

MARTA: LEVERAGING HUMAN RATIONALES FOR EXPLAINABLE TEXT CLASSIFICATION

INES AROUS¹, LJILJANA DOLAMIC², JIE YANG³, AKANSHA BHARDWAJ¹, GIUSEPPE CUCCU¹, PHILIPPE CUDRÉ-MAUROUX¹

¹{INES.AROUS, AKANSHA.BHARDWAJ, GIUSEPPE.CUCCU, PCM}@UNIFR.CH, ²LJILJANA.DOLAMIC@ARMASUISSE.CH, ³J.YANG-3@TUDELFT.NL

1. MOTIVATION AND CHALLENGES

Motivation: Explainability is a key requirement for text classification in many application domains including medical diagnosis and legal reviews.

Challenges: Existing explainable models rely on attention mechanisms, but:

- Attention-based models often provide an **inaccurate explanation** of text classification.

Using human rationales is challenging because:

- Workers provide **different rationales** for the same text classification.
- **Quality** of rationales provided by workers depend on their **reliability**.

2. PROBLEM DEFINITION

Given

- A set of textual documents
- Labels and rationales provided from workers of some documents

Our goal is to:

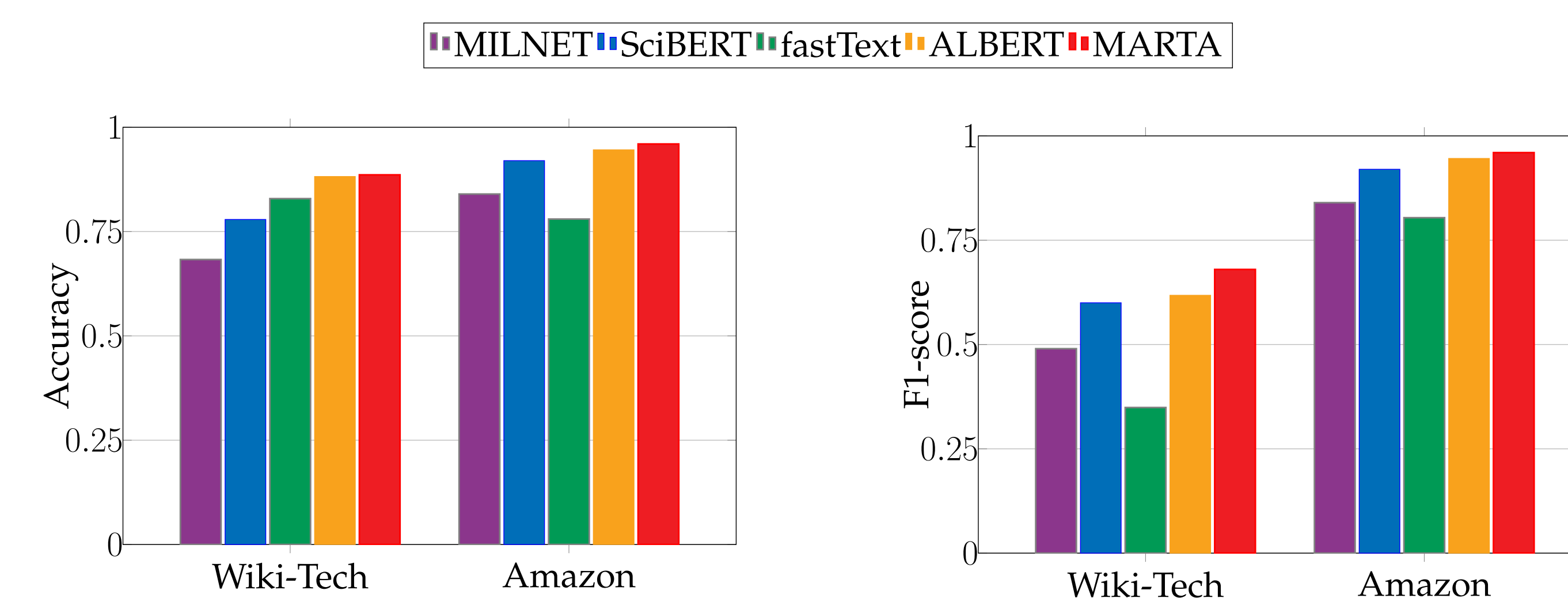
- Infer **the true label** of each document
- Estimate **the importance of each sentence** in a document.

3. CONTRIBUTIONS

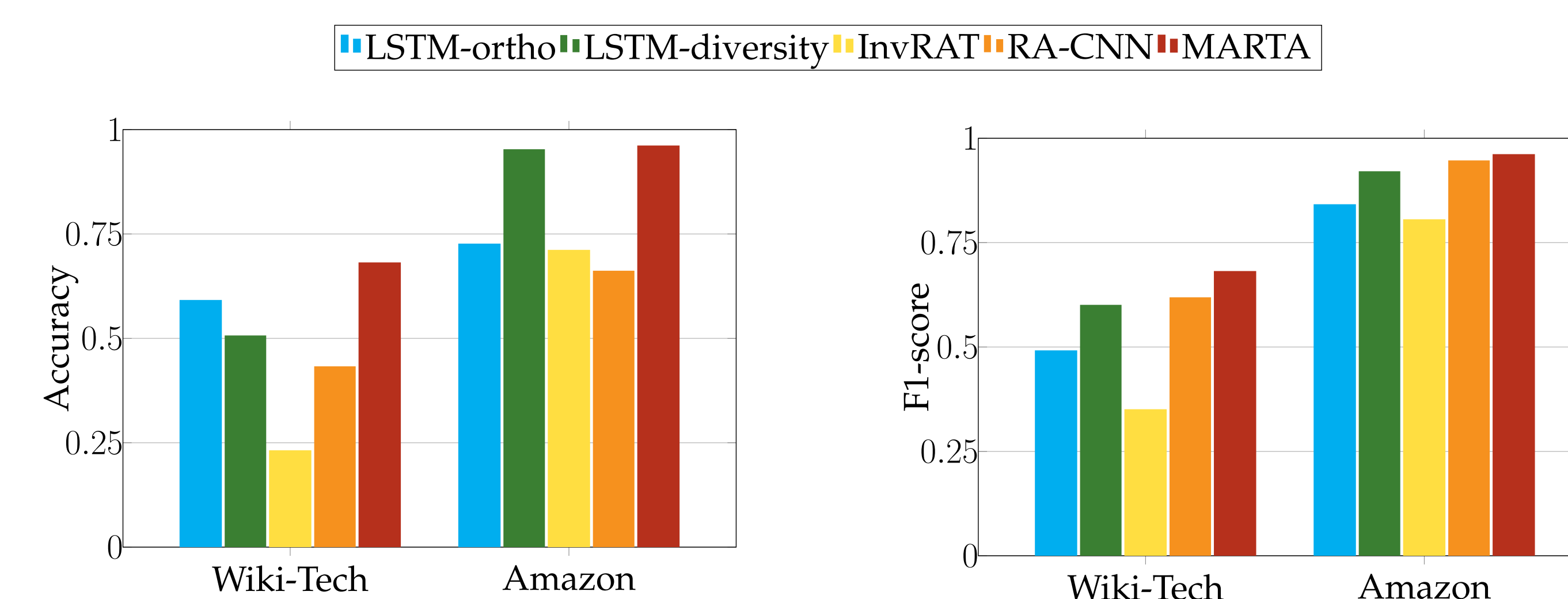
- A Bayesian framework called MARTA for Mapping human Rationales To Attention.
- A Variational Inference algorithm with incremental updating rules.

5. RESULTS

- MARTA improves **text classification methods** by **0.97% accuracy** and **5.76% F1-score**.

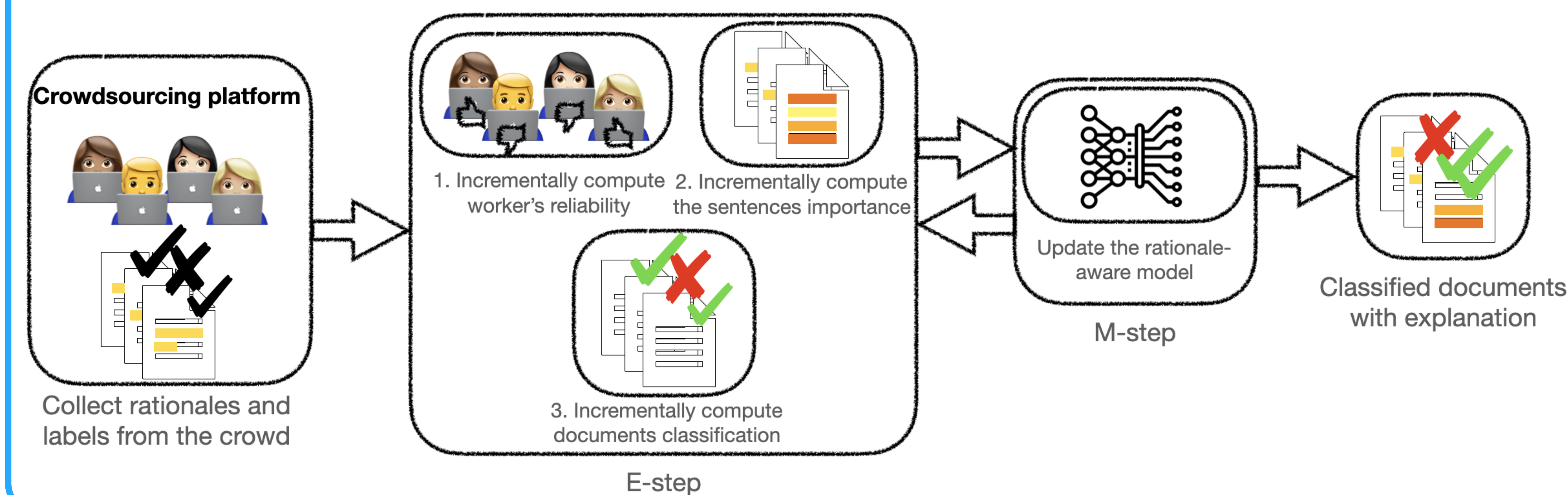


- MARTA improves **rationale-aware methods** by **18.68% accuracy** and **17.61% F1-score**.



- MARTA is better at leveraging rationale by weighing workers reliability.

4. THE MARTA FRAMEWORK



6. PROPERTIES

- **Explainability:** The overlap between the rationales chosen by workers and those highlighted by MARTA is **66%** on average on two datasets.
- **Adjustable Supervision Degree:** MARTA requires a small amount (max 30%) of ground truth labels for training.
- **Robustness:** The standard deviation is small (~ 0.03) compared to the absolute accuracy which demonstrates MARTA's robustness.

7. CONCLUSION

- We proposed MARTA for explainable text classification that integrates **human rationales** into an **attention-based model**.
- As future work, we plan to take into account worker's rationales expressed in a syntax different from the original text.