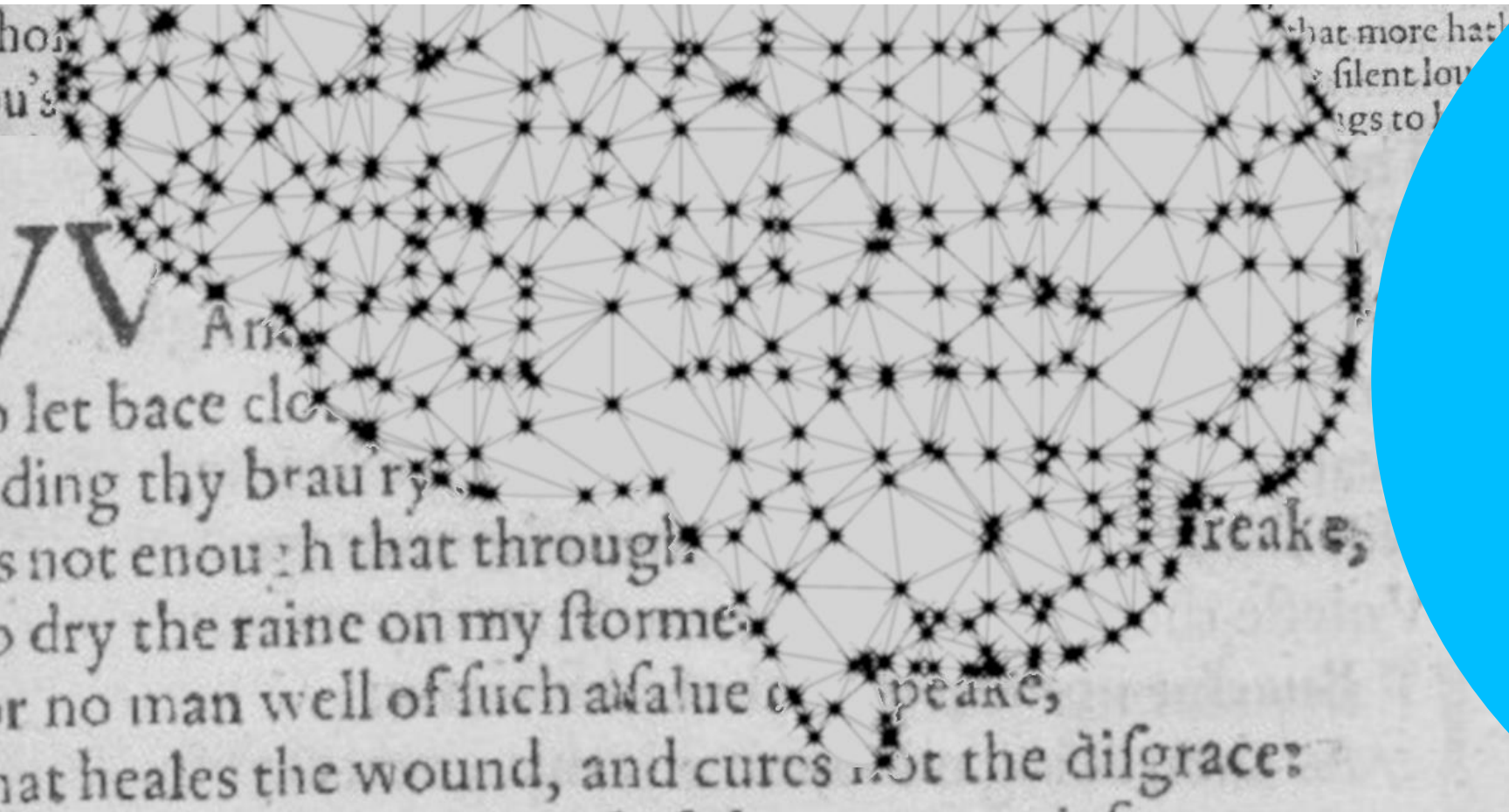**University of Stuttgart**
Institut für Maschinelle Sprachverarbeitung

# Author Classification in Poetry

Katrin Schmidt & Carlotta Quensel

# Author classification

## Goal

- Stylometrics: what age, education, gender, social class has the author

  - Find specific values for one author

- Author prediction: classification task with known classes

  - Which of 5 known authors has written this text

- Our goal: find a method which works well with poetry

  - Style choices might reflect the medium more than the author

  - Special features (number of verses, meter, anaphora, ...)

  - Choice of data is limited

# Author classification
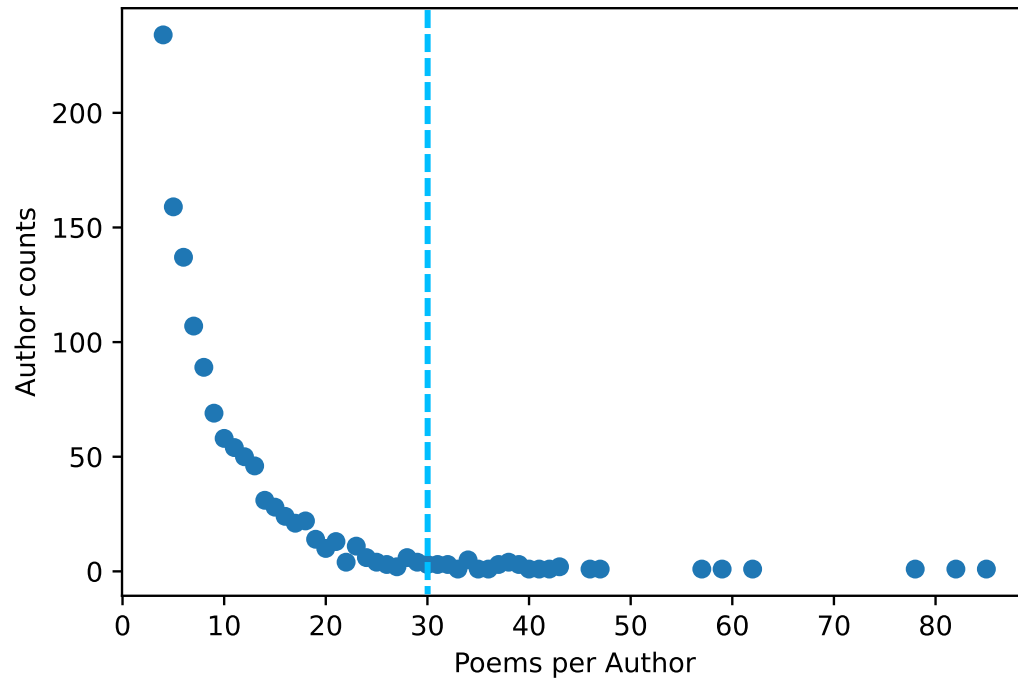
## Data



***Fig 1*** *– Poem distribution per author*

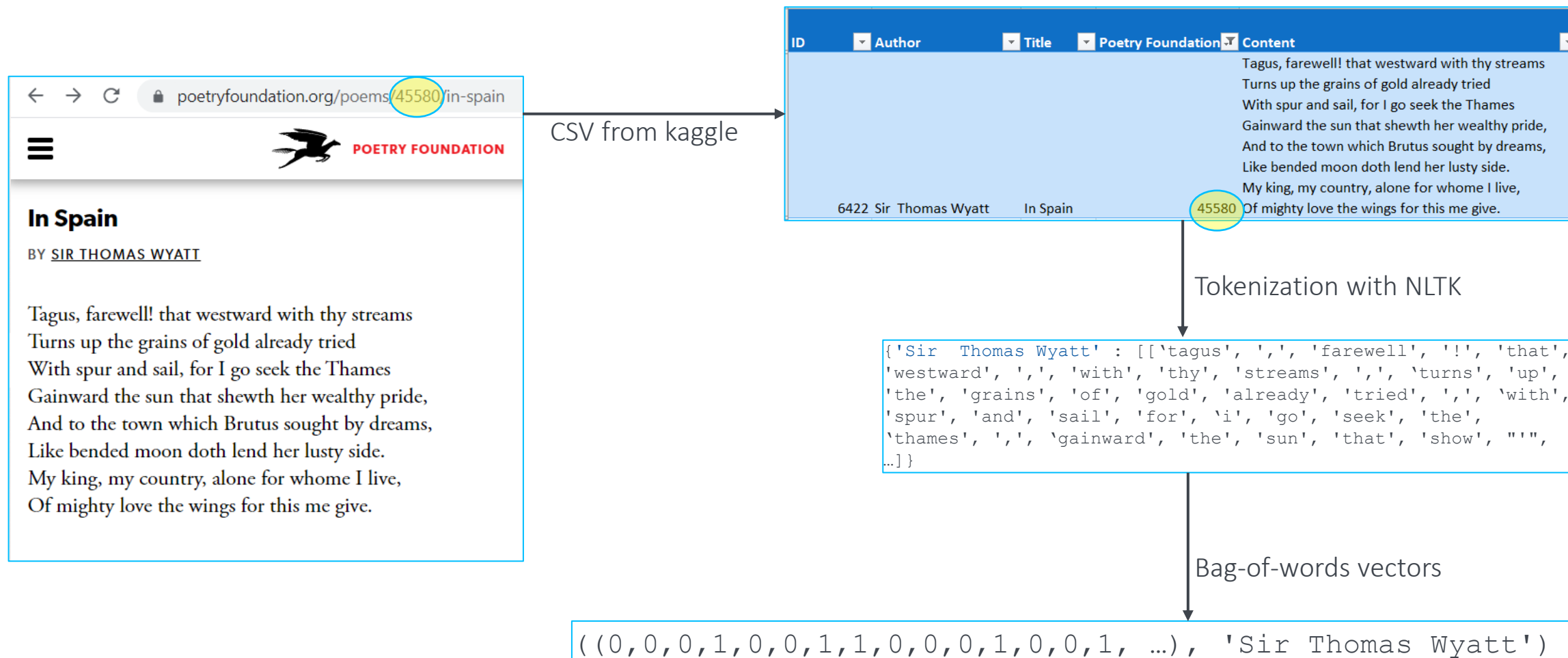Poetry Foundation (founded in 2003)

- 3 309 authors

- 15 567 poems


We decided on the 30 most prolific authors

- More datapoints per class

- Lower computational effort

- 1 569 poems

# Author classification

## Data

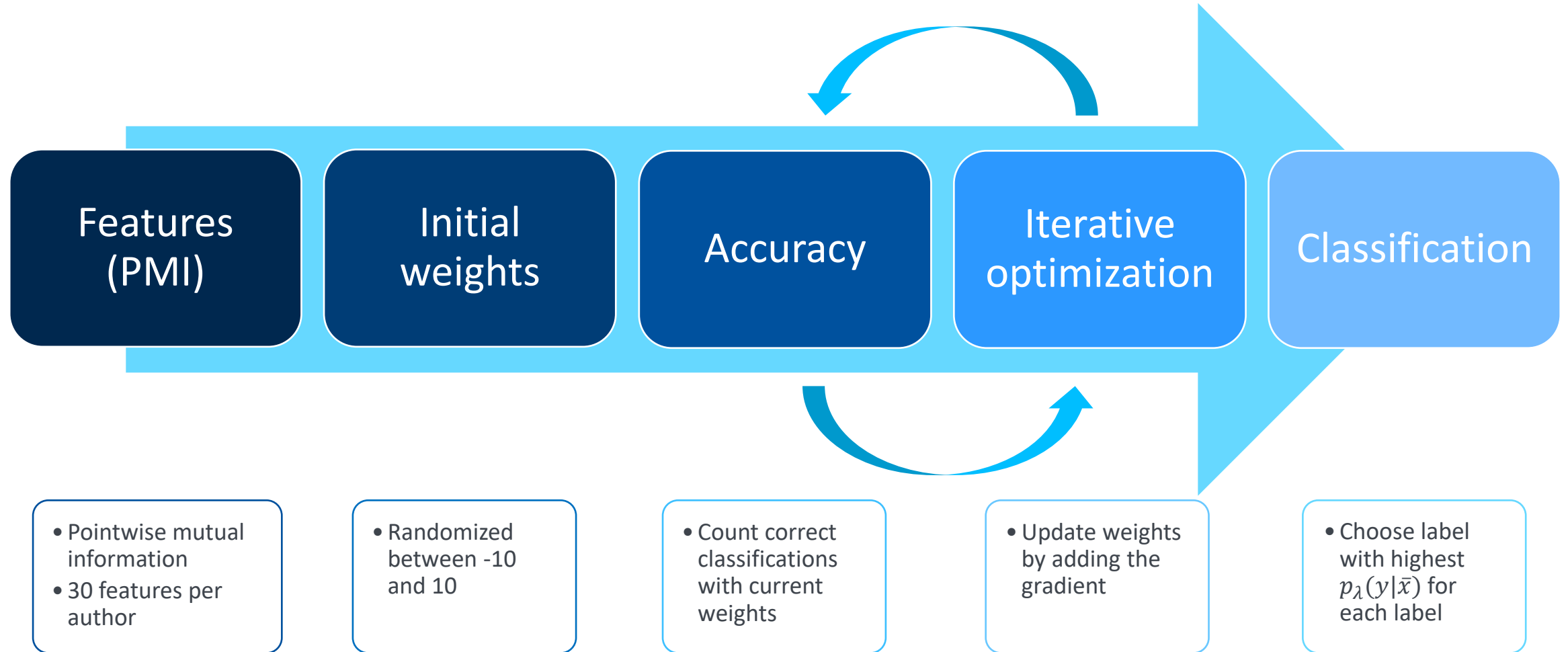# Approach

## Baseline method

- Maximum Entropy Classifier

$$p_\lambda(y|\bar{x}) = \frac{\exp \sum_i \lambda_i \cdot f_i(y, \bar{x})}{\sum_{y'} \exp \sum_i \lambda_i \cdot f_i(y', \bar{x})}$$

- Feature $f_i$: data properties paired with a label

  - e.g.: $f_1(y, \bar{x}) = \begin{cases} 1, & y = \text{Shakespeare} \land \text{thou} \in \bar{x} \\ 0, & \text{otherwise} \end{cases}$

- Weight $\lambda_i$ for each feature to represent the importance of the feature

⇨ Training by optimizing the weights

Baseline method



| Features (PMI) | Initial weights | Accuracy | Iterative optimization | Classification |

- Pointwise mutual information
- 30 features per author

- Randomized between -10 and 10

- Count correct classifications with current weights

- Update weights by adding the gradient

- Choose label with highest $p_\lambda(y|\bar{x})$ for each label
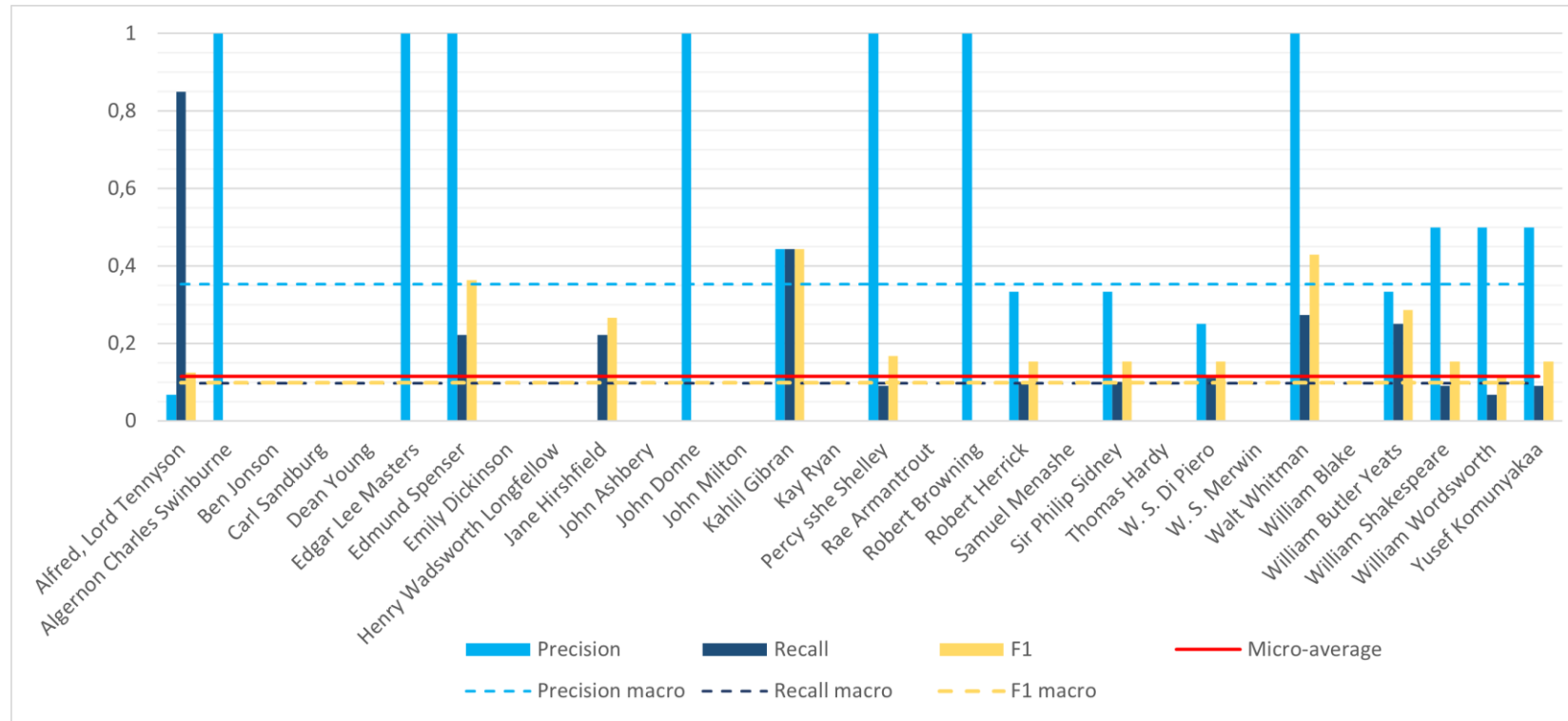
## Quantitative



*Fig 2 – Evaluation on test data (338 poems) for 30 authors, trained on 959 poems*

- Accuracy much worse on the test data
  - test: 0.11
  - training: 0.49
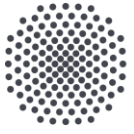
- Many authors never predicted

# Evaluation

## Qualitative

- Many authors never predicted

    **?** Bad features for this author

    **?** Data not evenly distributed


- Features are equally important (random weights approximate each other)

    **?** Different model or feature selection


- The alphabetically first author is chosen too often

    - Every author has the same probability

    **?** Choose most prolific author instead

# Next steps

- Research Question

    **?** Which features are inherent to poetry writers

    **?** Are there inter-dependencies between the features


- Advanced method

    - Extending the features                    ⇨ until 07/05

        - Number of stanzas and verses        ⇨ Katrin

        - Rhyming and rhyme schemes        ⇨ Carlotta

    - Non-linear model for sparse features and interaction

        - Multi-layer NN                            ⇨ until 07/12

**University of Stuttgart**
Germany

# Thank you!

**Katrin Schmidt**
katrin.schmidt@ims.uni-stuttgart.de

**Carlotta Quensel**
carlotta.quensel@ims.uni-stuttgart.de