# Explainable AI (XAI)

September 10th, 2024

*Regnskap Norge – OVERSKUDD 2024*

Claudia Battistin, PhD

AI Program Director, R&I – Telenor Group

# *Sometimes AI makes small mistakes...*



ChatGPT

Q: *Who is Claudia Battistin?*

A: *Claudia Battistin is an Italian academic and professional known for her expertise in international relations, with a particular focus on conflict and post-conflict situations, humanitarian interventions, and the impact of war on societies. She has been involved in various roles related to international cooperation, peacebuilding, and the protection of human rights in conflict zones.*

*Battistin has also contributed to the academic field through research and publications, particularly on issues related to conflict resolution, peace processes, and international law. Her work often intersects with the efforts of international organizations, NGOs, and other institutions dedicated to addressing the complex challenges of war and peace.*
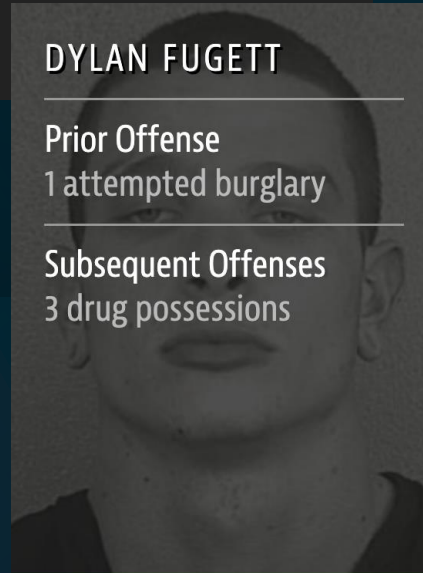
*… sometimes makes big mistakes*

# *Why does AI fail?*



**Training data**
- Poor quality (standard, size, ...)
- Skewed
- Not representative of production data

**AI design**
- Poor feature selection
- Spurious correlations

**User prompt**
- Limited instructions for use
- Absence of guardrails

# Explainable AI – trust enabler



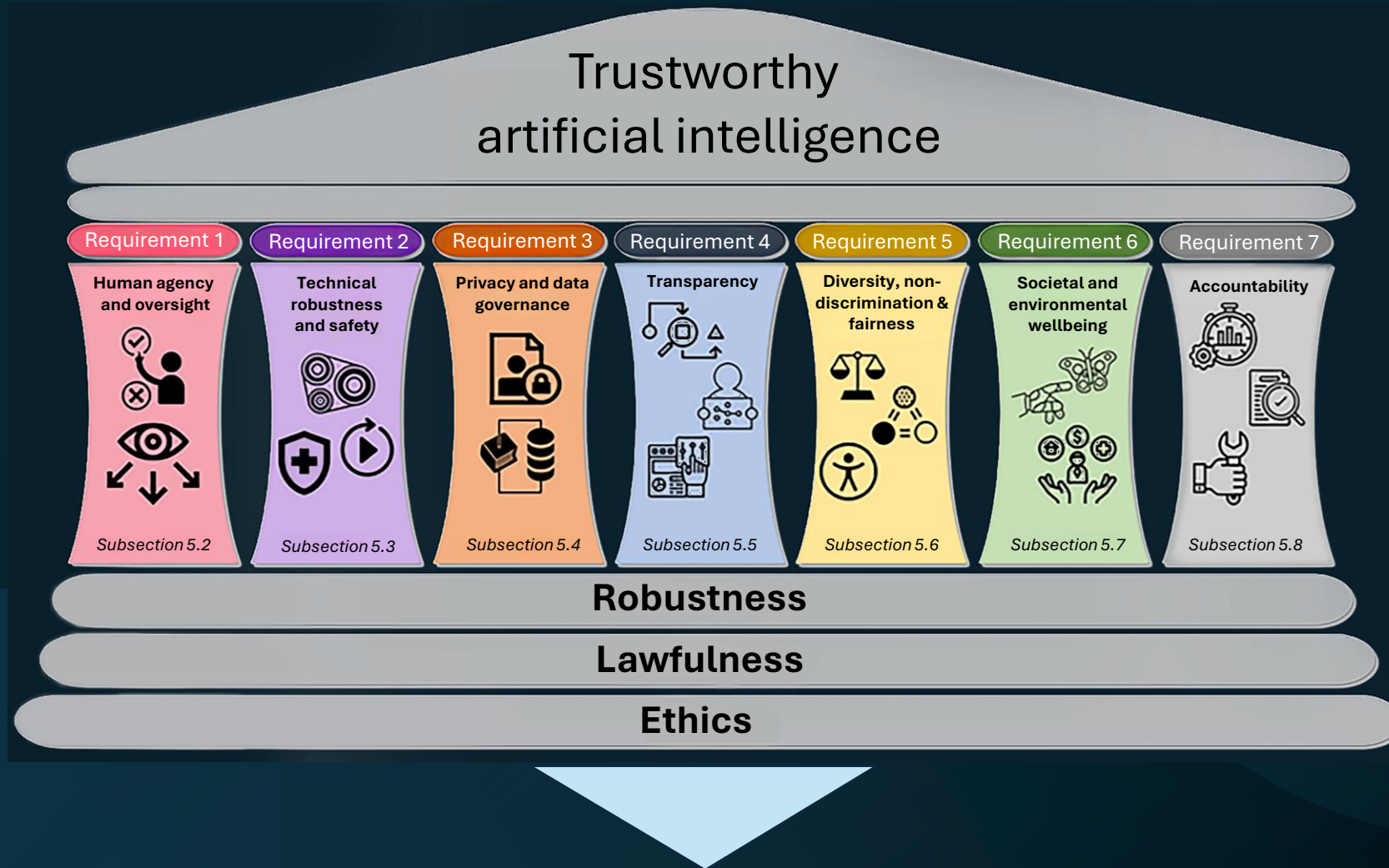Trustworthy artificial intelligence

| Requirement 1 | Requirement 2 | Requirement 3 | Requirement 4 | Requirement 5 | Requirement 6 | Requirement 7 |
|---|---|---|---|---|---|---|
| Human agency and oversight | Technical robustness and safety | Privacy and data governance | Transparency | Diversity, non-discrimination & fairness | Societal and environmental wellbeing | Accountability |
| Subsection 5.2 | Subsection 5.3 | Subsection 5.4 | Subsection 5.5 | Subsection 5.6 | Subsection 5.7 | Subsection 5.8 |

**Robustness**

**Lawfulness**

**Ethics**

Explaining model and predictions

# Explainable AI in the Nordics



Legend: ● 2021 ● 2022 ● 2023

| Category | 2021 | 2022 | 2023 |
|---|---|---|---|
| Computer vision | 16.76% | 40.34% | 42.86% |
| NLP | 18.5% | 42.02% | 60% |
| Recommendation | 12.72% | | |
| Intelligent search | 12.72% | 20.17% | 51.43% |
| Cognitive decision security | 12.14% | | |
| Optimization engines | 12.72% | 24.37% | 37.14% |
| Automated generator | 6.94% | | |
| Speech recognition | | 19.33% | 28.57% |
| Machine learning | | 74.79% | |
| Deep learning | | 51.26% | 42.86% |
| Explainable AI | | 54.62% | 17.14% |
| Time series forecastin | | 28.57% | 40% |
| Reinforcement learning | | 10.92% | 37.14% |
| Bayesian learning | | | 20% |
| Large language models | | | 71.43% |
| Generative AI | | | 77.14% |
| Synthethic data for AI training | | | 14.29% |
| Advanced analytics | | | 60% |
| None of the above | | 8.4% | 3% |
| Other | 7.51% | 7.56% | |
| Not sure | | 4.2% | 14.29% |

The Nordic State of AI, 2024 – Silo AI

AI vs XAI

AI vs XAI

Training data

AI model  +  explainability layer

Input data, prompt

Humans

Prediction, response

Politics

welfare

Mountain landscape

# Explainable AI approaches



Model agnostic

- Surrogate interpretable models
- Stats of local explanations

Feature attribution – Counterfactuals -

Local ←————————————————→ Global

Interpretable models
(linear, decision trees,...)

Model specific

*adapted from Zhang et al. 2023*

# Probability of material restatement



f(x)= -2.021

adapted from Zhang et al. 2023

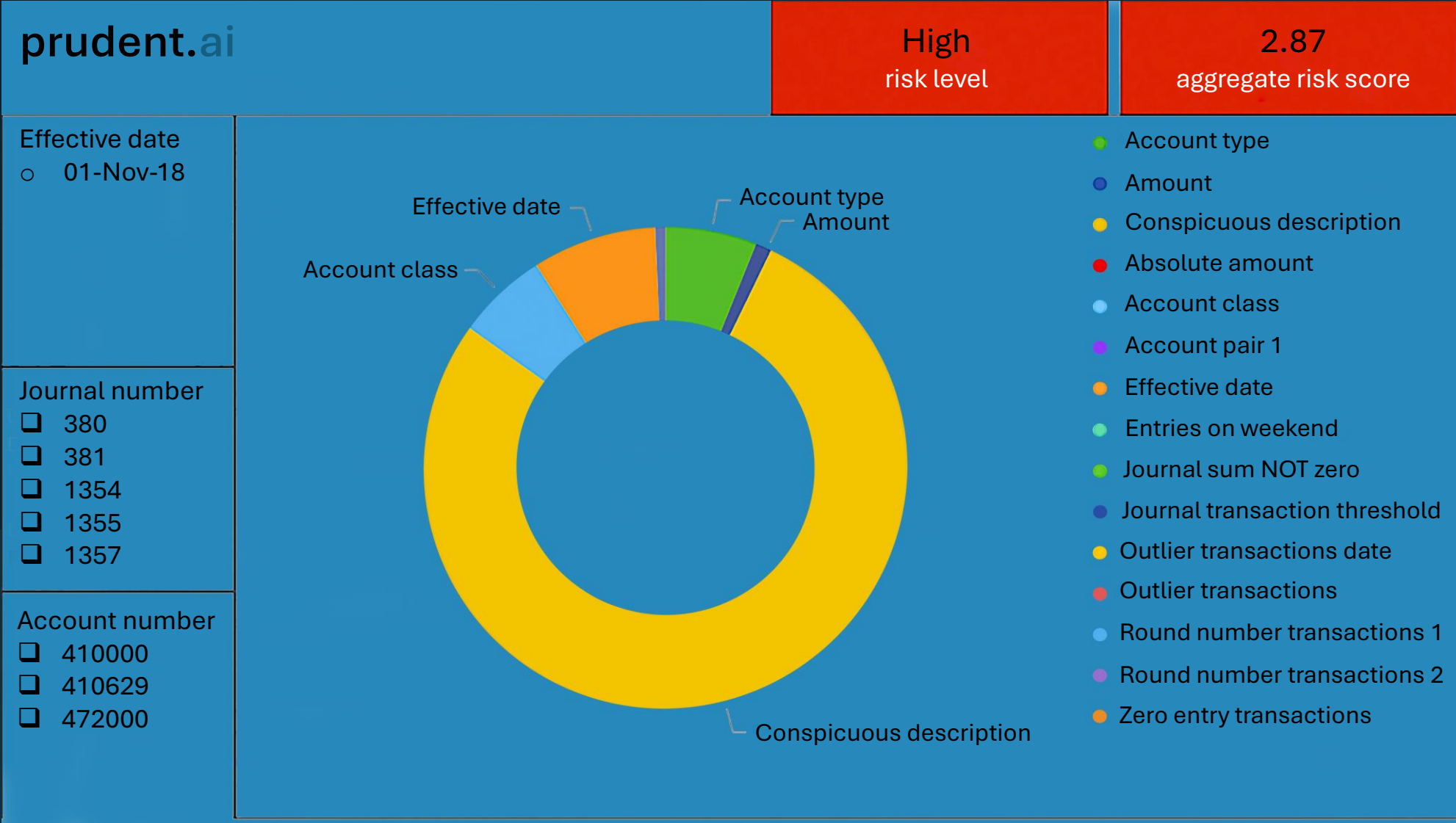0.01 = earnings to prices                                    -0.8

-0.203 = lag_Return                                         -0.41

0.002 = change in return of assets                          -0.27

0.324 = Soft Assets                                         +0.24

0.027 = change in cash margin                               +0.23

0.002 = performance matched discretionary accruals          -0.22

-0.179 = lag_mean adjusted absolute DD residual             -0.21

0.041 = deferred tax expense                                -0.21

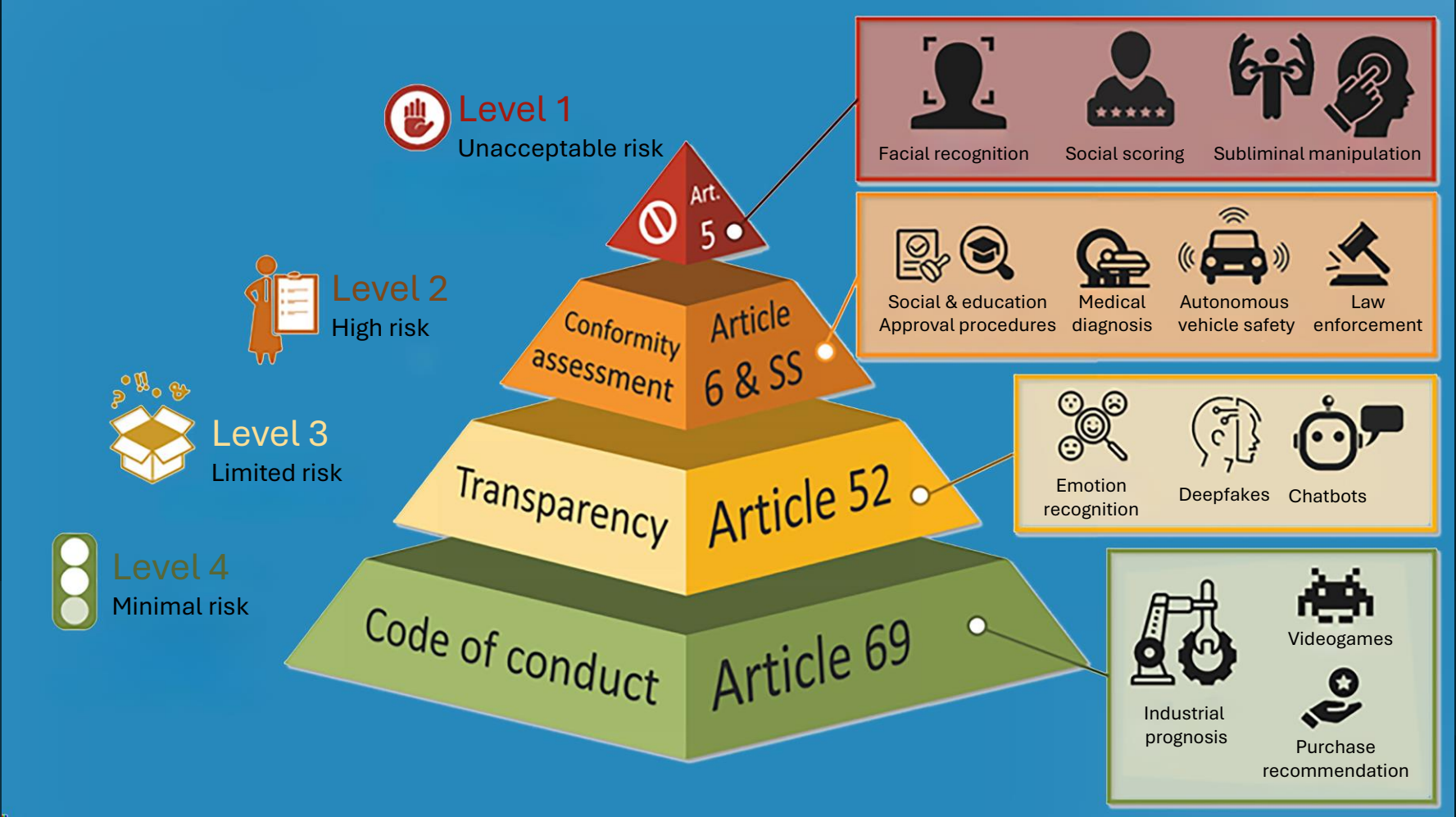-0.014 = change in cash sales                               -0.2
12 other features                                           -0.13

-2.0        -1.5        -1.0        -0.5

E[f(x)]= -2.021

# Risk of transaction



**prudent.**ai

**High**
risk level

**2.87**
aggregate risk score

**Effective date**
- 01-Nov-18

**Journal number**
- ☐ 380
- ☐ 381
- ☐ 1354
- ☐ 1355
- ☐ 1357

**Account number**
- ☐ 410000
- ☐ 410629
- ☐ 472000

Legend:
- Account type
- Amount
- Conspicuous description
- Absolute amount
- Account class
- Account pair 1
- Effective date
- Entries on weekend
- Journal sum NOT zero
- Journal transaction threshold
- Outlier transactions date
- Outlier transactions
- Round number transactions 1
- Round number transactions 2
- Zero entry transactions

Chart labels: Effective date, Account type, Amount, Account class, Conspicuous description

*adapted from Vaidyanathan, ACCA*

# Explainable AI – EU AI Act



adapted from Díaz-Rodríguez et al. 2023

# Explainable AI – GDPR

**Right to explanation**

**Recital 71 (not binding) –** LOCAL EXPLANATIONS

*The data subject should have the right [...] to obtain an explanation of the decision [solely on automated processing] reached after such assessment and to challenge the decision.*

**Article 13(2)(f) –** GLOBAL EXPLANATIONS

*In the event a controller relies on automated decision-making, [...] the controller is under the duty to provide the data subject with "meaningful information about the logic involved" [...] of such forms of processing.*

# Thanks