

On the Automatic Detection and Classification of Linguistic Bias

Richard Diehl Martinez, Sabri Eyuboglu



Dataset

Wikipedia Neutrality Corpus

50,000 biased sentences

Words responsible for bias are tagged

Schnabel himself did a **fantastic** reproduction of Basquiat's work.

No bias type labels!



Hand-labeled Data

500 biased sentences labeled as framing or epistemological

Schnabel himself did a **fantastic** reproduction of Basquiat's work. → **Framing**

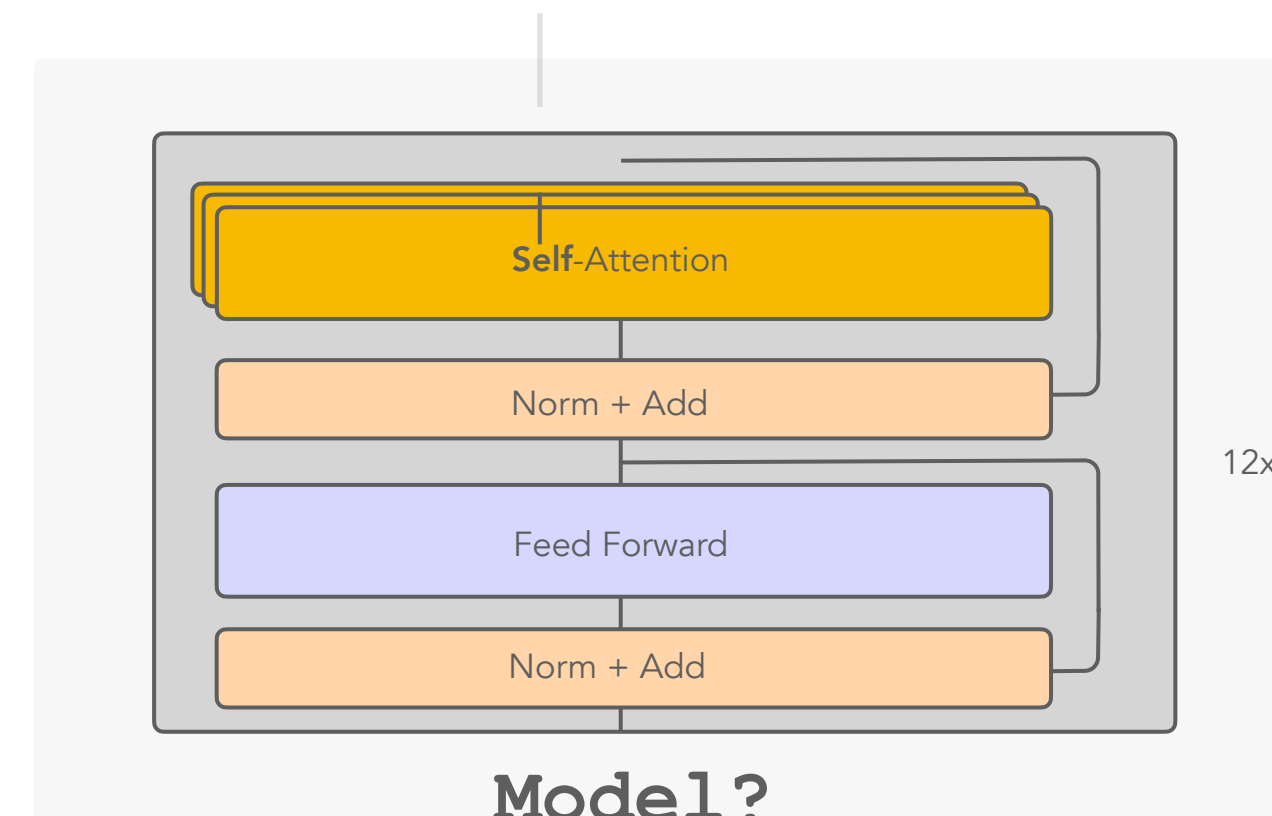
Lower cost of living means a **possibly** higher standard of living → **Epistemogical**

Challenge

✓ **Wikipedia Neutrality Corpus (WNC)**
n = 50,000

✓ **Bias Detection Model**
BERT Trained on WNC

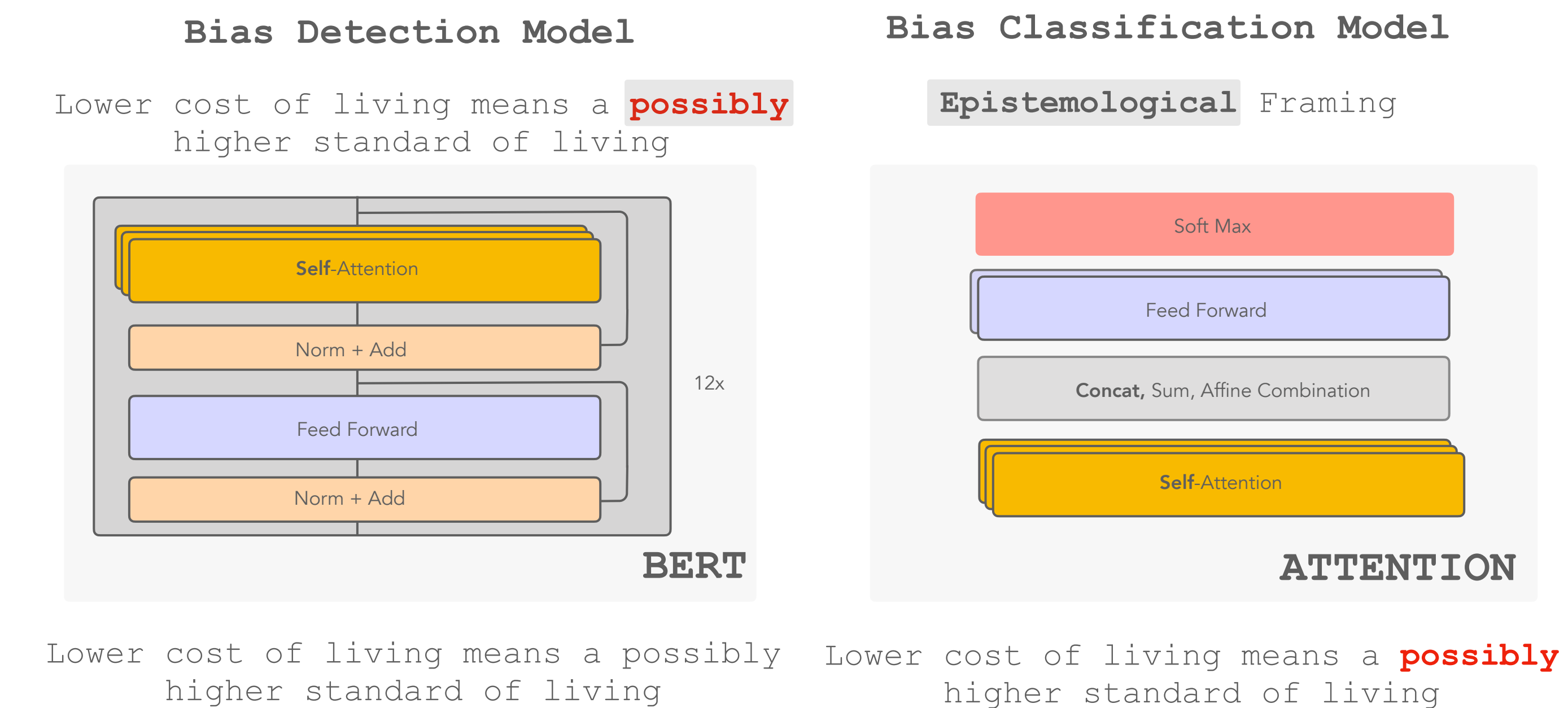
Schnabel himself did a **fantastic** reproduction of Basquiat's work.
Biased Sentence



✗ **Epistemological Bias** **Framing Bias**

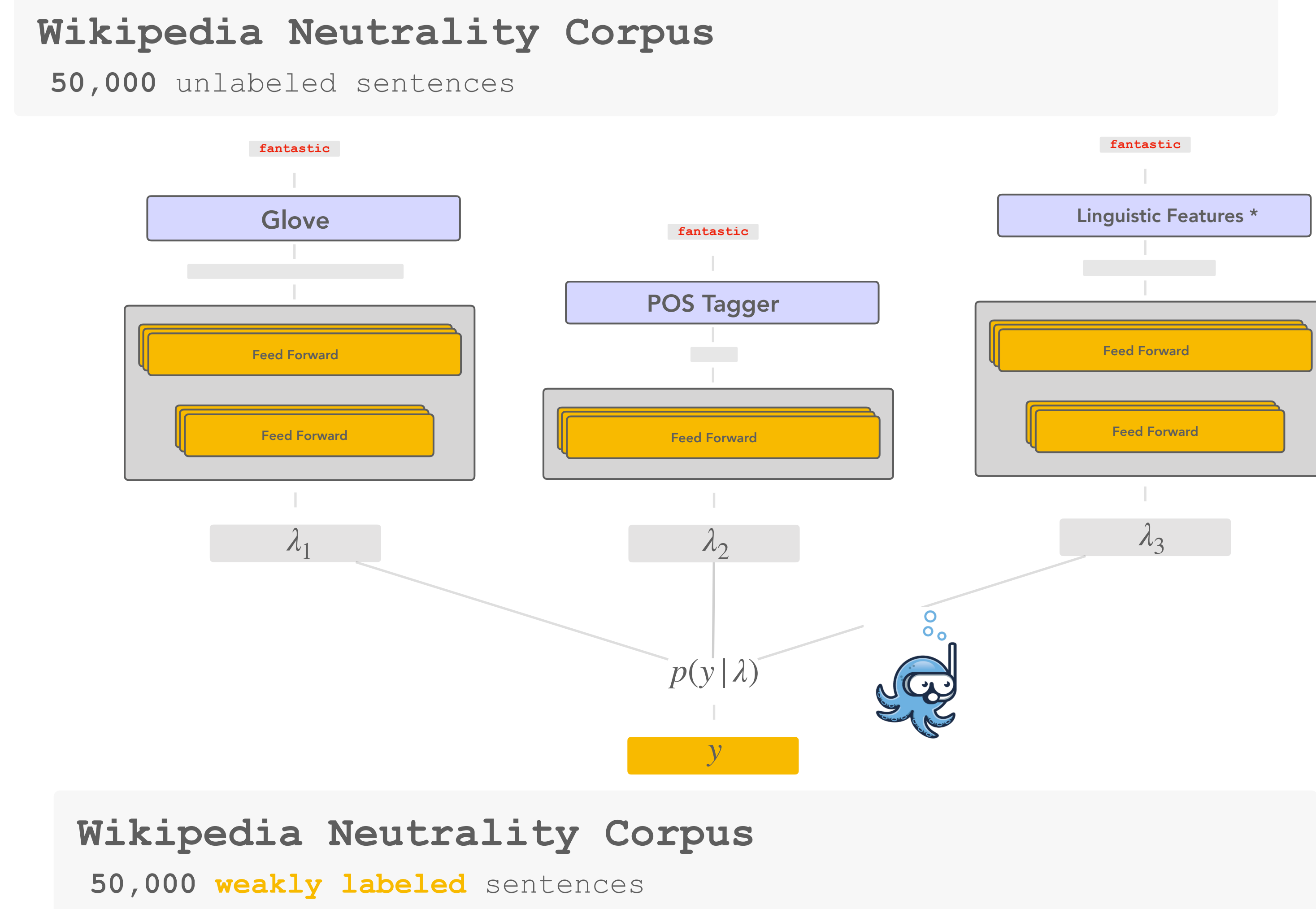
Principal Challenge: **LACK OF LABELS**

Training on Gold Labels



	BERT		Attention	
BERT out of Box Accuracy	0.573	0.606	0.679	
AUROC	0.618	0.608	0.756	

Training on Weak Labels



	BERT		Attention	
Accuracy	0.742	0.847	0.704	
AUROC	0.847		0.770	