# Text Annotation with OpenNLP and UIMA

Graham Wilcock University of Helsinki graham.wilcock@helsinki.fi

### Abstract

The tutorial presents a practical overview of automatic linguistic annotation of texts using freely available open source tools.

## 1 OpenNLP

Text annotation typically involves tasks at several linguistic levels, such as sentence boundary detection, tokenization, part-of-speech tagging, phrase chunking, syntactic parsing, named entity recognition, coreference resolution, and semantic role labelling. Most of these tasks can be done with appropriate combinations of OpenNLP tools (http://opennlp.sourceforge.net).

Practical examples will show annotations of a short English text. OpenNLP outputs annotations in a simple plain text format.

The OpenNLP tools do a good job of creating annotations automatically, but a number of issues arise. Although the OpenNLP tools themselves are open source Java and platform-independent, the annotation pipelines (where the output of one component is input to the next component) are created by Linux shell scripts and Windows .bat files that are platform-dependent and errorprone. Apache Ant can be used to gain platformindependence, but Ant requires technical skills.

#### 2 WordFreak

OpenNLP tools can also be used in WordFreak (http://wordfreak.sourceforge.net) as plugins. WordFreak provides an attractive, easy-to-use GUI for linguistic annotations. It is open source Java and platform-independent, and is convenient for manually correcting annotations made by the OpenNLP tools. However, Word-Freak creates annotations in its own specific XML stand-off annotation format.

This raises the issue of interoperability. How can annotations be interchanged between tools that use different annotation formats? This can be done by XSLT transformations, for example WordFreak XML format can be transformed by XSLT to OpenNLP plain text annotation format. However, writing such XSLT stylesheets requires specific technical skills.

# 3 UIMA

UIMA (Unstructured Information Management Architecture) provides solutions to many of the above issues. UIMA is open-source Java (http: //incubator.apache.org/uima). It aims to support interoperability and scalability.

In UIMA, annotators run in analysis engines. New annotators are written in Java, and existing annotation tools such as the OpenNLP tools are converted to UIMA annotators by Java wrappers. Pipelines of annotators run in aggregate analysis engines. Pipelines can be configured by writing XML descriptors (similar in some ways to Ant tasks), or by means of an easy-to-use graphical configuration tool in the Eclipse GUI (Figure 1).

UIMA supports interoperability at the level of annotation formats by adopting XML Metadata Interchange (XMI), which has been proposed as an interchange standard. Instead of having its own specific XML annotation format, the UIMA annotation format is XMI.

UIMA also supports interoperability at the level of annotation tools by means of a type system that defines annotation types and their features. Types are used to check that output from one component is the right type for input to the next component.

Practical examples will show how to configure and use pipelines of OpenNLP tools in UIMA, and how to view the annotations in UIMA (Figure 2).

# References

Graham Wilcock. 2009. Introduction to Linguistic Annotation and Text Analytics. Morgan and Claypool. Graham Wilcock

🗲 Java - uimaj-examples/opennlp_wrappers	:/descriptors/OpenNLPAggregate.xml - Eclipse Platform				
Eile Edit Navigate Search Project Run UIM	IA <u>W</u> indow <u>H</u> elp				
📬 • 🖫 👜   🏇 • 💽 • 💁   🍰	₩ @ • ] @ @ \$ ] 1 • 1 • 1 • 1 • 1 • •	🔛 🚭 Java			
📙 Package Explor 🕴 🍃 Hierarchy 🗖 🗖	OpenNLPAggregate.xml 🛛 🗳 OpenNLPParser.xml	- 0			
	OpenNLPAggregate,xml				
□ 🔐 uimaj-examples	Aggregate Delegates and Flows				
get opennip_wrappers/src	▼ Component Engines	▼ Component Engine Flow			
Berger Library (fre1.6.0_05] Berger Libraries Construct Const	The following engines are included in this aggregate.     Delegate   Key Name     OpenNLPSentenceDetector xml   OpenNLPSentenceDetector     OpenNLPOSTagger xml   OpenNLPOSTagger     OpenNLPTokenizer xml   OpenNLPTokenizer     Ø   OpenNLPPostager xml     Ø   OpenNLPTokenizer xml     Ø   OpenNLPPostager     Ø   OpenNLPPostager xml     Ø   OpenNLPPostager     Ø   OpenNLPPostager xml     Ø   OpenNLPPostager     Ø   Addmemote     Find AE   Add	Choose a flow type and describe the execution order of your engines. The totale shows the delegates using their key names. Flow Kind: Exed Flow F Flow Controller: From Search Key Name: Search GopenNLPSentenceDetector Up CopenNLPROSTagger Down CopenNLPPortser			
Overview Aggregate Parameters Parameters Settings Type System Capabilities Indexes Resources Source            Problems @ Javadoc @ Declaration @ Console 23           = 第 🦓 : 1 + 5           = 1					
	21				
	<u> </u>				
] []					

Figure 1: Configuring an OpenNLP annotation pipeline in UIMA

Annotation Resu Sonnet 130 by William Shakespear My mistress' eyes are it Coral is far more red th If snow be white, why If hairs be wires, black I have seen roses dan But no such roses see And in some perfumes Than in the breath the I love to hear her spec That music hath a far i I grant I never saw a My mistress when she And yet, by Heaven, 1 As any she belied with	Its for sonnet130.4 e (1564-1616) nothing like the sun, han her lips red. then her breasts are (wires grow on her he asked, red and white I in her cheeks. is there more delight if from ny mistress re ak, yet well I know more pleasing sound. goddess go, walks, treads on the (think my love as rare false compare.	txt.xmi in C:\Docu dun, aadj eks. ground.	ments and Settings\gw\	,My Documents∖Pr	ojects\UJMA\data\processed   ×     Click In Text to See Annotation Detail   Annotations     Annotations   •     • begin = 493   • end = 501     • componentId = OpenNLP Parser   •     • begin = 493   • end = 510     • componentId = OpenNLP Parser   •     • NP   •     • Dypl"   •     • end = 510   •     • componentId = OpenNLP Parser   •     • NP   •     • begin = 491   •     • end = 516   •     • componentId = OpenNLP Parser   •     • SBAR   •     • SBAR   •     • YP   •     • WP   •     • WP   •     • SBAR   •     • YP   •     • YP   •     • YP   •     • SUP   •     • SBAR   •     • YP   •     • SUP   •
			Dogument Appo		
SourceDocume	Token	VP		1_ Jondence	

Figure 2: Viewing annotations by OpenNLP Parser in UIMA