# Technical standards

# for AI governance

อาทิตย์ สุริยะวงศ์กุล
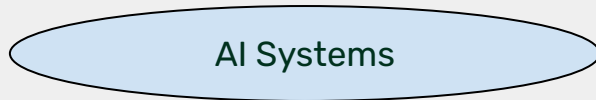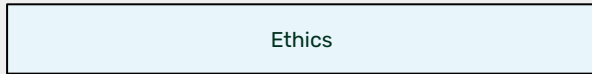Arthit Suriyawongkul, ADAPT Centre, Trinity College Dublin
SFI Centre for Research Training in Digitally-Enhanced Reality (d-real)

*20 September 2024 - Office of the Council of State International Symposium 2024, Bangkok*

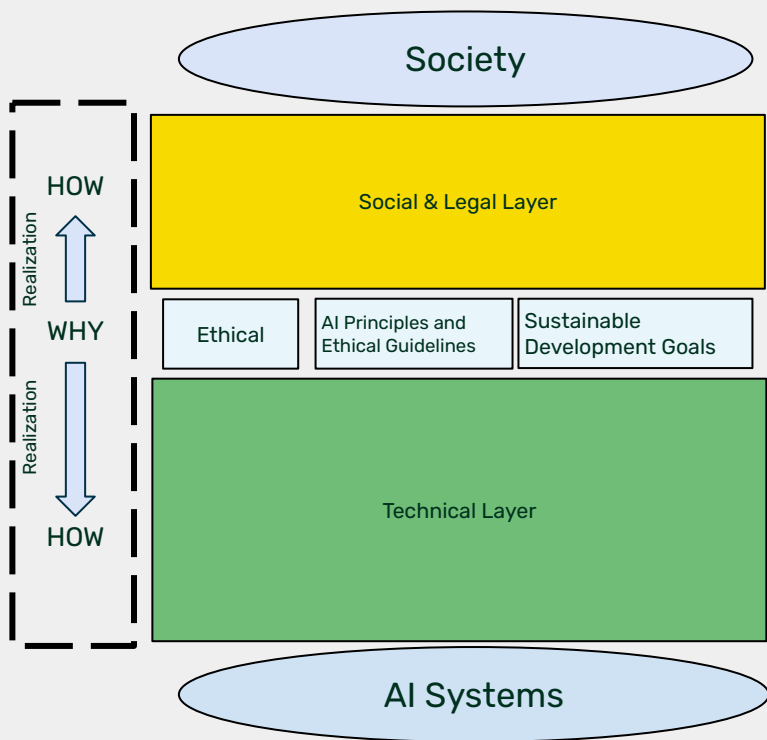Society

AI Systems

Society

Ethics

AI Systems

**Society**

**Social & Legal Layer**

| Ethical | AI Principles and Ethical Guidelines | Sustainable Development Goals |

**Technical Layer**

**AI Systems**

HOW

Realization

WHY

Realization

HOW

# A Layered Model for AI Governance



**HOW** · **WHY** · **HOW**

Realization · Realization

Society/Sector A

Society/Sector B

**Social & Legal Layer**
- EU AI Act
- Charter of Fundamental Rights of EU
- Local Norms
- UNESCO's Recommendation on the Ethics of AI

**Ethical**
- AI Principles and Ethical Guidelines
- Sustainable Development Goals

**Technical Layer**

Technical Standards
- NIST AI Risk Management Framework 1.0
- AI Life Cycle Processes (ISO/IEC 5338:2023)
- Data Privacy Vocabulary
- Data quality for analytics and ML (ISO/IEC FDIS 5259)

Software
- AI Verify toolkit
- Logging
- Inspect AI

AI System 1 · AI System 2

Adaptable and exchangeable across different jurisdictions

Technology to support the realization of AI principles

# AI Accountability

**Intergovernmental**

### ASEAN Guide on AI Governance and Ethics

1. Transparency and Explainability

2. Fairness and Equity

3. Security and Safety

4. Human-centricity

5. Privacy and Data Governance

6. Accountability and Integrity

7. Robustness and Reliability

**Governmental**

### Thailand AI Ethics Principles (MDES)

1. Competitiveness and Sustainability Development
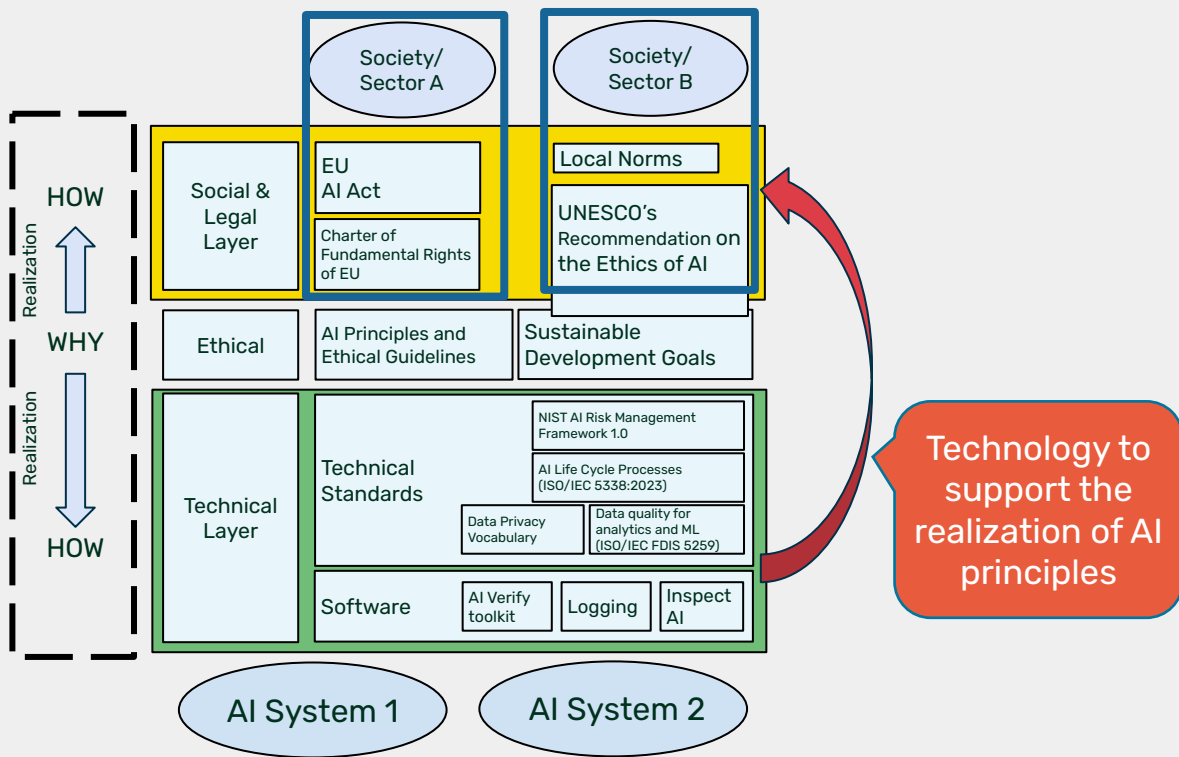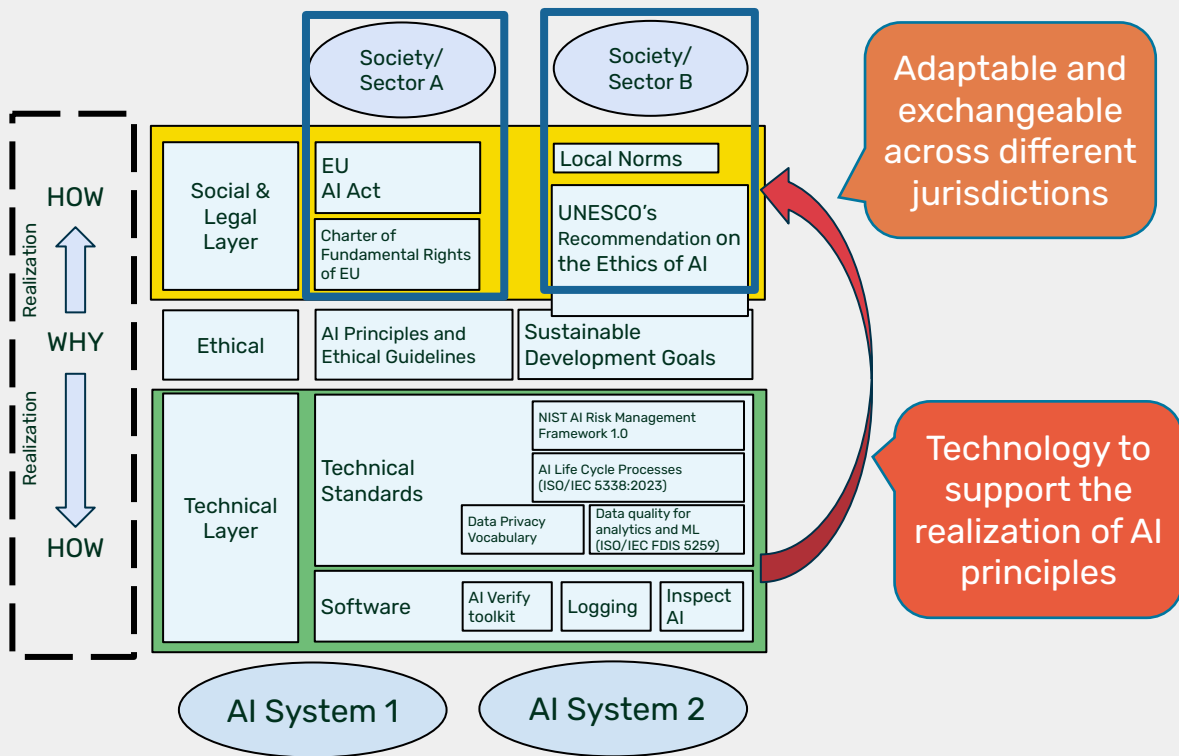
2. Laws, Ethics, and International Standards

3. Transparency and Accountability

4. Security and Privacy

5. Fairness

6. Reliability

**National research body / Grant-making agency**

### NSTDA AI Ethics Principles

1. Privacy

2. Security and Safety

3. Reliability

4. Fairness and non-discrimination

5. Transparency and Explainability

6. Accountability

7. Human Oversight and Human Agency

**Non-profit organization / Technical community**

### LF AI & Data's Principles for Trusted AI

1. Reproducibility

2. Robustness

3. Equitability

4. Privacy

5. Explainability

6. Accountability

7. Transparency

8. Security

**Intergovernmental**

### ASEAN Guide on AI Governance and Ethics

1. Transparency and Explainability

2. Fairness and Equity

3. Security and Safety

4. Human-centricity

5. Privacy and Data Governance

6. **Accountability** and Integrity

7. Robustness and Reliability

**National research body / Grant-making agency**

### NSTDA AI Ethics Principles

1. Privacy

2. Security and Safety

3. Reliability

4. Fairness and non-discrimination

5. Transparency and Explainability

6. **Accountability**

7. Human Oversight and Human Agency

**Governmental**

### Thailand AI Ethics Principles (MDES)

1. Competitiveness and Sustainability Development

2. Laws, Ethics, and International Standards

3. Transparency and **Accountability**

4. Security and Privacy

5. Fairness

6. Reliability

**Non-profit organization / Technical community**

### LF AI & Data's Principles for Trusted AI

1. Reproducibility

2. Robustness

3. Equitability

4. Privacy

5. Explainability

6. **Accountability**

7. Transparency

8. Security

**Accountability**
the fact of being responsible for what you do and <u>able to give a satisfactory reason for it</u>

# 3 Levels of AI Transparency

These three levels of AI transparency are working together and impact accountability and human oversight.

## Algorithmic transparency

- The ability to access and scrutinise code, data sets, and accompanying systems.

- Output like probabilities and charts from AI explainability methods (like LIME* and SHAP**) may be relevant to domain-experts and auditors/regulators, but not accessible to a person without background in AI or in the domain.

**Risk**

## Interaction transparency

- The ability to understand the strengths and limitations of an AI system, through the knowledge exchange between the AI system and its users.

- Tangibility, relevant metaphors to make sense of the environment, and the design paradigm that knowledge (transparency) is **co-created during an interaction, form a compelling basis for interaction transpar**

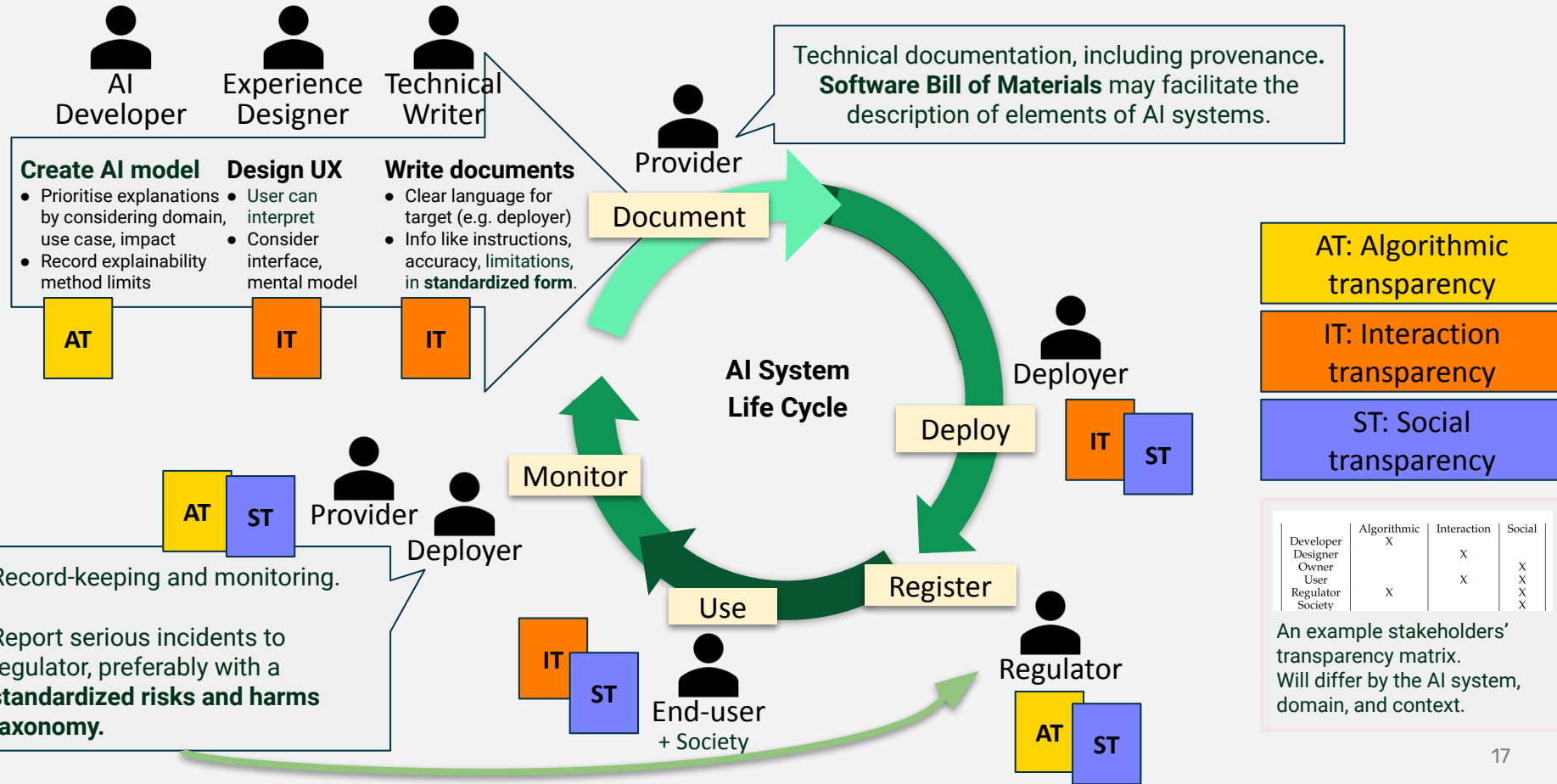**Risk**

## Social transparency

- The legal and cultural ability of the society/social institutions to understand and formulate responses to the use of an AI system.

- Institutionalised solutions that will **not overload information to users.**

**Risk**

* Local Interpretable Model-agnostic Explanations
** SHapley Additive exPlanations

# Transparency for Accountability in AI Life-Cycle

**AI Developer**

**Experience Designer**

**Technical Writer**

**Create AI model**
- Prioritise explanations by considering domain, use case, impact
- Record explainability method limits

**AT**

**Design UX**
- User can interpret
- Consider interface, mental model

**IT**

**Write documents**
- Clear language for target (e.g. deployer)
- Info like instructions, accuracy, limitations, in **standardized form**.

**IT**

**Provider**

Technical documentation, including provenance. **Software Bill of Materials** may facilitate the description of elements of AI systems.

**Document**

**AI System Life Cycle**

**Deployer**

**IT** **ST**

**Deploy**

**Monitor**

**AT** **ST** **Provider** **Deployer**

Record-keeping and monitoring.

Report serious incidents to regulator, preferably with a **standardized risks and harms taxonomy.**

**Use**

**IT** **ST**

**End-user + Society**

**Register**

**Regulator**

**AT** **ST**

AT: Algorithmic transparency

IT: Interaction transparency

ST: Social transparency

| | Algorithmic | Interaction | Social |
|---|---|---|---|
| Developer | X | | |
| Designer | | X | |
| Owner | | | X |
| User | X | X | X |
| Regulator | | X | X |
| Society | | | X |

An example stakeholders' transparency matrix.
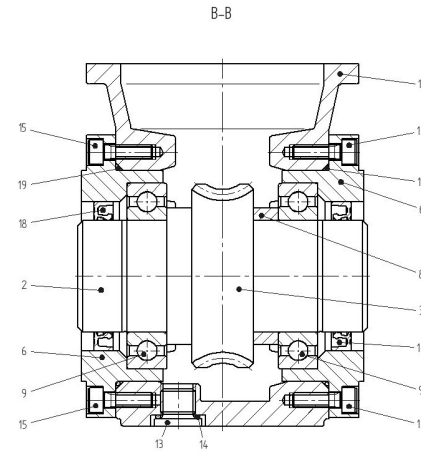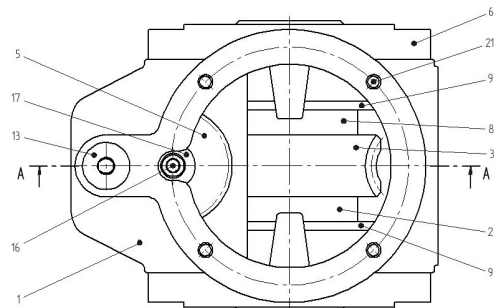Will differ by the AI system, domain, and context.

17

# Information obligations in EU AI Act that can support accountability (partial)

| For high-risk AI systems |
| --- |
| Provider name, registered trade name |
| Intended purpose |
| Instruction for use |
| Design choices |
| Standards applicable |
| Data origin, Collection original purpose |
| Possible biases, Measures to detect |

| For general purpose AI models |
| --- |
| Intended tasks, Limitations |
| Instruction for use |
| Model design specification |
| Training process, Testing process |
| Information on the data used |
| Copyright protection policy |
| Acceptable use policies applicable |

# Standards and Tools

A–A

B–B



| | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Pos. | Menge | Einheit | Benennung | Sachnummer / Norm – Kurzbezeichnung | | Bemerkung |
| 1 | 1 | Stck. | Gehäuse | | | G — AlSi10Mg |
| 2 | 1 | Stck. | Hohlwelle | | | 34CrMo4 |
| 3 | 1 | Stck. | Schneckenrad | | | G — CuSn12Ni |
| 4 | 1 | Stck. | Schneckenwelle | | | 16MnCr5 |
| 5 | 1 | Stck. | Zahnrad | | | 16MnCr5 |
| 6 | 2 | Stck. | Lagerdeckel groß | | | S235JR |
| 7 | 1 | Stck. | Lagerdeckel klein | | | S235JR |
| 8 | 1 | Stck. | Distanzring | | | S235JR |
| 9 | 2 | Stck. | Rillenkugellager | DIN 625 — 6009 | | |
| 10 | 2 | Stck. | Kegelrollenlager | DIN 720 — 30203 | | |
| 11 | 1 | Stck. | Passfeder groß | DIN 6885 — B 12 x 8 x 22 | | |
| 12 | 1 | Stck. | Passfeder klein | DIN 6885 — B 5 x 5 x 10 | | |
| 13 | 2 | Stck. | Verschlussschraube | DIN 908 — M14 x 1,5 — St | | |
| 14 | 2 | Stck. | Dichtring | DIN 7603 — A 14 x 18 Vf | | |
| 15 | 15 | Stck. | Zylinderschraube mit Innensechskant | ISO 4762 — M6 x 20 — 8.8 | | |
| 16 | 1 | Stck. | Zylinderschraube mit Innensechskant | ISO 4762 — M6 x 16 — 8.8 | | |
| 17 | 1 | Stck. | Scheibe | DIN 9021 — B 6,4 | | |
| 18 | 2 | Stck. | Radial-Wellendichtring | DIN 3760 — AS 45 x 60 x 8 | | |
| 19 | 2 | Stck. | O-Ring | DIN 3771-85x3,55-N-NBR 70 | | |
| 20 | 1 | Stck. | O-Ring | DIN 3771-40x3,55-N-NBR 70 | | |
| 21 | 4 | Stck. | Stiftschraube | Kaufteil gemäß Zeichnung | | S235JR |

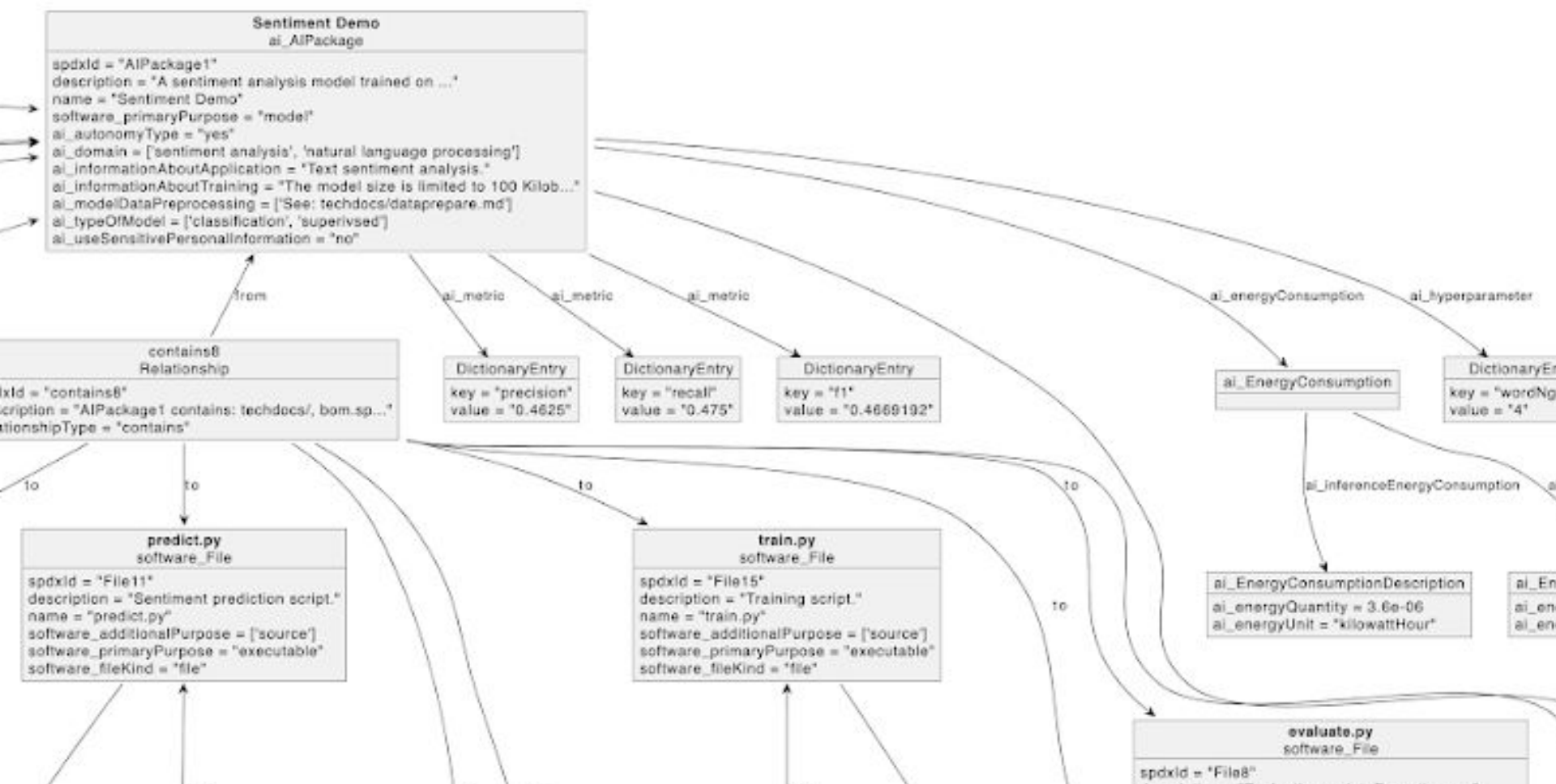| | Datum | Name | | |
|---|---|---|---|---|
| Bearb. | 20.03.2006 | | | |
| Gepr. | | | Schneckengetriebe | |
| Norm | | | | |
| | | | | Blatt 1 |
| | | | | 10 Bl. |
| Zust. | Änderung | Datum | Name | |

# Software Bill of Materials

- "formal record containing the details and supply chain relationships of various components used in building software" – Executive Order on Improving the Nation's Cybersecurity (EO 14028)

- "analogous to a list of ingredients" "can help organisations or persons avoid consumption of software that could harm them." – Wikipedia

- "communicating a release: name, version, components, licenses, copyrights, and useful security references." – SPDX

- ISO/IEC 5962:2021 Software Package Data Exchange (SPDX) Specification V2.2.1

profile AI

| AIPackage |
|---|
| + Core/Artifact/suppliedBy: Agent[1] |
| + Software/Package/downloadLocation: anyURI[1] |
| + Software/Package/packageVersion: String[1] |
| + Software/SoftwareArtifact/primaryPurpose: SoftwarePurpose[1] |
| + Core/Artifact/releaseTime: DateTime[1] |
| + energyConsumption: String[0..1] |
| + standardCompliance: String[0..*] |
| + limitation: String[0..1] |
| + typeOfModel: String[0..*] |
| + informationAboutTraining: String[0..1] |
| + informationAboutApplication: String[0..1] |
| + hyperparameter: DictionaryEntry[0..*] |
| + modelDataPreprocessing: String[0..*] |
| + modelExplainability: String[0..*] |
| + sensitivePersonalInformation:PresenceType[0..1] |
| + metricDecisionThreshold: DictionaryEntry[0..*] |
| + metric: DictionaryEntry[0..*] |
| + domain: String[0..*] |
| + autonomyType: PresenceType[0..1] |
| + safetyRiskAssessment: SafetyRiskAssessmentType[0..1] |

https://github.com/bact/sentimentdemo

25 June 2024 06:00

22

**Sentiment Demo**
ai_AIPackage

spdxId = "AIPackage1"
description = "A sentiment analysis model trained on ..."
name = "Sentiment Demo"
software_primaryPurpose = "model"
ai_autonomyType = "yes"
ai_domain = ['sentiment analysis', 'natural language processing']
ai_informationAboutApplication = "Text sentiment analysis."
ai_informationAboutTraining = "The model size is limited to 100 Kilob..."
ai_modelDataPreprocessing = ['See: techdocs/dataprepare.md']
ai_typeOfModel = ['classification', 'superivsed']
ai_useSensitivePersonalInformation = "no"

from

**contains8**
Relationship

dxId = "contains8"
cription = "AIPackage1 contains: techdocs/, bom.sp..."
ationshipType = "contains"

ai_metric

**DictionaryEntry**
key = "precision"
value = "0.4625"

ai_metric

**DictionaryEntry**
key = "recall"
value = "0.475"

ai_metric

**DictionaryEntry**
key = "f1"
value = "0.4669192"

ai_energyConsumption

**ai_EnergyConsumption**

ai_hyperparameter

**DictionaryEn**
key = "wordNg
value = "4"

to

to

to

to

to

**predict.py**
software_File

spdxId = "File11"
description = "Sentiment prediction script."
name = "predict.py"
software_additionalPurpose = ['source']
software_primaryPurpose = "executable"
software_fileKind = "file"

**train.py**
software_File

spdxId = "File15"
description = "Training script."
name = "train.py"
software_additionalPurpose = ['source']
software_primaryPurpose = "executable"
software_fileKind = "file"

ai_inferenceEnergyConsumption

**ai_EnergyConsumptionDescription**
ai_energyQuantity = 3.6e-06
ai_energyUnit = "kilowattHour"

**ai_E**
ai_en
ai_en

**evaluate.py**
software_File

spdxId = "File8"

# Transparency - Content labelling

- TikTok and Meta labeling of AI-generated images is based on C2PA & IPTC content metadata standards
- C2PA (Adobe, BBC, Google, Microsoft, Sony, etc) built upon Content Authenticity Initiative & Project Origin
- IPTC Photo Metadata (International Press Telecommunications Council)

# Standards (organisational/operational)

- ISO/IEC 23894:2023 Information technology — Artificial intelligence — Guidance on risk management

- Artificial Intelligence Risk Management Framework (AI RMF 1.0) from National Institute of Standards and Technology

# Standards (evaluation framework)

- **ONDE AI Ethics (Thailand)**
  ai-ethics.onde.go.th

- **AI Verify (Singapore)**
  aiverifyfoundation.sg

- **Inspect AI (UK)**
  inspect.ai-safety-institute.org.uk

## CONFIGURATION (.ENV)

Model  Logging

### Model

OpenAI

gpt-4-0125-preview

| Connections | Retries | Timeout |
|---|---|---|
| 20 | default | default |

### TASKS

- benchmarks
  - arc.py
    - ⊙ arc_challenge
    - ⊙ arc_easy
  - gpqa.py
  - gsm8k.py
  - hellaswag.py
  - mathematics.py
  - mmlu.py
- examples
  - agents
    - langchain
      - wikipedia.py
      - biology_qa.py
      - popularity.py

### TASK

benchmarks > arc.py > arc_challenge

```python
14
15  from inspect_ai import Task, task
16  from inspect_ai.dataset import Sample, hf_dataset
17  from inspect_ai.scorer import answer
18  from inspect_ai.solver import multiple_choice
19
20  def arc_task(dataset_name):
21      return Task(
22          dataset=hf_dataset(
23              path="allenai/ai2_arc",
24              name=dataset_name,
25              split="test",
26              sample_fields=record_to_sample,
27              shuffle=True,
28          ),
29          plan=multiple_choice(),
30          scorer=answer("letter"),
31      )
32
33
34  @task
35  def arc_easy():
36      return arc_task("ARC-Easy")
37
38
39  @task
40  def arc_challenge():
41      return arc_task("ARC-Challenge")
42
43
```

Debug Task | ▷ Run Task

Debug Task | ▷ Run Task

Inspect View   2024-05-09T17-19-51_arc-challenge_beSzAz3bBgHuWk37r2...

**arc_challenge** openai/gpt-4-0125-preview

5/9/2024, 5:19:51 PM— 1 min 22 sec

| accuracy | bootstrap_std |
|---|---|
| **0.953** | **0.007** |

| DATASET | PLAN | SCORER |
|---|---|---|
| allenai/ai2_arc — 1000 samples | multiple_choice | answer |

Samples   Info   Logging   JSON

Scores: All

Sort: sample asc    ⊹ Open All

| | Input | Target | Answer | Score |
|---|---|---|---|---|
| 1 | An astronomer observes that a planet rotates faster after a meteorite... | C | C | C |
| 2 | A group of engineers wanted to know how different building... | B | B | C |
| 3 | The end result in the process of photosynthesis is the... | C | C | C |
| 4 | A physicist wants to determine the speed a car must reach to jump... | D | D | C |
| 5 | An astronaut drops a 1.0 kg object and a 5.0 kg object on the Moon.... | D | C | I |

# theory_of_mind openai/gpt-4

4/28/2024, 8:43:13 PM— 1 min 50 sec

accuracy **0.81**

bootstrap_std **0.040**

| DATASET | PLAN | SCORER |
|---------|------|--------|
| theory_of_mind — 100 samples | chain_of_thought → generate | model_graded_fact |

**Samples**  Info  Logging  JSON

Scores: All ▾  Sort: sample asc ▾  ⬍ Open All

| | Input | Target | Answer | Score | |
|---|-------|--------|--------|-------|---|
| 1 | Jackson entered the hall. Chloe entered the hall. The boots is in the bathtub. Jackson exited the hall. Jackson entered the dining_room. Chloe moved the boots to the pantry.... | bathtub | First, we see Jackson and Chloe entering the hall but no... | C | ▾ |
| 2 | Jackson entered the hall. Chloe entered the hall. The boots is in the bathtub. Jackson exited the hall. Jackson entered the dining_room. Chloe moved the boots to the pantry.... | pantry | 1. The narrative starts with Jackson and Chloe both entering... | C | ▾ |
| 3 | Jackson entered the hall. Chloe entered the hall. The boots is in the bathtub. Jackson exited the hall. Jackson entered the dining_room. Chloe moved the boots to the pantry.... | bathtub | Jackson initially entered the hall. Chloe also entered... | I | ▾ |
| 4 | Jackson entered the hall. Chloe entered the hall. The boots is in the bathtub. Jackson exited the hall. Jackson entered the dining_room. Chloe moved the boots to the pantry.... | pantry | First, Jackson and Chloe entered the hall. This provides... | C | ▾ |
| 5 | Jackson entered the hall. Chloe entered the hall. The boots is in the bathtub. Jackson exited the hall. Jackson entered the dining_room. Chloe moved the boots to the pantry.... | bathtub | Firstly, Jackson was present in the hall when Chloe entered... | C | ▾ |
| 6 | Jackson entered the hall. Chloe entered the hall. The boots is in the bathtub. Jackson exited the hall. Jackson entered the dining_room. Chloe moved the boots to the pantry.... | bathtub | First, Jackson entered the hall, then Chloe also did so. At this... | C | ▾ |
| 7 | Hannah entered the patio. Noah entered the patio. The sweater is in the bucket. Noah exited the patio. Ethan entered the study. Ethan exited the study. Hannah moved... | bucket | Step 1: Hannah entered the patio Step 2: Noah entere... | C | ▾ |

28

# Standards (incident report)

- **AIAAIC harm taxonomy**



**AI, algorithmic, and automation harms taxonomy**

**Autonomy**
- Autonomy/agency loss
- Impersonation/identity theft
- IP/copyright loss
- Personality rights loss

**Physical**
- Bodily injury
- Loss of life
- Personal health deterioration
- Property damage

**Psychological**
- Addiction
- Alienation/isolation
- Anxiety/distress
- Coercion/manipulation
- Dehumanisation/objectification
- Harassment/abuse/intimidation
- Over-reliance
- Radicalisation
- Self-harm
- Sexualisation
- Trauma

**Reputational**
- Defamation/libel/slander
- Loss of confidence/trust

**Business & financial**
- Business operations/infrastructure damage
- Confidentiality loss
- Financial/earnings loss
- Livelihood loss
- Monopolisation
- Opportunity loss

**Human rights & civil liberties**
- Benefits/entitlements loss
- Dignity loss
- Discrimination
- Loss of freedom of speech/expression
- Loss of freedom of assembly/association
- Loss of social rights and access to public services
- Loss of right to information
- Loss of right to free elections
- Loss of right to liberty and security
- Loss of right to due process
- Privacy loss

**Societal & cultural**
- Breach of ethics/values/norms
- Cheating/plagiarism
- Chilling effect
- Cultural dispossession
- Damage to public health
- Historical revisionism
- Information degradation
- Job loss/losses
- Labour exploitation
- Loss of creativity/critical thinking
- Stereotyping
- Public service delivery deterioration
- Societal destabilisation
- Societal inequality
- Violence/armed conflict

**Political & economic**
- Critical infrastructure damage
- Economic instability
- Power concentration
- Electoral interference
- Institutional trust loss
- Political instability
- Political manipulation

**Environmental**
- Biodiversity loss
- Carbon emissions
- Electronic waste
- Excessive energy consumption
- Excessive landfill
- Excessive water consumption
- Natural resources extraction
- Pollution

29

# Purposes of Public Accountability
**(adapted from Bovens et al. 2010)**

- **Democratic perspective**
    - Popular control

- **Constitutional perspective**
    - Prevention of corruption
    and abuse of power

- **Learning perspective**
    - Maximising public value

# Purposes of Public Accountability

**(adapted from Bovens et al. 2010)**

*Related measures*

- **Democratic perspective**
  - Popular control

  ☐ *Explainability (legitimacy) + Human oversight (lawful + ethical)*

- **Constitutional perspective**
  - Prevention of corruption and abuse of power

  ☐ *Bias and drift detection (technically robust + ethical)*

- **Learning perspective**
  - Maximising public value

  ☐ *Information that allow the improvement     of the system (technically robust,    organizational learning)*

# Taxonomies for AI Accountability

Society/ Sector A

Society/ Sector B

**HOW**

Realization

**WHY**

Realization

**HOW**

Social & Legal Layer

EU AI Act

Charter of Fundamental Rights of EU

Local Norms

UNESCO's Recommendation on the Ethics of AI

Ethical

AI Principles and Ethical Guidelines

Sustainable Development Goals

Technical Layer

Technical Standards

NIST AI Risk Management Framework 1.0

AI Life Cycle Processes (ISO/IEC 5338:2023)

Data Privacy Vocabulary

Data quality for analytics and ML (ISO/IEC FDIS 5259)

Software

AI Verify toolkit

Logging

Inspect AI

AI Systems

Adaptable and exchangeable across different jurisdictions

Technology to support the realization of AI principles

Standard taxonomy to serve three accountability purposes:

- Democratic _Technical documentation_ for informed popular control

- Constitutional _Record keeping_ to minimize corruption or abuse of power

- Learning _Incident reporting_ to maximize public value and safety

Gasser, U., Almeida, V.A.F., 2017. A Layered Model for AI Governance. IEEE Internet Computing. https://doi.org/10.1109/MIC.2017.4180835

Bovens, M., Schillemans, T., Goodin, R.E., 2014. Public Accountability, in: The Oxford Handbook of Public Accountability. Oxford University Press, Oxford, New York, pp. 1–20. https://doi.org/10.1093/oxfordhb/9780199641253.013.0012

# Thank you

อาทิตย์ สุริยะวงศ์กุล
Arthit Suriyawongkul
suriyawa@tcd.ie

# ความโปร่งใสของ  AI 3 ระดับ

ความโปร่งใส 3 ระดับนี้ ทำงานร่วมกัน และส่งผลต่อภาระความรับผิดและการควบคุมโดยมนุษย์

## ความโปร่งใสเชิงอัลกอริทึม

- ความสามารถในการเข้าถึงและตรวจสอบ-ตั้งคำถามต่อโค้ด, ชุดข้อมูล, และระบบที่ประกอบเข้าด้วยกัน

- ค่าความน่าจะเป็น แผนภูมิ หรือสิ่งที่ได้จากวิธีในการอธิบาย AI (เช่น LIME* และ SHAP**) อาจถูกอ่านเข้าใจได้โดยผู้เชี่ยวชาญเฉพาะเรื่อง ผู้ตรวจสอบ และผู้กำกับกิจการ **แต่อาจเป็นการลำบากสำหรับผู้ที่ไม่มีภูมิหลังทาง AI หรือความรู้ในกิจการดังกล่าว**

## ความโปร่งใสเชิงปฏิสัมพันธ์

- ความสามารถในการเข้าใจสิ่งที่ระบบ AI ทำได้ดีและสิ่งที่ทำได้จำกัด ซึ่งได้มาจากการแลกเปลี่ยนความรู้ระหว่างตัวระบบ AI และผู้ใช้

- อุปลักษณ์ (metaphor) ที่จับต้องได้-ฝังอยู่ในประสบการณ์การใช้งาน เป็นอุปลักษณ์ที่สามารถทำให้เข้าใจสภาพแวดล้อมและวิธีคิดของการออกแบบระบบ **ความรู้หรือคำอธิบายนี้ เป็นสิ่งที่ระบบและผู้ใช้สร้างขึ้นร่วมกันในระหว่างที่มีปฏิสัมพันธ์กัน**

## ความโปร่งใสเชิงสังคม

- ความสามารถทางกฎหมายและทางวัฒนธรรม ของ(สถาบันทาง)สังคมในการเข้าใจและหาหนทางตอบสนองกับการใช้งานระบบ AI

- วิธีการ**ที่ไม่เสนอข้อมูลหรือ "ทางเลือก" ให้กับผู้ใช้จนเกินรับไหว** (เช่น กล่องข้อความขอความยินยอมเก็บคุกกี้) วิธีการควรถูกผนวกเข้าไปในการทำงานของสถาบัน (เช่น มาตรการความปลอดภัยในอุตสาหกรรมอาหาร การบิน)

เสี่ยง

เสี่ยง

เสี่ยง