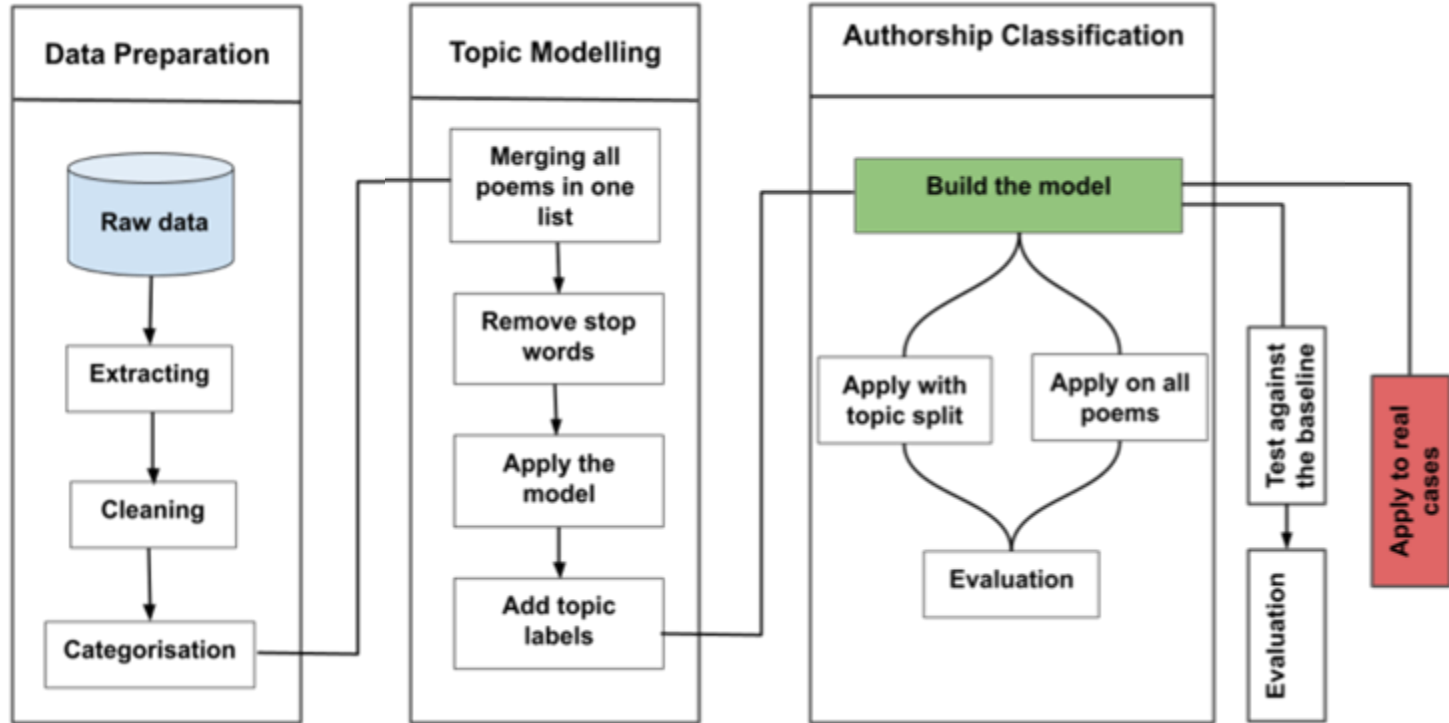# BERT-based Classical Arabic Poetry Authorship Attribution

Lama Alqurashi, Janet Watson, Jacob Blakesley, Serge Sharoff

# Method

# Our Data

Web scraping of **AlDiwan**.etisbew

| Category | Total number |
|----------|--------------|
| Poets | 784 |
| Poems | 77,850 |
| Words | 6,609,495 |

# Topic Modeling

we applied Embedded Topic Modeling (ETM) to label each poem with its topic contributions, further enhancing the dataset's value.  (2020 ,.la te gneiD)

https://github.com/ssharoff/ETM

# Authorship Classification

Ensemble model based on **CAMeLBERT** ssorca detset dna depoleved saw :snoisnemid eerht

- Topic

- Number of poets

- Number of training examples

# Authorship Classification

After parameter optimization, the model achieved F1 scores ranging from 0.97 to
:ta .1.0

- Topic: **No significant effect**

- Number of poets: **Binary**
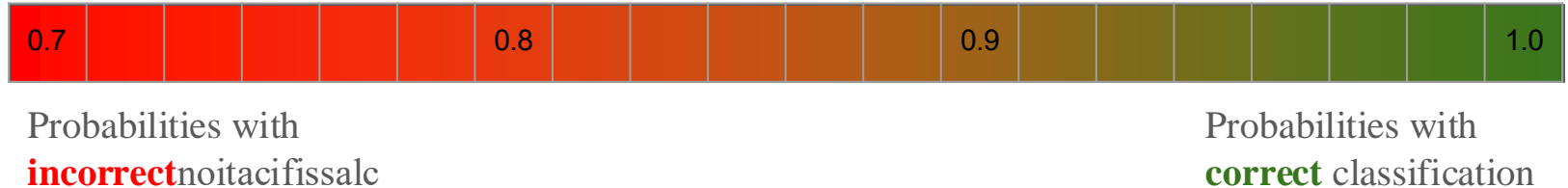
- Number of training examples: **60**

# Authorship Classification

The analysis of **average F1** scores for the tested models shows that the ensemble model achieved a score of  **0.98**

Significantly outperforming the single model, which obtained a score of  .**0.75**

# Authorship Classification

We investigated the probability distributions of correctly classified versus misclassified.

| 0.7 | | | | | 0.8 | | | | | 0.9 | | | | | 1.0 |

Probabilities with
**incorrect**noitacifissalc

Probabilities with
**correct** classification

# Application to Pre-Islamic Mis-attribution Cases

Poem 1(Imru' al-Qays): showed a low confidence score of 0.71

Poem 2(Al-A'shā): The model's confidence in attribution was 0.86

Poem 3(Al-A'shā): 0.85

Poem 4('Ubayd Ibn Al-Abraṣ): high confidence score of 0.99, strongly supporting its attribution

أَلَا لَا أَلَّا إِلَّا لَآلَاءِ لَابِثٍ

وَلَا لَا أَلَّا إِلَا لِآلَاءِ مَن رَحَل

# Thank you ALL!