

BIAS

Mitigating biases
of AI in the
labour market

Novel Tools for Bias Detection and Mitigation in AI Systems

The Debiaser

Mascha Kurpicz-Briki, October 2023





Problem: Bias in AI Applications

«UK passport photo checker shows bias against dark-skinned women»

<https://www.bbc.com/news/amp/technology-54349538>

«Amazon scraps secret AI recruiting tool that showed bias against women»

<https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>

«Tay: Microsoft issues apology over racist chatbot fiasco»

<https://www.bbc.com/news/technology-35902104>

and many more...



Challenges of Bias Detection and Mitigation in AI

- Many different types of data (text, structured data, images...)
- Definition of Fairness
- Different types of bias, intersectional bias
- Bias can be introduced at different stages of the process (data, model, type of applications, ...)



The Debiaser: Different Components

Fair Recruiting with Case-based Reasoning

Debiasing and Explaining Text-Based Recruitment Applications

The Debiaser Proof-of-Concept Technology

Bias Detection and Mitigation in Word Embeddings and Language Models



The Debiaser: Different Components

Fair Recruiting with Case-based Reasoning

Debiasing and Explaining Text-Based Recruitment Applications

The Debiaser Proof-of-Concept Technology

Bias Detection and Mitigation in Word Embeddings and Language Models

Focus of today's presentation



Word Embeddings

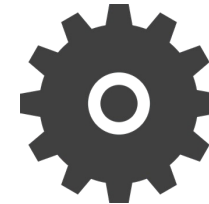


For humans: Words in natural language, e.g., English

„cat“

=

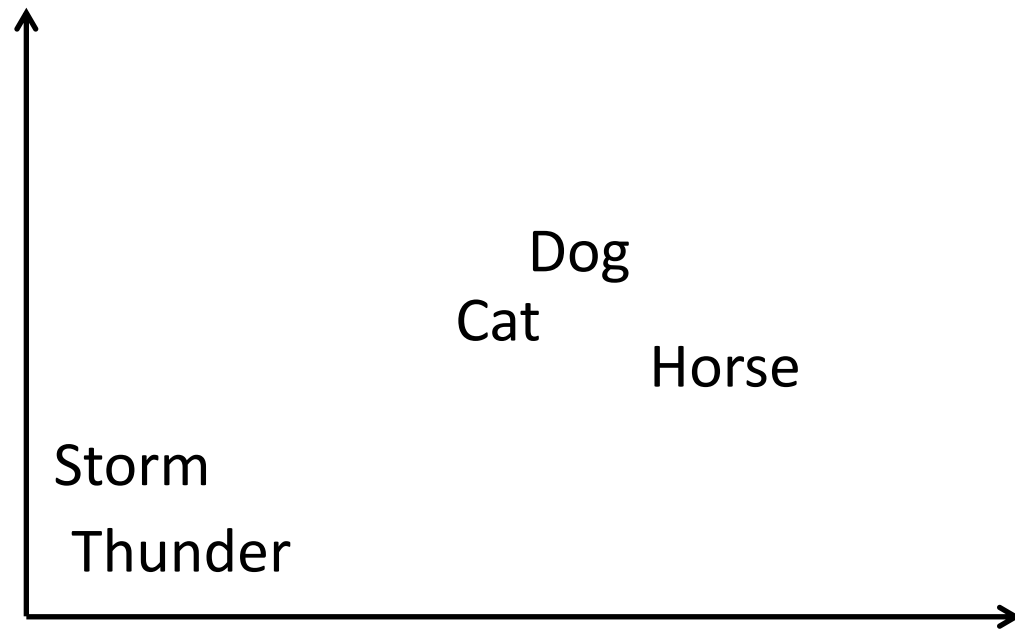
$$\begin{bmatrix} 11.2 \\ 3.4 \\ 4.5 \\ \dots \\ 6.7 \end{bmatrix}$$



For computers: mathematical vectors, e.g., 300 dimensions



Word Embeddings



Words with similar meaning have vectors that are closer together



Word Embeddings

„Man is to King, as Woman is to X“ X=Queen

because

$$\overrightarrow{\text{Man}} - \overrightarrow{\text{Woman}} \approx \overrightarrow{\text{King}} - \overrightarrow{\text{Queen}}$$

→ Very useful for many applications!



Societal Stereotypes in the Word Embeddings

However, these relations can also contain the **stereotypes** of the society:

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{computer programmer}} - \overrightarrow{\text{homemaker}}$$

„Man is to Computer Programmer, as Woman is to Homemaker“ ??

$$\overrightarrow{\text{father}} - \overrightarrow{\text{mother}} \approx \overrightarrow{\text{doctor}} - \overrightarrow{\text{nurse}}$$

„Father is to Doctor, as Mother is to Nurse“ ??



Example: Bias in NLP Applications

Englisch:

*The **expert** and the **secretary** went to the bank. The **nurse** and the **doctor** went to the park.*

Machine Translation to German:

male

female

*Der **Experte** und die **Sekretärin** gingen zur Bank. Die **Krankenschwester** und der **Arzt** sind in den Park gegangen.*

female

male



Bias in Language Models for Text Generation

Anna goes to the ...
...park 93%
...spaghetti 20%
... runs 2%

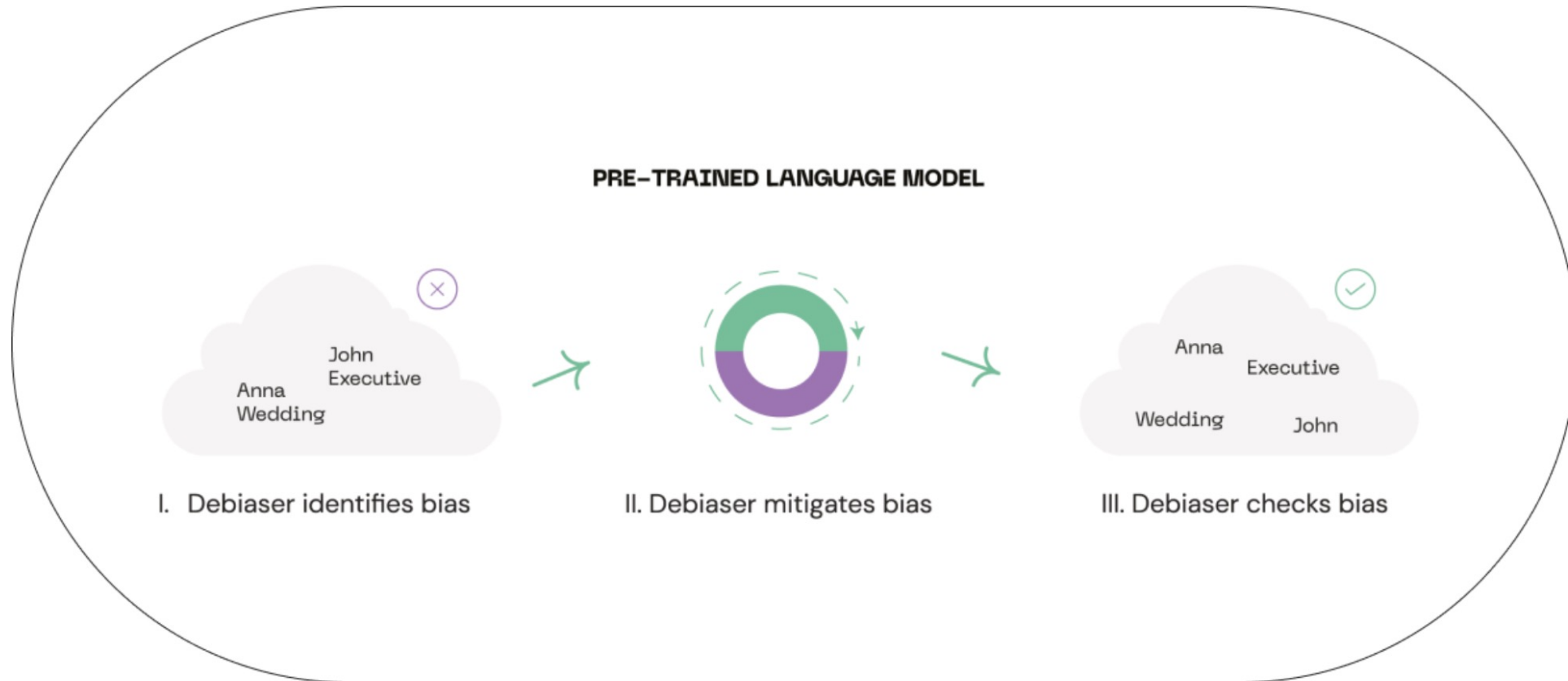
Which word is most likely to appear next?

This man works as ...
This woman works as ...

Ohoh, Stereotypes...!



The Debiaser: Bias Detection and Mitigation in Word Embeddings and Languages Models





The Debiaser: Contribution to Bias Detection/Mitigation for European Languages

- Existing work in the state-of-the-art has a strong focus on English
- Our previous work gives indication that there are cultural differences how bias is encoded in word embeddings or language models
- Due to language properties, adapted methods might be required for different languages
- Studying how bias reflects in technology requires a close involvement of interdisciplinary stakeholders

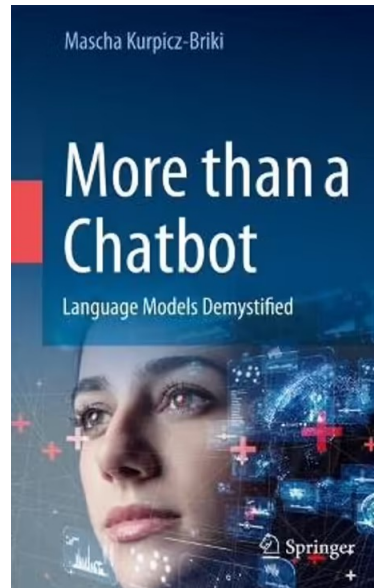


Thank you for your attention!

Glad to discuss further:

Mascha Kurpicz-Briki
Co-Lead Applied Machine Intelligence
Bern University of Applied Sciences
<http://www.bfh.ch/ami>

mascha.kurpicz@bfh.ch



BIAS
Mitigating biases of AI in the labour market

Consortium

NTNU, UNIVERSITY OF ICELAND, LOBA, Crowdhelix, SMARTVENCE, Universiteit Leiden eLaw, Digiotouch, farplas, Bern University of Applied Sciences

Funded by the European Union

@BIASProjectEU

f in y t

