# Project Title: Intent Based Duplicate Question Removal

Team Number: 14

Team Members:

Nikitha Rao (01FB15ECS364)

Rahul Ragesh (01FB15EEC303)

## Abstract:

Quora, a well known social media platform, allows people from all around the world to ask questions and enables them to connect with various domain experts who provide quality answers and several new insights. With over 100 million people visiting Quora every month, the user base is humungous. It is not difficult to inference that several people may end us asking similarly worded questions. When there are multiple questions being asked with the same intent, a reader would have to spend more time in order to find the most suitable answer to their question. At the same time a writer may feel that they are answering multiple versions of the same question which may seem like a waste of time. This in turn leads us to our problem. Given any two questions, identify whether they have the same intent i.e. is one question a duplicate of the other or not. We will be using the dataset provided by Quora for the same. However, we aim to generalise our approach to be able to tackle the problem for any Question - Answer Forum, such as Quora itself or Stack Overflow.

## Libraries Used:

* numpy

* Tensorflow

* Keras

* scikit-learn

* h5py

* hdf5

## Dataset Used:

Quora Question Pairs:

Link - https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs

There are over 400,000 lines of potential question duplicate pairs. Each line contains IDs for each question in the pair, the full text for each question, and a binary value that indicates whether the line truly contains a duplicate pair.

| id | qid1 | qid2 | question1 | question2 | is_duplicate |
|---|---|---|---|---|---|
| 0 | 1 | 2 | What is the step by step guide to invest in share market in india? | What is the step by step guide to invest in share market? | 0 |
| 1 | 3 | 4 | What is the story of Kohinoor (Koh-i-Noor) Diamond? | What would happen if the Indian government stole the Kohinoor (Koh-i-Noor) diamond back? | 0 |
| 2 | 5 | 6 | How can I increase the speed of my internet connection while using a VPN? | How can Internet speed be increased by hacking through DNS? | 0 |
| 3 | 7 | 8 | Why am I mentally very lonely? How can I solve it? | Find the remainder when [math]23^{24}[/math] is divided by 24,23? | 0 |
| 4 | 9 | 10 | Which one dissolve in water quikly sugar, salt, methane and carbon di oxide? | Which fish would survive in salt water? | 0 |
| 5 | 11 | 12 | Astrology: I am a Capricorn Sun Cap moon and cap rising...what does that say about me? | I'm a triple Capricorn (Sun, Moon and ascendant in Capricorn) What does this say about me? | 1 |
| 6 | 13 | 14 | Should I buy tiago? | What keeps childern active and far from phone and video games? | 0 |
| 7 | 15 | 16 | How can I be a good geologist? | What should I do to be a great geologist? | 1 |
| 8 | 17 | 18 | When do you use ⊱ instead of ∟? | When do you use "&" instead of "and"? | 0 |
| 9 | 19 | 20 | Motorola (company): Can I hack my Charter Motorolla DCX3400? | How do I hack Motorola DCX3400 for free internet? | 0 |

## GloVe: Global Vectors for Word Representation

GloVe is an unsupervised learning algorithm for obtaining vector representations for words. Training is performed on aggregated global word-word co-occurrence statistics from a corpus, and the resulting representations showcase interesting linear substructures of the word vector space.

## Project Files and their Descriptions:

main.py – The main function. Calls the rest of the functions in the order of the files mentioned below.

get_dataset.py – The preprocessing function. It checks if the dataset(Quora Questions) is present, if not it downloads the necessary data from the links specified. The downloaded data is extracted, tokenized and encodes it. We then check if the GloVe data is present and download it if not. This has information about the word embeddings. Using the GloVe embeddings we create a new embedding matrix to denote the word similarity for all the words present in the Quora dataset. This is followed by saving the preprocessed data in the form of files – q1_train.npy, q2_train.npy, label_train.npy, word_embedding_matrix.npy.

model_def.py -  The preprocessed data is split in the form of test data and training data (1:9) ratio. We then define the model architecture for the feature extraction and classification problem. Hyper parameters are tuned and the various layers are defined in the neural network. The output is in the form of 0 or 1 representing true or false depending on whether the two questions are convey the same intent or not. The loss function used is binary cross-entropy.

training.py – Trains the model based on the training data for 50 epochs. After each epoch, there is an early stopping criteria through cross validation is used. The weights for the epoch having the best validation accuracy  is saved in a file.

testing.py – The model is tested by using the test data.


## Instructions to run:

python3 main.py
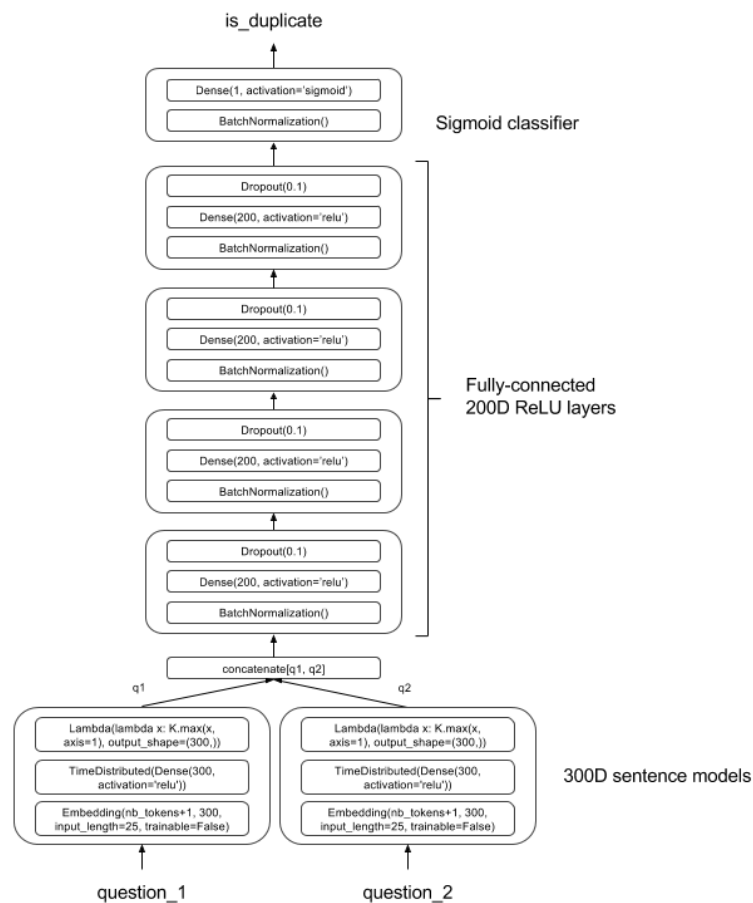

## Control Flow of the Program:

1. Pre-process the Dataset
2. Extract relevant information from the dataset.
3. Tokenize all the questions.
4. Encode text to sequence.
5. Generate word embedding matrix using Stanford GLoVe model.
6. Split the dataset into training, test and validation set.
7. Build the network model as per the given architecture.
8. Adam optimizer with Binary Cross Entropy Loss function was employed.
9. Save the configuration files.
10. Train the model.
11. Model evaluation is done to generate cross validation accuracy.

## Training Time:

Training takes approximately 300 secs/epoch on an average, using Tensorflow as a backend for Keras on an intel core i7 8<sup>th</sup> gen, 8GB RAM. No GPU used.

For 50 epocs, 250 mins ~ 4 to 5 hours.

## Network Architecture: