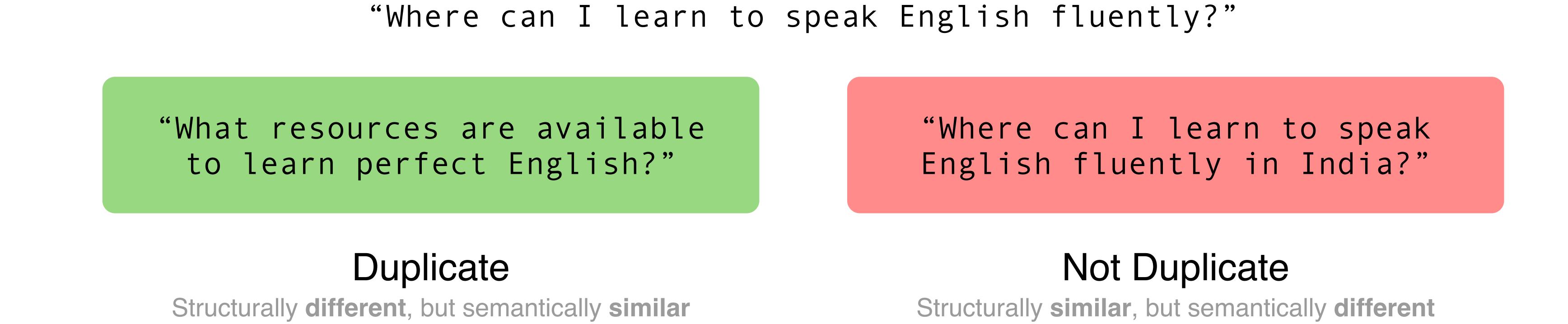
Identifying QuOTA Question Pair Duplicates Enzo Blindow, Tom Dop (Quora-The-Explorer)

1. The Problem

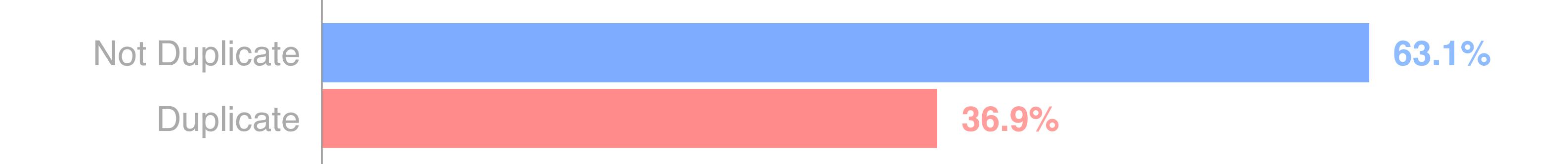


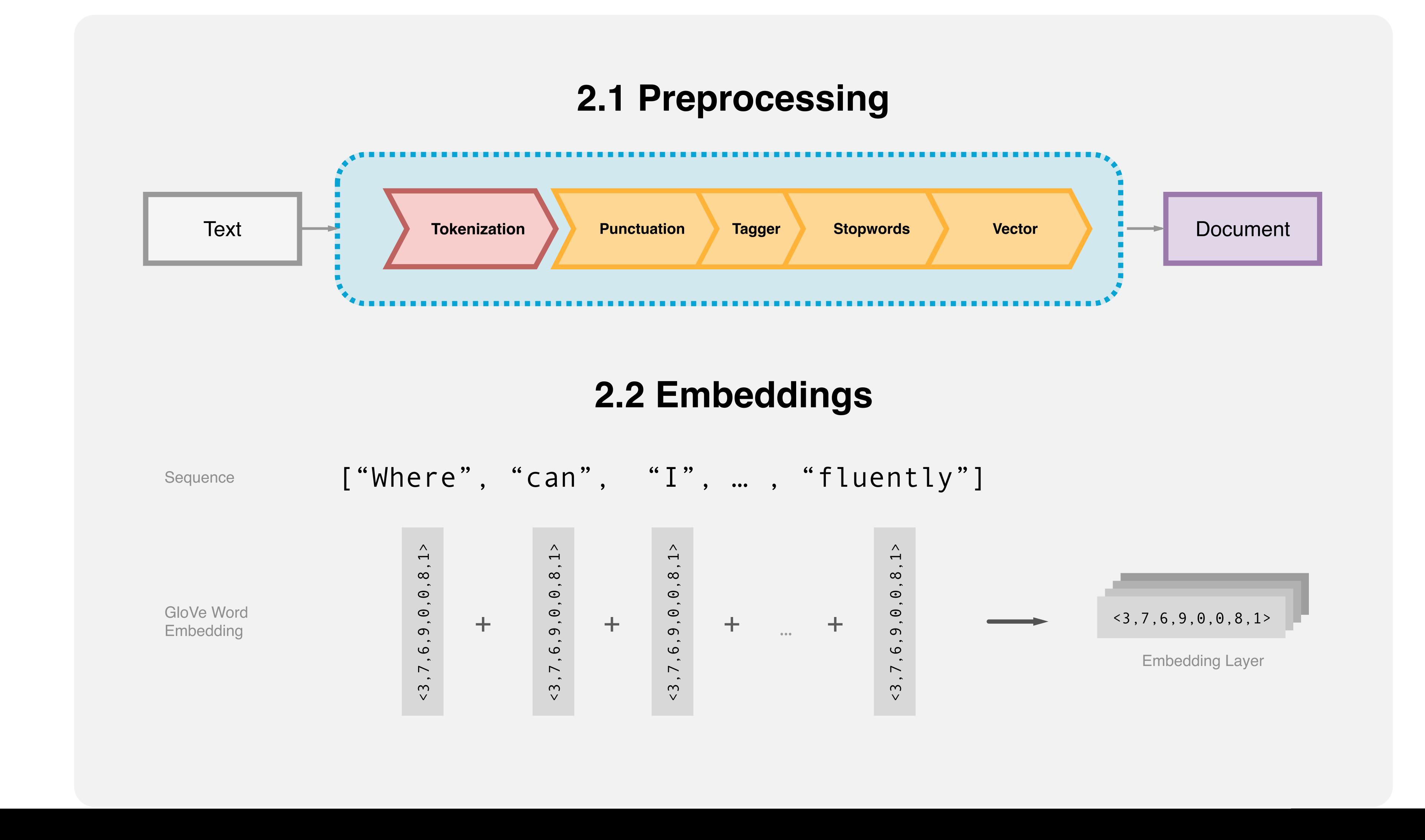
Our Ensemble Approach

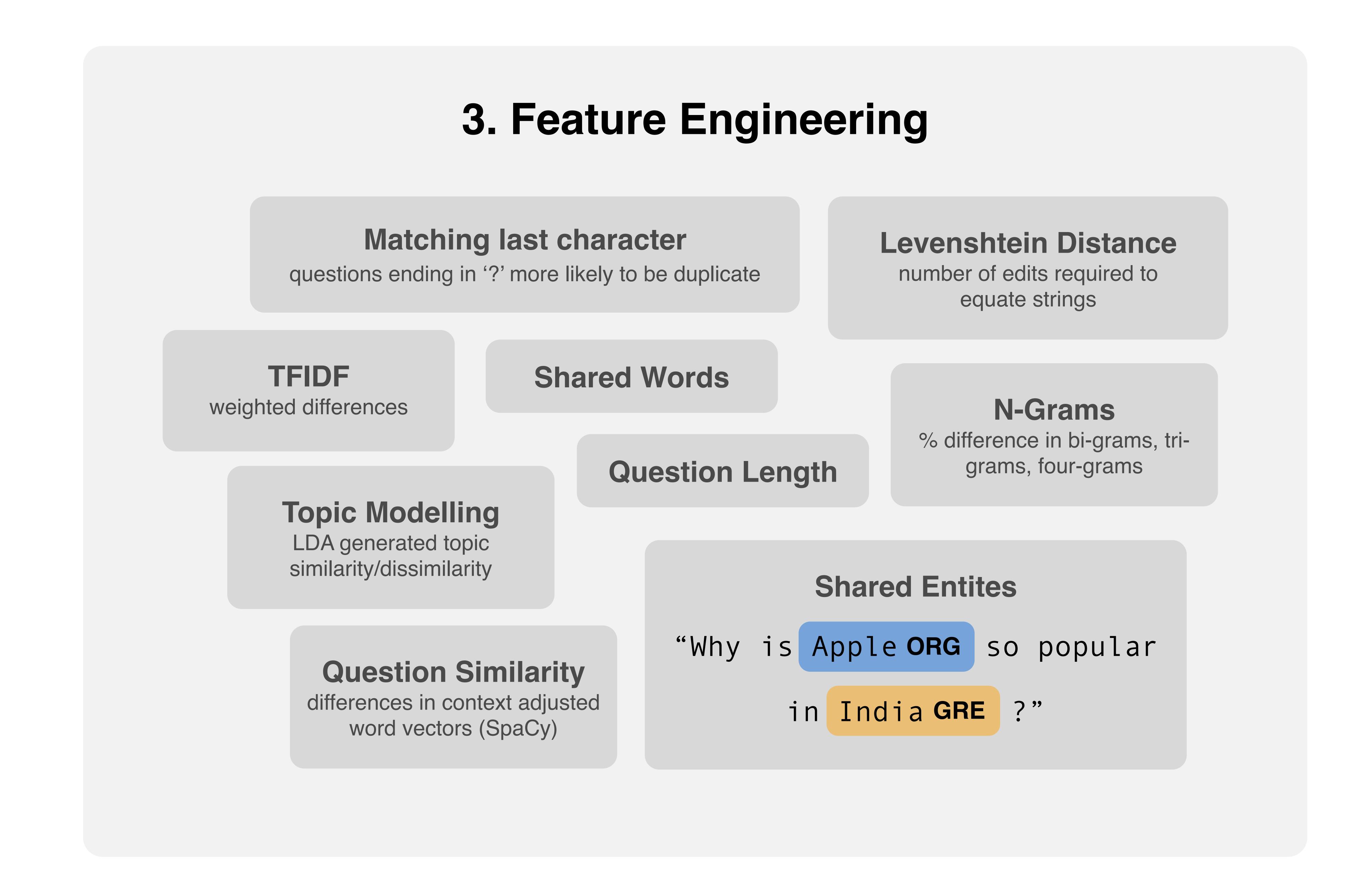
Utilizing a combination of LSTMs (Long-Short-Term-Memory) and NLP features, feeding into various classification algorithms to predict duplicate questions.

2 Data

Binary classification problem with imbalanced classes







4. RNN Network with LSTMs

While traditional NLP approaches only tend to take **lexical differences** into account, LSTM cells can account for the **ordering** of words and may be better for measuring **meaning**.

- Train / Crossvalidation Split: 20%
- Rebalancing of classes by undersampling
- · Siamese network (two inputs) seem to perform well on semantic similarity tasks
- We tried both subtracting and concatenating the outputs from the two LSTM layers

4.1 Hyperparameter Tuning

After 5 Epochs

Loss		OSS	Accuracy		
Description	Training	Validation	Training	Validation	
Base	0.19	0.61	92.0%	80.3%	
More Nodes	0.18	0.61	92.4%	81.9%	
Bidirectional	0.15	0.68	93.8%	81.5%	
Sigmoid	0.27	0.47	88.1%	80.9%	
No dropout	0.11	0.53	95.6%	83.7%	
Reweighted	0.33	0.59	88.1%	80.6%	
Reweighted, Low dropout	0.17	0.88	94.4%	81.6%	
Reweighted, No dropout	0.06	1.33	98.2%	82.2%	
Subtract	0.07	0.65	97.4%	84.1%	

Base Model 2 parallel LSTM layers (300 nodes), 2 Dense Layers (200 nodes), ReLu activation, 20% dropout

4.2 Model Performance 1 0.8 0.6 0.4 0.2 0.2 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30

5. Classification Results

Epochs

——Train Loss ——Validation Loss ——Train Accuracy ——Validation Accuracy

After combining outputs from LSTM and engineered features, we tried several different algorithms for classifying duplicate questions and compared results (Accuracy).

Logistic Regression 0.84585 Random Forest max depth 2 0.82639 SVM rbf kernel, c=1, gamma=1 (different set) 0.81001	Neural Network	3 Dense Layers 200 nodes each, 20% dropout	0.84515
	Logistic Regress	sion	0.84585
SVM rbf kernel, c=1, gamma=1 (different set) 0.81001	Random Forest	max depth 2	0.82639
	SVM rbf kernel, c=	:1, gamma=1	ferent set) 0.81001

6. Summary

- Although we achieved a reasonable final results (84.6% accuracy), our LSTM performance graph, suggests the model was overfitting
- While accuracy was used to judge the final result, precision and recall may have been a better measure of performance due to imbalanced classes
- Although we attempted to tune our model hyperparameters, we were limited by time and computational power and so were only able to test each model for 5 epochs