

Given a news article, after installing opennlp, the following steps should be completed using opennlp:

1. *Detect sentences in the given news article (2 points)*
2. *Tokenize each sentence into words (2 points)*
3. *Perform Part-of-Speech (POS) on each sentence (3 points)*
4. *Find name entities including person's name entities and locations (3 points)*

How this report is organized:

- The first section, "Results" describe the results from the experiments. There is one sub-section for each requested task. Results are abbreviated to make the report more legible. The full results are in the accompanying file, assignment1-full-results-cgarbin.txt.
- The appendix has the code and screenshots.
- The source code is available as accompanying files. The appendix section describes the organization of the code.

Results

Part 1 - Sentence Detection

1. Detect sentences in the given news article (2 points)

The code detected 16 sentences, listed below, in the order they appear in the document. Sentences are abbreviated to make the results more legible. Full results are available in the attached file assignment1-full-results-cgarbin.txt.

- [00] Anglo-French Channel Tunnel operator Eurotunnel Monday ... one billion pounds (\$1.56 billion) of its debt.
- [01] The long-awaited restructuring brings to an end months ... nearly nine billion pounds (\$14.1 billion).
- [02] The deal, announced simultaneously in Paris and London,... only 54.5 percent of the company.
- [03] "The restructuring plan provides Eurotunnel with the ... " Eurotunnel co-chairman Alastair Morton said.
- [04] The firm was now making a profit before interest, he added.
- [05] Although shareholders will see their interests diluted,... negotiated during the tunnel's construction.
- [06] Eurotunnel, which has taken around half the cross-Channel ... dividend within the next 10 years.
- [07] French co-chairman Patrick Ponsolle said shareholders ... the benefits of the company's success.
- [08] He called the debt restructuring plan "an acceptable compromise" for holders of Eurotunnel shares.
- [09] The company said there was still considerable work to be ... for approval, probably early in 1997.
- [10] Monday's announcement followed two weeks of highly secretive ... Eurotunnel and its six leading banks.
- [11] This was extended to the 24 "instructing banks" at a meeting late last week in London.
- [12] Eurotunnel said the debt-for-equity swap would be at 130 pence, or 10.40 francs, per share.
- [13] That is considerably below the level of around 160 pence ... billion) by 1.0 billion (\$1.56 billion).
- [14] The company said a further 3.7 billion pounds (\$5.8 billion)... be able to participate in this issue.
- [15] If they choose not to take up free warrants entitling them... the company by the end of December 2003.
- [16] Eurotunnel's shares, which were suspended last week at 113.5... trading on Tuesday, the company said.

Parte 2 - Tokenization

2. Tokenize each sentence into words (2 points)

OpenNLP has three tokenizers:

- SimpleTokenizer: based on character classes - *"sequences of the same character class are tokens"*.
- WhitespaceTokenizer: *"non whitespace sequences are identified as tokens"*
- TokenizerME: a probabilistic tokenizer - *"detects token boundaries based on probability model"*

The table below shows the difference between the tokenizers. The table has the results for each tokenizer on the first sentence of the document. Full results are available in the attached file `assignment1-full-results-cgarbin.txt`.

It is clear that TokenizerME does a better job of identifying tokens as humans would do. For example, it performs better when extracting numbers from the surrounding text, as shown in the highlighted tokens in the table. For this reason, TokenizerME was used in the subsequent tasks (part-of-speech and named entities).

TokenizerME	WhitespaceTokenizer	SimpleTokenizer
[00] Anglo-French	[00] Anglo-French	[00] Anglo
[01] Channel	[01] Channel	[01] -
[02] Tunnel	[02] Tunnel	[02] French
[03] operator	[03] operator	[03] Channel
[04] Eurotunnel	[04] Eurotunnel	[04] Tunnel
[05] Monday	[05] Monday	[05] operator
[06] announced	[06] announced	[06] Eurotunnel
[07] a	[07] a	[07] Monday
[08] deal	[08] deal	[08] announced
[09] giving	[09] giving	[09] a
[10] its	[10] its	[10] deal
[11] creditor	[11] creditor	[11] giving
[12] banks	[12] banks	[12] its
[13] 45.5	[13] 45.5	[13] creditor
[14] percent	[14] percent	[14] banks
[15] of	[15] of	[15] 45
[16] the	[16] the	[16] .
[17] company	[17] company	[17] 5
[18] in	[18] in	[18] percent
[19] return	[19] return	[19] of
[20] for	[20] for	[20] the
[21] wiping	[21] wiping	[21] company
[22] out	[22] out	[22] in
[23] one	[23] one	[23] return
[24] billion	[24] billion	[24] for
[25] pounds	[25] pounds	[25] wiping
[26] ([26] (\$1.56	[26] out
[27] \$	[27] billion)	[27] one
[28] 1.56	[28] of	[28] billion
[29] billion	[29] its	[29] pounds
[30])	[30] debt.	[30] (
[31] of		[31] \$
[32] its		[32] 1
[33] debt		[33] .
[34] .		[34] 56
		[35] billion
		[36])
		[37] of
		[38] its
		[39] debt
		[40] .

Part 3 - Part-of-Speech (POS)

Part-of-speech in OpenNLP requires tokenization first. Based on the results of the previous section, TokenizerME was used as the tokenizer.

The results for the first three sentences in the document are shown in the table below. Full results are available in the attached file `assignment1-full-results-cgarbin.txt`.

Sentence 1			Sentence 2			Sentence 3		
[00]	Anglo-French	JJ	[00]	The	DT	[00]	The	DT
[01]	Channel	NNP	[01]	long-awaited	JJ	[01]	deal	NN
[02]	Tunnel	NNP	[02]	restructuring	NN	[02]	,	,
[03]	operator	NN	[03]	brings	VBZ	[03]	announced	VBD
[04]	Eurotunnel	NNP	[04]	to	TO	[04]	simultaneously	RB
[05]	Monday	NNP	[05]	an	DT	[05]	in	IN
[06]	announced	VBD	[06]	end	NN	[06]	Paris	NNP
[07]	a	DT	[07]	months	NNS	[07]	and	CC
[08]	deal	NN	[08]	of	IN	[08]	London	NNP
[09]	giving	VBG	[09]	wrangling	VBG	[09]	,	,
[10]	its	PRP\$	[10]	between	IN	[10]	brings	VBZ
[11]	creditor	NN	[11]	Eurotunnel	NNP	[11]	the	DT
[12]	banks	NNS	[12]	and	CC	[12]	company	NN
[13]	45.5	CD	[13]	the	DT	[13]	back	RB
[14]	percent	NN	[14]	225	CD	[14]	from	IN
[15]	of	IN	[15]	banks	NNS	[15]	the	DT
[16]	the	DT	[16]	to	TO	[16]	brink	NN
[17]	company	NN	[17]	which	WDT	[17]	of	IN
[18]	in	IN	[18]	it	PRP	[18]	insolvency	NN
[19]	return	NN	[19]	owes	VBZ	[19]	but	CC
[20]	for	IN	[20]	nearly	RB	[20]	leaves	VBZ
[21]	wiping	VBG	[21]	nine	CD	[21]	shareholders	NNS
[22]	out	RP	[22]	billion	CD	[22]	owning	VBG
[23]	one	CD	[23]	pounds	NNS	[23]	only	RB
[24]	billion	CD	[24]	(-LRB-	[24]	54.5	CD
[25]	pounds	NNS	[25]	\$	\$	[25]	percent	NN
[26]	(-LRB-	[26]	14.1	CD	[26]	of	IN
[27]	\$	\$	[27]	billion	CD	[27]	the	DT
[28]	1.56	CD	[28])	-RRB-	[28]	company	NN
[29]	billion	CD	[29]	.	.	[29]	.	.
[30])	-RRB-						
[31]	of	IN						
[32]	its	PRP\$						
[33]	debt	NN						
[34]	.	.						

Part 4 - Named Entities

Named entities in OpenNLP require tokenization first. Based on the results of the previous section, TokenizerME was used as the tokenizer.

The results for the first four sentences in the document are shown below, with the identified spans highlighted. Full results are available in the attached file assignment1-full-results-cgarbin.txt.

Anglo-French Channel Tunnel operator Eurotunnel Monday announced a deal giving its creditor banks **<PERCENTAGE> 45.5 percent </PERCENTAGE>** of the company in return for wiping out one billion pounds (**<MONEY> \$ 1.56 billion </MONEY>**) of its debt .

The long-awaited restructuring brings to an end months of wrangling between Eurotunnel and the 225 banks to which it owes nearly nine billion pounds (**<MONEY> \$ 14.1 billion </MONEY>**) .

The deal , announced simultaneously in **<LOCATION> Paris </LOCATION>** and **<LOCATION> London </LOCATION>** , brings the company back from the brink of insolvency but leaves shareholders owning only **<PERCENTAGE> 54.5 percent </PERCENTAGE>** of the company .

" The restructuring plan provides Eurotunnel with the medium-term financial stability to allow it to consolidate its substantial commercial achievements to date and to develop its operations , " Eurotunnel co-chairman **<PERSON> Alastair Morton </PERSON>** said .

References

The following references were consulted to write the code and to produce this report:

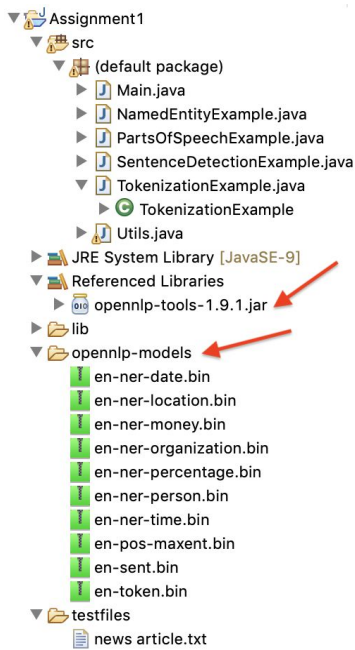
- <https://opennlp.apache.org/docs/1.5.3/manual/opennlp.html>
- <https://opennlp.apache.org/models.html>
- <https://www.tutorialkart.com/opennlp/how-to-setup-opennlp-java-project/>
- <http://www.sfs.uni-tuebingen.de/~keberle/NLPTools/presentations/OpenNLP/OpenNLP%20BorenshteinDaeubler.pdf>

Appendix

Source code overview

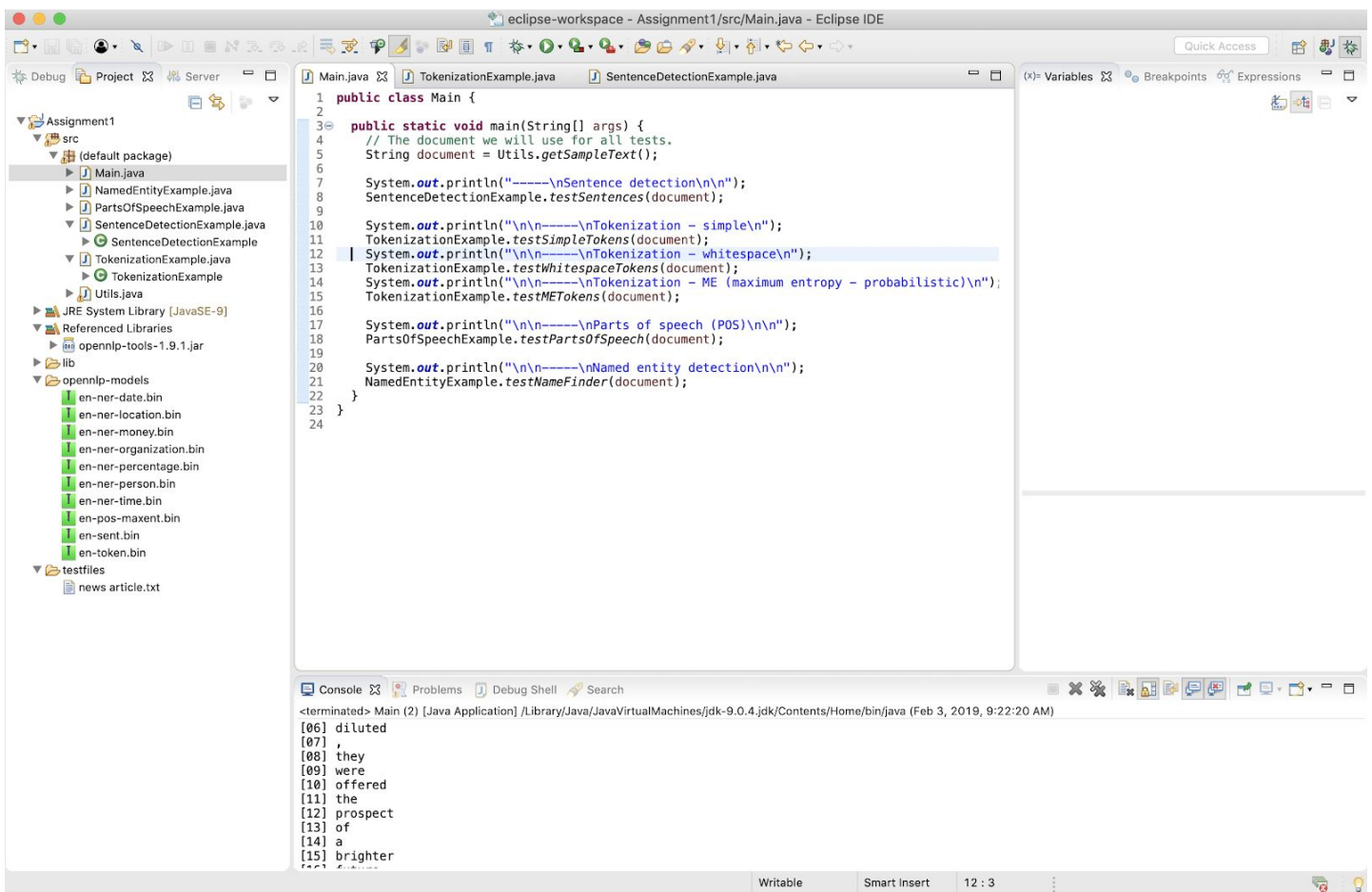
The code is an eclipse project, structured as in the picture below. Each task is in a separate file named accordingly, e.g. NamedEntityExamples.java has the code for the named entity code.

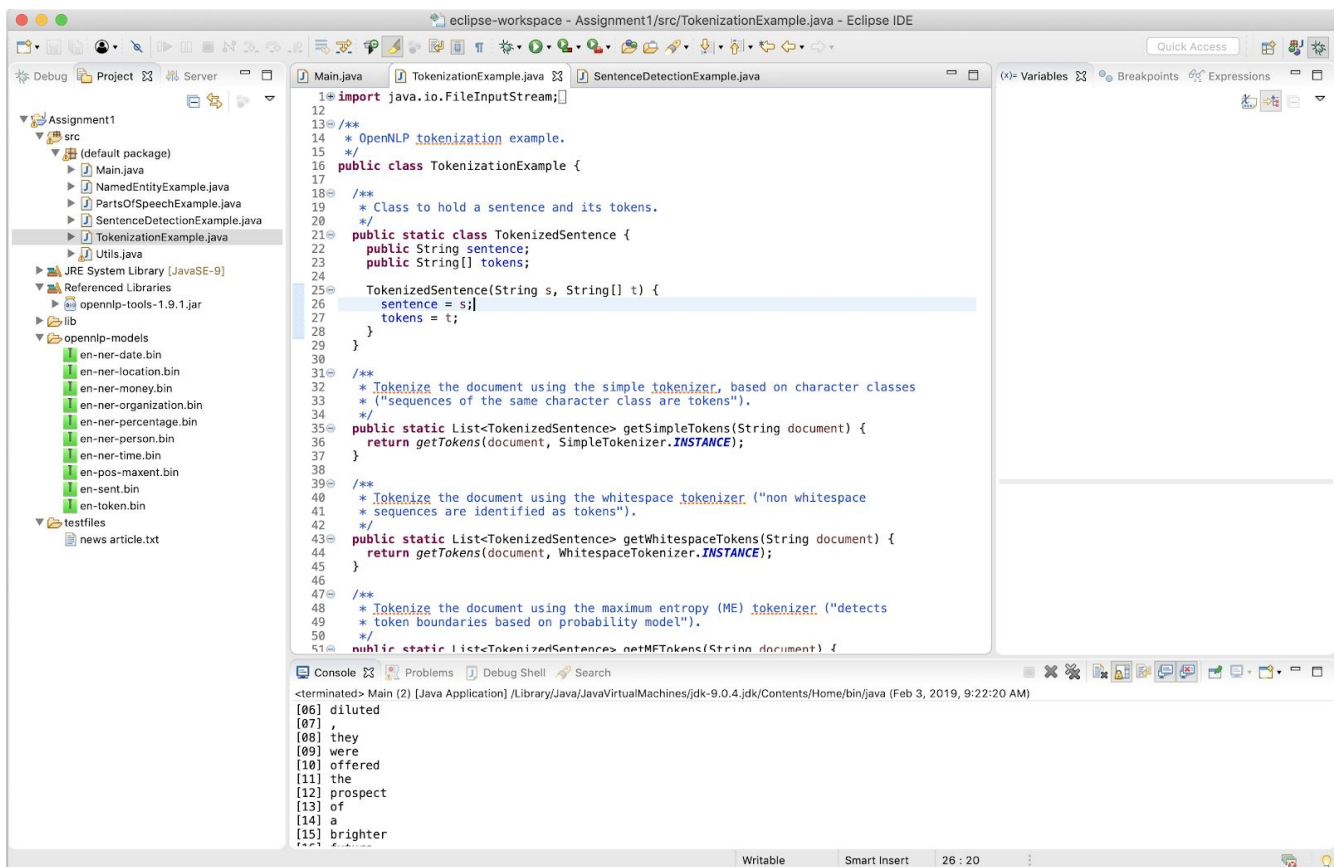
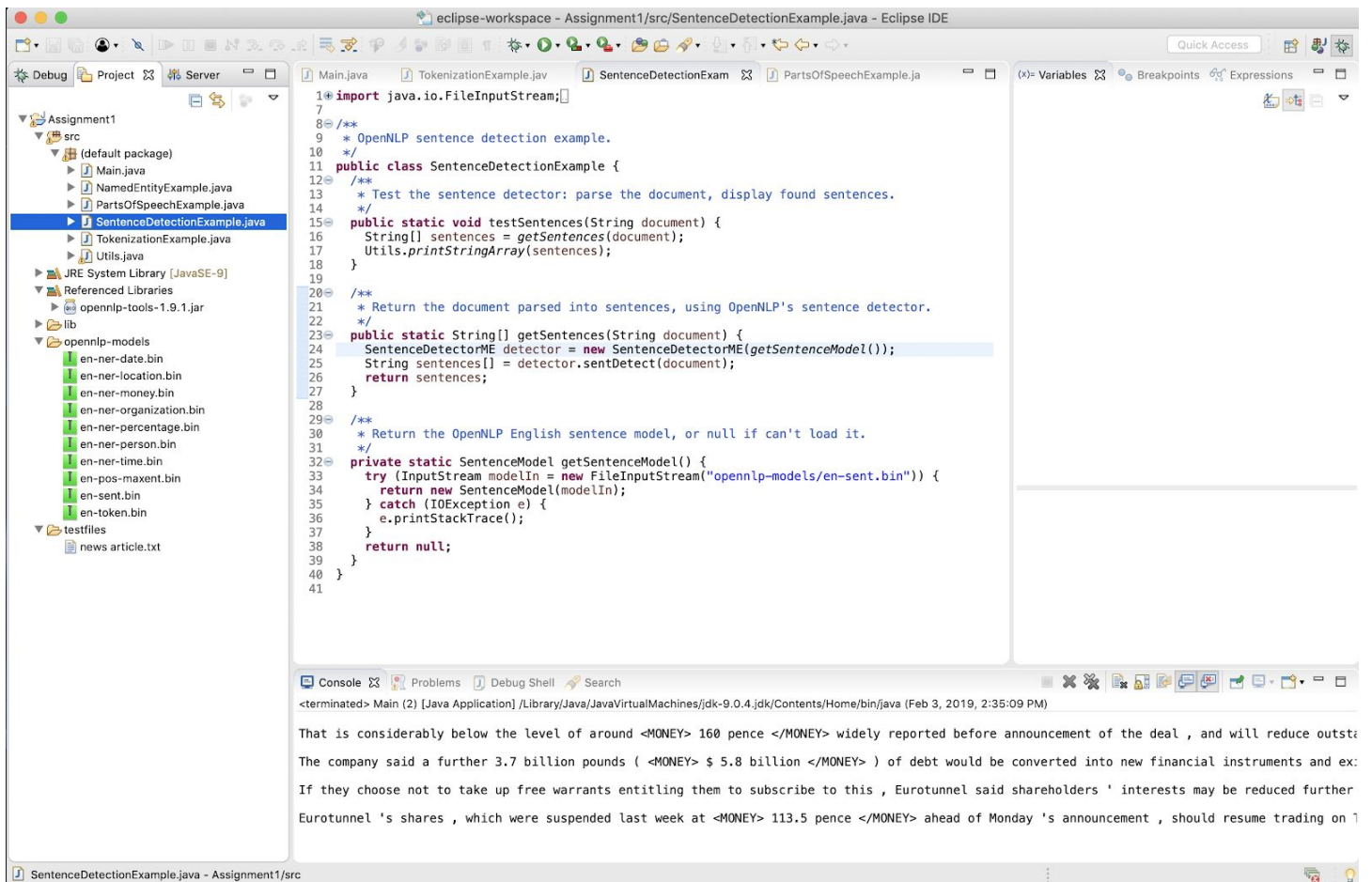
The Java code used for this report is submitted in a separate file. It does not contain OpenNLP components. The arrows indicate OpenNLP libraries and models that have to be downloaded and added to the project.



Source code for the tasks is listed below.

Screenshots






```

1 import java.io.FileInputStream;
9 /**
10 * OpenNLP parts of speech (POS) example.
11 */
12 public class PartsOfSpeechExample {
13
14     public static void testPartsOfSpeech(String document) {
15         POSModelME tagger = new POSModelME(getPOSModel());
16
17         // Use ME tokenization for POS because ME tokenization does a better job of
18         // returning tokens that are meaningful for humans.
19         List<TokenizationExample.TokenizedSentence> tokens = TokenizationExample.getMETokens(document);
20         for (TokenizationExample.TokenizedSentence ts : tokens) {
21             String[] tags = tagger.tag(ts.tokens);
22             printTags(ts.sentence, ts.tokens, tags);
23         }
24     }
25
26     /**
27     * Return the OpenNLP English maximum entropy POS model, or null if can't load
28     * it.
29     */
30     private static POSModel getPOSModel() {
31         try (InputStream modelIn = new FileInputStream("opennlp-models/en-pos-maxent.bin")) {
32             return new POSModel(modelIn);
33         } catch (IOException e) {
34             e.printStackTrace();
35         }
36         return null;
37     }
38
39     private static void printTags(String sentence, String[] tokens, String[] tags) {
40         System.out.println("\n" + sentence + "\n");
41         for (int i = 0; i < tokens.length; i++) {
42             System.out.print(String.format("%02d %-15s %s", i, tokens[i], tags[i]));
43         }
44     }
45 }
46
47

```

Console Output:

```

[06] diluted
[07] ,
[08] they
[09] were
[10] offered
[11] the
[12] prospect
[13] of
[14] a
[15] brighter

```

```

1 import java.io.FileInputStream;
9
10 public class NamedEntityExample {
11
12     public static void testNameFinder(String document) {
13         // Get the different named entities
14         NameFinderME personFinder = new NameFinderME(getNameFinderModel("en-ner-person.bin"));
15         NameFinderME locationFinder = new NameFinderME(getNameFinderModel("en-ner-location.bin"));
16         NameFinderME moneyFinder = new NameFinderME(getNameFinderModel("en-ner-money.bin"));
17         NameFinderME percentageFinder = new NameFinderME(getNameFinderModel("en-ner-percentage.bin"));
18
19         List<TokenizationExample.TokenizedSentence> tokens = TokenizationExample.getMETokens(document);
20         for (TokenizationExample.TokenizedSentence ts : tokens) {
21             // Get all different types of entities
22             Span personSpans[] = personFinder.find(ts.tokens);
23             Span locationSpans[] = locationFinder.find(ts.tokens);
24             Span moneySpans[] = moneyFinder.find(ts.tokens);
25             Span percentageSpans[] = percentageFinder.find(ts.tokens);
26
27             // Combine all entities into one set
28             Span[] entities = Utils.concatAll(personSpans, locationSpans, moneySpans, percentageSpans);
29
30             // Remove overlapping spans to simplify this code - in real life we would
31             // probably want to know about overlapping spans of different types.
32             entities = NameFinderME.dropOverlappingSpans(entities);
33
34             // Show tokens tagged with the named entity, if there is one for that
35             // token. This code is not particularly efficient, but it's easy to follow.
36             for (int t = 0; t < ts.tokens.length; t++) {
37                 // Check if this token is the start of a named entity
38                 for (Span span : entities) {
39                     if (span.getStart() == t) {
40                         System.out.print("<" + span.getType().toUpperCase() + "> ");
41                     }
42                 }
43
44                 // Show the token
45                 System.out.print(ts.tokens[t] + " ");
46
47                 // Check if this token is the end of a named entity
48                 for (Span span : entities) {
49                     if (span.getEnd() == t) {
50                         System.out.print("> ");
51                     }
52                 }
53             }
54         }
55     }
56 }
57
58

```

Console Output:

```

That is considerably below the level of around <MONEY> 160 pence </MONEY> widely reported before announcement of the deal , and will reduce outsta
The company said a further 3.7 billion pounds ( <MONEY> $ 5.8 billion </MONEY> ) of debt would be converted into new financial instruments and ex
If they choose not to take up free warrants entitling them to subscribe to this , Eurotunnel said shareholders ' interests may be reduced further
Eurotunnel 's shares , which were suspended last week at <MONEY> 113.5 pence </MONEY> ahead of Monday 's announcement , should resume trading on 1

```