

Programming with Maltcms 0.95

Nils Hoffmann

Nils.Hoffmann@CeBiTec.Uni-Bielefeld.DE

12. November 2009

Inhaltsverzeichnis

1	Introduction	1
1.1	Audience	2
2	Components	2
2.1	A document tree datastructure: IFileFragment and FileFragment	2
2.1.1	Referencing previous IFileFragment Objects	2
2.2	Manipulating document trees: AFragmentCommand	2
2.3	Building a pipeline with CommandSequence	5
2.4	Manipulating Arrays: ArrayCommand	5
3	Tools	5
3.1	ArrayTools	5
3.2	MaltcmsTools	5
3.3	FileTools	5
4	CLI, Configuration and Logging Infrastructure	5
4.1	Command Line Interface with Jakarta Commons CLI	6
4.2	Configuration with Jakarta Commons	6
4.3	Logging with Simple Logging Facade for Java: SLF4J	6
5	Supporting Scripts	6
5.1	scripts/maltcms.sh	6
5.2	scripts/buildCP.sh	8
5.3	scripts/buildAntCP.sh	9
5.4	scripts/setJavaHome.sh	9
5.5	scripts/updateBuildNumber.sh	10

1 Introduction

The following Tutorial was written as an introduction to the Java based framework Maltcms (Modular Application Toolkit for Chromatography Mass Spectrometry). It tries to give an overview of key classes and concepts, as well as in depth information, where necessary.

1.1 Audience

The target audience of this tutorial are typically programmers or scientists, wishing to extend the functionality of Maltcms, or wishing to incorporate Maltcms functionality into their own projects.

2 Components

The following section gives a short overview of the core components, mainly situated in the generic package `cross.*` (Common Runtime Object Support System). These components provide basic datastructures and infrastructure for logging, IO and other things.

2.1 A document tree datastructure: `IFileFragment` and `FileFragment`

¹ Maltcms makes use of the netcdf² file format which is based on an abstract model for scientific data. This model allows for naming of variables, attributes and dimensions within one file, thus creating self-describing data. There is also a possibility for using structures and groups within such files. However, in Maltcms only a subset of this model is used, and mirrored as an extended, internal data structure, which allows for serialization to XML and lazy (on demand) loading of array data, as well as indexed reading in row compressed storage (rcs) format of one array with indices stored in another array.³ The internal structure can be imagined as a document tree, where the root of the tree is an instance of `IFileFragment`. It is possible to add instances of `IVariableFragment` to the root, as well as removing them. Each `IVariableFragment` has different attributes, which can be altered.

Arbitrary input file formats are mapped to the internal array style format by corresponding `cross.io.IDataSource` implementations. Such a data file must provide methods for reading single array data, as well as for reading scan oriented indexed arrays. The implementation is required to react on passed in references to objects of type `cross.datastructures.fragments.IVariableFragment`, which provide a name for the data part associated with that name. The names of variables mimic those defined in the AIA/ANDI-MS standard, so that once the internal data structure has been established from a different input file via an appropriate `cross.io.IDataSource` implementation, the data can be saved to an AIA/ANDI-MS compatible netcdf file.

2.1.1 Referencing previous `IFileFragment` Objects

2.2 Manipulating document trees: `AFragmentCommand`

The classes extending `cross.commands.fragments.AFragmentCommand` build the core classes of a processing pipeline. Each such class must provide a method

¹found in package `cross.datastructures.fragments`

²<http://www.unidata.ucar.edu/software/netcdf/>

³The associated classes can be found in the package `maltcms.datastructures.fragments`. The package name will be omitted from now on, unless necessary for unique identification of classes.

public TupleND<IFileFragment> apply(TupleND<IFileFragment> t) to receive a tuple of IFileFragment objects and return them possibly processed or replaced.

As an example of how such a AFragmentCommand is created, the class Normalizer will be used as an example. Let's focus on the necessary imports first:

```

1  /* The package declaration is placed at the top of the File
2  * Since this example will be application domain specific, we
3  * place it under maltcms.commands.fragments
4  */
5  package maltcms.commands.fragments;
6  // The most basic imports are
7  import cross.datastructures.fragments.IFileFragment;
8  import cross.datastructures.fragments.FileFragment;
9  import cross.datastructures.tuple.TupleND;
10 // In order to use the logging and configuration infrastructure of Maltcms,
11 // we also need the following imports
12 import cross.Logging;
13 import org.slf4j.Logger;
14 import org.apache.commons.configuration.Configuration;
15 // We also add an import to throw Exceptions from not yet implemented methods
16 import maltcms.exception.NotImplementedException;

```

After having imported the basic dependencies necessary for our Normalizer, we can begin describing our class.

```

1  // Now we can start describing Normalizer, which extends
2  // FragmentCommand
3  public class Normalizer extends AFragmentCommand {
4  // Define some private variables
5      private boolean normalize_mass_channels = false;
6      private boolean normalize_scan_intensities = false;
7      private Logger log = Logging.getLogger(this.getClass());
8  // Implement the apply method of FragmentCommand
9      public TupleND<IFileFragment> apply(TupleND<IFileFragment> t) {
10 // Let it throw a NotImplementedException for now:
11          throw new NotImplementedException();
12      }
13      ...

```

```

1  public class Normalizer extends AFragmentCommand {
2      ...
3  // Next, we should override the configure method of
4  // Configurable to allow the Factory to configure objects of this class
5  // at runtime. Since Java version 1.5, overridden methods should be annotated
6  // with @Override.
7      @Override
8      public void configure(Configuration cfg) {
9  // We will set the private variables in here and publish log messages in
10 // info-mode, the second parameter in cfg.getBoolean(String s, boolean b)
11 // is the default value, if s is not found within the configuration.
12          this.normalize_mass_channels=cfg.getBoolean("normalize.
13              mass_channels",true);
14  // The braces within the logging statement are used for pattern
15 // substitution, to prevent costly variable expansion/evaluation if we are in a
16 // different logging mode.
17          log.info("Normalizing_mass_channels=_{}",this.
18              normalize_mass_channels);
19          this.normalize_scan_intensities=cfg.getBoolean("normalize.
20              scan_intensities",true);
21          log.info("Normalizing_scan_intensities=_{}",this.
22              normalize_scan_intensities);

```

```

19     }
20 }

```

In order to put some functionality into our class, we will begin implementing the `apply` method.

```

1 public TupleND<IFileFragment> apply(TupleND<IFileFragment> t)
2 {
3     // Iterate over all IFileFragment objects in t
4     for(IFileFragment ff:t){
5         // Create a IFileFragment with default name to take the
6         // results of operation
7         IFileFragment work = FragmentTools.create(this.
8             getClass());
9         // Query for a specific variable
10        if(ff.hasChildren("intensity_values","scan_index")){
11            // intensity_values is an indexed variable
12            VariableFragment intensity_values = ff.getChild("
13                intensity_values");
14            VariableFragment scan_index = ff.getChild("
15                scan_index");
16            intensity_values.setIndex(scan_index);
17        }
18    }
19 }

```

Now that we have some `VariableFragments` ready, we can try to read their associated arrays.

```

1 // Returns the scans contained in intensity_values.
2 // Read access uses the index variable scan_index
3 ArrayList<Array> intens_scan_arrays = intensity_values.
4     getIndexArray();
5 // Traversing the scans
6 for(Array arr:intens_scan_arrays) {
7     //we can try to cast an array to a more concrete type
8     if(arr instanceof ArrayInt.D1){
9         ArrayInt.D1 scan_i = (ArrayInt.D1) arr;
10        //read value 100 of scan
11        int value = scan_i.get(100);
12    }
13 // otherwise, we can always use an IndexIterator to traverse all values
14 IndexIterator iter = arr.getIndexIterator();
15 while(iter.hasNext()) {
16     double d = iter.getDoubleNext();
17 }
18 // We can also read the array of intensity_values as a whole,
19 // in it's row compressed storage representation
20 Array intens_values_array = intensity_values.getArray();
21 // We can of course read scan_index as well
22 Array a = scan_index.getArray();
23 // If the type of an array and it's shape are known in advance,
24 // we can directly cast to the appropriate type.
25 // This allows for direct element access, without using an Index.

```

```
26 | ArrayInt.D1 scan_index_array = ((ArrayInt.D1)a);
```

2.3 Building a pipeline with CommandSequence

Command sequence is a collection type class, allowing a special kind of iteration. Initialization is performed by passing in a List of `FragmentCommand` objects and input `IFileFragment` objects. Then, upon each invocation of the usual

```
if(hasNext)->next()
```

paradigm on iterable objects, the output of the last active `FragmentCommand` is pushed as input to the next `FragmentCommand`. It is thus possible to build processing chains⁴, which start for example with loading of default values, then preprocessing, then continue applying a feature detection algorithm and end with a visualization. The file `cfg/default_pipeline.properties` contains a key to set the elements of the pipeline used in left to right order called `pipeline`.

The simple pipeline scheme has some limitations though, if for example a `FragmentCommand` creates additional data, which is needed by `FragmentCommand` later in the pipeline. For this case, we currently use the configuration, obtainable via `Factory.getConfiguration()`, adding a new key-value mapping pointing e.g. to the file representation of the data.

```
1 //The following line would load FragmentCommand objects
2 //as given under property anchors.class, loading possibly
  available retention
3 //indices, then would run the default.varloader, which loads
  all variables listed under
4 //default.variables. Then, the dense arrays are produced from
  row compressed storage arrays
5 //and finally, pairwise distances between the arrays are
  calculated
6 pipeline = ${anchors.class},${default.varloader},${dense.
  arrays},${pairwise.distances}
```

2.4 Manipulating Arrays: ArrayCommand

3 Tools

3.1 ArrayTools

3.2 MaltcmsTools

3.3 FileTools

4 CLI, Configuration and Logging Infrastructure

Maltcms uses pre-existing software solutions to provide a common logging, command line and configuration system. Those systems are in wide use and thus rather bug free.

⁴Even for batch processing, since n-Tuples of `IFileFragment` serve as input.

4.1 Command Line Interface with Jakarta Commons CLI

The class `apps.Maltcms` uses the command line interface provided by Apache Jakarta Commons CLI and defines custom command line parameters, which can be inspected by typing the usual `-?`, `-help`, `-h` options after the program name.

4.2 Configuration with Jakarta Commons

Configuration in Maltcms is usually done in java properties file format, e.g. `PROPERTYNAME = VALUE`. The main configuration options are stored in `cfg/default.properties`. Configuration files can be aggregated by including other configuration files: `include = other.properties`.

4.3 Logging with Simple Logging Facade for Java: SLF4J

The simple logging facade allows easy exchange of the underlying logging system. Currently, apache log4j is used, but potentially any other compatible logging framework, such as the jdk's own logging could be used.

5 Supporting Scripts

5.1 scripts/maltcms.sh

```
1 #!/bin/bash
2 export JARCH=""
3 export JBIN="\vol\java-1.6.0\bin"
4 export EXEC="apps.Maltcms"
5 export MXSIZE="2G"
6 export MSSIZE="256M"
7 export MALTCMSARGS=""
8 export PROFILE=""
9 export USRCLSPATH=""
10
11 if [ -z "$MALTCMSDIR" ]; then
12     echo "Please_enter_path_to_Maltcms_installation_or_add_
13         MALTCMSDIR_to_your_bash_profile:"
14     read MALTCMSUSRDIR;
15     if [ -z "$MALTCMSUSRDIR" ]; then
16         echo "No_user_defined_directory_for_Maltcms_
17             installation_entered,no_default_given,
18             exiting!";
19         exit 1;
20     fi
21 else
22     MALTCMSUSRDIR=$MALTCMSDIR;
23 fi
24 LOG4J_LOCATION="-Dlog4j.configuration=file://$MALTCMSUSRDIR/
25     cfg/log4j.properties"
26
27 #Check if javahome exists => contains path to java
28 if [ -f $MALTCMSUSRDIR/javahome ]; then
```



```

69         -mx)
70             shift
71             export MXSIZE="$1"
72             ;;
73         -ms)
74             shift
75             export MSSIZE="$1"
76             ;;
77         --)
78             echo "Running a_${JARCH}_VM_ with -Xmx_${MXSIZE}"
79             shift
80             #Check for clspath file
81             if [ -f $MALTCMSUSRDIR/clspath ]; then
82                 echo "File clspath exists";
83             else
84                 $MALTCMSUSRDIR/scripts/buildCP.sh
85             fi
86             USRCLSPATH="$CLASSPATH:$(cat _$MALTCMSUSRDIR/
87                 clspath)"
88             echo -e "Passing args to _$EXEC"
89             echo -e "$@"
90             sleep 1
91             $JAVA_HOME/bin/java -cp $USRCLSPATH $PROFILE -
92                 Xms$MSSIZE -Xmx$MXSIZE $JARCH $LOG4J_LOCATION
93                 $EXEC "$@"
94             exit $?
95             ;;
96         -"?"|--help)
97             printHelp $0
98             ;;
99     esac
100     shift
101 done
102 exit 1

```

5.2 scripts/buildCP.sh

```

1  #!/bin/bash
2  if [ "$1" != "" ]; then
3      FPATH=$1
4  else
5      FPATH=$(pwd)/scripts
6  fi
7  if [ -e clspath ]; then
8      rm clspath
9  fi
10 for i in `ls $(echo -e $FPATH | sed -e s/scripts/lib/g)/*.jar`
11 do
12     echo -n ":$i" >> clspath;
13 done

```



```

14 #for i in `ls $(echo -e $FPATH | sed -e s/scripts/cfg/g)/*.
    properties `
15 #do
16 #    echo -n ":$i" >> clspath;
17 #done
18 #echo -e $FPATH
19 PATH=`echo -e $FPATH | sed -e s/scripts/bin/g`
20 echo -e ":$PATH" >> clspath
21 #export CLASSPATH="$CLASSPATH$(cat clspath)"
22 #echo $CLASSPATH

```

5.3 scripts/buildAntCP.sh

```

1 #!/bin/bash
2 if [ "$1" != "" ]; then
3     FPATH=$1
4 else
5     FPATH=$(pwd)/scripts
6 fi
7 if [ -e antclspath ]; then
8     rm antclspath
9 fi
10 for i in `ls $(echo -e $FPATH | sed -e s/scripts/lib/g)/*.jar `
11 do
12     echo -en "<pathelement location=\"$i\"/>\n" >> antclspath;
13 done

```

5.4 scripts/setJavaHome.sh

```

1 #!/bin/bash
2 ARCH=`uname -a`
3 OSTYPE=`uname -s`
4 #echo -e "Running on_$ARCH\n"
5 JAVA_HOME=""
6 case "$OSTYPE" in
7     SunOS)
8         JAVA_HOME="/vol/java-1.6/"
9         ;;
10     Darwin)
11         JAVA_HOME="/System/Library/Frameworks/JavaVM.framework
            /Versions/1.6.0/Home/"
12         ;;
13     Linux)
14         JAVA_HOME="/usr/lib/jvm/java-6-sun/"
15         ;;
16     *)
17         echo "Unknown_OS_$OSTYPE, please set the path to your
            Java-Installation by hand."
18         echo "Save it in a file called javahome within maltcms
            _basedir"
19         ;;

```

```

20 esac
21 echo -e "$JAVA_HOME" > javahome
22 echo $JAVA_HOME

```

5.5 scripts/updateBuildNumber.sh

```

1 #!/bin/bash
2 REVISION='svn info | grep -e Revision | sed -e "s/Revision: _//
   "'
3 DATE='date +%Y_%b_%d'
4 sed -i -e "s/application.version.revision=r[0-9]*/application
   .version.revision=r$REVISION/" cfg/application.properties
5 sed -i -e "s/application.build.date=-[0-9A-Za-z_]*/application
   .build.date=$DATE/" cfg/application.properties
6 cut -f 2 -d= cfg/application.properties > versionfile;
7 grep -v "\\$" versionfile > tmpversionfile;
8 paste -s -d\\0 tmpversionfile > versionfile;
9 rm tmpversionfile;
10 sed -e "s/PROJECT_NUMBER\\W*=.* /PROJECT_NUMBER=$(cat _
   versionfile)/g" -i Doxyfile
11 cat versionfile;

```