

Quora Question Pairs

Sukanya Behera [160714733008]

Chandrabhatta Sriram [160713733042]

Abul Faiz Mohammed Abdul Hadi Mouzzam [160714733108]

Project Guide: E. Shailaja [Asst. Prof. Dept. of CSE]

Project Co-ordinator: R. Sandeep [Asst. Prof. Dept. of CSE]

Introduction [Problem Definition]

- More formally, the duplicate question detection problem can be defined as follows: given a pair of questions $q1$ and $q2$, train a model that learns the function:
 - $f(q1, q2) \rightarrow 0 \text{ or } 1$.
 - where 1 represents that $q1$ and $q2$ have the same intent and 0 otherwise.
 - NOTE: $q1$ and $q2$ are given as string data.

Problem Input [test.csv]

id	qid1	qid2	question1	question2
447	895	896	What are natural numbers?	What is a least natural number?
1518	3037	3038	Which pizzas are the most popularly ordered pizzas on Domino's menu?	How many calories does a Dominos pizza have?
3272	6542	6543	How do you start a bakery?	How can one start a bakery business?
3362	6722	6723	Should I learn python or Java first?	If I had to choose between learning Java and Python, what should I choose to learn first?

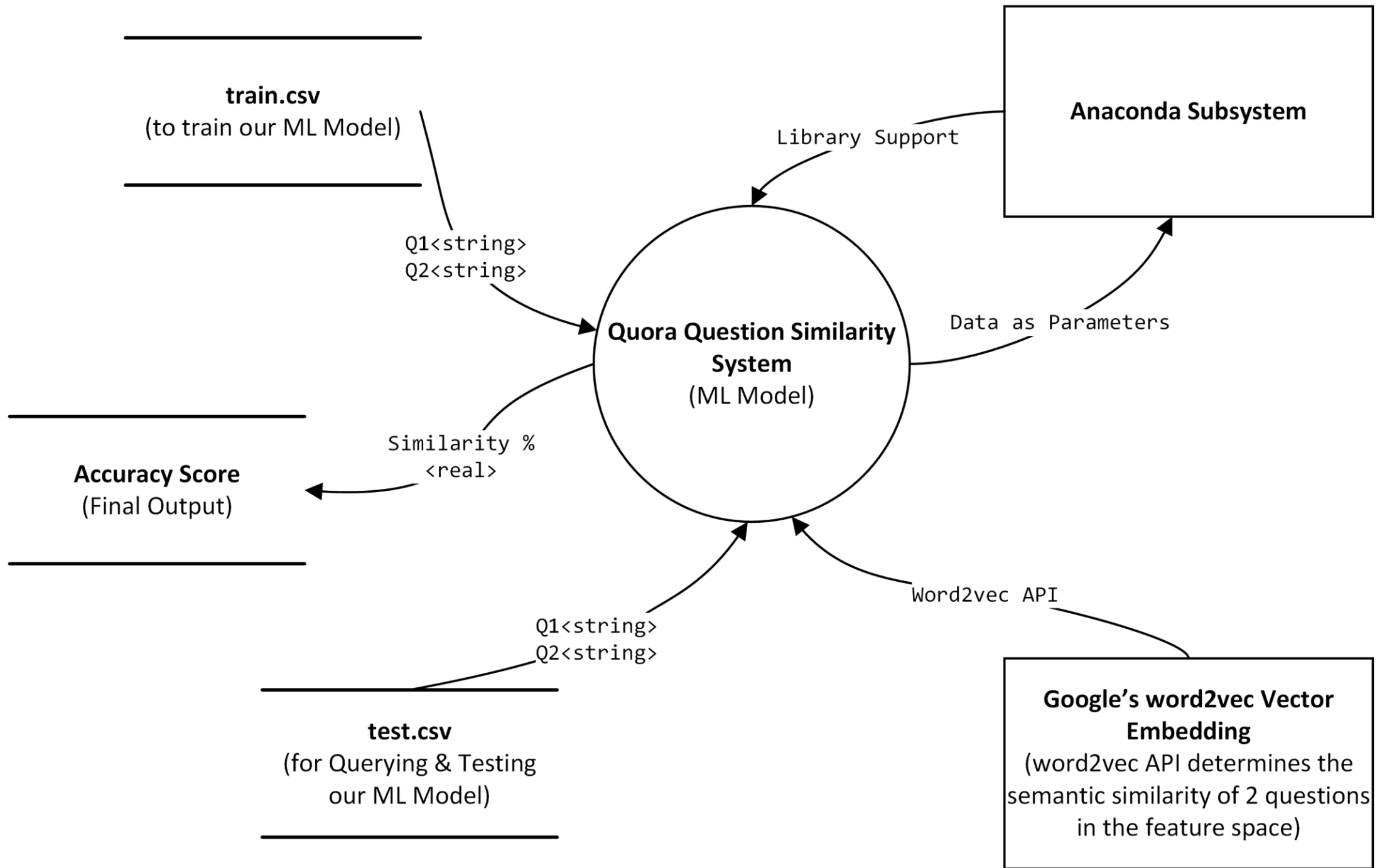
Problem Output [Expected]

id	qid1	qid2	question1	question2	is_duplicate
447	895	896	What are natural numbers?	What is a least natural number?	0
1518	3037	3038	Which pizzas are the most popularly ordered pizzas on Domino's menu?	How many calories does a Dominos pizza have?	0
3272	6542	6543	How do you start a bakery?	How can one start a bakery business?	1
3362	6722	6723	Should I learn python or Java first?	If I had to choose between learning Java and Python, what should I choose to learn first?	1

Modules

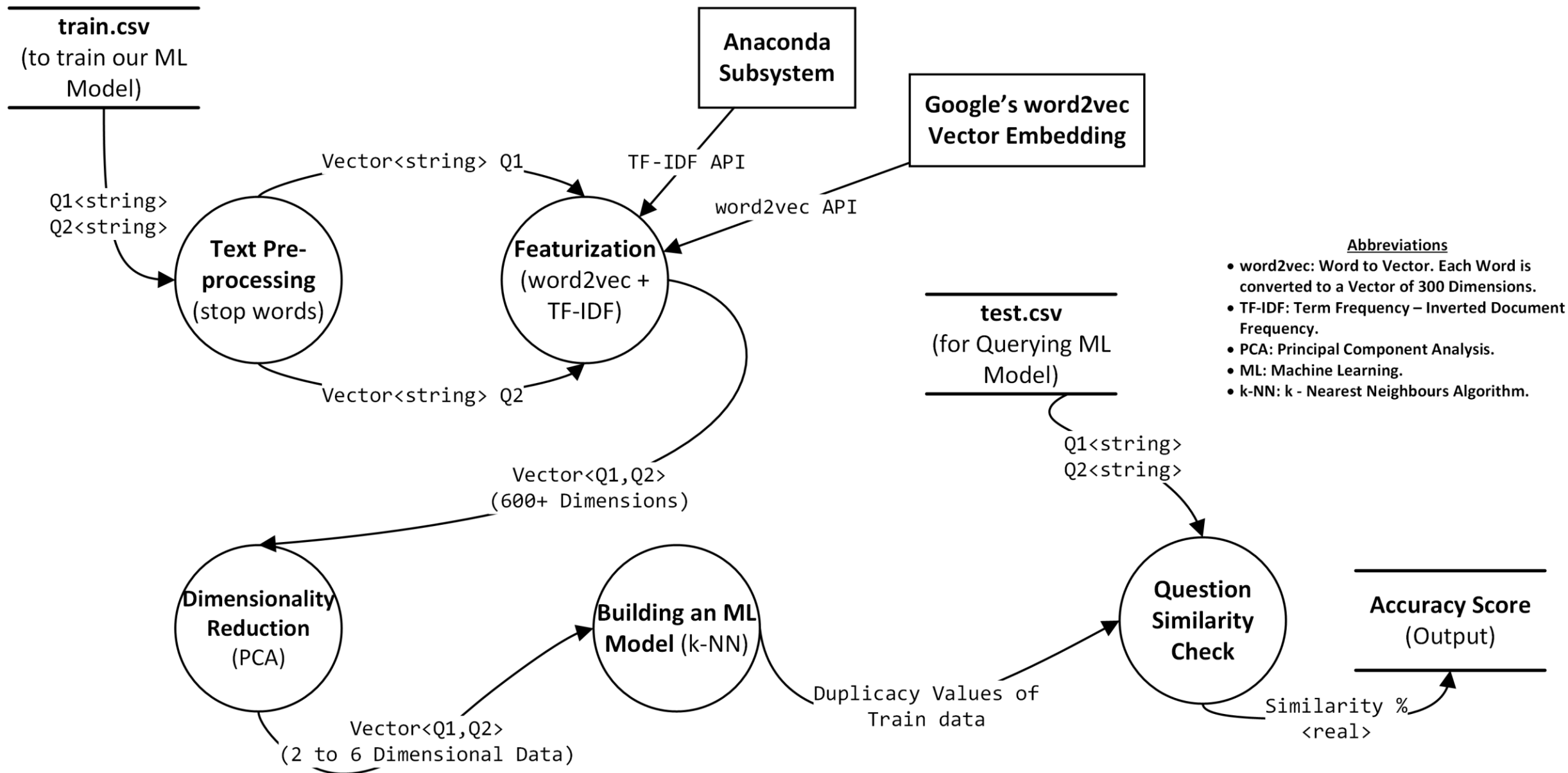
- Text Pre-processing: *text tokenization, uniform text casing, stop word removal, lemmatization, stemming*, etc.
- Featurization / Vectorization of the texts [$q1$ and $q2$]: Convert the given $q1$ and $q2$ into Word Vectors using word2vec. Word vectors are positioned in the vector space such that words that share common contexts in the corpus are located in close proximity to one another in the space.
- Apply a Machine Learning Model [using k – Nearest Neighbours algorithm] to get the Expected Output.

Data Flow Diagram – Level 0



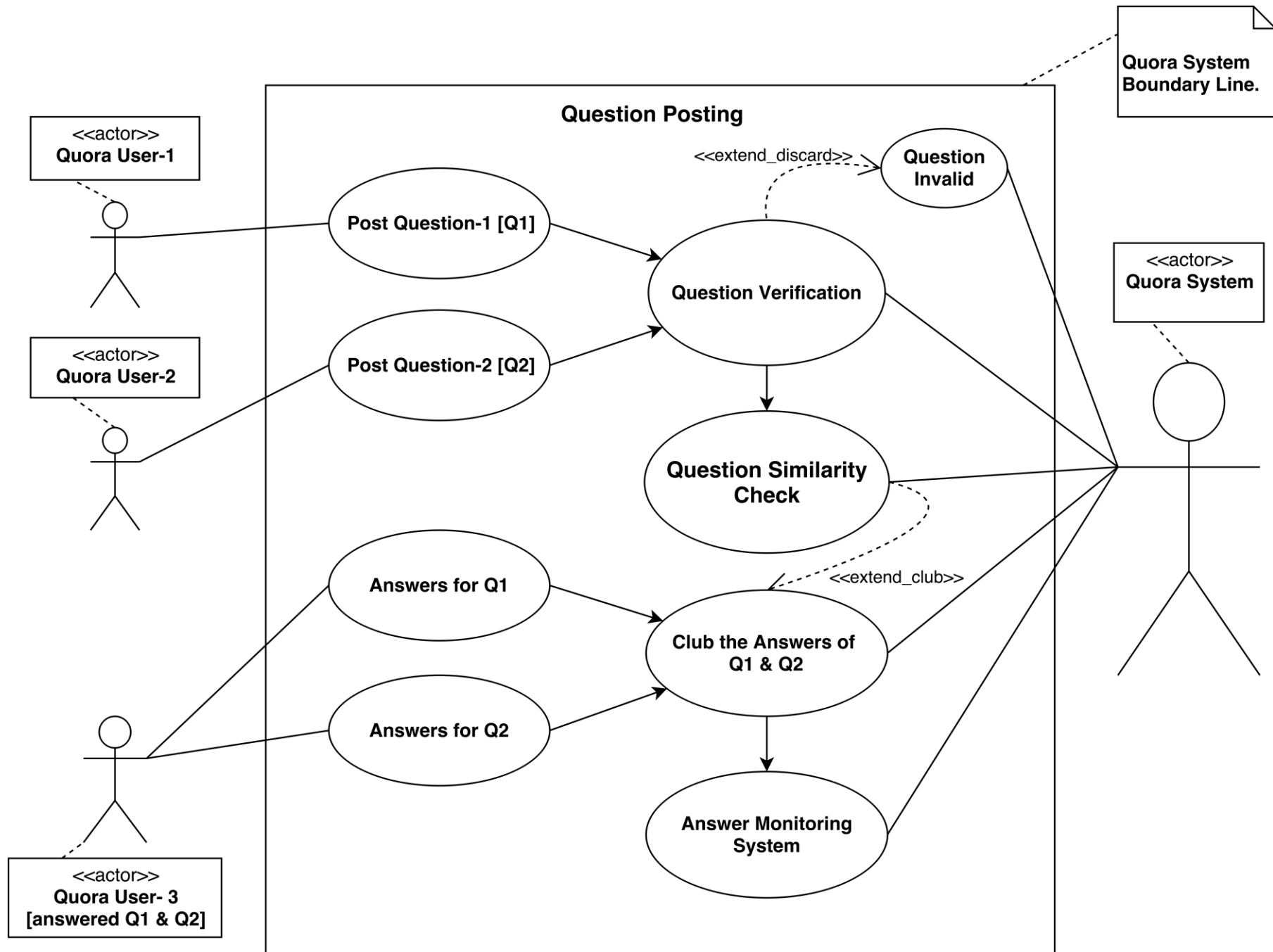
Data Flow Diagram – Level 1

Quora Question Similarity System



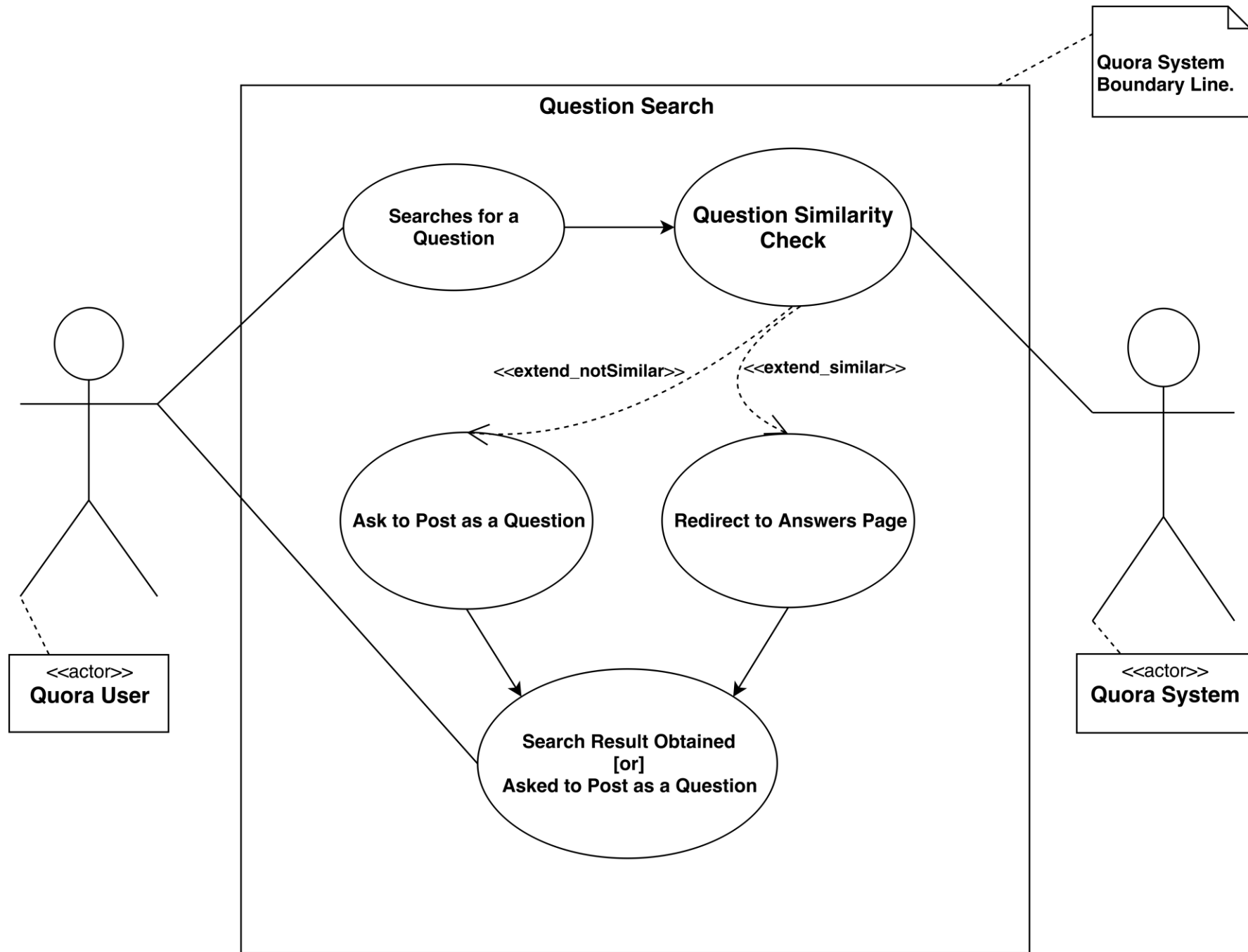
Use Case Diagram – 1

[Question Posting]

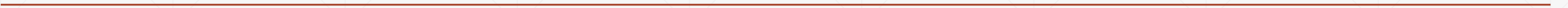


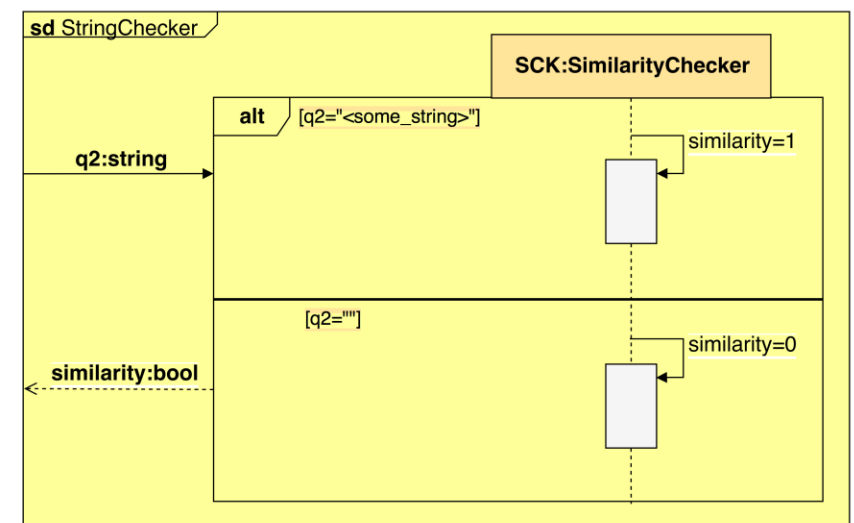
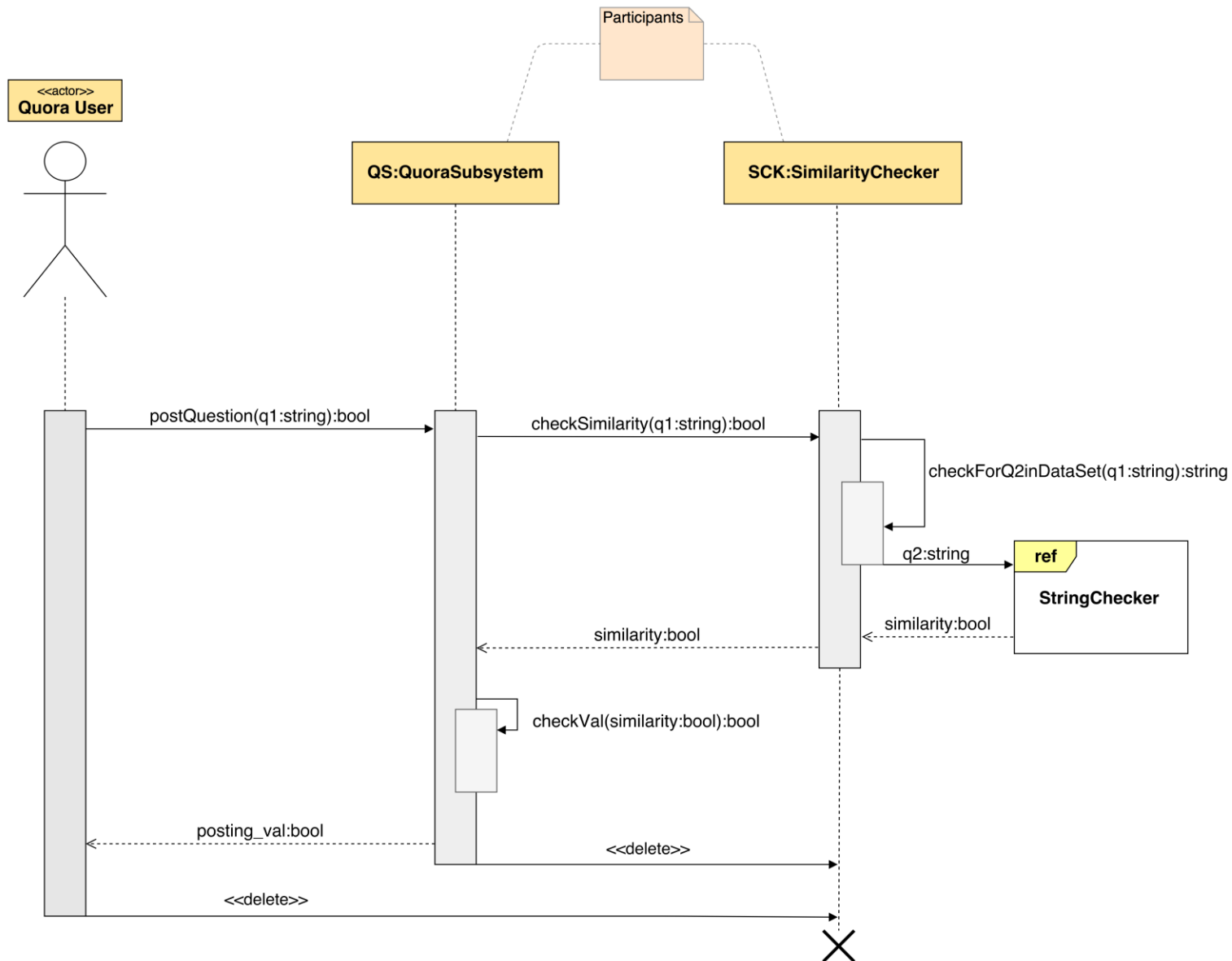
Use Case Diagram – 2

[Question Search]

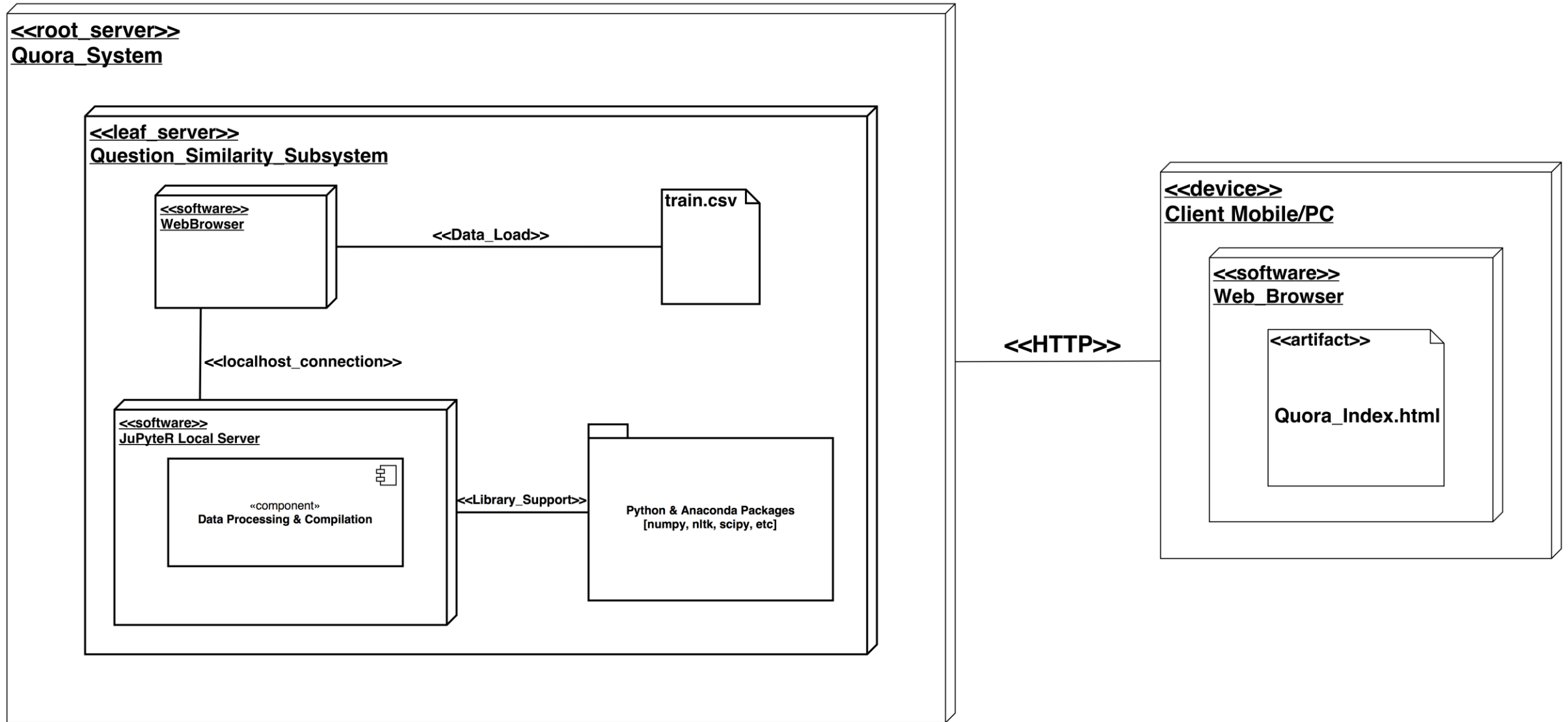


Sequence Diagram





Deployment Diagram



Algorithm Information

- In ***k*-NN** classification, the output is a class membership. An object is classified by a majority vote of its neighbours, with the object being assigned to the class most common among its ***k*** nearest neighbours (***k*** is a positive integer, typically small). If ***k*** = 1, then the object is simply assigned to the class of that single nearest neighbour.
- Distance between 2 points in the feature space, can be found by taking Cosine Distance (= 1 - Cosine Similarity) between those 2 points.

Thank You!

**Any
Questions/Queries?**
