# Task 3: Implementing A Named Entity Recognizer
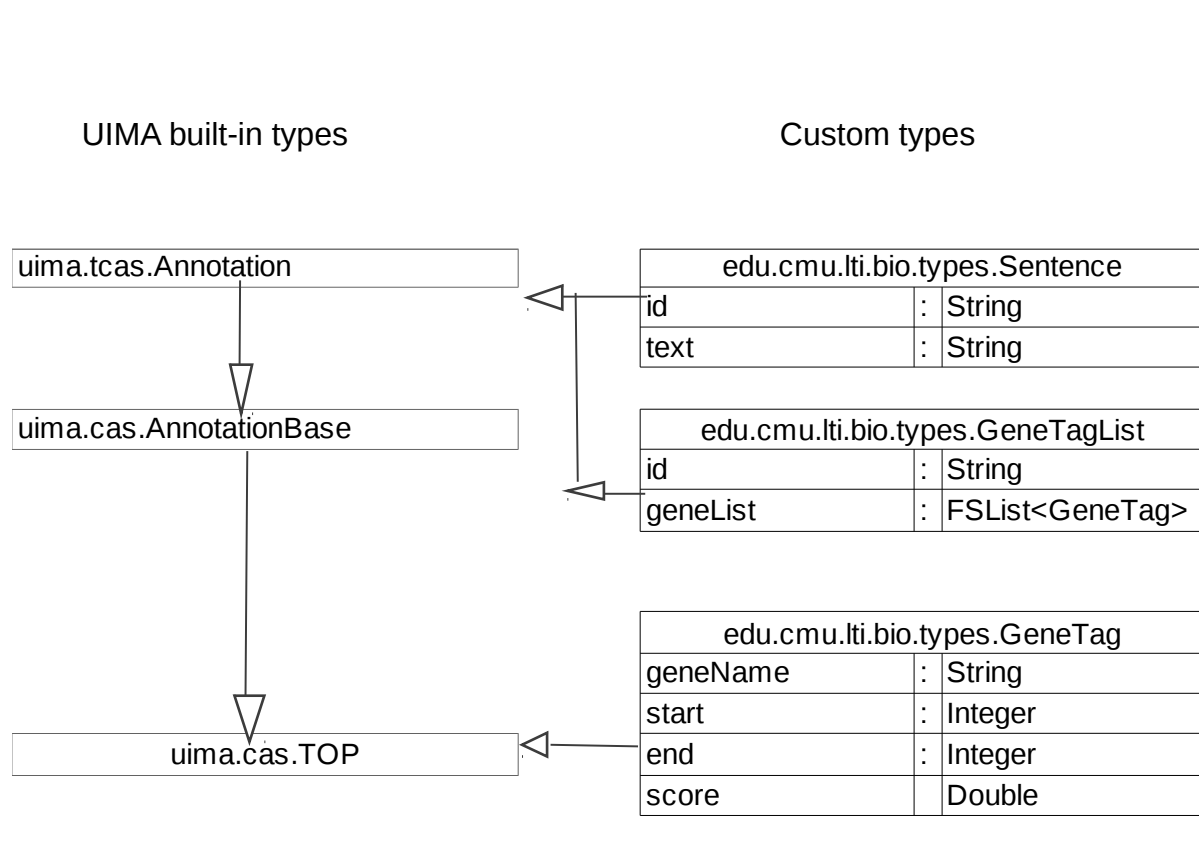
## 1. High-level architecture



## 2. UML Class Diagram

## 3. Type System

| UIMA built-in types | Custom types |
|---|---|



## 4. NLP Tools Used
* Stanford-Core-NLP toolkit was used for following tasks in the UIMA pipeline:
  – Tokenizer
  – Part-Of-Speech Tagger
  – Name Entity Recognizer
* Lucene was used to index N-Grams extracted from Corpus.

## 5. External Resource
GENETAG Training data was used for creating dictionary.
Genomics TREC -2006 data was used for extracting 5-Grams.

## 6. N-Gram Model
Genomics Track data was pulled from BIO-OAQA repository and upto 5-gram analysis was done. All N-Grams were indexed using Lucene to create index for faster search. Following steps were followed for determining whether given noun-phrase qualify for gene name:

  – Extract noun-phrases using Standford-CoreNLP toolkit
  – Search each noun-phrase in Lucene Idex
  – Apply threshold on relevance score returned by Lucene for a noun-phrase in order to

select the best gene name based on highest relevance score

## 7. Future Work
The threshold used for filtering best gene names is by gutt feeling. In future, I would look for detailed evaluation mechanism to figure out where system performed good and where it should have perform better. Would also look for better noun-phrase extraction strategy.