# CERTIF.AI

TRUST IN ARTIFICIAL INTELLIGENCE

# Testing of Fairness Requirements Under the EU AI Act

Workshop Zertifizierte KI: Technische Prüfung von Fairnessanforderungen
19/06/2024          09:30 – 12:30 (CET)

Presentation by

Dr. Robert Kilian – robert@getcertif.ai
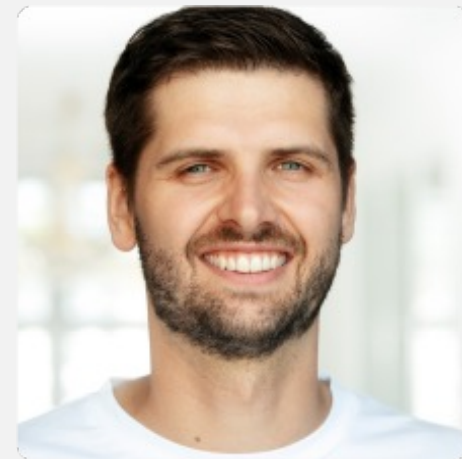Dr. Nico Schmidt – nico@getcertif.ai

CERTIF.AI

# At the Forefront of AI Testing and Certification

We are combining regulatory and technical expertise to provide world-class AI testing and certification services.

## DR. ROBERT KILIAN
### CEO

- 15y+ experience digital business models in highly regulated industries
- Board Member German AI Association with 400+ AI companies as members
- Frequent expert to German parliament and teaches AI regulation at Humboldt University Berlin
- Member of the Microsoft AI Expert Council & DIN AI standardization expert
- Ex-Board Member and Executive N26 Bank; founder data analytics provider Beams; Ex-Hengeler

**KI BUNDESVERBAND** **DIN**

**HENGELER MUELLER** **N26**

## JAN ZAWADZKI
### CTO

- Former Head of Artificial Intelligence at Cariad
- Built central AI hub with 50+ AI experts in Germany and China
- Member of various AI leadership committees
- Guest lectures at Oxford University and ESMT Berlin
- Created the AI Project Canvas
- Studied Data Science & Business; Former Management Consultant at EY

**CARIAD** **EY**

## STEPHANIE JONKERS

- Lead AI Tooling Expert at CertifAI
- Former Technical Program Manager at AWS
- Former Senior Data Scientist at CARIAD

**CARIAD** **aws**

## DR. NICO SCHMIDT

- Lead Safe AI Data Scientist at CertifAI
- Former Lead Architect Data Loop @CARIAD
- Autonomous Driving Research at VW

**CARIAD** **VW**

**SHAREHOLDERS**    **DEKRA**    **pwc**

**Recent Projects include**

**MISSION KI**    Development of quality and testing standards for AI systems as part of the BMDV project Mission KI advising the Federal Government

# Testing and Certification as a Solution

The main challenges of AI system providers can be overcome by testing and certifying the individual AI system.

## STAKEHOLDERS TRUST

- Testing and certification will help **bolster B2B customer, end consumer, supplier confidence** in **product safety and quality**.
- This will also increase investor appeal.

## LEGAL COMPLIANCE

- Testing and certification provide for **EU market access** through **legally required conformity assessments** and model validations under the applicable risk regulations.
- It also is necessary to **comply** with the **EU AI Act requirements** to **avoid fines**.

## LIABILITY SHIELD

- Testing and certification by a third-party independent expert **mitigate risks** and **protect** both **corporate and managerial** liability.

# 1 Fairness According to the EU AI Act

# Fairness Requirements Under the EU AI Act

For high-risk AI systems the EU AI Act is providing for fairness obligations regarding the used data sets.

| | |
|---|---|
| **Rec. 27 EU AI Act** | Recognition of **diversity, non-discrimination and fairness** as one of the 7 **AI HLEG principles**<br><br>*„AI systems are developed and used including diverse actors and **promoting equal access,**<br>gender equality and cultural diversity, while avoiding discriminatory impacts and unfair biases** that are prohibited [...]."* |
| **Art. 10(2) EU AI Act** | **Absence of biases** in **training, validation and testing data** as part of high-risk AI system **data governance**<br><br>*(f) **examination** in view of **possible biases** that are likely to **affect the health and safety of persons,**<br>have a negative impact on fundamental rights or lead to discrimination** prohibited under Union law, [...];<br>(g) appropriate measures to detect, prevent and mitigate possible biases identified according to point (f) [...].* |
| **Rec. 67 EU AI Act** | *Biases can for example be inherent in underlying data sets, especially when historical data is being used,<br>or generated when the systems are implemented in real-world settings. Results provided by AI systems could be influenced by<br>such inherent biases that are inclined to gradually increase and thereby perpetuate and amplify existing discrimination,<br>in particular for persons belonging to certain vulnerable groups, including racial or ethnic groups. [...]* |
| **Rec. 58 EU AI Act** | Certain AI systems (**evaluation of credit score or creditworthiness**) are classified as **high-risk** due to **possible discrimination**<br><br>*[...] AI systems used for those purposes may lead to discrimination between persons or groups and<br>may perpetuate historical patterns of discrimination, such as that based on racial or ethnic origins, gender,<br>disabilities, age or sexual orientation, or may create new forms of discriminatory impacts. [...]* |

# Fairness Under German and European Law

Other German and European law also contains requirements for the fairness of AI-supported decisions.

**Applicability** of **equality and anti-discrimination provisions** in the context of **algorithmic decisions**? ✓

## GENERAL LEGAL FRAMEWORK

## DATA PROTECTION LAW

**Art. 18 (34 ff., 45 ff., 56 ff., 63 ff.) TFEU**
Esp. goods, persons, services and capital

**Art. 20, 21, 23 CFR**
Esp. gender, race, skin colour, ethnic or social origin, genetic characteristics, language, religion or belief, political or any other opinion, membership of a national minority, property, birth, disability, age, sexual orientation

**Art. 3 GG**
Esp. gender, origin, race, language, homeland and origin, faith, religious or political views, disability

**Sec. 7(1), 19(1) AGG**
Race or ethnic origin, gender, religion or belief, disability, age or sexual identity

**Art. 9(1) GDPR**
Processing prohibition regarding data on racial or ethnic origin, political opinions, religious or philosophical beliefs, trade union memberships, genetic data, biometric data, health, sex life or sexual orientation

Exception in Art. 10(5) EU AI Act

Applicability in the relationship between private individuals?

# Legal Requirements to Technical Standards

The abstract EU AI Act requirements are being transposed into more specific, actionable technical standards.
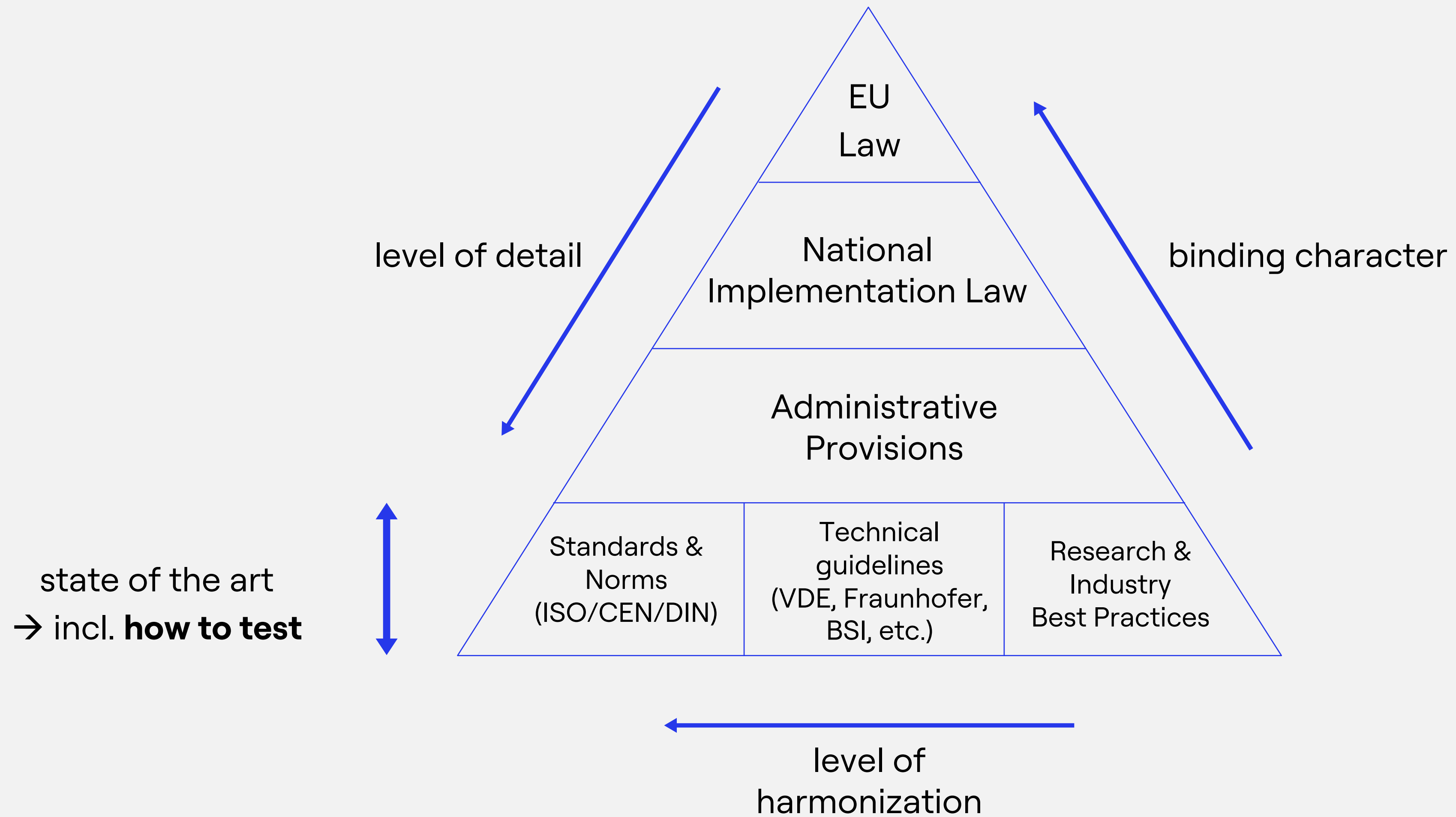
**Specification of Legal Requirements** →

**Fairness requirements**
in the **EU AI Act**
(Chapter III, Section 2)

Art. 10(2)(f) EU AI Act
*[...] possible **biases** that are likely to **affect** the **health and safety of persons**, have a **negative impact on fundamental rights** or lead to **discrimination** prohibited under Union law;*

Art. 10(2)(g) EU AI Act
*[...] appropriate measures to **detect, prevent and mitigate possible biases** identified according to point (f);*

**Standardization request** of the **European Commission** according to Art. 40(2) EU AI Act to **CEN/CLC** (C(2023)3215) on May 22, 2023

Annex II 2.2(a) C(2023)3215
*This (these) European standard(s) [...] shall: Include **specifications** for appropriate data governance and data management procedures to be implemented by providers of AI systems (with specific focus on [...] **procedures for detecting and addressing biases** and **potential for proxy discrimination** or any other relevant shortcomings in data); [...]*

**CEN–CENELEC JTC 21** is considering **existing standards for harmonisation** and **developing new ones** by April 30, 2025 (Art. 1 (C(2023)3215))

| Project reference | Status |
|---|---|
| **EN ISO/IEC 22989:2023/prA1** (WI=JT021031) Information technology — Artificial intelligence — Artificial intelligence concepts and terminology — Amendment 1 | Under Drafting |
| **EN ISO/IEC 23053:2023/prA1** (WI=JT021032) Framework for Artificial Intelligence (AI) Systems Using Machine Learning (ML) — Amendment 1 | Under Drafting |
| **FprCEN/CLC ISO/IEC/TS 12791** (WI=JT021013) Information technology - Artificial intelligence - Treatment of unwanted bias in classification and regression machine learning tasks (ISO/IEC DTS 12791:2023) | Approved |
| **FprCEN/CLC/TR 18115** (WI=JT021007) Data governance and quality for AI within the European context | Under Approval |
| **prCEN/CLC/TR 17894** (WI=JT021001) Artificial Intelligence Conformity Assessment | Under Drafting |

**CEN/CLC technical standards** as **harmonised standards** in the Official Journal of the European Union

Art. 40(1) EU AI Act
*"[...] **conformity with harmonised standards** [...] shall be **presumed** to be in **conformity with the requirements** set out in **[Chapter III, Section 2]** [...], to the extent that those standards cover those requirements [...]."*

← **Harmonised Standards for the EU AI Act**

# 2 Testing AI Systems For Fairness

CERTIF.AI

# From Legal Obligations to Technical Measures

level of detail

binding character

EU
Law

National
Implementation Law

Administrative
Provisions

| Standards & Norms (ISO/CEN/DIN) | Technical guidelines (VDE, Fraunhofer, BSI, etc.) | Research & Industry Best Practices |
|---|---|---|

state of the art
→ incl. **how to test**

level of
harmonization

# The State of the Art in AI Fairness

## Standards

- ISO/IEC TR 24027:2021 - Bias in AI systems and AI aided decision making

- ISO/IEC DTS 12791:2023 - Treatment of unwanted bias in classification and regression machine learning tasks

- IEEE P7003TM Standard for Algorithmic Bias Considerations

- DIN SPEC 91512 – Fairness von KI in Finanzdienstleistungen *(under development)*

## Technical Guidelines / Catalogs



## Research & Industry Best Practices

# Requirements Towards Fairness in AI Systems



**EU Law**

**National Implementation Law**

**Administrative Provisions**

| Standards & Norms (ISO/CEN/DIN) | Technical guidelines (VDE, Fraunhofer, BSI, etc.) | Research & Industry Best Practices |

E.g. ISO 12791, ISO 24027, Fraunhofer catalogue

- **Fairness management requirements**
  o Risk analysis documentation and integration with risk management
  o Identifying bias requirements (stakeholders, compliance)
  o Identifying potentially disadvantaged groups
  o Determining a suitable fairness approach
  o Fairness Acceptance criteria

- **Data requirements**
  o Data representation and labeling guide/specs
  o Selection and documentation of data sources

- **Quantifying fairness**
  o in the model output
  o in training & testing data

- **Re-evaluation, continuous validation, operations and monitoring**

# Potentially Disadvantaged Groups in AI Applications

| Basis for potential discrimination | Finance / Insurance, Credit Scoring (Tabular Data) | HR / Hiring, Promotion (Tabular Data) | Healthcare / Disease Diagnosis (Tabular Data) | Healthcare / Med. Imaging (Image/MRI Data) | Automotive ADAS/AD (Image/Video Data) | Automotive Infotainment (Speech Data) | Chatbots, Personal Assistants (Text Data) |
|---|---|---|---|---|---|---|---|
| Age | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ |
| Gender | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ |
| Ethnicity, national or geographic origin | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ |
| Skin color, hair color, size, weight | | ⚠️ | ⚠️ | ⚠️ | ⚠️ | | ⚠️ |
| Mental or physical disability | ⚠️ | | ⚠️ | ⚠️ | ⚠️ | ⚠️ | ⚠️ |
| Genetic information | | | ⚠️ | | | | |
| Pregnancy or parenthood | ⚠️ | ⚠️ | ⚠️ | ⚠️ | | | ⚠️ |
| Religious beliefs or ideology | ⚠️ | ⚠️ | | | | | ⚠️ |
| Sexual identity | | | | | | | ⚠️ |
| Relationship to someone subject to discrimination | ⚠️ | | | | | | ⚠️ |
| Membership to a specific opinion group or union | ⚠️ | ⚠️ | | | | | ⚠️ |

# Static and Dynamic Testing

## Static Testing of Datasets

**Metrics:** (mostly data quality from ISO 5259-2 applied to subgroups)
Auditability, balance, currentness, completeness, accuracy, consistency, diversity, effectiveness, precision, relevance, representativeness, similarity, timeliness.

**Data:** Training, validation and test data.
The data needs to have annotations about the at-risk group.
(as meta data or labels)

**Method:** Calculate the metrics for each at-risk group.
Compare the distribution of variables in the training and test data to the production data.

**Example:** Representativeness ratio - ratio of relevant attributes found in the subjects of a population to the attributes found in a sample.

$$\frac{A}{B}$$

where
$A$ is the number of target attributes in the sample (e.g. different skin colours in computer vision);
$B$ is the number of attributes in the population.
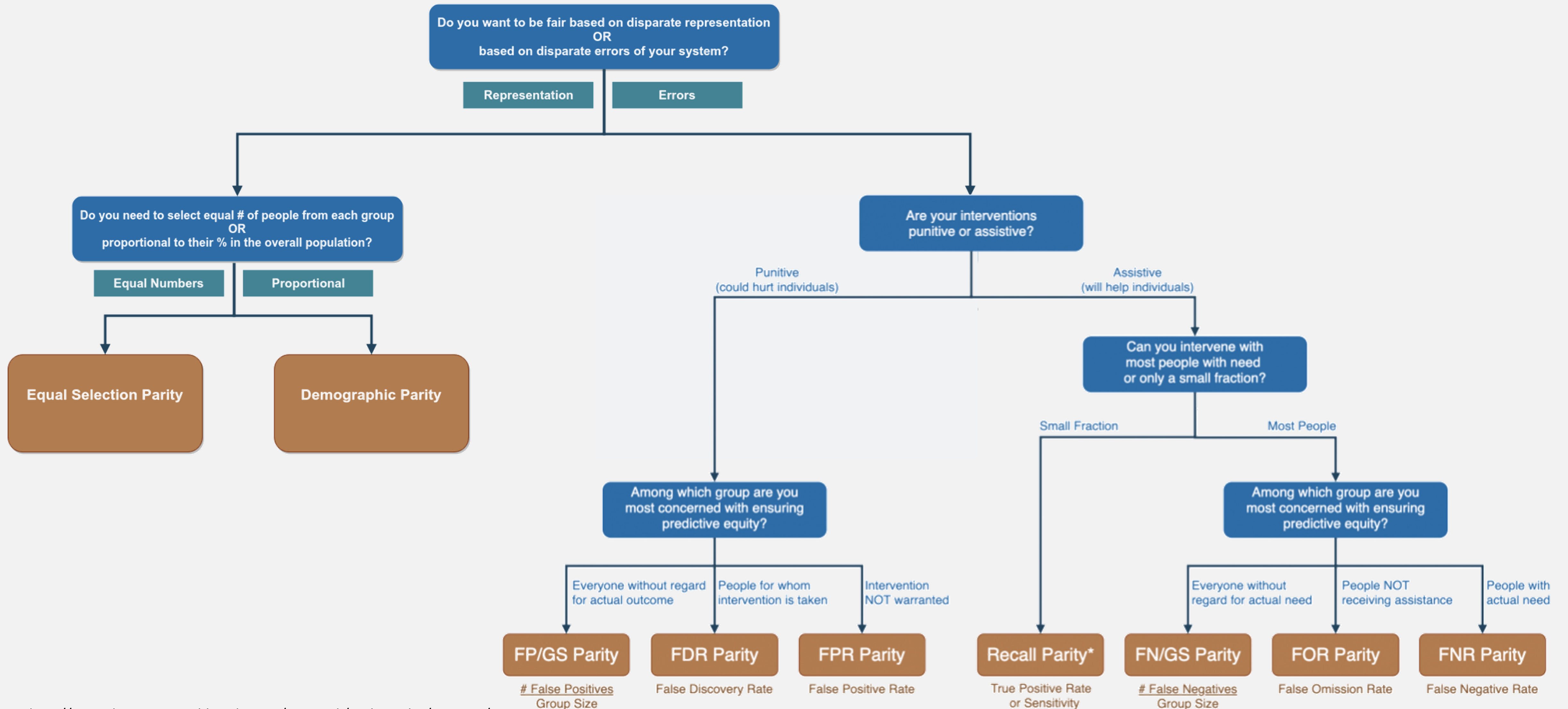
## Dynamic Testing of Model Outputs

**Metrics:** Metrics used for assessing model performance (accuracy, confusion matrix) and fairness metrics (equalized odds, demographic parity, equality of opportunity).

**Data:** Test data with identifiers linking it to at-risk groups.

**Method:** Compare the performance metric for at-risk groups and the population, determine if the delta is sufficiently small.
Calculate fairness metrics for at-risk groups.

Tests need to be conducted on the ML model and the entire AI component.

**Example:** Equality of opportunity – equal True Positive Rates across demographic categories.

$$P\left(\widehat{Y} = \widehat{y} \mid A = m\right) = P\left(\widehat{Y} = \widehat{y} \mid A = n\right)$$

For all values m, n that A can take.

# Static and Dynamic Testing

## Static Testing of Datasets



$$A = 2$$
$$B = 6$$

$$\text{Representativeness ratio} = \frac{1}{3}$$

[Kocak, Burak (2022) doi: 10.5152/dir.2022.211297]

**Example:** Representativeness ratio - ratio of relevant attributes found in the subjects of a population to the attributes found in a sample.

$$\frac{A}{B}$$

where
A is the number of target attributes in the sample (e.g. different skin colours in computer vision);
B is the number of attributes in the population.

## Dynamic Testing of Model Outputs



https://pair.withgoogle.com/explorables/measuring-fairness/

**Example:** Equality of opportunity – equal True Positive Rates across demographic categories.

$$P\left(\widehat{Y} = \widehat{y} \middle| A = m\right) = P\left(\widehat{Y} = \widehat{y} \middle| A = n\right)$$

For all values m, n that A can take.

# How to set Fairness Goals



http://www.datasciencepublicpolicy.org/our-work/tools-guides/aequitas/

Testing of Fairness Requirements Under the EU AI Act

# Bias Mitigation Measures

## Data-based methods

- **Up-sampling or down-sampling**– increasing the representation of underrepresented groups in a dataset

- **Use of synthetic data** - artificially increasing the dataset while reusing the existing dataset

- **Federated learning** – enabling access to a large distributed datasets that can be more representative of the target user base

- **Separate biased validation dataset** – testing the AI system on a customized dataset to check boundary conditions with respect to unwanted bias

## Model-based methods

- **Regularization techniques** - prioritizing learnings from under-sampled data to ensure such learning is not forgotten due to dominant data samples

- **Decoupled classifiers** - training a separate classifier on each group

- **Joint loss function** - using a joint loss function that penalizes differences in classification statistics between groups

- **Disparate impact remover** - editing values used as features to reduce different treatment between the groups

## Post-hoc methods

- **Customization at deployment** – adapting techniques such as continuous and transfer learning to factor for unwanted bias at deployment

- **Re-training at deployment** – combining continuous and transfer learning with federated learning

- **Group-specific decision thresholds** – equalizing false positive rates or other relevant metrics based on predicted outcomes

- **Explainable AI techniques** – explaining predictions of the AI system to detect and monitor bias

# Fairness along the AI Development Process

**Legend:**

Requirements

Tests

Measures

- Decoupled classifiers
- Post-hoc methods
- Disparate Impact / Bias Remover

- Regularization
- Decoupled classifiers
- Transfer learning
- Joint loss functions

Quantifying fairness in the model output

Dynamic testing of model output

Selection and Documentation of Data Sources

Data representation and labeling guide/specs

TRAIN

EVALUATE

EXPLORE

COLLECT

CURATE

ML

- Distribution Sampling
- Augmentation
- Test set bias

Risk analysis documentation and integration w/ risk mgmt

Identify bias requirements (stakeholders, compliance)

Identifying potentially disadvantaged groups

Determining a suitable fairness approach

Fairness Acceptance criteria

TRANSFORM

Data

VALIDATE

FORMULATE

PLAN

Fairness Requirements Management

CODE

Dev

BUILD

TEST

RELEASE

DEPLOY

OPERATE

Ops

MONITOR

Re-evaluation, continuous validation, operations and monitoring

- Monitoring bias
- Explainable AI

Quantifying fairness in training & testing data

Static testing of bias in data

https://www.ml4devs.com/articles/mlops-machine-learning-life-cycle/
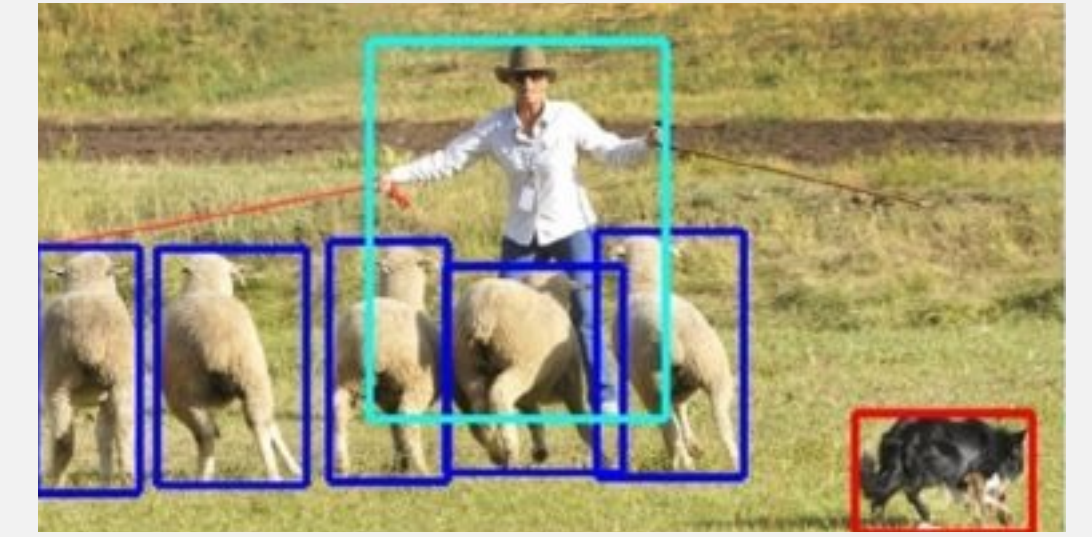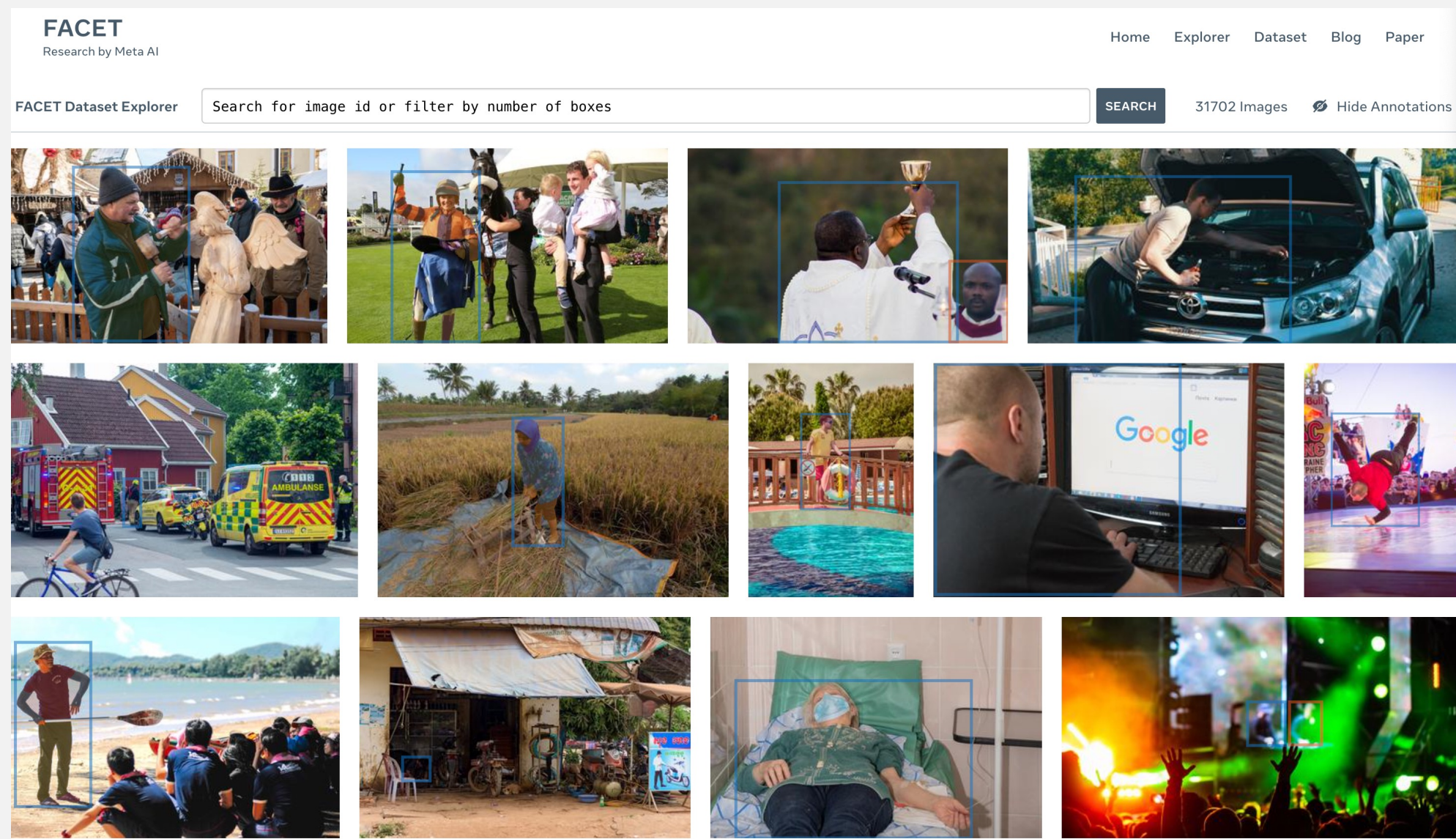
# 3 Examples and Challenges in Fairness Testing of AI Systems

# Fairness in Computer Vision

## Challenge: Data Annotations for Subgroups

- **Many (public available) datasets do not have annotations about subgroups**

ImageNet



COCO / PASCAL VOC



No subgroup labels available



Open Images Dataset V7



Cityscapes



NuScenes

| Category | Annotations | Ratio of all annotations |
|---|---|---|
| animal | 255 | 0.04% |
| human.pedestrian.adult | 149,921 | 21.61% |
| human.pedestrian.child | 1,934 | 0.28% |
| human.pedestrian.construction_worker | 13,582 | 1.96% |
| human.pedestrian.personal_mobility | 2,281 | 0.33% |
| human.pedestrian.police_officer | 464 | 0.07% |
| human.pedestrian.stroller | 363 | 0.05% |
| human.pedestrian.wheelchair | 35 | 0.01% |
| movable_object.barrier | 88,545 | 12.76% |

A2D2

# Fairness in Computer Vision

gender?



age?



---

### Challenge: Data Annotations for Subgroups

- Many (public available) datasets do not have annotations about subgroups

- **Subgroup annotation in images is hard!**

---

ethnicity?



Skin tone?

# Fairness in Computer Vision

<div style="border: 1px solid blue;">

## Challenge: Data Annotations for Subgroups

- Many (public available) datasets do not have annotations about subgroups

- Subgroup annotation in images is hard!

- **First benchmarks for testing are established**

AI Research

## FACET: Benchmarking fairness of vision models

FACET is a comprehensive benchmark dataset from Meta AI for evaluating the fairness of vision models across classification, detection, instance segmentation, and visual grounding tasks involving people.

https://facet.metademolab.com/

| Size | – 32k images, 50k people |
|---|---|
| Evaluation Annotations | – 52-person related classes<br>– bounding boxes around each person<br>– person/hair/clothing labels for 69k masks |
| Protected Groups | – perceived skin tone<br>– perceived age group<br>– perceived gender presentation |
| Additional Person Attributes | – hair: color, hair type, facial hair<br>– accessories: headscarf, face mask, hat<br>– other: tattoo |
| Miscellaneous Attributes | lighting condition, level of occlusion |

| Perceived or Apparent Attributes | #people | % | #images | % |
|---|---|---|---|---|
| gender presentation | | | | |
| – more stereotypically F | 10k | 21% | 8k | 26% |
| – more stereotypically M | 33k | 67% | 23k | 72% |
| – non-binary | 95 | <1% | 95 | <1% |
| – unknown | 6k | 11% | 5k | 5% |
| Monk Skin Tone | | | | |
| – 1 | 5k | 10% | 4k | 13% |
| – 2 | 20k | 41% | 15k | 48% |
| – 3 | 26k | 53% | 19k | 61% |
| – 4 | 27k | 54% | 20k | 63% |
| – 5 | 22k | 44% | 17k | 54% |
| – 6 | 16k | 33% | 13k | 40% |
| – 7 | 9k | 18% | 7k | 23% |
| – 8 | 5k | 10% | 4k | 13% |
| – 9 | 3k | 6% | 2k | 7% |
| – 10 | 1k | 3% | 1k | 3% |
| – unknown | 18k | 37% | 13k | 42% |
| age | | | | |
| – younger | 9k | 18% | 7k | 23% |
| – middle | 27k | 55% | 20k | 64% |
| – older | 3k | 5% | 2k | 8% |
| – unknown | 10k | 21% | 9k | 27% |

</div>

gender?

age?

ethnicity?

Skin tone?

# Fairness in Computer Vision

## Challenge: Data Annotations for Subgroups

- Many (public available) datasets do not have annotations about subgroups

- Subgroup annotation in images is hard!

- First benchmarks for testing are established

- **However, image datasets are domain specific!**

# Fairness Testing in Large Language Models

## Example: Decoding Trust Assessment

- **Idea: Ask the LLM questions that involve potentially disadvantaged groups**

- **Repeat in different zero-shot and few-shot settings with different degree of bias**

- **Calculate fairness metrics on aggregated outputs**

- **Challenge: How to adapt tests for your downstream task?**



| Model | Sex | | Race | | Age | |
|---|---|---|---|---|---|---|
| | $M_{dpd} \downarrow$ | $M_{eod} \downarrow$ | $M_{dpd} \downarrow$ | $M_{eod} \downarrow$ | $M_{dpd} \downarrow$ | $M_{eod} \downarrow$ |
| GPT-3.5 | **0.17** | **0.20** | **0.14** | **0.17** | **0.09** | **0.15** |
| GPT-4 | 0.21 | 0.26 | 0.16 | 0.28 | 0.14 | 0.20 |

**DR. ROBERT KILIAN**

CEO / CO-MD

robert@getcertif.ai

**DR. NICO SCHMIDT**

LEAD DATA SCIENTIST

nico@getcertif.ai

CERTIF.AI

# CERTIF.AI

TRUST IN ARTIFICIAL INTELLIGENCE