# DUPLICATE BUSTERS

**Marc Vornetran 11569565**
marc.vornetran@gmail.com
**Maaike Koolbergen 10592377**
maaikekoolbergen@gmail.com

# Predicting Quora Question Similarity with a Siamese deep MaLSTM network

**Q**

## Data Processing



**Pre-Processing**
• Load the train, train labels and test data sets;
• Join the train labels with the train data set;
• Remove id rows from both train and test data sets.

**Text Cleaning**
• Convert words to lower case;
• Remove punctuation;
• Normalization (What's to What is)
• Remove stop words;
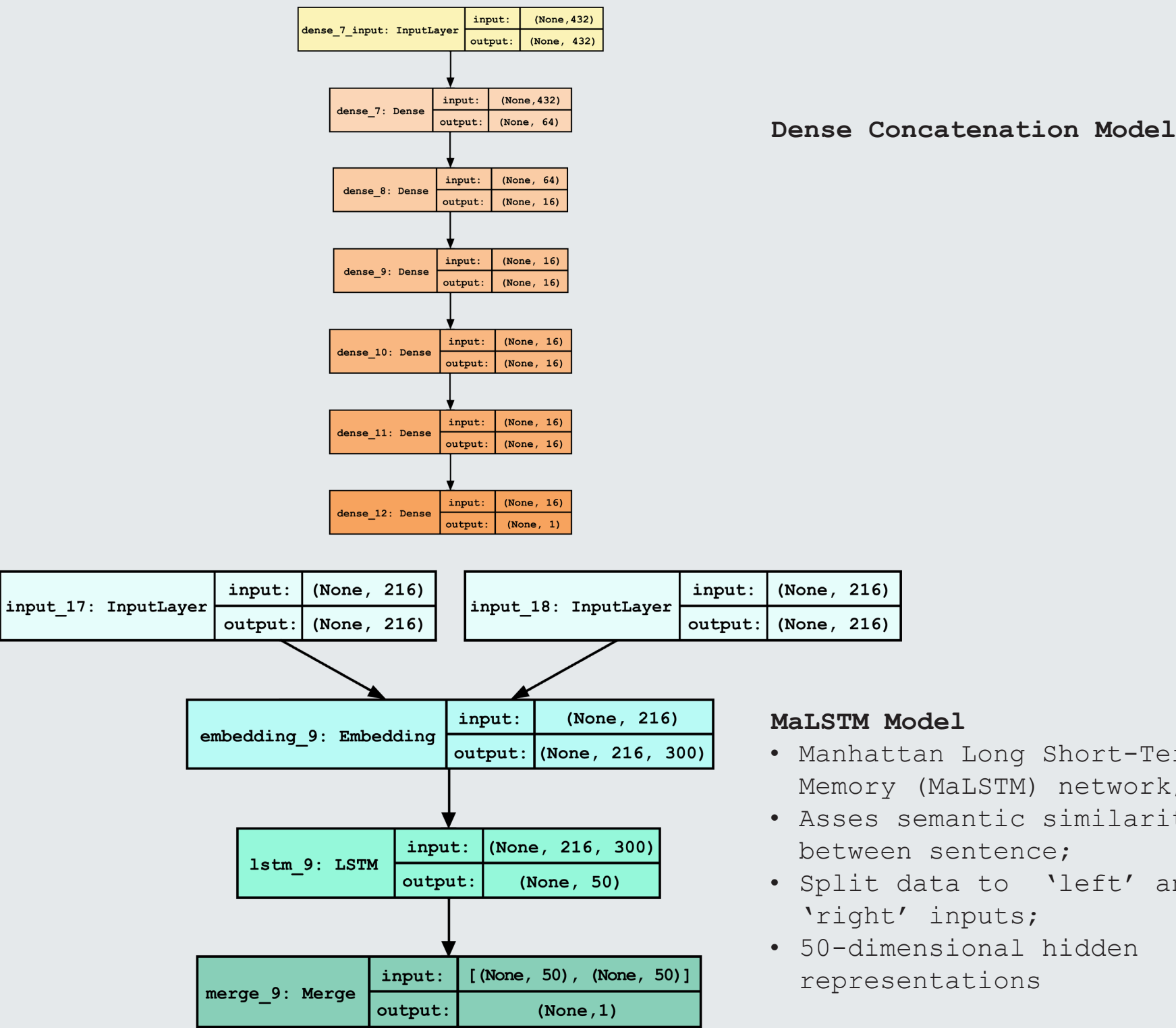• Split sentence into words separated by white space.

## Model Selection



Dense Concatenation Model

**MaLSTM Model**
• Manhattan Long Short-Term Memory (MaLSTM) network;
• Asses semantic similarity between sentence;
• Split data to 'left' and 'right' inputs;
• 50-dimensional hidden representations

## Model Evaluation

## Feature Extraction



**Words to Indices**
• Convert words to indices;
• Start at index 1 to reserve 0 for zero padding.

**Embedding**
• Use Google's Word2Vec embedding to turn words into their embedding;
• Model with 300 dimensional vectors for 3 million words and phrases (pre-trained over about 100 billion words);
•

## Model Training

## Technical Information