



Securing Artificial Intelligence (SAI); Problem Statement

(<https://standards.iteh.ai>)
Document Preview

[ETSI TR 104 221 V1.1.1 \(2025-01\)](https://standards.iteh.ai/catalog/standards/etsi/47df8e66-7723-478b-93fc-473fc359f252/etsi-tr-104-221-v1-1-1-2025-01)

<https://standards.iteh.ai/catalog/standards/etsi/47df8e66-7723-478b-93fc-473fc359f252/etsi-tr-104-221-v1-1-1-2025-01>

Reference
DTR/SAI-008
Keywords
artificial intelligence, security

ETSI

650 Route des Lucioles
F-06921 Sophia Antipolis Cedex - FRANCE

Tel.: +33 4 92 94 42 00 Fax: +33 4 93 65 47 16

Siret N° 348 623 562 00017 - APE 7112B
Association à but non lucratif enregistrée à la
Sous-Préfecture de Grasse (06) N° w061004871

Important notice

The present document can be downloaded from the
[ETSI Search & Browse Standards](#) application.

The present document may be made available in electronic versions and/or in print. The content of any electronic and/or print versions of the present document shall not be modified without the prior written authorization of ETSI. In case of any existing or perceived difference in contents between such versions and/or in print, the prevailing version of an ETSI deliverable is the one made publicly available in PDF format on [ETSI deliver](#) repository.

Users should be aware that the present document may be revised or have its status changed,
this information is available in the [Milestones listing](#).

If you find errors in the present document, please send your comments to
the relevant service listed under [Committee Support Staff](#).

If you find a security vulnerability in the present document, please report it through our
[Coordinated Vulnerability Disclosure \(CVD\)](#) program.

<https://standards.iteh.ai/catalog/standards/etsi/47d8-66-7723-478b-92fc-4726-259f252/etsi-tr-104-221-v1-1-1-2025-01>

Notice of disclaimer & limitation of liability

The information provided in the present deliverable is directed solely to professionals who have the appropriate degree of experience to understand and interpret its content in accordance with generally accepted engineering or other professional standard and applicable regulations.

No recommendation as to products and services or vendors is made or should be implied.

No representation or warranty is made that this deliverable is technically accurate or sufficient or conforms to any law and/or governmental rule and/or regulation and further, no representation or warranty is made of merchantability or fitness for any particular purpose or against infringement of intellectual property rights.

In no event shall ETSI be held liable for loss of profits or any other incidental or consequential damages.

Any software contained in this deliverable is provided "AS IS" with no warranties, express or implied, including but not limited to, the warranties of merchantability, fitness for a particular purpose and non-infringement of intellectual property rights and ETSI shall not be held liable in any event for any damages whatsoever (including, without limitation, damages for loss of profits, business interruption, loss of information, or any other pecuniary loss) arising out of or related to the use of or inability to use the software.

Copyright Notification

No part may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm except as authorized by written permission of ETSI.

The content of the PDF version shall not be modified without the written authorization of ETSI.

The copyright and the foregoing restriction extend to reproduction in all media.

© ETSI 2025.
All rights reserved.

Contents

Intellectual Property Rights	5
Foreword.....	5
Modal verbs terminology.....	5
1 Scope	6
2 References	6
2.1 Normative references	6
2.2 Informative references.....	6
3 Definition of terms, symbols and abbreviations.....	8
3.1 Terms.....	8
3.2 Symbols.....	8
3.3 Abbreviations	8
4 Context	9
4.1 History	9
4.2 AI and machine learning	9
4.3 Data processing chain (machine learning).....	10
4.3.1 Overview	10
4.3.2 Data Acquisition	11
4.3.2.1 Description.....	11
4.3.2.2 Integrity challenges	11
4.3.3 Data Curation.....	12
4.3.3.1 Description.....	12
4.3.3.2 Integrity challenges	12
4.3.4 Model Design.....	12
4.3.5 Software Build	12
4.3.6 Training	12
4.3.6.1 Description	12
4.3.6.2 Confidentiality challenges.....	12
4.3.6.3 Integrity challenges	13
4.3.6.4 Availability challenges	13
4.3.7 Testing	13
4.3.7.1 Description	13
4.3.7.2 Availability challenges	14
4.3.8 Deployment and Inference	14
4.3.8.1 Description	14
4.3.8.2 Confidentiality challenges	14
4.3.8.3 Integrity challenges	14
4.3.8.4 Availability challenges	15
4.3.9 Upgrades	15
4.3.9.1 Description	15
4.3.9.2 Integrity challenges	15
4.3.9.3 Availability challenges	15
5 Design challenges and unintentional factors	15
5.1 Introduction	15
5.2 Bias.....	15
5.3 Ethics	16
5.3.1 Introduction.....	16
5.3.2 Ethics and security challenges	16
5.3.2.1 Access to data.....	16
5.3.2.2 Decision-making	17
5.3.2.3 Obscurity	17
5.3.2.4 Summary	17
5.3.3 Ethics guidelines	17
5.4 Explainability (explicability).....	18
5.5 Software and hardware	18

6	Attack types.....	19
6.1	Poisoning.....	19
6.2	Input attack and evasion	19
6.3	Backdoor Attacks	19
6.4	Reverse Engineering.....	20
7	Misuse of AI.....	20
8	Real world use cases and attacks.....	20
8.1	Overview	20
8.2	Ad-blocker attacks.....	20
8.3	Malware Obfuscation	21
8.4	Deepfakes	21
8.5	Handwriting reproduction	21
8.6	Human voice	21
8.7	Fake conversation.....	22
Annex A:	Bibliography	23
History	24	

i T h S t a n d a r d s

(h t t p s : / / s t a n d a r d s . i t e h . a i / c a t a l o g u e s)

D o c u m e n t i e P w r

E TTSRI 1 V0 14. 12.211 (2 0 2 5 - 0 1)

h t t p s : / / s t a n d a r d s . i t e h . a i / c a t a l o g u e s)

Intellectual Property Rights

Essential patents

IPRs essential or potentially essential to normative deliverables may have been declared to ETSI. The declarations pertaining to these essential IPRs, if any, are publicly available for **ETSI members and non-members**, and can be found in ETSI SR 000 314: *"Intellectual Property Rights (IPRs); Essential, or potentially Essential, IPRs notified to ETSI in respect of ETSI standards"*, which is available from the ETSI Secretariat. Latest updates are available on the [ETSI IPR online database](#).

Pursuant to the ETSI Directives including the ETSI IPR Policy, no investigation regarding the essentiality of IPRs, including IPR searches, has been carried out by ETSI. No guarantee can be given as to the existence of other IPRs not referenced in ETSI SR 000 314 (or the updates on the ETSI Web server) which are, or may be, or may become, essential to the present document.

Trademarks

The present document may include trademarks and/or tradenames which are asserted and/or registered by their owners. ETSI claims no ownership of these except for any which are indicated as being the property of ETSI, and conveys no right to use or reproduce any trademark and/or tradename. Mention of those trademarks in the present document does not constitute an endorsement by ETSI of products, services or organizations associated with those trademarks.

DECT™, PLUGTESTS™, UMTS™ and the ETSI logo are trademarks of ETSI registered for the benefit of its Members. **3GPP™, LTE™** and **5G™** logo are trademarks of ETSI registered for the benefit of its Members and of the 3GPP Organizational Partners. **oneM2M™** logo is a trademark of ETSI registered for the benefit of its Members and of the oneM2M Partners. **GSM®** and the **GSM** logo are trademarks registered and owned by the GSM Association.

Foreword

<https://standards.iten.ai>

Document Preview

This Technical Report (TR) has been produced by ETSI Technical Committee Securing Artificial Intelligence (SAI).

In the present document "**should**", "**should not**", "**may**", "**need not**", "**will**", "**will not**", "**can**" and "**cannot**" are to be interpreted as described in clause 3.2 of the [ETSI Drafting Rules](#) (Verbal forms for the expression of provisions).

"**must**" and "**must not**" are **NOT** allowed in ETSI deliverables except when used in direct citation.

1 Scope

The present document describes the problem of securing AI-based systems and solutions, with a focus on machine learning, and the challenges relating to confidentiality, integrity and availability at each stage of the machine learning lifecycle. It also describes some of the broader challenges of AI systems including bias, ethics and explainability. A number of different attack vectors are described, as well as several real-world use cases and attacks.

NOTE: The present document updates and replaces ETSI GR SAI 004 [i.32].

2 References

2.1 Normative references

Normative references are not applicable in the present document.

2.2 Informative references

References are either specific (identified by date of publication and/or edition number or version number) or non-specific. For specific references, only the cited version applies. For non-specific references, the latest version of the referenced document (including any amendments) applies.

NOTE: While any hyperlinks included in this clause were valid at the time of publication, ETSI cannot guarantee their long term validity.

The following referenced documents are not necessary for the application of the present document but they assist the user with regard to a particular subject area.

- [i.1] Florian Tramèr, Pascal Dupré, Gili Rusak, Giancarlo Pellegrino, Dan Boneh: "[Adversarial Perceptual Ad Blocking meets Adversarial Machine Learning](#)", in Proceedings of the 2019, ACM SIGSAC Conference on Computer and Communications Security, pp. 2005-2021, November 2019.
- [i.2] Stuart Millar, Niall McLaughlin, Jesus Martinez del Rincon, Paul Miller, Ziming Zhao: "[DANDROID: A Multi-View Discriminative Adversarial Network for Obfuscated Android Malware Detection](#)", in Proceedings of the 10th ACM Conference on Data and Application Security and Privacy, 2019.
- [i.3] Leslie D.: "[Understanding artificial intelligence ethics and safety: A guide for the responsible design and implementation of AI systems in the public sector](#)", The Alan Turing Institute, 2019.
- [i.4] High Level Expert Group on Artificial Intelligence, European Commission: "Ethics Guidelines for Trustworthy AI", April 2019.
- [i.5] UK Department for Digital, Culture, Media & Sport: "Data Ethics Framework", August 2018.
- [i.6] Song C., Ristenpart T., and Shmatikov V.: "Machine Learning Models that Remember Too Much", ACM CCS 17, Dallas, TX, USA.
- [i.7] Finn C., Abbeel P., Levine S.: "[Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks](#)".
- [i.8] Chen X., Liu C., Li B., Lu K., Song D.: "[Targeted Backdoor Attacks on Deep Learning Systems Using Data Poisoning](#)".
- [i.9] Tom S. F. Haines, Oisin Mac Aodha, and Gabriel J. Brostow: "[My Text in Your Handwriting](#)", ACM Trans. Graph. 35, 3, Article 26 (June 2016), 18 pages.
- [i.10] K. Eykholt et al.: "[Robust Physical-World Attacks on Deep Learning Visual Classification](#)", 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, 2018, pp. 1625-1634.

[i.11] Florian Tramèr, Fan Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart: "Stealing machine learning models via prediction APIs", in Proceedings of the 25th USENIX Conference on Security Symposium (SEC'16). USENIX Association, USA, pp. 601-618, 2016.

[i.12] Seong Joon Oh, Max Augustin, Bernt Schiele, Mario Fritz: "Towards reverse-engineering black-box neural networks Max-Planck Institute for Informatics", Saarland Informatics Campus, Saarbrücken, Germany Published as a conference paper at ICLR 2018.

[i.13] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, Koray Kavukcuoglu: "[WaveNet: A Generative Model for Raw Audio](#)", September 2016.

[i.14] Miles Brundage, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allan Dafoe, Paul Scharre, Thomas Zeitzoff, Bobby Filar, Hyrum Anderson, Heather Roff, Gregory C. Allen, Jacob Steinhardt, Carrick Flynn, Seán Ó hÉigeartaigh, Simon Beard, Haydn Belfield, Sebastian Farquhar, Clare Lyle, Rebecca Crootof, Owain Evans, Michael Page, Joanna Bryson, Roman Yampolskiy, Dario Amodei: "[The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation](#)".

[i.15] Oscar Schwarz, IEEE Tech Talk: "[Artificial Intelligence, Machine Learning](#)", November 2019.

[i.16] Haberer, J. et al.: "Gutachten der Datenethikkommission", 2019.

[i.17] Hagendorff T.: "[The Ethics of AI Ethics: An Evaluation of Guidelines](#)", Minds & Machines, vol. 30, pp. 99-120, 2020.

[i.18] Uesato J., Kumar A., Szepesvari C., Erez T., Ruderman A., Anderson K., Heess N. and Kohli P.: "Rigorous agent evaluation: An adversarial approach to uncover catastrophic failures", 2018, arXiv preprint arXiv:1812.01647.

[i.19] Weng T.W., Zhang H., Chen H., Song Z., Hsieh C.J., Boning D., Dhillon I.S. and Daniel L.: "Towards fast computation of certified robustness for relu networks", 2018, arXiv preprint arXiv:1804.09699.

[i.20] Kingston J. K. C.: "[Artificial Intelligence and Legal Liability](#)", 2018.

[i.21] Won-Suk Lee, Sung Min Ahn, Jun-Won Chung, Kyoung Oh Kim, Kwang An Kwon, Yoonjae Kim, Sunjin Sym, Dongbok Shin, Inkeun Park, Uhn Lee, and Jeong-Heum Baek.: "[Assessing Concordance with Watson for Oncology, a Cognitive Computing Decision Support System for Colon Cancer Treatment in Korea](#)", JCO Clinical Cancer Informatics 2018.

[i.22] Pr. Ronald C. Arkin: "[The Case for Ethical Autonomy in Unmanned Systems, Journal of Military Ethics](#)", 9:4, pp. 332-341, 2010.

[i.23] Pega Systems: "[What Consumers Really Think About AI: A Global Study](#)", 2017.

[i.24] Reza Shokri, Marco Stronati, Congzheng Song, Vitaly Shmatikov: "Membership Inference Attacks Against Machine Learning Models", IEEE security and privacy, 2017.

[i.25] Matt Fredrikson, Somesh Jha, Thomas Ristenpart: "Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures", ACM CCS 2015.

[i.26] Association for Computing Machinery: "[Top Two Levels of The ACM Computing Classification System \(1998\)](#)".

[i.27] Yim J., Chopra R., Spitz T., Winkens J., Obika A., Kelly C., Askham H., Lukic M., Huemer J., Fasler K. and Moraes G.: "[Predicting conversion to wet age-related macular degeneration using deep learning](#)", Nature Medicine, pp.1-8, 2020.

[i.28] McKinney S.M., Sieniek M., Godbole V., Godwin J., Antropova N., Ashrafian H., Back T., Chesus M., Corrado G.C., Darzi A. and Etemadi M.: "[International evaluation of an AI system for breast cancer screening](#)", Nature, 577(7788), pp. 89-94, 2020.

[i.29] Massachusetts Institute of Technology (MIT): "[Moral Machine](#)".

- [i.30] Organisation for Economic Co-operation and Development (OECD): "[Council recommendation on Artificial Intelligence](#)".
- [i.31] [Someone built chatbots that talk like the characters from HBO's Silicon Valley](#).
- [i.32] ETSI GR SAI 004 (2020-12): "Securing Artificial Intelligence (SAI); Problem Statement".
- [i.33] [Regulation \(EU\) 2024/1689](#) of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828 (Artificial Intelligence Act).
- [i.34] ENISA: "[Foresight 2030 Threats](#)".

3 Definition of terms, symbols and abbreviations

3.1 Terms

For the purposes of the present document, the following terms apply:

artificial intelligence: ability of a system to handle representations, both explicit and implicit, and procedures to perform tasks that would be considered intelligent if performed by a human

availability: property of being accessible and usable on demand by an authorized entity

confidentiality: assurance that information is accessible only to those authorized to have access

full knowledge attack: attack carried out by an attacker who has full knowledge of the system inputs and outputs and its internal design and operations

integrity: assurance of the accuracy and completeness of information and processing methods

opaque system: system or object which can be viewed solely in terms of its input, output and transfer characteristics without any knowledge of its internal workings

partial knowledge attack: attack carried out by an attacker who has full knowledge of the system inputs and outputs, but only a limited understanding of its internal design and operations

zero knowledge attack: attack carried out by an attacker who has knowledge of the system inputs and outputs, but no knowledge about its internal design or operations

3.2 Symbols

Void.

3.3 Abbreviations

For the purposes of the present document, the following abbreviations apply:

ACM	Association for Computing Machinery
AI	Artificial Intelligence
ASIC	Application Specific Integrated Circuit
CCS	Cascading Style Sheet
CCTV	Closed Circuit Television
CNN	Convolutional Neural Network
CVF	Computer Vision Foundation
DAN	Discriminative Adversarial Network
EPFL	École Polytechnique Fédérale de Lausanne
FG	Focus Group
FPGA	Field Programmable Gate Array

GPU	Graphics Processing Unit
HTML	Hyper Text Markup Language
IEEE	Institute of Electrical and Electronics Engineers
ITU	International Telecommunications Union
MIT	Massachusetts Institute of Technology
ML5G	Machine Learning for Future Networks including 5G
NIST	National Institute of Standards and Technology (USA)
OECD	Organisation for Economic Co-operation and Development
OS	Operating System
RNN	Recurrent Neural Network
TEE	Trusted Execution Environment
UN	United Nations
URL	Uniform Resource Locator

4 Context

4.1 History

The term 'artificial intelligence' originated at a conference in the 1950s at Dartmouth College in Hanover, New Hampshire, USA. At that time, it was suggested that true artificial intelligence could be created within a generation. By the early 1970s, despite millions of dollars of investment, it became clear that the complexity of creating true artificial intelligence was much greater than anticipated, and investment began to drop off. The years that followed are often referred to as an 'AI winter' which saw little interest or investment in the field, until the early 1980s when another wave of investment kicked off. By the late 1980s, interest had again waned, largely due to the absence of sufficient computing capacity to implement systems, and there followed a second AI winter.

In recent years, interest and investment in AI has once again surfaced, due to the implementation of some practical AI systems enabled by:

- The evolution of advanced techniques in machine learning, neural networks and deep learning.
- The availability of significant data sets to enable robust training.
- Advances in high performance computing enabling rapid training and development.
- Advances in high-performance devices enabling practical implementation.

After the emergence of practical AI systems, suggested theoretical attacks on such systems have become plentiful. Increasingly, real-world practical attacks with sufficient motivation and impact are being observed, particularly those using forms of generative AI and for masquerade (e.g. deep fakes) and are forecast as increasing in impact and viability over the coming years [i.34].

4.2 AI and machine learning

The field of artificial intelligence is broad, so in order to identify the issues in securing AI, the first step is to define what AI means.

The breadth of the field creates a challenge when trying to create accurate definitions.

EXAMPLE 1: The Association for Computing Machinery (ACM) Computing Classification System [i.26] breaks down Artificial Intelligence into eleven different categories, each of which has multiple sub-categories.

EXAMPLE 2: The AI Act [i.33] does not directly define Artificial Intelligence but rather defines an 'AI system' as a machine-based system that is designed to operate with varying levels of autonomy and that may exhibit adaptiveness after deployment, and that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments.

This represents a complex classification system with a large group of technology areas at varying stages of maturity, some of which have not yet seen real implementations, but does not serve as a useful concise definition. For the purposes of the present document, the following outline definition is used:

- **Artificial intelligence** is the ability of a system to handle representations, both explicit and implicit, and procedures to perform tasks that would be considered intelligent if performed by a human.

NOTE: The above definition is consistent with the definition of AI system found in the AI Act [i.33] on the understanding that the tasks performed can influence physical or virtual environments.

This definition still represents a broad spectrum of possibilities. However, there are a limited set of technologies which are now becoming realizable, largely driven by the evolution of machine learning and deep learning techniques. Therefore, the present document focusses on the discipline of machine learning and some of its variants, including:

- **Supervised learning** - where all the training data is labelled, and the model can be trained to predict the output based on a new set of inputs.
- **Semi-supervised learning** - where the data set is partially labelled. In this case, even the unlabelled data can be used to improve the quality of the model.
- **Unsupervised learning** - where the data set is unlabelled, and the model looks for structure in the data, including grouping and clustering.
- **Reinforcement learning** - where a policy defining how to act is learned by agents through experience to maximize their reward; and agents gain experience by interacting in an environment through state transitions.

Within each of these machine learning paradigms, there are various model structures that might be used, with one of the most common approaches being the use of deep neural networks, where learning is carried out over a series of hierarchical layers that mimic the behaviour of the human brain.

There are also a number of different training techniques which can be used, including adversarial learning, where the training set contains not only samples which reflect the desired outcomes, but also adversarial samples, which are intended to challenge or disrupt the expected behaviour.

Document Preview

4.3 Data processing chain (machine learning)

[ETSI TR 104 221 V1.1.1 \(2025-01\)](https://standards.iteh.it/octa/1/e/standards/etsi/47df8e66-7723-478b-93fc-473fc359f252/etsi-tr-104-221-v1-1-1-2025-01)

4.3.1 Overview

The question of securing AI systems can be simply stated as ensuring the confidentiality, integrity and availability of those systems throughout their lifecycle. The lifecycle for machine learning can be considered to have the following stages, as shown in Figure 1:

- 1) Data acquisition
- 2) Data curation
- 3) Model design
- 4) Build
- 5) Train
- 6) Test
- 7) Deployment
- 8) Results
- 9) Updates

Stages 4), 5) and 6) (Build, Train, Test) can together be considered as an iterative implementation cycle.

In the machine learning lifecycle, the training phase can be considered as the most critical, since it is this stage that establishes the baseline behaviour of the system.