# Trustworthy AI – Obligation or Entrepreneurial Opportunity?

EAA e-Conference on
Data Science & Data Ethics

29 June 2021

*Dr. Maximilian Poretschkin*

***Fraunhofer Institute for Intelligent Analysis and Information Systems***

**Claims processing**

**Automated estimation of damage amount**

**Personalized insurance pricing**

**Assessment of the probability of default**

**Detection of fraudulent claims**

**Document creation**

**Responding to customer queries**

**Churn prediction**

Image sources: © welcomeinside - stock.adobe.com; © Anthony Brown - stock.adobe.com; © Cozine - stock.adobe.com; © samrit - stock.adobe.com; ©WrightStudio - stock.adobe.com; © Zerbor - stock.adobe.com; © REDPIXEL - stock.adobe.com; © Tom Bayer - stock.adobe.com

! New providers with expertise in data-driven business models can **occupy the customer interface**

! Insurtechs **unbundle the insurance value chain** by using data and AI use cases

! AI-based **loss prevention**, **differentiated pricing** and **new insurance models** (insurance of AI)

„3 crashes, 3 deaths raise questions about Tesla's autopilot"
AP News, January 3, 2020

„Apple Card Investigated After Gender Discrimination Complaints"
New York Times, November 10, 2019

„Amazon scraps secret AI recruiting tool that showed bias against women"
Reuters, October 11, 2018

„Alipay responds to risks of facial recognition payment"
China Daily, September 9, 2019

Image source: ©Zerbor - stock.adobe.com

# GUARDRAILS FOR ARTIFICIAL INTELLIGENCE

01

- Embedding BDAI within **a proper business organization**
- **No black box excuses –** explainability/traceability of models is necessary
- Continuing to develop existing **governance concepts**
- Defining **supervisory requirements for the explainability** and effectiveness of compliance processes
- Defining **prerequisites for BDAI use in models** requiring **supervisory approval**
- Addressing increased **information security risks** and using BDAI to combat them

Source: Big Data meets Artificial Intelligence, Challenges and implications for the supervision and regulation of financial services, BaFin, 2018

## Human agency and oversight

- Fundamental rights, human agency and human oversight

## Technical robustness and safety

- Resilience to attack and security, fall back plan and general safety, accuracy, reliability and reproducibility

## Privacy and data governance

- Respect for privacy, quality and integrity of data, and access to data

## Transparency

- Traceability, explainability and communication

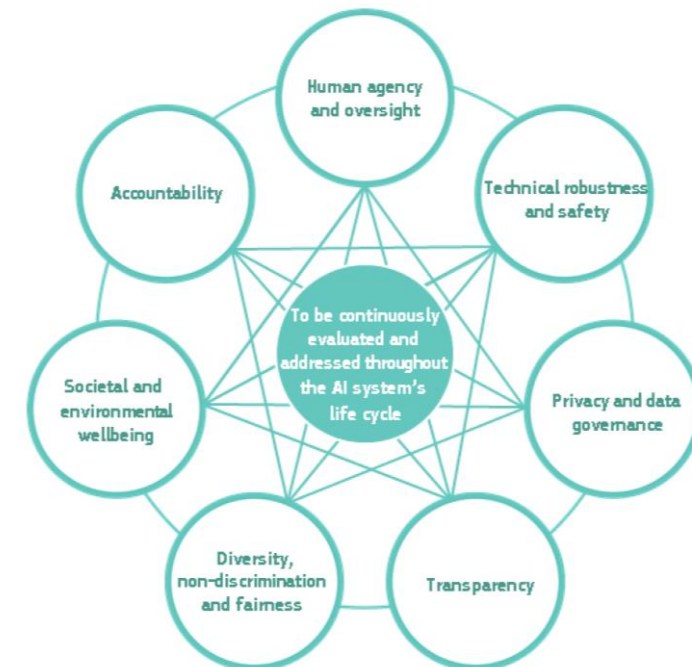## Diversity, non-discrimination and fairness

- Avoidance of unfair bias, accessibility and universal design, and stakeholder participation

## Societal and environmental wellbeing

- Sustainability and environmental friendliness, social impact, society and democracy

## Accountability

- Auditability, minimisation and reporting of negative impact, trade-offs and redress.

"The list of requirements is non-exhaustive."

Source: ETHICS GUIDELINES FOR TRUSTWORTHY AI, High-Level Expert Group on Artificial Intelligence, 2019

## *WORKING GROUP QUALITY, CONFORMITY ASSESSMENT AND CERTIFICATION*

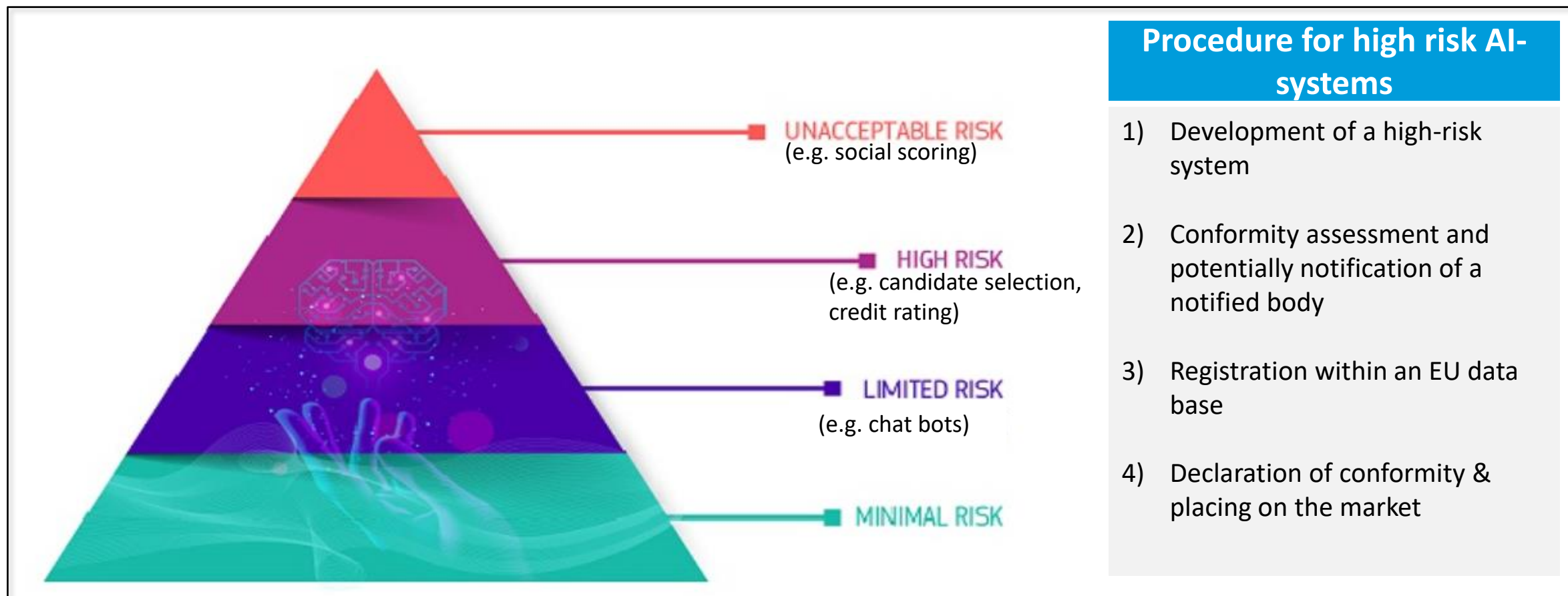**Two deliverables for establishing a testing procedure**:

1) Testing framework that guarantees comparability of tests (and is compatible with existing IT testing procedures!).

   - Process testing (standards for the development and operation of AI systems)

   - Product testing (verification of assured properties)

   - Differentiated assurance levels / testing depths

2) Criteria frameworks that operationalize trustworthiness requirements and map AI-specific challenges.

   - Use case dependency in formulation is challenge (metrics, thresholds)

   - Completely new testing tools and methods

Source: GERMAN STANDARDIZATION ROADMAP ON ARTIFICIAL INTELLIGENCE, DIN & DKE, 2019

DEUTSCHE NORMUNGSROADMAP
**KÜNSTLICHE INTELLIGENZ**

# MANY INSURANCE USE CASES ARE AFFECTED



**UNACCEPTABLE RISK**
(e.g. social scoring)

**HIGH RISK**
(e.g. candidate selection, credit rating)

**LIMITED RISK**
(e.g. chat bots)

**MINIMAL RISK**

## Procedure for high risk AI-systems

1) Development of a high-risk system

2) Conformity assessment and potentially notification of a notified body

3) Registration within an EU data base

4) Declaration of conformity & placing on the market

➡ **Procedure to systematically evaluate risks of AI-Systems is crucial!**

Image source: https://ec.europa.eu/info/strategy/priorities-2019-2024/europe-fit-digital-age/excellence-trust-artificial-intelligence

# AN UNFAIR RACE? WE HAVE TO FIND A MIDDLE GROUND



**European AI company**

Data access

Privacy

Regulation

Taxes

**American & Chinese hyperscaler**
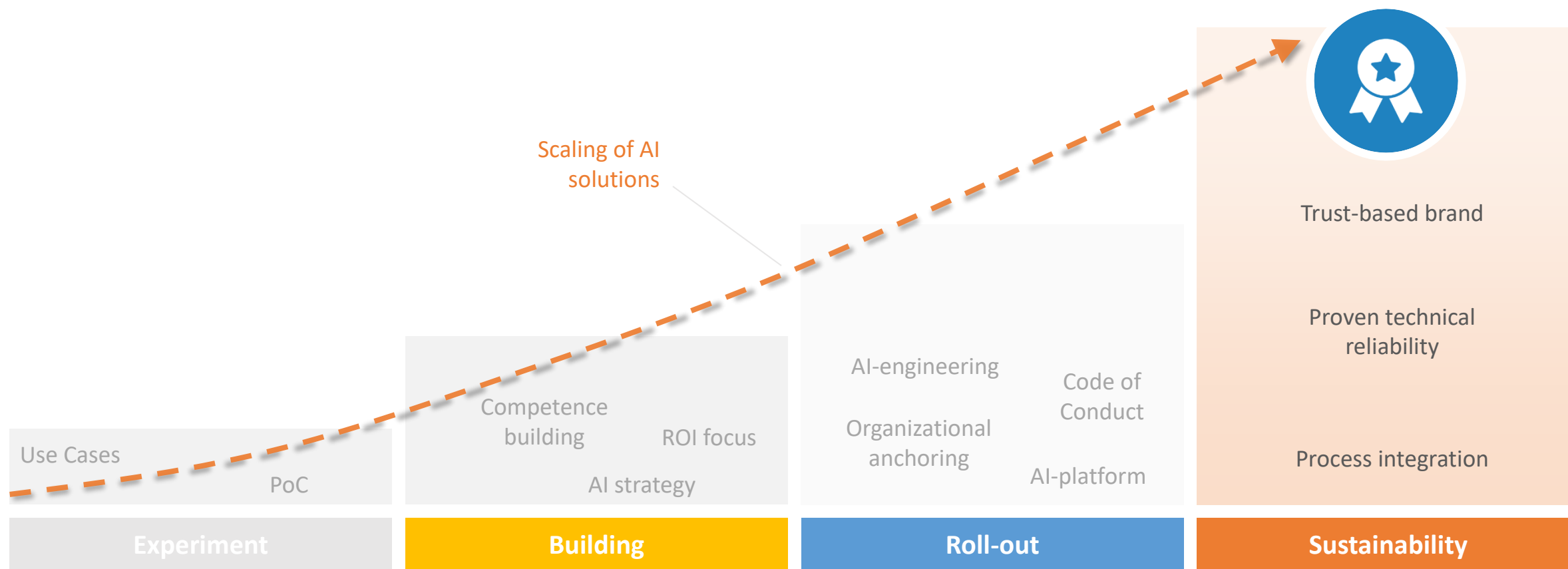
Image source: Adopted from Prof. Dr. Dr. Wolfgang Wahlster, Plattform Lernende Systeme

# ENTRE-PRENEURIAL OPPORTUNITY

02

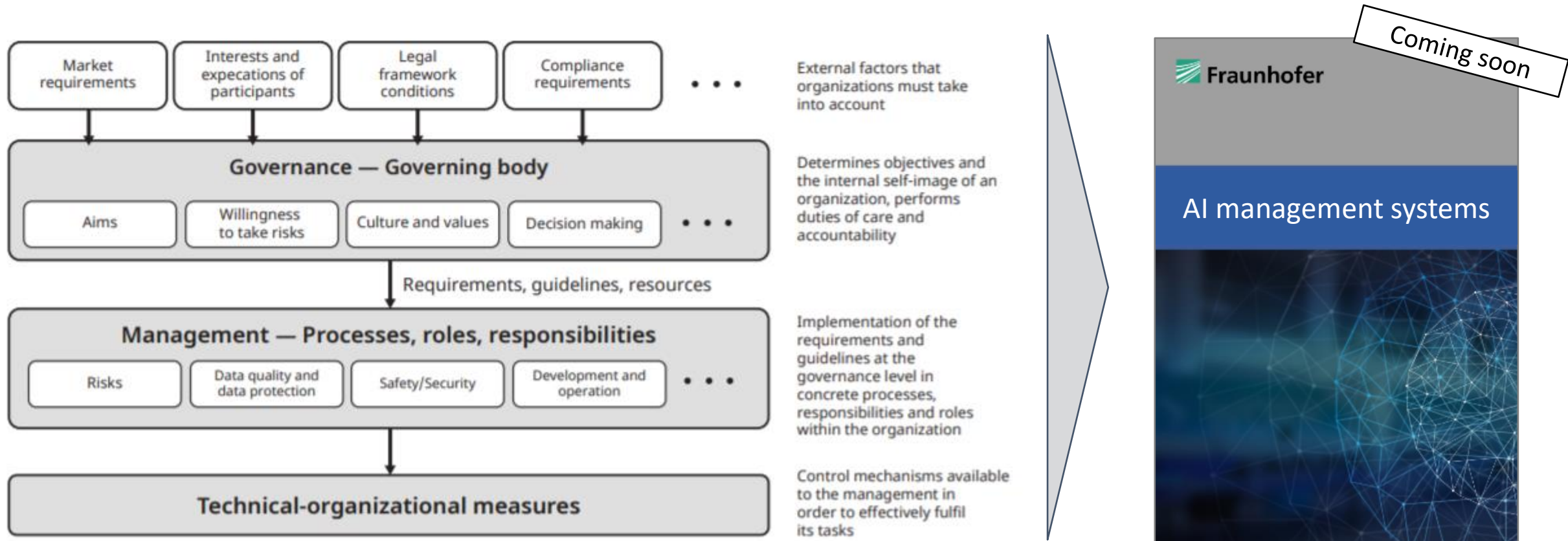## FROM AI EXPERIMENT TO SUSTAINABLY SCALABLE AI SOLUTION



Scaling of AI solutions

Trust-based brand

Proven technical reliability

Process integration

AI-engineering

Code of Conduct

Organizational anchoring

AI-platform

Competence building

ROI focus

AI strategy

Use Cases

PoC

| Experiment | Building | Roll-out | Sustainability |
|------------|----------|----------|----------------|

# GOVERNANCE, MANAGEMENT & TECHNICAL-ORGANIZATIONAL MEASURES



Coming soon

Fraunhofer

AI management systems

Image source: GERMAN STANDARDIZATION ROADMAP ON ARTIFICIAL INTELLIGENCE, DIN & DKE, 2019

## Building internal trust

**Business-critical decisions**

Are the AI-based recommendations comprehensible and trustworthy?
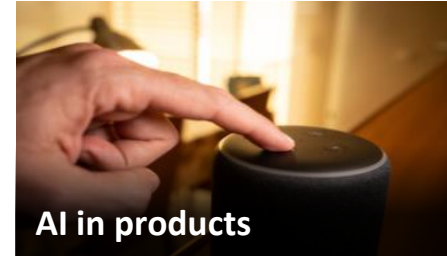
**AI in sensitive areas**

Can malfunctions cause significant (personal and/or financial) damage?

**Global deployment of AI systems**

Is the AI system reliable enough to be rolled out globally?

## Building external trust

**AI in products**

Can a competitive advantage be generated through proven technical reliability?

**Product brand**

How can a trusted brand be maintained for products with AI components?

## Understanding risks

**Acquisition of external AI solutions**

Does the purchased AI solution meet the required characteristics?

**Technical Due Diligence**

Does a company takeover entail technical risks? Does an acquired AI solution meet the expected requirements?
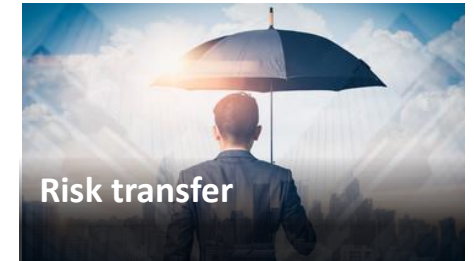
**Quality- and risk management**

Are AI risks recorded and assessed transparently? Are internal AI guidelines implemented?

## Risk transfer

**Risk premium**

Can proof of technical reliability reduce the insurance premium?

**Risk transfer**

Can the residual risk be covered by AI insurance?

Image sources: fizkes - stock.adobe.com, Nataliya Hora - stock.adobe.com, Fraunhofer IAIS, tanaonte - stock.adobe.com, Jacob Lund - stock.adobe.com, Looker_Studio - stock.adobe.com, amnaj - stock.adobe.com, Song_about_summer - stock.adobe.com, Zerbor - stock.adobe.com, TimeStopper - stock.adobe.com
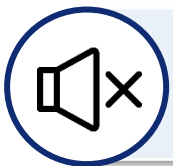
## CLASSICAL APPROACHES ARE NOT TRANSFERABLE

Data from past damage caused by AI applications is scarce

Risks are AI application-specific and can only be generalized to a limited extent

AI risks and damage scenarios are continuously changing. Approaches and competencies must be continuously updated

"Silent AI risks" in existing policies must be made transparent

Icon source: flaticon

# SYSTEMATIC EVALUATION OF AI-RISKS

03

**Ethics & Law**

Key questions concerning ethical issues

**Fairness**

Historically unbalanced data

**Autonomy & Control**

Appropriate degree of autonomy

**Transparency**

Incomprehensibility of results from neural networks

**Reliability**

Robustness of results processed by AI-systems

**Safety & Security**

Safety risks due to probabilistic output from AI component

**Privacy**

New types of personal data through AI

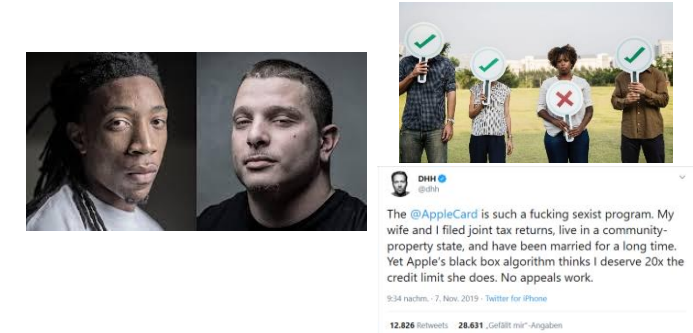| | | |
|---|---|---|
| **Design** | | ▪ The conception and architecture of the AI-system which ensures that certain characteristics are fulfilled „by design", like Privacy-by-Design, Safety-by-Design and Verifiability-by-Design. |
| **Develop-ment** | **Data** | ▪ The selection, augmentation, pre-processing of the training-, test- and input data of the AI–system as a key pre-requisite for a high quality of the AI-system. |
| | **AI-Com-ponent** | ▪ The selection of a method/algorithm, the training and test/validation of the model, aspects of transparency and explainability. The implementation into (standard) software. |
| | **Embed-ding** | ▪ The embedding of the AI-component into the AI-system with a focus on those aspects of the AI-system whose behaviour is based on the AI component. |
| **Operation** | | ▪ Application-related testing and assurance of the model's quality during operations. Verifiability and logging of the behaviour which is based on the AI-component. |

Icon source: https://iconmonstr.com/construction-35-png; https://iconmonstr.com/folder-20-png/

- Case 1: There is a **commonly preferred label**
  - We don't want to be refused that label due to
    - gender, ethnicity, …
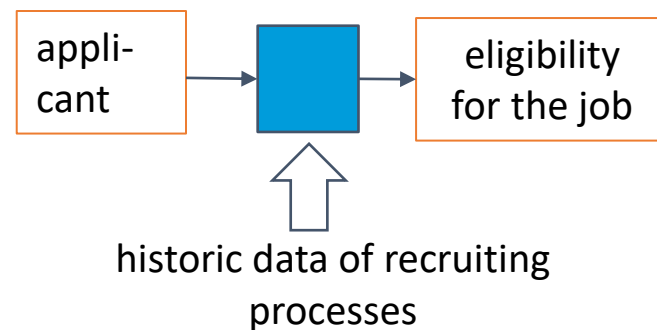  - Such attributes should not play a role in the decision process

- Case 2: There is **no commonly preferred label**
  - But we care whether we are **assigned the correct label**
  - We don't want to be treated with less care / worse service quality due to
    - gender, ethnicity, …
  - Such attributes should not influence the model performance



The @AppleCard is such a fucking sexist program. My wife and I filed joint tax returns, live in a community-property state, and have been married for a long time. Yet Apple's black box algorithm thinks I deserve 20x the credit limit she does. No appeals work.

Gender was misidentified in **up to 12 percent of darker-skinned males** in a set of 318 photos.

Gender was misidentified in **35 percent of darker-skinned females** in a set of 271 photos.

Icon sources: https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing; https://www.pexels.com/de-de/foto/afroamerikaner-betrubt-draussen-farbige-frau-1656594; Joy Buolamwini, M.I.T. Media Lab

## DATA IS A MAJOR CAUSE OF DISCRIMINATION BY AI-APPLICATIONS

Example: Recruiting Tool
**Historic bias**: Men were systematically preferred → ML-model could **learn negative correlation** of *woman* and *eligibility*

appli-cant → ▢ → eligibility for the job

historic data of recruiting processes

Example: Age Predictor
Certain group **underrepresented** in training data (or only **incomplete/ inaccurate data** available)
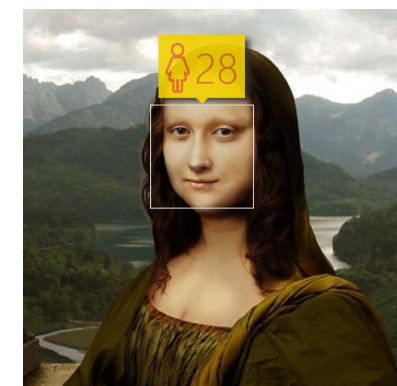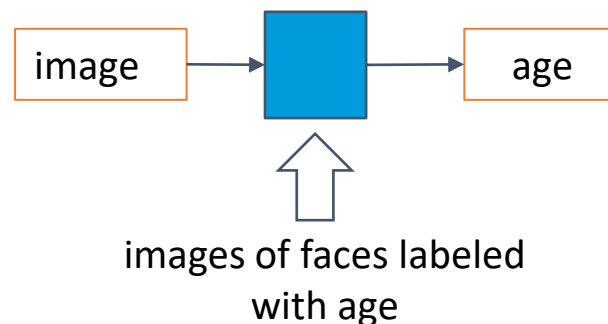→ **higher error rate** wrt this group

image → ▢ → age

images of faces labeled with age

## CHOICE OF FAIRNESS METRIC IS CRUCIAL FOR SAFEGUARDING

Risk analysis

What are the **risks** regarding fairness?

Which **concept of fairness** is appropriate in the given context?

To what extent is there a **trade-off** between fairness and utility of the application?

Choice of fairness metric **determines the procedure** of safeguarding and **enables objective evaluation** of measures

Icon sources: Flaticon.com icons by srip: Scale blue, risk blue, collaboration blue

## WIDE RANGE OF METRICS / FAIRNESS CONCEPTS

| Concepts of group fairness | | Other concepts |
| --- | --- | --- |
| Statistical/Demographic Parity | Overall Accuracy Equality | Individual Fairness |
| Predictive Rate Parity | Treatment Equality | Causal Discrimination |
| Equalized Odds | Well-Calibration | Counterfactual Fairness |
| Equal Opportunity | Test-fairness | ... |

# MITIGATE UNFAIRNESS BY MODIFYING DATASETS FOR TRAINING

## Pre-Processing

- Massaging
- Uniform Sampling
- Preferential Sampling
- Reweighing

- Unawareness
- Disparate Impact Remover
- Optimized Preprocessing
- Learning Fair Representations

## In-Processing

- Adversarial Debiasing
- Classifier without Disparate Mistreatment
- Prejudice Remover Regularizer

## Post-Processing

- Equalized Odds Threshold Predictor
- Reject Option Classification
- Calibrated Equalized Odds Postprocessing

## ASSESSMENT CATALOGUE PROVIDES GUIDANCE

| Dimension | Risk area |
|---|---|
| **Fairness** | Fairness |
| | Control of dynamics |
| **Autonomy and Control** | Distribution of tasks between human and AI-system |
| | Information and empowerment of users and stakeholders |
| | Control of dynamics |
| **Transparency** | Explainability to users |
| | Interpretability for experts |
| | Auditability |
| | Control of dynamics |

| Dimension | Risk area |
|---|---|
| **Privacy** | Protection of personal data |
| | Protection of business-relevant information |
| | Control of dynamics |
| **Reliability** | Reliability during regular operation |
| | Robustness |
| | Evasion strategies |
| | Estimation of uncertainty |
| | Control of dynamics |
| **Safety and Security** | Functional safety |
| | Integrity and confidentiality |
| | Availability |
| | Control of dynamics |

**Assessment catalogue is available here: www.iais.fraunhofer.de/de/forschung/kuenstliche-intelligenz/ki-pruefkatalog**

With regard to a specific operational context ⟹ Assessment of the trustworthiness of the AI function

Dimension 1 · Dimension 2 · … · Dimension 6

Protection needs analysis · … · …

low / Medium/high

Dimension 1 is not relevant

Risk area 1.1 · … · Risk area 1.n

The risk areas represent the possible malfunctions /unintended behaviour within a dimension.

Risk analysis

Definition of targets

Definition of criterion 1.1.1 · … · Definition of criterion 1.1.m

Set criteria that make the targets measurable

With regard to a specific operational context ⟹ **AI-function is trustworthy**

Cross-dimensional assessment

**AI-function is safeguarded w.r.t. dimension1** | …Dimension 2 | …Dimension 6

…

Summarizing dimension 1

*Dimension 1 is not relevant*

**AI-function is safeguarded w.r.t. risk area 1.1** … **…risk area 1.n**

The risk areas represent the possible malfunctions /unintended behaviour within a dimension.

Overall evaluation

…

**Criterion 1.1.1** … Criterion 1.1.m ⟸ The criteria can be used to verify the control of the identified risks.

Measures

Data | AI-com-ponent | Em-bedding | Operation ⟸ Documentation of measures related to these four categories that, taken together, justify meeting the criteria. Measures range from design decisions to testing.

## JOINT PROGRAMME TO DEVELOP TESTING METHODS FOR AI-SYSTEMS

**ZERTIFIZIERTE KI**

## Testing principles

- Testing scope
- Criteria
- Depth of testing
- Requirements for testing tools
- Concept for infrastructure

## Requirement Assessment

- Customer analysis
- Impact analysis
- Development of business models

## Use Cases

- Image recognition
- Natural Language Understanding
- Natural language processing
- Informed Machine Learning

## Testing ecosystem

- Platform for testing tools
- Testing labs
- Development of safeguarding methods

## Societal dialogue

- Considerations of ethical, legal and philosophical topics
- Public events

**Broad-based Participation Process**

**Key Partners:**

Fraunhofer IAIS

Bundesamt für Sicherheit in der Informationstechnik

UNIVERSITÄT BONN

Universität zu Köln

RWTH AACHEN UNIVERSITY

DIN

# AUTOMATED AI QUALITY ASSESSMENT

04

Quality estimation despite **restricted test coverage**

**Quantitative minimal requirements are highly use case specific**

Manual tests are **time consuming** or **uncomplete**

**Dynamics** of AI-systems
(continuous learning -> continuous assessment)

**Comprehensibility or explainability** for human auditors/ assessors

## Automation ⟷ Auditability

Icon source: https://iconmonstr.com/text-25-png; https://iconmonstr.com/school-7-png; https://iconmonstr.com/networking-7-png; https://iconmonstr.com/gear-11-png;
https://iconmonstr.com/search-thin-png

## OPEN-SOURCE PACKAGES FOR METRICS, ALGORITHMS AND LONG-TERM SIMULATION



AI Fairness 360 - Demo

Data — Check — Mitigate — Compare

### 2. Check bias metrics

Dataset: German credit scoring
Mitigation: none

**Protected Attribute: Sex**

Privileged Group: *Male*, Unprivileged Group: *Female*
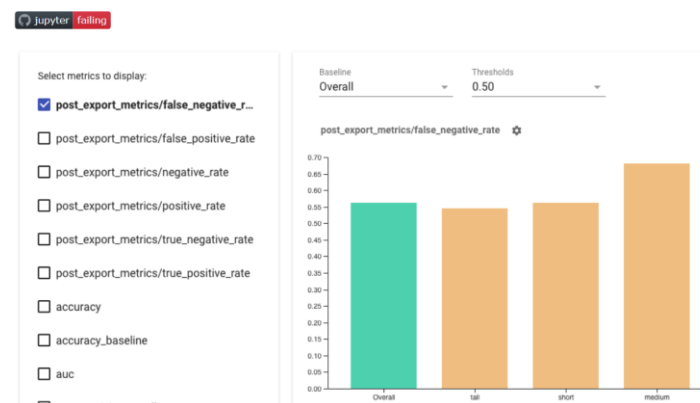Accuracy with no mitigation applied is 76%
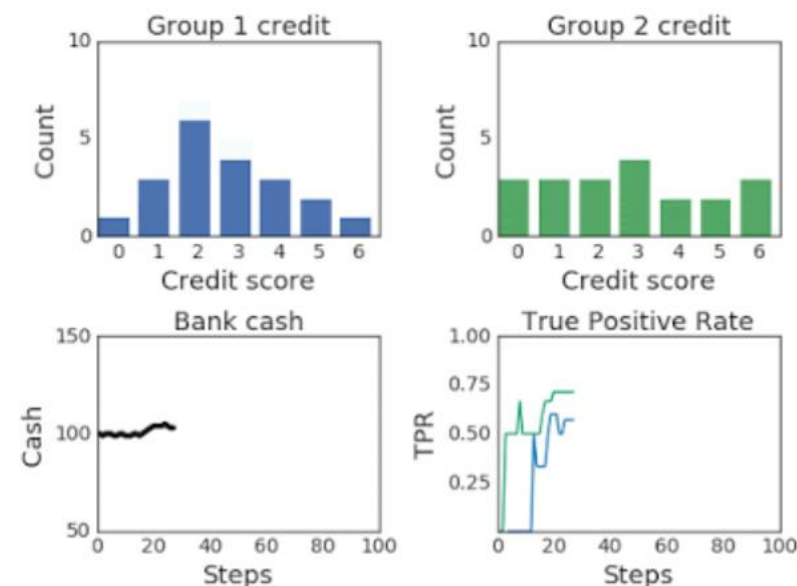With default thresholds, bias against unprivileged group

**Statistical Parity Difference**

**Equal Opportu Difference**



Aequitas
Bias & Fairness Audit

Upload Data → Select Protected Groups → Select Fairness Metrics → The Bias Report

Fairlearn

**Fairness Indicators BETA**

Select metrics to display:
- ☑ post_export_metrics/false_negative_r...
- ☐ post_export_metrics/false_positive_rate
- ☐ post_export_metrics/negative_rate
- ☐ post_export_metrics/positive_rate
- ☐ post_export_metrics/true_negative_rate
- ☐ post_export_metrics/true_positive_rate
- ☐ accuracy
- ☐ accuracy_baseline
- ☐ auc
- ☐ auc_precision_recall

Baseline: Overall    Thresholds: 0.50

post_export_metrics/false_negative_rate



algofairness / fairness-comparison

⟨⟩ Code    ⊙ Issues 1    ⎇ Pull requests 0    ▷ Actic

ML-fairness-gym: A Tool for Exploring Long-Term Impacts of Machine Learning Systems
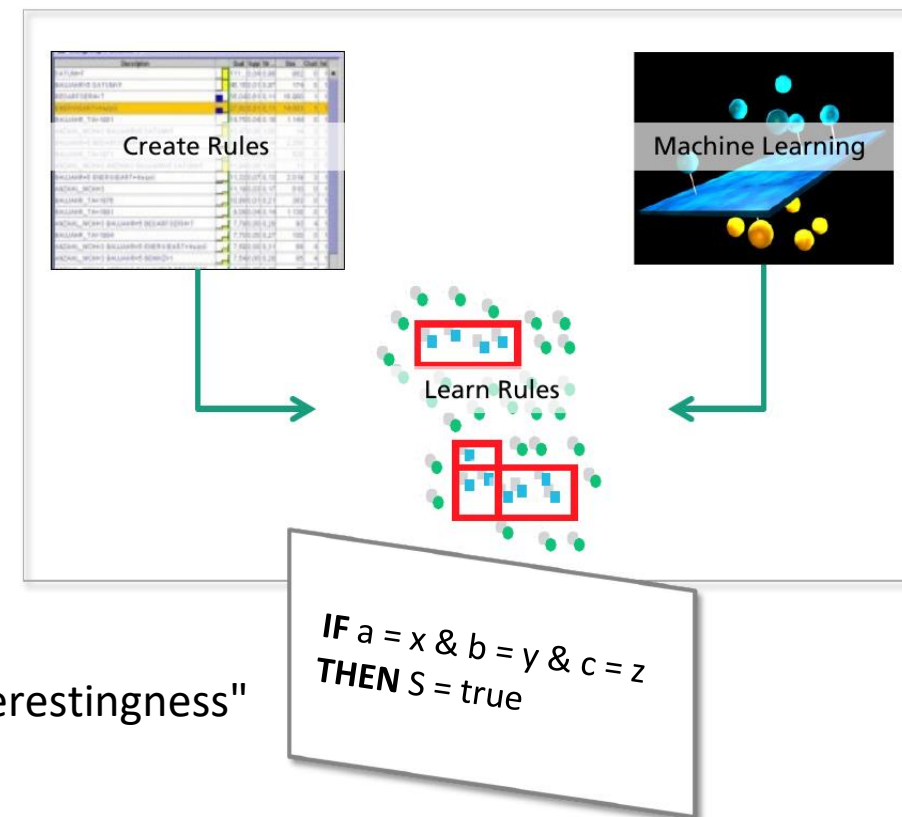Wednesday, February 5, 2020

# AUTOMATED RULE SEARCH FOR ERROR CONDITIONS

## Challenges

- Black-box models complex and powerful, but predictions **difficult or impossible to understand**

- **Huge search space** for potential failure modes

- **Reason for failure** in individual cases **difficult to explain** or generalize

## Approach

- Subgroup search finds rules for error cases with maximum "interestingness" (size and accuracy) regarding search criterion

- Prerequisite: meaningful **metadata** available or generatable



**IF** $a = x$ & $b = y$ & $c = z$
**THEN** $S = true$

## CONTENT ANALYSIS OF CUSTOMER REQUESTS

### Use Case

- Model classifies customer requests

### Meta data

- Lots of customer data & data concerning claim / request
- Extensible (manually or automatically generated attributes)

### Example

- Request not recognized as damage claim, if
  - Subject line is missing
  - Claim request < 10€



©StockPhotoPro - stock.adobe.com



```json
{
    "Documents": [
        {
            "Title": "De Finibus Bonorum et Malorum",
            "Type": "Article",
            "Author": "Cicero",
            "Date": "45 BC",
            "Keywords": [ "fake", "Latin", "filler", "text" ],
            ...
            "Content": "Lorem ipsum dolor sit amet, consectetur ad
        },
        {
            ...
```

**Insurers need to get ready for upcoming regulation and new business opportunities in AI**

**Therefore a systematic evaluation of AI risks is required**

**AI risks are use case specific – so are countermeasures**

**AI quality assessments require new methods and tools**

Image sources: Fraunhofer IAIS

**Contact**

*Dr. Maximilian Poretschkin*

*Fraunhofer Institute for Intelligent Analysis and Information Systems*

*+49 2241 14-1984*

*maximilian.poretschkin@iais.fraunhofer.de*

EAA e-Conference on
Data Science & Data Ethics

29 June 2021