# 11-791: Design and Engineering of Intelligent Information Systems
## Project Individual 2: UIMA Type System

Keith Maki

September 14, 2015

## 1 Introduction

In this report, I present my design and implementation for the logical data model pursuant to the PI2 sample information processing task outlined in figure 1. The report is structured as follows: section 2 describes the overall type system hierarchy and data class design; section 3 provides discussion of important design considerations and limitations to the selected approach; section 4 presents the UIMA type system descriptor and other relevant project files; and section 5 concludes the report.

## 2 Type System Design Process

In designing my first type system, I found it useful to focus initially on top-down organization. I began by describing the types of annotations produced for each test element by each of the analysis engine steps shown in Figure 1. Next, I identified primary data elements which would need to be represented in the type system (e.g. Question, Answer, Tokenization, Score). I then built a preliminary type system using the Eclipse interface. This type system centered around a generic datatype for each test element, which would be populated with additional annotations as the analysis pipeline continued to process the associated test element.

This solution likely could have been instantiated in an analysis framework, but it suffered from two major flaws. First, it resulted in many design choices which limited the ability of the type hierarchy to generalize to additional domains. My top-down approach
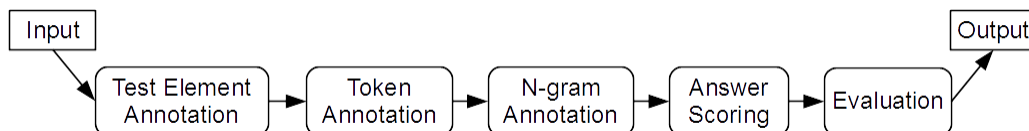


Figure 1: High-level architecture of the system scaffolded by this project.
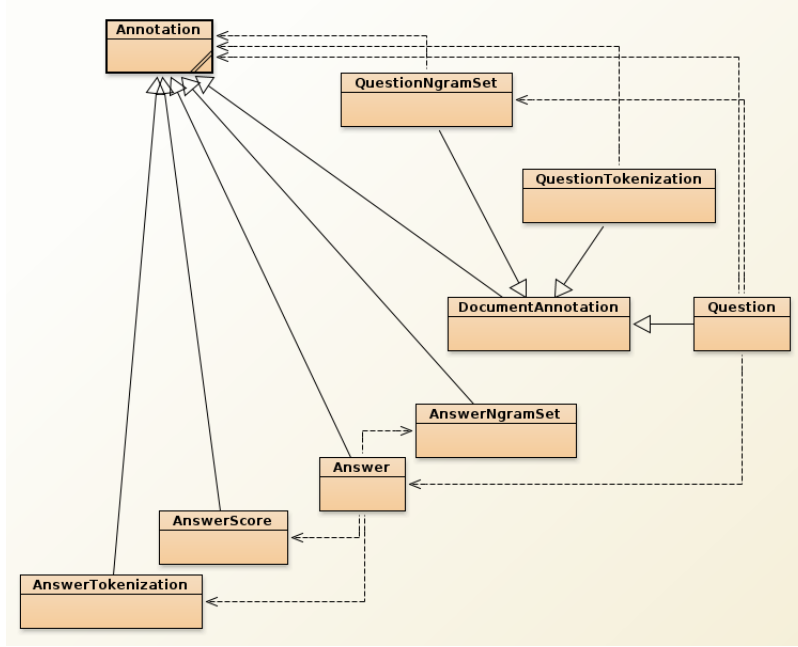
Figure 2: BlueJ Class Diagram for Preliminary Type System

inherently began with the sample task outline and refined the description from there, leaving it unable to easily capture broader notions. Second, it was difficult to take advantage of substructure present in the task, in that classes could not easily inherit from other classes, although many of my types began to look more similar as I continued to add detail to their descriptions.

An inheritance diagram of my preliminary type system is shown in Figure 2. Although I ultimately decided to scrap this version of my type system hierarchy, it was helpful to give myself an initial view of what was possible within the framework in order to inform a more intelligent design. It is important to note the weak notions of inheritance that plagued this design, with a heavy reliance on prebuilt types and an unfounded compulsion to separate Question and Answer-related data types. Notice also the unused QuestionTokenization class (which I believe was an error in this commit); it is easy to overlook missing dependencies like this in a poorly designed hierarchy.

Once I had gotten a feel for the kinds of elements which would be useful to complete the top-level hierarchy, I rebuilt my type system using a bottom-up perspective. This allowed me to ensure my design included all necessary subtypes before I extended or included them as fields in other types. During this step, I also completed the documentation strings for each field and type. In taking a broad view of the types of data that would be useful to annotate in general systems involving the analysis engine tasks outlined by Figure 1, I feel I was better able to conceptualize not only the aspects of each data element which were useful to extend and inherit within my type system design, but also those aspects which would enable them to be extensible and useful in more general analysis frameworks.
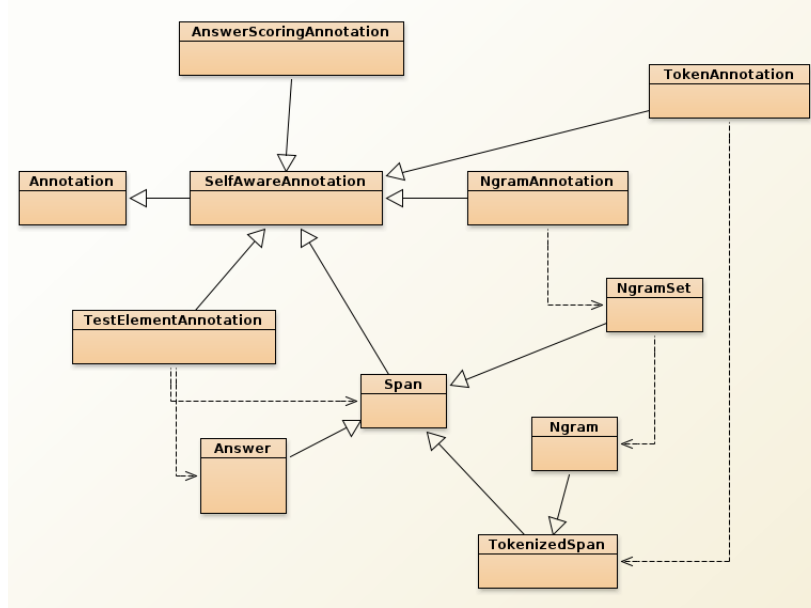
Figure 3: BlueJ Class Diagram for Final Type System

An inheritance diagram of my final type system is shown in Figure 3. This design has much deeper motivation behind its inheritance structure and avoids the field declaration redundancies that were so prevalent in the earlier design.

# 3    Design Considerations

In developing my final type system, I aimed to ensure that the machinery required to produce the related types was as generalizable as possible. Each of the analysis engines should not depend heavily on the structure of the system overall, and should only rely on behavior from upstream engines that can reasonably be expected to be part of such an upstream engine. For example, the scorer engine should not need to reproduce Answer information (e.g. the "correct" field from the gold data) even though it might be easier for downstream engines like the evaluator if these data were in the same place as the annotations produced by those engines. Tweaking an analysis engine's designs is only justified when it becomes clear that engine does not meet the specifications drafted for its *own* requirements.

I also tried to encourage flexibility in my type system design where possible with respect to configurability on the part of the analysis engines. For example, I provide for arbitrary numbers of ngram annotations, but I require each ngram annotation keep track of its order of n so that a scorer which should only consider certain ngrams (e.g. no ngrams of length $> 5$) can easily filter the annotations to suit its needs.

Of particular note is the affordance in my Span annotations of a String text field aside from the inherited begin and end Integer indices. This allows for more general anno-

tation engines can be easily implemented using the same underlying type system. For instance, a many-to-one tokenizer could be implemented, which maps alternate spellings onto the same tokenization (e.g. "don't" → "do not", "do not" → "do not"). Similarly, tokenizers which perform disambiguation using an ontology (e.g. "Booth shot Lincoln" → "John_Wilkes_Booth shot Abraham_Lincoln") can be easily implemented by annotating the appropriate text field in each Span.

## 4   UIMA Type System Descriptor

Figure 4 provides the complete UIMA type system descriptor for the described logical data model. The concise documentation provided in the form of class and field descriptions works alongside the UIMA JCasGen automatic code generation scripts to produce functional and easy-to-work-with Java class files from this xml. The full pi2-kmaki project repository including this report be obtained at https://github.com/kortemaki/pi2-kmaki.

## 5   Conclusion

In this report, I presented my UIMA Type system in pursuit of the requirements of the PI2 assignment. The techniques employed in the design of this type system will be foundational in designing and implementing more elaborate information systems in later units of this course and in the Team Project assignments in the latter half of the semester. Beyond this, the underlying design principles serve as important building blocks for simple yet effective software engineering which will continue to provide benefits long after this course has been completed.

```xml
1 <?xml version="1.0" encoding="UTF-8"?>
2 <typeSystemDescription xmlns="http://uima.apache.org/resourceSpecifier">
3     <name>pi2-kmaki-typesystem</name>
4     <description/>
5     <version>1.0</version>
6     <vendor/>
7   <types>
8     <typeDescription>
9       <name>SelfAwareAnnotation</name>
10      <description>Annotation subclass which is aware of the annotator
            that produced it.
11
12 Fields inherited from Annotation:   begin, end</description>
13      <supertypeName>uima.tcas.Annotation</supertypeName>
14      <features>
15        <featureDescription>
16          <name>annotator</name>
17          <description>The name of the annotator that produced this
                annotation.</description>
18          <rangeTypeName>uima.cas.String</rangeTypeName>
19        </featureDescription>
20      </features>
21    </typeDescription>
22    <typeDescription>
23      <name>Span</name>
24      <description>Subclass of SelfAwareAnnotation.  Annotates a span of
            text with a String that may encode information about the span of
             text.
25
26 Fields inherited from SelfAwareAnnotation:   begin, end, annotator</
      description>
27      <supertypeName>SelfAwareAnnotation</supertypeName>
28      <features>
29        <featureDescription>
30          <name>text</name>
31          <description>The text annotated by the span.
32
33 **Note that this text may differ from the portion of the annotated
      TestElement as indexed between begin and end!!</description>
34          <rangeTypeName>uima.cas.String</rangeTypeName>
35        </featureDescription>
36      </features>
37    </typeDescription>
38    <typeDescription>
39      <name>TokenizedSpan</name>
40      <description>Subclass of Span which annotates a tokenization for
            its span.
41
42 Fields inherited from Span:   begin, end, annotator, text</description>
43      <supertypeName>Span</supertypeName>
44      <features>
```

Figure 4:   The UIMA Type System Descriptor XML for the presented type system

```
45          <featureDescription>
46            <name>tokens</name>
47            <description>The tokens identified in the tokenization of this
                  Span.</description>
48            <rangeTypeName>uima.cas.FSList</rangeTypeName>
49            <elementType>Span</elementType>
50            <multipleReferencesAllowed>true</multipleReferencesAllowed>
51          </featureDescription>
52        </features>
53      </typeDescription>
54      <typeDescription>
55        <name>Ngram</name>
56        <description>Subclass of TokenizedSpan which annotates a subset of
            a tokenization corresponding to an ngram of finite length n.
57
58 Fields inherited from TokenizedSpan:   begin, end, annotator, text,
      tokens</description>
59        <supertypeName>TokenizedSpan</supertypeName>
60        <features>
61          <featureDescription>
62            <name>n</name>
63            <description>The cardinality of the ngram identified by this
                  TokenizedSpan.</description>
64            <rangeTypeName>uima.cas.Integer</rangeTypeName>
65          </featureDescription>
66        </features>
67      </typeDescription>
68      <typeDescription>
69        <name>NgramSet</name>
70        <description>Subclass of Span annotating selected ngrams for the
            given span.
71
72 Fields inherited from Span: begin, end, annotator, text</description>
73        <supertypeName>Span</supertypeName>
74        <features>
75          <featureDescription>
76            <name>ngrams</name>
77            <description>List of ngrams identified for this Span.</
                  description>
78            <rangeTypeName>uima.cas.FSList</rangeTypeName>
79            <elementType>Ngram</elementType>
80            <multipleReferencesAllowed>true</multipleReferencesAllowed>
81          </featureDescription>
82        </features>
83      </typeDescription>
84      <typeDescription>
85        <name>Answer</name>
86        <description>Subclass of span annotating an answer choice for a
            particular TestElement.  Is aware of whether it is a correct
            answer choice or not.
87
88 Fields inherited from Span:   begin, end, annotator, text</description>
```

Figure 5:  Continued from Figure 4

```xml
 89          <supertypeName>Span</supertypeName>
 90          <features>
 91            <featureDescription>
 92              <name>correct</name>
 93              <description>Indicates correctness of this Span as an answer
                      choice under "gold" labels.</description>
 94              <rangeTypeName>uima.cas.Boolean</rangeTypeName>
 95            </featureDescription>
 96          </features>
 97      </typeDescription>
 98      <typeDescription>
 99        <name>TestElementAnnotation</name>
100        <description>Subclass of SelfAwareAnnotation holding question and
              Answer Span annotations for a TestElement.
101
102 Fields inherited from SelfAwareAnnotation:   begin, end, annotator</
        description>
103        <supertypeName>SelfAwareAnnotation</supertypeName>
104        <features>
105          <featureDescription>
106            <name>question</name>
107            <description>Identifies the question for this TestElement.</
                  description>
108            <rangeTypeName>Span</rangeTypeName>
109          </featureDescription>
110          <featureDescription>
111            <name>answers</name>
112            <description>Array holding annotations identifying each answer
                    choice for this TestElement.</description>
113            <rangeTypeName>uima.cas.FSArray</rangeTypeName>
114            <elementType>Answer</elementType>
115            <multipleReferencesAllowed>true</multipleReferencesAllowed>
116          </featureDescription>
117        </features>
118      </typeDescription>
119      <typeDescription>
120        <name>TokenAnnotation</name>
121        <description>Subclass of SelfAwareAnnotation holding tokenization
              annotations for a TestElement.
122
123 Fields inherited from SelfAwareAnnotation:   begin, end, annotator</
        description>
124        <supertypeName>SelfAwareAnnotation</supertypeName>
125        <features>
126          <featureDescription>
127            <name>questionTokens</name>
128            <description>Tokenization of the question for this TestElement.
                  </description>
129            <rangeTypeName>TokenizedSpan</rangeTypeName>
130          </featureDescription>
131          <featureDescription>
132            <name>answersTokens</name>
```

Figure 6: Continued from Figure 5

7

```xml
133            <description>Array holding a tokenization for each answer
                   choice for this TestElement.</description>
134            <rangeTypeName>uima.cas.FSArray</rangeTypeName>
135            <elementType>TokenizedSpan</elementType>
136            <multipleReferencesAllowed>true</multipleReferencesAllowed>
137          </featureDescription>
138        </features>
139      </typeDescription>
140      <typeDescription>
141        <name>NgramAnnotation</name>
142        <description>Subclass of SelfAwareAnnotation holding ngram
                annotations for a TestElement.
143
144 Fields inherited from SelfAwareAnnotation:   begin, end, annotator</
        description>
145        <supertypeName>SelfAwareAnnotation</supertypeName>
146        <features>
147          <featureDescription>
148            <name>questionNgrams</name>
149            <description>The ngrams identified in the question for this
                   TestElement.</description>
150            <rangeTypeName>NgramSet</rangeTypeName>
151          </featureDescription>
152          <featureDescription>
153            <name>answersNgrams</name>
154            <description>Array of ngrams identified in each answer choice
                   for this TestElement.</description>
155            <rangeTypeName>uima.cas.FSArray</rangeTypeName>
156            <elementType>NgramSet</elementType>
157            <multipleReferencesAllowed>true</multipleReferencesAllowed>
158          </featureDescription>
159        </features>
160      </typeDescription>
161      <typeDescription>
162        <name>AnswerScoringAnnotation</name>
163        <description>Subclass of SelfAwareAnnotation holding answer scoring
                 annotations for a TestElement.
164
165 Fields inherited from SelfAwareAnnotation:   begin, end, annotator</
        description>
166        <supertypeName>SelfAwareAnnotation</supertypeName>
167        <features>
168          <featureDescription>
169            <name>scores</name>
170            <description>Array of scores assigned to each answer by the
                   analysis engine.</description>
171            <rangeTypeName>uima.cas.FloatArray</rangeTypeName>
172            <multipleReferencesAllowed>true</multipleReferencesAllowed>
173          </featureDescription>
174        </features>
175      </typeDescription>
176    </types>
177 </typeSystemDescription>
```

Figure 7:  Continued from Figure 6