

Introduction

Personal values are the abstract motivations that drive our opinions and actions. Using state-of-the-art NLP methods, we design a classifier to study their expression in text.

Moral Foundations Theory¹ (MFT) proposes five "irreducible basic elements" of morality, that we can frame our study in: *care/harm, authority/subversion, fairness/cheating, loyalty/betrayal, purity/degradation*.

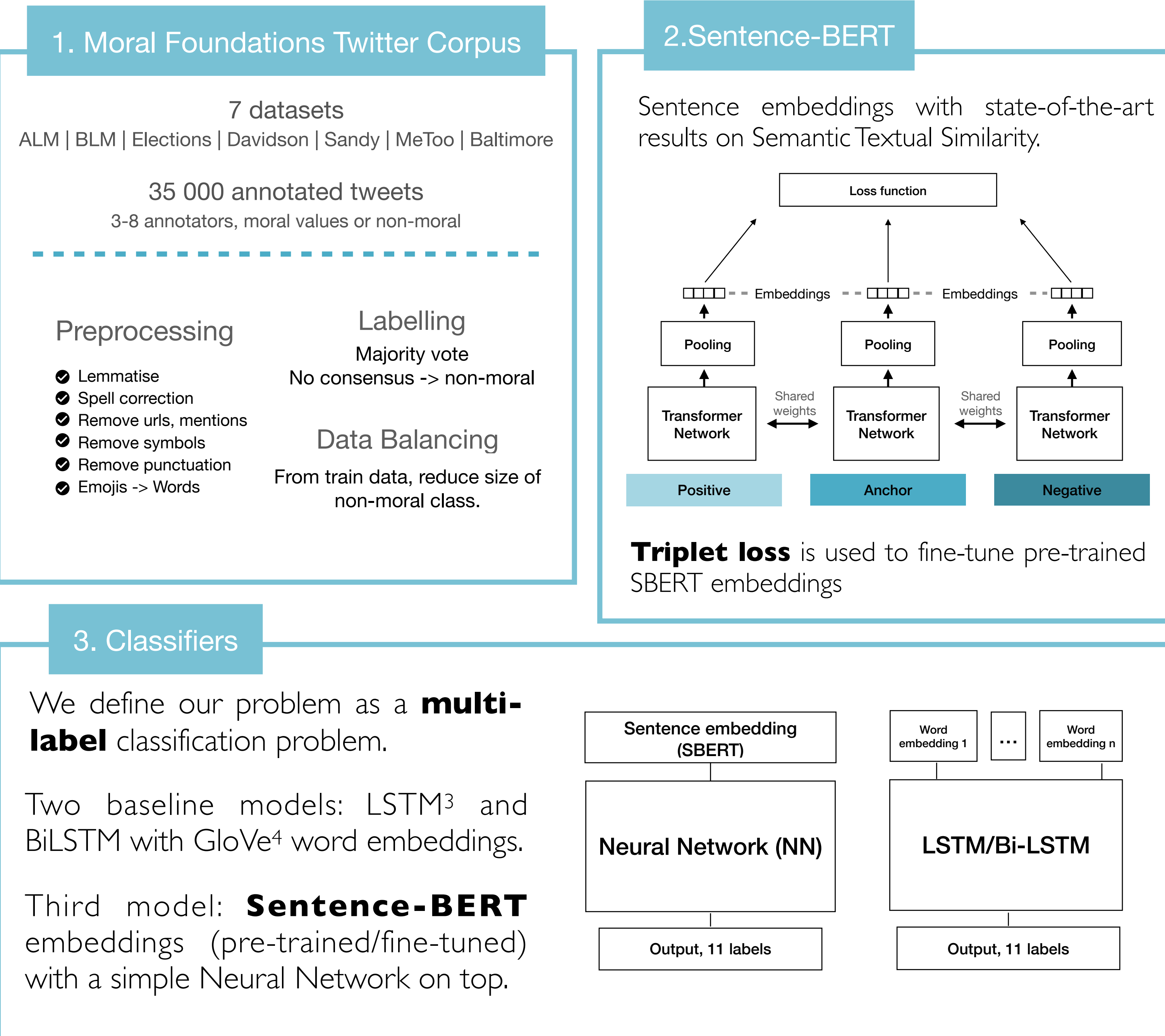
Embeddings convert word and sentences to meaningful vectors and they are an important step in a text classifier's pipeline. They can be domain-adapted to improve the model's performance.

Research Goal: Train embeddings to learn moral foundations and assess our method by answering three research questions:

- 1) Does our fine-tuning method increase the moral classifier's **performance**.
- 2) Do fine-tuned embeddings **generalise** across domains of discourse.
- 3) Are fine-tuned embeddings **transferable**.

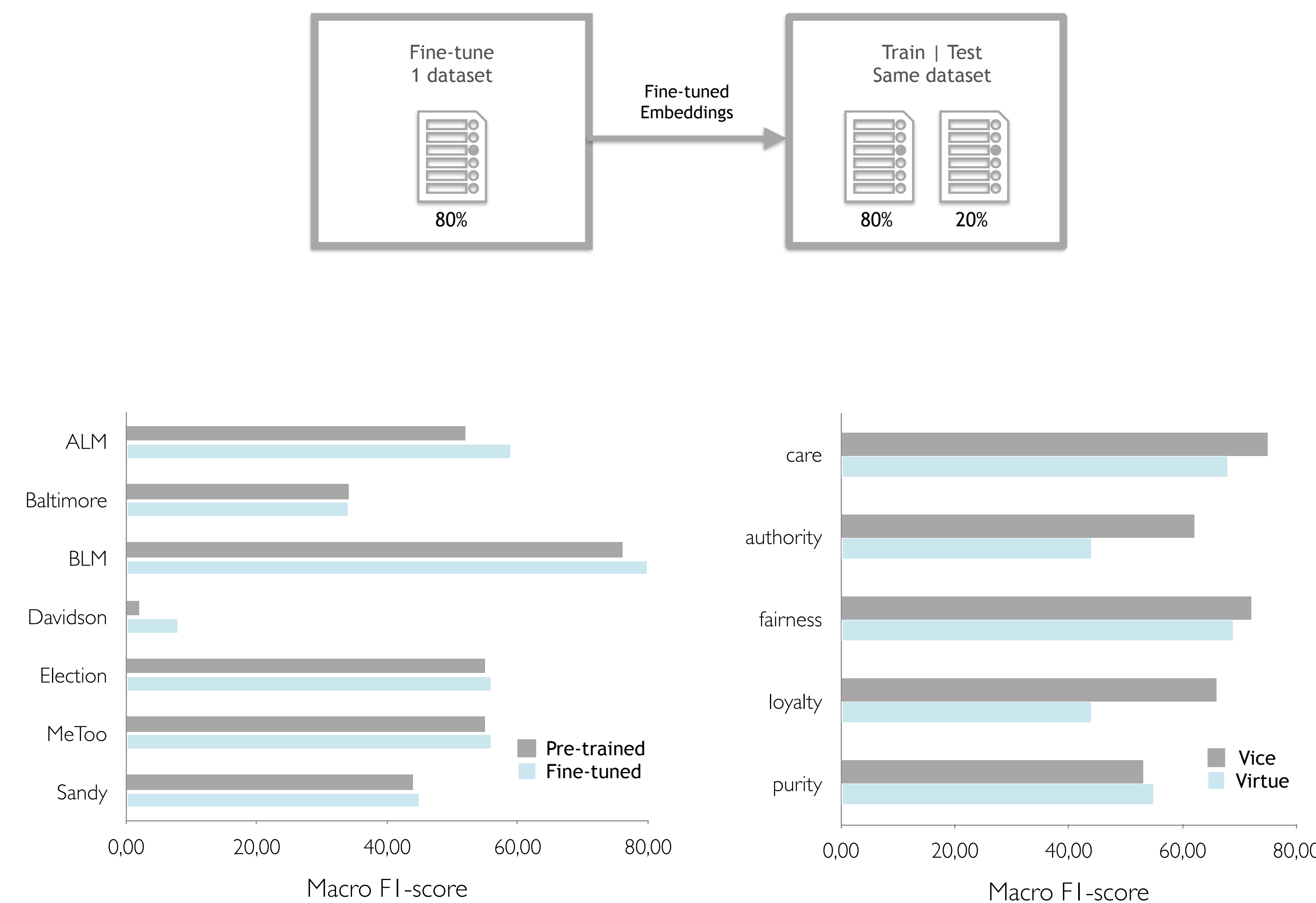
Motivation: no prior moral classifiers focus on fine-tuning state-of-the-art embeddings (Sentence-BERT²) to improve the model's performance. Moreover, after training, embeddings' utility is not limited to the classification task: Semantic Textual Similarity, clustering.

Methods



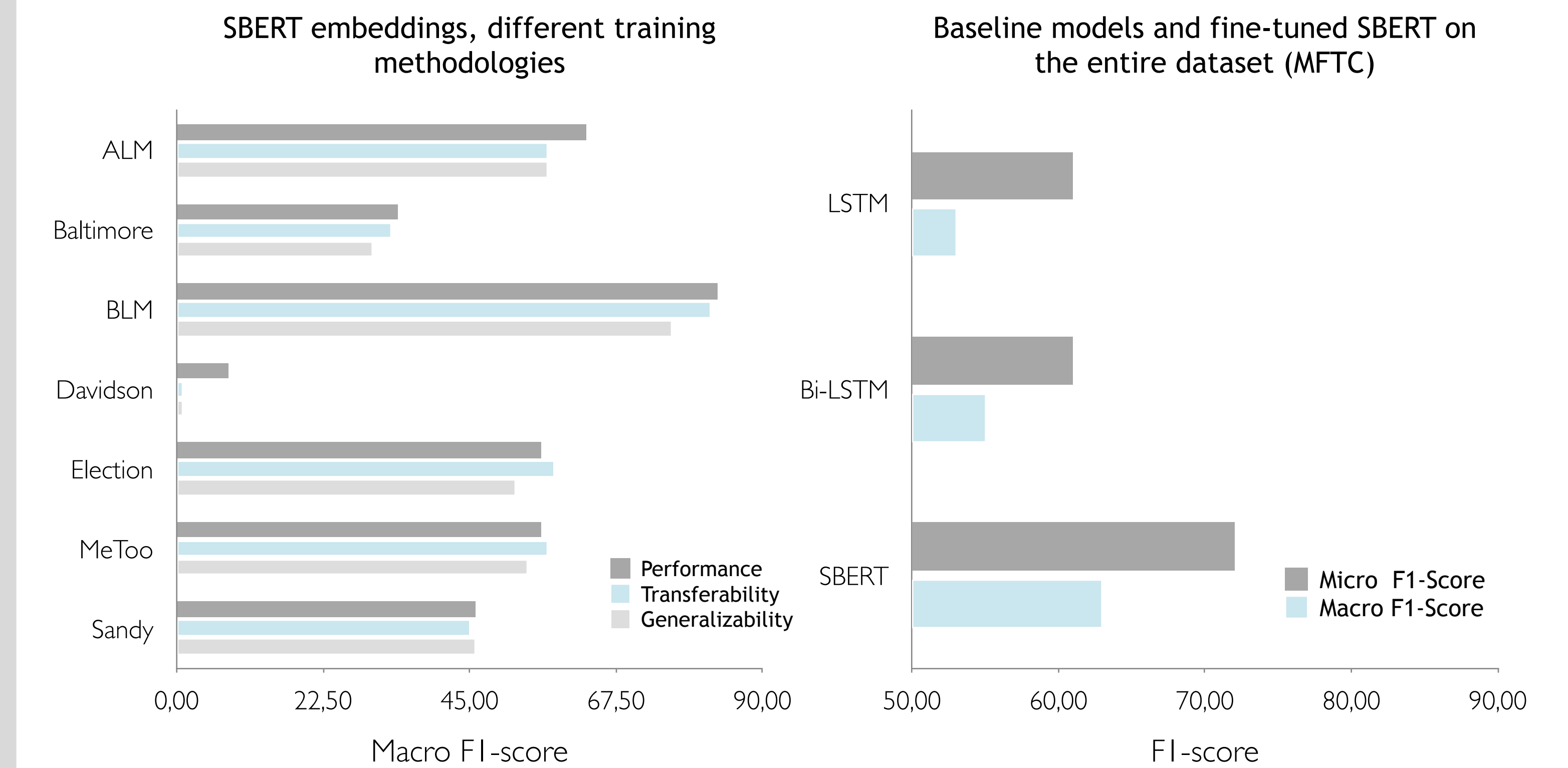
Performance

1) Does fine-tuning Sentence-BERT embeddings with triplet loss increase the moral classifier's **performance**.



The chart on the left shows how fine-tuning SBERT improves the Macro F1-Scores for the moral classifiers. On the right, we illustrate how embeddings trained on the entire MFTC recognise each moral value.

Results



Discussion

For the moral classification task, we proposed a method to fine-tune state-of-the-art embeddings. The resulting classifier achieves **72% Micro F1-score** on the MFTC dataset.

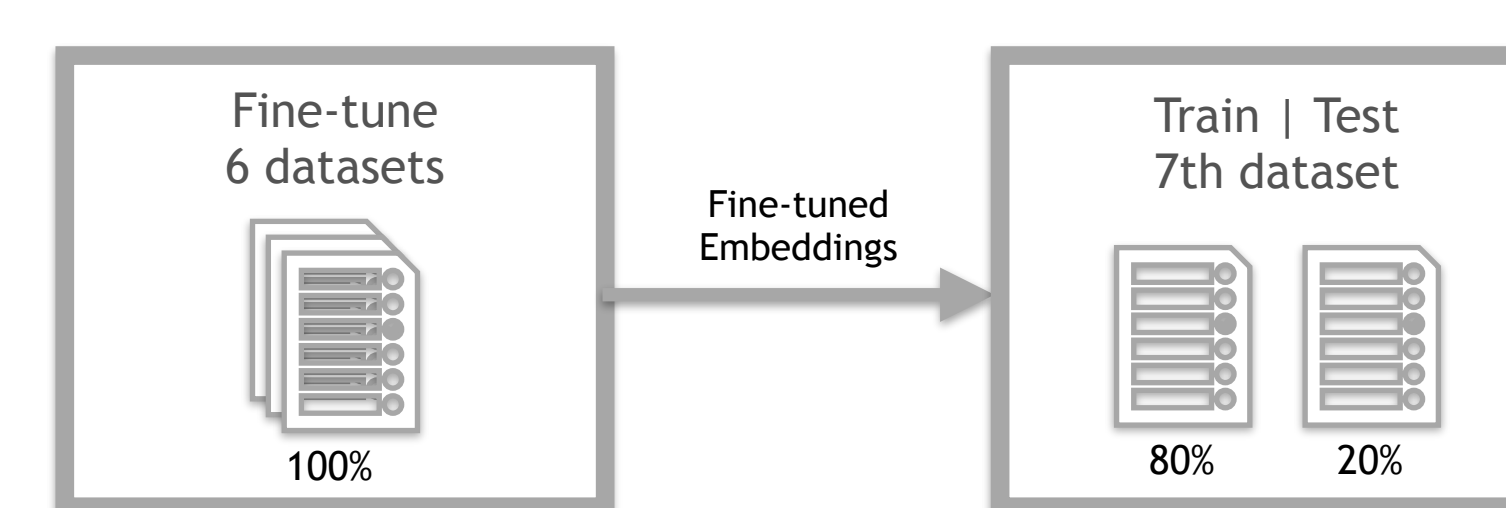
Future work

For a complete understanding of moral embedding's transferability, MFTC should be extended. As MFT annotating is labour intensive, we recommend experimenting with semi-supervised annotating methods⁵.

To better explain our method's success, it should be investigated if *semantically similar text expresses similar moral values*.

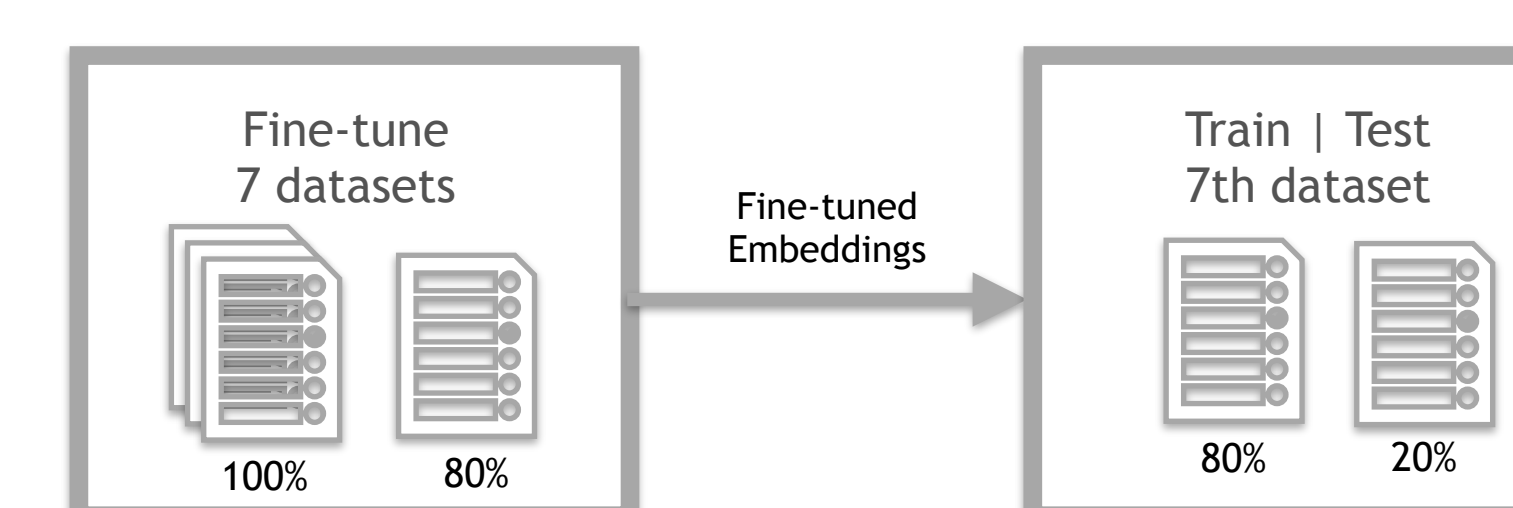
Generalisability

2) Do fine-tuned embeddings **generalise** across domains of discourse.



Transferability

3) Are fine-tuned embeddings **transferable**.



References

1. Graham, J. (2013). Moral Foundations Theory. *Advances in Experimental Social Psychology*, 47, 55-130.
2. Reimers, N. (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks
3. Hochreiter, S. (1997). Long short-term memory. *Neural Computation*, 9, 1735-1780
4. Pennington, J. (2014). Glove: Global vectors for word representation
5. Settles, B. (2011). Closing The Loop: Fast, Interactive Semi-Supervised Annotation with Queries on Features and Instances.