

Classification of Questions

Petr Lorenc 2017
FIT CTU Prague
12.5.2017

Why ?

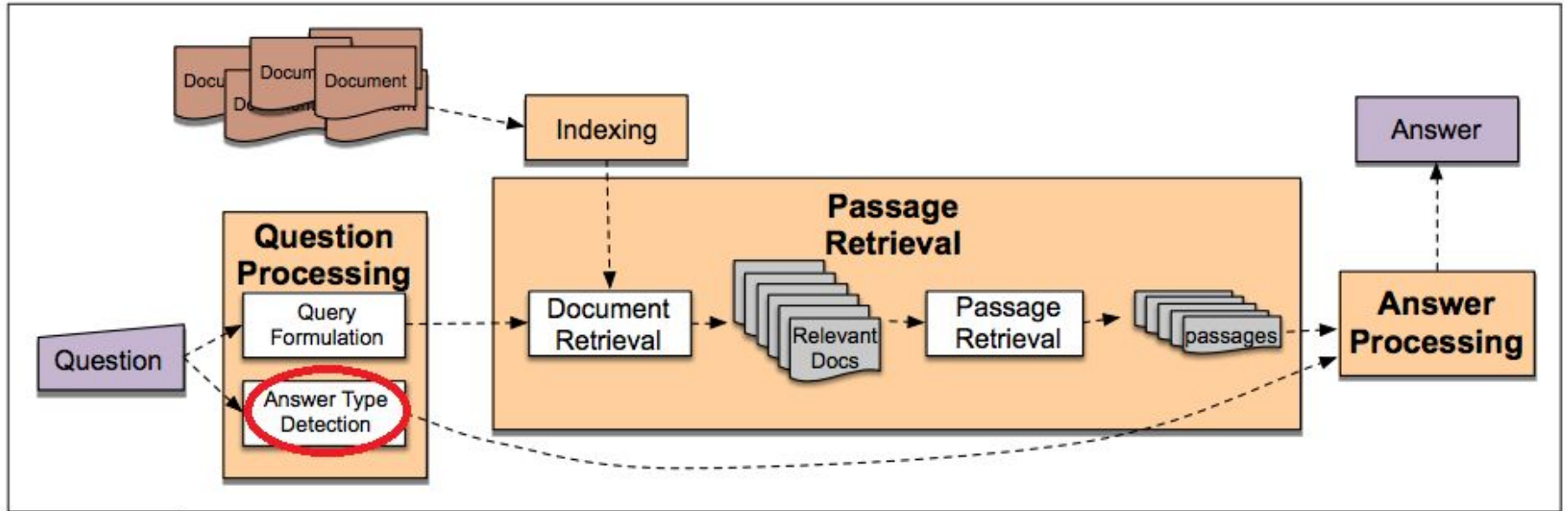
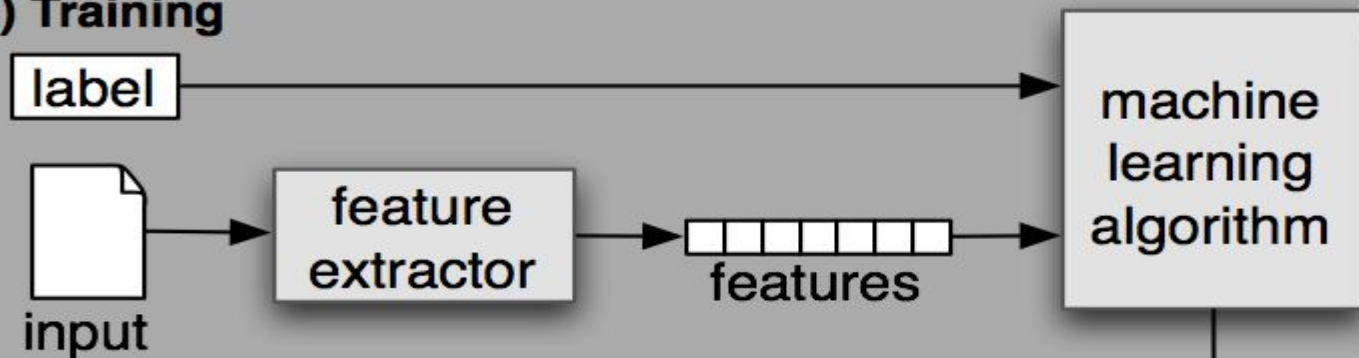


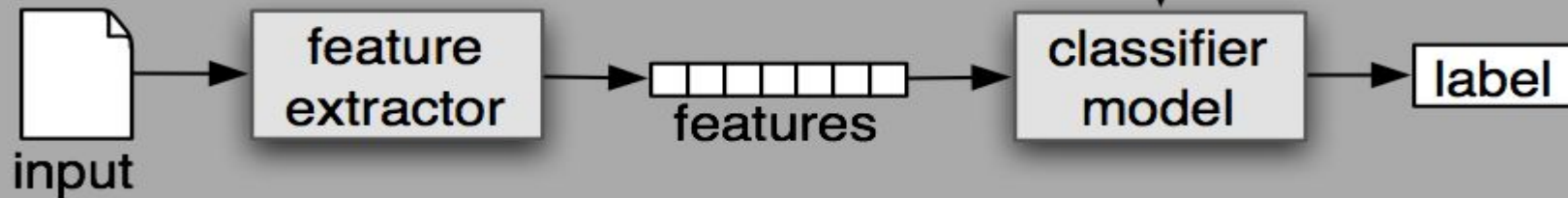
Figure 28.2 IR-based factoid question answering has three stages: question processing, passage retrieval, and answer processing.

How?

(a) Training



(b) Prediction



Representation of data

1. Tf-Idf
2. Bi-grams + Tf-Idf
3. Vector representation (with/without scale normalization)
 - a. Word2Vec
 - b. Glove
 - c. Doc2Vec
 - d. FastText

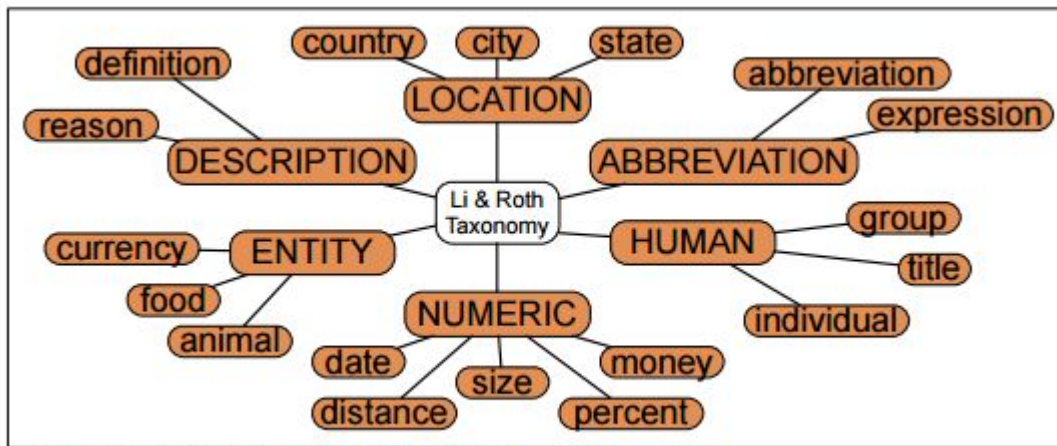
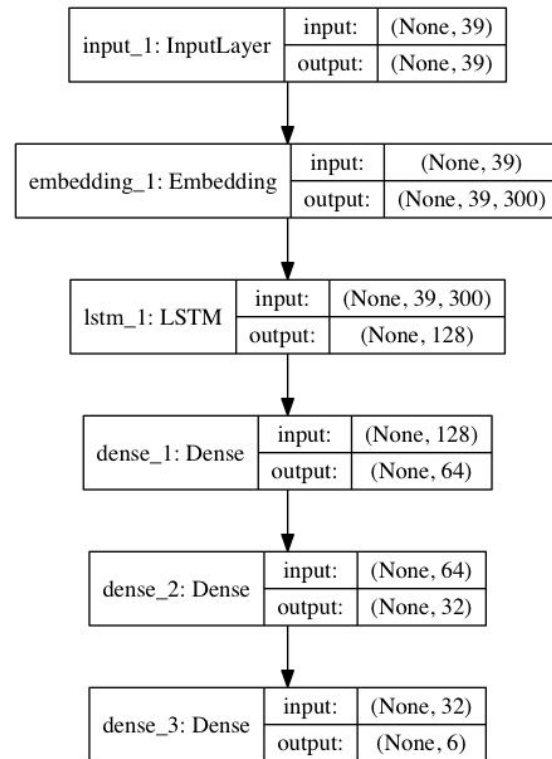
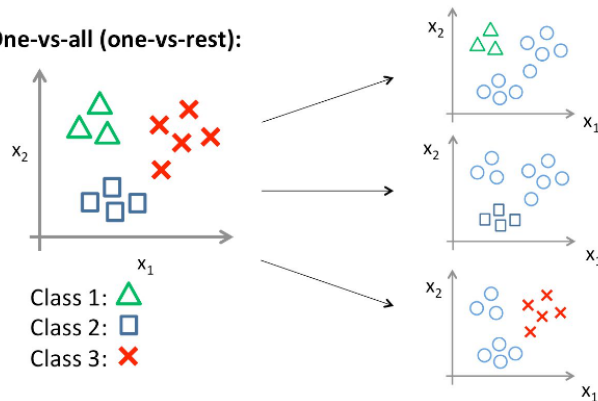


Figure 28.3 A subset of the Li and Roth (2005) answer types.

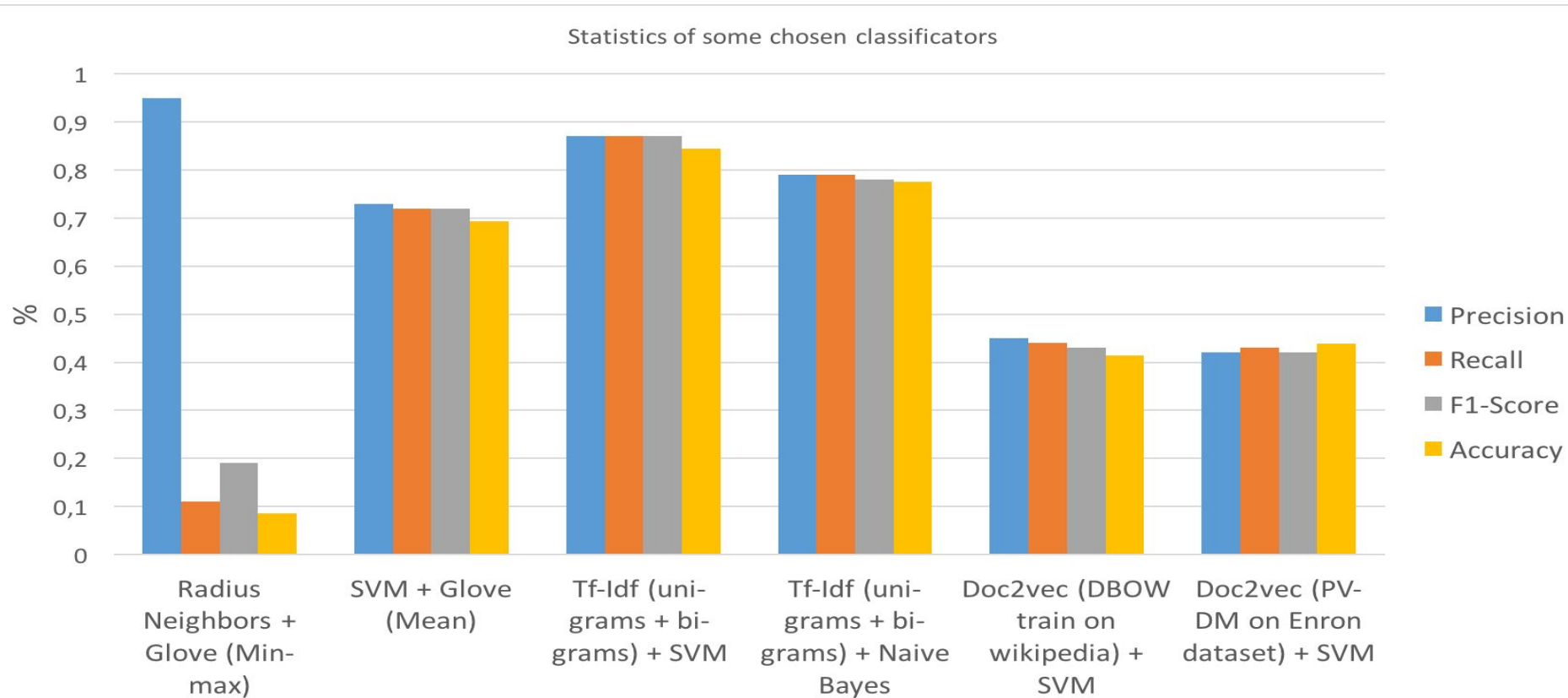
Algorithms

1. Naive Bayes
- 2. SVM (+ bigrams) - 83%**
3. kNN
4. Neural network
 - a. LSTM - 89 %**
 - b. CNN
 - c. LSTM + CNN

One-vs-all (one-vs-rest):

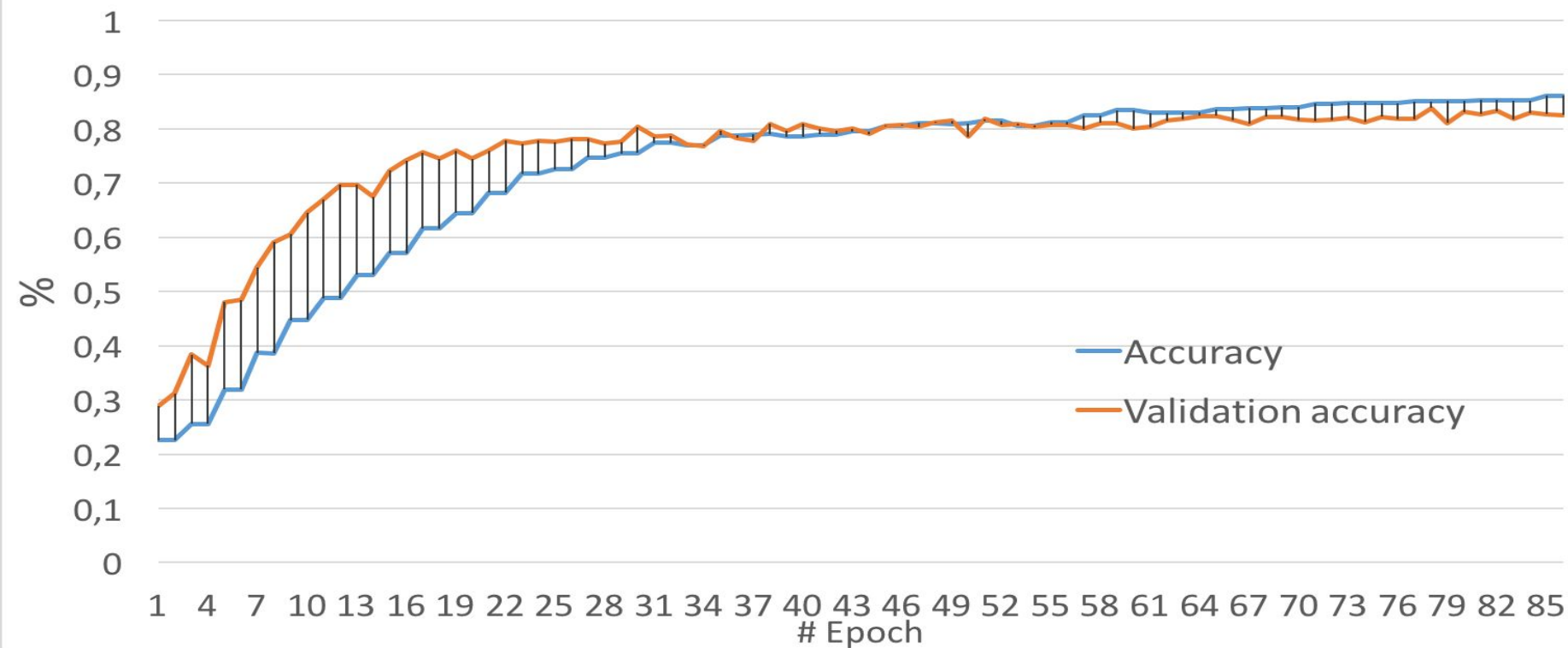


Results



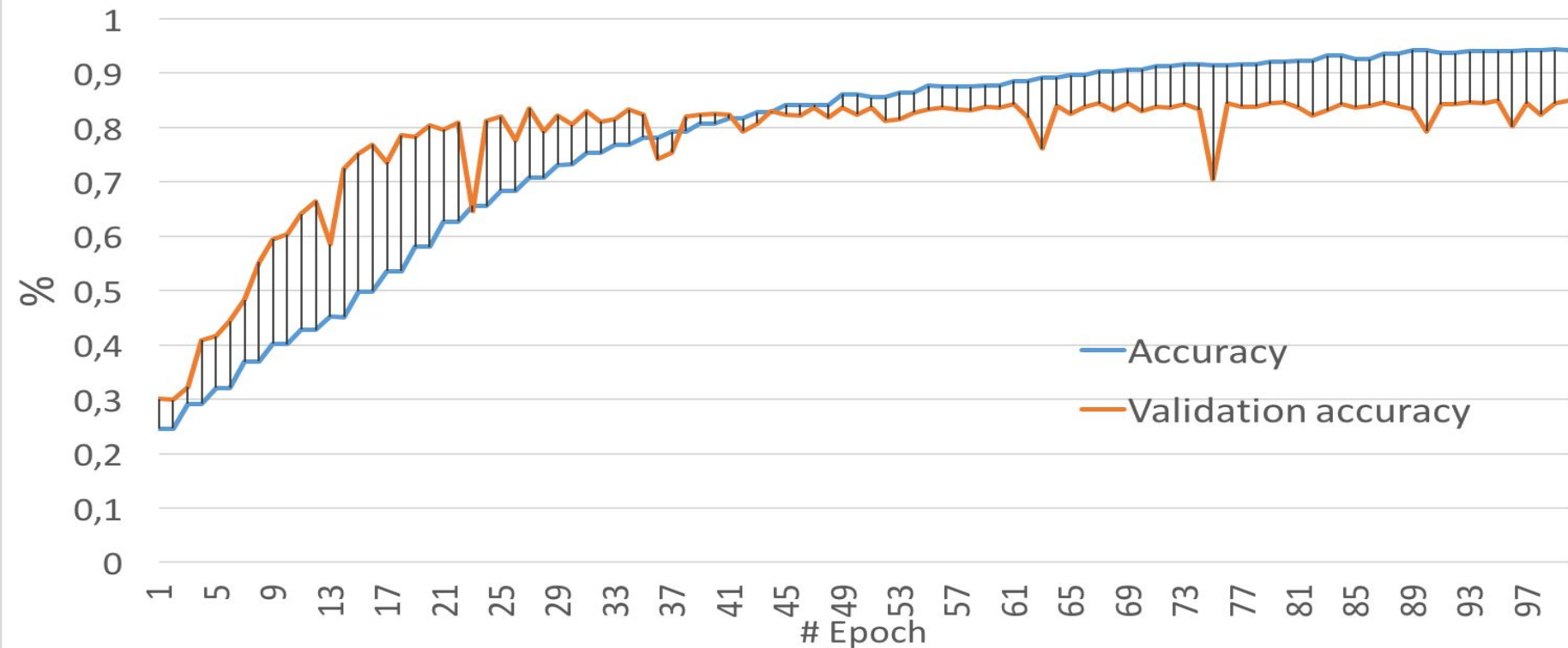
Results

Accuracy over time LSTM



Results

Accuracy over time LSTM + CNN



Results

- FastText is library from Facebook, which has build-in function to classify text (based on train data), it also provide similar functionality like word2vec

	# Unique Words	Precision	Recall	Time to train 5 Epoch
Train on Quora Test dataset + Retrain on Base dataset	136 665	0.862	0.862	7 min
Train on Quora Train dataset + Retrain on Base dataset	57 165	0.85	0.85	94 s
Train on Base dataset	7 111	0.818	0.818	8 s

Reference

<https://web.stanford.edu/~jurafsky/slp3/28.pdf>

Frameworks



Code

Available at: <https://github.com/petrLorenc/EmailReply>

Report at: <https://github.com/petrLorenc/EmailReply/blob/master/report.pdf>