

## Implementing a Simple Information Processing Task with UIMA SDK

In implementing this QA system, I used the example given in class as a point of reference, and made several modifications, taking several considerations into account.

Using the `deis_types.xml` file as a starting point, I first considered the requirements of the information pipeline, and then broke it into several stages (Analysis Engine components) which are described below.

Stage / Description	Input Types	Output Types
<b>1. TestElementAnnotator</b> <ul style="list-style-type: none"> <li>The TestElementAnnotator annotates each input document. For each line of the input file, it determines if that line is a Question, or an Answer, and annotates the span accordingly. <ul style="list-style-type: none"> <li>For each Answer annotation, the TestElementAnnotator will sets the <code>isCorrect</code> feature to yes / no (1/0), based upon the flag given in the input file.</li> </ul> </li> </ul>	N/a (text file as input)	<b>Question, Answer</b>
<b>2. TokenAnnotator</b> <ul style="list-style-type: none"> <li>The TokenAnnotator converts each Question and Answer annotation into a set of Token annotations. The Token annotations are then used in downstream processing.</li> <li><b>Please note:</b> TokenAnnotator.java utilizes the Stanford libraries utilized in HW0. Since they already exist on the <code>mu.its</code> repository, I assume there will be no problems with loading them on the repository side, but if you run into any issues, please feel free to contact me. Specifically, the script imports the following: <ul style="list-style-type: none"> <li><code>edu.stanford.nlp.ling.Word;</code></li> <li><code>edu.stanford.nlp.objectbank.TokenizerFactory;</code></li> <li><code>edu.stanford.nlp.process.Tokenizer;</code></li> <li><code>edu.stanford.nlp.process.PTBTTokenizer.PTBTTokenizerFactory;</code></li> </ul> </li> <li>Note: I had hoped to add the following features to the Token type – these would help build a better machine learning (or rules-based) model, however getting up to speed on Java and UIMA syntax and structures took longer than anticipated. If we do a future iteration of the QA system, I hope to add the following, using Stanford CoreNLP: <ul style="list-style-type: none"> <li>A <b>stem</b> (or <b>lemma</b>) feature containing the stem (or lemma) of the token. For n-gram models, using stemmed/lemmatized words can contribute to a smoother n-gram that can better generalize to text that contains new forms of previously seen words (e.g. “Does John <u>love</u> Mary?” → “John <u>loves</u> Mary.”)</li> <li>A <b>pos</b> feature containing the Part of Speech tag for the token. Introducing this as a feature in an n-gram model would better distinguish between sentences like “John loves Mary” vs. “Mary loves John”.</li> </ul> </li> </ul>	<b>Question, Answer</b>	<b>Token</b>
<b>3. NgramAnnotator</b> <ul style="list-style-type: none"> <li>The NgramAnnotator converts sequences of Token annotations into Ngrams.</li> <li>A parameter allows the user to specify the which NGram orders will be created. <ul style="list-style-type: none"> <li>Ex: 1: setting the parameter to “1,2,3” will create unigrams, bigrams, and trigrams.</li> <li>Ex: 2: setting the parameter to “2,4” would create bigrams and 4-grams.</li> </ul> </li> <li>Each time an <b>NGram</b> annotation is created: <ul style="list-style-type: none"> <li>The <b>start</b> feature is set equal to the <b>start</b> of the first Token in the NGram.</li> <li>The <b>end</b> feature is set equal to the <b>end</b> of the last Token in the NGram.</li> </ul> </li> </ul>	<b>Token</b>	<b>NGram</b>

<ul style="list-style-type: none"> <li>○ The <b>elements</b> feature is an array that contains pointers to each Token in the NGram.</li> <li>○ The <b>elementsType</b> feature is set to "Token". (Note: this could become parameterized in future releases, if types other than Tokens are used in NGram models)</li> <li>○ The <b>ngramOrder</b> feature tells whether the NGram is a unigram ("1"), bigram ("2"), etc – this may be used in future releases if I develop a machine learning model that weights the NGrams differently.</li> <li>○ The <b>source</b> feature indicates whether the NGram came from a Question or an Answer. This is used in downstream processing to compare Answer NGrams against their corresponding Question NGrams.</li> </ul>		
<b>4. Answer Scoring</b>	<b>Question, Answer, NGram</b>	<b>AnswerScore</b>
<ul style="list-style-type: none"> <li>• The AnswerScorer loops through each <b>Question</b> (there should only be one in each document). For each Question, the AnswerScorer loops through the candidate <b>Answer</b> annotations and scores them.</li> <li>• For this implementation, the scoring methodology is a simple NGram overlap rule: for each NGram contained in the candidate Answer, the AnswerScorer looks to see if that NGram is also contained in the Question. If so, it increments a "matching_ngrams" counter by 1. The final score feature is simply the count of "matching_ngrams" / the total number of n-grams contained in the candidate Answer.</li> <li>• Note: the approach above corresponds to a interpolated n-gram model, with even weights between each n-gram order. For this implementation, I used 1-grams, 2-grams, and 3-grams, so the weights are 1/3 unigram, 1/3 bigram, 1/3 trigram. This is a very simple model. Per the note above (re: adding features to the Token type), in a future release, I would like to have a richer feature set, and perhaps use machine learning techniques to build a more robust scoring mechanism.</li> <li>• When the score is returned, a new <b>AnswerScore</b> annotation is created. This annotation spans the text of the answer sentence. In addition, it has a <b>score</b> feature (that stores the score), as well as an <b>answer</b> feature (that has a pointer to the Answer object that was used during scoring).</li> </ul>		
<b>5. Evaluator</b>	<b>AnswerScore</b>	N/a (results print to screen)
<ul style="list-style-type: none"> <li>• The Evaluator evaluates the performance of the AAE by comparing the system outputs (<b>AnswerScore.score</b> values) with the gold standard outputs. It measures precision by doing the following: <ul style="list-style-type: none"> <li>○ 1) Ranks the Answers according to AnswerScore.score (descending).</li> <li>○ 2) Selects the top N Answers (where N = the total number of correct answers that were possible for that Question).</li> <li>○ 3) Measures precision by computing how many of the top N ranked Answers were actually correct (based upon the gold standard).</li> </ul> </li> <li>• After performing these calculations, it prints the results to screen.</li> </ul>		

I enhanced the type system as follows:

- For the **NGram** type:
  - Added the **ngramOrder** feature, which tells whether the NGram is a unigram ("1"), bigram ("2"), etc – this may be used in future releases if I develop a machine learning model that weights the NGrams differently.

- Added the **source** feature, which indicates whether the NGram came from a Question or an Answer. This is used in downstream processing to compare Answer NGrams against their corresponding Question NGrams.

**Figure 1: Types and Features**

Type Name or Feature Name	SuperType or Range
<input type="checkbox"/> edu.cmu.deiis.types.Annotation	uima.tcas.Annotation
casProcessorId	uima.cas.String
confidence	uima.cas.Double
<input type="checkbox"/> edu.cmu.deiis.types.Answer	edu.cmu.deiis.types.Annotation
isCorrect	uima.cas.Boolean
<input type="checkbox"/> edu.cmu.deiis.types.AnswerScore	edu.cmu.deiis.types.Annotation
score	uima.cas.Double
answer	edu.cmu.deiis.types.Answer
<input type="checkbox"/> edu.cmu.deiis.types.NGram	edu.cmu.deiis.types.Annotation
elements	uima.cas.FSArray
elementType	uima.cas.String
ngramOrder	uima.cas.Integer
source	uima.cas.String
edu.cmu.deiis.types.Question	edu.cmu.deiis.types.Annotation
edu.cmu.deiis.types.Token	edu.cmu.deiis.types.Annotation