# (Gen)AI security & privacy controls

Developments at ISO/IEC, CEN/CENELEC and OWASP

Rob van der Veer

**Rob van der Veer**

Senior director AI, security & privacy
Software Improvement Group

> 32 years experience AI, security & privacy

> ISO/IEC:

> > 5338 - AI lifecycle (Lead author)

> > 27090 - AI security

> > 27091 - AI privacy

> Advisor ENISA, Dutch NCSC, CIP

> OWASP: SAMM, AI Exchange, ML top 10, OpenCRE.org

> CEN/CENELEC JTC21/WG 5  (EU AI Act cybersec requirements)



**Software Improvement Group**

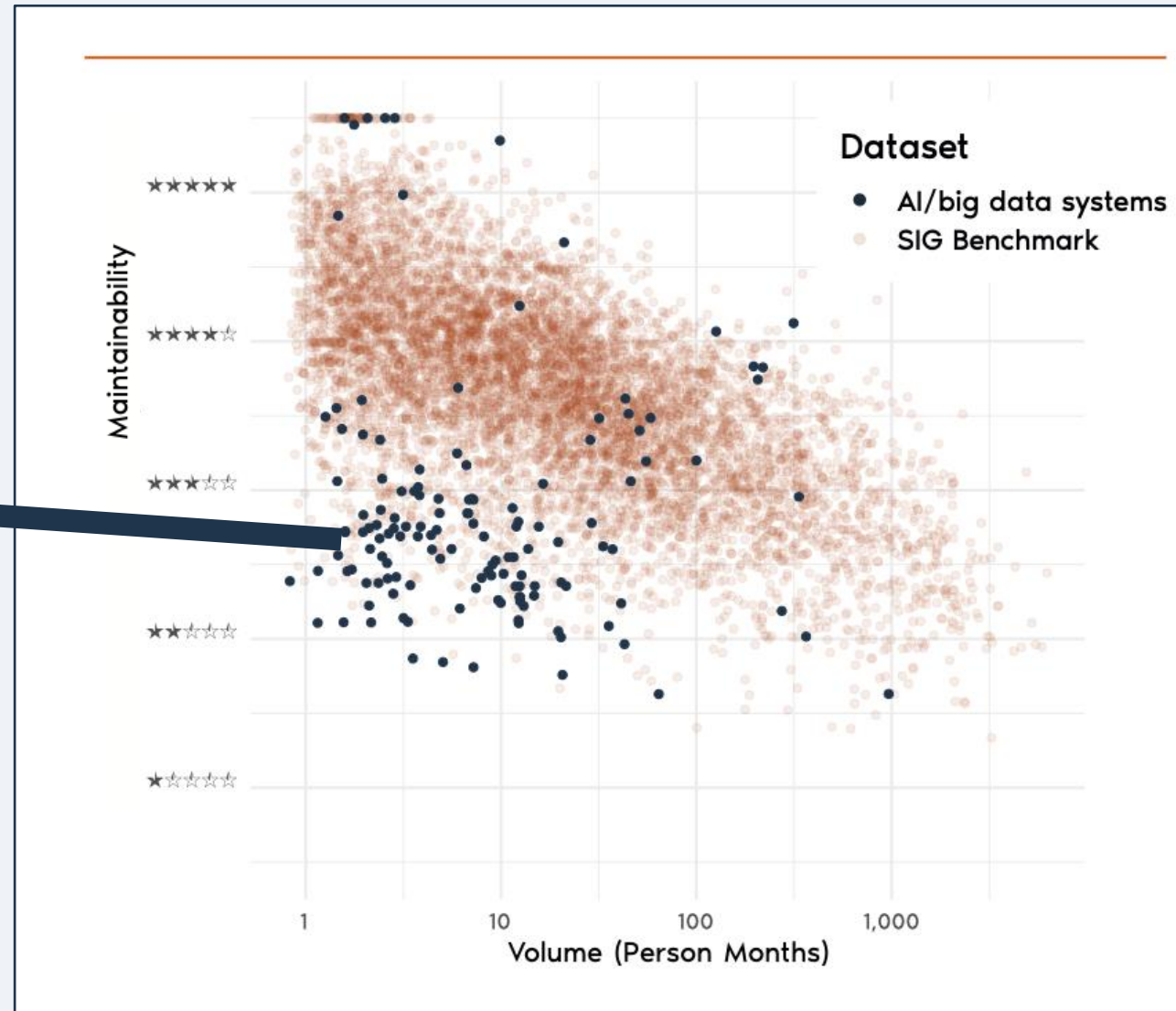rob.vanderveer@softwareimprovementgroup.com

@robvanderveer

+31 6 20437187

# Why ISO/IEC 5338: Software engineering best practices are often lacking in AI

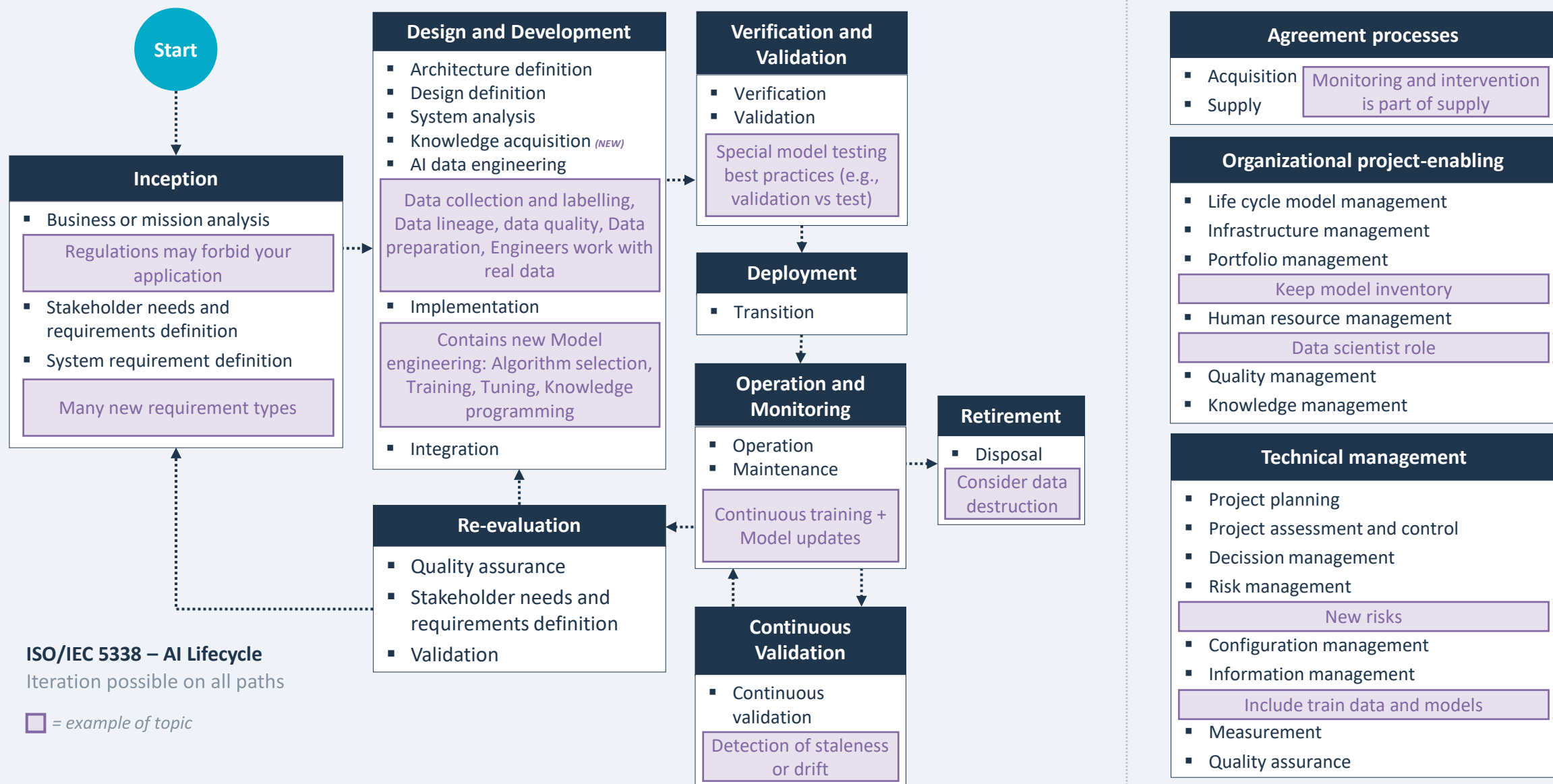**AI systems built by
data scientists in the lab**

- Bad maintainability
- Zero test code
- Lack of security
- Privacy issues
- Scalability issues
- Undocumented

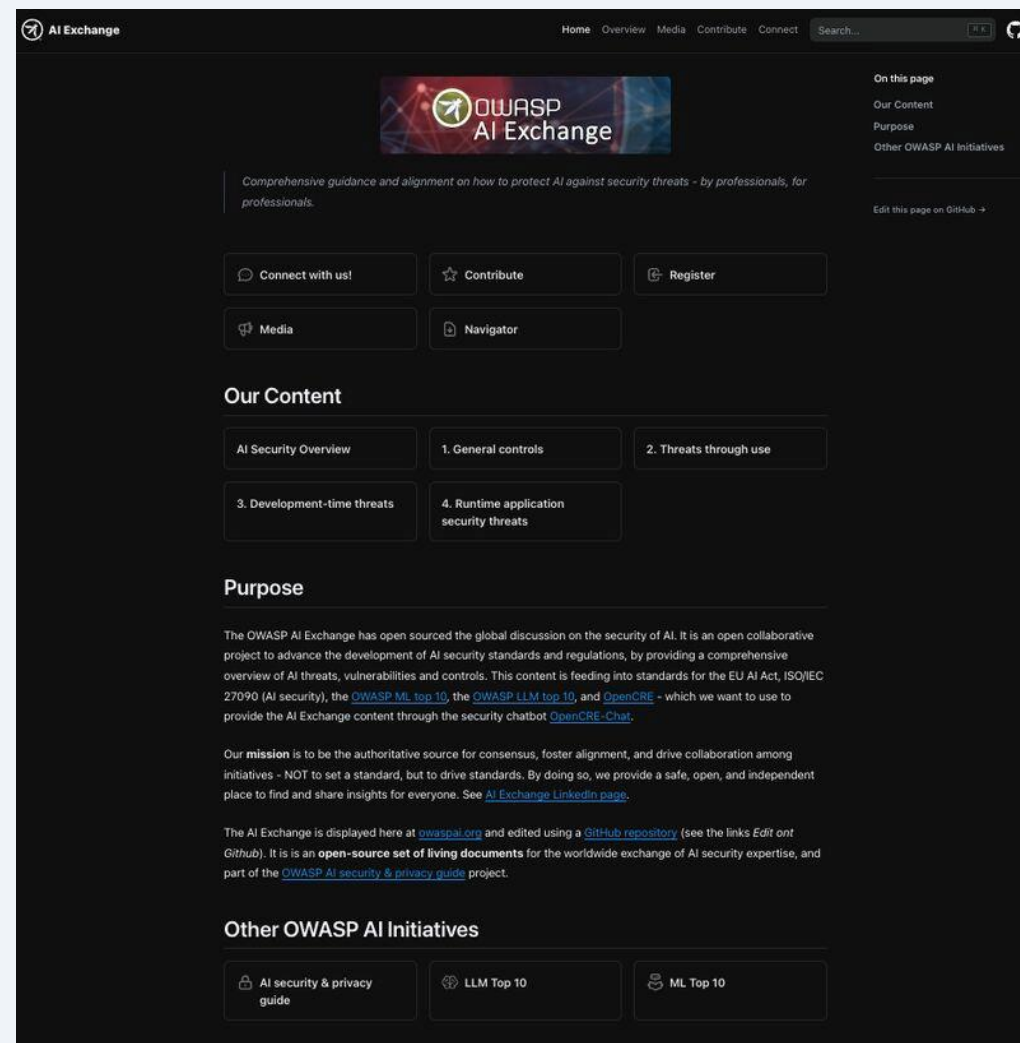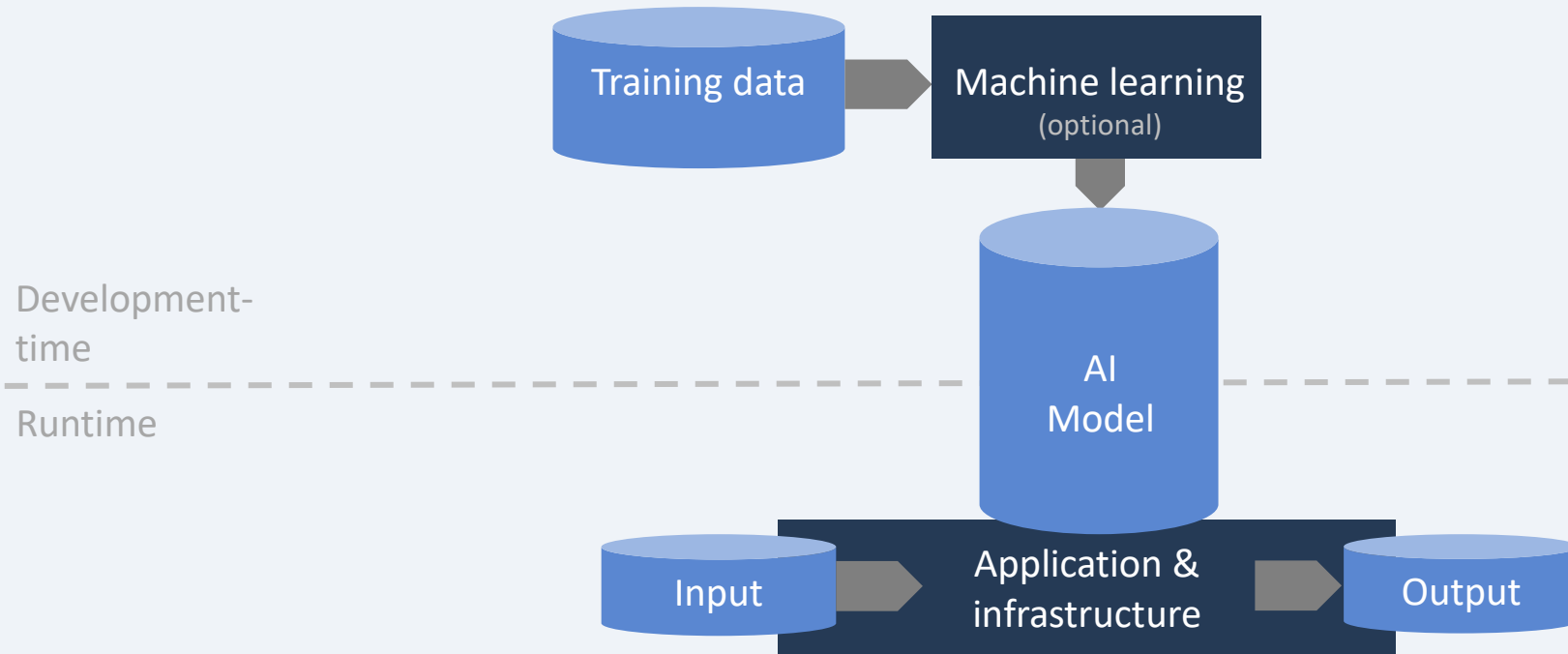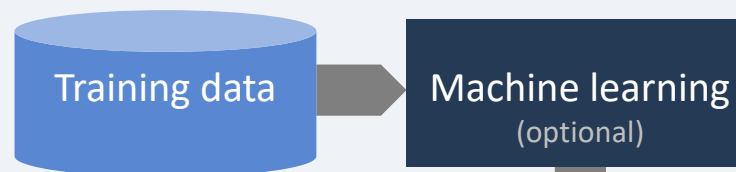**An accident or failure waiting to happen**

# ISO/IEC 5338 on AI system lifecycle: AI particularities for 12207 (software lifecycle)

**Start**

## Inception

- Business or mission analysis

  *Regulations may forbid your application*

- Stakeholder needs and requirements definition
- System requirement definition

  *Many new requirement types*

## Design and Development

- Architecture definition
- Design definition
- System analysis
- Knowledge acquisition *(NEW)*
- AI data engineering

  *Data collection and labelling, Data lineage, data quality, Data preparation, Engineers work with real data*

- Implementation

  *Contains new Model engineering: Algorithm selection, Training, Tuning, Knowledge programming*

- Integration

## Re-evaluation

- Quality assurance
- Stakeholder needs and requirements definition
- Validation

## Verification and Validation

- Verification
- Validation

  *Special model testing best practices (e.g., validation vs test)*

## Deployment

- Transition

## Operation and Monitoring

- Operation
- Maintenance

  *Continuous training + Model updates*

## Continuous Validation

- Continuous validation

  *Detection of staleness or drift*

## Retirement

- Disposal

  *Consider data destruction*

## Agreement processes

- Acquisition
- Supply

  *Monitoring and intervention is part of supply*

## Organizational project-enabling

- Life cycle model management
- Infrastructure management
- Portfolio management

  *Keep model inventory*

- Human resource management

  *Data scientist role*

- Quality management
- Knowledge management

## Technical management

- Project planning
- Project assessment and control
- Decission management
- Risk management

  *New risks*

- Configuration management
- Information management

  *Include train data and models*

- Measurement
- Quality assurance

**ISO/IEC 5338 – AI Lifecycle**
Iteration possible on all paths

☐ = *example of topic*

SIG

4

# OWASP AI Exchange at owaspai.org

- Well-received OWASP project since Q4 2022

- Open sources the AI security discussion

- Cross-geo, cross-domain, cross-industry

- Copyright and attribution-free

- Active liaisons:
  CEN/CENELEC, CSA, ISO 27090, NIST, MITRE, NCSCs, AISIC

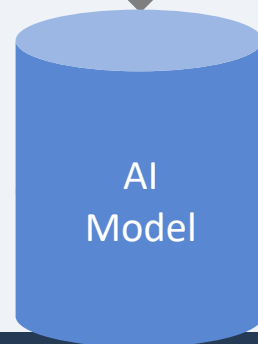- Content adopted verbatim by CEN/CENELEC JTC 21 for the AI Act

Training data → Machine learning (optional) → AI Model

Development-time

Runtime

Input → Application & infrastructure → Output

Training data → Machine learning (optional) → AI Model

Development-time

Runtime

Application & infrastructure

Input → Output

**Threats through use:**

- Evasion[B]
- Model theft[P]
- Model inversion [T]
- Data disclosure[T]
- Membership inference[T]
- Denial of model service[A]
- Prompt injection[B]

**Impact legend:**

(T) Train data confidentiality
(B) Model behaviour
(P) Intellectual property
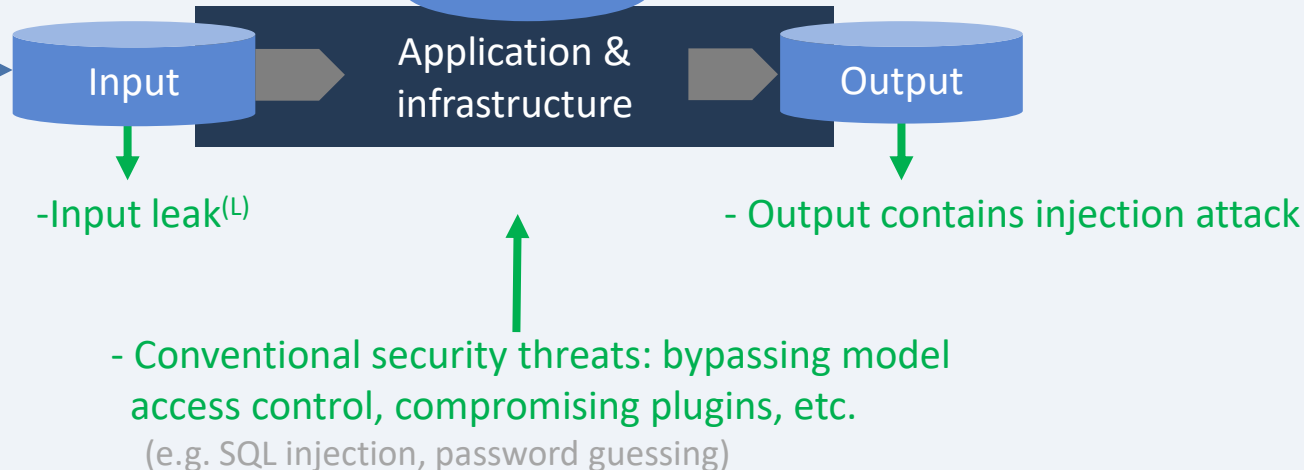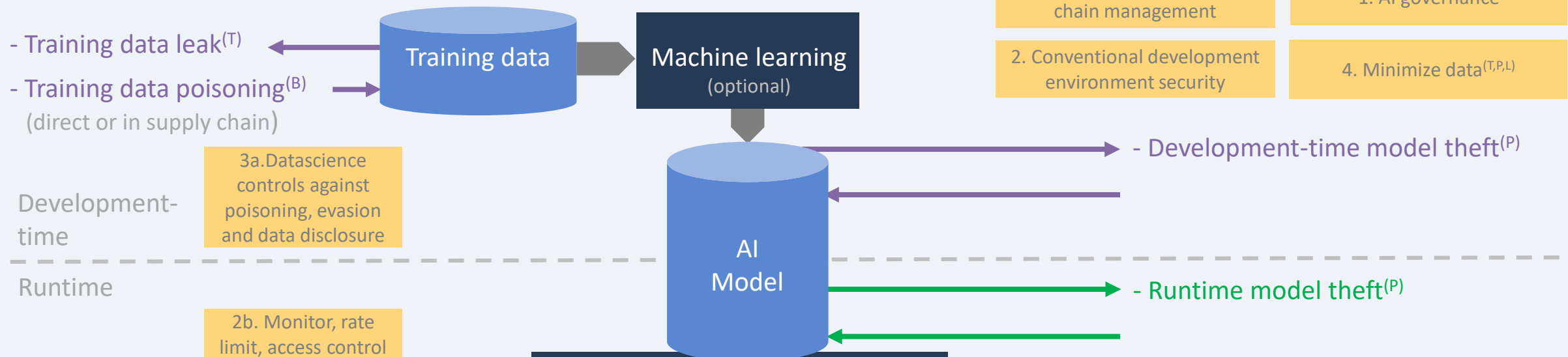(A) Availability
(L) Input confidentiality

➡ = threat

# Development-time threats

- Training data leak[T]

- Training data poisoning[B]
  (direct or in supply chain)

**Training data** → **Machine learning** (optional)

Development-time

Runtime

**AI Model**

- Development-time model theft[P]

- Development-time model poisoning[B]
  (direct or in supply chain)

**Impact legend:**

(T) Train data confidentiality
(B) Model behaviour
(P) Intellectual property
(A) Availability
(L) Input confidentiality

➡ = threat

Development-time

Runtime

Training data → Machine learning (optional)

AI Model

- Runtime model theft$^{(P)}$

- Runtime model poisoning$^{(B)}$

Input → Application & infrastructure → Output

-Input leak$^{(L)}$

- Output contains injection attack

- Conventional security threats: bypassing model access control, compromising plugins, etc.
(e.g. SQL injection, password guessing)

**Runtime security threats**

**Impact legend:**

(T) Train data confidentiality
(B) Model behaviour
(P) Intellectual property
(A) Availability
(L) Input confidentiality

➡ = threat

**Development-time threats**

- Training data leak$^{(T)}$

- Training data poisoning$^{(B)}$
 (direct or in supply chain)

Training data → Machine learning (optional)

- Development-time model theft$^{(P)}$

- Development-time model poisoning$^{(B)}$
 (direct or in supply chain)

Development-time

Runtime

AI Model

- Runtime model theft$^{(P)}$

- Runtime model poisoning$^{(B)}$

**Threats through use:**

- Evasion$^{(B)}$

- Model theft$^{(P)}$

- Model inversion $^{(T)}$

- Data disclosure$^{(T)}$

- Membership inference$^{(T)}$

- Denial of model service$^{(A)}$

- Prompt injection$^{(B)}$

Input → Application & infrastructure → Output

-Input leak$^{(L)}$

- Output contains injection attack

- Conventional security threats: bypassing model access control, compromising plugins, etc.
 (e.g. SQL injection, password guessing)

**Runtime security threats**

**Impact legend:**

(T) Train data confidentiality
(B) Model behaviour
(P) Intellectual property
(A) Availability
(L) Input confidentiality

➡ = threat

**Development-time threats**

- Training data leak$^{(T)}$
- Training data poisoning$^{(B)}$
  (direct or in supply chain)

Training data → Machine learning (optional)

2. Model and data supply chain management

1. AI governance

2. Conventional development environment security

4. Minimize data$^{(T,P,L)}$

3a. Datascience controls against poisoning, evasion and data disclosure

Development-time

- Development-time model theft$^{(P)}$

AI Model

Runtime

- Runtime model theft$^{(P)}$

2b. Monitor, rate limit, access control

**Threats through use:**

- Evasion$^{(B)}$
- Model theft$^{(P)}$
- Model inversion $^{(T)}$
- Data disclosure$^{(T)}$
- Membership inference$^{(T)}$
- Denial of model service$^{(A)}$
- Prompt injection$^{(B)}$

3b. Datascience input filtering and detection

Input → Application & infrastructure → Output

5. Control behaviour impact e.g. oversight, validation $^{(B)}$

-Input leak$^{(L)}$

- Output contains injection attack

2. Runtime technical security: conventional + new

4. Minimize data$^{(T,P,L)}$

- Conventional security threats: bypassing model access control, compromising plugins, etc.
  (e.g. SQL injection, password guessing)

**Runtime security threats**

**Impact legend:**

(T) Train data confidentiality
(B) Model behaviour
(P) Intellectual property
(A) Availability
(L) Input confidentiality

➡ = threat

= control group

# Control groups with individual controls - OWASP AI Exchange

1.**Governance processes**: AI risk management, Information security management, AI (secure) software engineering:
   (AIPROGRAM, SECPROGRAM, DEVPROGRAM, SECDEVPROGRAM, CHECKCOMPLIANCE, SECEDUCATE)

2.**Technical IT security controls**:
   a. Apply conventional IT security controls for **AI-specific parts**:
      • Development-time: model & data storage, model & data supply chain, data science documentation
        (DISCRETE, DEVDATAPROTECT, DEVSECURITY, SEGREGATEDATA, SUPPLYCHAINMANAGE)
      • Runtime: model storage, model use and model IO
        (RUNTIMEMODELINTEGRITY, RUNTIMEMODELIOINTEGRITY, RUNTIMEMODELCONFIDENTIALITY, MODELINPUTCONFIDENTIALITY, ENCODEMODELOUTPUT, LIMITRESOURCES)
   b. **Adapt conventional** IT security controls (rate limit, monitor)
      (MONITORUSE, MODELACCESSCONTROL, RATELIMIT)
   c. **Adopt new** IT security controls
      (CONFCOMPUTE, MODELOBFUSCATION, PROMPTINPUTVALIDATION, INPUTSEGREGATION)

3.**Datascience security controls:**
   a. **Development-time** controls when developing the model
      (DATAQUALITYCONTROL, FEDERATIVELEARNING, EVASIONROBUSTMODEL, POISIONROBUSTMODEL, TRAINADVERSARIAL, TRAINDATADISTORTION, ADVERSARIALROBUSTDISTILLATION, FILTERSENSITIVETRAINDATA, MODELENSEMBLE, MORETRAINDATA, SMALLMODEL)
   b. **Runtime** controls when running the model:
      (DETECTODDINPUT, DETECTADVERSARIALINPUT, DOSINPUTVALIDATION, INPUTDISTORTION, FILTERSENSITIVEMODELOUTPUT, OBSCURECONFIDENCE)

4.**Limit the amount of data** and the time it is stored
   (DATAMINIMIZE, ALLOWEDDATA, SHORTRETAIN, OBFUSCATETRAININGDATA)

5. **Limit the effect** of unwanted model behaviour (e.g. oversight, validation)
   (OVERSIGHT, LEASTMODELPRIVILEGE, CONTINUOUSVALIDATION, AITRANSPARENCY, EXPLAINABILITY, UNWANTEDBIASTESTING)

# How about privacy controls?

**Four types of privacy controls:**

1. AI security controls to **protect personal data**

2. AI security controls to **minimize data**

3. AI security controls to **limit impact** of model behaviour

4. **Additional controls to protect individual privacy rights:**

   ▪ Validate purpose (e.g. repurposing personal data)

   ▪ Arrange consent

   ▪ Control unwanted bias

   ▪ Provide transparency/explanation

   ▪ Attain data accuracy & updating

   ▪ Provide features to correct, access, erase, and object

Source: OWASP AI Security & privacy guide

# SIG Software Improvement Group .com

**GETTING SOFTWARE RIGHT** FOR A HEALTHIER DIGITAL WORLD