SetFit for Automated Essay Scoring

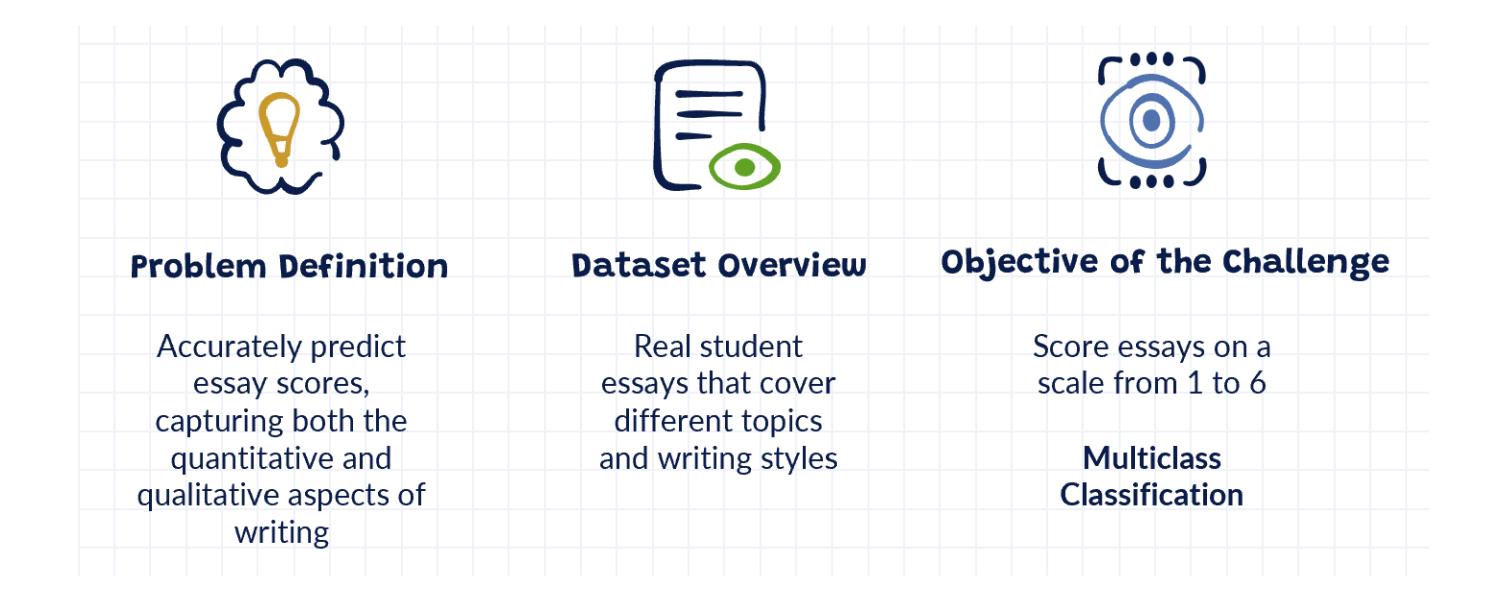
Leon Krug, Jannik Bundeli, Jannine Meier, Elena Nazarenko

Lucerne University of Applied Sciences and Arts

HOCHSCHULE LUZERN

Informatik

The Objective



What is Setfit?

SetFit (Set-Factory for Few-shot Text Classification) is a powerful text classification framework designed to work well even with very few labeled examples (few-shot learning). It combines Sentence Transformers with efficient fine-tuning techniques to deliver strong performance using minimal training data.

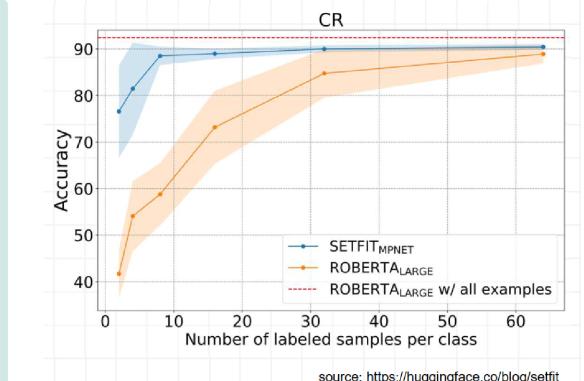
Key Features:

Few-shot capable: Trains models with as few as 8–64 labeled samples.

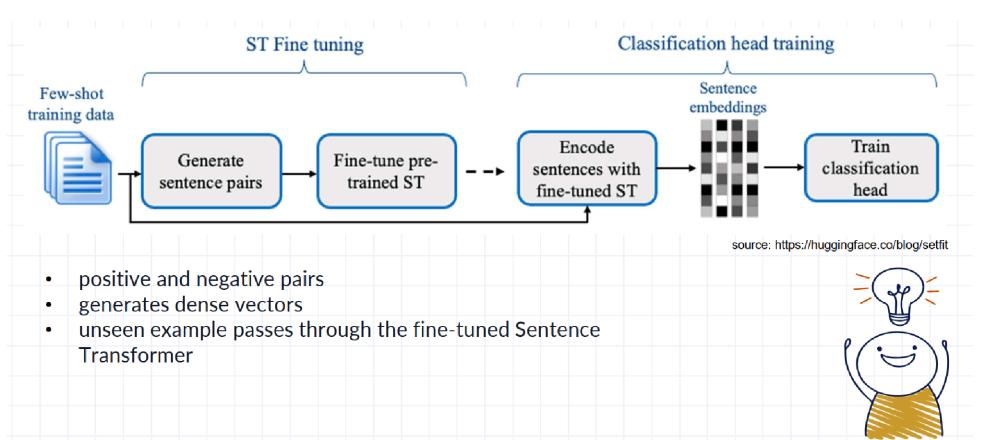
No need for large GPUs: Very efficient, often runs on CPU.

Fast fine-tuning: Trains in seconds to minutes.

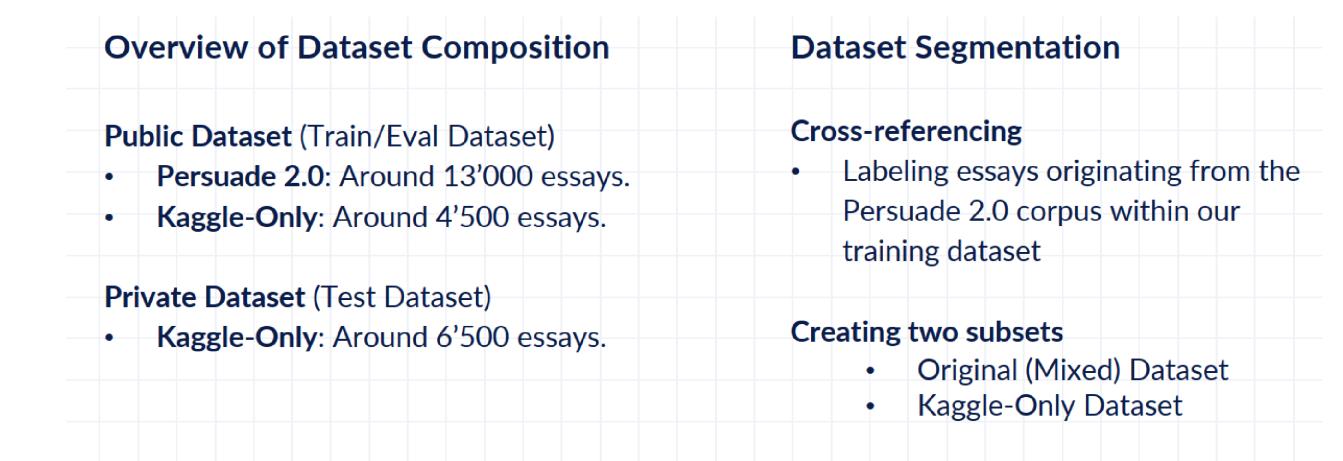
High accuracy: Competes with full fine-tuning methods on many benchmarks.

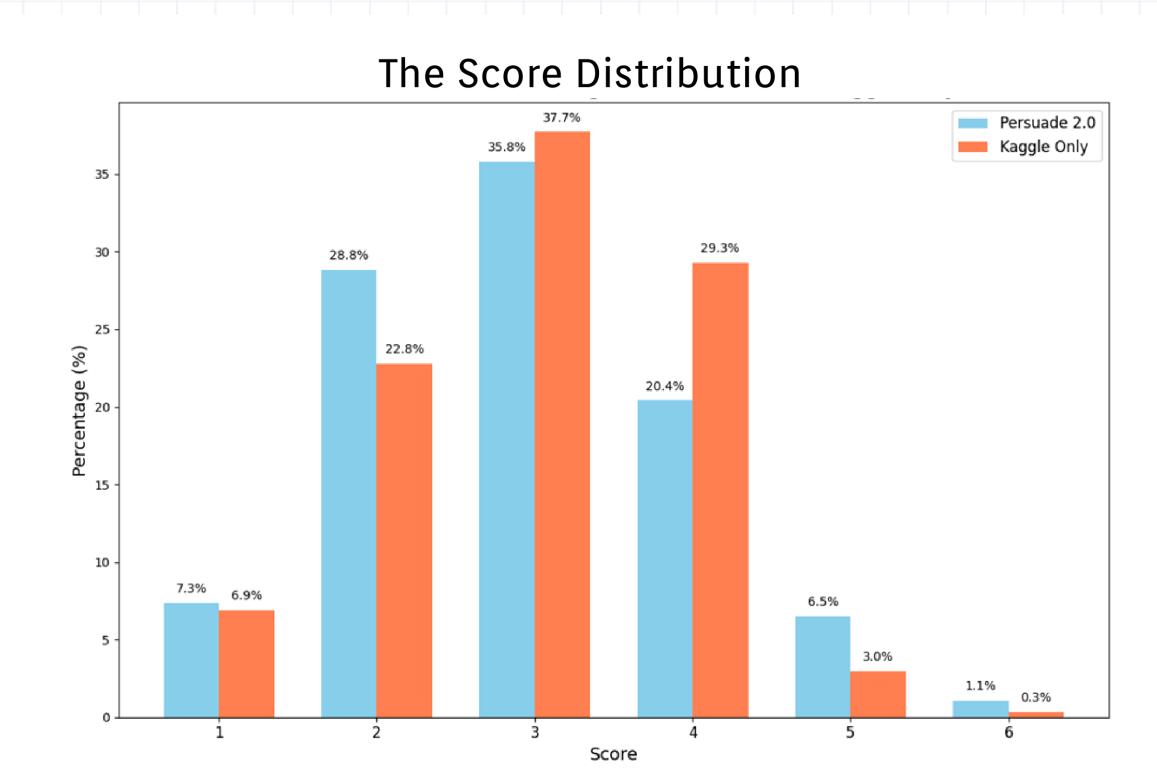


SetFit's two-stage Training



The Dataset





The Score Description

Clear mastery with few errors, outstanding critical thinking, appropriate evidence, well-organized, skilled language use.

Reasonable mastery with occasional errors, strong critical thinking, generally appropriate evidence, well-organized, good language use.

Adequate mastery with some lapses, competent critical thinking, adequate evidence, generally organized, fair language use.

Developing mastery with weaknesses, limited critical thinking, inconsistent evidence, limited organization, fair language use with weaknesses.

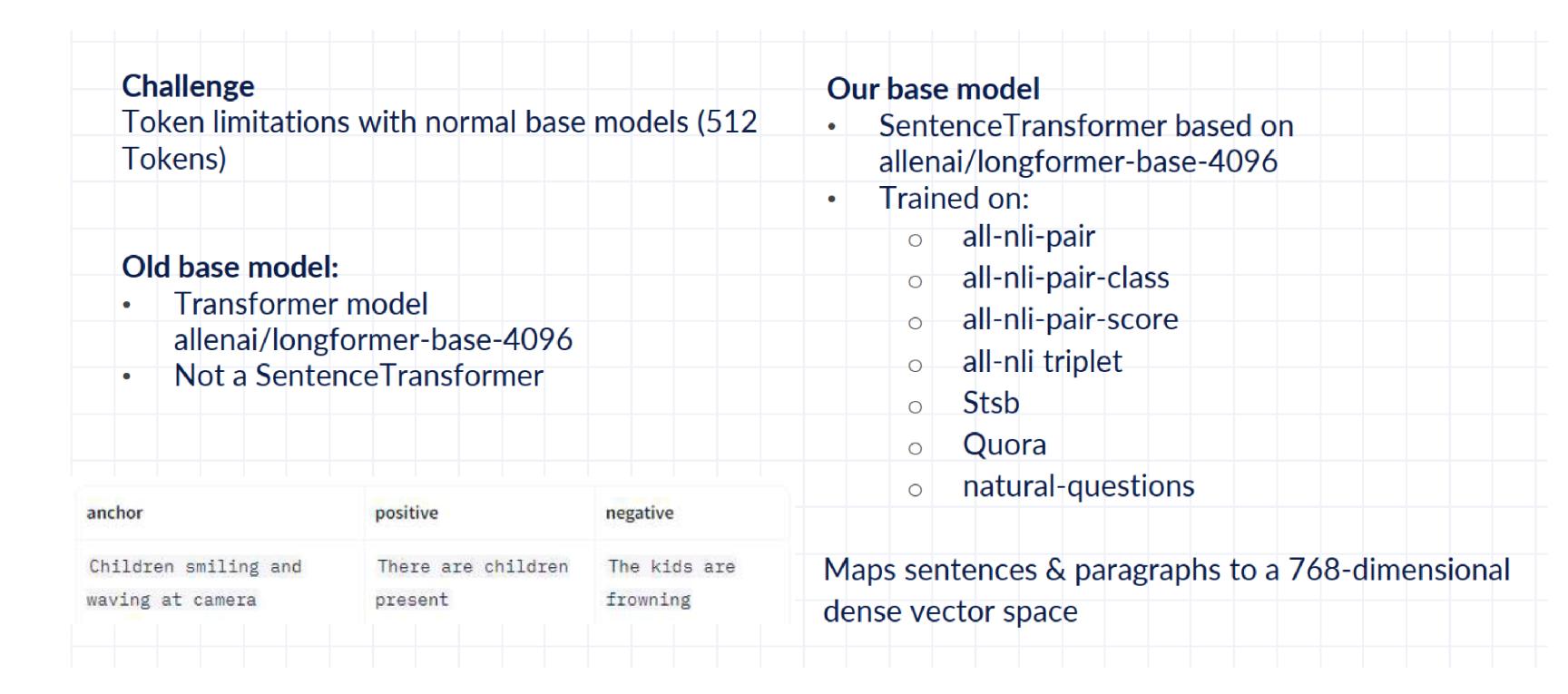
Little mastery with serious flaws, weak critical thinking, insufficient evidence, poor organization, limited language use with frequent errors.

Very little or no mastery, severely flawed, no viable point of view, disorganized,

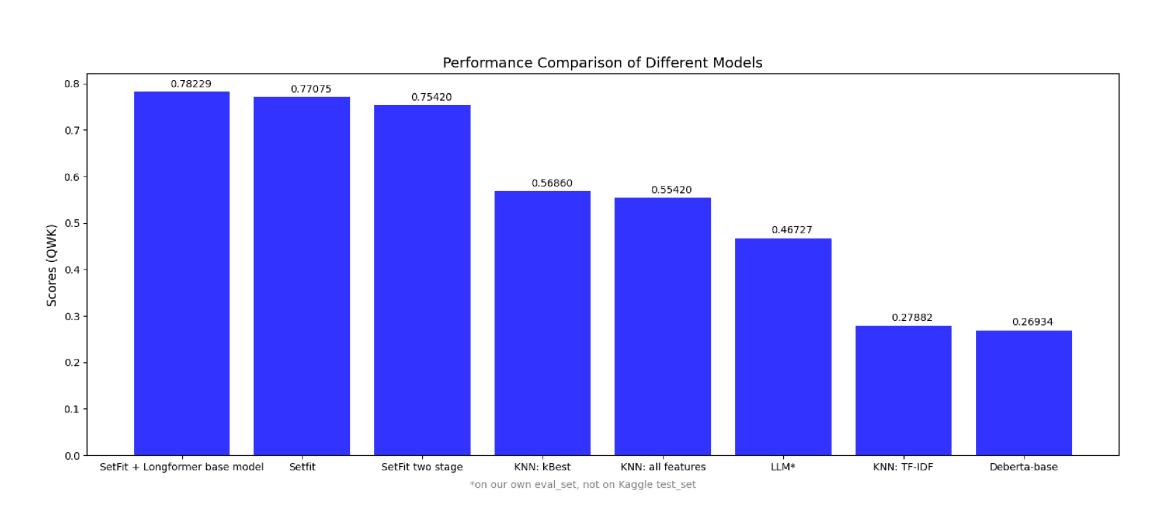
fundamental language flaws, pervasive grammar/mechanics errors.

Our Innovation

SetFit + Longformer Sentence Transformer



The Results



Quadratic Weighted Kappa Metrics (QWK)

- Measures agreement between predicted and actual scores
- Weighted: Penalizes larger disagreements

Range: -1 to 1

1: Perfect agreement
0: Random agreement
< 0: disagreement

- 0: disagreement

- 0.4-0.6: Moderate agreement

- 0.4: Poor agreement

Longformer as a Sentence Transformer
Old context size: 512 Tokens
New context size: 4096 Tokens
8x longer context





8 essay pages (4096 tokens)