



NIST Special Publication 800
NIST SP 800-218A ipd

Secure Software Development Practices for Generative AI and Dual-Use Foundation Models

An SSDF Community Profile

Initial Public Draft

Harold Booth
Murugiah Souppaya
Apostol Vassilev
Michael Ogata
Martin Stanley
Karen Scarfone

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.SP.800-218A.ipd>

NIST Special Publication 800
NIST SP 800-218A ipd

Secure Software Development Practices for Generative AI and Dual-Use Foundation Models

An SSDF Community Profile

Initial Public Draft

Harold Booth
Murugiah Souppaya
Apostol Vassilev
*Computer Security Division
Information Technology Laboratory*

Michael Ogata
*Applied Cybersecurity Division
Information Technology Laboratory*

Martin Stanley
*Cybersecurity and Infrastructure Security
Agency (CISA)*

Karen Scarfone
Scarfone Cybersecurity

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.SP.800-218A.ipd>

April 2024



U.S. Department of Commerce
Gina M. Raimondo, Secretary

National Institute of Standards and Technology
Laurie E. Locascio, NIST Director and Under Secretary of Commerce for Standards and Technology

Certain equipment, instruments, software, or materials, commercial or non-commercial, are identified in this paper in order to specify the experimental procedure adequately. Such identification does not imply recommendation or endorsement of any product or service by NIST, nor does it imply that the materials or equipment identified are necessarily the best available for the purpose.

There may be references in this publication to other publications currently under development by NIST in accordance with its assigned statutory responsibilities. The information in this publication, including concepts and methodologies, may be used by federal agencies even before the completion of such companion publications. Thus, until each publication is completed, current requirements, guidelines, and procedures, where they exist, remain operative. For planning and transition purposes, federal agencies may wish to closely follow the development of these new publications by NIST.

Organizations are encouraged to review all draft publications during public comment periods and provide feedback to NIST. Many NIST cybersecurity publications, other than the ones noted above, are available at <https://csrc.nist.gov/publications>.

Authority

This publication has been developed by NIST in accordance with its statutory responsibilities under the Federal Information Security Modernization Act (FISMA) of 2014, 44 U.S.C. § 3551 et seq., Public Law (P.L.) 113-283. NIST is responsible for developing information security standards and guidelines, including minimum requirements for federal information systems, but such standards and guidelines shall not apply to national security systems without the express approval of appropriate federal officials exercising policy authority over such systems. This guideline is consistent with the requirements of the Office of Management and Budget (OMB) Circular A-130.

Nothing in this publication should be taken to contradict the standards and guidelines made mandatory and binding on federal agencies by the Secretary of Commerce under statutory authority. Nor should these guidelines be interpreted as altering or superseding the existing authorities of the Secretary of Commerce, Director of the OMB, or any other federal official. This publication may be used by nongovernmental organizations on a voluntary basis and is not subject to copyright in the United States. Attribution would, however, be appreciated by NIST.

NIST Technical Series Policies

[Copyright, Use, and Licensing Statements](#)

[NIST Technical Series Publication Identifier Syntax](#)

How to Cite this NIST Technical Series Publication

Booth H, Souppaya M, Vassilev A, Ogata M, Stanley M, Scarfone K (2024) Secure Development Practices for Generative AI and Dual-Use Foundation AI Models: An SSDF Community Profile. (National Institute of Standards and Technology, Gaithersburg, MD), NIST Special Publication (SP) NIST SP 800-218A ipd.
<https://doi.org/10.6028/NIST.SP.800-218A.ipd>

Author ORCID iDs

Harold Booth: 0000-0003-0373-6219

Murugiah Souppaya: 0000-0002-8055-8527

Apostol Vassilev: 0000-0002-9081-3042

Michael Ogata: 0000-0002-8457-2430

Karen Scarfone: 0000-0001-6334-9486

Public Comment Period

April 29, 2024 – June 2, 2024

Submit Comments

Comments on NIST SP 800-218A may be sent electronically to SSDF@nist.gov with “NIST SP 800-218A, Secure Software Development Practices for Generative AI and Dual-Use Foundation Models” in the subject line.

Comments may also be submitted via www.regulations.gov: enter NIST-2024-0001 in the search field, click on the “Comment Now!” icon, complete the required fields, including “NIST SP 800-218A, Secure Software Development Practices for Generative AI and Dual-Use Foundation Models” in the subject field, and enter or attach your comments. Comments containing information in response to this notice must be received on or before **June 2, 2024, at 11:59 PM Eastern Time**.

Additional Information

Additional information about this publication is available at <https://csrc.nist.gov/pubs/sp/800/218/a/ipd>, including related content, potential updates, and document history.

All comments are subject to release under the Freedom of Information Act (FOIA).

1 **Abstract**

2 This document augments the secure software development practices and tasks defined in
3 Secure Software Development Framework (SSDF) version 1.1 by adding practices, tasks,
4 recommendations, considerations, notes, and informative references that are specific to AI
5 model development throughout the software development life cycle. These additions are
6 documented in the form of an SSDF Community Profile to support Executive Order (EO) 14110,
7 *Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence*, which tasked NIST
8 with “developing a companion resource to the [SSDF] to incorporate secure development
9 practices for generative AI and for dual-use foundation models.” This Community Profile is
10 intended to be useful to the producers of AI models, the producers of AI systems that use those
11 models, and the acquirers of those AI systems. This Profile should be used in conjunction with
12 NIST Special Publication (SP) 800-218, *Secure Software Development Framework (SSDF) Version*
13 *1.1: Recommendations for Mitigating the Risk of Software Vulnerabilities*.

14 **Keywords**

15 artificial intelligence; artificial intelligence model; cybersecurity risk management; generative
16 artificial intelligence; secure software development; Secure Software Development Framework
17 (SSDF); software acquisition; software development; software security.

18 **Reports on Computer Systems Technology**

19 The Information Technology Laboratory (ITL) at the National Institute of Standards and
20 Technology (NIST) promotes the U.S. economy and public welfare by providing technical
21 leadership for the Nation’s measurement and standards infrastructure. ITL develops tests, test
22 methods, reference data, proof of concept implementations, and technical analyses to advance
23 the development and productive use of information technology. ITL’s responsibilities include
24 the development of management, administrative, technical, and physical standards and
25 guidelines for the cost-effective security and privacy of other than national security-related
26 information in federal information systems. The Special Publication 800-series reports on ITL’s
27 research, guidelines, and outreach efforts in information system security, and its collaborative
28 activities with industry, government, and academic organizations.

29 **Audience**

30 There are three primary audiences for this document:

- 31 • *AI model producers* — Organizations that are developing their own generative AI and
32 dual-use foundation models, as defined in EO 14110
- 33 • *AI system producers* — Organizations that are developing software that leverages a
34 generative AI or dual-use foundation model
- 35 • *AI system acquirers* — Organizations that are acquiring a product or service that utilizes
36 one or more AI systems

37 Individuals who are interested in better understanding secure software development practices
38 for AI models may also benefit from this document.

39 Readers are not expected to be experts in secure software development or AI model
40 development, but such expertise may be needed to implement these recommended practices.

41 **Note to Reviewers**

42 NIST welcomes feedback on any part of this document but is particularly interested in the
43 following:

- 44 • The Profile suggests adding several practices and tasks to those defined in SSDF version
45 1.1. Are these additions reasonable? What other additions would help address secure
46 development practices for generative AI and dual-use foundation models?
- 47 • What changes should be made to the Profile's Recommendations, Considerations, and
48 Notes column to help address secure software development practices for generative AI
49 and dual-use foundation models?
- 50 • What additional cybersecurity, privacy, and/or reproducibility considerations should be
51 taken into account when selecting model training techniques (e.g., deterministic model
52 training)?
- 53 • What suggestions do you have for Implementation Examples and additional Informative
54 References for the Profile?
- 55 • Is this Profile fully applicable to the secure development of other types of AI models
56 besides generative and dual-use foundation models? If not, what changes could be
57 made to the Profile to accommodate other AI models?
- 58 • Is this Profile flexible enough to support AI model producers, the producers of AI
59 systems using those models, and the acquirers of those AI systems?
- 60 • What guidance, templates, or other resources on SSDF Community Profile use would
61 you find beneficial?

62 If you are from a standards developing organization (SDO) or another organization that is
63 defining a set of secure practices for AI model development and you would like to map your

64 standard or guidance to the SSDF profile, please contact the authors at ssdf@nist.gov. They will
65 introduce you to the [National Online Informative References Program \(OLIR\)](#), where you can
66 submit your mapping to augment the existing set of informative references.

67 **Trademark Information**

68 All registered trademarks belong to their respective organizations.

69 **Acknowledgments**

70 The authors thank all of the organizations and individuals who provided input for this
71 publication. In response to Executive Order (EO) 14110, [Safe, Secure, and Trustworthy](#)
72 [Development and Use of Artificial Intelligence](#), NIST held a [January 2024 workshop](#), where
73 speakers and attendees shared suggestions for adapting secure software development
74 practices and tasks to accommodate the unique aspects of AI model development and the
75 software that leverages them. The authors also thank all of their NIST colleagues and external
76 experts who provided suggestions and feedback that helped shape this publication.

77 **Call for Patent Claims**

78 This public review includes a call for information on essential patent claims (claims whose use
79 would be required for compliance with the guidance or requirements in this Information
80 Technology Laboratory (ITL) draft publication). Such guidance and/or requirements may be
81 directly stated in this ITL Publication or by reference to another publication. This call also
82 includes disclosure, where known, of the existence of pending U.S. or foreign patent
83 applications relating to this ITL draft publication and of any relevant unexpired U.S. or foreign
84 patents.

85 ITL may require from the patent holder, or a party authorized to make assurances on its behalf,
86 in written or electronic form, either:

- 87 a) assurance in the form of a general disclaimer to the effect that such party does not hold
88 and does not currently intend holding any essential patent claim(s); or
- 89 b) assurance that a license to such essential patent claim(s) will be made available to
90 applicants desiring to utilize the license for the purpose of complying with the guidance
91 or requirements in this ITL draft publication either:
 - 92 i. under reasonable terms and conditions that are demonstrably free of any unfair
93 discrimination; or
 - 94 ii. without compensation and under reasonable terms and conditions that are
95 demonstrably free of any unfair discrimination.

96 Such assurance shall indicate that the patent holder (or third party authorized to make
97 assurances on its behalf) will include in any documents transferring ownership of patents
98 subject to the assurance, provisions sufficient to ensure that the commitments in the assurance
99 are binding on the transferee, and that the transferee will similarly include appropriate
100 provisions in the event of future transfers with the goal of binding each successor-in-interest.

101 The assurance shall also indicate that it is intended to be binding on successors-in-interest
102 regardless of whether such provisions are included in the relevant transfer documents.

103 Such statements should be addressed to: ssdf@nist.gov

104	Table of Contents	
105	1. Introduction	1
106	1.1. Purpose	1
107	1.2. Scope	2
108	1.3. Sources of Expertise	2
109	1.4. Document Structure.....	2
110	2. Using the SSDF Community Profile	4
111	3. SSDF Community Profile for AI Model Development	6
112	References	18
113	Appendix A. Glossary	20
114	List of Tables	
115	Table 1. SSDF Community Profile for AI Model Development	8
116		

117 1. Introduction

118 Section 4.1.a of Executive Order (EO) 14110, *Safe, Secure, and Trustworthy Development and*
119 *Use of Artificial Intelligence* [1], tasked NIST with “developing a companion resource to the
120 Secure Software Development Framework to incorporate secure development practices for
121 [generative AI](#) and for [dual-use foundation models](#).” This document is that companion resource.

122 The software development and use of [AI models](#) and [AI systems](#) inherit much of the same risk
123 as any other digital system. A unique challenge for this community is the blurring of traditional
124 boundaries between system code and system data, as well as the use of plain human language
125 as the means of interaction with the systems. AI models and systems, their configuration
126 parameters (e.g., model weights), and the data they interact with (e.g., training data, user
127 queries, etc.) can form closed loops that can be manipulated for unintended functionality.

128 AI model and system development is still much more of an art than an exact science, requiring
129 developers to interact with model code, training data, and other parameters over multiple
130 iterations. Training datasets may be acquired from unknown, untrusted sources. Model weights
131 and other training parameters can be susceptible to malicious tampering. Some models may be
132 complex to the point that they cannot easily be thoroughly inspected, potentially allowing for
133 undetectable execution of arbitrary code. User queries can be crafted to produce undesirable
134 or objectionable output and — if not sanitized properly — can be leveraged for injection-style
135 attacks. The goal of this document is to identify the practices and tasks needed to address these
136 novel risks.

137 1.1. Purpose

138 The SSDF provides a common language for describing secure software development practices
139 throughout the software development life cycle. This document augments the practices and
140 tasks defined in SSDF version 1.1 by adding recommendations, considerations, notes, and
141 informative references that are specific to generative AI and dual-use foundation model
142 development. These additions are documented in the form of an *SSDF Community Profile*
143 (“Profile”), which is a baseline of SSDF practices and tasks that have been enhanced to address
144 a particular use case. An example of an addition is, “Secure code storage should include AI
145 models, model weights, pipelines, reward models, and any other AI model elements that need
146 their confidentiality, integrity, and/or availability protected.”

147 **This Profile supplements what SSDF version 1.1 already includes. The Profile is intended to be**
148 **used in conjunction with NIST Special Publication (SP) 800-218, *Secure Software Development***
149 ***Framework (SSDF) Version 1.1: Recommendations for Mitigating the Risk of Software***
150 ***Vulnerabilities*** [6]. Readers should also utilize the implementation examples and informative
151 references defined in SP 800-218 for additional information on how to perform each SSDF
152 practice and task for all types of software development, as they are also generally applicable to
153 AI model and AI system development.

154 **1.2. Scope**

155 This Profile's scope is *AI model development*, which includes data sourcing for, designing,
156 training, fine-tuning, and evaluating AI models, as well as incorporating and integrating AI
157 models into other software. Consistent with SSDF version 1.1 and EO 14110, practices for the
158 deployment and operation of AI systems with AI models are out of scope. Similarly, while
159 cybersecurity practices for training data and other forms of data being used for AI model
160 development are in scope, the rest of the data governance and management life cycle is out of
161 scope.

162 **1.3. Sources of Expertise**

163 This document leverages and integrates numerous sources of expertise, including:

- 164 • NIST research and publications on trustworthy and responsible AI, including the *Artificial*
165 *Intelligence Risk Management Framework (AI RMF 1.0)* [2], *Adversarial Machine*
166 *Learning: A Taxonomy and Terminology of Attacks and Mitigations* [3], *Towards a*
167 *Standard for Identifying and Managing Bias in Artificial Intelligence* [4], and the Dioptra
168 experimentation testbed for security evaluations of machine learning algorithms [5].
- 169 • NIST's *Secure Software Development Framework (SSDF) Version 1.1* [6], which is a set of
170 fundamental, sound, and secure software development practices. It provides a common
171 language to help facilitate communications among stakeholders, including software
172 producers and software acquirers. The SSDF has also been used in support of EO 14028,
173 *Improving the Nation's Cybersecurity* [7], to enhance software supply chain security.
- 174 • NIST general cybersecurity resources, including *The NIST Cybersecurity Framework (CSF)*
175 *2.0* [8], *Security and Privacy Controls for Information Systems and Organizations* [9], and
176 *Cybersecurity Supply Chain Risk Management Practices for Systems and Organizations*
177 [10].
- 178 • AI model developers, AI researchers, AI system developers, and secure software
179 practitioners from industry and government with expertise in the unique security
180 challenges of AI models and the practices for addressing those challenges. This expertise
181 was primarily captured through NIST's [January 2024 workshop](#), where speakers and
182 attendees shared suggestions for adapting secure software development practices and
183 tasks to accommodate the unique aspects of AI model development and the software
184 leveraging them.

185 **1.4. Document Structure**

186 This document is structured as follows:

- 187 • Section 2 provides additional background on the SSDF and explains what an SSDF
188 Community Profile is and how it can be used.
- 189 • Section 3 defines the SSDF Community Profile for AI Model Development.

- 190 • The References section lists all references cited in this document.
- 191 • Appendix A provides a glossary of selected terms used within this document.

192 2. Using the SSDF Community Profile

193 AI model producers, AI system producers, AI system acquirers, and others can use the SSDF to
194 foster their communications regarding secure AI model development throughout the
195 development life cycle. Following SSDF practices should help AI model producers reduce the
196 number of vulnerabilities in their AI models, reduce the potential impacts of the exploitation of
197 undetected or unaddressed vulnerabilities, and address the root causes of vulnerabilities to
198 prevent recurrences. AI system producers can use the SSDF's common vocabulary when
199 communicating with AI model producers regarding their security practices for AI model
200 development and when integrating AI models into the software they are developing. AI system
201 acquirers can also use SSDF terms to better communicate their cybersecurity requirements and
202 needs to AI model producers and AI system producers, such as during acquisition processes.

203 The Profile is not a checklist to follow, but rather a starting point for planning and implementing
204 a risk-based approach to adopting secure software development practices involving AI models.
205 The contents of the Profile are meant to be adapted and customized, as not all practices and
206 tasks are applicable to all use cases. Organizations should adopt a risk-based approach to
207 determine what practices and tasks are relevant, appropriate, and effective to mitigate the
208 threats to their software development practices. Factors such as risk, cost, feasibility, and
209 applicability should be considered when deciding which practices and tasks to use and how
210 much time and resources to devote to each one. Cost models may need to be updated to
211 effectively consider the costs inherent to AI model development. A risk-based approach to
212 secure software development may change over time as an organization responds to new or
213 elevated capabilities and risks associated with an AI model or system.

214 The SSDF Community Profile's practices, tasks, recommendations, and considerations can be
215 integrated into machine learning operations (MLOps) along with other software assets within a
216 continuous integration/continuous delivery (CI/CD) pipeline.

217 The responsibility for implementing SSDF practices in the Profile may be shared among multiple
218 organizations. For example, an AI model could be produced by one organization and executed
219 within an AI system hosted by a second organization, which is then used by other organizations.
220 In these situations, there is likely a shared responsibility model involving the AI model producer,
221 AI system producer, and AI system acquirer. An AI system acquirer can establish an agreement
222 with an AI system producer and/or AI model producer that specifies which party is responsible
223 for each practice and task and how each party will attest to its conformance with the
224 agreement.

225 A limitation of the SSDF and this Profile is that they only address cybersecurity risk
226 management. There are many other types of risks to AI systems that organizations should
227 manage in accordance with cybersecurity risk as part of a mature enterprise risk management
228 program. NIST resources on identifying and managing other types of risk include:

- 229 • *AI Risk Management Framework (AI RMF)* [2] and the *NIST AI RMF Playbook* [11]
- 230 • *Adversarial Machine Learning: A Taxonomy and Terminology of Attacks and Mitigations*
231 [3]

- 232 • *Towards a Standard for Identifying and Managing Bias in Artificial Intelligence* [4]
- 233 • *Cybersecurity Supply Chain Risk Management Practices for Systems and Organizations*
234 [10]
- 235 • *NIST Privacy Framework: A Tool for Improving Privacy Through Enterprise Risk*
236 *Management, Version 1.0* [12]
- 237 • *Integrating Cybersecurity and Enterprise Risk Management (ERM)* [13]

238 3. SSDF Community Profile for AI Model Development

239 Table 1 defines the SSDF Community Profile for AI Model Development. The meanings of each
240 column are as follows:

241 • **Practice** contains the name of the practice and a unique identifier, followed by a brief
242 explanation of what the practice is and why it is beneficial.

243 **Task** specifies one or more actions that may be needed to perform a practice. Each task
244 includes a unique identifier and a brief explanation.

245 All practices and tasks are unchanged from SSDF version 1.1 unless they are explicitly
246 tagged as “Modified from SSDF 1.1” or “Not part of SSDF 1.1.” An example is the PW.3
247 practice, “Confirm the Integrity of Training, Testing, Fine-Tuning, and Aligning Data
248 Before Use” and all of its tasks.

249 • **Priority** reflects the suggested relative importance of each task *within the context of the*
250 *profile* and is intended to be a starting point for organizations to assign their own
251 priorities:

252 ○ **High:** Critically important for AI model development security compared to other
253 tasks

254 ○ **Medium:** Directly supports AI model development security

255 ○ **Low:** Beneficial for secure software development but is generally not more
256 important than most other tasks

257 • **Recommendations, Considerations, and Notes Specific to AI Model Development** may
258 contain one or more items that recommend what to do or describe additional
259 considerations for a particular task. Each item has an ID starting with one of the
260 following:

261 ○ “R” (recommendation: something the organization should do)

262 ○ “C” (consideration: something the organization should consider doing)

263 ○ “N” (note: additional information besides recommendations and considerations)

264 An R, C, or N designation and its number can be appended to the task ID to create a
265 unique identifier (e.g., “PO.1.2.R1” is the first recommendation for task PO.1.2).

266 Note that a value of “none” in this column indicates that the Profile does not contain
267 recommendations, considerations, or notes specific to AI model development for the
268 task. Refer to SSDF version 1.1 [6] for baseline guidance on the secure development task
269 in question and to the other references in this document for additional information
270 related to the task.

271 • **Informative References** point (map) to parts of standards, guidance, and other content
272 containing requirements, recommendations, considerations, or other supporting

273 information on performing a particular task. In this draft, the Informative References
274 come from two sources:

- 275 ○ *AI Risk Management Framework 1.0* [2]. Several crosswalks have already been
276 defined between the AI RMF and other guidance and standards; see
277 https://airc.nist.gov/AI_RM_F_Knowledge_Base/Crosswalks for the current set.
- 278 ○ *OWASP Top 10 for LLM Applications Version 1.1* [14]. Each identifier indicates
279 one of the top 10 vulnerability types and might also refer to an individual
280 prevention and mitigation strategy for that vulnerability type.

281 NIST plans to include additional Informative References in the final version of the
282 Profile.

283 NIST is also considering adding a column for Implementation Examples in the final Profile. An
284 **Implementation Example** is a single sentence that suggests a way to accomplish part or all of a
285 task. While the Recommendations and Considerations column describes the “what,”
286 Implementation Examples would describe options for the “how.” Such examples added to this
287 Profile would supplement those already defined in SSDF version 1.1. See the [Note to Reviewers](#)
288 for more information on providing input on additional Informative References and
289 Implementation Examples.

290 **Note: This Profile supplements what SSDF version 1.1 [6] already includes and is intended to be used in conjunction with it.**

291 There are gaps in the numbering of some SSDF practices and tasks. For example, the PW.4 practice has three tasks: PW.4.1, PW.4.2, and PW.4.4.

292 PW.4.3 was a task in SSDF version 1.0 that was moved elsewhere for version 1.1, so its ID was not reused.

Table 1. SSDF Community Profile for AI Model Development

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
Prepare the Organization (PO)				
Define Security Requirements for Software Development (PO.1): Ensure that security requirements for software development are known at all times so that they can be taken into account throughout the software development life cycle (SDLC) and duplication of effort can be minimized because the requirements information can be collected once and shared. This includes requirements from internal sources (e.g., the organization’s policies, business objectives, and risk management strategy) and external sources (e.g., applicable laws and regulations).	PO.1.1: Identify and document all security requirements for the organization’s software development infrastructures and processes, and maintain the requirements over time.	High	R1: Include AI model development in the security requirements for software development infrastructure and processes.	AI RMF: Map 1.3, 1.5, 1.6
	PO.1.2: Identify and document all security requirements for organization-developed software to meet, and maintain the requirements over time.	High	R1: Organizational policies should support all current requirements specific to AI model development security for organization-developed software. These requirements should include the areas of AI model development, AI model operations, and data science. Requirements may come from many sources, including laws, regulations, contracts, and standards. C1: Consider reusing or expanding the organization’s existing data classification policy and processes.	AI RMF: Govern 1.1, 1.2, 3.2, 4.1, 5.1, 6.1; Map 1.1
	PO.1.3: Communicate requirements to all third parties who will provide commercial software components to the organization for reuse by the organization’s own software.	Medium	R1: Include AI model development security in the requirements being communicated for third-party software components.	AI RMF: Map 4.1, 4.2 OWASP: LLM05-1
Implement Roles and Responsibilities (PO.2): Ensure that everyone inside and outside of the organization involved in the SDLC is prepared to perform their SDLC-related roles and responsibilities throughout the SDLC.	PO.2.1: Create new roles and alter responsibilities for existing roles as needed to encompass all parts of the SDLC. Periodically review and maintain the defined roles and responsibilities, updating them as needed.	High	R1: Include AI model development security in SDLC-related roles and responsibilities throughout the SDLC. The roles and responsibilities should include AI model development, AI model operations, and data science.	AI RMF: Govern 2.1

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
			N1: Roles and responsibilities involving AI system producers, AI model producers, and other third-party providers can be documented in agreements.	
	PO.2.2: Provide role-based training for all personnel with responsibilities that contribute to secure development. Periodically review personnel proficiency and role-based training, and update the training as needed.	High	R1: Role-based training should include understanding cybersecurity vulnerabilities and threats to AI models and their possible mitigations.	AI RMF: Govern 2.2 OWASP: LLM04-7
	PO.2.3: Obtain upper management or authorizing official commitment to secure development, and convey that commitment to all with development-related roles and responsibilities.	Medium	R1: Leadership should commit to secure development practices involving AI models.	AI RMF: Govern 2.3
Implement Supporting Toolchains (PO.3): Use automation to reduce human effort and improve the accuracy, reproducibility, usability, and comprehensiveness of security practices throughout the SDLC, as well as provide a way to document and demonstrate the use of these practices. Toolchains and tools may be used at different levels of the organization, such as organization-wide or project-specific, and may address a particular part of the SDLC, like a build pipeline.	PO.3.1: Specify which tools or tool types must or should be included in each toolchain to mitigate identified risks, as well as how the toolchain components are to be integrated with each other.	High	R1: Plan to develop and implement automated toolchains that secure AI model development and reduce human effort, especially at the scale often used by AI models. N1: Ideally, automated toolchains will perform the vast majority of the work related to securing AI model development. N2: See PO.4, PO.5, PS, and PW for information on tool types.	AI RMF: Measure 2.1 OWASP: LLM08
	PO.3.2: Follow recommended security practices to deploy, operate, and maintain tools and toolchains.	High	R1: Execute the plan to develop and implement automated toolchains that secure AI model development and reduce human effort, especially at the scale often used by AI models. R2: Verify the security of toolchains at a frequency commensurate with risk.	AI RMF: Measure 2.1 OWASP: LLM05-3, LLM05-9, LLM08, LLM09
	PO.3.3: Configure tools to generate artifacts of their support of secure software development practices as defined by the organization.	Medium	N1: An <i>artifact</i> is “a piece of evidence” [15]. <i>Evidence</i> is “grounds for belief or disbelief; data on which to base proof or to establish truth or falsehood” [16]. Artifacts provide	AI RMF: Measure 2.1

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
			records of secure software development practices. Examples of artifacts specific to AI model development include attestations of training dataset integrity and provenance.	
Define and Use Criteria for Software Security Checks (PO.4): Help ensure that the software resulting from the SDLC meets the organization’s expectations by defining and using criteria for checking the software’s security during development.	PO.4.1: Define criteria for software security checks and track throughout the SDLC.	Medium	C1: Consider requiring review and approval from a human-in-the-loop for software security checks beyond risk-based thresholds.	AI RMF: Measure 2.3, 2.7; Manage 1.1 OWASP: LLM01-2
	PO.4.2: Implement processes, mechanisms, etc. to gather and safeguard the necessary information in support of the criteria.	Low	None	AI RMF: Measure 2.3, 2.7; Manage 1.1 OWASP: LLM01-2
Implement and Maintain Secure Environments for Software Development (PO.5): Ensure that all components of the environments for software development are strongly protected from internal and external threats to prevent compromises of the environments or the software being developed or maintained within them. Examples of environments for software development include development, AI model training, build, test, and distribution environments. [Modified from SSDF 1.1]	PO.5.1: Separate and protect each environment involved in software development.	High	C1: Consider separating execution environments from each other to the extent feasible, such as by using sandboxing or containers. R1: Monitor, track, and limit resource usage and rates for AI model users.	OWASP: LLM01-1, LLM01-4, LLM04, LLM08, LLM10
	PO.5.2: Secure and harden development endpoints (endpoints for software designers, developers, testers, builders, etc.) to perform development tasks using a risk-based approach.	Medium	None	OWASP: LLM01-1, LLM05-3, LLM05-9, LLM08
	PO.5.3: Continuously monitor software execution performance and behavior in software development environments to identify potential suspicious activity and other issues. [Not part of SSDF 1.1]	High	R1: Perform continuous security and performance monitoring for all development environment components that host an AI model or related resources (e.g., model APIs, weights, configuration parameters, training datasets). R2: Continuous monitoring and analysis tools should generate alerts when detected activity involving an AI model passes a risk threshold or otherwise merits additional investigation.	AI RMF: Measure 2.4 OWASP: LLM03-7, LLM04, LLM05-8, LLM09, LLM10
Protect Software (PS)				

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
<p>Protect All Forms of Code and Data from Unauthorized Access and Tampering (PS.1): Help prevent unauthorized changes to code and data, both inadvertent and intentional, which could circumvent or negate the intended security characteristics of the software. For code and data that are not intended to be publicly accessible, this helps prevent theft of the software and may make it more difficult or time-consuming for attackers to find vulnerabilities in the software. [Modified from SSDF 1.1]</p>	<p>PS.1.1: Store all forms of code – including source code, executable code, and configuration-as-code – based on the principle of least privilege so that only authorized personnel, tools, services, etc. have access.</p>	<p>High</p>	<p>R1: Secure code storage should include AI models, model weights, pipelines, reward models, and any other AI model elements that need their confidentiality, integrity, and/or availability protected. R2: Follow the principle of least privilege to minimize direct access to AI models and model elements regardless of where they are stored or executed. R3: Store reward models separately from AI models and data. C1: Consider preventing all human access to model weights. C2: Consider requiring all AI model development to be performed within organization-approved environments only.</p>	<p>OWASP: LLM10</p>
	<p>PS.1.2: Protect all training, testing, fine-tuning, and aligning data from unauthorized access and modification. [Not part of SSDF 1.1]</p>	<p>High</p>	<p>R1: Continuously monitor the confidentiality and integrity of training, testing, fine-tuning, and aligning data. C1: Consider securely storing training, testing, fine-tuning, and aligning data for future use and reference if feasible.</p>	<p>OWASP: LLM03, LLM06, LLM10</p>
	<p>PS.1.3: Protect all model weights and configuration parameter data from unauthorized access and modification. [Not part of SSDF 1.1]</p>	<p>High</p>	<p>R1: Keep model weights and configuration parameters separate from training, testing, fine-tuning, and aligning data. R2: Continuously monitor the confidentiality (for closed models only) and integrity of model weights and configuration parameters. R3: Follow the principle of least privilege to restrict access to AI model weights, configuration parameters, and services during development. R4: Specify and implement additional risk-proportionate cybersecurity practices around model weights, such as encryption, multi-</p>	<p>OWASP: LLM10</p>

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
			party authorization, and air-gapped environments.	
Provide a Mechanism for Verifying Software Release Integrity (PS.2): Help software acquirers ensure that the software they acquire is legitimate and has not been tampered with.	PS.2.1: Make software integrity verification information available to software acquirers.	Medium	R1: Generate and provide cryptographic hashes or digital signatures for an AI model and its components.	OWASP: LLM05-6
Archive and Protect Each Software Release (PS.3): Preserve software releases in order to help identify, analyze, and eliminate vulnerabilities discovered in the software after release.	PS.3.1: Securely archive the necessary files and supporting data (e.g., integrity verification information, provenance data) to be retained for each software release.	Low	R1: Perform versioning and tracking for infrastructure tools (e.g., pre-processing, transforms, collection) that support dataset creation and model training.	OWASP: LLM10
	PS.3.2: Collect, safeguard, maintain, and share provenance data for all components of each software release (e.g., in a software bill of materials [SBOM], through Supply-chain Levels for Software Artifacts [SLSA]). [Modified from SSDF 1.1]	Medium	R1: Track the provenance of an AI model and its components, including the training libraries and frameworks used to build the model. C1: Consider disclosing the provenance of the training, testing, fine-tuning, and aligning data used for an AI model.	OWASP: LLM03-1, LLM05-4, LLM05-5, LLM10
Produce Well-Secured Software (PW)				
Design Software to Meet Security Requirements and Mitigate Security Risks (PW.1): Identify and evaluate the security requirements for the software; determine what security risks the software is likely to face during operation and how the software’s design and architecture should mitigate those risks; and justify any cases where risk-based analysis indicates that security requirements should be relaxed or waived. Addressing security requirements and risks during software design (secure by design) is key for improving software security and also helps improve development efficiency.	PW.1.1: Use forms of risk modeling – such as threat modeling, attack modeling, or attack surface mapping – to help assess the security risk for the software.	High	R1: Incorporate relevant AI model-specific vulnerability and threat types in risk modeling. Examples of these vulnerability and threat types include poisoning of training data, malicious code or other unwanted content in inputs and outputs, denial-of-service conditions, supply chain attacks, unauthorized information disclosure, and theft of AI model weights.	AI RMF: Govern 4.1, 4.2; Map 5.1; Measure 1.1; Manage 1.2, 1.3 OWASP: LLM01, LLM02, LLM03, LLM04, LLM05, LLM06, LLM07, LLM08, LLM09, LLM10
	PW.1.2: Track and maintain the software’s security requirements, risks, and design decisions.	Medium	None	AI RMF: Govern 4.1, 4.2; Map 2.1, 2.2, 2.3, 3.2, 3.3, 4.1, 4.2, 5.2; Manage 1.2, 1.3, 1.4

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
	PW.1.3: Where appropriate, build in support for using standardized security features and services (e.g., enabling software to integrate with existing log management, identity management, access control, and vulnerability management systems) instead of creating proprietary implementations of security features and services.	Medium	None	
Review the Software Design to Verify Compliance with Security Requirements and Risk Information (PW.2): Help ensure that the software will meet the security requirements and satisfactorily address the identified risk information.	PW.2.1: Have 1) a qualified person (or people) who were not involved with the design and/or 2) automated processes instantiated in the toolchain review the software design to confirm and enforce that it meets all of the security requirements and satisfactorily addresses the identified risk information.	Medium	None	AI RMF: Measure 2.7; Manage 1.1
Confirm the Integrity of Training, Testing, Fine-Tuning, and Aligning Data Before Use (PW.3): Prevent data that is likely to negatively impact the cybersecurity of the AI model from being consumed as part of AI model training, testing, fine-tuning, and aligning. [Not part of SSDF 1.1]	PW.3.1: Analyze data for signs of data poisoning, bias, homogeneity, and tampering before using it for AI model training, testing, fine-tuning, or aligning purposes, and mitigate the risks as necessary. [Not part of SSDF 1.1]	High	R1: Verify the provenance and integrity of all training, testing, fine-tuning, and aligning data before use. R2: Select and apply appropriate methods for analyzing and altering the training, testing, fine-tuning, and aligning data for an AI model. Examples of methods include anomaly detection, bias detection, data cleaning, data curation, data filtering, data sanitization, fact-checking, and noise reduction. C1: Consider using a human-in-the-loop to examine data, such as with exploratory data analysis techniques [17].	AI RMF: Measure 2.1; Manage 1.2, 1.3 OWASP: LLM03, LLM06
	PW.3.2: Track the provenance of all training, testing, fine-tuning, and aligning data used for an AI model. [Not part of SSDF 1.1]	Medium	None	AI RMF: Measure 2.1 OWASP: LLM03-1

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
	PW.3.3: Include adversarial samples in the training and testing data to improve attack detection. [Not part of SSDF 1.1]	Medium	None	OWASP: LLM03-6, LLM05-7
Reuse Existing, Well-Secured Software When Feasible Instead of Duplicating Functionality (PW.4): Lower the costs of software development, expedite software development, and decrease the likelihood of introducing additional security vulnerabilities into the software by reusing software modules and services that have already had their security posture checked. This is particularly important for software that implements security functionality, such as cryptographic modules and protocols.	PW.4.1: Acquire and maintain well-secured software components (e.g., software libraries, modules, middleware, frameworks) from commercial, open-source, and other third-party developers for use by the organization’s software.	Medium	C1: Consider using an existing AI model instead of creating a new one.	OWASP: LLM05
	PW.4.2: Create and maintain well-secured software components in-house following SDLC processes to meet common internal software development needs that cannot be better met by third-party software components.	Low	None	
	PW.4.4: Verify that acquired commercial, open-source, and all other third-party software components comply with the requirements, as defined by the organization, throughout their life cycles.	High	R1: Verify the integrity, provenance, and security of an existing AI model or any other acquired AI components — including training, testing, fine-tuning, and aligning datasets; reward models; adaptation layers; and configuration parameters — before using them. R2: Scan and thoroughly test acquired AI models and their components for vulnerabilities before use.	OWASP: LLM05-2, LLM05-6
Create Source Code by Adhering to Secure Coding Practices (PW.5): Decrease the number of security vulnerabilities in the software, and reduce costs by minimizing vulnerabilities introduced during source code creation that meet or exceed organization-defined vulnerability severity criteria.	PW.5.1: Follow all secure coding practices that are appropriate to the development languages and environment to meet the organization’s requirements.	High	R1: Expand secure coding practices to include AI technology-specific considerations. R2: Code the handling of inputs (including prompts and user data) and outputs carefully. All inputs and outputs should be logged, analyzed, and validated within the context of the AI model, and those with issues should be sanitized or dropped. R3: Encode inputs and outputs to prevent the execution of unauthorized code.	AI RMF: Manage 1.2, 1.3, 1.4 OWASP: LLM01, LLM02, LLM04-1, LLM06, LLM07, LLM09-9, LLM10

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
Configure the Compilation, Interpreter, and Build Processes to Improve Executable Security (PW.6): Decrease the number of security vulnerabilities in the software and reduce costs by eliminating vulnerabilities before testing occurs.	PW.6.1: Use compiler, interpreter, and build tools that offer features to improve executable security.	Low	None	
	PW.6.2: Determine which compiler, interpreter, and build tool features should be used and how each should be configured, then implement and use the approved configurations.	Low	None	
Review and/or Analyze Human-Readable Code to Identify Vulnerabilities and Verify Compliance with Security Requirements (PW.7): Help identify vulnerabilities so that they can be corrected before the software is released to prevent exploitation. Using automated methods lowers the effort and resources needed to detect vulnerabilities. Human-readable code includes source code, scripts, and any other form of code that an organization deems human-readable.	PW.7.1: Determine whether code <i>review</i> (a person looks directly at the code to find issues) and/or code <i>analysis</i> (tools are used to find issues in code, either in a fully automated way or in conjunction with a person) should be used, as defined by the organization.	Medium	R1: Code review and analysis policies or guidelines should include code for AI models and other related components. C1: Consider performing scans of AI model code in addition to testing the AI models.	
	PW.7.2: Perform the code review and/or code analysis based on the organization’s secure coding standards, and record and triage all discovered issues and recommended remediations in the development team’s workflow or issue tracking system.	High	R1: Scan all AI models for malware, vulnerabilities, backdoors, and other security issues in accordance with the organization’s code review and analysis policies or guidelines.	AI RMF: Measure 2.3, 2.7; Manage 1.1, 1.2, 1.3, 1.4 OWASP: LLM03-7d, LLM07-4
Test Executable Code to Identify Vulnerabilities and Verify Compliance with Security Requirements (PW.8): Help identify vulnerabilities so that they can be corrected before the software is released in order to prevent exploitation. Using automated methods lowers the effort and resources needed to detect vulnerabilities and improves traceability and repeatability. Executable code includes binaries, directly executed bytecode and source code, and any other form of code that an organization deems executable.	PW.8.1: Determine whether executable code testing should be performed to find vulnerabilities not identified by previous reviews, analysis, or testing and, if so, which types of testing should be used.	High	R1: Include AI models in code testing policies and guidelines. Several forms of code testing can be used for AI models, including unit testing, integration testing, penetration testing, red teaming, and adversarial testing.	
	PW.8.2: Scope the testing, design the tests, perform the testing, and document the results, including recording and triaging all discovered issues and recommended remediations in the development team’s workflow or issue tracking system.	High	R1: Test all AI models for vulnerabilities in accordance with the organization’s code testing policies or guidelines.	AI RMF: Measure 2.2, 2.3, 2.7; Manage 1.1, 1.2, 1.3, 1.4 OWASP: LLM03-7d, LLM05-7, LLM07-4
Configure Software to Have Secure Settings by Default (PW.9): Help improve the security of	PW.9.1: Define a secure baseline by determining how to configure each setting	Medium	None	AI RMF: Measure 2.7

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
the software at the time of installation to reduce the likelihood of the software being deployed with weak security settings, putting it at greater risk of compromise.	that has an effect on security or a security-related setting so that the default settings are secure and do not weaken the security functions provided by the platform, network infrastructure, or services.			
	PW.9.2: Implement the default settings (or groups of default settings, if applicable), and document each setting for software administrators.	Medium	N1: Documenting settings can be performed earlier in the process, such as when defining a secure baseline (see PW.9.1).	AI RMF: Measure 2.7; Manage 1.2, 1.3, 1.4
Respond to Vulnerabilities (RV)				
Identify and Confirm Vulnerabilities on an Ongoing Basis (RV.1): Help ensure that vulnerabilities are identified more quickly so that they can be remediated more quickly in accordance with risk, reducing the window of opportunity for attackers.	RV.1.1: Gather information from software acquirers, users, and public sources on potential vulnerabilities in the software and third-party components that the software uses, and investigate all credible reports.	High	R1: Log, monitor, and analyze all inputs and outputs for AI models to detect possible security and performance issues (see PO.5.3). R2: Make the users of AI models aware of mechanisms for reporting potential security and performance issues. R3: Monitor vulnerability and incident databases for information on AI-related concerns, including the machine learning frameworks and libraries used to build AI models.	AI RMF: Govern 4.3, 5.1, 6.1, 6.2; Measure 1.2, 2.4, 2.5, 2.7, 3.1, 3.2, 3.3; Manage 4.1 OWASP: LLM03-7a, LLM09, LLM10
	RV.1.2: Review, analyze, and/or test the software’s code to identify or confirm the presence of previously undetected vulnerabilities.	Medium	R1: Scan and test AI models frequently to identify previously undetected vulnerabilities. R2: Rely mainly on automation for ongoing scanning and testing, and involve a human-in-the-loop as needed. R3: Conduct periodic audits of AI models.	AI RMF: Govern 4.3; Measure 1.3, 2.4, 2.7, 3.1; Manage 4.1 OWASP: LLM03-7b, LLM03-7d
	RV.1.3: Have a policy that addresses vulnerability disclosure and remediation, and implement the roles, responsibilities, and processes needed to support that policy.	Medium	R1: Include AI model vulnerabilities in organization vulnerability disclosure and remediation policies. R2: Make users of AI models aware of their inherent limitations and how to report any cybersecurity problems that they encounter.	AI RMF: Govern 4.3, 5.1, 6.1; Measure 3.1, 3.3; Manage 4.3
Assess, Prioritize, and Remediate Vulnerabilities (RV.2): Help ensure that vulnerabilities are remediated in accordance	RV.2.1: Analyze each vulnerability to gather sufficient information about risk to	Medium	None	AI RMF: Govern 4.3, 5.1, 6.1; Measure 2.7,

Practice	Task	Priority	Recommendations, Considerations, and Notes Specific to AI Model Development	Informative References
with risk to reduce the window of opportunity for attackers.	plan its remediation or other risk response.			3.1; Manage 1.2, 2.3, 4.1
	RV.2.2: Plan and implement risk responses for vulnerabilities.	High	R1: Risk responses for AI models should consider the time and expenses that may be associated with rebuilding them. R2: Be prepared to roll back to a previous AI model, since that may be the most feasible response in some cases. C1: Consider being prepared to stop using an AI model at any time and to continue operations through other means until the AI model's risks are sufficiently addressed.	AI RMF: Govern 5.1, 5.2, 6.1; Measure 3.3; Manage 1.3, 2.1, 2.3, 2.4, 4.1
Analyze Vulnerabilities to Identify Their Root Causes (RV.3): Help reduce the frequency of vulnerabilities in the future.	RV.3.1: Analyze identified vulnerabilities to determine their root causes.	Medium	N1: The ability to review training, testing, fine-tuning, and aligning data after the fact can help identify some root causes.	AI RMF: Govern 5.1, 6.1; Measure 2.7, 3.1; Manage 2.3, 4.1
	RV.3.2: Analyze the root causes over time to identify patterns, such as a particular secure coding practice not being followed consistently.	Medium	None	AI RMF: Govern 5.1, 6.1; Measure 2.7, 3.1; Manage 4.1, 4.3
	RV.3.3: Review the software for similar vulnerabilities to eradicate a class of vulnerabilities, and proactively fix them rather than waiting for external reports.	Medium	None	AI RMF: Govern 5.1, 5.2, 6.1; Measure 2.7, 3.1; Manage 4.1, 4.2, 4.3
	RV.3.4: Review the SDLC process, and update it if appropriate to prevent (or reduce the likelihood of) the root cause recurring in updates to the software or in new software that is created.	Medium	None	AI RMF: Govern 5.2, 6.1; Measure 2.7, 3.1; Manage 4.2, 4.3

293 References

- 294 [1] Executive Order 14110 (2023) Safe, Secure, and Trustworthy Development and Use of
295 Artificial Intelligence. (The White House, Washington, DC), DCPD-202300949, October
296 30, 2022. Available at <https://www.govinfo.gov/app/details/DCPD-202300949>
- 297 [2] National Institute of Standards and Technology (2023) Artificial Intelligence Risk
298 Management Framework (AI RMF 1.0). (National Institute of Standards and Technology,
299 Gaithersburg, MD) NIST Artificial Intelligence (AI) Report, NIST AI 100-1.
300 <https://doi.org/10.6028/NIST.AI.100-1>
- 301 [3] Vassilev A, Oprea A, Fordyce A, Anderson H (2024) Adversarial Machine Learning: A
302 Taxonomy and Terminology of Attacks and Mitigations. (National Institute of Standards
303 and Technology, Gaithersburg, MD) NIST Artificial Intelligence (AI) Report, NIST AI 100-
304 2e2023. <https://doi.org/10.6028/NIST.AI.100-2e2023>
- 305 [4] Schwartz R, Vassilev A, Greene K, Perine L, Burt A, Hall P (2022) Towards a Standard for
306 Identifying and Managing Bias in Artificial Intelligence. (National Institute of Standards
307 and Technology, Gaithersburg, MD), NIST Special Publication (SP) 1270.
308 <https://doi.org/10.6028/NIST.SP.1270>
- 309 [5] NIST (2024) Dioptra. (National Institute of Standards and Technology, Gaithersburg,
310 MD.) Available at <https://pages.nist.gov/dioptra/>
- 311 [6] Souppaya MP, Scarfone KA, Dodson DF (2022) Secure Software Development
312 Framework (SSDF) Version 1.1: Recommendations for Mitigating the Risk of Software
313 Vulnerabilities. (National Institute of Standards and Technology, Gaithersburg, MD),
314 NIST Special Publication (SP) 800-218. <https://doi.org/10.6028/NIST.SP.800-218>
- 315 [7] Executive Order 14028 (2021) Improving the Nation’s Cybersecurity. (The White House,
316 Washington, DC), DCPD-202100401, May 12, 2021. Available at
317 <https://www.govinfo.gov/app/details/DCPD-202100401>
- 318 [8] National Institute of Standards and Technology (2024) The NIST Cybersecurity
319 Framework (CSF) 2.0 (National Institute of Standards and Technology, Gaithersburg,
320 MD). <https://doi.org/10.6028/NIST.CSWP.29>
- 321 [9] Joint Task Force (2020) Security and Privacy Controls for Information Systems and
322 Organizations. (National Institute of Standards and Technology, Gaithersburg, MD), NIST
323 Special Publication (SP) 800-53, Rev. 5. Includes updates as of December 10, 2020.
324 <https://doi.org/10.6028/NIST.SP.800-53r5>
- 325 [10] Boyens JM, Smith AM, Bartol N, Winkler K, Holbrook A, Fallon M (2022) Cybersecurity
326 Supply Chain Risk Management Practices for Systems and Organizations. (National
327 Institute of Standards and Technology, Gaithersburg, MD), NIST Special Publication (SP)
328 800-161r1. <https://doi.org/10.6028/NIST.SP.800-161r1>
- 329 [11] NIST (2023) NIST AI RMF Playbook. (National Institute of Standards and Technology,
330 Gaithersburg, MD.) Available at
331 https://airc.nist.gov/AI_RM_F_Knowledge_Base/Playbook
- 332 [12] National Institute of Standards and Technology (2020) NIST Privacy Framework: A Tool
333 for Improving Privacy Through Enterprise Risk Management, Version 1.0. (National

- 334 Institute of Standards and Technology, Gaithersburg, MD), NIST Cybersecurity White
335 Paper (CSWP) NIST CSWP 10. <https://doi.org/10.6028/NIST.CSWP.10>
- 336 [13] Stine KM, Quinn SD, Witte GA, Gardner RK (2020) Integrating Cybersecurity and
337 Enterprise Risk Management (ERM). (National Institute of Standards and Technology,
338 Gaithersburg, MD), NIST Interagency or Internal Report (IR) 8286.
339 <https://doi.org/10.6028/NIST.IR.8286>
- 340 [14] OWASP (2023) OWASP Top 10 for LLM Applications Version 1.1. Available at
341 <https://llmtop10.com>
- 342 [15] Waltermire DA, Scarfone KA, Casipe M (2011) Specification for the Open Checklist
343 Interactive Language (OCIL) Version 2.0. (National Institute of Standards and
344 Technology, Gaithersburg, MD), NIST Interagency or Internal Report (IR) 7692.
345 <https://doi.org/10.6028/NIST.IR.7692>
- 346 [16] Ross RS, McEvilley M, Winstead M (2022) Engineering Trustworthy Secure Systems.
347 (National Institute of Standards and Technology, Gaithersburg, MD), NIST Special
348 Publication (SP) NIST SP 800-160v1r1. <https://doi.org/10.6028/NIST.SP.800-160v1r1>
- 349 [17] NIST/SEMATECH (2012) What is EDA? *Engineering Statistics Handbook*, eds Croarkin C,
350 Tobias P (National Institute of Standards and Technology, Gaithersburg, MD), Section
351 1.1.1. Available at <https://www.itl.nist.gov/div898/handbook/eda/section1/eda11.htm>

352 **Appendix A. Glossary**

353 **artificial intelligence**

354 A machine-based system that can, for a given set of human-defined objectives, make predictions,
355 recommendations, or decisions influencing real or virtual environments. [1]

356 **artificial intelligence model**

357 A component of an information system that implements AI technology and uses computational, statistical, or
358 machine-learning techniques to produce outputs from a given set of inputs. [1]

359 **artificial intelligence red-teaming**

360 A structured testing effort to find flaws and vulnerabilities in an AI system, often in a controlled environment and
361 in collaboration with developers of AI. [1]

362 **artificial intelligence system**

363 Any data system, software, hardware, application, tool, or utility that operates in whole or in part using AI. [1]

364 **dual-use foundation model**

365 An AI model that is trained on broad data; generally uses self-supervision; contains at least tens of billions of
366 parameters; is applicable across a wide range of contexts; and that exhibits, or could be easily modified to exhibit,
367 high levels of performance at tasks that pose a serious risk to security, national economic security, national public
368 health or safety, or any combination of those matters, such as by:

369 (i) substantially lowering the barrier of entry for non-experts to design, synthesize, acquire, or use
370 chemical, biological, radiological, or nuclear (CBRN) weapons;

371 (ii) enabling powerful offensive cyber operations through automated vulnerability discovery and
372 exploitation against a wide range of potential targets of cyber attacks; or

373 (iii) permitting the evasion of human control or oversight through means of deception or obfuscation.

374 Models meet this definition even if they are provided to end users with technical safeguards that attempt to
375 prevent users from taking advantage of the relevant unsafe capabilities. [1]

376 **generative artificial intelligence**

377 The class of AI models that emulate the structure and characteristics of input data in order to generate derived
378 synthetic content. This can include images, videos, audio, text, and other digital content. [1]

379 **model weight**

380 A numerical parameter within an AI model that helps determine the model's outputs in response to inputs. [1]

381 **provenance**

382 Metadata pertaining to the origination or source of specified data. [12]