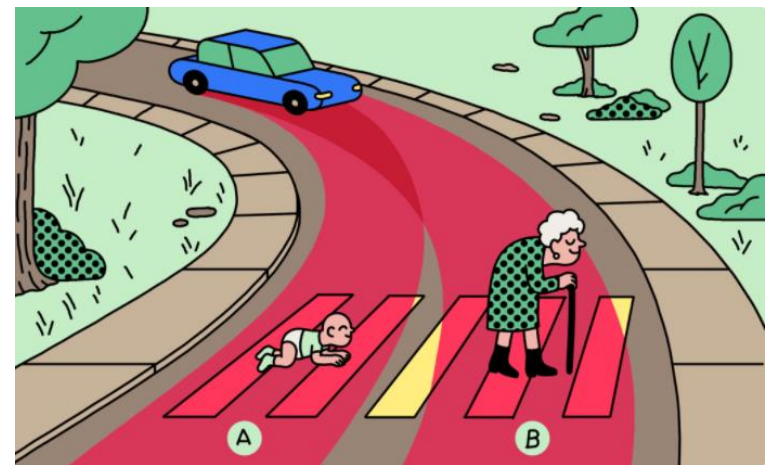# AI Ethics Issues in Real World: Evidence from AI Incident Database

Mengyi Wei

Chair of Cartography and Visual Analytics

Technical University of Munich

mengyi.wei@tum.de
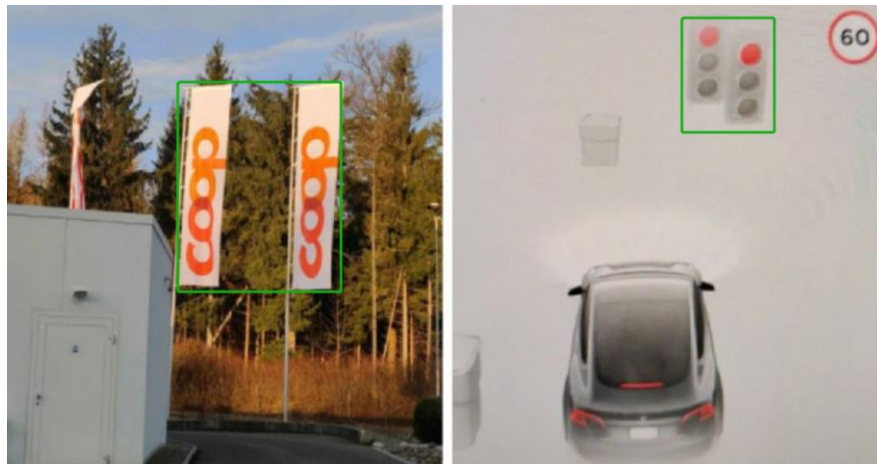
Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

# AI Ethics Issues in Real World: Evidence from AI Incident Database

- Introduction

- Method

- Results

- Limitations and Outlook

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

# AI Ethics Issues in Real World: Evidence from AI Incident Database

## ➢ **Introduction**

- AI technology has become a mega trend.

- AI guidelines are too theoretical and disjointed from practical problems.

- How AI ethics issues take place in real world and how repetitive AI failures can be mitigated?



Self-driving car mistakes red letters on flag for red traffic lights

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

ТШП

## ➢ **Method**

How to describe AI ethics incidents?　➡　Build a taxonomy of AI ethics incidents in real world

- **Data Collection:**

  150 AI ethics incidents from AI Incident Database ( https://incidentdatabase.ai/ )
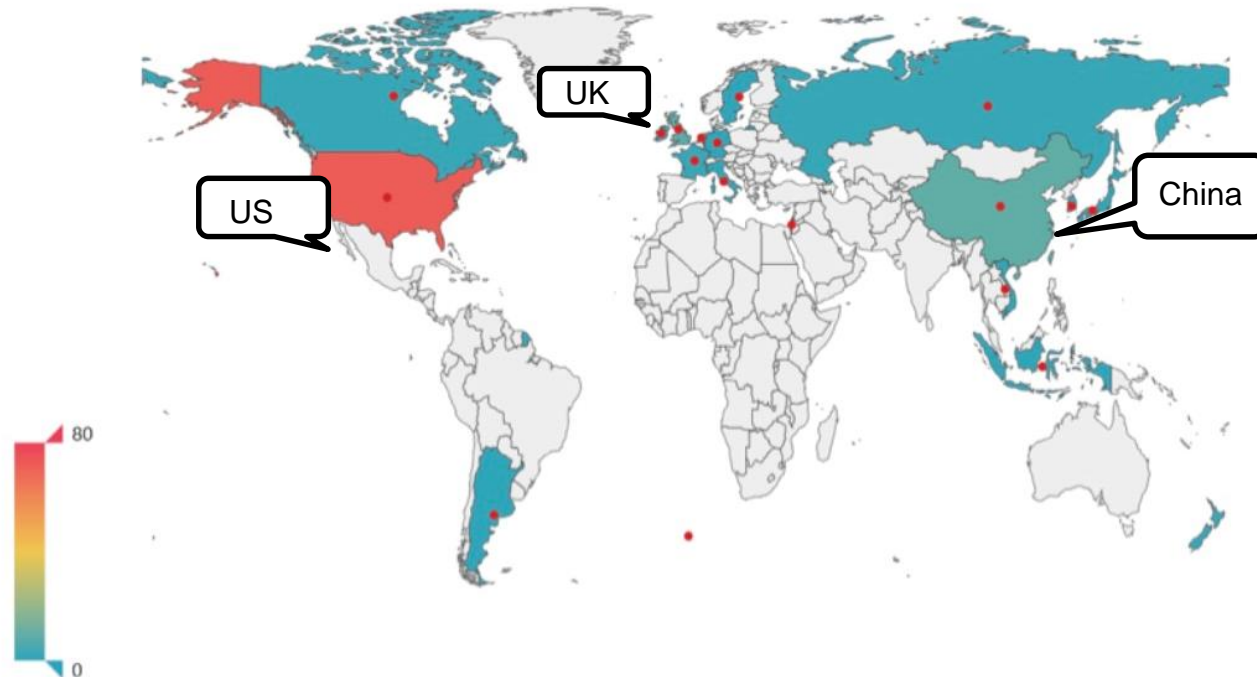
  4 attributes: Time, Geographic locations, Application areas, Taxonomy of AI ethics issues

| No | Title | Time | Location | Application areas | AI ethics issue |
|----|-------|------|----------|-------------------|-----------------|
| 46 | Robot passport checker rejects Asian man's photo for having his eyes closed | 2016.12.07 | New Zealand | Identity Authentication | Racial Discrimination |

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

- Temporal evolution of AI ethics incidents

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

Π∏

- Geographic distribution of AI ethics incidents



Cumulative number of AI ethics incidents from 2010 to 2021

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

TUM

## ➤ **Method**

How to describe AI ethics incidents?  ➡  Build a taxonomy of AI ethics incidents in real world

- **Content Analysis:**

    Krippendorff' s alpha is computed as:

    $$\kappa = \frac{P_A - P_c}{1 - P_c}$$

    where:

    $P_A$ = proportion of units on which the raters agree

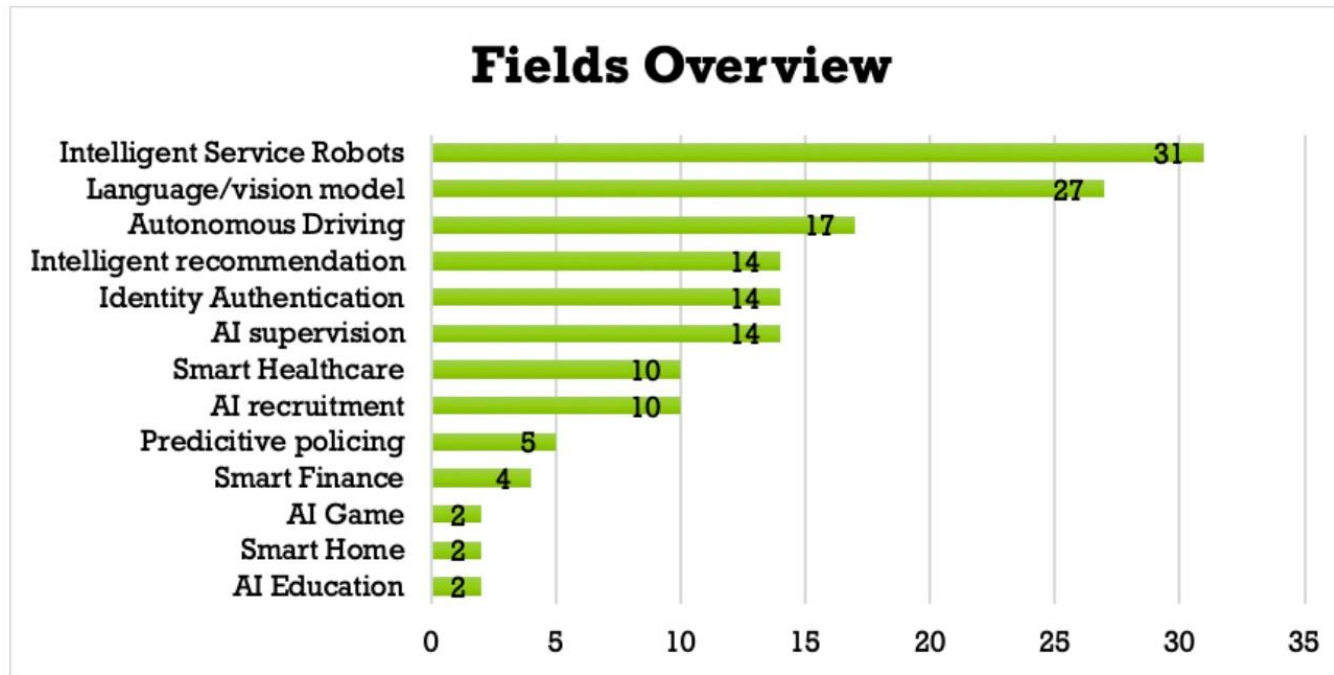    $P_c$ = the proportion of units for which agreement is expected by chance.

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

TUM

## ➢ **Method**

0.94 > 0.8

Almost Perfect

| Kappa Statistic | Strength of Agreement |
|---|---|
| <0.00 | Poor |
| 0.00– 0.20 | Slight |
| 0.21– 0.40 | Fair |
| 0.41– 0.60 | Moderate |
| 0.61– 0.80 | Substantial |
| 0.81– 1.00 | Almost Perfect |

Table 1  Krippendorff's alpha for each variable

| Content Category | Krippendorff's Alpha |
|---|---|
| AI supervision | 0.79 |
| AI recruitment | 0.44 |
| Identity Authentication | 1 |
| Language/vision model | 0.98 |
| Intelligent recommendation | 0.96 |
| Autonomous Driving | 1 |
| Intelligent Service Robots | 1 |
| Smart Healthcare | 1 |
| AI Education | 1 |
| Predicitive policing | 1 |
| Smart Home | 1 |
| AI Game | 1 |
| Smart Finance | 1 |
| Privacy | 1 |
| Inappropriate Use(Bad Performance) | 0.90 |
| Unethical Use(illeagal Use) | 0.97 |
| Racial Discrimination | 1 |
| Gender Discrimination | 0.98 |
| Unfair Algorithm (Evaluation) | 0.94 |
| Mental Health | 0.86 |
| Physical Safety | 1 |
| Average | 0.94 |

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

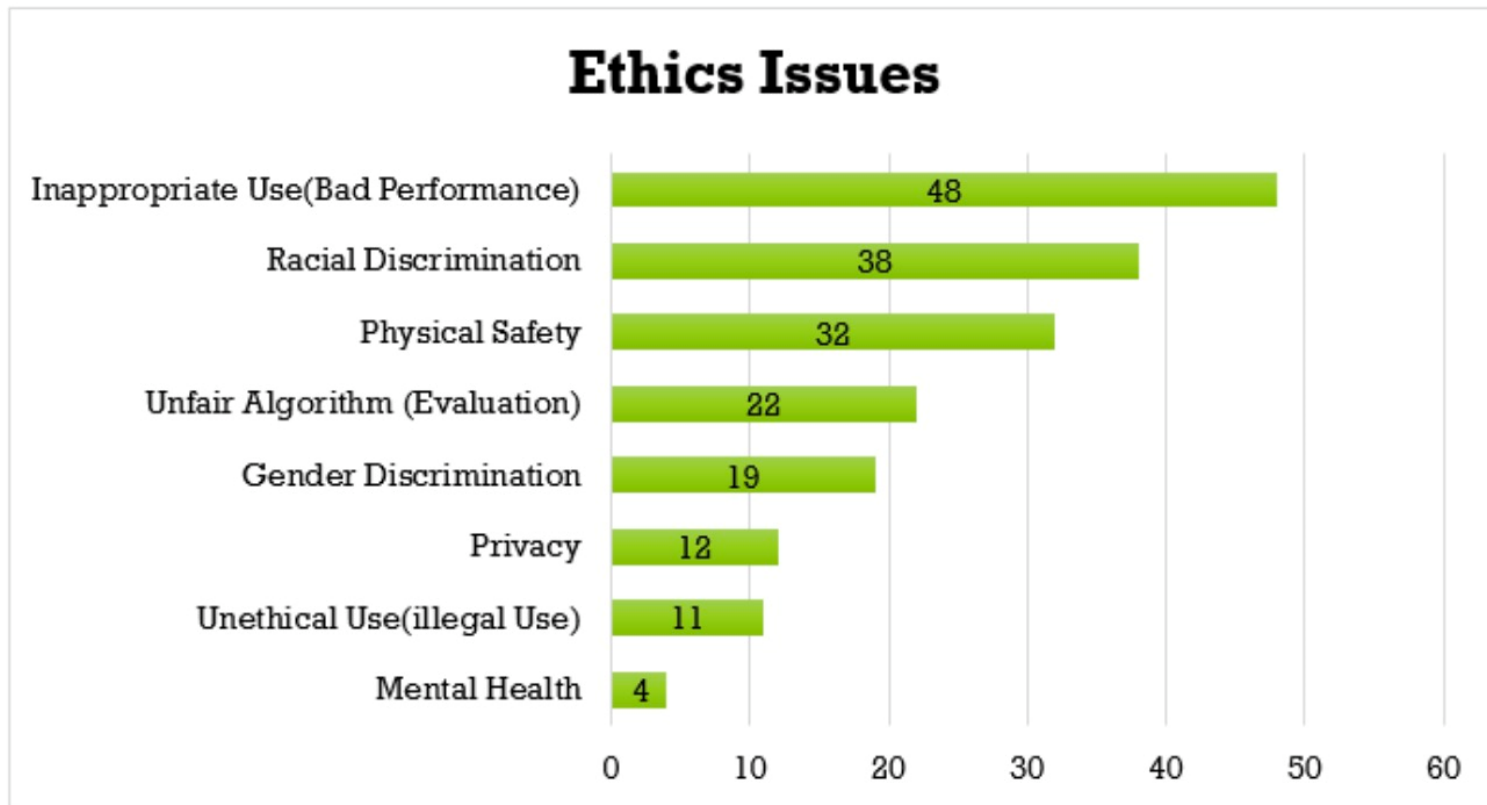AI Ethics Issues in Real World: Evidence from AI Incident Database

# ➢ **Results**

- Application areas of AI ethics incidents



**Fields Overview**

| Field | Value |
|---|---|
| Intelligent Service Robots | 31 |
| Language/vision model | 27 |
| Autonomous Driving | 17 |
| Intelligent recommendation | 14 |
| Identity Authentication | 14 |
| AI supervision | 14 |
| Smart Healthcare | 10 |
| AI recruitment | 10 |
| Predicitive policing | 5 |
| Smart Finance | 4 |
| AI Game | 2 |
| Smart Home | 2 |
| AI Education | 2 |

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

# ➢ **Results**

- Taxonomy of AI ethics issues



**Ethics Issues**

| Issue | Count |
|---|---|
| Inappropriate Use(Bad Performance) | 48 |
| Racial Discrimination | 38 |
| Physical Safety | 32 |
| Unfair Algorithm (Evaluation) | 22 |
| Gender Discrimination | 19 |
| Privacy | 12 |
| Unethical Use(illegal Use) | 11 |
| Mental Health | 4 |

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

ПЛП

## ➢ **Results**

- Distribution of AI ethics issues in different fields

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

ΠΠ

## ➢ **Limitations and Outlook**

- Limitations:

- The size and variety of data are limited.

- Only manually analyze the AI incident database, without applying NLP models to analyze topics and sentiments.

- Outlook:

- Expand the number of AI ethics incident

- Build NLP models to analyze topics and sentiments

- More work to refine the theoretical and operable parts of the guidelines

Assist principle makers in formulating more practical AI guidelines.

Chair of Cartography and Visual Analytics
TUM School of Engineering and Design
Technical University of Munich

AI Ethics Issues in Real World: Evidence from AI Incident Database

TUM

# Thank you !

Mengyi Wei

Chair of Cartography and Visual Analytics

Technical University of Munich

mengyi.wei@tum.de