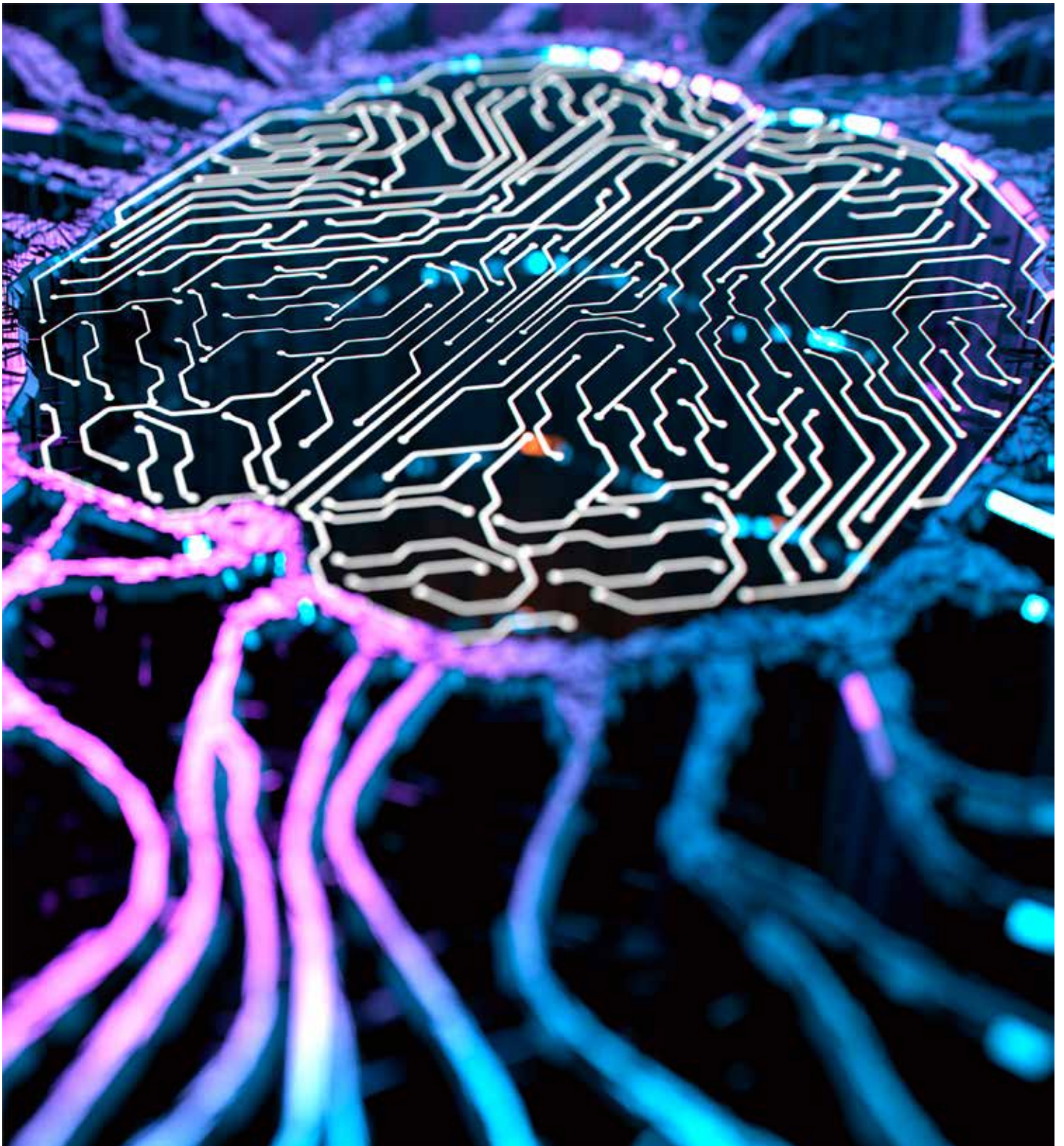


# Taiala mo le faalauteleina o le secure AI system





Communications Security Establishment  
Canadian Centre for Cyber Security

Centre de la sécurité des télécommunications  
Centre canadien pour la cybersécurité



National Cyber and Information Security Agency



REPUBLIC OF ESTONIA  
INFORMATION SYSTEM AUTHORITY



RÉPUBLIQUE FRANÇAISE  
Liberté  
Égalité  
Fraternité



Federal Office for Information Security



INCD Israel National Cyber Directorate



NISC 内閣サイバーセキュリティセンター  
National center of Incident readiness and Strategy for Cybersecurity

National Cyber Security Centre

Ni TDA



NSM  
NORWEGIAN NATIONAL  
CYBER SECURITY CENTRE



NASK



Ministerstwo  
Cyfryzacji

CSA  
SINGAPORE  
Cyber Security Agency of Singapore



## Faatatau i lenei pepa

O lenei pepa ua lomia e le UK National Cyber Security Centre (NCSC), le US Cybersecurity and Infrastructure Security Agency (CISA), ma paaga faavaomalo o loo taua i lalo:

- National Security Agency (NSA) (Ofisa Aoao tau le Puiipuiga)
- Federal Bureau of Investigations (FBI)
- Australian Signals Directorate's Australian Cyber Security Centre (ACSC)
- Canadian Centre for Cyber Security (CCCS)
- New Zealand National Cyber Security Centre (NCSC-NZ)
- Chile's Government CSIRT
- Czechia's National Cyber and Information Security Agency (NUKIB)
- Information System Authority of Estonia (RIA) and National Cyber Security Centre of Estonia (NCSC-EE)
- French Cybersecurity Agency (ANSSI)
- Germany's Federal Office for Information Security (BSI)
- Israeli National Cyber Directorate (INCD)
- Italian National Cybersecurity Agency (ACN)
- Japan's National center of Incident readiness and Strategy for Cybersecurity (NISC)
- Japan's Secretariat of Science, Technology and Innovation Policy, Cabinet Office
- Nigeria's National Information Technology Development Agency (NITDA)
- Norwegian National Cyber Security Centre (NCSC-NO)
- Poland Ministry of Digital Affairs
- Poland's NASK National Research Institute (NASK)
- Republic of Korea National Intelligence Service (NIS)
- Cyber Security Agency of Singapore (CSA)

## Faalauiloaina

O faalapopotoga o loo taua i lalo na fai o latou sao i le faalauteleina o nei taiala:

- Alan Turing Institute
- Anthropic
- Databricks
- Georgetown University's Center for Security and Emerging Technology
- Google
- Google DeepMind
- IBM
- ImBue
- Microsoft
- OpenAI
- Palantir
- RAND
- Scale AI
- Software Engineering Institute at Carnegie Mellon University
- Stanford Center for AI Safety
- Stanford Program on Geopolitics, Technology and Governance

## Tautinoga Patino

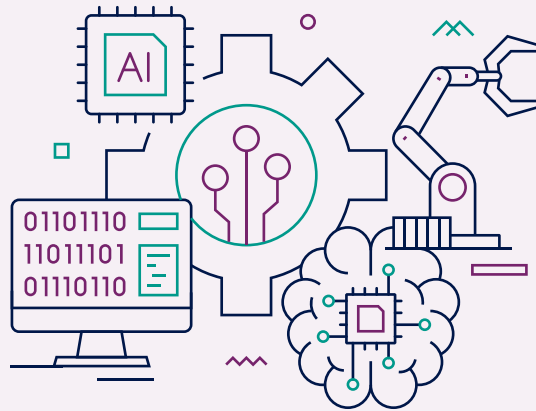
O le faamatalaga i lenei pepa ua saunia ile "tulaga o loo iai" e le NCSC ma ofisa ua tuufaatasia lenei tusitusiga, e le tataua ona agavaa i latou mo soo se mea e leiloa, se manuaga poo soo se ituaiga mea e faaleagaina e mafua mai i le faaogaina e pei ona manaomia i lalo o le tulafono. O le faamatalaga i lenei pepa e le tulai pe faamauina ai se ioega poo se fautuaga a se isi o isi faalapopotoga o loo aofia ai, oloa poo se auunaga a le NSC ma ofisa o loo tuufaatasia lenei tusitusiga. Sootaga ma faamatalaga e faatatau i upega tafailagi ma isi meafaitino o loo aofia ai ma saunia mo nao ni faamatalaga ae le avea o se ioega poo se fautuaina o ia meafaitino e sili atu i lo isi.

O lenei pepa o loo avanoa atu i se tulaga TLP:CLEAR (<https://www.first.org/tlp/>).



# Mea o loo aofia

Aotelega mai Pulega .....	5
Folasaga .....	6
Aisea e ese ai le AI security .....	6
O ai e tatau ona faitau i lenei pepa .....	7
O ai e feagai ma le faalauteleina o le secure AI .....	7
Taiala mo le faalauteleina o le system mauu a le AI?.....	8
1. Secure design .....	9
2. Secure development.....	12
3. Secure deployment.....	14
4. Secure operation and maintenance .....	16
Nisi faitauga.....	17



# Aotelega mai Pulega

O lenei pepa ua fautuaina ai taiala mo providers (i latou e saunia auunaga)- o soo se ituaiga system latou te faaogaina artificial intelligence (poto tau tekonoosi e tai tutusa ma le tomai o le tagata) tusa poo na systems ua fausia mai ni vaegamea pe fausia foi i luga o meafaigaluega ma auunaga e saunia e isi. O le faatinoina o nei taiala o le a fesoasoani i providers e fausia AI systems e gaii e pei ona fuafuaina, ma o loo avanoa pe a manaomia, ma e galue foi e aunoa ma le faailoaina o faamatalaga maaleale i pātī e le faatagaina.

O lenei pepa ua taula'i tonu i providers o AI systems o loo latou faaogaina models poo faataitaiga faatino e se faalapotopotoga, pe o faaoga foi se isi polokalame (external application programming interface (APIs). Matou te unia stakeholders **uma** (e aofia ai ma faamatalaga faasaienitisi, le 'aufaalautele, pulega, i latou e faia faaiuga ma umia tulaga lamatia) ina ia o latou faitau i nei taiala e fesoasoani ai ia i latou e faia faaiuga lelei e faatatau i le **faatulagana, faalateleina, faasoaina** ma le **faagaioia** o a latou AI systems.

## Faatatau i taiala

O AI systems e iai lo latou agavaa e aumai ai le tele o penefiti i le sosaiete. Peitai, ina ia matuai manino avanoa o le AI, e tatau ona faalateleina, faasoia ma faagaioia i se auuala e saogalemu ma lelei.

O AI systems e ono iai ni ona tulaga maaleale tau le puipuiga e manaomia ai ona iloilo lelei faatasi ma tulaga masani o faamata'u o cyber security. Afai e maualuga le faatinoina o le faalateleiga - e pei foi ona iai le tulaga tau le puipuiga o le AI- e masani ona silasila ile puipuiga o se iloilo lona lua. O le puipuiga lava e tatau ona avea ma manaomiaga autu, ae le nao le vaega o le faalateleiga ae ia faaoga ile faataamiloga atoa o le system.

Ona o lenei mafuaaga, ua vaevaeina ai taiala i ni vaega autu se fa i totonu o le faalateleina o le faataamiloga o le AI system: **faatulagana mautu, faalateleiga mautu, faasoaina mautu**, ma le **faatinoga ma le faatumauiina mautu**. Mo vaega taitasi, matou te fautuaina le faia o iloilo ma ni tulaga e fesoasoani i le faaitiitia o lamatiaga atoa i le faalateleina o se system a le AI.

### 1. Secure design (Faatulagana mautu)

O le vaega lenei e aofia ai taiala e apalai i le vaega o le faatulagaina o le faalateleina ole faataamiloga o le olaga o le AI system. E aofia ai le malamalama i lamatiaga ma faataitaiga o faamata'u, aemaise o nisi o ulutala patino ma ni feaatauaiga e ao ona iloilo ile faatulagana o le system ma faataitaiga (model).

### 2. Secure development (faalateleiga mautu)

O le vaega lenei e aofia ai taiala e apalai ile vaega o le faalateleiga o le olaga o le AI system, e aofia ai ma tulaga tau le puipuiga i le supply chain, faapepaina ma le puleaina o aseta ma aitalafu i tulaga tau masini.

### 3. Secure deployment (faasoaga mautu)

O lenei vaega e aofia ai taiala e apalai i le vaega o le faasoaga o le faalateleina o le olaga o le AI system, e aofia ai ma le puipua o meatotino i le siomaga ma models poo faataitaiga mai tulaga lamatia, faamata'u pe o le lusi foi o faiga tau le puleaina o tulaga ia atoa ma lona faatinoga talafeagai.

### 4. Secure operation and maintenance (Faagaioiga ma le faatumauiina mautu)

O le vaega lenei e aofia ai taiala e apalai i le faatumauiina o gaioiga ma le vaega e faatumauiina ai le faalateleiga o le faataamiloga o le olaga o le AI system. Na te saunia taiala e faatatau i gaioiga e faia fua agai i le faasoaina o se system, e aofia ai le tatalaina ma le mataituina, faafouina o pulega ma le faasoaina o faamatalaga.

O taiala e taumulimuli i le faiga a le 'secure by default', ma o loo faatulaga faalatalata i faiga o loo faamanino mai i le NCSC's [Secure development and deployment guidance](#), NIST's [Secure Software Development Framework](#), ma le '[secure by design principles](#)' lomina e le CISA, le NCSC ma faalapotopotoga faavaomalo o puipuiga tau initaneti. Latou te faamuamua:

- umia o iuga o tulaga tau le puipuiga mo tagata faatau
- opogi le faamaoni o le faamalalamalama ma mea ma le tali manino atu
- fausia o faatulagana o faalapotopotoga ma taitaiga e ave ai le faamuamua a le pisinisi i faiga a le secure by design



# Folasaga

O systems a le Artificial Intelligence (AI) e iai le agavaa e mafai ona o latou aumaia ni mau penefiti i le sosaiete. Peitai, ina ia matua'i iloa lelei avanoa ia o le AI, e tatau ona faalateleina, faasoa ma faatinoina i se auala e mautu ma lelei. O le Cyber security (puipuiga mai osofaiga i luga o le initaneti) o se tulaga e tatau ona muai faataatia mo le saogalemu, lavasa'ia, tulaga faalilolilo, sa'o, faatuatuaina ma le faamoemoeina o systems o le AI.

Peitai, o systems a le AI e fua lava i o latou tulaga fou maaleale (novel security vulnerabilities) e tatau ona silasila i ai faatasi ma ni faamata'u tau cyber security. Afai e maua tulaga le saosaoa o le faalateleina — e pei ona iai ile tulaga o le AI — e ono masani ona avea tulaga tau le puipuiga o se mataupu e le o autu i ai se silasila. O le tulaga lava tau le puipuiga e tatau ona avea o se manaomiaga autu, e le na'o le faa aofia i le vaega o le faalateleina, ae ia faatino lava i le faataamilosaga o le olaga atoa o le system.

**O lenei pepa e fautuaina taiala mo providers' o soo se system e faaoga ai le AI, pe o ituaiga systems ua fausia mai ni vaegamea pe fausia mai luga o ni meafaigaluega ma auaunaga e saunia mai e isi. O le faataunuina o nei taiala o le a fesoasoani ai i providers e fausia AI systems e gaiioi e pei ona fuafuaina, o loo avanoa ane pe a manaomia, ma e gaiioi e aunoa ma le faalauiloina o ni faamatalaga maaleale i ni pātī e le faatagaina.**

O nei taiala e tatau ona silasila i ai faatasi ma nisi faiga ua fautuina mo le puipuiga, faatonutonuina o lamatiaga ma faataitaiga aupito sili. I se tautalaga patino, matou te unaia ai providers ina ia mulimulita'i i faiga a le 'secure by design' o loo faalatele e le US Cybersecurity and Infrastructure Security Agency (CISA), le UK National Cyber Security Centre (NCSC), ma a tatou paaga uma faavaomalo. O manatu autu e faamuamua ai:

- le umia o iuga o tulaga tau le puipuiga mo tagata faatau
- opogi le faamaoni o le faamalamalamaina o mea ma le tali manino atu
- fausia o faatulagana o faalapopotoga ma taitaiga e ave ai le faamuamua a le pisinisi i faiga a le secure by design

O le mulimulita'i i faiga a le 'secure by design' e manaomia ai ni meafaitino taua ile faataamilosaga o le olaga o le system. O lona uiga o developers (tagata latou te faalateleina) e tatau ona inivesi i le ave o le faamuamua **features (mea o loo aofia ai), mechanisms (auala e faatino ai),** ma le **implementation (faatinoina)** o meafaigaluega e puipui ai tagata faatau i vaega taitasi o le faatulagana o le system, aemaise i vaega uma o le faalateleina o le olaga o le oloa. O le faia o lea tulaga o le a taofia ai le toe faia o nisi faatulagana taugata e mulimuli ane, aemaise o le faasaogalemuina o tagata faatau ma o latou faamaumauga i se taimi lata mai.

## Aisea e ese ai le AI security (tulaga tau le puipuiga a le AI)?

I lenei pepa o loo matou faaoga le 'AI' e faauiga patino lava i le a'oa'oina o talosaga tau masini (ML)<sup>3</sup>. O ituaiga uma o ML o loo suesueina. Matou te faamatalaina talosaga a le ML o talosaga e:

- aofia ai mea tau software components (models) e mafai ai e komipiuta ona mataitu ma aumai i totonu ni faasologa o ni faamaumauga e aunoa ma le faapolokalameina auilili o tulafono e se tagata
- ia maua ni matematega, fautuaga poo ni faaiuga e faia e fua i luga o finauga tau faamatalaga

E ese mai le iai o ni taufaamata'u i le tulaga tau le puipuiga, o AI systems e ono iai foi nisi ituaiga maaleale fou e ono aafia ai. O le faaupuga 'adversarial cyber machine learning' (AML) ua faaogaina e faamatala ai le faaogaina o tulaga maaleale i mea o loo aofia i le ML, e aofia ai hardware (meafaitino), software (polokalame tau le komipiuta), sologa lelei o galuega ma le faasologa o le sapalai o le oloa. O le AML e mafai ona ia faatagaina tagata osofai e faatupu nisi ituaiga amioga e lei fuafuaina i le ML systems ma e ono aafia ai:

- aafiaga i le faavasegaina o le model poo lana gaiioiga
- faataga ai users e faatino nisi gaiioiga lē faatagaina
- ina ia maua mai faamatalaga maaleale e faatatau i le model poo le faataitaiga

E tele auala e mafai ona ausia ai nei aafiaga, e pei o le faatino o le prompt injection attacks poo osofaiga faatotope i le (LLM) domain, pe fai ma le mautinoa le faaleagaina o faamatalaga mo toleniga, poo user feedback (e faaigoa foi o le 'data poisoning').



## O ai e tatau ona faitau i lenei pepa?

O lenei pepa o loo faatatau sa'o lava i providers o le AI systems, pe fua a latou auaunaga i models poo faataitaiga ua talimalo ai se faalapotopotoga, pe faaoga foi o nisi ituaiga polokalame, o external application programming interfaces (APIs). Peitai, matou te unaia stakeholders **uma** (e aofia ai saienitisi tau faamaumauga, developers, pule, le au faifaaiuga ma i latou e feagai ma tulaga maaleale) ina ia faitau i taiala nei e fesoasoani ai ia i latou ia faia faaiuga maiio e faatatau i le **faatulagaina**, **faasoaina** ma le **faagaioia** o a latou systems o machine learning o le AI.

Ina ua faailoa atu lena tulaga, e le o taiala uma e talafeagai tonu ma faalapotopotoga uma. O le ituaiga maualuga o le tulaga e iai ma metotia e osofaia ai e fesuisuia'i e fua lava i le fili o loo taulai mai lana osofaiga i le AI system, o lona uiga la o taiala e ao ona silasila i ai faatasi ma isi mataupu tau i lau faalapotopotoga ma ni faamata'u o iai.

## O ai e feagai ma le faalauteleina o le secure AI?

E masani ona toatele actors i le vaega o le sapalaia o modern AI. O se faiga faigofie e mafai ona faaoga ai ni pisinisi se lua:

- o le 'provider' lea e feagai ma le faia o faamatalaga, faalauteleina o faiga tau masini komipiuta, faatulagana, faasoaina ma le faatumauina
- o le 'user', lea na te saunia ni mea e tuu i totonu (input) ma ia mauina foi mea sau i fafo (output)

Ao faaoga ai lea ituaiga faiga o le provider-user approach i le tele o talosaga, ua faasolo atu ina faateteleina le lē taatele<sup>4</sup>, aemaise e ono silasila providers e tuu faatasi faamatalaga tau komipiuta, faamaumauga, faataitaiga aemaise o le faaoga o auaunaga e saunia mai e isi pātī agai i a latou systems. O nei supply chains faigata e faafaigata ai foi i le end user ona malamalama i le vaega o loo taatia ai le matafaioi a le AI.

O users (tusa poo 'end users', pe o se provider e latou te tuu faatasi ni mea e faatatau i AI components,<sup>5</sup>) e lē lava so latou silafia a'ia'i pe o se tomai foi e malamalama atoa ai, iloilo pe tagofia foi aafiaga e iai sootaga ma systems o loo latou faaogaina. O lea la, ina ia o gatasi ma manatu faavae a le 'secure by design', **e tatau i providers o mea o loo aofia ile AI ona faatino le latou matafaioi mo se iuga tau le puipuiga mo users ile faasolo atu agai i le supply chain (gaosiga ma le tufatufaina o oloa).**

E tatau i providers ona faatino ni tulaga tau le puipuiga ma mea e faaititia ai aafiaga i soo se vaega e mafai ai i a latou models, alapaipa ma/poo systems, ma mea o loo faaoga ai faatulagana, faatino le mea e aupito saogalemu e avea ma default. I le vaega e lē mafai ona faia i ai ni tulaga faaititia, e tatau ona feagai lea ma le provider:

- faailoa i users ile agai i le supply chain, lamatiaga e ono feagai ma latou (ma afai e apalai) taliaina foi e a latou lava users
- fautuaina i latou pe faapefea ona faaoga mea o loo aofia ai i se tulaga mautu

I vaega e ono lamatia ai, e ono o'o atu ai i ni tulaga e mafai ona tagofia pe faaleagaina atoa ai le talaaga, o se gau tele i faiga faapisinisi, liki i tua o faamaumauga maaleale ma faalilolilo aemaise ni aafiaga faaletulafono, lamatiaga tau AI cyber security e tatau ona faia faapei o se tulaga **tugā**.







# 1. Secure design (faatulagana mautu)

O le vaega lenei e aofia ai taiala e apalai ile **design** vaega o le faalauteleina o le olaga o le AI system. E aofia ai ma le malamalama i lamatiaga ma faataitaiga o faamata'u, aemaise o ulutala patino ma fefaatauaiga e ao ona iloilo i le faatulagana o le system ma le model.

## Siitia le silafia o le aufaigaluega i faamata'u ma lamatiaga



O i latou e umia systems ma taitai sinia latou te malamalama i faamata'u agai i le secure AI ma ni auala e tau tuuaititia ai. O au saienitisi tau faamatalaga ma developers latou te faamautuina ni faamata'u tau le puipuiga e ono aafia ai ma le faaletonu o nisi o modes ma fesoasoani i tagata o loo iai i tulaga lamatia ia faia ni faaiuga maioio. E te saunia faamatalaga i users ma taiala e faatatau i lamatiaga tau le puipuiga o loo feagai ma le AI system. (faataitaiga, o se vaega o le standard InfoSec training) ma toleni ai developers ina ia faamautu auala mo le makaina o codes ma metotia faanonaponei ma mautu ai foi ma lelei faiga ia a le AI.

## Faatitai faamata'u i lau system



O se vaega o lau fuafuaga mo le faatonutonuina o lamatiaga, e te apalai ai se fuafuaga e iloilo ai faamata'u i lau system, lea e aofia ai lou malamalama i ni aafiaga e ono iai ile system, users, faalapopotoga ma le sosaiete aao pe afai ua lamatia se meatotino o le AI pe ua faafuasei ona suia uiga<sup>7</sup>. O lenei fuafuaga e aofia ai le iloiloina o aafiaga o le AI-specific threats<sup>8</sup> ma le faamaumauina o au faaiuga fai.

E te maitauina o ituaiga faamatalaga maaleale ua faaoga i lau system e ono aafia ai lona taua o se taulaiga mo le au osofai. O lau iloiloga e tatau ona silasila ai i nisi o faamata'u e ono faatupulaia a'o faaauau pea ona silasila tagata i systems a le AI oni taulaiga taua, ma pe afai e faataga e le AI nisi osofaiga otomeki mai vectors.

## Faatulaga lau system mo le puipuiga aemaise le faagaioiina ma lana gaioi.



O loo e lototele o le galuega ua i ou lima ua matuai talafeagai ona tagofia faaoga ai le AI. Ina ua e faamautinoaina lenei, ona e iloilo lea o le sa'o atoatoa o lau filifiliga o lau AI design. E te manatu o lau threat model e iai sootaga ma faiga tau le puipuiga aemaise lona faagaioiina, o le silafia o le user, siosiomaga o loo faasoaina ai, gaioi, faamautinoaina, silasila mamao, mea manaomia faaletulafono ma le tagata aemaise isi vaega ua iloiloina. Mo se faataitaiga:

- e te manatu ile puipuiga o le supply chain pe a fai lau filifiliga pe faalautele i totonu o lou lava siomaga pe faaoga external components, mo se faataitaiga:
  - o lau faaiuga e toleni se model fou, faaoga se model o loo iai (pe iai pe leai le fine tuning) pe tagofia foi se model e ui atu i le external API e talafeagai i mea o loo e manaomia
  - lau filifiliga e galue ma se external model provider e aofia ai se iloiloga ua faataunuuina i le tulaga o lona provider ma lona puipuiga
  - afai o faaoga se faletusi i fafo atu (external library), e te faataunuuina se iloiloga (mo se faataitaiga, e faamautinoaina ai o loo iai mea faatonutonu a le faletusi e taofia ai le system mai le la'uina i luga o models e le faatuatuaina e aunoa ma le faailoaina o i latou agai i se faiga ua faaigoa arbitrary code execution<sup>9</sup>)
  - e te faatinoina le siakiina (scanning) ma le faanofoesea pe afai e aumaia models ia (third-party models) poo serialised weights, e tatau ona faia e pei o se isi pāti e lē faatuatuaina e ono mafai ai ona faaoga le remote code execution

- Afai o faaoga se APIs mai fafo (external) e te apalaiina siaki talafeagai ma faamama o faamatalaga ma mea e tuuina i totonu; e aofia ai le tuufaatasia o finagalo faaalua o tagata e faaogaina poo faamatalaga faaauau tau aoga i totonu o lau faataitaiga, e amanaia ai o faamatalaga mo aoga e fua i ai fausaga o amioga (system behaviour)
- e te apalai siaki talafeagai ma le faamamāina o faamatalaga ma mea e tuuina i totonu; e aofia ai ma le tuufaatasia o ni finagalo faaalua mai tagata e faaogaina le auunaga poo le tuuina o faamatalaga i lau faataitaiga, ma amanaia ai o faamatalaga mo toleniga ua fua agai iai le system o amioga
- e te tuufaatasia le AI software system ma lona faalauteleiga i totonu o faalauteleiga mautu ma auala aupito lelei mo le mautu o le faalauteleiga o ia faiga aupito lelei; o elemeni uma o le AI system o loo tusia i siosiomaga talafeagai e faaoga ai coding practices ma gagana e faaitiitia ai pe aveesea ai ituaiga tulaga maaleale
- Afai e manaomia AI components e faaosofia ai gaioga, mo se faataitaiga o le suia o faila poo le tuuina foi o mea e agai i fafo, e te talosagaina ni sa talafeagai i gaioga e pei ona tuuina mai (o lea e aofia ai le internal AI e tatau lava ona faailoa)
- o faaiuga e faatatau i le fesootaiga o loo faailoa mai e le AI-Specific risks, mo se faataitaiga:
  - o lau system na te saunia users e iai a latou usable outputs e aunoa ma se faailoaina o ituaiga e manaomia agai atu i se tagata e ono osofaia maia
  - afai e mafai, e saunia e lau system effective guardrails e faataaliolio i mea maua (outputs) mai au models
  - afai e ofoina se API i tagata faatau i fafo poo ni collaborators, e te apalai i ai faatonutonga talafeagai e faaitiitia ai osofaiga i le AI system e ala i le API
  - e te tuufaatasi ai ituaiga faatulagana aupito sili ona mautu i totonu o le system by default
  - e te apalai nai faamanuiaga mai manatu faavae e faaitiitia ai le avanoa i le faagaioia o le system
  - e te faamanino agavaa e iai lamatiaga agai i users ma manaomia users e tali atu i le mea lea ina ia mafai ai ona o latou faaogaina; e te fesootai agai i mataupu e faasaina, ma afai e mafai, faailoa i users o nisi tali

### Iloilo penefiti o tulaga tau le puipuiga ma fefaatauaiga ao filifili lau AI model



O lau filifiliga i le AI model o le a aofia ai le faapaleniina o nisi o mea manaomia. E aofia ai ma le filifiliga o le ituaiga architecture model, faatulagana, faamatalaga tau toleniga, polokalame tau aoga ma hyperparameters. O au faaiuga o le a faailoa e lau threat model, ma e fai lava ma toe iloilo ona o loo agai pea i luma soo se tasi o le AI security research advances ma le malamalama i faamata'u.

Afai e te filifilia se model a le AI, o au fuafuaga o le a ono aofia ai ae le faatapulaa ai:

- o le faigata o le model o loo e faaogaina, ona o le architecture ua faaogaina ma le aofai o parameters, o le chosen architecture ma le numera o parameters, o le a, ese mai isi mafuaaga, aafia ai le ituaiga faamatalaga o toleniga sa aoina ma e manaomia ai pe o le ā le tele o le faamatalaga e manaomia ae ole ā foi lona tulaga malosī o iai fua agai i suiga o faamatalaga tuuina i totonu pe a faaogaina
- O le sa'o o le model o loo e faaogaina ma le sutesuga o loo faatulaga agai i ai i ou manaoga (faapei o le fine tuning)
- o le agavaa e toe faatulaga ai, faaliliu ma faamalamalama le iuga o au models (faataitaiga, debugging, sueina poo le ioeina o le faatulafonoina); e ono iai ni penefiti ile faaogaina o models faigofie ma tonu nai lo models e faigata atu ona faaliliu
- o uiga o training dataset(s) e aofia ai le ituaiga tele o loo manaomia, le atamai, ituaiga tulaga maaleale, tausaga, ma tulaga e tali tutusa ma eseese ai



## 2. Secure development (faalautelega mautu)

O lenei vaega e aofia ai taiala e apalai i le **development** vaega o le faalautelega o le olaga o le AI system, e aofia ai le puipuiga o le supply chain, faapepaina ma aseta faapea le puleaina o aitalafu tau masini.

### Faamautu lau supply chain



E te iloilo ma mataitu le puipuiga o lau AI supply chains ile olaga o le system, ma e manaomia suppliers latou te usitaia le ituaiga tulaga lava lea o loo apalai i au faalapopotoga ma isi polokalame tau komipiuta (software). Afai e le usitaia e suppliers tulaga masani o lau faalapopotoga, e te gaiioi e fua i le faiga faavae tau pulega o loo iai lamatiaga.

I tulaga e lē o faia ai i totonu lava o le lotoifale, e te aumaia ma faamautuina meatotino o faamaumauga faamauina lelei o hardware ma software (mo se faaitaiga, models, faamaumauga, faletusi o polokalame tau komipiuta, modules, middleware, faavaa ma external APIs) mai isi vaega ua maea ona faamaoniaina faatasi ai ma isi third party developers e faamautuina ai le tulaga lelei o le puipuiga o au systems.

Ua e sauni e te faaaogaina nisi tali o avanoa mai mo systems tuga aafiaga, pe afai e lē ausia tulaga tau le puipuiga. E te faaoga punaoa e pei o le NCSC [Supply Chain Guidance](#) ma faavaa e pei o le Supply Chain Levels for Software Artifacts (SLSA)<sup>10</sup> mo le siakiina o tautinoga o le supply chain ma le faalauteleina o le software ma lona olaga.

### Mataitu, mulimulita'i ma puipui au aseta



E te malamalama i le tatau i lau faalapopotoga o au aseta e iai sootaga i le AI, e aofia ai models, faamaumauga, (e aofia ai ma finagalo faaalii mai tagata e faaaogaina le auunaga), faamanatu, software, faamaumauga, o api ma iloiloga (e aofia ai faamatalaga e faatatau i agavaa le talafeagai ma tulaga sa iai faaletonu) amanaia ai ma vaega o loo latou tulai mai ai o ni inivesi tatau ma poo fea tonu o le auala agai atu ia i latou e mafai ona sao ane ai se tagata osofai. E te tausia api o faamatalaga o ni faamaumauga maaleale ma faatino tulaga faatonutonu e puipui ai o latou tulaga faalilolilo, agavaa ma lo latou avanoa.

E te iloa le mea o loo iai au aseta ma ua mae'a ona e iloiloina ma taliaina nisi lamatiaga e iai ni ona sootaga. O loo iai au fuafuaga ma meafaigaluega e mulimulita'i ai, faamaonia, faatonutonu ma faamautu au aseta, ma e mafai foi ona toe faatumauiina se tulaga lelei i se taimi e tulai ai se mataupu e ono lamatia ai.

O loo iai au fuafuaga ma faatonutonu ua maea faatulaga e pulea ai poo a faamatalaga a le AI system e mafai ona tagofia, ma pulea mea o loo aofia ai e fua agai i tulaga maaleale o le AI (ma le tulaga maaleale o inputs lea na faaoga e faatino ai).

### Faapepa au faamaumauga, models ma faamanatu



E te faapepaina le faatuina, faagaioia, ma le faatonutonuina o le olaga o soo se model, seti o faamatalaga ma faamanatu o meta poo systems. O au faamaumauga e aofia ai faamatalaga e faatatau i le puipuiga e pei o le mafuaaga o faamatalaga mo toleniga (e aofia ai faamatalaga fine tune pe tagata foi poo isi foi finagalo faaalii), fuafua mo le faatapulaa, guardrails, faaoga o faailo tau faamaumauga poo ni saina e nana ai faamaumauga, faatumauiina o le taimi, fautuaga faatatau i se iloiloga ma models ua aafia pe iai ni faaletonu. O nisi faiga e fesoasoani ile faiga o lea gaiioiga, e aofia ai model cards, data cards ma software bills of materials (SBOMs). O le fausia o faamaumauga lavelave, e lagolago ai le faamalamalamaina o mea ma le mafai ona tali manino atu.<sup>11</sup>



## 3. Secure deployment (Faasoaina mautu)

O lenei vaega e aofia ai taiala e apalai i le **deployment** vaega o le faalauteleina o le olaga o le AI system, e aofia ai le puipuiga o lona siomaga ma models mai tulaga e ono afaina ai, faamata'u ma mea ua lusi/gau, faalauteleina o faiga o pulega o mea e tutupu ma lona faamatuuina atu.

### Faamautu lou siomaga



E te apalaia ni manatu faavae tau le puipuiga i le siomaga o loo faaogaina i soo se vaega o le olaga o lau system. E te apalaia le avanoa e talafeagai e faatonutonu ai au APIs, models, ma faamatalaga, ma le latou toleniga ma alapaipa o loo faasolosolo atu, ia faia ai suesuega ma faalautele aemaise o le faasoaina. E aofia 'i i le tuueseseina talafeagai o siosiomaga o loo umia ni faamatalaga maaleale. O le a fesoasoani foi lenei e faaititia osofaiga masani i tulaga tau le puipuiga o loo taulai e gaioia se model poo le faamoemoe e faaleaga lana gaioi.

### Ia puipui faauau lava lau model



E ono mafai e le 'auosofai ona toe fausia le faatulagana o se model<sup>13</sup> poo le faamatalaga o loo faia ai toleniga<sup>14</sup>, e ala i le tagofia tonu lava o se model (e ala i le faaogaina o model weights) pe (fesiligia foi le model e taula ane i se talosaga poo se auunaga). E ono mafai foi ona lotea e le 'auosofai o models, faamatalaga poo ni faamanatu ile taimi poo le mae'a o toleniga, ma faapea ai o le output ua le faatuatuaina.

E te puipuia le model ma faamatalaga mai auala sa'o ma lē sa'o foi , e ala i le:

- faatinoina o faiga masani tau cyber security
- faatinoina o ni mea faatonutonu e tusa o le faafesili ina ia faailoa pe taofia foi nisi taumafaiga e agai i totonu, fesuia'i ma faailoa faamatalaga faalilolilo

Ina ia faamautinoa e mafai e systems (consuming) ona faaogaina models, e te faasoaina codes i komipiuta ma faailoga (hash) pe o saina foi a faila o models (mo se faataitaiga, model weights) ma seti o faamatalaga (e aofia ai checkpoints) i le taimi lava e a'oa'oina ai le model. E masani ai lava i le cryptography (faaogaina o se code e nana ai faamatalaga) e taua lava le lelei o le faatonutonuina o le ki<sup>15</sup>.

O lau auala mo le faaititia o aafiaga faalilolilo o ni lamatiaga o le a faalagolago lava i le faaogaina o le case ma le threat model. O nisi talosaga, mo se faataitaiga e aofia ai nisi o faamatalaga maaleale, e ono manaomia ai theoretical guarantees (ni mea e faamaonia ai) ma e mafai ona faigata pe taugata foi ona apalai. Afai e talafeagai, o privacy-enhancing technologies (tai pei o differential privacy ma homomorphic encryption) e mafai ona faaogaina e asiasi ai pe faamautinoa ai ituaiga o lamatiaga o iai sootaga ma tagata e faaogaina auunaga, ma i latou e osofaia ma maua le avanoa e tagofia ai models ma outputs.

### Faalautele le faasologa o le puleaina o mataupu



O le tulaga mautinoa o mataupu tau le puipuiga e aafia ai AI systems e atagia lea i lau tali atu i le mataupu, faagasologa ma fuafuaga mo se toe faaleleiga. O au fuafuaga e atagia ai vaega eseese ma e matele ina toe iloilo pe a faasolo mai nisi o systems ma suesuega lautele. E te teuina punaoa tuga o faamaumauga a le kamupani i ni faila e lē o i luga o le initaneti. O i latou e tali atu i ni faalavelave (responders) ua mae'a ona toleniina ina ia iloilo ma tagofia mataupu faapea e afua mai i ni sootaga i le AI. E te saunia ni api o faamaumauga tulaga lelei ma isi ituaiga features tau le puipuiga poo ni faamatalaga agai i tagata faatau ma users e aunoa ma se totogi faaopopo, e faataga ai ni a latou tali atu i ni mataupu e tulai mai.

### Faamatuu atu ma le toto'a le AI



E te faamatuu atu models, talosaga ma systems pe afai nao le iloiloga o le tulaga tau le puipuiga o loo talafeagai ma outou e pei o le benchmarking ma le red teaming (faapea foi ma isi suesuega sa faia ona o nei taiala e pei o le saogalemu poo le amio sa'o), ma o loo e malamalama i au users ma faatapulaa poo ni modes ua faaletonu. O faamatalaga e faatatau i faletusi ma faataitaiga o tulaga tau le puipuiga, o loo tuuina atu i le [isi tusitusiga ile](#) faaiuga o lenei pepa.

### Ia faafaigofie mo users e faia le mea sa'o



O loo e maitauina o ituaiga faatulagana taitasi poo le configuration option e tatau ona iloiloina faatasi ma penefiti e maua ai a le pisinisi, ma nisi ituaiga lamatiaga e afua atu ai. O le mea e lelei, o le ituaiga faatulagana aupito sili ona mautu o le a tuuina faatasi i le system e nao le pau lea o le option. Afai e manaomia ona faatino le configuration, e tatau lava ona matuai faamautuina le default option agai i faamata'u taatele (e aofia ai le secure by default). E te apalaia tulaga faatonutonu e taofia ai le faaaogaina poo le faasoaina o lau system i ni auala le lelei.

E te sauniaina mo users ni taiala e faatatau i le faaaogaina talafeagai o lau model poo le system, lea e aofia ai le faamatamata tetele o tapulaa ma nisi o modes ua faaletonu. E te faailoa manino i users poo fea tonu vaega o le tulaga tau le puipuiga o loo feagai ma latou, ma o loo faamaoni foi poo fea (ae pe faapefea) ona faaaogaina ona faamaumauga, tagofia ma teuina (mo se faataitaiga, pe afai e faaaoga mo le faatumauina o le model, poo le iloiloina o tagata faigaluega poo ni paaga).

## 4. Secure operation and maintenance (faagaioiga ma le faatumauina mautu)

O lenei vaega e aofia ai taiala e apalai ile **secure operation and maintenance** vaega o le faalateleina o le olaga o le AI system. E saunia atu ai taiala i gaioiga e talafeagai ma patino tonu pe a faasoaina se system, e aofia ai ma le tatalaina ma mataituina, faafou faamatalaga faasoa tau pulega.

### Mata'itu le amio a lau system



E te fuaina le output (mea e ave i fafo) ma le gaioi a lau model ma le system, ina ia mafai ona e maitau ni suiga vave i le amio e aafia ai ma le puipuiga. E mafai ona e tautala pe iloa nisi o faalavelave ma lamatiaga e ono o'o atu, aemaise le faagaioiga o faamatalaga masani.

### Mata'itu inputs a lau system



Ina ia o gatasi ma tulaga faalilolilo ma le puipuiga manaomia mo faamatalaga. e te mata'ituina ma faamau inputs i lau system (e pei o faaiuga ua faia fua i faamaoniga, fesili poo ni faamanatu) e faataga ai le usitaia o matafaioi, suetusi, suesuega ma toe fetuunaiga pe afai e iai se mataupu i se tulaga lamatia pe faaoga sese foi. E ono aofia ai ii se iloa maumaututu o out-of-distribution poo adversarial inputs, e aofia ai ma i latou e fuafua e faaoga faamatalaga mo laasaga o tapenaga (e pei o le toe fetuunaiga o le tetele o ata).

### Mulimulita'i i faiga a le secure by design



E te faaafua faafouga e otomeki (by default) i soo se oloa ma faaoga tulaga mautu, fuafuaga o faafouga o modules mo lo latou tufatufaina. O au fuafuaga aupito lata mai (e aofia ai faataitaiga ma iloiloga) e atagia ai le mea moni o le suiga i faamatalaga, models poo faamanatu e ono o'o ai i ni suiga i amioga a le system (mo se faataitaiga, e te faia faafouga tetele e pei o ni versions fou). E te lagolagoina users ina ia toe siaki ma tali atu i suiga i models (mo se faataitaiga e ala i le saunia o se auala e muai silasila ai ma isi ituaiga APIs).

### Ao ma faasoa lesona ua a'o'ina



E te auai i le faasoaina o faamatalaga i le komiuniti, galulue faatasi i le siomaga o fale gaosi oloa faavaomalo, aoga ma faigamalo e faasoa ituaiga faiga aupito sili ma talafeagai. E te faatumauina le tatala o laina o fesootaiga mo le faailoa o finagalo e faatatau i le puipuiga o systems, e lē gata i totonu o lona lotoifale ao fafo foi i lau faalapotopotoga, e aofia ai le tuuina o le ioega i le 'ausuesue i mataupu tau puipuiga e suesue ai ma lipoti tulaga maaleale. Afai ae manaomia, e te faasolo atu ni mataupu i le komiuniti lautele, mo se faataitaiga, lomina o puletini e tali atu ai i mataupu tau tulaga maaleale, e aofia ai le auiliiliga ma le faamauina o tulaga maaleale taatele. E te gaioi e faaititia ma toe faalelei mataupu i se tulaga vave ma talafeagai.



# Nisi faitauga

## AI development (faalateleina o le AI)

[O manatu faavae mo le puipuiga o machine learning](#)

O le taiala a le NCSC e faatatau i le faalateleina, faasoaina ma le faagaioia o se system e iai sona ML component.

[Secure by design - Fetu'una'iga o le Paleni i Lamatiaga mai Osofaiga ile Initaneti: Manatu faavae ma Auala e faatino ai le polokalame o le Secure by Design](#)

Tusia faatasi e le CISA, NCSC ma isi ofisa, o leni taiala e faamatala atu ai pe faapefea ona faia laasaga e manufacturers o software systems, e aofia ai ma le AI, ina ia tuuina ai puipuiga i vaega o faatulagana o le faalateleina o le oloa, ma lafoina oloa e ave atu o loo saogalemu i fafo o a latou pusa.

[AI Security Concerns in a Nutshell](#)

Saunia mai e le German Federal Office for Information Security (BSI), o leni pepa e saunia atu ai se folasaga i ni osofaiga e ono tutupu i systems o machine learning ma ni puipuiga e ono faatino e tetee atu ai i na osofaiga.

[O le Hiroshima Process International Guiding Principles for Organizations Developing Advanced AI Systems ma le Hiroshima Process International Code of Conduct for Organizations Developing Advanced AI Systems](#)

O nei pepa, ua saunia o se vaega o le G7 Hiroshima AI Process, e saunia atu ai taiala mo faalapotopotoga o loo latou faalateleina ituaiga systems aupito lelei o le AI, e aofia ai ma faataitaiga (models) aupito lelei ma AI systems ma le faamoemoe e faalauiloa lona saogalemu, mautu ma le faamoemoeina o le AI ile lalolagi atoa.

[AI Verify](#)

O le AI Governance Testing Framework and Software toolkit a Singapore e faataga ai le gaiio o AI systems faasagatau atu i se seti o manatu faavae e iloa tele faavaomalo e ala i faataitaiga masani.

[Multilayer Framework for Good Cybersecurity Practices for AI – ENISA \(europa.eu\)](#)

O se faavaa e taitai ai National Competent Authorities ma AI stakeholders i laasaga e manaomia ona o latou mulimulita'i i ai ina ia mautu ai a latou AI systems, faatinoga ma fuafuaga.

[ISO 5338: O AI system life cycle processes \(Lea o loo iloiloaina\)](#)

O se seti o fuafuaga e iai sootaga mo le faamatalaina o le olaga o le AI system e fua i machine learning ma ituaiga o systems e faaigoa heuristics systems.

[AI Cloud Service Compliance Criteria Catalogue \(AIC4\)](#)

O le BSI's AI Cloud Service Compliance Criteria Catalogue na te saunia ni aiaiga patino o le AI, e mafai ai ona faatino le iloiloa o le tulaga tau le puipuiga o se auunaga tau AI i lona olaga atoa.

[NIST IR 8269 \(Draft\) A Taxonomy and Terminology of Adversarial Machine Learning](#)

O se seti o fuafuaga e iai sootaga i ni faiga e faamatala ai le olaga o le AI system e fua i luga o machine learning ma heuristics systems.

[MITRE ATLAS](#)

O se faavae o malamalamaaga e faatatau i metotia a le fili, o lona iloa, ma nisi o suesuega sa faia faatatau i machine learning (ML) ua faataitaia i le maea o sootaga ma le faavaa o le MITRE ATT&CK.

[An Overview of Catastrophic AI Risks \(2023\)](#)

Saunia e le Centre for AI Safety, o leni pepa ua faatulaga atu ai vaega o lamatiaga mo le AI.

[Gagana o Faata'ita'iga Tetele: Opportunities and Risks for Industry and Authorities](#)

O le pepa na saunia e le BSI mo kamupani, pulega ma developers e mananao e fia silafia atili i nisi avanoa ma lamatiaga mo le faalateleina, faasoaina ma le faaogaina o LLMs.

Nisi o poloketi o loo tatala mea na aumai ai e fesoasoani i users e faata'ita'i AI models e aafia ai:

- [Adversarial Robustness Toolbox](#) (IBM)
- [CleverHans](#) (University of Toronto)
- [TextAttack](#) (University of Virginia)
- [Prompt Bench](#) (Microsoft)
- [Counterfit](#) (Microsoft)
- [AI Verify](#) (Infocomm Media Development Authority, Singapore)

## Cyber security

[O le CISA's Cybersecurity Performance Goals](#)

O se seti taatele o puipuiga e tatau ona faatino e pisinisi uma ina ia uiga ai le faaitiitia o le ono aafia i lamatiaga o loo iloa ma metotia a le fili.

[NCSC CAF Framework](#)

The Cyber Assessment Framework (CAF) na te saunia taiala mo faalapotopotoga e feagai ma auaunaga aupito taua ma mea fai.

[MITRE's Supply Chain Security Framework](#)

O se faavaa mo le iloiloina o suppliers ma i latou e saunia auaunaga i totonu o le supply chain.

## Puleaina o lamatiaga

[NIST AI Risk Management Framework \(AI RMF\)](#)

O le AI RMF e faaata mai ai pe faapefea ona faatonutonu lamatiaga e taua o socio-technical risks i tagata taitoatasi, faalapotopotoga ma le sosaiete o tagata e iai sootaga i le AI.

[ISO 27001: Faamatalaga tau le puipuiga, cybersecurity ma le puipuiga o tulaga faailolilo](#)

O lenei puipuiga e saunia ai faalapotopotoga i ni taiala e faatatau i le faatuina, faatinoina ma le faatumauiina o se system e pulea ai faamatalaga tau le puipuiga.

[ISO 31000: O le faatonutonuina o lamatiaga](#)

O se faatulagana faavaomalo e saunia ai i faalapotopotoga ni taiala ma ni manatu faavae mo le faatonutonuina o lamatiaga i totonu o faalapotopotoga.

[NCSC Risk Management Guidance \(Taiala a le NCSC mo le faatonutonuina o lamatiaga\)](#)

O lenei taiala e fesoasoani ai i practitioners ia malamalama lelei ma pulea ituaiga lamatiaga o loo aafia ai a latou faalapotopotoga.

# Faamaumauga

1. O loo faamatala atu 'i i o se tagata, se pulega lautele, ofisa poo soo se isi totino latou te faalauteleina le AI system (pe ua iai se AI system ua uma ona faalautele) ma tuuina ai le system i luga o le maketi pe tuu foi i se auunaga i lalo o lona ia lava igoa poo se faamaufaaailoga
2. Mo nisi faamatalaga e faatatau i le secure by design, silasila i le upega tafailagi a le CISA's [Secure by Design](#) ma taiala [Fetu'una'iga o le Paleni i lamatiaga mai osofaiga i luga o initaneti: Manatu faavae ma Auala e faatino ai le polokalame o le Secure by Design](#)
3. E ese mai i auala faatino a le non-ML AI e pei o systems e faavae mai tulafono
4. Ua faamatala e le CEPS ituaiga e fitu o fesootaiga tau AI ma lona faalauteleina i a latou lomiga ['Reconciling the AI Value Chain with the EU's Artificial Intelligence Act'](#)
5. [ISO/IEC 22989:2022\(en\)](#) e faamatala ai o se 'elemeni lenei o loo gaioi e faatulaga e pei o se AI system'.
6. Ua tuuina le galuega i le NIST latou te saunia ai taiala (ma le faia o isi gaioga) ina ia agai i luma saogalemu, mautu ma faatuatuaina le faalauteleina ma le faaogaina o Artificial Intelligence (AI). [Silasila i Matafaioi a le NIST's i lalo o le Poloaga \(Executive Order\) o le aso 30 o Oketopa, 2023](#)
7. O nisi faamatalaga e faatatau i faaitaiga o faamata'u (threat modelling) o loo avanoa mai mai le [OWASP Foundation](#)
8. Silasila i le MITRE ATLAS [Adversarial Machine Learning 101](#)
9. GitHub: [RCE PoC for Tensorflow using a malicious Lambda layer](#)
10. SLSA: ['Malupuipua o le tulaga lelei o artifacts i soo se vaega o le supply chain'](#)
11. O le METI (Japanese Ministry of Economy, Trade and Industry, 2023), ['Guide of Introduction of Software Bill of Materials \(SBOM\) for Software Management'](#)
12. Sailiga a le Google: [Machine Learning \(Aoaoga o le Masini\): The High Interest Credit Card of Technical Debt](#)
13. Tramèr et al 2016, [Stealing Machine Learning Models via Prediction APIs](#)
14. Boenisch, 2020, [Attacks against Machine Learning Privacy \(Part 1\): Model Inversion Attacks with the IBM-ART Framework](#)
15. National Cyber Security Centre, 2020, [Design and build a privately hosted Public Key Infrastructure](#)

---

© Crown copyright 2023. O ata ma ata o faamatalaga e ono aofia ai ma nisi o mea i lalo o le laisene mai isi pāti ma e lē avanoa mo le toe faaaogaina. O faamatalaga o loo aofia o loo laiseneina mo le toe faaaogaina i lalo o le Open Government Licence v3.0  
(<https://www.nationalarchives.gov.uk/doc/open-government-licence/version/3/>)

