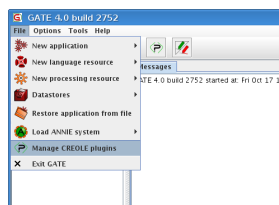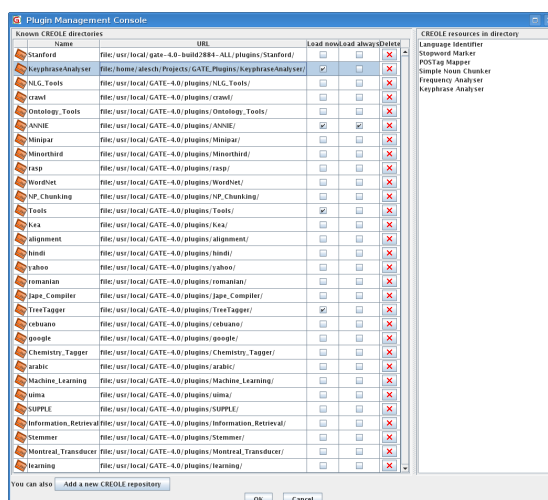# Usage as GATE plugin

Alexander Schutz

November 11, 2008

The keyphrase extraction has been developed as a number of GATE plugins and it makes use of some of the nice features the GATE framework has to offer. Therefore, it is possible to run the keyphrase extraction as a pipeline inside the GATE gui, for testing or evaluation purposes, or for prototyping. The following provides a complete walk through of what has to be done in order to run the keyphrase extraction as GATE plugin, and how results can be inspected.

1. Download the GATE plugins at `http://resources.smile.deri.ie/nlp/ keyphrase-extraction/current/keyphrase-extraction-1.0.0-gate. zip`.

2. Unzip in a folder of your choice, in the remainder denoted as keyphrase-extraction-folder.

3. Start GATE.

4. When the user interface comes up, select *File > Manage CREOLE plugins*.
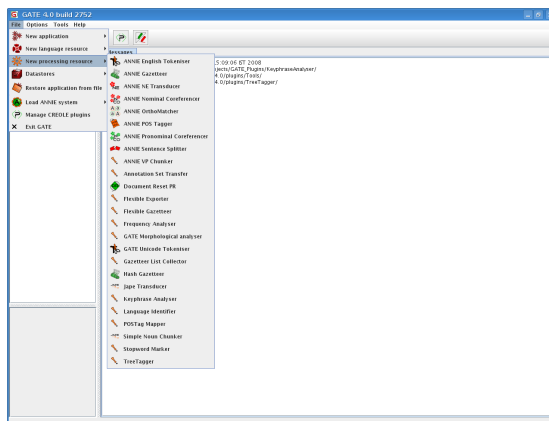


5. In the following dialogue, click on *Add a new CREOLE repository*, then select the directory *keyphrase-extraction-folder* (this is where you unzipped the plugins), confirm with OK in the file chooser. Now, an additional entry should appear in the Plugin Management Console, named *KeyphraseAnalyser*. By clicking on it, a number of CREOLE resources should be listed on the right hand side.

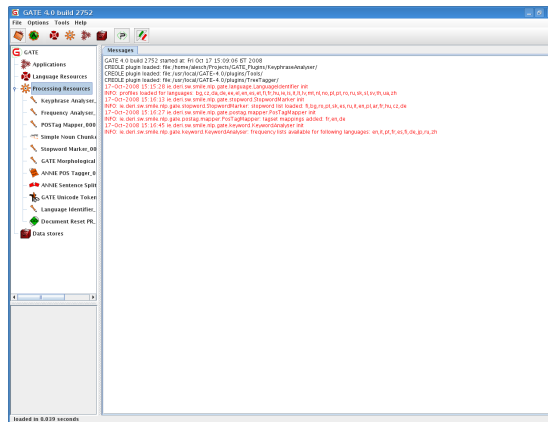6. Click on *Load Now* (or *Load Always*).



1

Now is also a good time to enable/load a number of different GATE plugins (namely *ANNIE* and *Tools*, and for German and French texts *TreeTagger*), as they are needed for the Keyphrase Extraction pipeline. Mark the mentioned plugins as *Load now* or *Load always* and return to the GATE main interface by clicking *OK* in the Plugin Management Console. The next steps show how to construct the keyphrase extraction processing pipeline

7. Load the required Processing Resources. To do so, right click on the *Processing Resources* field in the panel on left hand side. A menu item *New* should appear, which, when hovered over by the mouse, should trigger the display of enabled Processing Resources. Alternatively, you can click on the Menu *File > New Processing Resource.*
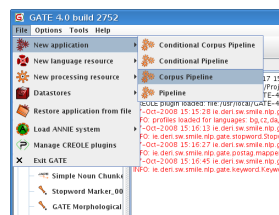


8. Choose the following Processing Resources (you can blindly accept/confirm the upcoming dialogue, there is no need to adjust the init-time parameters):

   - Document Reset PR
   - Language Identifier
   - GATE Unicode Tokeniser
   - ANNIE Sentence Splitter
   - ANNIE POS Tagger
   - GATE Morphological Analyser
   - Stopword Marker
   - Simple Noun Chunker
   - POSTag Mapper
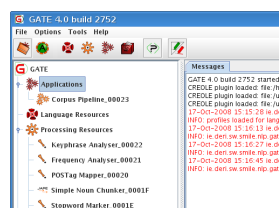   - Frequency Analyser
   - Keyphrase Analyser

   While loading these Processing Resources, you should be able to observe some informative output in the Messages tab.
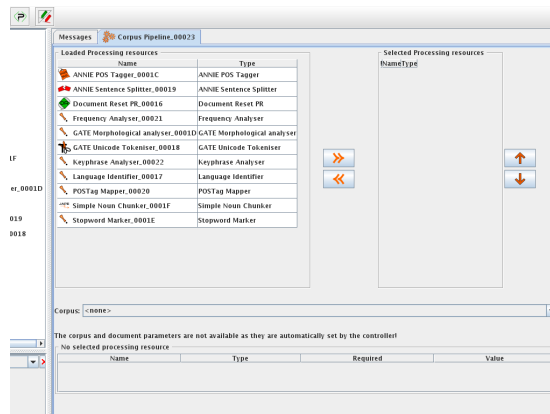
9. Now, right click on *Applications* in the panel on left hand side, and select *New > Corpus Pipeline.* Alternatively, you can also select from the *Menu File > New Application > Corpus Pipeline.*
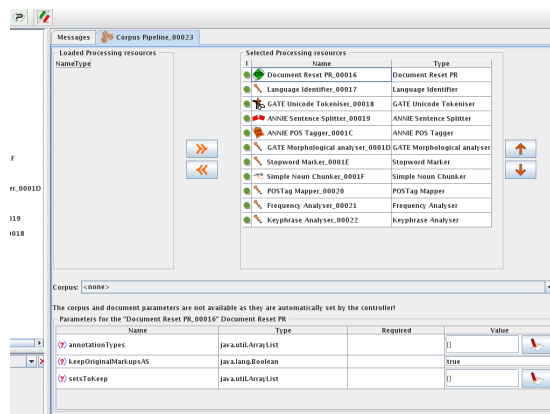


You may name the pipeline and confirm, after that an instance of the Corpus Pipeline with the assigned name should appear in the panel on the left hand side.
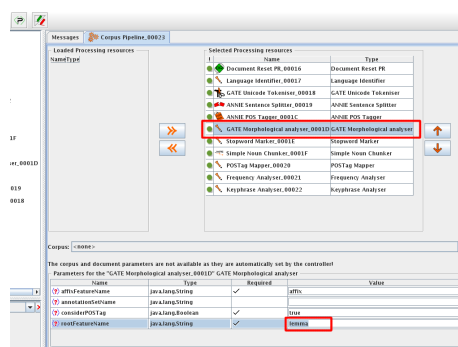


10. Double click the newly created instance of Corpus Pipeline, and a tab opens in the central area of the screen, which provides a view of the *Loaded Processing Resources* (left hand side) and the *Selected Processing Resources* (right hand side). Currently, no Processing Resource is selected.
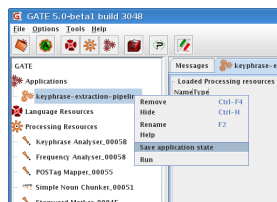
Mark all of the previously loaded Processing Resources, and use the arrow button pointing to the right hand side to assign the selected Processing Resources to the Processing Pipeline. Now, with the up and down arrows, assemble the Processing Resources to the order given at point 8.
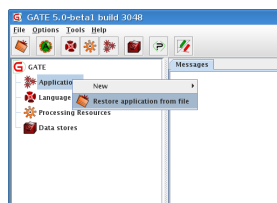


11. Still in the Corpus Pipeline view, click on the *GATE Morphological Analyser*. In the lower part of the view, you can see the default runtime parameters. Click on the value *rootFeatureName*, which is set to "root" by default. You need to change this value to "lemma".
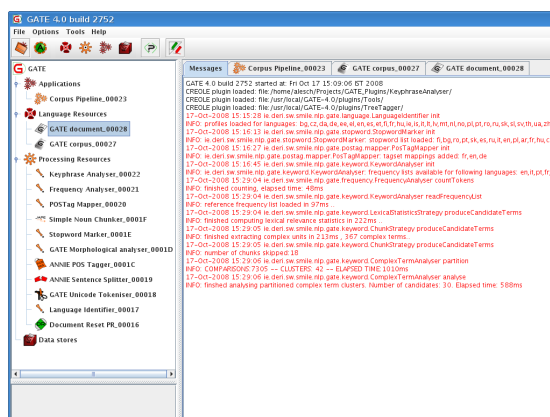
12. For convenience reasons, you can now save the application state to a file, in order to easily load the pipeline at a later point in time. This way you can avoid having to load every processing resource and the configuration process. To do so, simply right-click on the corpus pipeline created in the previous step, select *Save application state* and enter a file name, followed by the ".gapp" (or ".xgapp") extension.



To load a previously saved application pipeline, select *File > Restore application from state*, and choose the desired *.gapp (or *.xgapp) file.



13. Create a corpus, populate it with some documents, and set it as the corpus to be processed by the Processing Pipeline.

14. When you run the previously constructed Processing Pipeline over your corpus, each document will be annotated with a number of additional features, which you can observe in the lower, left-hand corner after double-clicking on a document.

The additional feature names are as follows, and they can be accessed programmatically via API as well:

**keywords:** a sorted map of <keyphrase:confidence > pairs, where confidence is an double value between 0 and 1. High confidence keyphrases have a value closer to 1.

**language:** the language of the document

**document_size:** the number of words/tokens in the document

**lexicon_size:** the number of different lemmas in the document

**overall-token_frequency:** the frequencies of the tokens in the document, implemented as a sorted map

**overall-token_frequency:** the frequencies of the tokens in the document, implemented as a sorted map

**noun-lemma_frequency:** the lemma frequencies of the nouns in the document, implemented as a sorted map

**verb-lemma_frequency:** the lemma frequencies of the verbs in the document, implemented as a sorted map

**adj-lemma_frequency:** the lemma frequencies of the adjectives in the document, implemented as a sorted map



15. The results of the keyphrase extraction can be inspected by looking at the keywords feature of the document.