

# Dissecting Biases in Relation Extraction

## A Cross-Dataset Analysis on People's Gender and Origin

Marco Antonio Stranisci, Pere-Lluís Huguet Cabot, Elisa Bassignana, Roberto Navigli



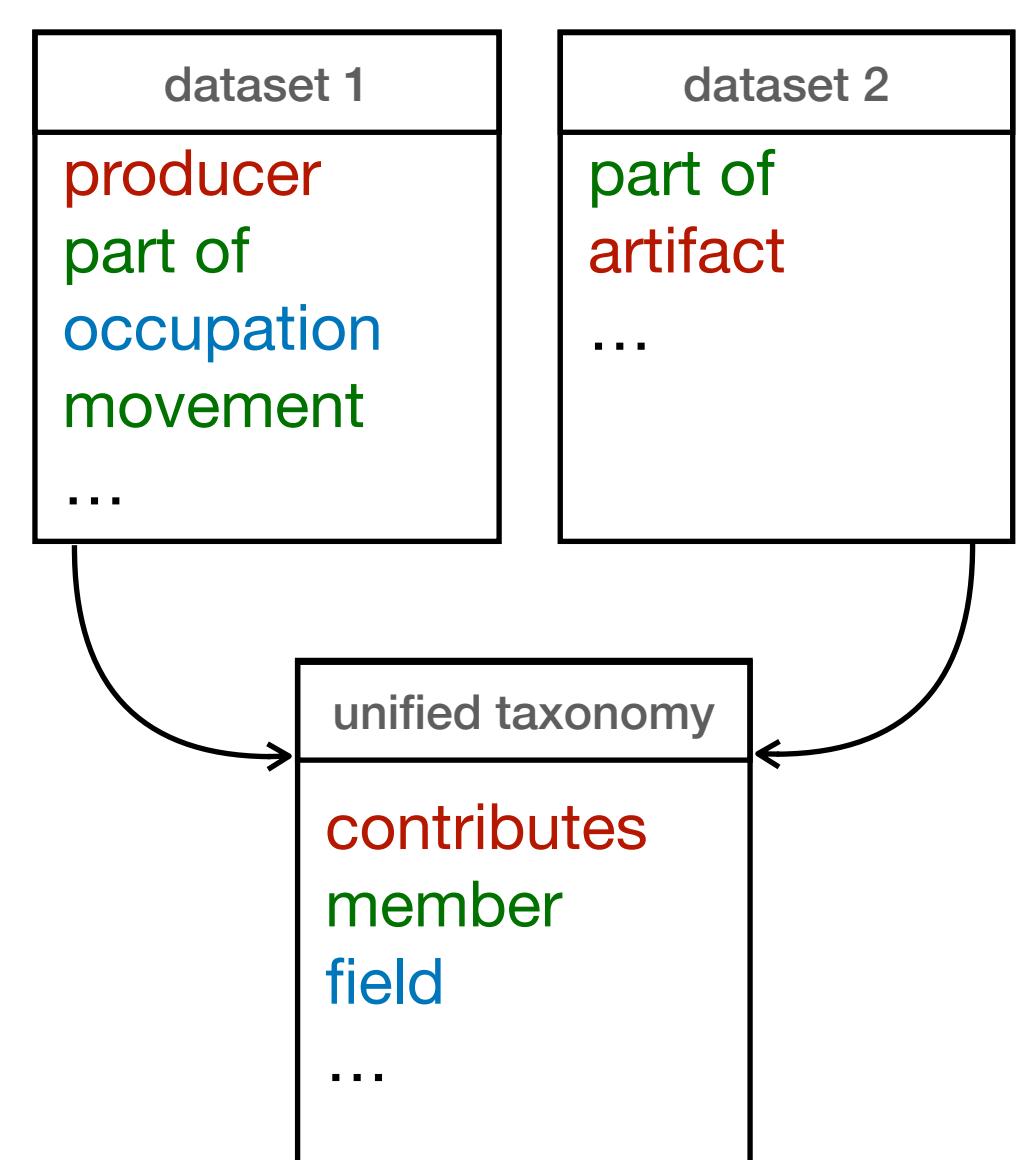
### What?

We propose a methodology for the analysis of socio-demographic biases in the Relation Extraction pipeline, which is completely transparent in terms of interpretability.

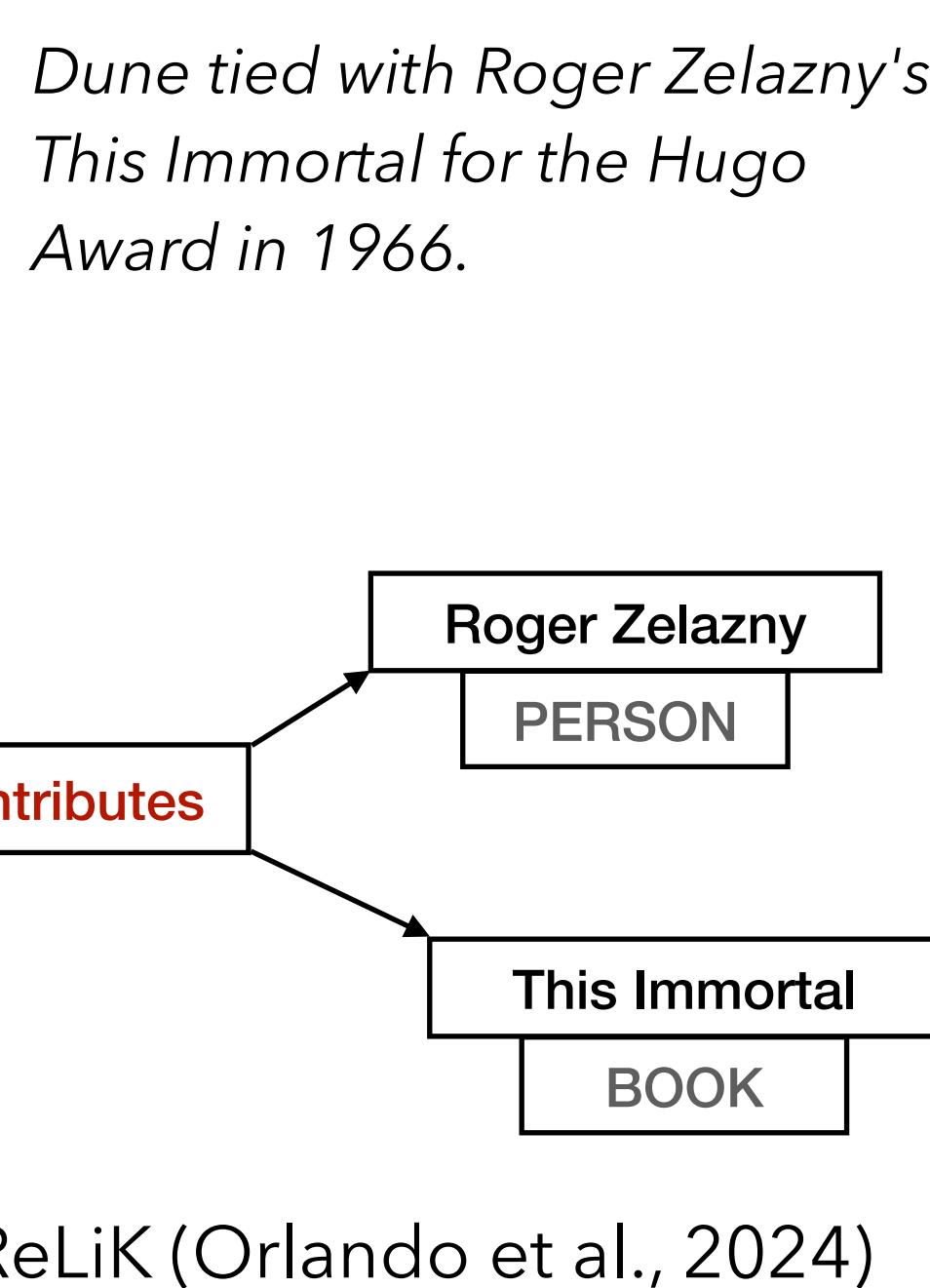


### How?

#### 1 Label alignment

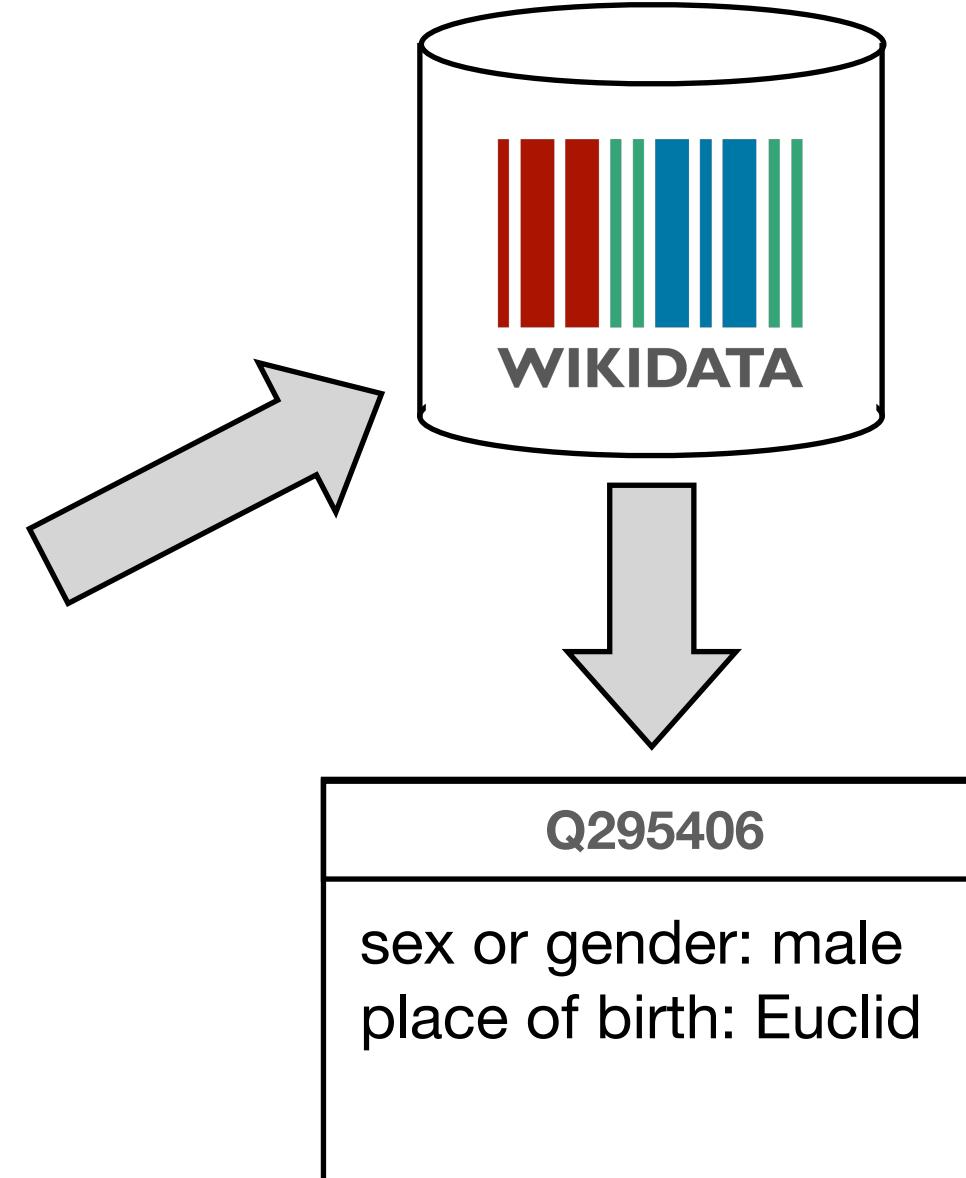


#### 2 Relation Extraction



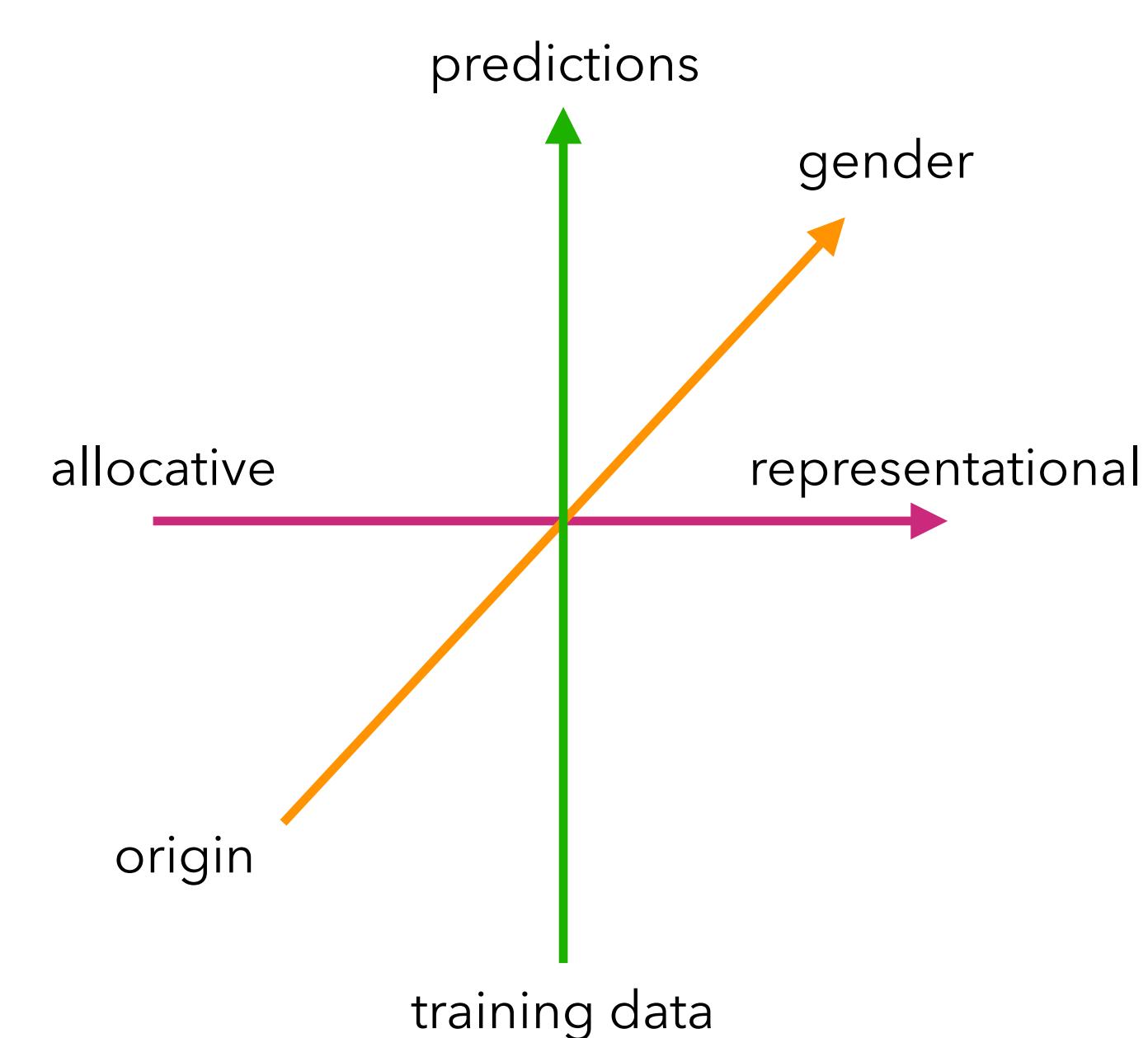
ReLiK (Orlando et al., 2024)

#### 3 Entity Linking



EntQA (Zhang et al., 2022)

#### 4 Bias analysis



## Experiments and Results

### Datasets:

- SRED<sup>FM</sup> (Huguet Cabot et al., 2023)
- CrossRE (Bassignana and Plank, 2022)
- NYT (Riedel et al., 2010)

|               | SRED <sup>FM</sup> |              | CrossRE      |       | NYT          |              |
|---------------|--------------------|--------------|--------------|-------|--------------|--------------|
|               | M                  | W            | M            | W     | N            | S            |
| contributes   | 0.28               | 0.475        | <b>0.407</b> | 0.291 | <b>0.758</b> | 0.162        |
| date          | <b>1.038</b>       | 0.926        |              |       | <b>1.07</b>  | 0.993        |
| field         | <b>0.388</b>       | 0.291        |              |       | 0.394        | 0.451        |
| geographical  | <b>0.469</b>       | 0.368        | 0.218        | 0.218 | 0.501        | <b>0.64</b>  |
| language      | 0.013              | 0.006        |              |       | 0.025        | 0.024        |
| member        | 0.21               | 0.164        | 0.229        | 0.218 | 0.252        | 0.201        |
| participated  | 0.088              | 0.049        | <b>0.278</b> | 0.145 | 0.052        | 0.08         |
| position held | 0.091              | 0.038        | 0.745        | 0.727 | <b>0.085</b> | 0.012        |
| relationship  | 0.124              | <b>0.215</b> | <b>0.098</b> | 0.036 | 0.144        | <b>0.196</b> |
| topic         | 0.001              | 0.001        | 0.018        | 0.018 | 0.002        | 0.002        |

### Allocative bias in training data (% of entities)

|                     | SRED <sup>FM</sup> | CrossRE | NYT   |
|---------------------|--------------------|---------|-------|
| Women               | 20.0%              | 11.8%   | 17.3% |
| Global South        | 18.9%              | 10.0%   | 12.2% |
|                     |                    |         |       |
| Women               | -3.5%              | -2.2%   | +5.6% |
| +SRED <sup>FM</sup> |                    | -5.8%   | +0.6% |
| + gen. balanced     | -2.9%              | -4.4%   | 0.0%  |
|                     |                    |         |       |
| Global South        | -1.7%              | -8.3%   | -2.1% |
| +SRED <sup>FM</sup> |                    | -6.7%   | -1.6% |
| + gen. balanced     | -0.3%              | -9.9%   | -5.9% |

### Allocative bias in predictions

(False Positive Balance score)

|               | SRED <sup>FM</sup> |              | CrossRE      |       | NYT          |              |
|---------------|--------------------|--------------|--------------|-------|--------------|--------------|
|               | M                  | W            | M            | W     | M            | W            |
| contributes   | 0.28               | 0.475        | <b>0.407</b> | 0.291 | <b>0.758</b> | 0.162        |
| date          | <b>1.038</b>       | 0.926        |              |       | <b>1.07</b>  | 0.993        |
| field         | <b>0.388</b>       | 0.291        |              |       | 0.394        | 0.451        |
| geographical  | <b>0.469</b>       | 0.368        | 0.218        | 0.218 | 0.501        | <b>0.64</b>  |
| language      | 0.013              | 0.006        |              |       | 0.025        | 0.024        |
| member        | 0.21               | 0.164        | 0.229        | 0.218 | 0.252        | 0.201        |
| participated  | 0.088              | 0.049        | <b>0.278</b> | 0.145 | 0.052        | 0.08         |
| position held | 0.091              | 0.038        | 0.745        | 0.727 | <b>0.085</b> | 0.012        |
| relationship  | 0.124              | <b>0.215</b> | <b>0.098</b> | 0.036 | 0.144        | <b>0.196</b> |
| topic         | 0.001              | 0.001        | 0.018        | 0.018 | 0.002        | 0.002        |

### Representational bias in training data (t-test statistics)