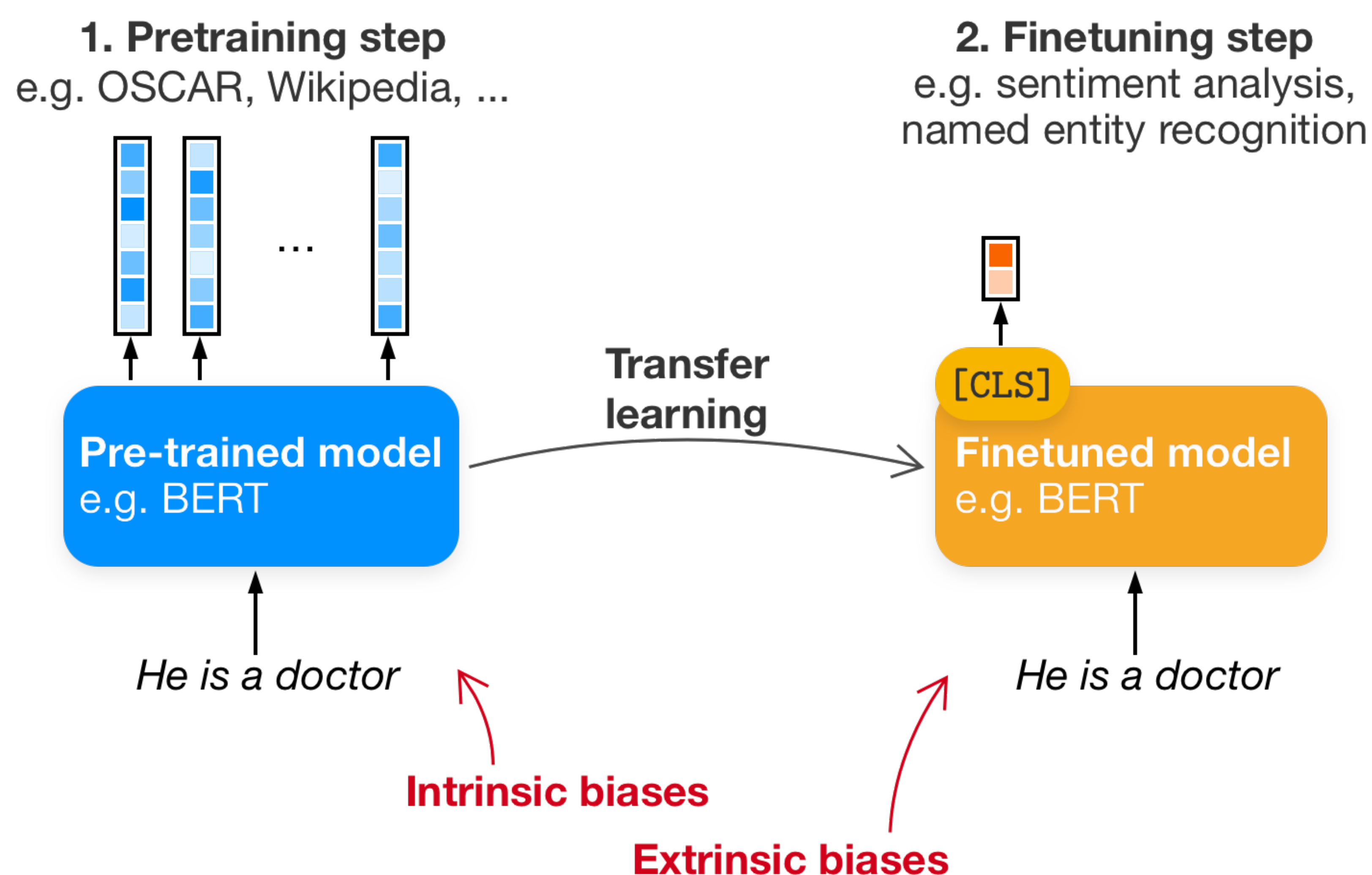


# Measuring Fairness with Biased Rulers: A Comparative Study on Bias Metrics for Pre-trained Language Models

Pieter Delobelle, Ewoenam Kwaku Tokpo, Toon Calders, Bettina Berendt

[pieterdelobelle](#)

## What is *fairness* for language models?

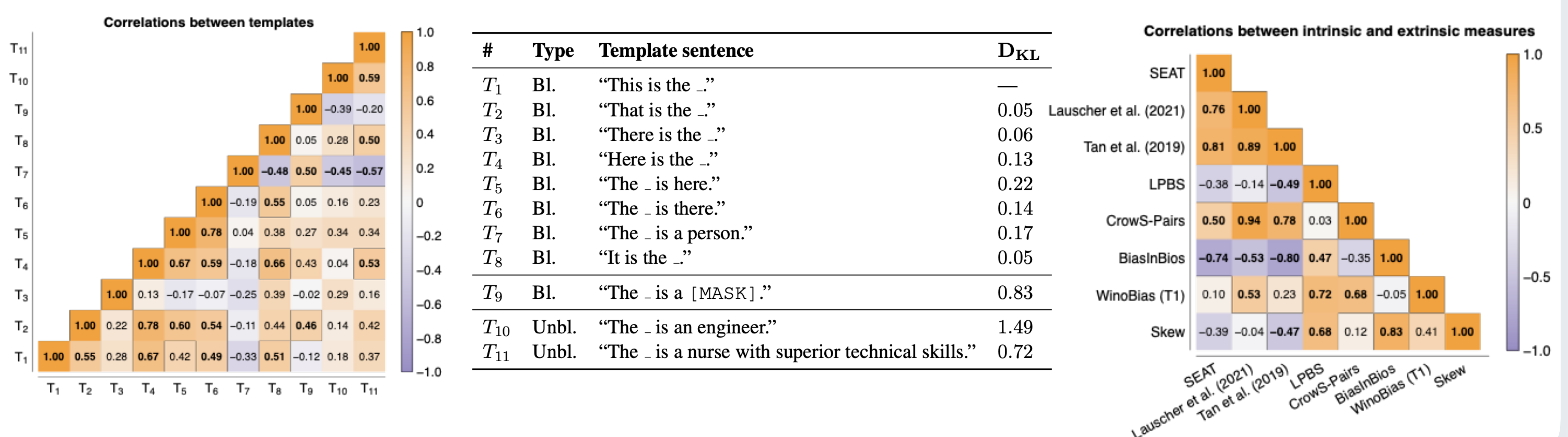


## Fairness metrics

DisCo (Webster et al., 2020)  
Lauscher et al. (2021)  
LPBS (Kurita et al., 2019)  
BEC-Pro (Bartl et al., 2020)  
**Based on WEAT**  
SEAT (May et al., 2019)  
Lauscher et al. (2021)  
Tan and Celis (2019)  
CEAT (Guo and Caliskan, 2021)  
CAT (Nadeem et al., 2021)  
CrowS-Pairs (Nangia et al., 2020)  
Basta et al. (2019)  
Zhao et al. (2019)  
Sedoc and Ungar (2019)

## Fairness metrics don't correlate

We found that metrics are not compatible with each other and highly depend on (i) templates, (ii) attribute and target seeds and (iii) the choice of embeddings.



## Our advice

Use a mix of some intrinsic measures of fairness that don't use embeddings directly and eliminate one source of variance, for example DisCo or LPBS. However, we also recommend to **perform extrinsic fairness evaluations on downstream tasks.**



YouTube

KU LEUVEN

LEUVEN.AI

University of Antwerp

AI FLANDERS  
BUILDING OUR DIGITAL FUTURE

weizenbaum  
institut

DTAI  
DECLARATIVE LANGUAGES &  
ARTIFICIAL INTELLIGENCE



blog post