

- The retrieved documents given by the system depend on how the query is being entered by the user

- If there are any spelling errors or incomplete sentences, the system will not be able to give the correct relevant documents
 - Eg. Query: Mone ban
 Relevant Doc: Money is stored in the bank.
 This relevant document will not be retrieved as the words don't match due to the spelling errors.
 - **Title:** Titles of documents are not considered and hence their information is lost
2. State your hypothesis for addressing the above limitation(s). Describe how you would realize the above hypotheses in your search engine.

Answer:

- **High dimensionality:**
 - We will use LSA or other methods to reduce the dimensionality while preserving enough information
 - Embeddings like Word2Vec and BERT embeddings will be tried instead of TF-IDF
 - **Realization:** The tf-idf calculation step will be replaced with word2vec / bert embeddings in informationRetrieval.py, LSA will be used after to reduce dimensions
- **Word Order:**
 - We will be using n-gram (mostly n=2) for preserving the word order and co-occurrence information
 - **Realization:** In the tokenization step (tokenizer.py), we will add a n-gram tokenizer function to the Tokenizer class and implement the n-gram tokenizer using nltk library
- **Synonymy:** Embeddings like Word2Vec/ BERT will address this as synonymous word vectors will be similar. Also LSA helps in capturing this information (2 words meaning same thing).
- **Polysemy:** This is partially addressed by using LSA and n-gram (since neighborhood information is captured)
- **Human Error:**
 - We will use a spell checker on the query to prevent misspelled words from giving wrong results
 - **Realization:** In the lemmatization step (inflectionReduction.py), we will add a spell checker and corrector (from library). Hence all queries and documents will be spelling checked and corrected before they are tokenized
- **Title:** We will add the terms in the title also to the document and weight their vector values by a $K > 1$ (around 2 or 3)

3. Describe how you would evaluate your system.

Answer: To evaluate the system,

We will continue to use the Cranfield dataset.

We will use MAP @ k and Mean nDCG measures since they provide a better evaluation metric than others. Since in qrels.json, we are provided the relevant scores of the relevant documents, we can make use of Mean nDCG to get a good evaluation of our system. MAP @ k is also calculated for comparing it with nDCG.