*Execution Architecture with CPE and Deployment*
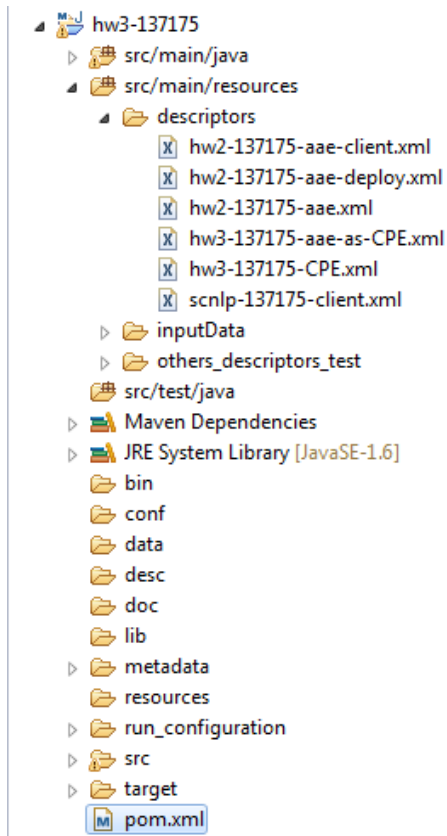*Architecture with UIMA-AS*

hw3-137175

Tania Patiño

28-03-2014

*Organization for this project is the next one:*

```
hw3-ID
|- pom.xml
'- src
   '- main
      |- java
      |  '- **/*.java      /* Java classes generated by JCasGen
      |                       and your UIMA annotators        */
      '- resources
         |- hw2-ID-aae.xml        /* your aggregate analysis engine of homework 2*/
         |- hw2-ID-aae-client.xml /* clinet descriptor of your AAE */
         |- hw2-ID-aae-deploy.xml /* deployment descriptor of your AAE */
         |- hw3-ID-aae-as-CPE.xml /* CPE descriptor to test your AAE service */
         |- hw3-ID-CPE.xml        /* CPE descriptor of your homework 2 pipeline */
         |- scnlp-ID-client.xml   /* clinet descriptor for the remote UIMA-AS service
         |- **/*.*                /* analysis engine and other resources */
         '- docs
            '- hw3-ID-report.pdf  /* your report for design */
```

- hw3-137175
  - src/main/java
  - src/main/resources
    - descriptors
      - hw2-137175-aae-client.xml
      - hw2-137175-aae-deploy.xml
      - hw2-137175-aae.xml
      - hw3-137175-aae-as-CPE.xml
      - hw3-137175-CPE.xml
      - scnlp-137175-client.xml
    - inputData
    - others_descriptors_test
  - src/test/java
  - Maven Dependencies
  - JRE System Library [JavaSE-1.6]
  - bin
  - conf
  - data
  - desc
  - doc
  - lib
  - metadata
  - resources
  - run_configuration
  - src
  - target
  - pom.xml

On the folder others_descriptiors_test, there are some files that were tested when the experimentation of the environment UIMA-Eclipse and UIMA-AS was done.

In this task, was created a Collection Processing Engine (CPE) and also the task was to understand how to use it. Requirements were, run the pipeline with a CPE instead of the UIMA Document Analyzer.

**Task 1.1 Learning CPE**

*Resources CPE*

**Basic concepts and usage about CPE from** *Chapter 2. Collection Processing Engine Developer's Guide* **(http://uima.apache.org/d/uimaj-2.4.0/tutorials_and_users_guides.html#ugr.tug.cpe).**

**Manual for CPE GUI (http://uima.apache.org/d/uimaj-2.4.0/tools.html#ugr.tools.cpe).**

The acronym CPE, corresponds to a *Collection Processing Engine,* which one ilustrates the management of the data flow that goes between different types of components that make up a CPE. Each component form part of this data flow. There are many components, for example:

1.**CAS Initializer:** is a component to populate a CAS from a document. This one is focus on the population of a document, how a document could be filled by data and what kind of data.

2.**CAS Consumer:** there could be more than one component like this one, the main idea about this, is to consume the enriched CAS that was generated by the sequence of Analysis Engines. Types of CAS Consumers are search engine like Google, index (this one is related to how to index the data in a database or what?) or a database like Oracle.

3.**Collection Reader:** interfaces to a collection of documents that should be analyzed. Maybe this is related to XML files and the analysis about every structure of the data like tags <building></building>, labels and keywords, etc.

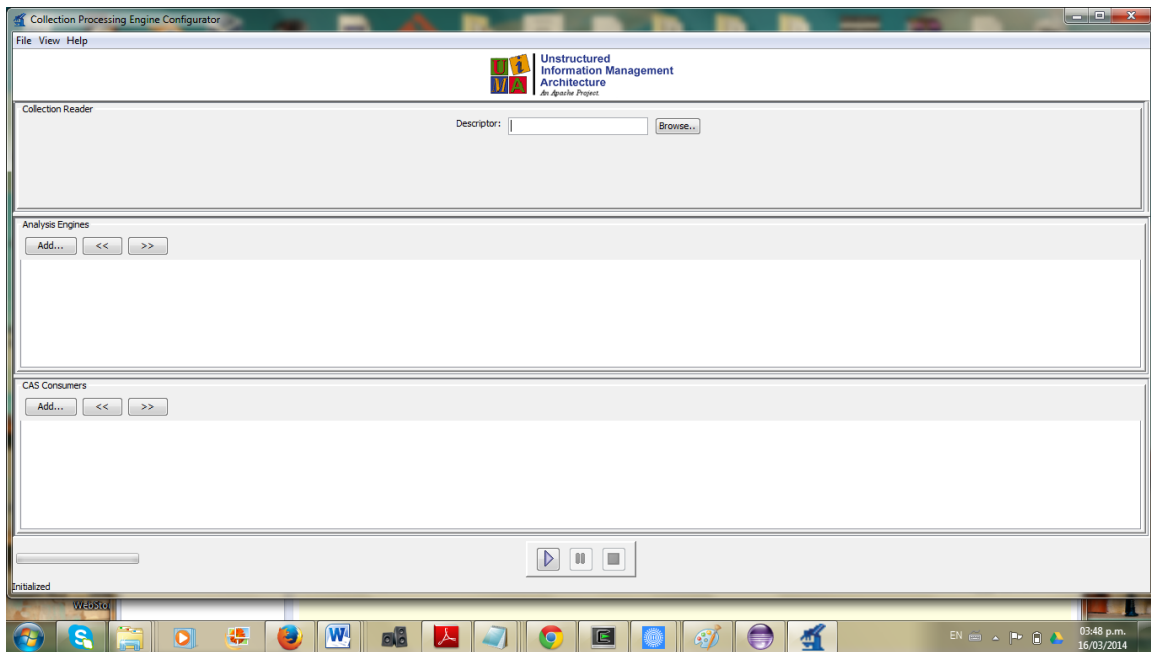4.**Analysis Engine:** this one takes a CAS analyses takes its contents and produce enriched CAS.

Secondly, following the next link related to how to make a CPE step by step, and also there is a explanation about exploration and applications of this tool: http://uima.apache.org/d/uimaj-2.4.0/tools.html#ugr.tools.cde

In the next description, there is  an explanation related to the CPE Configuration steps and CAS viewer.

1.  Open Eclipse-IDE then go to Menu option ->
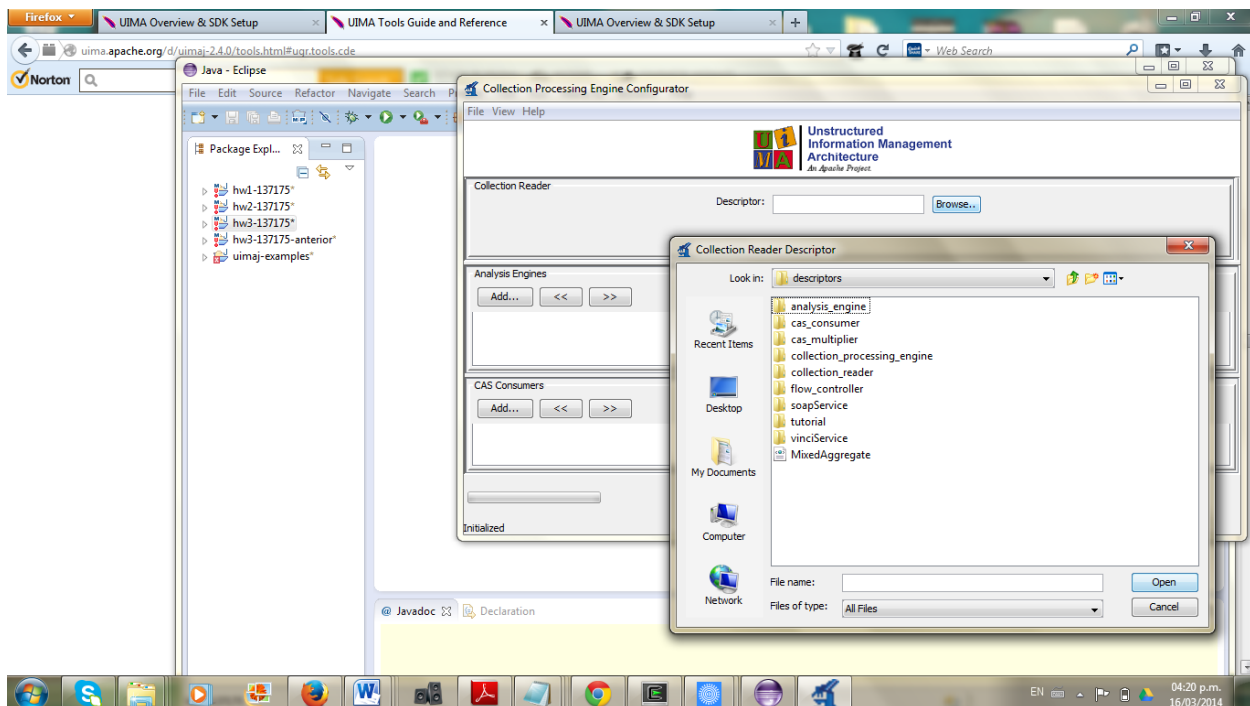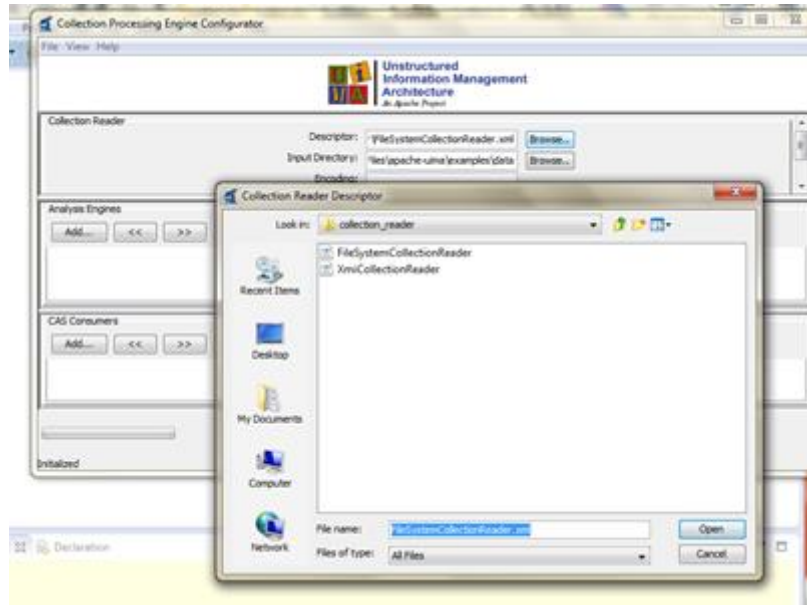2.  Run -> Run Configurations -> Select UIMA CPE GUI, then click on Run

Then the next window related to Collection Processing Engine Configurator is open, where a descriptor have to be select the descriptor and more options are activated.



Then making some experiments with the files that come by default on the descriptors folder. Finally make a selection following this path:

Selecting Collection Reader Descriptor-> collection_reader->FileSystemCollectionReader.xml

The collection reader needs to be general enough to establish a connection to the file source and open a stream to read the content. Creating a org.apache.uima.tools.components.FileSystemCollectionReader directly, and only need to write a Collection Reader descriptor.

*Descriptor:C:\ProgramFiles\uimaj-2.5.0-bin\apache uima\examples\descriptors\collection_reader\FileSystemCollectionReader.xml*

*Input Directory: C:\Program Files\apache-uima\examples\data*

*Analysis Engines:*

*Aggregate TAE- Name Recognizer and Person Title Annotator*

*Cas Consumer: XCAS Writer Cas Consumer*

Following the first steps, appears many issues, some were created because the wrong configuration paths or the wrong calling of a file for a descriptor option.



Document Analyzer with UIMA making test and experiments that ran all right:

With the Document Analyzer with UIMA  i made tests and experiments, that ran all right: testElementAnnotator.xml example was running fine and after some analysis of some files that came by default like q001.txt.xmi and q002.txt.xmi the entire test works.



Task 1.2 Creating and Running your CPE (25 pts)

1. (10 pts) The collection reader needs to be general enough to establish a connection to the file source and open a stream to read the content.
2. Task, consider a folder of files located on the file system, where apply org.apache.uima.tools.components.FileSystemCollectionReader directly, and write a Collection Reader descriptor appropied.

   Some lines about the Collection Reader Descriptor that was created (CollectionReaderDescriptor.xml):

```xml
<?xml version="1.0" encoding="UTF-8"?>
<!-- Task 1
     Section 1. Consider the folder of files located on the file system, wich means you can use:
     org.apache.uima.tools.components.FileSystemCollectionReader.
     Then write a Collection Reader Descriptor to fit the needs -->
<collectionReaderDescription xmlns="http://uima.apache.org/resourceSpecifier">
  <frameworkImplementation>org.apache.uima.java</frameworkImplementation>
  <implementationName>org.apache.uima.examples.cpe.FileSystemCollectionReader</implementationName>
  <processingResourceMetaData>
    <name>FileSystemCollectionReaderDescriptor</name>
    <description>FileSystem Collection Reader Descriptor that uses or not CAS Initializer.</description>
    <version>1.0</version>
    <vendor>Apache Enterprise Foundation</vendor>
    <configurationParameters searchStrategy="none">
      <configurationParameter>
        <name>InputDirectory</name>
        <description>this directory contains inputs files</description>
        <type>String</type>
        <multiValued>false</multiValued>
        <mandatory>true</mandatory>
      </configurationParameter>
    <configurationParameter>
        <name>Encoding</name>
        <description>Encoding parameter. The idea is that CAS is responsible to deal with character
encoding issues.</description>
        <type>String</type>
        <multiValued>false</multiValued>
        <mandatory>false</mandatory>
      </configurationParameter>
```
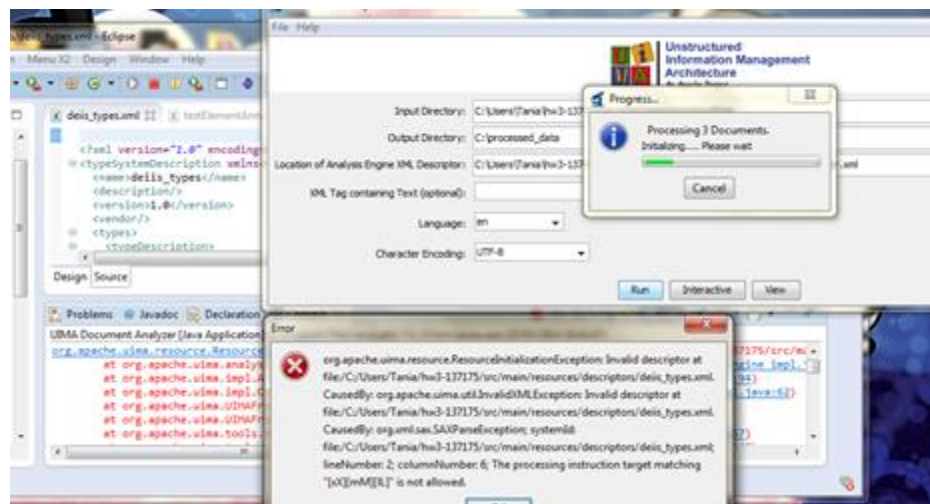
```xml
<configurationParameter>
  <name>InputDir</name>
  <description>C:\Users\Tania\hw3-137175\src\main\resources\inputData</description>
```

**Note:** the Java class called (FileSytemCollectionReader.java), was implemented  and refactored and it is related to FileSytemCollectionReader which analyze the file CollectionReaderDescriptor.xml



Error message that appears, this one occurred because there was a white space at the beginning of the xml code.



*Error message: org.apache.uima.ResourceInitializationException: Invalid description at file c:/users/Tania/hw3-137175/src/main/resources/descriptors/deiis_types.xml*
*Caused by: org.xml.sax.SAXParseException.*

*Delete White space and runs all right!!*

3. (10 pts) You are required to create a Cas Consumer based on the Evaluator component of homework 2, and include it in your CPE pipeline.

```xml
<?xml version="1.0" encoding="UTF-8"?>
<casConsumerDescription xmlns="http://uima.apache.org/resourceSpecifier">
  <frameworkImplementation>org.apache.uima.java</frameworkImplementation>

  <implementationName>edu.cmu.deiis.types.FileSystemCollectionReader</implementationName>
  <processingResourceMetaData>
    <name>casConsumerDescriptor</name>
    <description>This is related to print answerScores</description>
    <version>1.0</version>
    <vendor>Apache Enterprise Foundation</vendor>
    <configurationParameters>
      <configurationParameter>
        <name>Outputdirectory</name>
        <description>This is the directory where XMI files should be located</description>
        <type>String</type>
```

| Node | Content |
|---|---|
| ?-? xml | version="1.0" encoding="UTF-8" |
| ▲ e casConsumerDescription | |
|     ⓐ xmlns | http://uima.apache.org/resourceSpecifier |
|     e frameworkImplementation | org.apache.uima.java |
|     e implementationName | edu.cmu.deiis.types.FileSystemCollectionReader |
|    ▲ e processingResourceMetaData | |
|      e name | casConsumerDescriptor |
|      e description | This is related to print answerScores |
|      e version | 1.0 |
|      e vendor | Apache Enterprise Foundation |
|     ▲ e configurationParameters | |
|      ▷ e configurationParameter | |
|     ▷ e configurationParameterSettings | |
|     ▷ e typeSystemDescription | |
|     ▷ e capabilities | |
|     ▷ e operationalProperties | |
|     e resourceManagerConfiguration | |

4. (5 pts) Please name your CPE descriptor as hw3-ID-CPE.xml and put it under src/main/resources/, so that we could easily find the entry point of your pipeline.

| Node | Content |
|---|---|
| ?-? xml | version="1.0" encoding="UTF-8" |
| ▲ e cpeDescription | |
|     ⓐ xmlns | http://uima.apache.org/resourceSpecifier |
|    ▷ e collectionReader | |
|    ▷ e casProcessors | |
|    ▷ e casProcessor | |
|    ⌐-- | Configuration CPE |
|    ▷ e cpeConfig | |

```xml
<?xml version="1.0" encoding="UTF-8"?>
<cpeDescription xmlns="http://uima.apache.org/resourceSpecifier">
    <collectionReader>
        <collectionIterator>
            <descriptor>
                <import location="../../../../../../Program Files/uimaj-2.5.0-bin/apache-uima/examples/descriptors/collection_reader/FileSystemCollectionReader.xml"/>
            </descriptor>
            <configurationParameterSettings>
                <nameValuePair>
```

```
            <name>InputDirectory</name>
            <value>
                <string>C:\</string>
            </value>
        </nameValuePair>
    </configurationParameterSettings>
  </collectionIterator>
</collectionReader>

<casProcessors casPoolSize="3" processingUnitThreadCount="1"></casProcessors>

 <casProcessor name="hw2-137175-aae" deployment="integrated">
   <descriptor>
      <import location="/hw3-137175/src/main/resources/hw2-137175-aae.xml"/>
   </descriptor>
<deploymentParameters/>
```

Task 2

Deployment Architecture with UIMA-AS

Integrate a remote UIMA-AS service (Stanford CoreNLP) into the CPE pipeline, and deploy the aggregate
analysis engine in homework 2 as an UIMA-AS service.


Task 2.1 Learning UIMA-AS,

I read about the information that came on each link, the first one describes the concepts about UIMA-AS
the other two makes a description about UIMA-AS and how to apply it for making some experiments with data.

Task 2.2 Creating an UIMA-AS client (25 pts)

    Checking that UIMA-AS was installed on Eclipse graphically by the Available Software Tool:

1.  UIMA AS binary package installation.

Finally, i realise that UIMA AS was installed in another moment, so is ready, UIMA AS version is 2.4.2



To configure UIMA-AS was the next one that is explained step by step in the link: http://svn.apache.org/viewvc/uima/uima-as/tags/uima-as-2.4.0/README?view=co, following the steps that are explained on README about UIMA-as.

1. Set the environment variable MAVEN_OPTS to -Xmx800m -XX:MaxPerSize-256m.

Then build from the directory containing this README by using the command for MAVEN environment.

➢ mvn clean install

Many problems appears. issue: cmd is not recognized as an internal or external command. searching on internet then one option is to reinstall Java, i did that and magically all workjust fine!!



UIMA_HOME variables

In the next path is located the UIMA-AS, C:\Users\Tania\apache-uima-as-2.4.2\bin

Running the file $startBroker, then appears this execution.



Listening for conections at: tcp://Tania-PC:61616

Connector openwire Started

ActiveMQ JMS Message Broker <localhost, ID:Tania-PC-1403-1395983594388-0:1> started

Then the amq folder is created after  $ startBroker file was executed.
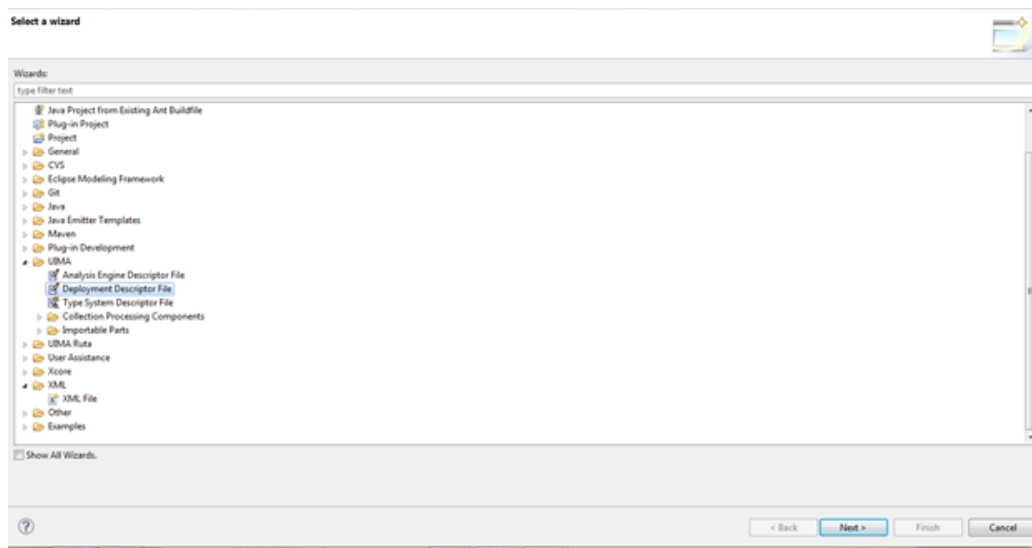
The next step was: run the runRemoteAsyncAE.



Task 2.2 Creating an UIMA-AS client

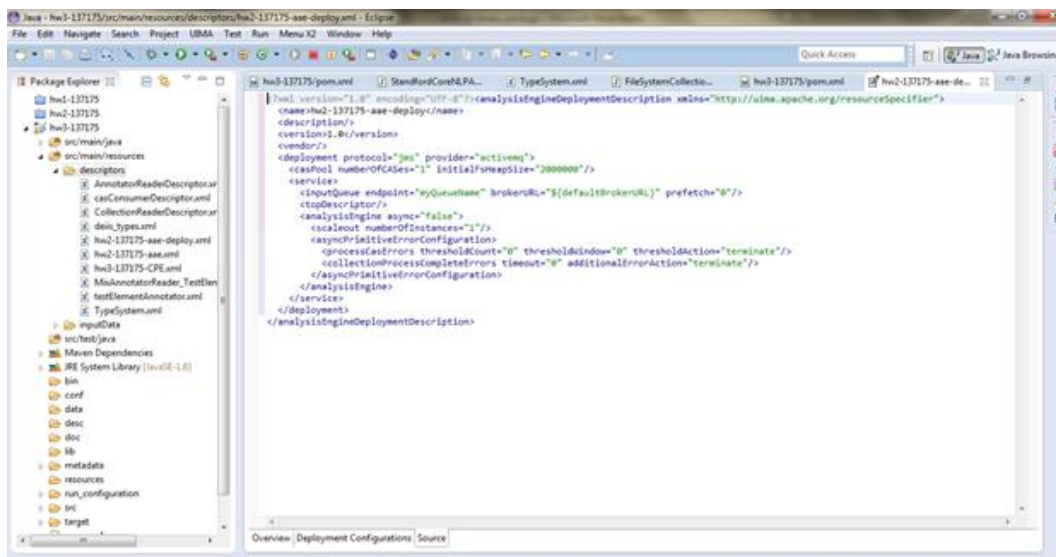Create a UIMA-AS client descriptor (scnlp-137175-client.xml) for a remote UIMA-AS service (Stanford Core NLP), and integrate your client with your CPE pipeline.

UIMA-AS service provided for this homework is the Stanford CoreNLP. Annotator from Clear TK toolkit. This annotator reads the DocumentText form JCas and do tokenization, sentence splitting, POS tagging, lemmatization, NER, syntactic parsing, and coreference resolution.
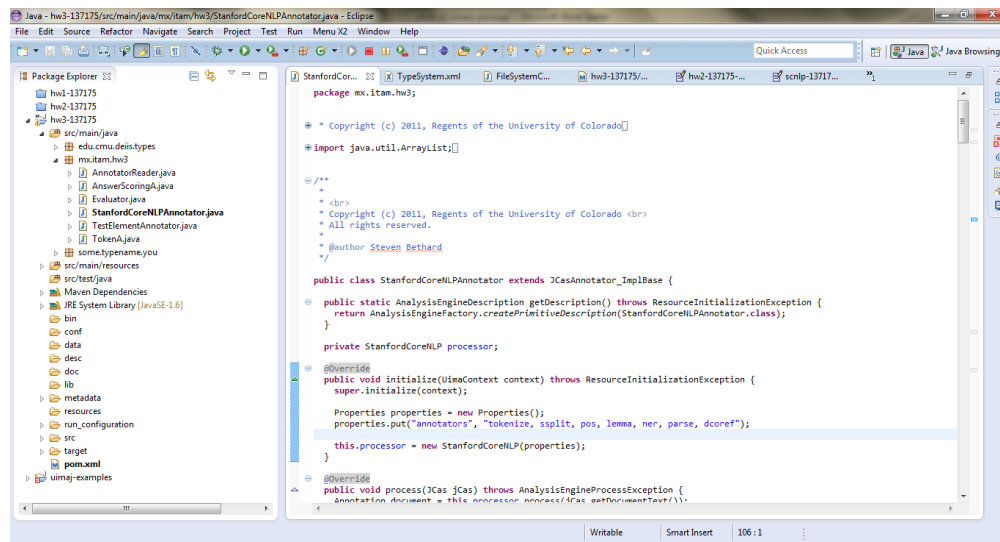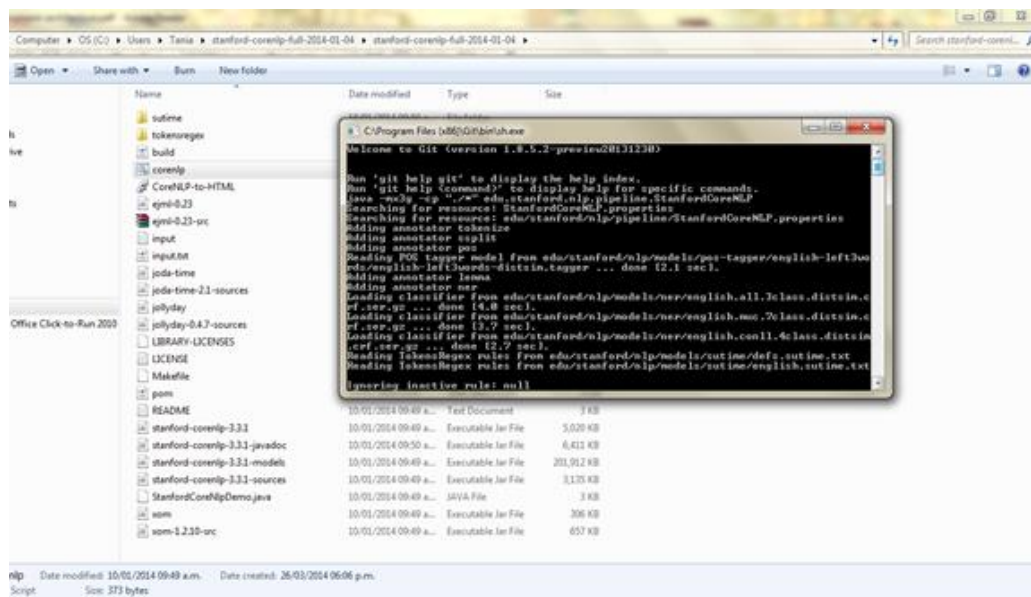
Creating the Deployment Descriptor File



Creating the next files hw2-137175-aae-deploy.xml and scnlp-137175-client.xml
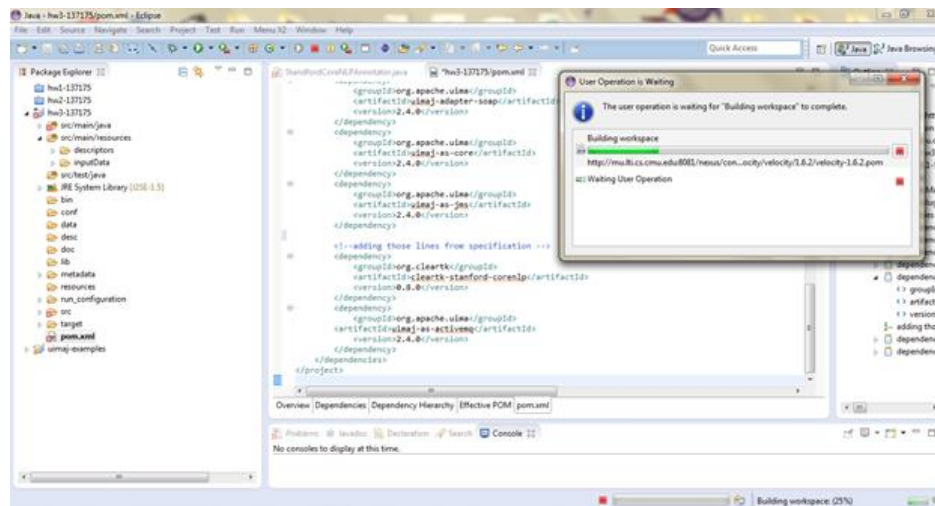
This one is the StanfordCoreNLPAnnotator, (http://nlp.stanford.edu/software/corenlp.shtml)



Testing and Running coreNLP:

Adding those lines into pom.xml file:



## IMPORTING DEPENDENCES INTO POM.XML

```xml
<dependency>
        <groupId>org.apache.uima</groupId>
            <artifactId>uimaj-as-jms</artifactId>
                    <version>2.4.0</version>
                </dependency>

        <!--adding those lines from specification -->
            <dependency>
                <groupId>org.cleartk</groupId>
                    <artifactId>cleartk-stanford-corenlp</artifactId>
                      <version>0.8.0</version>
                </dependency>

          <dependency>
                    <groupId>org.apache.uima</groupId>
                      <artifactId>uimaj-as-activemq</artifactId>
                          <version>2.4.0</version>
                </dependency>
            </dependencies>
</project>
```
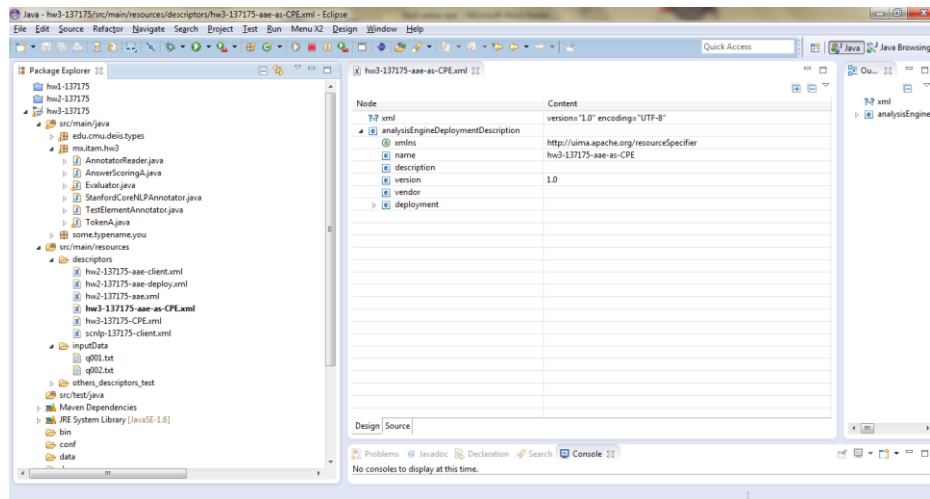
Exploring the tool ClearTK and checking the TypeSystem.xml reference that is here:
https://code.google.com/p/cleartk/source/browse/cleartk-typesystem/src/main/resources/org/cleartk/TypeSystem.xml

Then the file $runRemoteAsyncAE throws some messages. Also the installation about ClearTK is in progress.

Right now the project about the hw3-137175 UIMA and UIMA-AS has the next structure and organization:



**Note:** now I am in the integration the Name Entity annotations from *StanfordCoreNLP.java* into the answer scoring component but it is not running yet.