

ETSI TS 104 050 V1.1.1 (2025-03)



Securing Artificial Intelligence (SAI); AI Threat Ontology and definitions

Reference

RTS/SAI-005

Keywords

artificial intelligence

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.

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	4
Foreword.....	4
Modal verbs terminology.....	4
1 Scope	5
2 References	5
2.1 Normative references	5
2.2 Informative references.....	5
3 Definition of terms, symbols and abbreviations.....	6
3.1 Terms.....	6
3.2 Symbols.....	7
3.3 Abbreviations	7
4 From taxonomy to an ontology for secure AI	7
4.1 Overview	7
4.2 Formal expression of an ontology	10
4.3 Relationship to other work	11
5 Threat landscape.....	13
5.1 Threat dimensions	13
5.2 Attacks as instance of threat agent	14
5.3 Adversarial Goals	14
5.3.1 Violation of Confidentiality.....	14
5.3.2 Violation of Integrity and Availability.....	14
5.4 Threat modelling	15
5.4.1 Attacker objectives	15
5.4.2 Attack surface	15
5.4.2.1 AI effect on impact and likelihood.....	15
5.4.2.2 Data acquisition and curation.....	16
5.4.2.3 Implementation	16
5.4.2.4 Deployment.....	16
5.4.2.5 Humans	16
5.4.3 Trust model.....	17
5.5 Statistics in AI and ML	17
6 AI and SAI ontology	18
6.1 Nouns, verbs, adverbs and adjectives.....	18
6.2 Taxonomy and ontology.....	18
6.3 Core SAI ontology relationships	19
Annex A (informative): Cultural origins of ICT based intelligence.....	22
Annex B (informative): Machine processing to simulate intelligence.....	25
B.1 Overview of the machine intelligence continuum.....	25
B.2 Expert systems.....	25
B.3 Data mining and pattern extraction	25
Annex C (informative): Bibliography.....	26
C.1 AGI analysis.....	26
C.2 AI in the context of threat analysis.....	27
C.3 Societal and cultural references to AI	27
History	28

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

This Technical Specification (TS) has been produced by ETSI Technical Committee Securing Artificial Intelligence (SAI).

NOTE: The present document updates and extends ETSI GR SAI 001 [i.20] prepared by ISG SAI.

Modal verbs terminology

In the present document "**shall**", "**shall not**", "**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 defines what an Artificial Intelligence (AI) threat is and defines how it can be distinguished from any non-AI threat. The model of an AI threat is presented in the form of an ontology to give a view of the relationships between actors representing threats, threat agents, assets and so forth and defines those terms (see also [1]). The ontology in the present document extends from the base taxonomy of threats and threat agents described in ETSI TS 102 165-1 [2] and addresses the overall problem statement for SAI presented in ETSI TR 104 221 [i.21] and the mitigation strategies described in ETSI TR 104 222 [i.22]. Note that, although both technical reports are listed in clause 2.2, they are indeed essential for understanding the scope of the present document.

NOTE 1: The ontology described in the present document applies to AI both as a threat agent and as an attack target.

NOTE 2: The present document extends the content of ETSI GR SAI 001 [i.20], and retains significant elements of its content where relevant for clarity.

2 References

2.1 Normative 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.

Referenced documents which are not found to be publicly available in the expected location might be found in the [ETSI docbox](#).

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 necessary for the application of the present document.

[1] [ISO/IEC 22989:2022](#): "Information technology - Artificial intelligence - Artificial intelligence concepts and terminology".

NOTE: Many of the terms defined in the cited document above are also visible on the ISO Online Browsing Platform: <https://www.iso.org/obp/>.

[2] [ETSI TS 102 165-1](#): "Cyber Security (CYBER); Methods and protocols; Part 1: Method and pro forma for Threat, Vulnerability, Risk Analysis (TVRA)".

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] Alan Turing: "On computable numbers, with an application to the Entscheidungsproblem".

[i.2] Alan Turing: "Computing Machinery and Intelligence".

[i.3] Philip K. Dick: "Do androids dream of electric sheep?" (ISBN-13: 978-0575094185).

[i.4] Isaac Asimov: "I, robot" (ISBN-13: 978-0008279554).

- [i.5] W3C[®] Recommendation 11 December 2012: "OWL: OWL 2 Web Ontology Language Document Overview (Second Edition)".
- [i.6] RDF: RDF 1.1 Primer; W3C[®] Working Group Note; 24 June 2014.
- [i.7] Cohen, Jacob (1960): "A coefficient of agreement for nominal scales". Educational and Psychological Measurement. 20 (1): 37-46. doi:10.1177/001316446002000104. hdl:1942/28116. S2CID 15926286.
- [i.8] W3C[®] Recommendation 16 July 2020: "JSON-LD 1.1: A JSON-based Serialization for Linked Data".
- [i.9] ETSI GS CIM 009 (V1.2.2): "Context Information Management (CIM); NGSI-LD API".
- [i.10] ["The Emergence Of Offensive AI"](#).
- [i.11] ["Weaponizing Data Science for Social Engineering: Automated E2E Spear Phishing on Twitter"](#).
- [i.12] Li Chen, Chih-Yuan Yang, Anindya Paul, Ravi Sahita: ["Towards resilient machine learning for ransomware detection"](#).
- [i.13] Alejandro Correa Bahnsen, Ivan Torroledo, Luis David Camacho and Sergio Villegas: ["DeepPhish: Simulating Malicious AI"](#).
- [i.14] [Common Weakness Enumeration Project](#).
- [i.15] ETSI TS 118 112: "oneM2M; Base Ontology".
- [i.16] [The Smart Appliances REference \(SAREF\) ontology](#).
- [i.17] ETSI TS 102 165-2: "CYBER; Methods and protocols; Part 2: Protocol Framework Definition; Security Counter Measures".
- [i.18] ETSI TR 104 048: "Securing Artificial Intelligence (SAI); Data Supply Chain Security".
- [i.19] Andrew Marshall, Jugal Parikh, Emre Kiciman and Ram Shankar Siva Kumar: ["Threat Modeling AI/ML Systems and Dependencies"](#).
- [i.20] ETSI GR SAI 001: "Securing Artificial Intelligence (SAI); AI Threat Ontology".
- [i.21] ETSI TR 104 221: "Securing Artificial Intelligence (SAI); Problem Statement".
- [i.22] ETSI TR 104 222: "Securing Artificial Intelligence; Mitigation Strategy Report".

3 Definition of terms, symbols and abbreviations

3.1 Terms

For the purposes of the present document, the terms given in ETSI TR 104 221 [i.21], ISO/IEC 22989 [1] and the following apply:

Artificial General Intelligence (AGI): applying intelligence to any intellectual task, at a level equivalent to a human

NOTE: AGI is also termed Strong AI.

Artificial Intelligence (AI): 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: From ETSI TR 104 221 [i.21].

Artificial Narrow Intelligence (ANI): applying intelligence to only one context

EXAMPLE: Autonomous driving, speech recognition.

NOTE: ANI is also termed Weak AI.

Artificial Super Intelligence (ASI): extending beyond AGI to apply intelligence to a level significantly beyond those of humans across a comprehensive range of categories and fields of endeavour

cognition: mental action or process of acquiring knowledge and understanding through thought, experience, and the senses

predicate: part of a sentence or clause containing a verb and stating something about the subject

NOTE: In the context of the present document as applied to RDF statements the predicate illustrates the nature of the relationship between two objects or concepts.

reasoning: application of learned strategies in order to solve puzzles, and make judgments where there is uncertainty in either the input or the expected outcome

3.2 Symbols

Void.

3.3 Abbreviations

For the purposes of the present document, the abbreviations given in ETSI TR 104 221 [i.21] and the following apply:

AGI	Artificial General Intelligence
AI	Artificial Intelligence
ANI	Artificial Narrow Intelligence
ASI	Artificial Super Intelligence
CAV	Connected and Autonomous Vehicles
CIA	Confidentiality Integrity Availability
CVE	Common Vulnerabilities and Exposures
CVSS	Common Vulnerability Scoring System
CWE	Common Weakness Enumeration
GAN	Generative Adversarial Networks
ICT	Information Communications Technology
IQ	Intelligence Quotient
IT	Information Technology
ITS	Intelligent Transport Systems
JSON	JavaScript Object Notation
LD	Linked Data
ML	Machine Learning
NGSI	Next Generation
NGSI-LD	Next Generation Service Interface - Linked Data
OWL	Ontology Web Language
RDF	Resource Description Framework
RL	Reinforcement Learning
SAI	Securing Artificial Intelligence
TVRA	Threat Vulnerability Risk Analysis
UML	Unified Modelling Language
XSS	Cross Site Scripting

4 From taxonomy to an ontology for secure AI

4.1 Overview

An ontology in information science identifies a set of concepts and categories within a particular field of knowledge that shows the properties of the concepts and categories and the relations between them.

This overview illustrates and demonstrates how the various concepts that are taken for granted in the security standards space are implicit as taxonomies. The overview extends to illustrate that by adopting a broader understanding of these implicit taxonomies in the form of an ontology, in which concepts are related, will help in making systems more resilient against AI attackers, or which make better use of AI in defence.

NOTE 1: The model of ontology from philosophy is the study of being, and addresses concepts such as becoming, existence and reality. For many, the ultimate aim of AI is general intelligence i.e. the ability of a single machine agent able to learn or understand any task, covering the range of human cognition. If and when AI moves closer to any concept of independent sentience, there will be increasing overlap between the worlds of information science and philosophy. However, this is likely to be decades away at least, and so the present document focusses on so-called weak AI: the use of software to perform specific, pre-defined, reasoning tasks. Also, in the philosophical domain there is a degree of crossover in the role of intelligence and the role of ethics. The present document does not attempt to define the role of ethics other than to reflect that in an ontology of intelligence that there are various schools of ethics that apply. So, an intelligence framework is influenced by its ethical framework, where the impact of the ethical framework can be realized in various ways.

In many domains that apply some form of AI, the core data model is presented in an ontological form and from that it is possible to apply more sophisticated search algorithms to allow for semantic reasoning. The technical presentation of an ontology is therefore significant of itself as it can pre-determine the way in which the programming logic is able to express intelligence. Ontologies, in the context of a semantic web, are often designed for re-use. In addition to conventional ontologies and the use of Resource Description Framework (RDF) [i.6] notations, there is growth in the use of Linked Data extensions to data passing mechanisms used widely on the internet.

EXAMPLE 1: JSON-LD [i.8] has been designed around the concept of a "context" to provide additional mappings from JSON to an RDF model. The context links object properties in a JSON document to concepts in an ontology.

EXAMPLE 2: NGSI-LD [i.9]. The term NGSI (Next Generation Service Interfaces) was first developed in work by the Open Mobile Alliance and has been extended using concepts of Linked Data to allow for wider adoption of ontologies and semantic as well as contextual information in data-driven systems.

As a pre-cursor to the development of a threat ontology for AI based threats, there are a number of threat taxonomies, some found in ETSI TS 102 165-1 [2] and in ETSI TS 102 165-2 [i.17]. These can serve as a starting point for the definition of a threat ontology, and more specifically of an AI threat ontology.

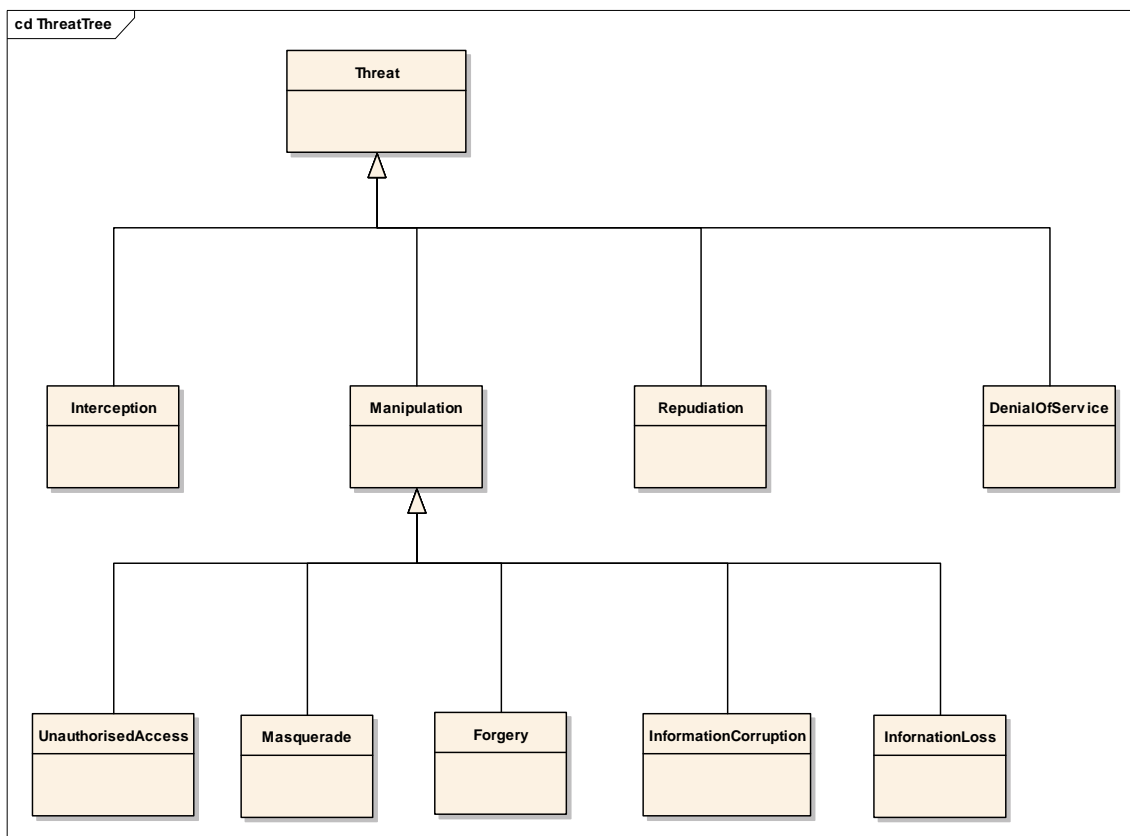


Figure 1: Threat tree (from ETSI TS 102 165-1 [2]) as a taxonomy

In the conventional taxonomy, as in Figure 1 for threats, the core relationship between entities is of type "is a", thus Forgery "is a" Manipulation, "is a" Threat. The relationships in a conventional taxonomy are often unidirectional, whereas in an ontology the normal expectation is that relationships are bidirectional and asymmetric.

EXAMPLE 3: Trust is asymmetric, a pupil is expected to trust a teacher, whereas the teacher is not expected to trust the child.

A simple taxonomy such as in Figure 1 does not easily express side channel attacks, or composite attacks, nor does it capture the asymmetric relationships of things like trust.

EXAMPLE 4: In order to perform a masquerade attack it is often necessary to first have intercepted data, or in order to corrupt data it can be necessary to first have masqueraded as an authorized entity.

Many of the forms of attack on AI that are described in the SAI Problem Statement (ETSI TR 104 221 [i.21]) are in the manipulation tree: data poisoning is a form of information corruption; incomplete data is a form of information loss. The relationship in these cases are "modifies", and "is modified by". Similarly, the terms "threat" and "vulnerability" as defined in ETSI TS 102 165-1 [2] are loosely expressed in the form of ontological relationships. Thus, threat is defined as the potential cause of an incident that may result in harm to a system or organization, where a threat consists of an asset, a threat agent and an adverse action of that threat agent on that asset, and a threat is enacted by a threat agent, and may lead to an unwanted incident breaking certain pre-defined security objectives.

NOTE 2: The nature of data poisoning is complex to clearly identify. The consequence of data poisoning are to limit the ability of the reasoning element of AI to reason towards the "right" solution, but rather to lead the reasoning algorithms to an invalid answer, often in favour of the attacker.

NOTE 3: The deliberate introduction of poisoned data may lead to the AI system exhibiting bias, where the original (non-poisoned) data may have zero biases.

NOTE 4: Whilst it is suggested that incomplete data is a form of information loss the attack vector in an AI system may be quite different than that from non-AI systems.

The structure of the term vulnerability has a similar ontological grouping of relationships, being modelled as the combination of a weakness that can be exploited by one or more threats. A more in-depth examination of the problems of and from AI is found in the SAI Problem Statement [i.21], and in the SAI report on mitigation strategies [i.22].

4.2 Formal expression of an ontology

There are many ways to express an ontology in information science. The most common are:

- OWL - Ontology Web Language [i.5]
- RDF - Resource Description Framework [i.6]

It should be noted, however, that OWL and RDF, whilst common when referring to ontologies, are not equivalent but are mutually supportive.

A simple model that underpins both OWL and RDF is the subject-predicate-object grammar structure (see Figure 2). However, there is also a more complex set of data structures that also look like the object-oriented design concepts (e.g. inheritance, overloading) underpinning design languages such as UML, and coding languages such as C++, Swift and Java. Such taxonomical classifications are also common in science, particularly in the biological sciences.

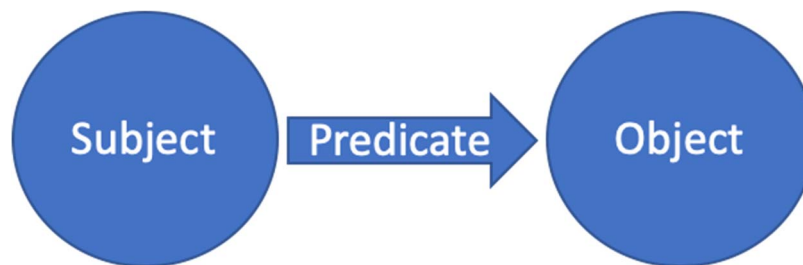


Figure 2: Simplified model of grammar underpinning Ontology

An ontology is expected to consist of the following elements:

- Classes, also known as type, sort or category.
- Attributes, which describe object instances, such as "has name", "has colour", "by definition has a".

EXAMPLE 1: A *protected object* belongs to class *network object*, of sub-type *router*, with name "Router-1" and, by definition, has 1 or more Ethernet ports.

EXAMPLE 2: *Ransomware* belongs to class *threat*, of subclass *denial-of-service*, with attribute *file-encryption*.

- Relationships, as outlined in 4.1 identifies how one class is associated to another.

Expanding from the taxonomy in [i.5], *threat* is modelled as one class, with *threat agent* modelled as another. This is then consistent with the definitions given for the terms "threat" and "vulnerability", and for the relationship to assets as the subject or object in the simplified grammar of ontology.

In the gap between an ontology and natural language, the present document classifies concepts around intelligence as nouns, and relationships as verbs, adverbs, adjectives.

NOTE: Whilst there is a risk in trying to explain AI only by mapping to programming constructs (e.g. objects and classes), or only from data modelling (e.g. tables, lists, numbers, strings and the relationships or type constraints a data model can impose) it is not addressed by the present document but is considered in ETSI TS 102 165-1 [2].

As stated above, an ontology is often described as a specification of a conceptualization of a domain. The result of such an approach to an ontology is to provide standardized definitions for the concepts of a specific domain. In the structure of a technical standard the ontology defines classes (concepts) for sets of objects in the domain that have common characteristics. The objects include specific events, actions, procedures, ideas, and so forth in addition to physical objects. In addition to the concepts, the ontology describes their characteristics or attributes, and defines typed relationships that may hold between actual objects that belong to one or more concepts.

Information in an ontology is conventionally encoded in languages such as RDF and in representations such as OWL, as a list of triplets (the "subject-predicate-object" concept), where the subject is the domain under analysis, the objects are all relevant concepts affecting the subject and the predicate defines the relationship between them.

For illustrative purposes, this can be expressed as a mathematical representation of a linear system:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_n X_n + \mu$$

$$\beta_x \neq 0$$

where Y is the domain variable to be explained by the ontology, X are the concepts as explicative variables, and coefficient β represents the relationship of the explicative variables over the variable Y , and μ is the error factor for the "unknown" concepts in Y .

In regard to standardization, ontologies, when formally modelled, provide explicit knowledge models for particular domains that can assist in both structuring the problem and in identifying where standards can assist in specifying the nature of the domain in such a way that it becomes known.

4.3 Relationship to other work

In the scope of the present document an ontology is also developed to assist in the development of strategies in securing AI. This addresses the modes in which AI can exist in a system, shown figuratively in Figure 3 below.

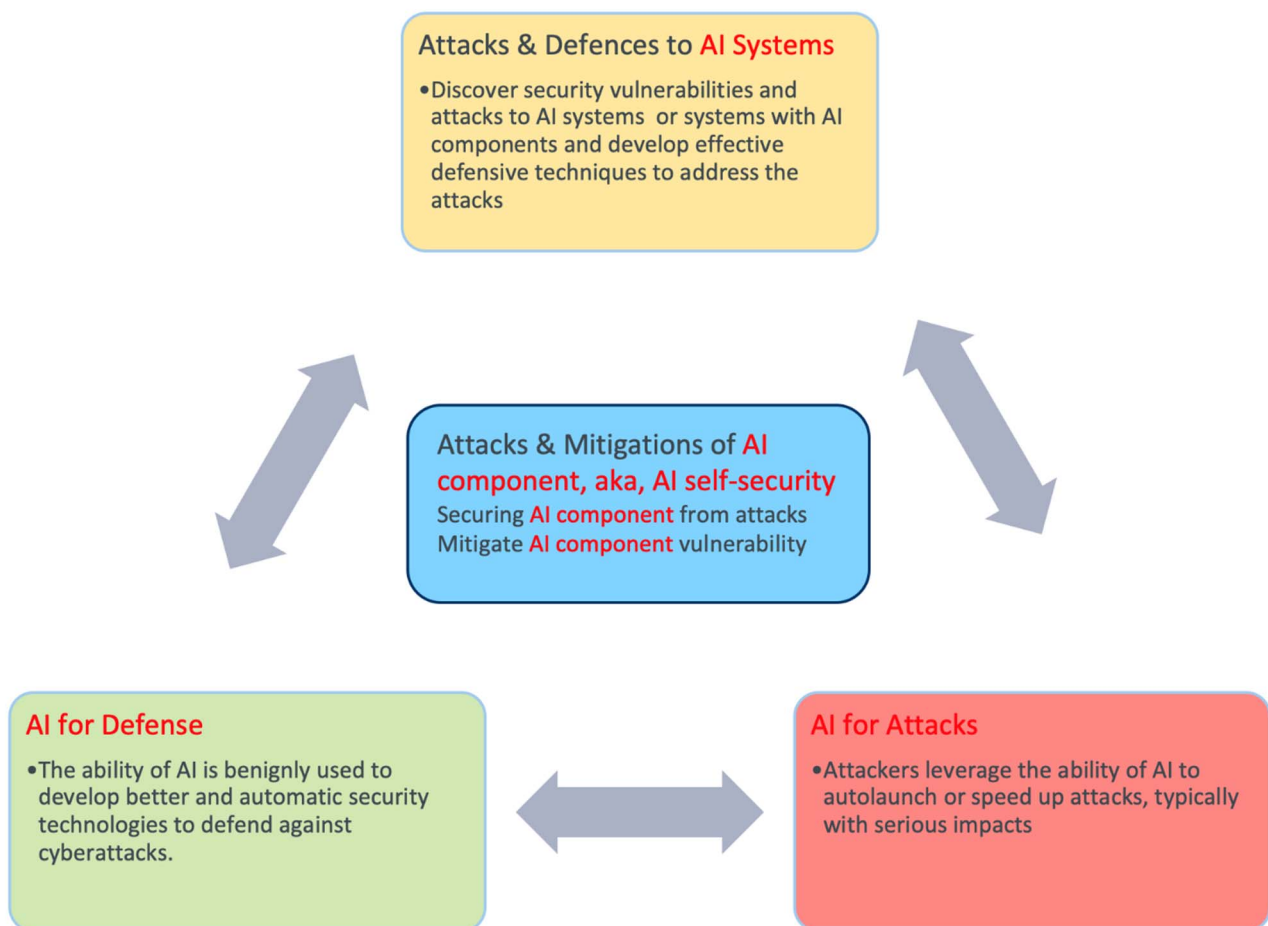


Figure 3: Modes of application of AI in networks and services

AI can be deployed in attack mode against components or systems, defence mode in countering attacks on components or systems, and proactively in understanding attacks on components and systems. Underpinning both attack and defence modes is the goal of understanding the problem associated to AI and the risk to the system of AI.

In ETSI TR 104 221 [i.21], "SAI Problem Statement", the following definition of AI is offered:

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.

In ETSI TR 104 221 [i.21], the intent is to describe the challenges of securing AI-based systems and solutions, in an environment where data, algorithms and models (in both training and implementation environments) are portrayed as challenges to the overall systems analysis and understanding in respect of AI as a system threat. In [i.21] and in ETSI TR 104 222 [i.22] there is a subtle change of emphasis with regards to data that is input to a system where the dominant AI mechanism is machine learning. In the specific subset of AI that is Machine Learning (ML), there are further modes of learning that can be defined. As noted above, the role of ontologies is implied in most ML systems as a means of structuring the input. In practice most ontologies are incomplete: they tend to be domain specific, whereas in practical systems an entity can exist in more than one domain.

EXAMPLE: A potato will naturally exist in an ontology describing tubers and variants of the nightshades, but will also naturally exist in an ontology describing foodstuffs and diet. Domain specific knowledge, knowing that a potato is related to a tomato in the family of nightshades, would not necessarily reveal the existence of poutine. Or in other words there may not be an obvious link between domains.

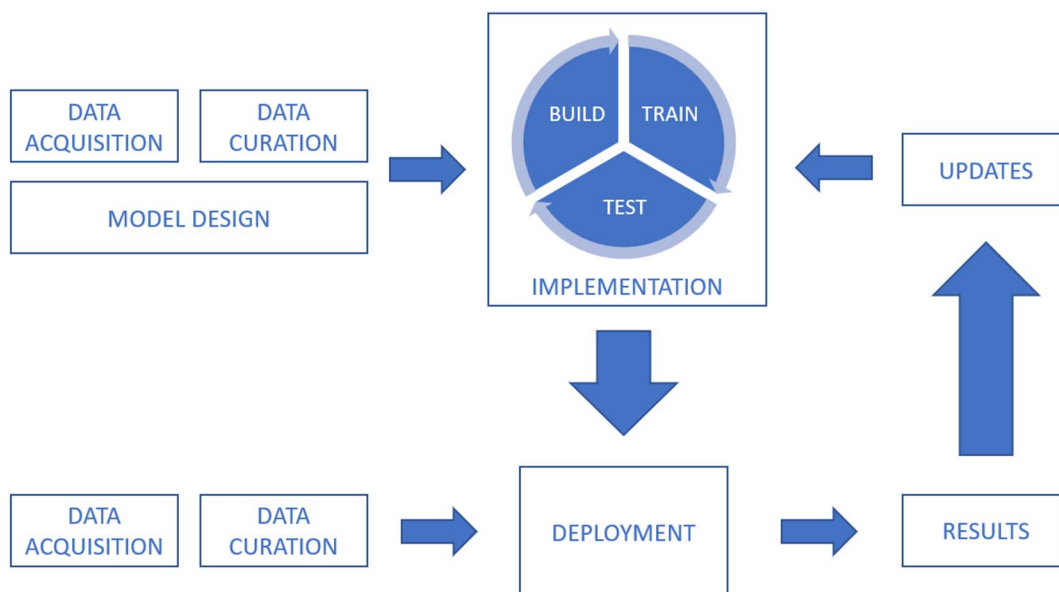


Figure 4: Typical machine learning lifecycle from ETSI TR 104 221 [i.21]

The ML lifecycle considered in ETSI TR 104 221 [i.21] identifies a number of ML strategies:

- **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.

In each case ML can be used for both classification problems, and for prediction problems.

In applying annotations or labels to data, it is common practice to first define the domain data using an ontology, in its simplest form of a semantic data set. As an ontological or semantic data definition is unlikely to be complete, one role often assigned to AI/ML is to further develop the data description, by identifying additional patterns through finding new correlations and asserting new causations.

Attack strategies against the ML workflow (illustrated in Figure 5) may apply to each stage or process flow, and the risk associated to each is assessed independently. The role of risk management is covered in more detail in clause 5 of the present document.

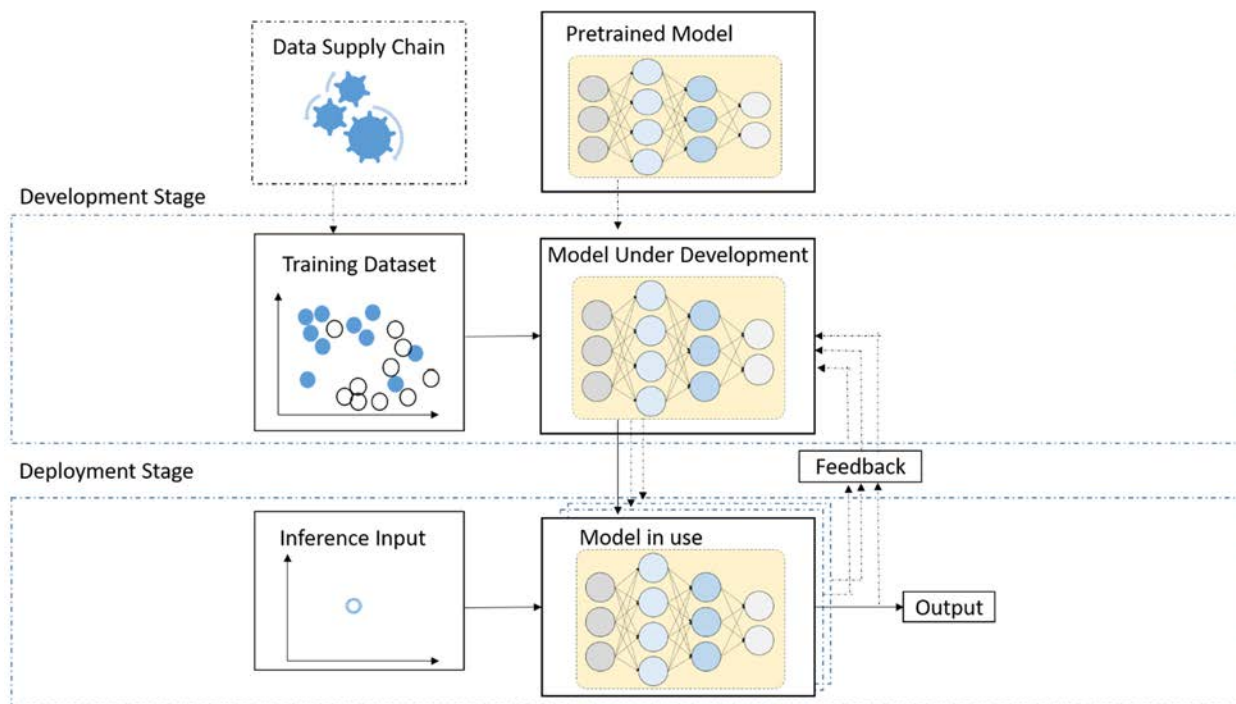


Figure 5: Machine learning model workflow from ETSI TR 104 222 [i.22]

5 Threat landscape

5.1 Threat dimensions

The TVRA model in ETSI TS 102 165-1 [2] states that "A **threat agent** enacts a specific **attack** against a system **weakness** to exploit a **vulnerability**". AI shall be considered in the context of each of the items in this statement identified with **bold text**, in both offensive and defensive contexts. The SAI problem statement in ETSI TR 104 221 [i.21] identifies ways in which a threat agent can invoke particular forms of attack, while in ETSI TR 104 222 [i.22] a number of specific mitigations are identified. In each of [i.21] and [1] the focus is on attack and defence of AI-inspired attacks on AI systems, whereas in the present document the focus is on the understanding of what AI means and of defining the AI domain itself.

Acquired intelligence, i.e. intelligence from learning, requires knowledge of data semantics (i.e. what data elements mean) and data context (i.e. how data elements are related), and conventional domain ontologies offer this form of data labelling. The richer the ontology of the input data, i.e. the more that data is labelled, the closer the ontology is to the world model required by the AI to represent the world view for the intelligence in the machine. In other words, semantic labelling is a major step forward in gaining an understanding of data.

EXAMPLE: The numerical value 42 can be syntactically represented as a signed integer which in computing terms means certain functions can be applied to it (arithmetic functions say) and a compiler will be able to warn if functions will fail based on knowledge of the syntax. However, of itself the value of the integer does not confer knowledge whereas adding a semantic label to it allows reasoning to be applied. Simple semantic labels can be seen in the names given by programmers to constants and variables, but in the wider context semantics have to be transferred with the data in order that the receiver has knowledge of what the value means, or is associated to.

5.2 Attacks as instance of threat agent

According to a 2019 report by Forrester, 86 % of cybersecurity decision makers are concerned about the offensive use of AI by threat actors [i.10]. As with many other organizations, adversaries are increasingly looking to AI to automate, scale and speed up activities which are currently conducted manually. This is particularly where malicious actors target indiscriminately, where the lower likelihood of a successful single attack is offset by the opportunity to attack many more targets. As such, although the impact of individual threats is unlikely to change, the scale of attacks, the likelihood of attacks being successfully carried out, and the difficulty in responding or remediating in a timely manner may all increase dramatically.

Opportunities range across the phases of the attack chain, from reconnaissance to exfiltration and impact. Public data, particularly social media, presents a significant opportunity for AI-based exploitation. For example, supervised and unsupervised learning can be used as tools to mine data for large-scale target discovery and spearphishing email generation, and reinforcement learning for conducting automated phishing [i.11].

AI can be used to evade defence mechanisms, including those defences based themselves on AI techniques:

EXAMPLE 1: Generative Adversarial Networks (GANs) can be used to evade ransomware detection [i.12].

EXAMPLE 2: The use of deep neural networks may be used to evade phishing detection [i.13].

AI tools can also be used to enable other kinds of attack. ML approaches are increasingly being demonstrated in side-channel analysis as an alternative to traditional statistical tools. The increasing sophistication of AI-based techniques for generating fake biometric data, and creating falsified images (so-called DeepFakes), can also pose a threat to, and undermine, trust in the integrity and authenticity of data in all aspects of a business:

EXAMPLE 3: In one case, AI-based voice-mimicking software was used to defraud an energy company of \$240 000, by convincing an employee to make a fraudulent bank transfer.

EXAMPLE 4: Deepfakes, a term applied to synthetic media in which a person in an existing image or video is replaced with someone else's likeness, has been used in a number of situations both positive and negative, including replacement of images to commit fraud or misdirection, and in more positive environments to place a dead actor in a film.

In all these examples, the use of AI is unlikely to change the impact of a successful exploitation. However, it can increase the likelihood of an organization being targeted and/or of attack attempts being successful, and hence can increase the overall risk.

5.3 Adversarial Goals

5.3.1 Violation of Confidentiality

As mentioned above, data is a crucial asset in an AI-based system. By definition, information about training data is encoded in a model itself: a model can be considered the aggregated understanding of a scenario or task derived from analysis of many examples of that scenario. Techniques exist whereby an adversary can infer aspects of this information, for example reconstructing training data examples (so-called *model inversion*) which represents a violation of confidentiality, and potentially privacy where a model has been trained on personal data. Similarly, interrogation of a model can leak information about the model itself, which can represent a leak of proprietary information. This can be particularly damaging where a business model is based around a well-trained model.

5.3.2 Violation of Integrity and Availability

As described in clause 5, the role of AI in a system is usually to make inferences about data to enable downstream decision making (human or machine), or to carry out actions based on input data. Compromise of the AI component can hence lead to a violation of integrity of the system, in that inferences and decisions will be inaccurate, and overall system performance will be degraded. If significant enough, this degradation can constitute a loss of availability, either of the component itself or of the whole system, depending on the AI's role in that system. The various modes of compromise are described in ETSI TR 104 221 [i.21].

One variant of availability compromise that particularly applies in an AI context is "reputational compromise". AI is not well understood in many domains, and one of the purposes of the present document is to offer a wider understanding of AI. Users may be wary of trusting model outputs, especially when model reasoning is difficult to explain or is sufficiently different from human reasoning. An attacker can target an AI system in an attempt to damage trust in AI itself and disrupt or damage an organization by preventing or reversing the adoption of AI technologies.

A key aspect of the integrity definition of the CIA model is the ability to reverse the effects of attacks. The relationship between data and a model, as well as the probabilistic nature of models themselves, make understanding and reversing adversarial action more challenging for AI systems than for traditional software. Once a model is trained, it is extremely difficult to undo the effect of malicious datapoints without significant retraining, which requires large amounts of reliable training data and compute. Techniques exist to harden some types of AI against adversarial input, either as a preventative or remediation measure; these are areas of active research.

5.4 Threat modelling

5.4.1 Attacker objectives

As with any cyberattack, an adversary will ultimately be aiming to extract information from a system or affect its operation in some way. The adversary can choose to do so using AI, or by attacking AI components, but ultimately the objective will be to affect a system or the information within it. Where an AI component of a system is used to provide context for decision making within a system, or make decisions itself, then compromise of the AI component will affect those decisions. The effect of compromising decisions will depend on the design and purpose of the system.

On a more granular level, attacks on AI systems can be thought of as aiming to force a model to do, learn or reveal the wrong things:

- **Do** - the actor aims to engineer an input to a model such that the output will be incorrect. The actor has control over the input but not the model itself. This class of attack is known as *evasion*. An example would be a malware author manipulating an executable binary so that an ML-based security product classifies that binary as benign software.
- **Learn** - the actor wishes to *poison* a model such that it will fail to operate as the intended, in a targeted or indiscriminate way. The actor has control over the data and/or model. The actor may be looking to degrade the overall performance of a model (functionally a denial-of-service attack), or to introduce a backdoor or trojan. In the former case, the degradation or disruption can be the actor's ultimate aim, or they wish to use the reliably poor performance of a model to achieve a downstream effect. In the backdoor case, while overall model performance will remain consistent, the actor will be able to reliably conduct an evasion attack as above.
- **Reveal** - the actor aims to uncover information about the model and/or the data used to train it. This can be for espionage or theft purposes: actors wish to steal the model itself, reveal sensitive training data or learn if particular examples have been used for training. Many of the evasion and poisoning attacks above are enabled by access to the model or an approximation of it: as such, an actor may wish to steal a model as a preparatory step in conducting another type of attack.

As already discussed, an adversary's objective in using AI for offensive purposes is likely to be similar to any organization's: increased efficiency, speed and scale; automation; and easier exploitation of large amounts of data. If attacks using AI fail, a sufficiently motivated attacker may still be able to perform the attack manually.

5.4.2 Attack surface

5.4.2.1 AI effect on impact and likelihood

As discussed in clause 6.2, the use of AI in an offensive context is not likely to increase the attack surface, or if it does, only because AI is being used to attack vulnerabilities in AI components. It can, however, increase the likelihood of attacks being prosecuted at all and/or successfully. This restates the assertion that the impact of an attack is immutable irrespective of the mode of attack, but the likelihood of an attack will change over time. AI is a vector to modify, mostly by increasing, the likelihood of an attack.

EXAMPLE: The well-known existence of deepfakes, e.g. a video, can affect trust models in that the video as a form of validation may not be as trustworthy as in the past.

The typical machine learning lifecycle is given in Figure 4. All elements of the pipeline and their dependencies should be considered in threat modelling. Some of these dependencies will be specific to AI. However, traditional dependencies, for example the code in popular ML libraries and model serving frameworks, should also be considered.

5.4.2.2 Data acquisition and curation

As previously described, compromise of training data represents a compromise of the system itself as the training data represents a core asset of the AI's performance. This is the case regardless of whether by "training data" it is meant the data used to create the initial model, or any additional data used to fine-tune or retrain a model after deployment.

Compromise of data can be achieved via manipulation of data points themselves and/or their labels (and potentially their ordering, see ETSI TR 104 048 [i.18]), or by injection of extra data samples, and can be carried out in a targeted or indiscriminate way. Dataset modification in any form can be carried out at any part of the data lifecycle: as data is collected (for sensor input can be manipulated), stored, processed or transferred. Further discussion of data supply chains, their vulnerabilities and methods to detect and mitigate attacks are given in ETSI TR 104 048 [i.18].

5.4.2.3 Implementation

The training phase is where any offline manipulation of training data becomes encoded into a model i.e. where dataset compromise becomes model compromise. In reinforcement learning, an adversary can modify the environment in which the RL agent is learning, in order to cause it to learn an incorrect or suboptimal policy.

Logic corruption attacks target the code used to generate or train a model, maliciously editing it to change the way that the algorithm learns and hence behaves. This mode of behaviour is not unique to AI: any software can be targeted to change the way a downstream system behaves, and the threat can be mitigated by following standard software supply chain security practises. However, it should be emphasized that the effect of the attack will not be felt just within the system where the code is run, rather wherever the poisoned artifacts of that code (i.e. the model) will be deployed.

Similarly, models are increasingly trained in one environment and deployed in another, either as is or with some modifications. This approach, known as transfer learning, is particularly common in applications requiring vast amounts of data and/or compute to generate, for example natural language processing. If an upstream or pre-trained model is poisoned, the subsequent effect can be transferred into different environments and is not necessarily removed by fine tuning with local data. As such, the model supply chain should be considered and assured alongside the supply chains for data and other software.

Once a model is deployed, it can be refined or retrained using new data points, labels, or other feedback from the environment, including user interactions. This data should be handled with as much care and caution as the original training data.

5.4.2.4 Deployment

Ultimately, compromise of a model cannot cause harm until that model, or one based upon it, is deployed. If an actor can trigger an incorrect response from a model (e.g. misclassification by inputting adversarial examples or by taking advantage of a backdoor), then they may be able to achieve an effect on the system downstream. In assessing threats to a system, one identifies how changing model outputs cause changes in behaviour within a system, and what effects those different behaviours ultimately have on users and the physical world.

As described in clause 5.4.1, a malicious actor can infer information about training data or the model itself by interacting with that model once deployed. Malicious user interactions can also be incorporated into a model, as described above.

5.4.2.5 Humans

The human interaction with, and understanding of, AI is a factor that is worth emphasizing in this context. As previously mentioned, AI is a growing field, and understanding of AI is likely to be less widespread than for other software paradigms. This lack of understanding can affect threats and potential outcomes.

As previously noted, an adversary can attempt to damage an AI component not because of any downstream effect per se, but to damage trust in the overall system or technology.

EXAMPLE: An attacker wishes to flood a human security analyst with false positive results, damaging their trust in AI analytics.

Conversely, another potential threat is that of overconfidence in predictions generated by any model. As previously described, a weak AI model is effective only within the bounds of the problem for which it has been trained, and its outputs are based on probabilities. For example, an AI-based antivirus only classes as malicious malware from families appearing in its training set. A human user, not realizing this limitation, can then assume that novel malware will also be caught.

5.4.3 Trust model

Trust is at the root of security in ICT systems. In symmetric relationships for example, there is trust that Alice does not share with Eve any secrets at the heart of her relationship with Bob.

Several classes of actors are relevant to a deployed ML-based system as shown in Table 1.

Table 1: Classes of actor in deployed ML systems

Actor	Role	Role in an example system (face recognition)
Data owners	Owners or trustees of the environment in which a system is deployed and the data which that system stores and use	An IT organization deploying a face recognition authentication service
System/model providers	Constructors of the system, models and/or algorithms that are used in an AI system	Authentication service software vendors
Training data providers	Provider of any labelled datasets used to (re)train the model at any stage of its development or deployment	Face dataset labelling service
Consumers	Consumers of the service the system provides	Enterprise users in the IT organization
Outsiders	Anyone capable of influencing system inputs at any stage, whether by explicit or incidental access	Insider threats, organization guests, external threat actors

It should be noted that there can be multiple examples of each class of user involved in a given deployment, and that individuals can fulfil multiple roles (a consumer can provide training data, for example). The roles of the data and model providers also should be considered explicitly, given the novel role of data and models within an AI system compared to a traditional software system. As previously noted, interactions with the system may (depending on its function and architecture) significantly impact the system itself. Further discussion of trust boundaries in AI systems is given in [i.19].

5.5 Statistics in AI and ML

In many ML systems, such as those offered in packaged programming suites and covered in ML practitioner education, the underlying approach is based on forms of statistical analysis of large data sets. In ontological terms the role of statistical analysis in ML can be addressed by defining the relationship "is enabled by" from ML to the statistics group. It is possible to dismiss all ML as effectively being statistical analysis, although that misses the intent of ML as a waypoint on the continuum from ANI to AGI and thence to ASI.

As mentioned in clause 6.3, it should be noted that, unlike traditional programming, AI systems are inherently probabilistic as opposed to deterministic. The same inputs will not result in a single, testable output; for example, two models trained on the same data may not be identical, in part this is because randomness factors are introduced to the training or learning algorithms. This makes detecting deliberate misuse more difficult. Similarly, human users may not correctly interpret the results of an ML model if they do not understand the reliance on probability.

Where statistical models are used in assisting classification modelling (e.g. pattern recognition), primary, secondary or tertiary, judgement schemes can be applied. In any judgement-based system independent observers can diverge on how they classify items. Statistical measures such as Cohen's kappa [i.7] can be used to determine the degree of agreement between two independently developed AI when classifying N items into C mutually exclusive categories. Such statistical models can be used as a tool in identifying forms of attack that lead to AI based uncertainty, however detecting misuse when the model is addressed as a non-clear-box is generally infeasible.

6 AI and SAI ontology

6.1 Nouns, verbs, adverbs and adjectives

An ontology is fundamentally a language to describe a domain. For intelligence, and for the particular domain of AI, and then for the action of securing AI, it is useful to move from natural language to the ontological expression in small steps (see Table 2 for a mapping of terms in natural language to their model form).

Table 2: Use of natural language terms

Term	Modal form	Definition
Intelligence	Mass noun	The ability to acquire and apply knowledge and skills
Security	Mass noun	The state of being free from danger or threat, i.e. a state with attributes
Learn	Verb	Gain or acquire knowledge of or skill in (something) by study, experience, or being taught
Knowledge	Mass noun	Facts, information, and skills acquired through experience or education
Reason	Verb	Think, understand, and form judgements logically
Think	Verb	Within the act of reasoning the role of thinking may require to offer scenarios in order to form connected ideas
Thought	Noun	An idea or opinion produced by thinking, or occurring suddenly in the mind
Logic	Mass noun	Reasoning conducted or assessed according to strict principles of validity

The distinction of mass noun is important for intelligence and for AI, as the definition states that a mass noun is a noun denoting something that cannot be counted (e.g. a substance or quality), in English usually a noun which lacks a plural in ordinary usage and is not used with the indefinite article (one can say "the knowledge", but not "a knowledge"). In the context of the present document intelligence is not directly quantified which is consistent with the measure of intelligence used in many human societies, where intelligence is measured relative to the wider population, often referred to as Intelligence Quotient (IQ). The IQ term introduces measurement of intelligence as a measure of reasoning. There is, from this definition, no direct measure of intelligence, rather it is measurable only by its relation to another intelligent entity.

NOTE: For the purposes of the present document, IQ is viewed as a number representing a person's reasoning ability (measured using problem-solving tests) as compared to the statistical norm or average for their age, taken as 100.

Learning is one means of achieving intelligence. Intelligence is a pre-requisite of learning. This may appear to be a circular dependency, but it is one that underpins machine learning. Learning builds the knowledge base; intelligence refines the knowledge gained. Systems apply the knowledge using rules and mores also achieved through learning. In normal human discourse, there are multiple, shared, knowledge bases. Contamination of any one base has an indirect impact on other bases. If AI systems mimic human intelligence structures, then this creates side-channel attack vectors.

6.2 Taxonomy and ontology

A taxonomy is a means of classification, and an ontology is an extension of a taxonomy by the inclusion of how things are related across classifications. Thus, in the security domain a taxonomy will often identify authentication methods within a taxonomical class, or when reporting vulnerabilities using the Common Vulnerabilities and Exposures (CVE) or Common Vulnerability Scoring System (CVSS) approaches, or when expressing weaknesses using the Common Weakness Enumeration (CWE) approach [i.14].

EXAMPLE: The CWE project maintains a long list of weaknesses and illustrates their relationships (parents, children, peers and so on). The CWE project is an extended taxonomy with several types of classification.

A more detail summary of the CWE project examining a single weakness, the "*CWE-79: Improper Neutralization of Input During Web Page Generation ('Cross-site Scripting')*" variant follows (from <http://cwe.mitre.org/data/definitions/79.html>) and is illustrated in Figure 6:

- CWE-Class → An abstract view of the weakness lying between a Pillar Weakness (more abstract) and a Base Weakness (more specific).
- CWE-Variant → Specific to a particular technology, e.g. an XSS weakness that is specific to a particular browser or server.
- CWE-Chain → An element that links weaknesses from different classes together in an OR logical arrangement (at least one element in the chain has to be present).
- CWE-Composite → A join of two or more weaknesses in an AND logical arrangement (all weaknesses have to be present).
- CWE-Base → A weakness mostly independent of technology/implementation.
- CWE-Category → A container class with a set of other entries showing a common characteristic.

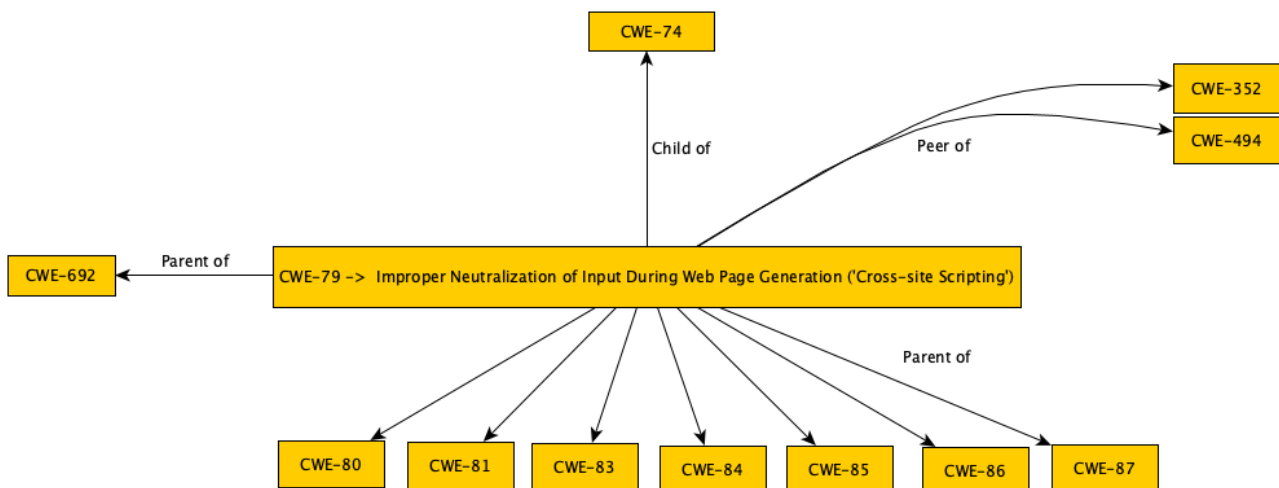


Figure 6: Illustration of CWE-79 in the form of a taxonomy or hierarchy

In the CWE model the child inherits attributes of the parent, and peers share those attributes.

6.3 Core SAI ontology relationships

Security can be understood within a wider ontology representing the state of relations between objects that symbolize, in broad terms, the Confidentiality Integrity Availability (CIA) paradigm. As such, attacks that undermine the relationships inherent in that paradigm are at the heart of understanding the role of securing AI.

The core concepts of a security ontology and taxonomy already exist in ETSI TS 102 165-1 [2] and to a lesser extent in ETSI TS 102 165-2 [i.17]. This has already been captured in Figure 8, although not expressed as an ontology. It is also recognized that ETSI TS 102 165-1 [2] and ETSI TS 102 165-2 [i.17] identify taxonomies for many forms of countermeasure.

EXAMPLE: In ETSI TS 102 165-2 [i.17] the authentication countermeasure has specializations for cryptographic authentication using a challenge response protocol to prove knowledge of a shared secret, and for digital signature approaches when an asymmetric system is used.

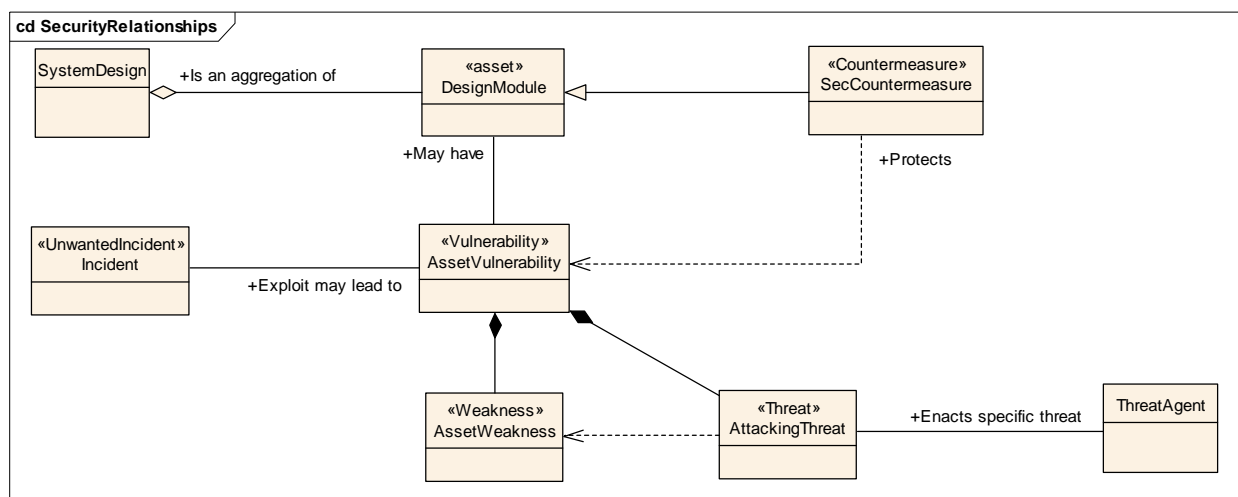


Figure 7: Model of threat from ETSI TS 102 165-1 [2]

The approach of RDF [i.6] requires statements in the form (see also Table 3):

<subject> <predicate> <object>

This is given in the example of Figure 7 as:

ThreatAgent → Enacts → Threat

or in reverse:

Threat → is enacted by → ThreatAgent.

Table 3: RDF language constructs

Construct	Syntactic form	Description
Class (a class)	C rdf:type rdfs:Class	C (a resource) is an RDF class
Property (a class)	P rdf:type rdf:Property	P (a resource) is an RDF property
type (a property)	I rdf: type C	I (a resource) is an instance of C (a class)

In addressing these RDF constructs Confidentiality, Integrity, Availability (including Authenticity) (the CIA paradigm) are properties of systems or of objects in the systems. AI is in this model also a characteristic of the attacker, i.e. a property of the attack, where an attack is an instance of a threat agent enacting a threat.

An AI agent shall be modelled as a specialisation of the Threat Agent when in adversarial mode. The classes of the AI enhanced threat model are then as shown in Figure 7 (threat, threat agent, countermeasure, and system asset are all classes):

- An AI threat agent is a specialisation of a threat agent.

In a defensive mode an AI-enabled countermeasure is modelled as a specialisation of a countermeasure:

- An AI-enabled countermeasure is a specialisation of a countermeasure.

The significant difference to the simple model given in Figure 7 is the learning and reactivity of the AI-enabled threat agent, or the AI-enabled countermeasure. Thus, the AI-enabled entity can observe the impact of the attack and modify the attack or the defence accordingly. This does not alter the core semantic relationship of threat agent to threat, or of a vulnerability being identified by a threat/threat-agent.

Observation is then added as a property of the AI-enabled threat agent and of the AI-enabled countermeasure. This expands upon the model of data learning given in clause 5:

- 1) data gathering;

- 2) data processing;
- 3) applying insights gained from the data processing.

In adversarial systems the data gathering can be applied to the rate of success of attacks and the way that attacks are countered. Data gathering can also be used to develop the ability of the threat agent by examination of the system as a whole.

NOTE: In an AI-enabled system for either adversarial or defensive application the core threat/threat-agent and countermeasure have to be adaptable. In other words, if there is only one way to apply the attack and it can be successively defended then no amount of intelligence can overcome the rigidity of the attack in that scenario.



Figure 8: AI threat agent class as a specialisation of threat agent

As shown in Figure 8 the AI Threat agent extends the core Threat Agent class by adding AI capabilities of Observe, React and Learn as methods, and its data-model as an attribute.

Annex A (informative): Cultural origins of ICT based intelligence

There is some debate regarding who first considered the role of machines as intelligent entities, but the broad consensus is that Alan Turing [i.1] is the father of machine intelligence from his paper "On computable numbers, with an application to the Entscheidungsproblem" [i.1] and its demonstration of an abstract logic that has become known as a "Turing machine", and later in his paper "Computing Machinery and Intelligence" [i.2]. The test of machine intelligence, also from Turing (the Turing Test), is one in which a machine's ability to exhibit intelligent behaviour equivalent to, or indistinguishable from, that of a human, suggests that the machine is intelligent. Turing's experiment was based on a game, the "Imitation game" in which the computer is open to lie or tell the truth, a human (the opponent), can only tell the truth, and a 3rd party asks questions to them both and based on the responses attempts to determine which is the human and which the computer. There are various rules to inhibit data leakage such that the determination can only be based on the answers to the questions. The same basic form was used in Philip K. Dick's "*Do androids dream of electric sheep?*" [i.3] in which the determination of human versus machine was based on the way in which the subject responded to questions.

NOTE 1: It is recognized that the fictional Voight-Kampff test of Philip K. Dick's vision is not directly comparable to Turing's original test but the premise that the evaluator, Deckard in the book, equivalent to party C of Turing's test, has the experience and skill to make the distinction is a failure of both tests. In practical application the machine only has to be better than the evaluator which is a lower barrier to success for the machine to be recognized as human.

The general definition of intelligence is somewhat simpler, "the ability to acquire and apply knowledge and skills", and is not comparative. However, it is also normal to consider cognition, "*the mental action or process of acquiring knowledge and understanding through thought, experience, and the senses*", as a synonym of intelligence. How much is understood regarding intelligence in humans is often difficult to gauge but the actions of neurons and their interconnection through synapses is increasingly being understood with their ability to change the weight of connection over time. Thus, whilst Turing's Universal Computing Engine is a binary system the system of neurons and synapses is more akin to an analogue system, or certainly non-binary.

NOTE 2: The connection linking neuron to neuron is termed the synapse with signals flowing in one direction, from the presynaptic neuron to the postsynaptic neuron via the synapse which acts as a variable attenuator giving weight to the signal (either chemical or electronic).

Turing's view was that a machine did not need to exist but that it could exist. It may be reasonably stated that the existence of neural processors, neural networks, and the near infinite resources of memory and processing available through cloud computing systems that such machines do now exist.

In contrast Ada Lovelace, who along with Charles Babbage, developed the concepts behind modern computer programming, stated that "*The Analytical Engine has no pretensions whatever to originate anything. It can do whatever we know how to order it to perform. It can follow analysis; but it has no power of anticipating any analytical relations or truth.*" where the Analytical Engine is one example of an early computing engine. In some respects, Ada Lovelace's concern regarding machine intelligence fits into a placement of computing as a machine able to implement combinational logic, or to act as a finite state machine, in the automata theory (see figure A.1), where such levels of computation do not anticipate the potential of the Turing Machine.

Automata theory

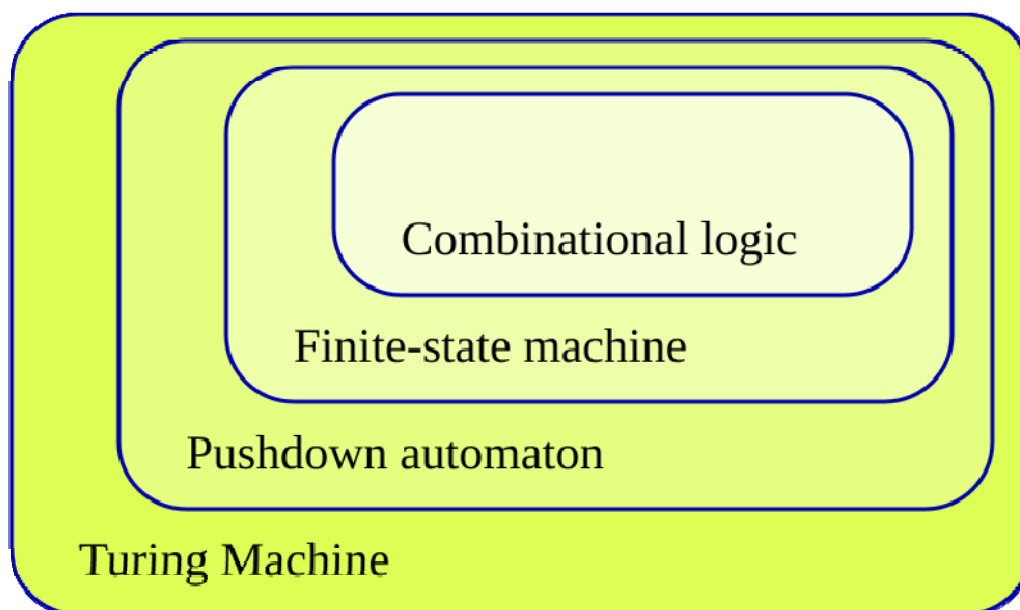


Figure A.1: Automata theory of machine logic (diagram released under Creative Commons license)

There are a number of ways of interpreting how computing engines may exhibit or gain "the ability to acquire and apply knowledge and skills". In normal human development the acquisition of knowledge or skill requires data, test and application.

In other parts of popular culture Asimov proposed three rules of robotics that may be assumed to apply to any form of near sentient AI. To avoid any issues of copying copyrighted material readers are asked to refer to the text of *i-Robot* [i.4], although the rules can be found cited in many locations.

It is possibly useful to note that Asimov proposed the rule of robotics as a moral compass, and that the rules are akin to an ethics framework (do no harm to others, do not let others come to harm, do not allow yourself to be harmed), where the ordering of the rules is critical and "others" is meant to mean humans but is generalized here.

NOTE 3: If Asimov's robotics rules are inversed such that the self-preservation rule is dominant it allows the machine to allow harm to be caused to others. In the "correct" order self-sacrifice is allowed to prevent harm to others.

In the forming of many ethical or moral rulesets this form of ordering is critical.

An initial problem with AI, and security aspects of AI, is that the domain does not appear to be well bounded, and the level of uncertainty is high. With respect to the Rumsfeld statement quoted, below the domain of AI has many unknown unknowns (*the ones we don't know we don't know*), the most pressing of which is a definitive view of intelligence.

QUOTE: "*Reports that say that something hasn't happened are always interesting to me, because as we know, there are known knowns; there are things we know we know. We also know there are known unknowns; that is to say we know there are some things we do not know. But there are also unknown unknowns - the ones we don't know we don't know. And if one looks throughout the history of our country and other free countries, it is the latter category that tend to be the difficult ones. Attributed to Donald Rumsfeld on 12-February-2002.*"

However, in addressing the AI problem in the present document, whilst there is a degree of uncertainty regarding Artificial General Intelligence (AGI), it has, for the purposes of the present document, a low likelihood of actualisation and therefore the focus of the present document is on Artificial Narrow Intelligence (ANI). See also the discussion in clause 4.5.

Intelligence and intellect are two additional terms that are often confused. An ANI will not be considered as intellectual although a AGI will be, where the definition of intellect is given as having the faculty of reasoning and understanding objectively, especially with regard to abstract matters, and a machine having this faculty can be described as intellectual. An ANI will in many cases be designed in such a way that it cannot be intellectual - rather it is designed to be good at a single function. In contrast a AGI will be able to apply learning or knowledge from one field to another and to abstract knowledge to multiple fields.

Annex B (informative): Machine processing to simulate intelligence

B.1 Overview of the machine intelligence continuum

NOTE: This annex is a guide and readers are also advised to consult the sources in the Bibliography and references of the present document to gather a full picture of the development and application of machine intelligence.

The use of data to build up intelligence can be represented as a continuum in terms of how data is interpreted.

At the far-left of the continuum are expert systems: Guided paths through data. Early forays in AI were often in the form of expert systems, and can be typified as a tree of questions of an "if ... then ... else" structure, where the enquirer was led to an answer based on their responses to specific questions. The future stages in the continuum are presented in figure B.1.

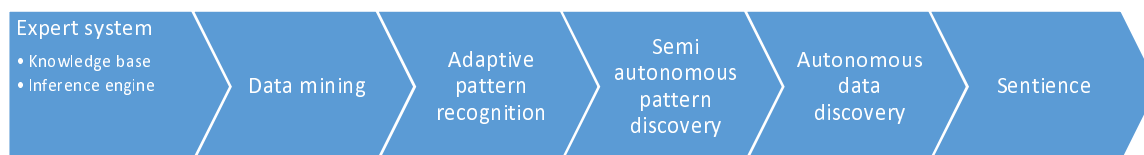


Figure B.1: Continuum of intelligence, one possible model

In figure B.1 it is noted that sentience is not a simple achievement that would be a trivial next step from autonomous data discovery, rather it has to be recognized that there are substantial intellectual and engineering barriers to overcome between each stage, and that all of these will co-exist for all time.

B.2 Expert systems

The founder of much of the early work on expert systems, Edward Feigenbaum, postulated "*intelligent systems derive their power from the knowledge they possess rather than from the specific formalisms and inference schemes they use*".

Expert systems are not autonomously intelligent, rather their observed intelligence is captured in the design of the inference systems. A crude simplification of an expert system is one in which a question is offered with a finite set of possible answers, such as in fault diagnosis, and based on the answer to the first question a second and then a third and further layers of questions and answers lead to a diagnosis. The "intelligence" is in the construct of the questions and answers and there is no processed intelligence.

B.3 Data mining and pattern extraction

The statistical analysis of large volumes of data to match patterns, or to discover patterns is the forefather of much of ML (see also clause 5.5 of the present document). The application of significant levels of processing power and memory to such statistical analysis led to the ability to sift through data quickly, in addition the introduction of alternative database models allowed changes in data analysis to be developed. Conventional database wisdom of relational databases with highly structured search patterns was challenged by offering pooled or unstructured data to be searched using statistical approaches. Adding self-described data into the mix, using semantic labelling and similar advances, allowed data resources to be searched or analysed without the structural rigidity of conventional relational databases.

Annex C (informative): Bibliography

C.1 AGI analysis

- [A] L Chen, Z Yi, X Chen: "[Research on Network Security Technology Based on Artificial Intelligence](#)", - Recent Trends in Intelligent Computing, Communication and Devices, 2020 - Springer.
- [B] W Wu, T Huang, K Gong: "[Ethical Principles and Governance Technology Development of AI in China](#)"; Engineering, 2020 - Elsevier.
- [C] [R Girasa](#): "[International Initiatives in AI](#)", Artificial Intelligence as a Disruptive Technology, 2020 - Springer.
- [D] FB Pizzini, [F Pesapane](#), [W Niessen](#): "[ESMRMB Round Table report on "Can Europe Lead in Machine Learning of MRI-Data?"](#)", 2020 - Springer.
- [E] [A Abuarqoub](#): "[D-FAP: Dual-Factor Authentication Protocol for Mobile Cloud Connected Devices](#)", Journal of Sensor and Actuator Networks, 2020 - mdpi.com.
- [F] W Hoffmann-Riem: "[Artificial Intelligence as a Challenge for Law and Regulation](#)", Regulating Artificial Intelligence, 2020 - Springer.
- [G] D Feldner: "[Designing a Future Europe](#)", Redesigning Organizations, 2020 - Springer.
- [H] SM Lee, [SC Jeong](#): "[A study on strategy for invigorating utilization of HPC in industry based on business building blocks model](#)", Nonlinear Theory and Its Applications, IEICE, 2020 - jstage.jst.go.jp.
- [I] [R Girasa](#): "[Bias, Jobs, and Fake News](#)", Artificial Intelligence as a Disruptive Technology, 2020 - Springer.
- [J] Jakob Schemmel: "[Artificial Intelligence and the Financial Markets: Business as Usual?](#)", Regulating Artificial Intelligence, 2020 - Springer.
- [K] RK Saini, C Prakash, [A Dua](#): "[A Survey on Artificial Intelligence Techniques for Cybersecurity](#)", Current Trends in Information Technology, 2020 - computerjournals.stmjournals.in.
- [L] RT Kreutzer, M Sirrenberg: "Fields of Application of Artificial Intelligence - Security Sector and Military Sector", Understanding Artificial Intelligence, 2020 - Springer.
- [M] S Saif, S Biswas: "[On the Implementation and Performance Evaluation of Security Algorithms for Healthcare](#)", Proceedings of the 2nd International Conference on Communication, Devices and Computing, 2020 - Springer.
- [N] X Liu, Y Lin, [H Li](#), J Zhang: "[A novel method for malware detection on ML-based visualization technique](#)", Computers & Security, 2020 - Elsevier.
- [O] T Chen, [J Liu](#), Y Xiang, W Niu, E Tong, Z Han: "[Adversarial attack and defense in reinforcement learning-from AI security view](#)", Cybersecurity, 2019 - Springer.
- [P] Russell, Stuart J.; Norvig, Peter (2003): Artificial Intelligence: A Modern Approach (2nd ed.), Upper Saddle River, New Jersey: Prentice Hall, ISBN 0-13-790395-2.
- [Q] Luger, George; Stubblefield, William (2004): Artificial Intelligence: Structures and Strategies for Complex Problem Solving (5th ed.), The Benjamin/Cummings Publishing Company, Inc., p. 720, ISBN 978-0-8053-4780-7.

C.2 AI in the context of threat analysis

- [R] ["Threat Modeling AI/ML Systems and Dependencies"](#), Andrew Marshall, Jugal Parikh, Emre Kiciman and Ram Shankar Siva Kumar; November 2019.
- [S] N. Papernot, P. McDaniel, A. Sinha and M. P. Wellman, "[SoK: Security and Privacy in Machine Learning](#)", 2018 IEEE™ European Symposium on Security and Privacy (EuroS&P), 2018, pp. 399-414, doi: 10.1109/EuroSP.2018.00035.

C.3 Societal and cultural references to AI

- [T] Alan Turing: "On computable numbers, with an application to the Entscheidungsproblem".
- [U] Alan Turing: "Computing Machinery and Intelligence".
- [V] Philip K. Dick: "Do androids dream of electric sheep?" (ISBN-13: 978-0575094185).
- [W] Isaac Asimov: "I, robot" (ISBN-13: 978-0008279554).

History

Document history		
V1.1.1	January 2022	Publication as ETSI GR SAI 001
V1.1.1	March 2025	Publication