ZHAO JIN

KANGJING FAN

# QUOTA QUESTION PAIRS

# OVERVIEW

▸ Task: Identify duplicate questions.

▸ Dataset size:

  ▸ Training:

    ▸ Size (404290), Field (question pairs+gold label)

  ▸ Test:

    ▸ Size (2345796), Field (question pairs)

▸ Character: Range:(15,150), Mean(~60)

▸ Word: Range:(3,30), Mean(~11)

# PREPROCESSING

▸ Use re package separating punctuation from the  word, restore abbreviation

   ▸ 's —- is

▸ Remove stop-words

▸ Correct spelling: Spell checker & word2vec

▸ Lemmatize word (WordNetLemmatizer())
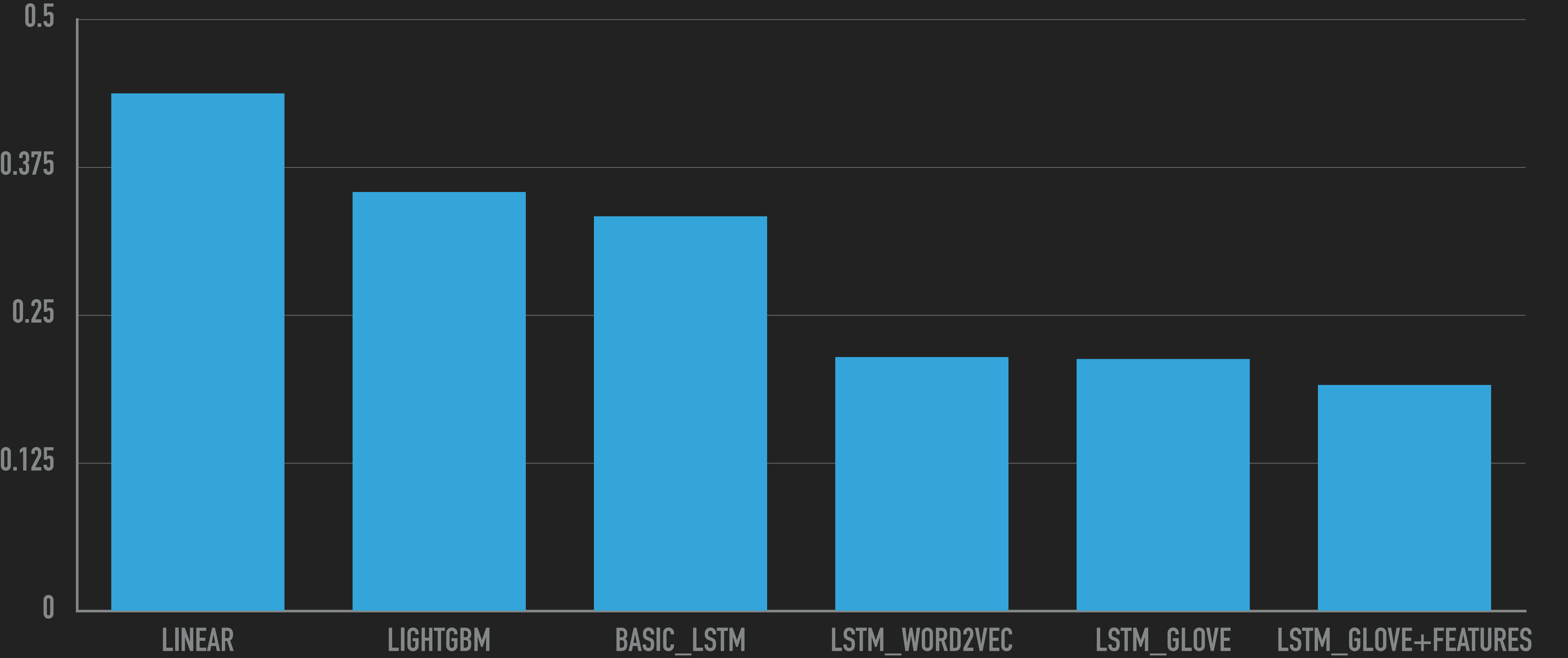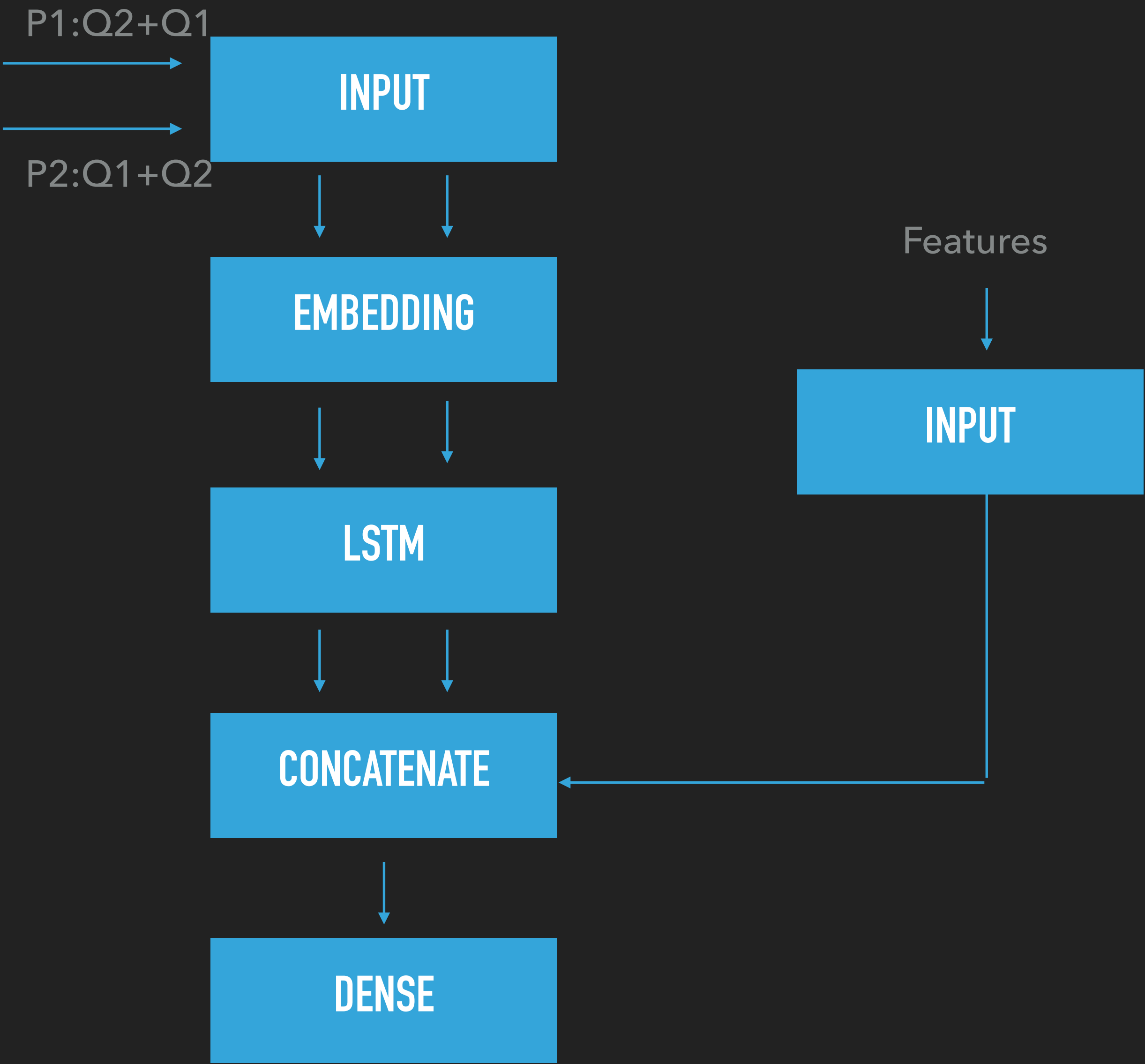
▸ Get word stem (SnowballStemmer('english'))

# METHOD

▸ Common word + linear model len (Q1 ∩ Q2) / len (Q1 ∪ Q2)

▸ Magic features + LIGHTGBM

▸ Basic LSTM

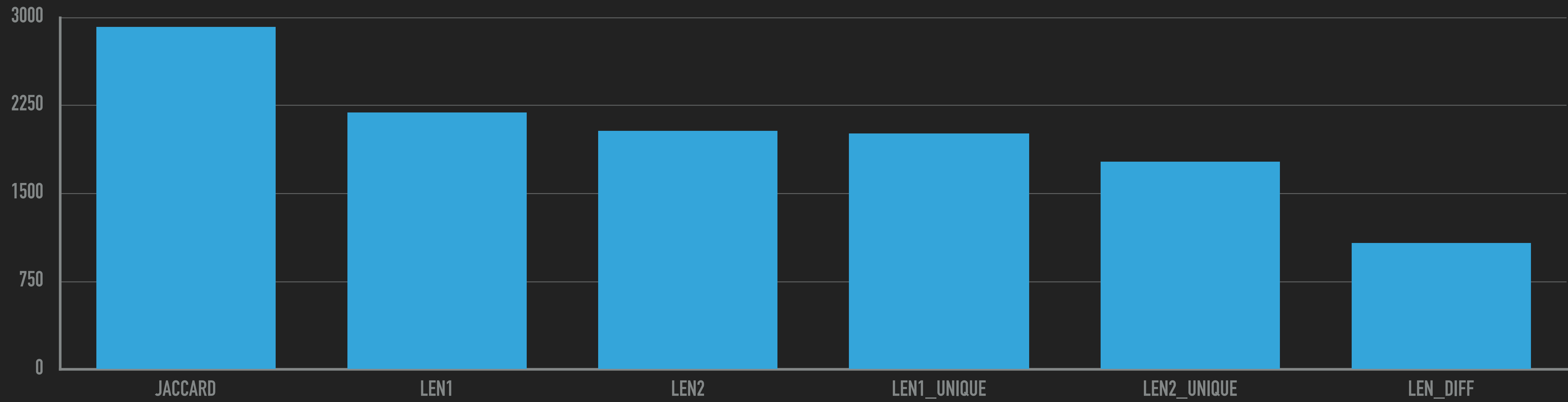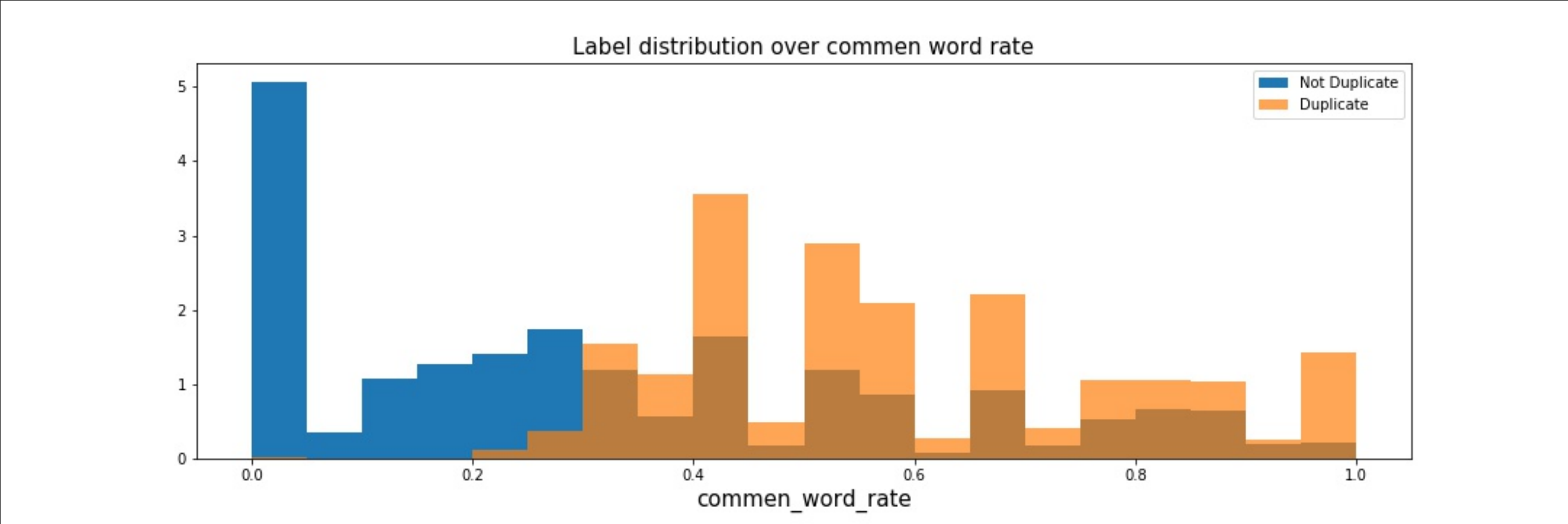▸ Boost LSTM (with pre-trained word embedding)

▸ Boost LSTM + magic features

# RESULT

# BEST MODEL



P1:Q2+Q1

P2:Q1+Q2

INPUT

EMBEDDING

LSTM

CONCATENATE

DENSE

Features

INPUT

# ANALYSIS



Label distribution over commen word rate

# CONCLUSION

▸ For the data pre-processing

  ▸ Spell check is necessary

  ▸ Spell check + lemmatize + stem

▸ For the model

  ▸ LSTM > Tree base model

  ▸ Ensemble model is the best choice

# THANK YOU

# FEATURES

▸ Jacquard: len(Q1 ∩ Q2) /len(Q1 ∪ Q2)

▸ Len1: len(Q1)

▸ Len2: len(Q2)

▸ Len1_unique: len(Q1- Q1 ∩ Q2)

▸ Len2_unique: len(Q2- Q1 ∩ Q2)

▸ Len_diff: |len(Q1) - len(Q2)|