# Singularity Score for Evaluating Topic Relevance in Tiny Text
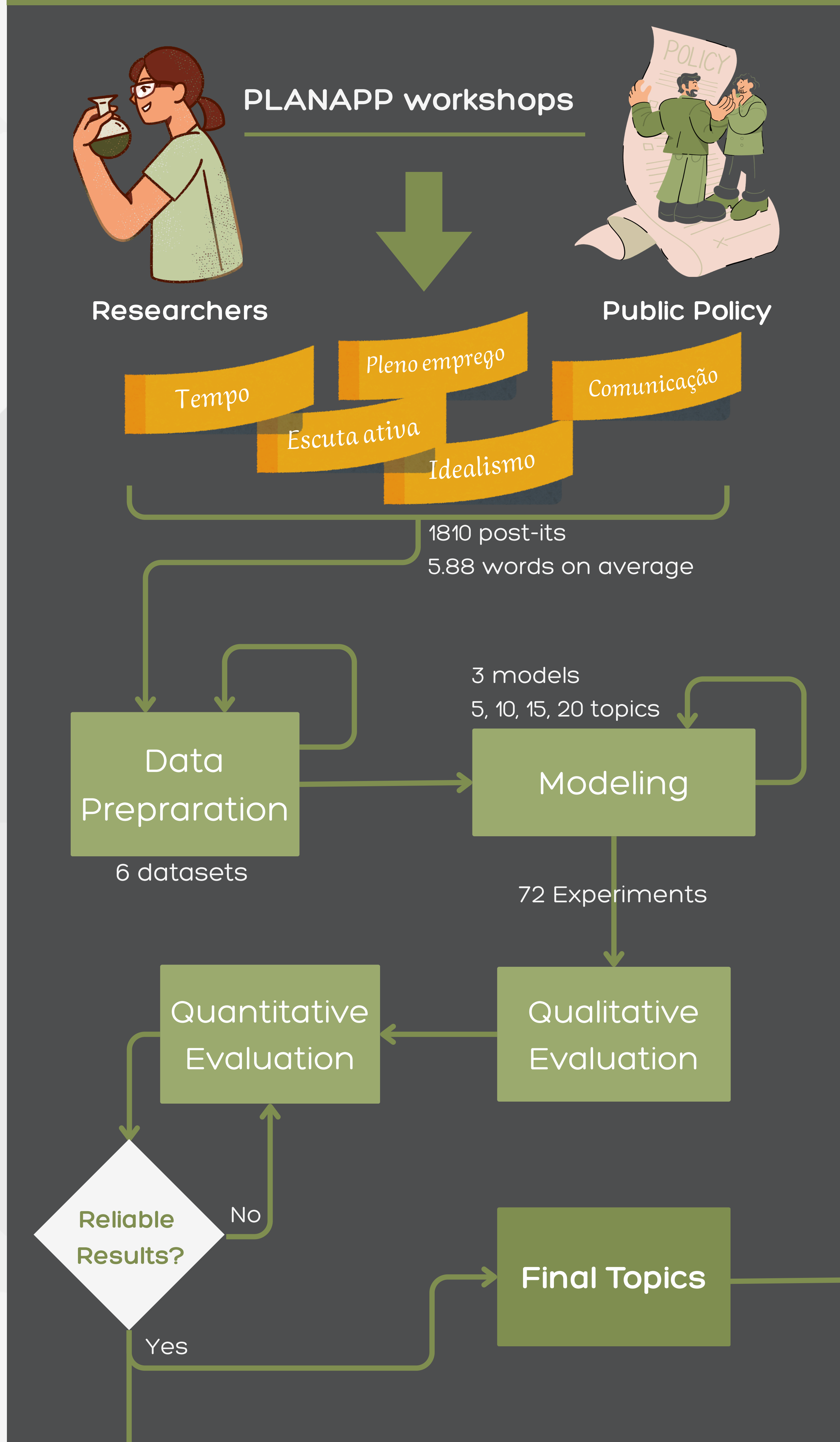
This study explores topic modeling on very short texts, a challenging task by the lack of reliable evaluation metrics. It tests various preprocessing strategies and modeling techniques to identify the most effective approach.
A new metric, the singularity score, is proposed to assess topic quality.

**Student**
Nicole Nunes - Nicole_Nunes@iscte-iul.pt

**Supervisor**
Ana Maria de Almeida - Ana.Almeida@iscte-iul.pt

**Supervisor**
Ana Rita Peixoto - Rita_Peixoto@iscte-iul.pt

## Project Workflow



PLANAPP workshops

Researchers → Public Policy

Tempo · Pleno emprego · Escuta ativa · Comunicação · Idealismo

1810 post-its
5.88 words on average

3 models
5, 10, 15, 20 topics

**Data Prepraration** → **Modeling**

6 datasets

72 Experiments

**Quantitative Evaluation** ← **Qualitative Evaluation**
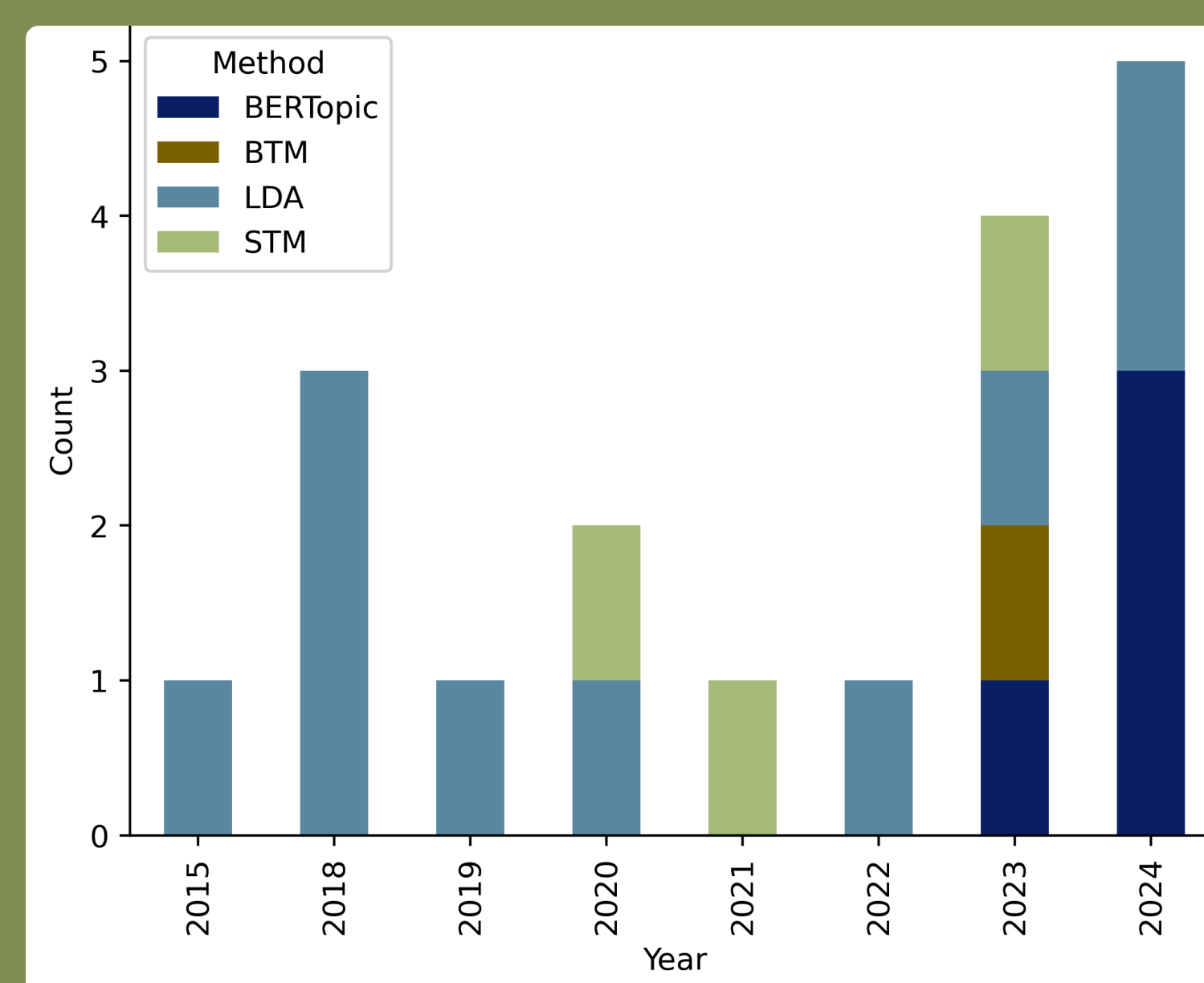
Reliable Results? — No → 
Yes → **Final Topics**

## Related Work

How do other authors deal with having a two-language corpus?
Either **translating** the documents written in the language with the least occurrence OR employing **models designed for multilingual texts** such as BERTopic multilingual.
How do other authors deal with suggestion text (tiny and informal)?
Avoid the problem by **removing documents with less than a minimum** of words OR apply **models more appropriate to this text length**, such as Biterm Topic Model, ST-LDA, Latent Dirichlet Allocation and BERTopic.



Topic modeling method evolution by year

## Experiment with highest SS – 0.998

**Dataset used**: Translated with text normalization and removal of stop words
**Model**: BERTopic with AlBERTina sentence-transformer



Language and communication · Politics · Research and projects · Goals and deadlines · Science, politics and gaps in knowledge transfer

## Singularity Score

- E**mulate the behaviour of annotators**.
- Based on the **stem of the top 10 words** of each topic.
- Significant Word (SW) $\in \{0, 1\}$;
- Count of Unique Words (UW) $\in [0, 10]$;
- Count of Non-Unique Words (NUW) $\in [0, 10]$.

$$tu_i = w_{SW}SW_i + w_{UW}\frac{UW_i}{10} + w_{NUW}\left(1 - \frac{NUW_i}{10}\right) \quad \text{Where } w_{SW} + w_{UW} + w_{NUW} = 1$$

- For topics with ST (strong topics) a reward is applied.
- Tetha is the threshold and beta the bonus.

$$ST = \frac{number\ of\ topics\ with\ tu \geq \theta}{N}$$

$$f(ST) = \begin{cases} 0 & if\ ST \leq 0.5 \\ \beta ST & if\ ST > 0.5 \end{cases}$$

Singularity Score is given by:

$$SS = TU(1 + ST)$$

## Conclusions

- Dealing with tiny text is extreamely challenging
- Traditional metrics are insufficient
- Singularity Score is proposed
- Future work: Validate Singularity Score with classical datasets

**Affiliations**
Instituto Universitário de Lisboa (ISCTE-IUL), ISTAR, Lisboa, Portugal

istar _iscte

PRR Plano de Recuperação e Resiliência · REPÚBLICA PORTUGUESA · Financiado pela União Europeia NextGenerationEU