

Apache UIMA and Academic Research

Grupo de PLN FAMAF

Pablo Ariel Duboue

Les Laboratoires Foulab
999 Rue du College
Montreal, H4C 2S3, Quebec

December 17th 2010



Outline

Ramblings about Frameworks

- Frameworks

- Academic Environments

Intro and Tutorial

- What is UIMA

- Mini-Tutorial

W3C Corpus Processing

- TREC Enterprise Track

Advanced Topics

- Custom Flow Controllers

- UIMA Asynchronous Scale-out

UIMA Scoreboard for Academic Research



Outline

Ramblings about Frameworks

Frameworks

Academic Environments

Intro and Tutorial

What is UIMA

Mini-Tutorial

W3C Corpus Processing

TREC Enterprise Track

Advanced Topics

Custom Flow Controllers

UIMA Asynchronous Scale-out

UIMA Scoreboard for Academic Research



Frameworks

- ▶ Frameworks enable:
 - ▶ Sharing and Collaboration
 - ▶ Growth
 - ▶ Deployment and Large scale implementations
 - ▶ Adoption
- ▶ Frameworks need:
 - ▶ Maintenance (no software is ever “completed”)
 - ▶ Documentation (to further collaboration / adoption)
 - ▶ Neutrality (w.r.t. applications being implemented)
 - ▶ Ownership (on behalf of their developers / maintainers)
 - ▶ Publicity (for widespread adoption)



Enabling Sharing and Collaboration

- ▶ Sharing with the world
 - ▶ Enabling the greater good, without paying a high price (support time, spoiling potential ventures)
- ▶ Sharing with partners
 - ▶ Collaborating with other groups
 - ▶ Commercial impact (with business partners)
- ▶ Sharing with new / potential students
 - ▶ Bringing new people up to speed
 - ▶ Attracting talent



Enabling Growth

- ▶ New phenomena
 - ▶ From syntactic parsing to semantic parsing
 - ▶ From parsing sentences to parsing USB traffic data
- ▶ New artifacts
 - ▶ From text to speech
- ▶ New architectures
 - ▶ From Understanding to Generation

Enable Deployment and Large scale implementations

- ▶ Multiple architectures
 - ▶ Windows, Linux
- ▶ On-line vs. off-line
 - ▶ Batch corpus processing vs. user-oriented Web services
- ▶ New programming languages (and old, efficient ones)
- ▶ New human languages

Enabling adoption

- ▶ Re-using of research components as black-boxes
- ▶ Re-productibility of results
- ▶ Comparison of approaches done by other groups

Outline

Ramblings about Frameworks

Frameworks

Academic Environments

Intro and Tutorial

What is UIMA

Mini-Tutorial

W3C Corpus Processing

TREC Enterprise Track

Advanced Topics

Custom Flow Controllers

UIMA Asynchronous Scale-out

UIMA Scoreboard for Academic Research



Challenges in academic environments

- ▶ High turnover of collaborators
 - ▶ Students
 - ▶ Contractors
 - ▶ Visiting researchers
- ▶ High level of cramming
 - ▶ Paper deadlines
 - ▶ Thesis writing
- ▶ Collaboration at higher levels (sharing full components) rather than finer levels
 - ▶ Easier ownership for paper and thesis writing
- ▶ Focus on *smallish* problems, less efficient programming languages.

Wishlist for a framework for academic research

- ▶ Self-documenting
 - ▶ In the sense of literate programming, e.g., CWEB
 - ▶ Hopefully to the point it won't run without adequate meta-data
- ▶ Allows agile experimentation
 - ▶ Natural bindings for quick-prototyping languages
 - ▶ Ease of exploration of *what-if* scenarios
- ▶ *Post-and-forget* sharing with the world
 - ▶ Communication with users is overhead
 - ▶ Person behind the component might have left the group already
 - ▶ Reproducibility of experiments is key to improve the impact of the group
- ▶ Reasonable transition to products
 - ▶ Reuse prototypes (some code?) with industrial partners
 - ▶ Elicit respect on potential industrial partners



Wishlist for a framework for academic research

- ▶ Simplifying transitions
 - ▶ New students (go read the framework documentation before arriving, helps reduce stress and focus energies)
 - ▶ Departing students (document before you leave, but document what?)
- ▶ A shared effort between research groups
 - ▶ Groups are in the business of doing **research** not frameworks
 - ▶ Clear opportunity to cooperate with other groups
- ▶ Your framework should enable uses of your work beyond your own imagination (i.e., not only a *tagging* framework)
 - ▶ New media
 - ▶ New languages
 - ▶ New problems

Outline

Ramblings about Frameworks

Frameworks

Academic Environments

Intro and Tutorial

What is UIMA

Mini-Tutorial

W3C Corpus Processing

TREC Enterprise Track

Advanced Topics

Custom Flow Controllers

UIMA Asynchronous Scale-out

UIMA Scoreboard for Academic Research



What is UIMA

- ▶ UIMA is a framework, a means to integrate text or other unstructured information analytics.
- ▶ Reference implementations available for Java, C++ and others.
- ▶ An Open Source project under the umbrella of the Apache Foundation.

Analytics Frameworks

- ▶ Find all telephone numbers in running text
 - ▶ $((\backslash([0-9]\{3}\backslash))|[0-9]\{3})-?[0-9]\{3}-?[0-9]\{4}$
- ▶ Nice but...
 - ▶ How are you going to feed this result for further processing?
 - ▶ What about finding non-standard proper names in text?
 - ▶ Acquiring technology from external vendors, free software projects, etc?



In-line Annotations

- ▶ Modify text to include annotations
 - ▶ This/**DET** happy/**ADJ** puppy/**N**
- ▶ It gets very messy very quickly
 - ▶ (S (NP (This/DET happy/ADJ puppy/N) (VP eats/V (NP (the/DET bone/N))))
- ▶ Annotations can easily cross boundaries of other annotations
 - ▶ He said **<confidential>**the project can't go on. The funding is lacking.**</confidential>**



Standoff Annotations

- ▶ Standoff annotations
 - ▶ Do not modify the text
 - ▶ Keep the annotations as offsets within the original text
- ▶ Most analytics frameworks support standoff annotations.
- ▶ UIMA is built with standoff annotations at its core.
- ▶ Example:

He said the project can't go on. The funding is lacking.

012345678901234567890123567890123456789012345678901234567

- ▶ Sentence Annotation: 0-32, 35-57.
- ▶ Confidential Annotation: 8-57.

Type Systems

- ▶ Key to integrating analytic packages developed by independent vendors.
- ▶ Clear metadata about
 - ▶ Expected Inputs
 - ▶ Tokens, sentences, proper names, etc
 - ▶ Produced Outputs
 - ▶ Parse trees, opinions, etc
- ▶ The framework creates an **unified** typesystem for a given set of annotators being run.

Many frameworks

- ▶ Besides UIMA
 - ▶ <http://incubator.apache.org/uima>
- ▶ LingPipe
 - ▶ <http://alias-i.com/lingpipe/>
- ▶ Gate
 - ▶ <http://gate.ac.uk/>
- ▶ NLTK
 - ▶ <http://www.nltk.org/>



UIMA Advantages

- ▶ Apache Licensed
- ▶ Enterprise-ready code quality
- ▶ Demonstrated scalability
- ▶ Developed by experts in building frameworks
 - ▶ Not domain (e.g., NLP) experts
- ▶ Interoperable (C++, Java, others)



Outline

Ramblings about Frameworks

Frameworks

Academic Environments

Intro and Tutorial

What is UIMA

Mini-Tutorial

W3C Corpus Processing

TREC Enterprise Track

Advanced Topics

Custom Flow Controllers

UIMA Asynchronous Scale-out

UIMA Scoreboard for Academic Research



UIMA Concepts

- ▶ Common Annotation Structure or CAS
 - ▶ Subject of Analysis (SofA or View)
 - ▶ JCas
- ▶ Feature Structures
 - ▶ Annotations
- ▶ Indices and Iterators
- ▶ Analysis Engines (AEs)
 - ▶ AEs descriptors



Room annotator

- ▶ From the UIMA tutorial, write an Analysis Engine that identifies room numbers in text.

Yorktown patterns: 20-001, 31-206, 04-123 (Regular Expression Pattern: `[0-9][0-9]-[0-2][0-9][0-9]`)

Hawthorne patterns: GN-K35, 1S-L07, 4N-B21 (Regular Expression Pattern: `[G1-4][NS]-[A-Z][0-9]`)

- ▶ Steps:
 1. Define the CAS types that the annotator will use.
 2. Generate the Java classes for these types.
 3. Write the actual annotator Java code.
 4. Create the Analysis Engine descriptor.
 5. Test the annotator.



Editing a Type System

TutorialTypeSystem.xml

Type System Definition

Types (or Classes)

The following types (classes) are defined in this analysis engine descriptor.
The grayed out items are imported or merged from other descriptors, and cannot be edited here. (To edit them, edit their source files).

Type Name or Feature Name	SuperType or Range	Element Type
org.apache.uima.tutorial.RoomNumber	uima.tcas.Annotation	
building	uima.cas.String	

Imported Type Systems

The following type systems are included as part of this one.

Buttons: Add..., Remove, Set DataPath

Kind	Location/Name

Buttons: Add Type, Add..., Edit..., Remove, Export..., JCasGen

Overview | Type System | Source



The XML descriptor

```
<?xml version="1.0" encoding="UTF-8" ?>
<typeSystemDescription xmlns="http://uima.apache.org/resourceSpecifier">
  <name>TutorialTypeSystem</name>
  <description>Type System Definition for the tutorial examples –
    as of Exercise 1</description>
  <vendor>Apache Software Foundation</vendor>
  <version>1.0</version>
  <types>
    <typeDescription>
      <name>org.apache.uima.tutorial.RoomNumber</name>
      <description></description>
      <supertypeName>uima.tcas.Annotation</supertypeName>
      <features>
        <featureDescription>
          <name>building</name>
          <description>Building containing this room</description>
          <rangeTypeName>uima.cas.String</rangeTypeName>
        </featureDescription>
      </features>
    </typeDescription>
  </types>
</typeSystemDescription>
```



The AE code

```

package org.apache.uima.tutorial.ex1;

import java.util.regex.Matcher;
import java.util.regex.Pattern;

import org.apache.uima.analysis_component.JCasAnnotator_ImplBase;
import org.apache.uima.jcas.JCas;
import org.apache.uima.tutorial.RoomNumber;

/**
 * Example annotator that detects room numbers using
 * Java 1.4 regular expressions.
 */
public class RoomNumberAnnotator extends JCasAnnotator_ImplBase {
    private Pattern mYorktownPattern =
        Pattern.compile("\\b[0-4]\\d-[0-2]\\d\\d\\d\\b");

    private Pattern mHawthornePattern =
        Pattern.compile("\\b[G1-4][NS]-[A-Z]\\d\\d\\d\\b");

    public void process(JCas aJCas) {
        // next slide
    }
}

```



The AE code (cont.)

```

public void process(JCas aJCas) {
    // get document text
    String docText = aJCas.getDocumentText();
    // search for Yorktown room numbers
    Matcher matcher = mYorktownPattern.matcher(docText);
    int pos = 0;
    while (matcher.find(pos)) {
        // found one – create annotation
        RoomNumber annotation = new RoomNumber(aJCas);
        annotation.setBegin(matcher.start());
        annotation.setEnd(matcher.end());
        annotation.setBuilding("Yorktown");
        annotation.addToIndexes();
        pos = matcher.end();
    }
    // search for Hawthorne room numbers
    // ..
}

```

UIMA Document Analyzer





UIMA Document Analyzer (cont)

The screenshot shows the UIMA Document Analyzer interface. The main window displays a list of events:

- UIMA Summer School
- August 26, 2003
UIMA 101 - The New UIMA Introduction
(Hands-on Tutorial)
9:00AM-5:00PM in HAW GN-K35
- August 28, 2003
FROST Tutorial
9:00AM-5:00PM in HAW GN-K35
- September 15, 2003
UIMA 201: UIMA Advanced Topics
(Hands-on Tutorial)
9:00AM-5:00PM in HAW 1S-F53
- September 17, 2003
The UIMA System Integration Test and Hardening Service
The "SITH"
3:00PM-4:30PM in HAW GN-K35

Below the list is a Legend section with two entries:

- DocumentA...
- RoomNumber

At the bottom of the window, there are buttons for "Select All" and "Deselect All", and a "Viewer Mode:" section with radio buttons for "Annotations" (selected) and "Entities".

Outline

Ramblings about Frameworks

Frameworks

Academic Environments

Intro and Tutorial

What is UIMA

Mini-Tutorial

W3C Corpus Processing

TREC Enterprise Track

Advanced Topics

Custom Flow Controllers

UIMA Asynchronous Scale-out

UIMA Scoreboard for Academic Research



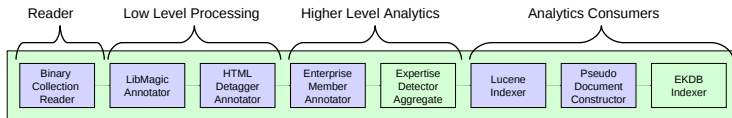
An Example

- ▶ TREC 2006 enterprise track
- ▶ Search for experts in W3C Website
 - ▶ Given topic, find expert in topic
- ▶ The IBM Enterprise Track 2006 Team
 - ▶ Guillermo Averboch, **Jennifer Chu-Carroll**, Pablo A Duboue, David Gondek, J William Murdock and John Prager .



Corpus Processing: Generalities

A simplified view of our TREC 2006 Enterprise Track Indexing Pipeline (actual pipeline has 23 components)



Pipeline

- ▶ Reader
 - ▶ Binary Collection Reader
- ▶ Low Level Processing
 - ▶ LibMagic Annotator
 - ▶ HTML Detagger Annotator
- ▶ Higher Level Analytics
 - ▶ Enterprise Member Annotator
 - ▶ Expertise Detector Aggregate
- ▶ Analytics Consumers
 - ▶ Lucene Indexer
 - ▶ Pseudo Document Constructor
 - ▶ EKDB Indexer



Binary Collection Reader

- ▶ Reads the TREC XML format
- ▶ 300,000+ documents (a full crawl of the w3.org site)
- ▶ Binary format, to allow auto-detection of file type, encoding, etc.

LibMagic Annotator

- ▶ Uses “magic” numbers to heuristically guess the file type.
- ▶ JNI wrapper to libmagic in Linux.
- ▶ Non-supported file types are dropped.
- ▶ UIMA can run this remotely from a Windows machine.



HTML Detagger Annotator

- ▶ For documents identified as HTML, parse them and extract the text.
- ▶ Perform also encoding detection (utf-8 vs. iso-8859-1).
- ▶ Other detaggers (not shown) are applied to other file formats.

Enterprise Member Annotator

- ▶ Detects inside running text the occurrence of any variant of the 1,000+ experts for the Enterprise Track.
- ▶ Dictionary extended with name variants.
- ▶ Simple TRIE-based implementation.



Expertise Detector Aggregate

- ▶ Hierarchical aggregate of 16 annotators leveraging existing technology into a new “expertise detection” annotator.
- ▶ Includes a named-entity detector and a relation detector for semantic patterns.

Lucene Indexer

- ▶ Integration with Open Source technology.
- ▶ Indexes the tokens from the text.
- ▶ The UIMA framework also contains JuruXML, an indexer for semantic information.

Pseudo Document Constructor

- ▶ Uses the name occurrences to create a “pseudo” document with all text surrounding each expert name.
- ▶ The pseudo documents are indexed off-line.

EKDB Indexer

- ▶ Stores extracted entities and relations in a relational database.
- ▶ Standards-based (JDBC, RDF).
- ▶ Employed in a variety of research applications for search and inference.

Outline

Ramblings about Frameworks

Frameworks

Academic Environments

Intro and Tutorial

What is UIMA

Mini-Tutorial

W3C Corpus Processing

TREC Enterprise Track

Advanced Topics

Custom Flow Controllers

UIMA Asynchronous Scale-out

UIMA Scoreboard for Academic Research



Custom Flow Controllers

- ▶ UIMA allows you to specify which AE will receive the CAS next, based on all the annotations on the CAS.
- ▶ `examples/descriptors/flow_controller/WhiteboardFlowController.xml`
 - ▶ FlowController implementing a simple version of the “whiteboard” flow model. Each time a CAS is received, it looks at the pool of available AEs that have not yet run on that CAS, and picks one whose input requirements are satisfied. Limitations: only looks at types, not features. Does not handle multiple Sofas or CasMultipliers.

Outline

Ramblings about Frameworks

Frameworks

Academic Environments

Intro and Tutorial

What is UIMA

Mini-Tutorial

W3C Corpus Processing

TREC Enterprise Track

Advanced Topics

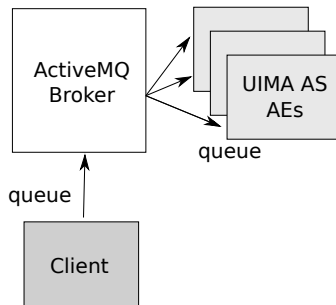
Custom Flow Controllers

UIMA Asynchronous Scale-out

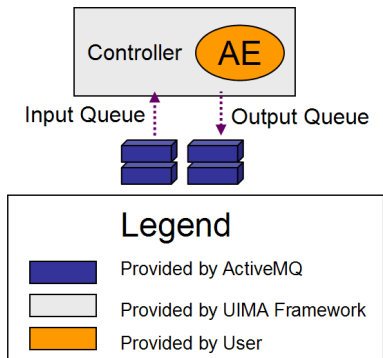
UIMA Scoreboard for Academic Research



UIMA AS: ActiveMQ



UIMA AS: Wrapping Primitive AEs



UIMA AS: More information

- ▶ <http://incubator.apache.org/uima/doc-uimaas-what.html>
- ▶ <http://svn.apache.org/viewvc/incubator/uima/uima\discretionary{-}{ }as/trunk/uima-as-distr/src/main/readme/README?view=markup>
- ▶ http://incubator.apache.org/uima/downloads/releaseDocs/2.3.0\discretionary{-}{ }incubating/docs-uima-as/html/uima_async_scaleout/uima_async_scaleout.html

How UIMA fares w.r.t. Academic Research

Self-documenting	Doing well
Agile experimentation	So, so
Post-and-forget	Yes
Transition to products	Yes
Students transitions	Could be
Shared efforts	Industrial at the moment
Beyond your own work	Yes



