



Paraphrase Detection in NLP



DataRobot

Yuriy Guts
ML Engineer

Paraphrasing

A composite image featuring a Tuscan vineyard landscape in the background. In the foreground, there is a wooden barrel with a bottle of wine and a glass of wine. A bunch of green grapes is hanging from the top of the barrel, and another bunch is on the table in front of the glass. The scene is set against a warm, golden sky.

Any trip to Italy should include a visit to Tuscany to sample their exquisite wines.

Be sure to include a Tuscan wine-tasting experience when visiting Italy.

Paraphrase Identification



Where can I get very **professional and reliable** envelope printing service in Sydney?

Where can I get very **affordable** branded envelope printing service in Sydney?



Why are doctors always late?

Why doctors always make you wait for 15-20 minutes before they see you?

The Quora logo is displayed in white text on a dark red background.

 Featured Prediction Competition

Quora Question Pairs

Can you identify question pairs that have the same intent?

\$25,000 · 3,307 teams · 2 months ago

\$25,000

Prize Money

Data

Q1 (ID, Title) Q2 (ID, Title)

404,290 pairs

Training set

Target

Binary (metric: log-loss)

2,345,796 pairs

Test set

Challenges

Sometimes a single word matters



What are the best books on **IT** leadership?

What are the best books about leadership?



Will the Miami Heat win the NBA championship in **2011**?

Will the Miami Heat win the NBA championship in **2012**?



How do I lose **20** pounds?

How do I lose **15** pounds?

Sometimes there are almost no shared words



I am unable to talk to girls, leave being friendly with them. Why?

I am shy to talk to any woman because i get nervous and freaked out around them. What is the solution?



Is there a Quora user who have seen a UFO?

Have you seen an alien?

Sometimes ALL the words are shared



Is the Government of **Pakistan** encouraging **India** by not taking any real action against ceasefire violations?

Is the Government of **India** encouraging **Pakistan** by not taking any real action against ceasefire violations?



What is the most interesting thing we learned about **Portugal's** World Cup team in their match against **Germany**?

What is the most interesting thing we learned about **Germany's** World Cup team in their match against **Portugal**?

Approaches

1. Counting Stuff.
2. Information Theory.
3. Linguistics.
4. Deep Learning.

BoW Representations

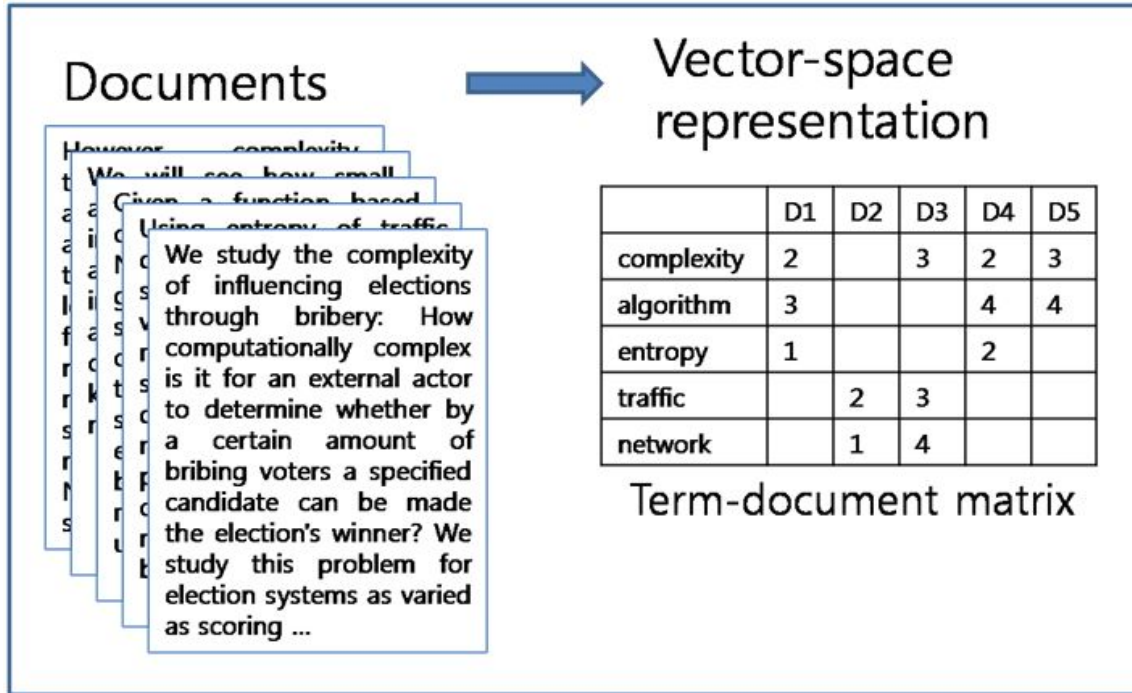


Image copyright © S. Mukherjee

TFIDF

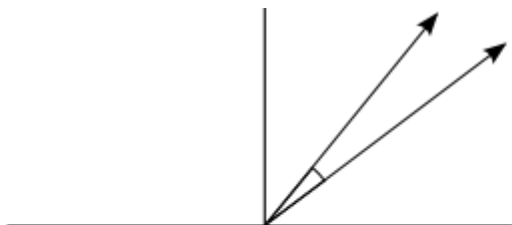
For a term i in document j :

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

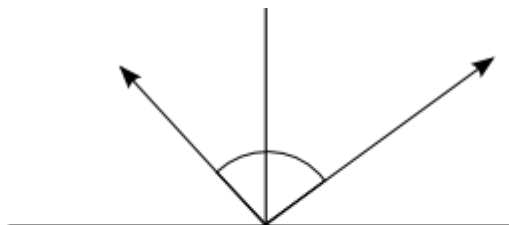
$tf_{i,j}$ = number of occurrences of i in j
 df_i = number of documents containing i
 N = total number of documents

Cosine Similarity

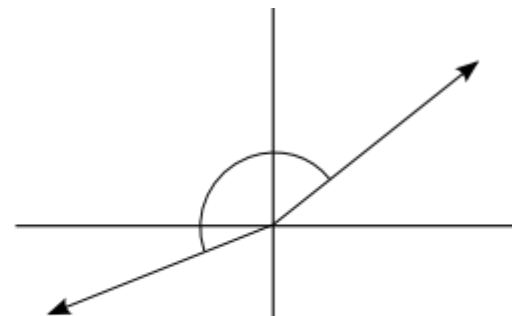
$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



Similar scores
Score Vectors in same direction
Angle between them is near 0 deg.
Cosine of angle is near 1 i.e. 100%



Unrelated scores
Score Vectors are nearly orthogonal
Angle between them is near 90 deg.
Cosine of angle is near 0 i.e. 0%

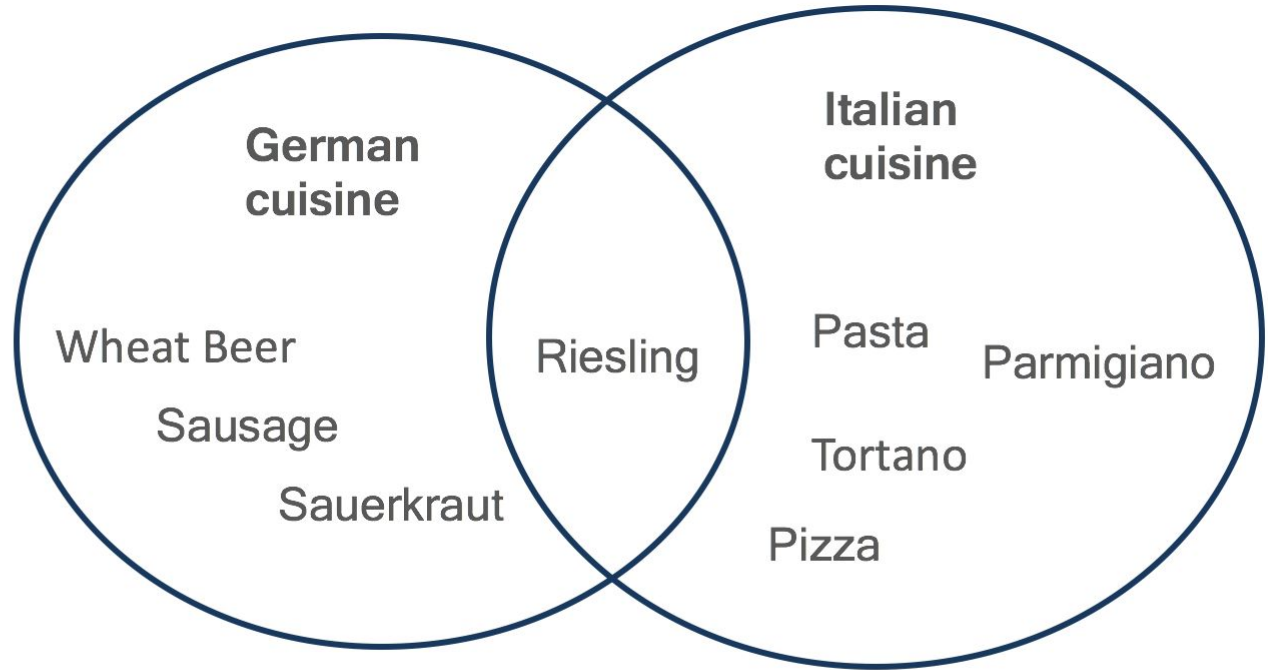


Opposite scores
Score Vectors in opposite direction
Angle between them is near 180 deg.
Cosine of angle is near -1 i.e. -100%

Jaccard Similarity

$$\text{sim}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Jaccard
similarity



$$\text{sim}\left(\text{🇩🇪}, \text{🇮🇹}\right) = \frac{1}{8}$$

Edit Distances

How similar are the names
“Peter” and “Piotr”?

Assume the following cost function

<i>Substitution</i>	1 Unit
<i>Insertion</i>	1 Unit
<i>Deletion</i>	1 Unit

$D(\text{Peter}, \text{Piotr})$ is 3

Peter



Substitution (i for e)

Piter



Insertion (o)

Pioter

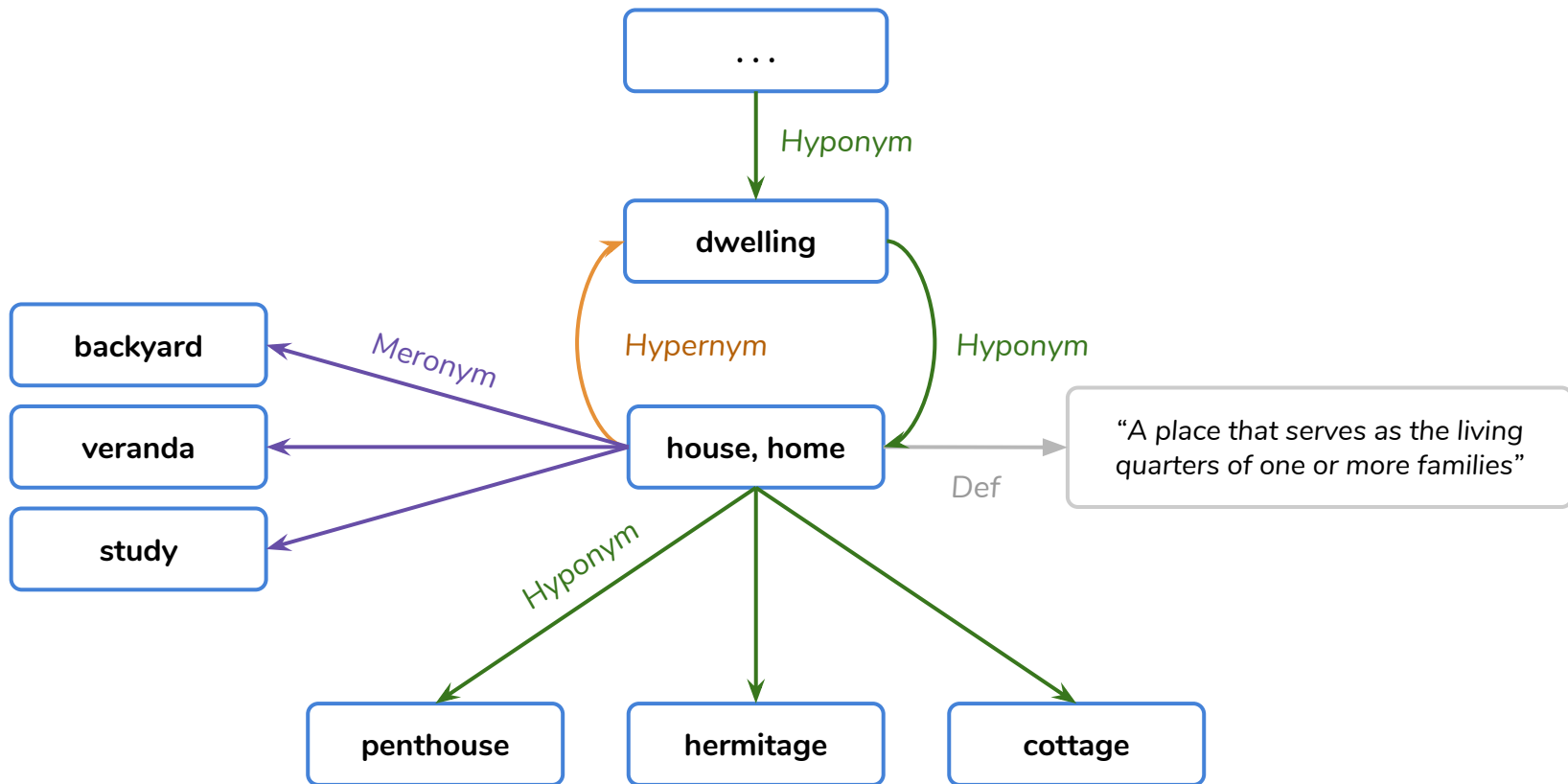


Deletion (e)

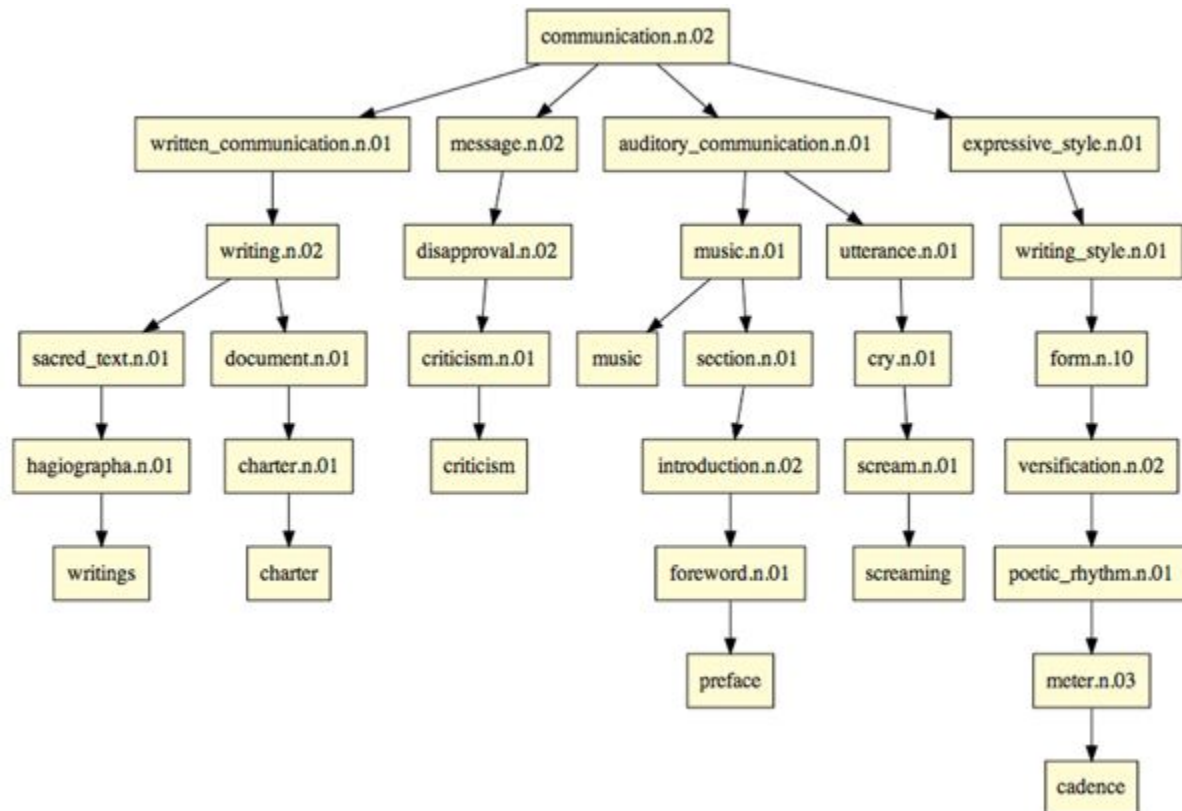
Piotr

1. Levenshtein
 2. Damerau-Levenshtein
 3. Jaro
 4. Jaro-Winkler
-

Lexical Databases and Ontologies



WordNet



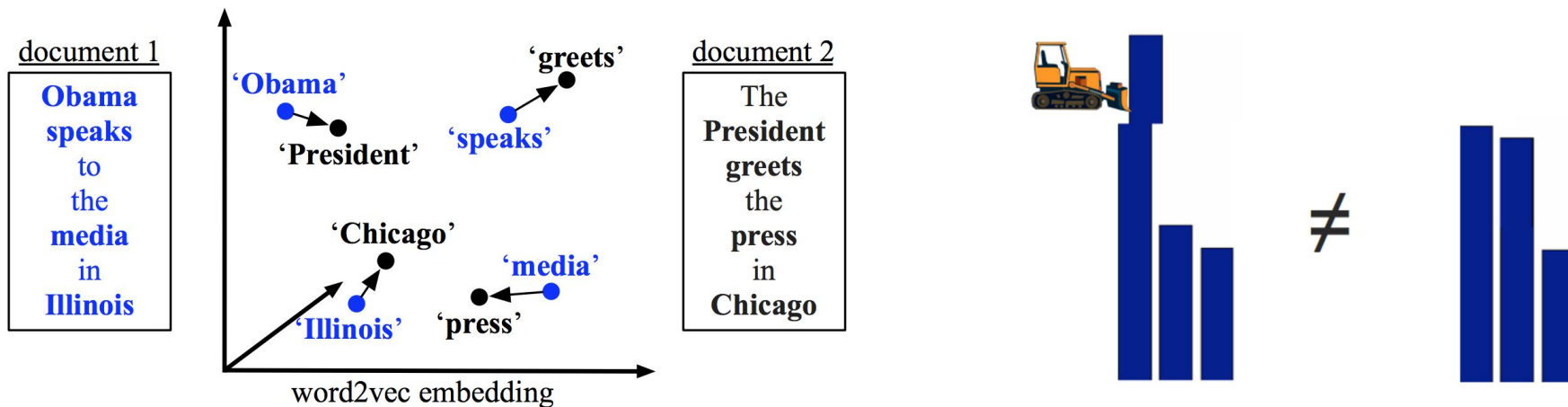
Distributed Hypothesis of Language

“The complete meaning of a word is always contextual, and no study of meaning apart from context can be taken seriously”

John R. Firth. The technique of semantics.

Transactions of the Philological Society, 1935.

Word Mover's Distance (WMD)

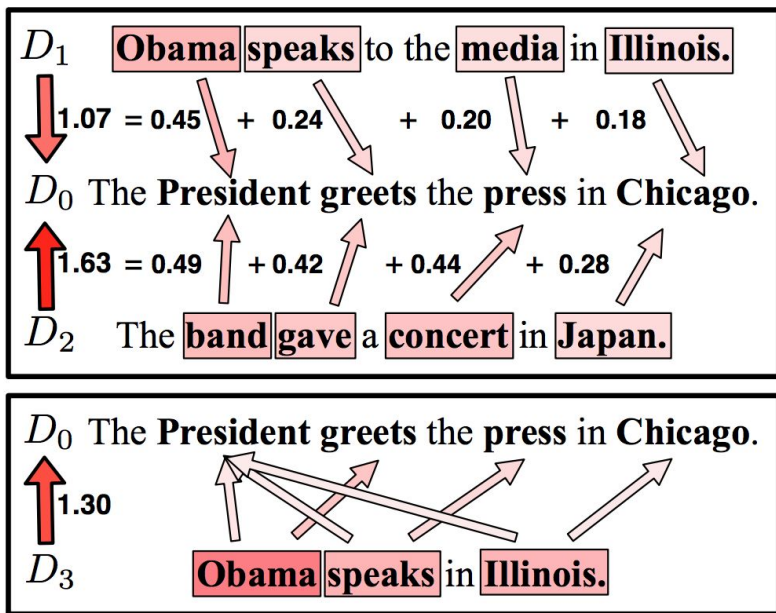


The minimum amount of “work” needed to transform document **1** to document **2**. Inspired by Earth Mover’s Distance, a well-studied transportation problem.

M. Kusner et al. “From Word Embeddings to Document Distances”, 2015.

<http://proceedings.mlr.press/v37/kusnerb15.pdf>

WMD: Linear Optimization Problem



$$d_i = \frac{c_i}{\sum_{j=1}^n c_j} \quad \text{nBOW frequency of the } i\text{-th word in the document}$$

$$\mathbf{T}_{ij} \geq 0 \quad \text{"How much" of word } i \text{ travels to word } j$$

$$\min_{\mathbf{T} \geq 0} \sum_{i,j=1}^n \mathbf{T}_{ij} c(i,j)$$

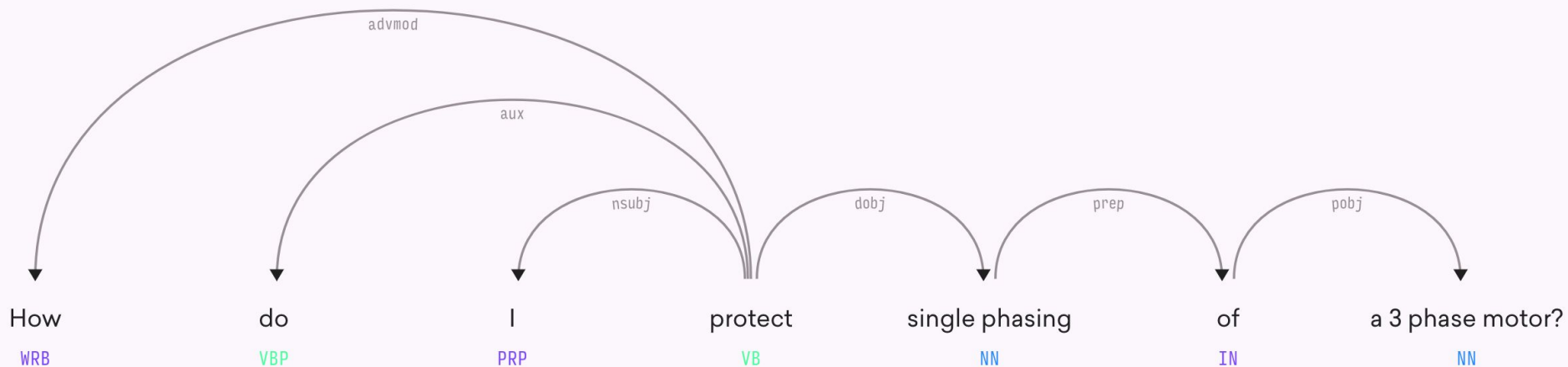
$$\text{subject to: } \sum_{j=1}^n \mathbf{T}_{ij} = d_i \quad \forall i \in \{1, \dots, n\}$$

$$\sum_{i=1}^n \mathbf{T}_{ij} = d'_j \quad \forall j \in \{1, \dots, n\}.$$

M. Kusner et al. "From Word Embeddings to Document Distances", 2015.

<http://proceedings.mlr.press/v37/kusnerb15.pdf>

Morphology & Syntax Features



What is the most interesting thing we learned about **Portugal** **GPE** 's **World Cup** **EVENT** team in their match against **Germany** **GPE** ?

Architectural Principles for State-of-the-Art Neural NLP

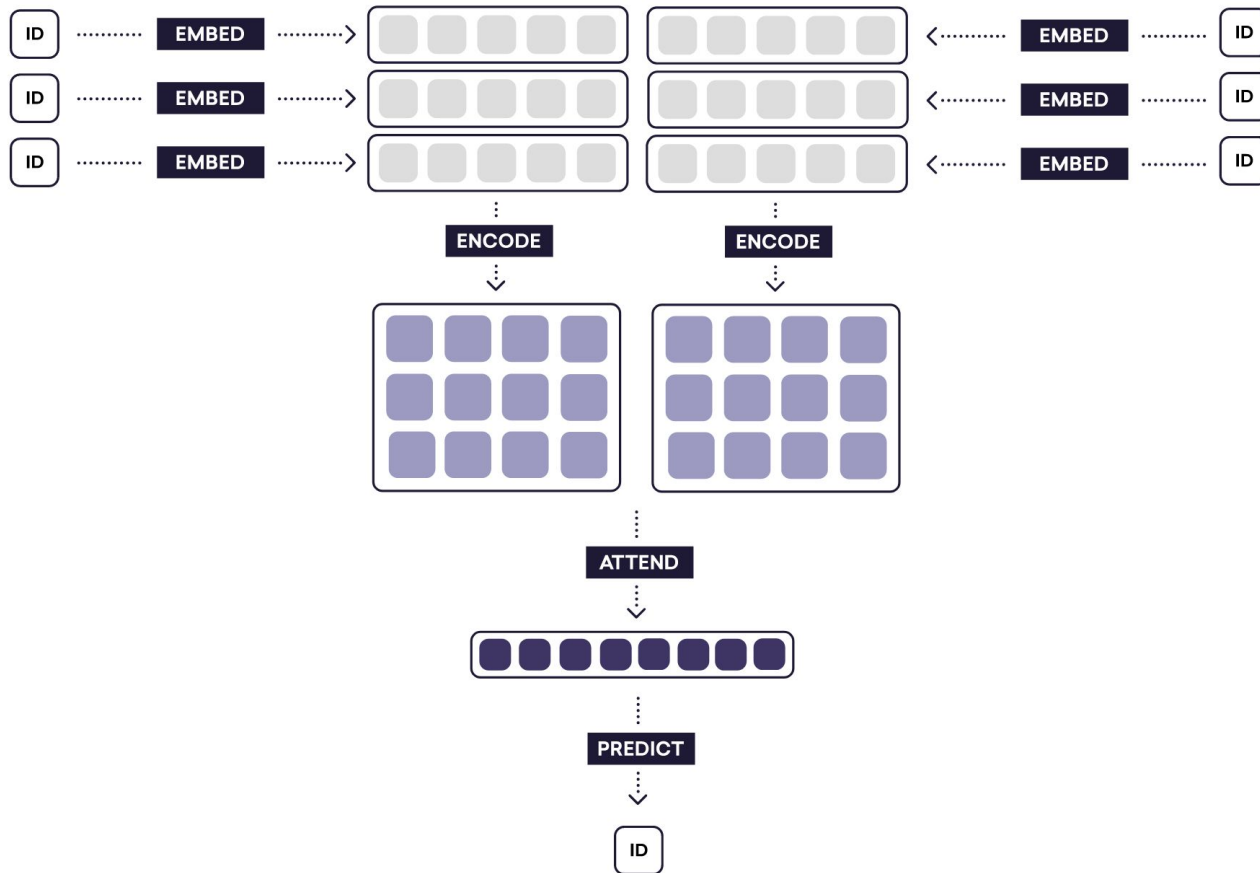
Embed

Encode

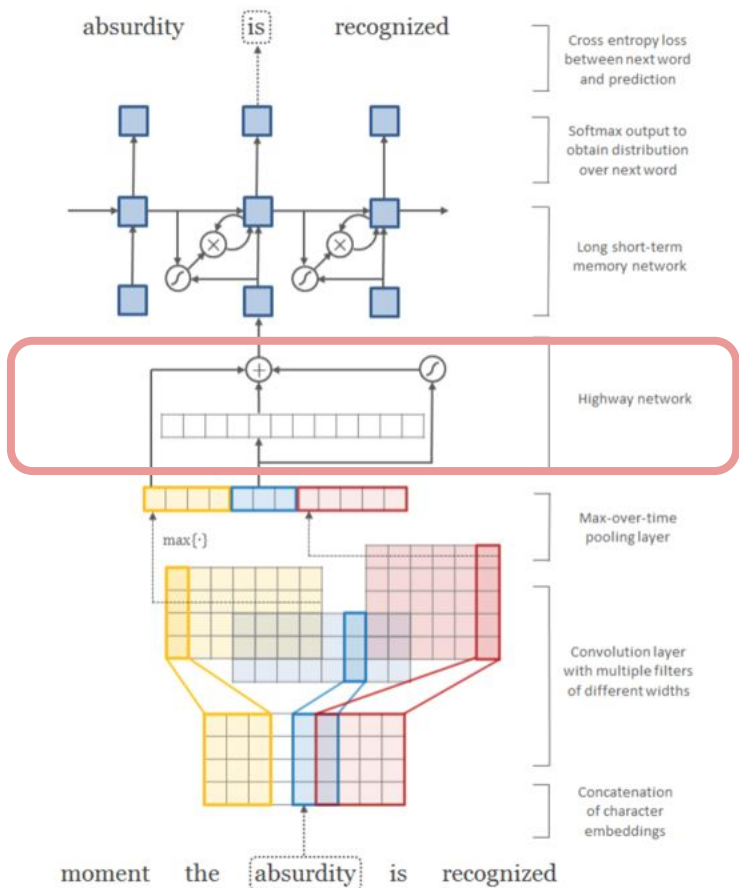
Attend

Predict

State-of-the-Art NLP Pipeline



Recurrent Highway Networks



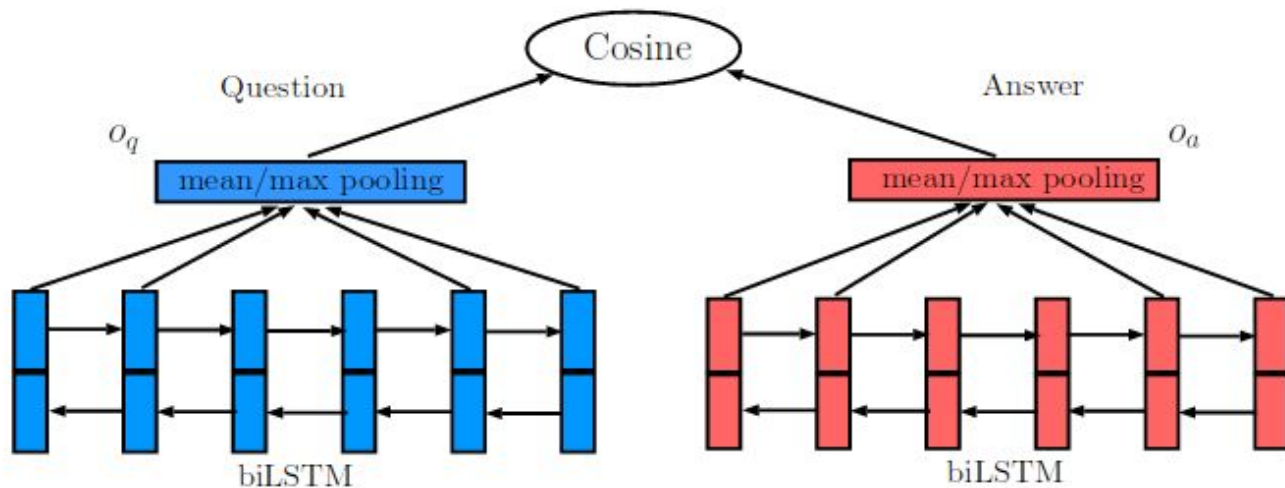
Highway layer:

$$\mathbf{y} = H(\mathbf{x}, \mathbf{W}_H) \cdot T(\mathbf{x}, \mathbf{W}_T) + \mathbf{x} \cdot (1 - T(\mathbf{x}, \mathbf{W}_T))$$

Srivastava et al. "Recurrent Highway Networks",
2015

<https://arxiv.org/pdf/1607.03474.pdf>

Two Input Documents? Siamese Network



Run the same encoding step for every input.

Share encoder weights, don't learn W_{E1} and W_{E2} separately.

Distance Learning, Contrastive Loss

$$L(W, Y, \vec{X}_1, \vec{X}_2) = (1 - Y) \frac{1}{2} (D_W)^2 + (Y) \frac{1}{2} \{ \max(0, m - D_W) \}^2$$

Penalize similar pairs by a monotonically **increasing** function of their learned distance.

Penalize different pairs by a monotonically **decreasing** function of their learned distance.

Hadsell, Chopra, LeCun “Dimensionality Reduction by Learning an Invariant Mapping”, 2006.

<http://yann.lecun.com/exdb/publis/pdf/hadsell-chopra-lecun-06.pdf>

The Quora logo is displayed in white text on a dark red background.

 Featured Prediction Competition

Quora Question Pairs

Can you identify question pairs that have the same intent?

\$25,000

Prize Money

\$25,000 · 3,307 teams · 2 months ago

Data

Q1 (ID, Title) Q2 (ID, Title)

404,290 pairs

Training set

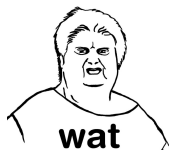
Target

Binary (metric: log-loss)

2,345,796 pairs

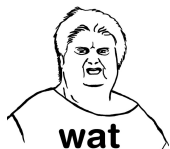
Test set

Sometimes the labeling is just plain wrong



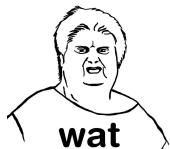
How can I prevent pneumonia?

How do you prevent pneumonia?



What is the car in this picture?

What car is in this picture?



How do I protect single phasing of a 3 phase motor?

What is single phase and 3 phase?

People misspell. A lot.

demonitization

demonitization

masturbation

masturbation

demonitizing

demonitization

masturbation

masturbation

demonitising

demonitization

masturbation

masturbation

demonitisation

demonitization

masturbation

masturbation

demonitize

demonitization

masturbation

masturbation

demonitisation

demonitization

masturbation

masturbation

demonitization

demonitization

masturbation

masturbation

demonitisation

demonitization

masturbation

masturbation

demonitisation

demonitization

masturbation

masturbation

demonitisation

demonitization

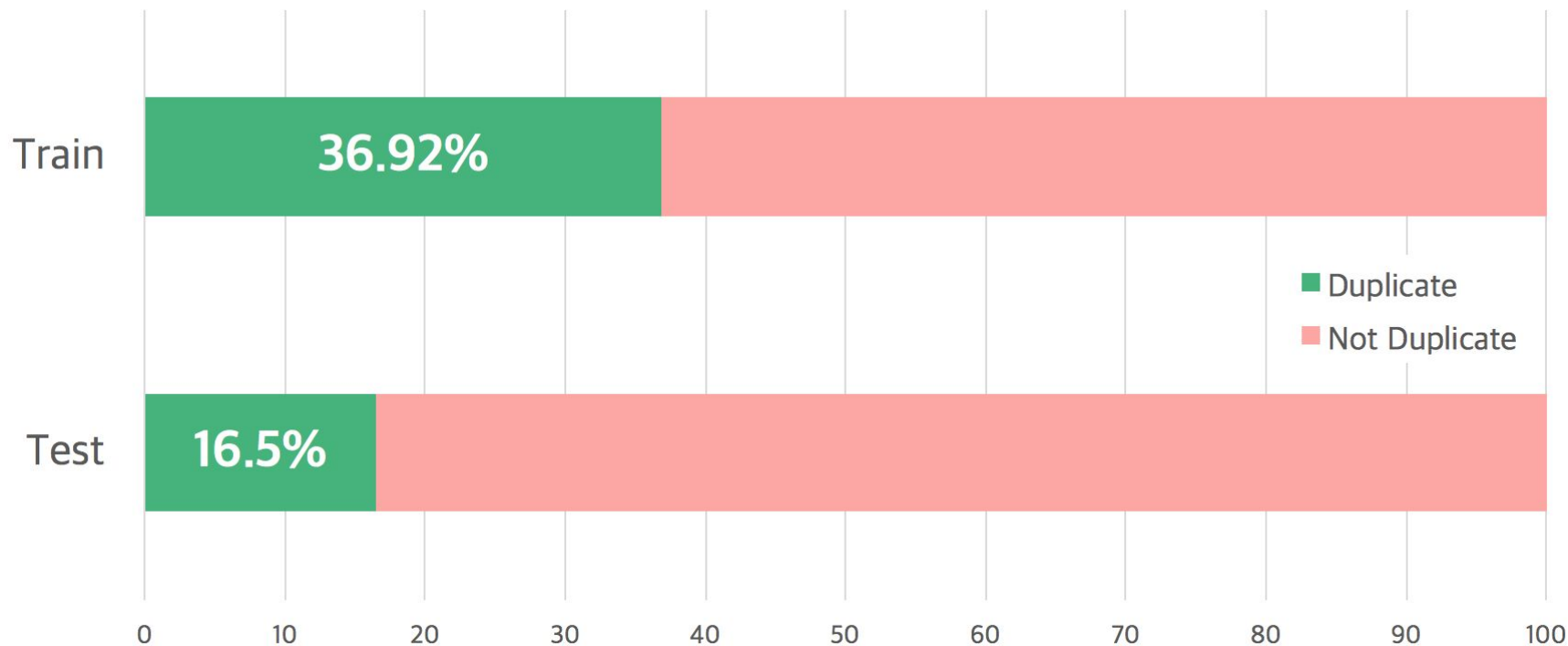
masturbation

demonitising

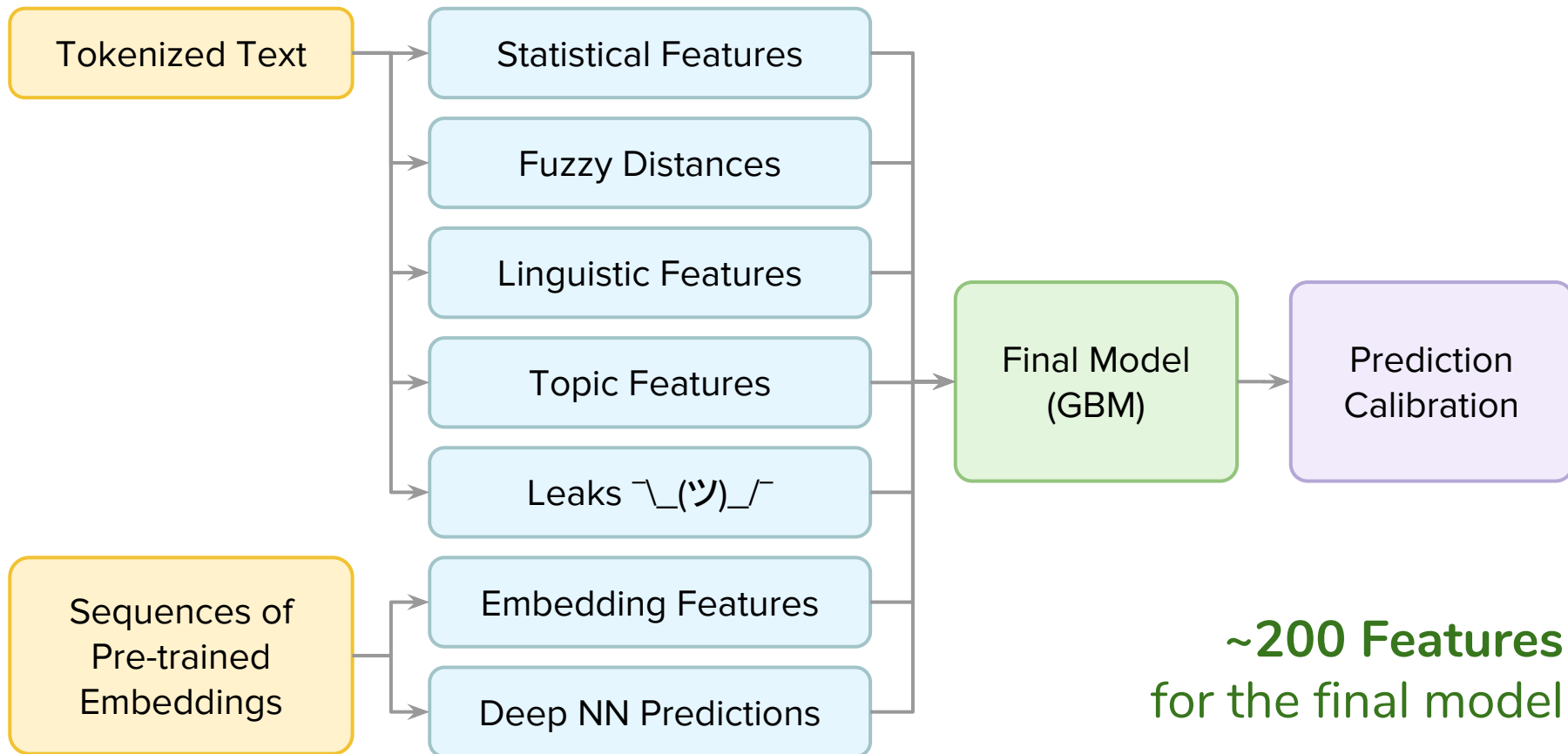
demonitization

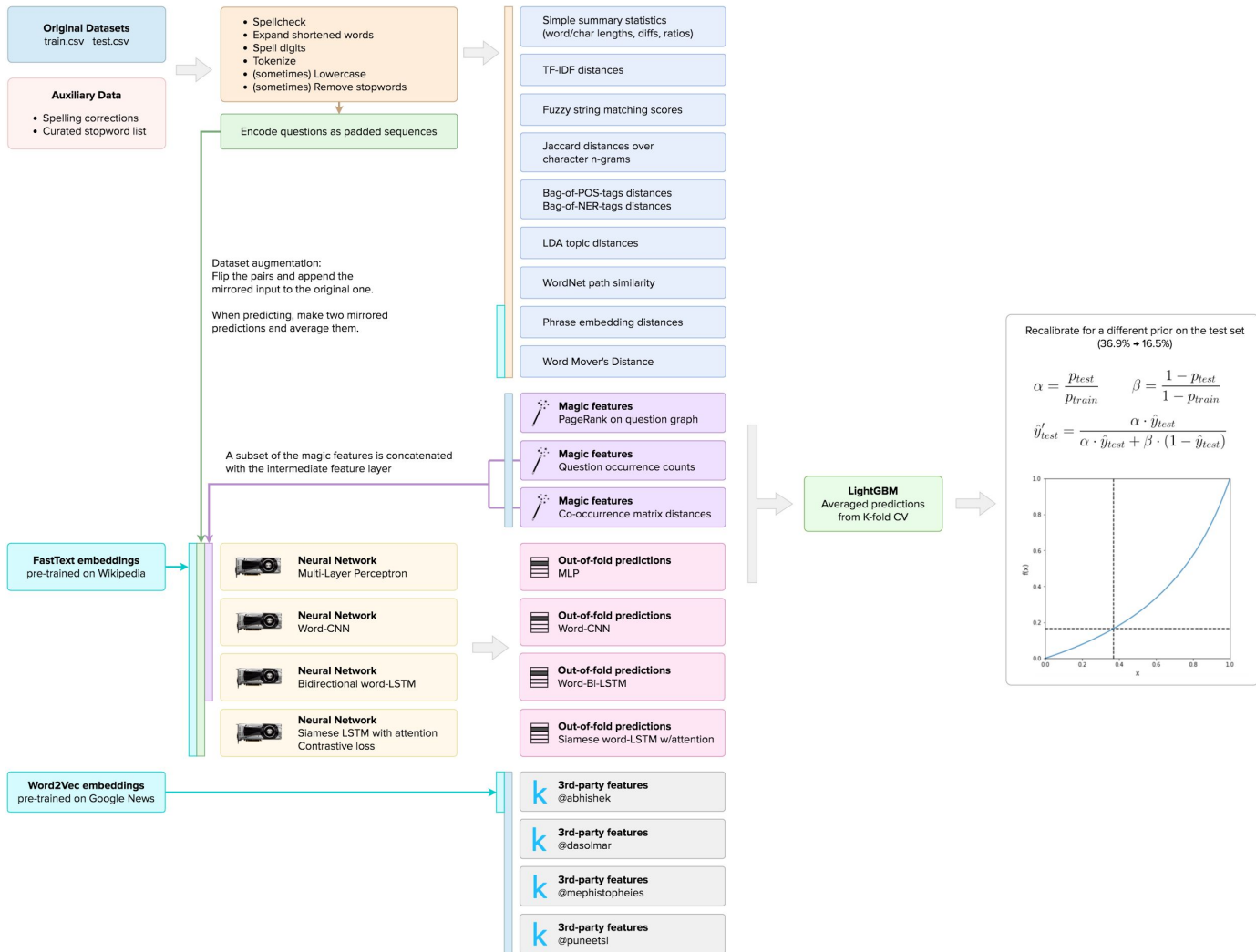
100,000+
spelling corrections

Different target balance for Train and Test



Final Solution (Simplified)





F1 Score

0.8829

True Positive Rate
(Sensitivity)

0.89

False Positive Rate
(Fallout)

0.0738

True Negative Rate
(Specificity)

0.9262

Positive Predictive
Value (Precision)

0.876

Negative Predictive
Value

0.935

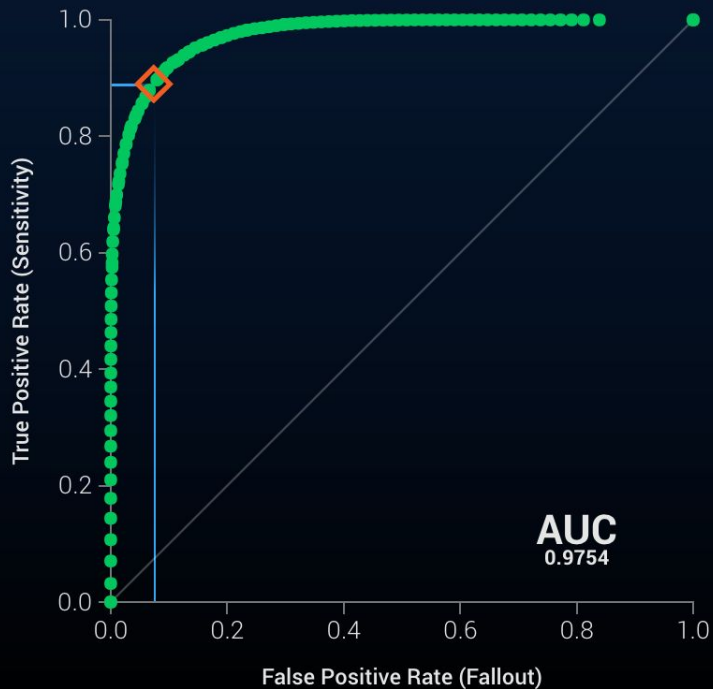
Accuracy

0.9129

ROC Curve

Data Source:

Holdout

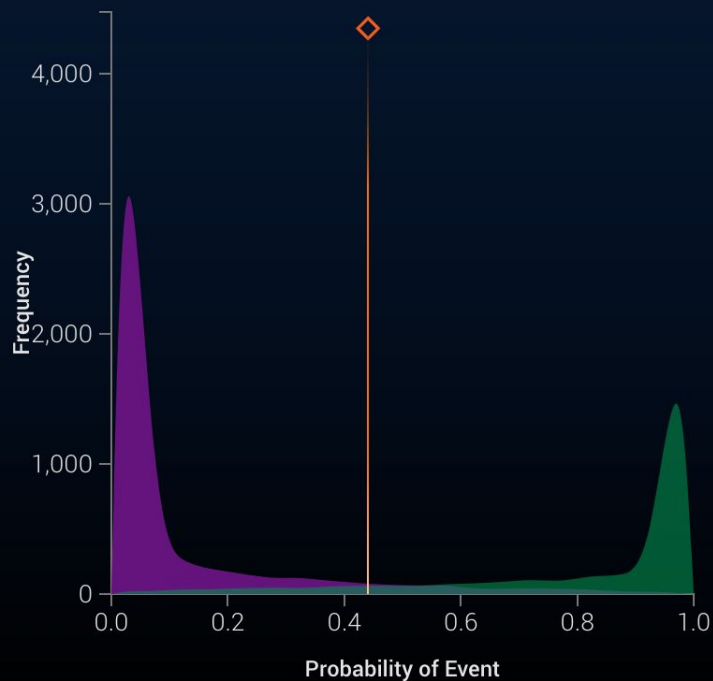


Prediction Distribution

Threshold (0-1):

0.4409

Frequency





facebook.com/yuriy.guts



github.com/YuriyGuts/kaggle-quora-question-pairs

