
808 validate risk controls that are in place to address risks identified in Sections 4.2.1 and 4.2.2. Usability  
809 validation can involve structured strategies, such as silent testing, which evaluates how well the  
810 device integrates into real-world workflows while blinding clinicians to the outputs of the model so that  
811 the existing workflow is not disrupted. It is also essential to test the device's usability in the context of  
812 its intended purpose, validating that users can understand its inputs and outputs, recognize its  
813 limitations, and use it effectively for decision-making. This may include, for example, testing the ability  
814 of users to understand labelling, as further described in Section 6. It is also important for  
815 manufacturers to consider reasonably foreseeable misuse and whether any potential user errors  
816 could lead to unintended consequences.

817 *Determining adequacy of measured performance:* After conducting the V&V described above, the  
818 next step is to decide whether the measured performance is adequate. Judgements about whether or  
819 not medical device performance in V&V activities (or any other performance) is acceptable are  
820 influenced by several factors:

- 821 • the nature of the condition of interest, including whether the condition is relatively minor, is  
822 potentially deadly or debilitating, or is somewhere in between;
- 823 • the clinical workflow that the AI-enabled medical device is intended to be used within,  
824 particularly whether it functions within a system that includes prior screening or subsequent  
825 confirmations that can catch potential errors; and
- 826 • the usefulness, importance, and autonomy of the AI-enabled medical device outputs in the  
827 clinical environment, especially regarding whether a specific clinical action or intervention can  
828 be initiated based on the result and whether such interventions occur automatically without  
829 additional clinical oversight.

830 These judgements can become even more complex when an AI-enabled medical device is intended  
831 to be used repeatedly for the same individual over time. It is important to consider the limitations,  
832 uncertainty, and performance (including acceptance criteria) and whether they are acceptable for a  
833 given clinical use. Manufacturers may benefit from crafting a clear justification to both articulate to  
834 users on the benefits and risks of an AI-enabled medical device and as part of regulatory approval on  
835 why the V&V activities and outcomes are relevant to the device's intended use.

#### 836 **Additional Standards and References to be Considered**

- 837 • ISO/IEC 23053:2022 *Framework for Artificial Intelligence (AI) Systems Using Machine*  
838 *Learning (ML)*

---

<sup>20</sup> See GMLP Guiding Principle 7.

## 839 5.5. Deployment

840 Deployment describes the process of integrating the AI-enabled medical device into its intended use  
841 environment, including “delivery, installation, setup, and configuration that support a controlled and  
842 effective distribution”<sup>21</sup> of the device, making it accessible for real-world use. This step demarks the  
843 transition into post-market and encompasses the processes between regulatory market authorization  
844 and first use of the device.

845 There are standards and resources that manufacturers can consider when planning for deployment  
846 requirements and processes, such as those described in ISO 13485, which emphasizes that  
847 deployment planning must be documented, risk-based, and include provisions for training, installation  
848 verification, and post-delivery activities to ensure safe and effective use of medical devices. It is  
849 important to note that the deployment process is highly context-specific to the device’s intended use.  
850 For example, an AI-enabled medical device may be deployed in a consumer product for direct use by  
851 patients and caregivers or, conversely, it may be deployed into a complex hospital Information  
852 Technology (IT) network and interact with multiple stakeholders. Additionally, AI-enabled medical  
853 device performance can be significantly impacted by variations in IT infrastructure (e.g. scanner  
854 models and operating systems), patient demographics and local health policies, requiring deployment  
855 to be managed carefully and proportionate to the risks of the intended use.

856 It is important for manufacturers to review the relevant legal and regulatory requirements, industry  
857 standards, and organizational policies in the jurisdictions and sites in which the device is being  
858 deployed. Such considerations could include data privacy, security, and ethical considerations to  
859 maintain trust and mitigate risks in the deployment stage, which are not extensively covered in this  
860 document.

861 This section focuses on deployment considerations that are unique to or particularly important for AI-  
862 enabled medical devices.

863 *Planning and Infrastructure:* As part of the Planning and Design life cycle step as well as Model  
864 Building and Tuning, deployment is planned and infrastructure needed for deployment, operations  
865 and monitoring and real-world performance evaluation is designed (see Sections 5.1 and 5.3). The  
866 deployment step is meant to ensure that assumptions made about the device and its infrastructure  
867 during development were accurate in the real-world environment and provides an opportunity for  
868 adjustments to be made so that the device can be used safely and effectively. For instance, when a  
869 device is designed to integrate with an EHR system, deployment involves establishing and verifying  
870 this connection and, if data fields are labelled differently than anticipated during development, the  
871 system can be reconfigured to properly map inputs per the device’s requirements. This can help to  
872 control for risks, for example those outlined in Section 4.2.4 related to interoperability issues.

873 As AI-enabled medical devices can have different failure modes than other medical devices and  
874 medical device software, it is important for this plan and associated infrastructure to include ways to  
875 determine whether any issues that arise can be attributed to AI model(s) or to medical device  
876 software generally, so that manufacturers can properly evaluate a failure and address it adequately.  
877 Additionally, it is important for the infrastructure to have mechanisms in place to determine whether  
878 the quality and quantity of inputs to and data requirements of the AI-enabled medical device meet  
879 specifications and account for variation across sites, enabling the ability to monitor for data drift and  
880 performance degradation, potentially through Human Oversight features (as described in Section 4.3).  
881 Furthermore, where applicable, it is important for the infrastructure to ensure that the clinician retains  
882 autonomy over the clinical decision-making including the ability to override AI-enabled medical device  
883 and its model(s)’s output where needed. If there is a plan to use real-world data to improve the  
884 model(s) in the future, it is important for that to be considered when building and implementing the  
885 infrastructure and principles and practices from the Data Collection and Management life cycle step  
886 can be employed. Finally, it is essential for the deployment plan and infrastructure to account for  
887 cybersecurity concerns detailed in Section 4.4.

---

<sup>21</sup> IMDRF/SaMD WG/N23 FINAL:2015

888 *Deployment Verification and Validation:* Deployment verification and validation ensures that the  
889 deployment plan was executed correctly and any anomalies that may have occurred during  
890 deployment are appropriately addressed before widespread use of the AI-enabled medical device.  
891 This can include ensuring that training for site staff and healthcare professionals was adequate and  
892 appropriate for the environment and users. Depending on the risk of the device, deployment V&V  
893 activities can range in complexity and timespan. For example, utilizing a combination of site-specific  
894 retrospective data and control data may be sufficient to verify that the deployment process was  
895 successful. However, complex workflows, which may contain highly variable data inputs, can require  
896 more longitudinal activities such as trial periods with sufficient human oversight to validate the  
897 intended purpose of the AI-enabled medical device is being achieved.

898 *Deployment Approach:* Defining the implementation approach in deployment plans is a critical step,  
899 whether that is deployment in phases, such as a preview release followed by a final release, or  
900 through phased release across multiple sites. A phased release across multiple sites may be  
901 beneficial for AI-enabled medical devices because it enables manufacturers to identify potential safety  
902 issues, performance variations, or integration challenges in a controlled manner before full-scale  
903 deployment, allowing them to understand the AI-enabled medical device's performance across patient  
904 populations, clinical workflows, and technical environments that may differ from the controlled  
905 conditions used during V&V.

906 *Site-Specific Customization and Localization*<sup>22</sup>: AI-enabled medical devices may have the ability to be  
907 customized and/or localized to achieve expected performance levels at a particular site. Depending  
908 on the risks and context of the deployment environment, various approaches may be utilized. For  
909 example, parameters may need to be adjusted from a range of previously validated options until the  
910 most appropriate configuration for a deployment site is found (customization). Alternatively, or  
911 potentially in addition, local data may be used to retrain the AI model contained within an AI-enabled  
912 medical device to more appropriately align to site-specific criteria (localization). For devices that are  
913 localized, deployment includes localization steps and the development of training and labelling  
914 materials specific to each deployment. Predetermined Change Control Plans (PCCPs) can be  
915 leveraged and applied to AI-enabled medical devices as a way for regulatory authorities within a  
916 specific jurisdiction to prospectively review and approve certain changes to an AI-enabled medical  
917 device, including customization and localization, where such an approach is available.

918 *Traceability and Version Control:* For AI-enabled medical devices, especially those that are  
919 customized or localized as described above, it is important for traceability and version control  
920 mechanisms to be in place so that manufacturers have the ability to identify failure events and root  
921 causes and ultimately intervene as necessary to minimize risk to patients, such as those identified in  
922 Section 4.2.4 Risks Related to Deployment and Post-Marketing Monitoring and Performance. Such  
923 interventions may include retraining, which can be enabled by a PCCP, or rolling back to prior  
924 versions. Unique device identifiers (UDI), as recommended in IMDRF/UDI (Unique Device Identifier)  
925 WG/N48 FINAL:2019 *Unique Device Identification system (UDI system) Application Guide*, can be  
926 employed for software/model version traceability, where available. If UDIs are not applicable in a  
927 particular jurisdiction, manufacturers may want to consider alternative auditable identifiers.

## 928 **Additional Standards and References to be Considered**

- 929 • IMDRF/SaMD WG/N23 FINAL:2015 *Software as a Medical Device (SaMD): Application of*  
930 *Quality Management System*
- 931 • ISO/IEC 23053:2022 *Framework for Artificial Intelligence (AI) Systems Using Machine*  
932 *Learning (ML)*

---

<sup>22</sup> Please note that related changes to MLMD are described in Section 7.1.1 of IMDRF/AIMD WG/N67 (Edition 1):2022 *Machine Learning-enabled Medical Devices: Key Terms and Definitions* as “[h]eterogeneous changes [that] are non-uniform changes that can be specific to one clinic, region, demographic, etc. (sometimes referred to as local adaptations)”. Both terms used in this section (localization and customization) are examples of these types of changes and the methods used to facilitate such changes.

## 933 5.6. Operations and Monitoring

934 As with all medical devices, monitoring of deployed AI-enabled medical devices can help detect and  
935 address any performance<sup>23</sup>, technical, security and operational issues arising in the intended use  
936 environment.

937 IMDRF/SaMD WG/N23 FINAL:2015 recommends that “post market surveillance including monitoring,  
938 measurement and analysis of quality data can include logging and tracking of complaints, clearing  
939 technical issues, determining problem causes and actions to address, identify, collect, analyse, and  
940 report on critical quality characteristics of products developed”. ISO 13485:2016 includes similar  
941 provisions for post-market surveillance and feedback.

942 The legal responsibility to meet medical device requirements, including post-market monitoring and  
943 surveillance, falls on the manufacturer of the AI-enabled medical device; in practice, manufacturers  
944 may depend upon the deploying entity/site to carry out post-market responsibilities, including  
945 collection of much of the data needed to track the performance of their device. Throughout this step of  
946 the life cycle (and Real-World Performance Evaluation in Section 5.7), it is important for  
947 manufacturers to consider how responsibilities may need to be shared and communicated to the user.

948 The rest of this section covers aspects of operations and monitoring that are specific or unique to AI-  
949 enabled medical devices, many of which enable the real-world performance evaluation described in  
950 the section below. Many, if not all, of the systems that facilitate operations and monitoring help to  
951 facilitate real-world performance evaluation. In this document, “Operations and Monitoring” refers  
952 primarily to continuous, often automated activities embedded within the quality management system  
953 (for example, logging, alerting, and incident handling), whereas “Real-World Performance Evaluation”  
954 (Section 5.8) refers to structured analyses or studies that use these data, and where appropriate  
955 additional data collection, to answer specific questions about safety, effectiveness, and clinical utility.  
956 Although concepts may overlap between these two sections, this distinction can help clarify roles and  
957 expectations across the AI-enabled medical device life cycle.

958 *Infrastructure, Data and Logging:* As discussed in Section 5.5 above, the deployed AI-enabled  
959 medical device’s infrastructure facilitates the collection of data needed to fulfill monitoring  
960 requirements, and can incorporate data logging and performance tracking capabilities, as well as  
961 anomaly detection algorithms, incident response procedures and troubleshooting workflows. For AI-  
962 enabled medical devices, it may be worthwhile for the infrastructure to have both local monitoring  
963 capabilities as well as the ability to contribute accumulated monitoring data to central repositories,  
964 ensuring a comprehensive and uninterrupted record of the AI-enabled medical device’s operational  
965 performance across deployments for the manufacturer. Where possible, the process may integrate  
966 advanced technologies that enable continual assessment of the AI-enabled medical device’s  
967 accuracy and responsiveness, supporting rapid intervention, adaptive recalibration, and maintenance  
968 of the AI-enabled medical device’s reliability and effectiveness across its operational life cycle. In  
969 addition, logging and auditing are critical components of an AI-enabled medical device’s operational  
970 monitoring activities and maintenance. For AI-enabled devices specifically, it may be beneficial for the  
971 logging mechanisms to capture detailed logs and audit trails (as allowed within jurisdictional  
972 regulations) of model predictions and any explainability elements (if available) such as heatmaps or  
973 top predictive features, data inputs, model outputs, and metadata such as timestamps, user  
974 interactions, model versions, and environmental context for traceability, accountability, and  
975 compliance purposes.

---

<sup>23</sup> See GMLP Guiding Principle 10.

976 *Performance Degradation and Drift Detection:* Performance of the models within the AI-enabled  
977 medical devices can degrade over time as the statistical patterns that informed the original model(s')  
978 training may become less representative at one or more deployment sites, causing the AI-enabled  
979 medical devices to generate increasingly inaccurate or inappropriate outputs, even though the  
980 model(s) have not changed. Drift can emerge from changes in input data characteristics, shifts in  
981 underlying population dynamics, changes in context of use, such as clinical behaviour or patient  
982 management, or gradual accumulation of data biases that were not apparent during initial training.  
983 Continuous post-market monitoring and surveillance enables early detection of these performance  
984 variations, yielding better visibility to manufacturer and supporting the implementation of timely  
985 recalibration, retraining, or intervention strategies when necessary and with appropriate regulatory  
986 approval. If drift is detected, it may be important to consider additional investigations, such as  
987 retrospective analysis. Monitoring drift in generative AI-enabled medical devices is especially critical  
988 as these systems can rapidly and unpredictably deviate from their original training objectives.  
989 Effective drift monitoring ensures that generative AI-enabled devices maintain their core functional  
990 integrity and alignment with the original design. The non-deterministic nature of some AI models  
991 creates unique challenges for post-market performance assessment, as outputs may vary from those  
992 observed during pre-market testing without necessarily indicating improved or degraded performance.  
993 In these cases, manufacturers can consider developing post-market metrics that account for output  
994 variability and assess whether the model's performance remains within acceptable bounds and  
995 implement monitoring strategies that distinguish between expected variation and performance drift.

996 *Alerting and Reporting:* The use of continuous monitoring and responsive alerting and reporting  
997 systems helps to ensure reliability, stability and effectiveness of AI-enabled medical devices at  
998 deployed sites. For example, use of thresholds for critical performance metrics, automated  
999 notifications and alerts are some of the ways manufacturers can continuously monitor their deployed  
1000 AI-enabled medical devices and promptly address risks. In addition to individual incidents, it is  
1001 important to monitor for changes in data distribution and performance over time, e.g., drift as  
1002 described above. Jurisdictional requirements may dictate relevant alerting and reporting needs to  
1003 ensure timelines and other regulatory requirements for reporting are met.

1004 *Monitoring for Advanced AI Models:* Operations monitoring for customizable and adaptive AI-enabled  
1005 devices may necessitate fundamentally different approaches due to their inherent architectural  
1006 differences. Monitoring customizable AI-enabled medical devices, which operate within predefined  
1007 parameters defined as part of regulatory authorization and set during deployment, focuses on  
1008 verifying consistent performance within established boundaries and detecting any deviations from  
1009 original specifications. Adaptable AI-enabled medical devices, which can have their algorithms and  
1010 learning processes modified, necessitate more dynamic and sophisticated monitoring strategies that  
1011 track not only performance metrics but also the AI-enabled medical device's own learning and  
1012 transformation mechanisms. This may include continuous surveillance of the adaptable device's  
1013 decision-making evolution, potential drift in core algorithmic logic, and ongoing assessment of  
1014 whether the AI-enabled medical device's modifications maintain alignment with the approved design.  
1015 The monitoring approach for adaptable AI-enabled medical devices is inherently more complex,  
1016 incorporating real-time analysis of the AI-enabled medical device's learning trajectories, unexpected  
1017 behaviour emergence, and potential unintended consequences of autonomous algorithmic  
1018 modifications.

1019 In addition, operational monitoring of AI-enabled medical devices such as generative AI-enabled  
1020 medical devices may necessitate different post-market monitoring due to their emergent behaviour,  
1021 user-driven variability, and context-sensitive outputs. To effectively capture unexpected uses,  
1022 dynamic, usage-aware, and risk-responsive monitoring approaches may be beneficial, including  
1023 approaches that monitor how the AI-enabled medical device is used due to the possibility of the AI-  
1024 enabled medical device being repurposed, misused, or used beyond its intended use. Tracking the  
1025 deployment context and collecting interaction logs and prompt-response histories from LLM's (with  
1026 privacy safeguards) as well as employing mechanisms to detect patterns of prompt chaining and  
1027 workarounds may be helpful to flag unsafe deviations. Prompt classification and risk-based escalation  
1028 to overseeing AI-enabled medical devices or manual review may be helpful mechanisms to implement  
1029 in order to monitor and address developing unsafe uses.

## 1030 **Additional Standards and References to be Considered**

- 1031 • ISO/IEC 23053:2022 *Framework for Artificial Intelligence (AI) Systems Using Machine*  
1032 *Learning (ML)*

## 1033 5.7. Real-World Performance Evaluation

1034 Understanding how an AI-enabled medical device performs in the real world<sup>24</sup> can be beneficial to  
1035 supporting clinical utility and better understanding of the risks and benefits of the AI-enabled medical  
1036 device. It can also facilitate the collection of data at scale as compared to pre-market evaluation, and  
1037 help manufacturers understand how the device is performing in subpopulations in detail. It is  
1038 important for manufacturers to consider implementing a structured and strategic approach to  
1039 establishing and tracking the real-world performance of AI-enabled medical devices.

1040 Monitoring infrastructure and data described in the sections above enables real-world performance  
1041 evaluation described in this section (e.g., detailed logs and audit trails of model predictions and  
1042 explainability elements). Drift monitoring in this step goes beyond the individual model's performance  
1043 to understand how the AI-enabled medical device and its model(s) are performing in the context of the  
1044 clinical workflow. For example, if clinicians are frequently choosing to override aspects of the AI-  
1045 enabled medical device's outputs or the output entirely (although the Operations and Monitoring  
1046 infrastructure indicates that the model is performing as expected), it may be worth the manufacturer  
1047 looking into this further as part of the performance evaluation in its clinical context.

1048 The rest of this section provides an overview of considerations with respect to real-world performance  
1049 evaluation of AI-enabled medical devices.

1050 *Define Device Performance Indicators:* It is important for manufacturers to establish device  
1051 performance indicators that are framed by the intended use of the device/model and that are both  
1052 informative and scalable. Ideal device performance indicators reflect the safety and effectiveness of  
1053 the AI-enabled medical device in its real-world setting and allow for the detection of post-market  
1054 safety signals. Depending on the AI-enabled medical device and model, manufacturers may wish to  
1055 consider multiple different device performance indicators to capture different failure modes and  
1056 different aspects of safety and performance, some of which could be the evaluation metrics used in  
1057 prior steps of the life cycle when feasible. For example, device performance indicators can include  
1058 technical measures of performance such as accuracy, precision, recall, and latency or other  
1059 measures that capture the performance of the human-AI team such as user satisfaction.

1060 Where feasible, it is important for device performance indicators to be defined not only in aggregate  
1061 but also for clinically relevant subgroups, such as demographics, disease subtypes, and deployment  
1062 sites, to support the identification of differential performance that could impact safety, effectiveness, or  
1063 equity. It is also helpful to pre-specify risk-based thresholds, trends, or trigger conditions for these  
1064 indicators (for example, sustained degradation in a safety-critical metric or emergence of significant  
1065 performance gaps between subgroups) that may prompt investigation, mitigation, or model updates,  
1066 in line with applicable regulatory requirements and change control mechanisms.

1067 For AI-enabled medical devices enabled by adaptable and/or autonomous models, it is important for  
1068 manufacturers to specify device performance indicators that measure whether the model operation  
1069 remains within an acceptable performance range. Manufacturers may also consider defining  
1070 governance procedures for remediation when model performance falls outside of the specified  
1071 acceptable range according to predetermined criteria (e.g., sustained performance degradation  
1072 exhibited with a particular frequency over a specified time span).

1073 *Data Collection:* The infrastructure described in Section 5.5. and 5.6 can enable collection of  
1074 necessary and representative data to measure the device performance indicators, with efforts made  
1075 to ensure that data is of sufficient quality and integrity and is anonymized and protected, where  
1076 applicable as required by jurisdictional regulations. Additionally, proactive data collection can occur  
1077 outside of this infrastructure as well, such as interviews with clinicians and other users and surveys to  
1078 collate actionable information.

---

<sup>24</sup> See GMLP Guiding Principle 10.

1079 *Optimizing Performance*: While performance optimization can occur based on signals that surface  
1080 during Operations and Monitoring, insights from real-world performance evaluation can also drive  
1081 continuous improvement in the AI-enabled medical device, including when device performance  
1082 indicators demonstrate consistent decline, when safety or effectiveness of the AI-enabled medical  
1083 device is compromised, when significant technological or methodological improvements emerge, or  
1084 from user feedback. This may lead to, for example, performance optimization by retraining the  
1085 model(s), and considerations throughout this document can inform that activity, along with an  
1086 assessment of the need for new regulatory approval.

1087 **Additional Standards and References to be Considered**

- 1088 • ISO/IEC 23053:2022 Framework for Artificial Intelligence (AI) Systems Using Machine  
1089 Learning (ML)

1090 **5.8. Sunsetting**

1091 Sometimes called decommissioning, retirement, or end-of-life, as discussed in IMDRF/SaMD  
1092 WG/N23 FINAL:2015 for a SaMD, sunsetting involves the termination of “maintenance, support, and  
1093 distribution” of the AI-enabled medical device “in a controlled and a managed fashion.” In common  
1094 with SaMD, “this process indicates an end to active support, and may entail deactivation and/or  
1095 removal of [the AI-enabled medical device] and its supporting data.” This can include, for example,  
1096 technical decommissioning, which can involve gradually scaling down and eventually deactivating  
1097 system components, ensuring a controlled shutdown to prevent data loss or service disruptions. It can  
1098 also include decommissioning and repurposing of hardware and cloud resources used by the AI-  
1099 enabled medical device. For AI-enabled medical devices specifically, it is important to communicate  
1100 sunsetting plans to stakeholders, both internal teams and external users and sites, with clear  
1101 timelines and reasons for the decision, as well as alternative solutions or support to help them  
1102 transition, as applicable. Furthermore, as described in the Data Collection and Management life cycle  
1103 step, retention of data for regulatory, legal and business reasons is critical, ensuring compliance with  
1104 jurisdictional data protection regulations. It is important for manufacturers to maintain detailed  
1105 documentation of the sunsetting process for audit purposes and to address any future legal or  
1106 compliance inquiries.

1107 **Additional Standards and References to be Considered**

- 1108 • ISO/IEC 23053:2022 *Framework for Artificial Intelligence (AI) Systems Using Machine*  
1109 *Learning (ML)*

1110

1111

1112

## 6. Transparency and Labelling

1113 Transparency, or the communication of clear, essential information<sup>25</sup> about an AI-enabled medical  
1114 device to relevant audiences, is a key principle throughout the life cycle of an AI-enabled medical  
1115 device and can help address risks related to information (see Section 4.2.1) and risks related to the  
1116 human-AI interaction (see Section 4.2.2).

1117 IMDRF/GRRP WG/N52: FINAL:2024 (Edition 2) *Principles of Labelling for Medical Devices and IVD*  
1118 *Medical Devices* provides “general labelling principles, including specific sections on the label,  
1119 instructions for use, and information intended for the patient” for all medical devices. These apply  
1120 when developing labelling for AI-enabled medical devices as well, in addition to jurisdictional  
1121 requirements for medical device labelling. As with all medical devices, AI-enabled medical devices  
1122 should be accompanied by information that is truthful, accurate, and sufficient to support their use by  
1123 the intended user and to ensure that they are safe and effective for their intended use or purpose.  
1124 This information helps ensure, for example, users have the information they need about both the  
1125 characteristics of the AI-enabled medical device, such as its model(s), and how those characteristics  
1126 were assessed. See Appendix C for the types of information to be included in labelling to support safe  
1127 and effective use.

1128 With AI-enabled medical devices, it is particularly important that necessary information is presented to  
1129 the user at the most appropriate time and in the most appropriate location and manner during medical  
1130 device procurement, implementation, use, modifications, or at the AI-enabled medical device’s end-of-  
1131 life. When appropriate, labelling may sometimes be implemented through the medical devices’ user  
1132 interface (UI), with additional information included in the UI that can help further achieve  
1133 transparency. This can include information related to the explainability and interpretability of the AI-  
1134 enabled medical device’s model behaviour and its output(s). For example, it is important for users  
1135 who will rely on the output of the AI-enabled medical device in time-critical decision making to have a  
1136 clear understanding of the relationship between how the safety and effectiveness of the device was  
1137 demonstrated with respect to the type of patient and condition being treated (including regulatory  
1138 authorization information as applicable) and when it is appropriate to rely on the output in these  
1139 situations. It may impede clinical decision-making, if the user needs to interpret complex model  
1140 outputs without accompanying clarifying information, especially in emergency situations, and  
1141 overburden these users. Manufacturers may also wish to consider electronic labelling, which can  
1142 support increased availability, utility, interactivity, and accessibility of labelling.

1143 For AI-enabled medical device models with no planned updates, along with the types of information  
1144 recommended within Appendix C, it is important that transparency efforts also focus on the model’s  
1145 initial technological and performance characteristics. For models that are expected to undergo  
1146 discrete updates, tailored transparency information can be provided at the time of each update,  
1147 including whether the update required regulatory approval, for example as part of PCCP. For  
1148 autonomously adaptable models designed to perform self-implementing updates, transparency may  
1149 also include details on the update process and the parameters or objectives guiding those updates.  
1150 When general-purpose, off-the-shelf models are incorporated (as described in Section 5.3), it is  
1151 important for users to have information on these models and their limitations. In all cases, it is  
1152 important for the manufacturer to ensure the intended user can understand what, if any, risk controls  
1153 and monitoring are in place and what, if anything, is expected of them to ensure safe and effective  
1154 use on an ongoing basis.

### 1155 **Additional Standards and References to be Considered**

- 1156 • Transparency for Machine Learning-Enabled Medical Devices: Guiding Principles (Published  
1157 June 2024): [https://www.fda.gov/medical-devices/software-medical-device-  
1158 \[samd/transparency-machine-learning-enabled-medical-devices-guiding-principles\]\(https://www.fda.gov/medical-devices/software-medical-device-samd/transparency-machine-learning-enabled-medical-devices-guiding-principles\)](https://www.fda.gov/medical-devices/software-medical-device-samd/transparency-machine-learning-enabled-medical-devices-guiding-principles)

---

<sup>25</sup> See GMLP Guiding Principle 9.

1159

## 7. Conclusion

1160 In summary, this document provides a harmonized framework to address the challenges and  
1161 opportunities of AI-enabled medical devices, ensuring safety, effectiveness, and patient-focused  
1162 innovation. The universal concepts of quality and risk management, cybersecurity, human oversight  
1163 (see Section 4) are essential for their responsible development and use. The Life Cycle Steps (see  
1164 Section 5) help manufacturers to ensure that AI-enabled medical devices and their model(s) are  
1165 designed, developed, and deployed with GMLP, patient safety and clinical effectiveness in mind.

1166 This document complements existing IMDRF publications and regulatory requirements, offering a  
1167 common language and principles adaptable to local contexts. Jurisdictions may integrate these  
1168 concepts to promote global consistency and alignment.

1169 The goal is to foster innovation while ensuring AI-enabled medical devices are safe, effective, and  
1170 capable of delivering meaningful benefits to patients and healthcare systems. By adopting the  
1171 considerations and principles outlined in this document, stakeholders can contribute to the  
1172 development of AI-enabled medical devices that are trustworthy, transparent, and aligned with the  
1173 needs of patients, healthcare providers, and requirements of regulatory authorities. This collaborative  
1174 approach will help ensure that AI-enabled medical devices continue to advance healthcare while  
1175 maintaining the highest standards of safety and effectiveness.

1176 **Appendix A: Traceability**  
 1177 **between GMLP and Life Cycle**  
 1178 **Steps/Document Sections**

AI Life Cycle Step and/or Document Section	Relevant GMLP Principle
<b>Planning and Design</b>	All
<b>Data Collection and Management</b>	Guiding Principle 3: Clinical evaluation includes the use of datasets that are representative of the intended patient population Guiding Principle 4: Training datasets are independent of test sets Guiding Principle 5: Selected reference standards are fit-for-purpose
<b>Model Building and Tuning</b>	Guiding Principle 2: Good software engineering, medical device design, and security practices are implemented throughout the total product life cycle Guiding Principle 6: Model choice and design are tailored to the available data and the intended use/ intended purpose of the device
<b>Verification and Validation, including Clinical Evaluation</b>	Guiding Principle 3: Clinical evaluation includes the use of datasets that are representative of the intended patient population Guiding Principle 4: Training datasets are independent of test sets Guiding Principle 5: Selected reference standards are fit-for-purpose Guiding Principle 7: The device is assessed with a focus on human-AI interactions in the intended use environment, including the performance of the human-AI team, rather than just the device in isolation Guiding Principle 8: Testing demonstrates device performance during clinically relevant conditions
<b>Deployment</b>	Guiding Principle 10: Deployed models are monitored for performance and re-training risks are managed
<b>Operations and Monitoring</b>	Guiding Principle 10: Deployed models are monitored for performance and re-training risks are managed
<b>Real-World Performance Evaluation</b>	Guiding Principle 10: Deployed models are monitored for performance and re-training risks are managed
<b>Sunsetting</b>	Guiding Principle 9: Users are provided clear, essential information

1179

DRAFT

# Appendix B: Examples of Common Evaluation Metrics

Examples of common evaluation metrics may include, but are not limited to, the following:

- Accuracy (measures the overall correctness of the predictions);
- Precision (measures the proportion of true positive predictions out of all positive predictions);
- Recall (measures the ability to correctly identify actual positives);
- Model performance, including different classification thresholds and discrimination ability;
- F1 score (measures precision and recall);
- Receiver Operating Characteristic – Area Under Curve (ROC-AUC) (evaluates model performance different classification thresholds and discrimination ability); and
- Mean squared error (MSE) (measures differences between predicted and actual values).

Other metrics like perplexity (metric to measure the uncertainty of a language model when predicting the next word in the sequence) and coherence (assesses the contextual relevance and logical structure of generated text) may be used for evaluation of LLM-based models.

# Appendix C: Labelling Elements

For AI-enabled medical devices, information that may be needed to support safe and effective use includes, but is not limited to:

- the intended use or purpose and intended user;
- any professional or other qualifications that are necessary to use the device safely and effectively;
- clear information on compatibility requirements for inputs (if applicable);
- how the device should be used, including any prerequisite steps or necessary preparations;
- the populations in which the device can be used and the conditions for which it can be used (including known limitations);
- measured performance and subgroup performance;
- when and for how long the device should be used;
- pre-market testing conditions and datasets;
- how to interpret the output of the device, known limitations, model description or characteristics.

See also IMDRF/GRRP WG/N52 FINAL:2024 (Edition 2) *Principles of Labeling for Medical Devices and IVD Medical Devices* for additional information to consider.

**Please visit our website  
for more details.**

[www.imdrf.org](http://www.imdrf.org)

#### **Disclaimer**

© Copyright 202X by the International Medical Device Regulators Forum.

This work is copyright. Subject to these Terms and Conditions, you may download, display, print, translate, modify and reproduce the whole or part of this work for your own personal use, for research, for educational purposes or, if you are part of an organisation, for internal use within your organisation, but only if you or your organisation do not use the reproduction for any commercial purpose and retain all disclaimer notices as part of that reproduction. If you use any part of this work, you must include the following acknowledgement (delete inapplicable):

"[Translated or adapted] from [insert name of publication], [year of publication], International Medical Device Regulators Forum, used with the permission of the International Medical Device Regulators Forum. The International Medical Device Regulators Forum is not responsible for the content or accuracy of this [adaption/translation]."

All other rights are reserved, and you are not allowed to reproduce the whole or any part of this work in any way (electronic or otherwise) without first being given specific written permission from IMDRF to do so. Requests and inquiries concerning reproduction and rights are to be sent to the IMDRF Secretariat.

Incorporation of this document, in part or in whole, into another document, or its translation into languages other than English, does not convey or represent an endorsement of any kind by the IMDRF.