# Apache UIMA, Part 2
## PU Tools, Ressourcen, Infrastruktur

Nils Reiter,
nils.reiter@uni-koeln.de

January 13, 2022
(Winter term 2020/21)

# Section 1

## Recap

## Exercise 10

https://github.com/idh-cologne-tools-ressourcen-infra/exercise-10

# Section 2

# Apache UIMA, Part 2

# Writing UIMA Components
Select Interface

▶ UIMA indexes all feature structures, such that we can efficiently access them

```
1 SelectFSs<Token> selector = jcas.select(Token.class);
2 selector.asList(); // a list of all feature structures
3 selector.count(); // number of feature structures
4 selector.get(index); // get nth token
```

# Writing UIMA Components
Select Interface

▶ UIMA indexes all feature structures, such that we can efficiently access them

```
1 SelectFSs<Token> selector = jcas.select(Token.class);
2 selector.asList(); // a list of all feature structures
3 selector.count(); // number of feature structures
4 selector.get(index); // get nth token
5 selector.coveredBy(anotherAnnotation); // feature structures
6                                         // covered by another
```

# Writing UIMA Components
Select Interface

▶ UIMA indexes all feature structures, such that we can efficiently access them

```
1 SelectFSs<Token> selector = jcas.select(Token.class);
2 selector.asList(); // a list of all feature structures
3 selector.count(); // number of feature structures
4 selector.get(index); // get nth token
5 selector.coveredBy(anotherAnnotation); // feature structures
6                                         // covered by another
7 selector.following(position); // first token after position
```

# Writing UIMA Components
Select Interface

▶ UIMA indexes all feature structures, such that we can efficiently access them

```
1 SelectFSs<Token> selector = jcas.select(Token.class);
2 selector.asList(); // a list of all feature structures
3 selector.count(); // number of feature structures
4 selector.get(index); // get nth token
5 selector.coveredBy(anotherAnnotation); // feature structures
6                                        // covered by another
7 selector.following(position); // first token after position
8 selector.allMatch(predicate); // get tokens for which predicate
9                               // is fulfilled
```

# Writing UIMA Components
Select Interface: coveredBy()

```java
1  // iterate over sentences
2  for (Sentence sentence : jcas.select(Sentence.class)) {
3    // iterate over tokens in a sentence
4    for (Token token : jcas.select(Token.class).coveredBy(sentence)) {
5      // do stuff with each token
6    }
7  }
```

# Writing UIMA Components

Select Interface: allMatch()

$$\text{Predicate} \langle Tok \rangle \quad prd = t \rightarrow t.getId() > 15;$$

```
1  // iterate over sentences that start with "T"
2  for (Sentence sentence : jcas.select(Sentence.class)
3       .allMatch(s -> s.getCoveredText().startsWith("T"))) {
4    // iterate over tokens in sentence with id larger than 15
5    for (Token token : jcas.select(Token.class).coveredBy(sentence)
6         .allMatch(t -> t.getId() > 15)) {
7      // do stuff with token
8    }
9  }
```

prd

# Lambda Expressions in Java

$t \rightarrow t.getId() > 15$

```
1 // long form
2 Predicate<Token> predicate = new Predicate<Token>() {
3   public boolean test(Token t) {
4     return t.getId() > 15;
5   }
6 };
7
```

## Lambda Expressions in Java

```
1 // long form
2 Predicate<Token> predicate = new Predicate<Token>() {
3    public boolean test(Token t) {
4      return t.getId() > 15;
5    }
6 };
7
8 // short form ("syntactic sugar")
9 Predicate<Token> predicate = t -> t.getId() > 15;
```

## Lambda Expressions in Java

```java
1 // long form
2 Predicate<Token> predicate = new Predicate<Token>() {
3   public boolean test(Token t) {
4     return t.getId() > 15;
5   }
6 };
7
8 // short form ("syntactic sugar")
9 Predicate<Token> predicate = t -> t.getId() > 15;
```

▶ java.util.function package (since 1.8)
▶ Functional interface: Classes/interfaces with a single method
  ▶ Predicate: $T \rightarrow$ boolean (public boolean test(T))
  ▶ BiPredicate: $T, S \rightarrow$ boolean (public boolean test(T, S))
  ▶ Function: $T \rightarrow S$
  ▶ …

# Section 3

# Resources and Arguments

# Resources and Arguments

- Many components need resources and/or parameters
- Define fields as usual
- Components may implement a method `public void initialize(UimaContext)`
  - $\simeq$ constructor, but frameworks handles construction as needed
- Arguments need to be passed as atomic values (or read from file)

# Resources and Arguments

- ▶ Many components need resources and/or parameters
- ▶ Define fields as usual
- ▶ Components may implement a method `public void initialize(UimaContext)`
  - ▶ $\simeq$ constructor, but frameworks handles construction as needed
- ▶ Arguments need to be passed as atomic values (or read from file)

```
 1 public class MyComponent extends JCasAnnotator_ImplBase {
 2
 3   @Override
 4   public void initialize(final UimaContext context)
 5     throws ResourceInitializationException {
 6
 7     // call super class method
 8     super.initialize(context);
 9
10   }
11
12   @Override
13   public void process(JCas jcas) throws AnalysisEngineProcessException { }
14 }
```

# Resources and Arguments

- ▶ Declare parameters as class fields
- ▶ Use @ConfigurationParameter() annotation to mark as configuration parameter
  - ▶ Javadoc:
    https://javadoc.io/static/org.apache.uima/uimafit-core/3.1.0/org/apache/uima/fit/descriptor/ConfigurationParameter.html
  - ▶ Arguments
    - ▶ name = ... Define a name for the parameter
    - ▶ defaultValue = ... Default value (as a string array)
    - ▶ description = ... Human-readable description
    - ▶ mandatory = ... Defines whether it can be omitted, boolean value

# Resources and Arguments

- ▶ Declare parameters as class fields
- ▶ Use @ConfigurationParameter() annotation to mark as configuration parameter
    - ▶ Javadoc:
        https://javadoc.io/static/org.apache.uima/uimafit-core/3.1.0/org/apache/uima/fit/descriptor/ConfigurationParameter.html
    - ▶ Arguments
        - ▶ name = ... Define a name for the parameter
        - ▶ defaultValue = ... Default value (as a string array)
        - ▶ description = ... Human-readable description
        - ▶ mandatory = ... Defines whether it can be omitted, boolean value
- ▶ Do I need to define initialize(UimaContext)?
    - ▶ Not for purely filling given argument values into fields variables
    - ▶ Only if there is something to be done with the arguments
        - ▶ E.g., loading a file

# Resources and Arguments: Best Practice

```
1  public class MyComponent extends JCasAnnotator_ImplBase {
2    // Define a static public variable that is used instead of the name
3    public static final String PARAM_WINDOW_SIZE = "window size";
4
5    @ConfigurationParameter(name = PARAM_WINDOW_SIZE, defaultValue="10", mandatory=false)
6    int windowSize = 10;
7
8    @Override
9    public void process(JCas jcas) throws AnalysisEngineProcessException {
10     for (Token token : jcas.select(Token.class)) {
11       // get all tokens
12       List<Token> nextTokens = jcas.select(Token.class)
13         // ... that follow this token
14         .following(token)
15         // ... and only the first windowSize, as list
16         .limit(windowSize).asList();
17       // do something
18     }
19   }
20 }
```
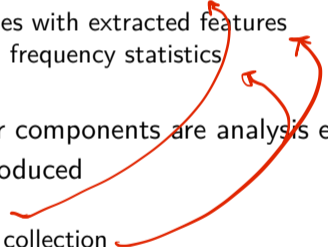
# Section 4

# Writer Components

# Writer Components

▶ We often need at least one component that produces output
  ▶ HTML files with the annotations
  ▶ CSV/ARFF files with extracted features
  ▶ CSV files with frequency statistics
  ▶ …

# Writer Components

- ▶ We often need at least one component that produces output
  - ▶ HTML files with the annotations
  - ▶ CSV/ARFF files with extracted features
  - ▶ CSV files with frequency statistics
  - ▶ …
- ▶ Technically, writer components are analysis engines like any other component
- ▶ Output can be produced
  - ▶ Per document
  - ▶ For the entire collection

# Writer Components
Writing to Files in Java

▶ Output writing in Java uses streams
▶ Stream: A stream of characters that flows from your code to somewhere else

## Example

```java
1 File file = new File("PATH/TO/MY/FILE");
2 FileOutputStream os = new FileOutputStream(file);
3 OutputStreamWriter writer = new OutputStreamWriter(os);
4 writer.write("This is the file content I want to write");
5 writer.flush();
6 writer.close();
```

# Writer Components

## Output per Document

▶ Done in `public void process(JCas)` method

▶ Usual Java mechanisms for opening files and writing to them

▶ Helpful: Inherit from `org.dkpro.core.api.io.JCasFileWriter_ImplBase`
  ▶ New method `getOutputStream(JCas, String)` can be used
  ▶ Creates a file with the same name as the input file (except extension)

# Writer Components

## Output per Document

- Done in `public void process(JCas)` method
- Usual Java mechanisms for opening files and writing to them
- Helpful: Inherit from `org.dkpro.core.api.io.JCasFileWriter_ImplBase`
  - New method `getOutputStream(JCas, String)` can be used
  - Creates a file with the same name as the input file (except extension)

## Output for entire Collection

- Implement method `public void collectionProcessComplete()`
- Called after all documents have been processed

# Section 5

## Next Exercise

https://github.com/idh-cologne-tools-ressourcen-infra/exercise-11