# CS6370
# NLP

Course Project
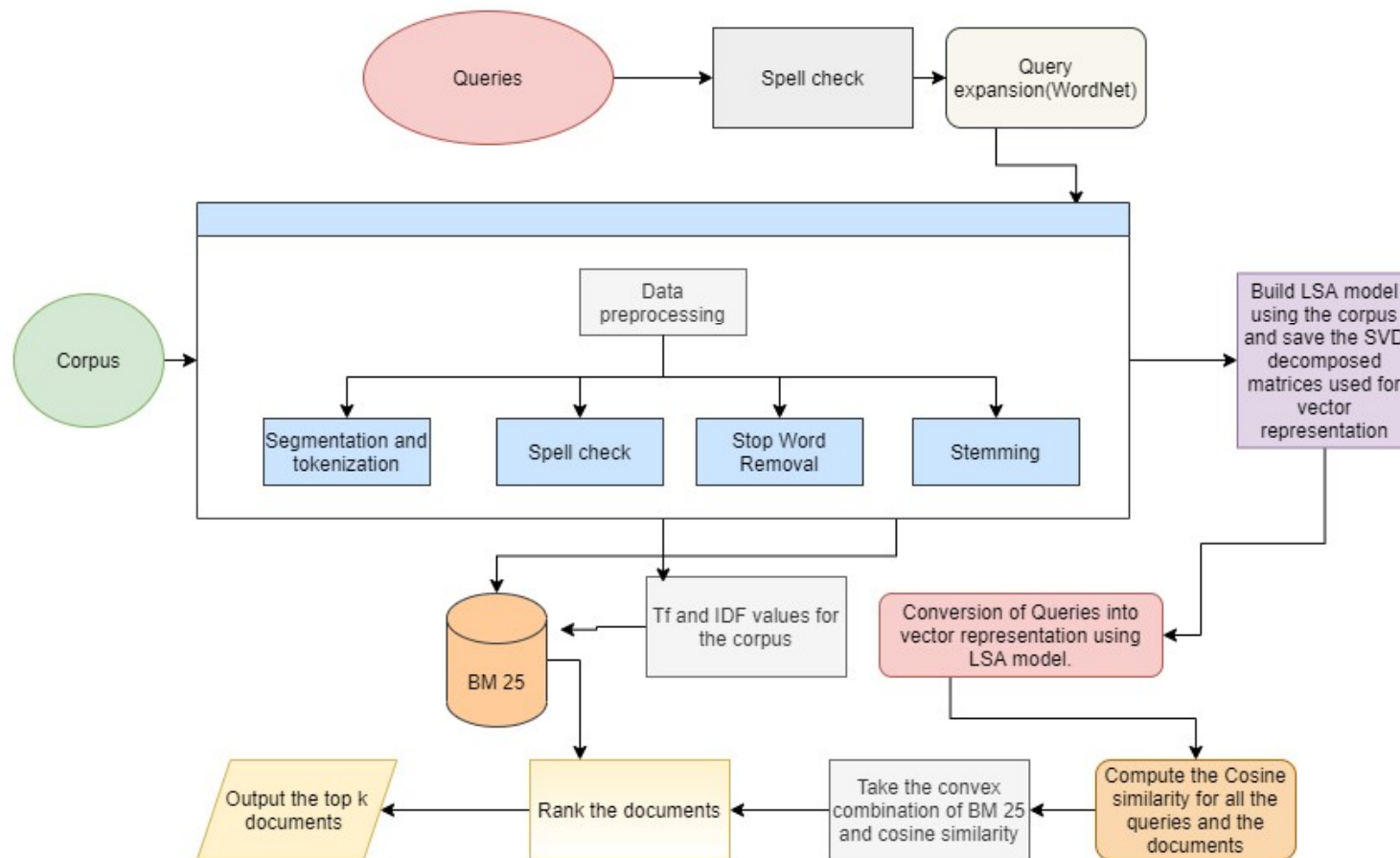
- Naveen Vakada (CS20S012)

- Arup Das (CS20S016)

# Limitations of the current search engine

Fails to classify documents containing synonymous words or documents with the same context as relevant documents.

The title of the document was not used.

Does not consider the co-occurrence relation between the terms in similar documents.

# Proposed methodology

# Completed tasks

Implemented BM25

$$\sum_{t \in q} IDF(t) * \frac{(k_1 + 1) * tf(t, d)}{k_1[(1-b) + \frac{b \cdot dl}{avdl}] + tf(t,d)} * \frac{(k_2 + 1) * tf(t, q)}{k_2 + tf(t, q)} \text{ where } IDF(t) = log\left(\frac{N}{df(t)}\right)$$

In the presence of full relevance judgement, $IDF(t) = log\left(\frac{\frac{r_t + 0.5}{R - r_t + 0.5}}{\frac{df(t) - r_t + 0.5}{N - df(t) - R + r_t + 0.5}}\right)$

Used titles of the document to improve the performance.

Completed reading about how to implement LSA in python.

# Pending Work

Implementing LSA

Implementing Statistical hypothesis testing to compare the mean nDCG of the old method and the new proposed method.

Implementing query expansion.

Implementing Spell Check.