



university of  
groningen

faculty of science  
and engineering

# Dependency Graph Creator Engine

Software Maintenance and Evolution

*<https://github.com/richardswesterhof/pyne>*

## Authors:

Job Heersink (s3364321)  
Richard Westerhof (s3479692)

## Supervisors:

Mohamed Soliman  
Filipe Capela

University of Groningen  
The Netherlands  
January 12, 2021



# CONTENTS

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>User Guide</b>	<b>4</b>
2.1	Running pyne . . . . .	4
2.2	Running Structure101 . . . . .	4
2.3	Running the dependency checker . . . . .	6
<b>3</b>	<b>Design</b>	<b>7</b>
3.1	Requirements . . . . .	7
3.2	Implementation . . . . .	8
3.2.1	Pyne-demo . . . . .	9
3.2.2	Pyne-cli . . . . .	9
3.2.3	Pyne-api . . . . .	10
3.3	External dependencies . . . . .	11
<b>4</b>	<b>Pyne vs Structure101</b>	<b>13</b>
4.1	Running Structure101 on Tajo . . . . .	13
4.2	Design of Tajo . . . . .	14
4.2.1	components . . . . .	14
4.2.2	External dependencies . . . . .	17
4.3	Running Pyne on Tajo . . . . .	18
4.4	Comparing the feature sets . . . . .	20
4.5	Comparing the Results of Pyne to Structure101 . . . . .	20
<b>5</b>	<b>Comparing results programmatically</b>	<b>21</b>
5.1	Creating a dependency checker program . . . . .	21
5.1.1	Design of the dependency checker . . . . .	21
5.2	Running our dependency checker program . . . . .	23
<b>6</b>	<b>Improving Pyne</b>	<b>24</b>
6.1	Faults in Pyne . . . . .	24
6.2	Changes to Pyne . . . . .	27
6.3	Result . . . . .	29
6.3.1	Approach . . . . .	29

---

6.3.2	Analyses . . . . .	31
<b>7</b>	<b>Discussion</b>	<b>34</b>
<b>8</b>	<b>Conclusion and Future work</b>	<b>35</b>
	<b>References</b>	<b>36</b>
<b>A</b>	<b>Troubleshooting</b>	<b>38</b>
<b>B</b>	<b>Change Log</b>	<b>40</b>
<b>C</b>	<b>Dependency checker output</b>	<b>42</b>
C.1	Output Template . . . . .	42
C.2	original version of Pyne . . . . .	44
C.2.1	classes . . . . .	44
C.2.2	packages . . . . .	45
C.3	improved version of Pyne . . . . .	46
C.3.1	classes . . . . .	46
C.3.2	packages . . . . .	47

# CHAPTER 1

## INTRODUCTION

Software maintenance has always been of great concern to software engineers. An application is almost never perfect on deployment and needs regular bug-fixes and enhancements to the system. To help with the software maintenance of an application, several different techniques have been developed.

One of these techniques is a dependency graph. A dependency graph is a directed graph that represents dependencies between objects of some application domain, where the nodes represent packages or classes of a software system and the edges the relation between them.

In this document we will compare an incomplete existing dependency graph creator called Pyne[1] with another widely used dependency graph creator engine called

The document is structured as follows: Chapter 2 will show how Pyne should be installed and used. Chapter 3 will discuss the design of Pyne itself. In Chapter 4 we will apply Pyne and Structure101 to a large java application called Apache Tajo[2]. In Chapter 5 we will discuss our dependency checker to evaluate the output of pyne. In chapter 6 we will discuss the problems of pyne and how we tried to fix them and finally in Chapter 7 and 8 we will discuss and conclude the project.

# CHAPTER 2

## USER GUIDE

### 2.1 RUNNING PYNE

#### Installing

To install Pyne, Java JDK 11+ and Maven is required. The installation itself is rather straightforward: Clone or download our fork of the Pyne repository:

```
1 git clone https://github.com/richardswesterhof/pyne.git
```

Install Pyne's dependencies and build an executable jar by running the following command inside Pyne's root directory:

```
1 mvn install
```

#### Executing

After the application has been built, one can use this application via CLI. This can be done by running the following command from the root directory:

```
1 java -jar <path-to-jar>/pyne-cli-1.0-SNAPSHOT-jar-with-dependencies.jar  
  ↪ <github-url>
```

Where the `<path-to-jar>` is the path to the built jar file (by default this is in `<pyne_root>/pyne-cli/target/`) and `<github-url>` is the url of the project you want to create a dependency graph for. For example, this would be the command to create a yearly dependency graph for the commits of Apache Tajo from 2015 to November 2020:

```
1 java -jar pyne-cli/target/pyne-cli-1.0-SNAPSHOT-jar-with-dependencies.jar  
  ↪ https://github.com/apache/tajo -s "2015-01-01" -e "2020-11-21" -p  
  ↪ "YEAR"
```

### 2.2 RUNNING STRUCTURE101

To generate the output from Structure101, follow the steps below:

1. Make sure Structure101 is installed.
2. Create a new project, select "Maven" as the option for discovering bytecode, select the pom file of the project you want to analyse, and make sure to check the "Parse Tests" and "Parse Profiles" boxes to include as many sources as possible.
3. Now you should already be able to select the way you want everything to be organized, i.e. whether you want packages, leaf packages or classes. If you only want to analyse one of the two, you can select this now (note that for analysing on a package level, make sure to select "Leaf packages" instead of just "package"). However, if you want to analyse both, it is easier to select "Packages" here.
4. Select the "Detail" granularity, and check "Included injected dependencies" to again cover as many dependencies as possible.
5. Select "Show externals" so that external packages are shown as well as internal packages, and check "Parse archives contained in classpath archives".
6. Leave the list of exclusions empty, and select the project you want to analyse's folder as the sources.

Now the project is set up, and the results can be exported to CSV. If there are any issues during the setup of the project, there is a Structure101 project for Tajo available on our github repository<sup>1</sup>.

7. To generate the expected matrix from Structure101, right-click anywhere outside the packages in the main graph view, hover over the "Flatten" option, and select "To leaf packages" (for comparing on package level) or "To classes" (for comparing on class level) from the sub-menu (skip this step if you already did this during the project setup).
8. Go to the "View" tab and select the composition view. You should now see the matrix that shows the amount of dependencies between each class/package.
9. A CSV file must be exported by clicking the export icon in the top right toolbar in the matrix panel. Make sure that you select "file" for the "Export to" option, then select "Matrix as CSV" in the "Export as" dropdown. Finally, select a file location to save the file to in the "Target file" field, and click "Ok" to export.
10. (optional) If you want to do an analysis on both class and package level, go back to the main view and use control+z to undo the flattening, then return to step 7.

Now you have the output file necessary from Structure101. This is the file you should use when using the dependency checker.

---

<sup>1</sup>[https://github.com/richardswesterhof/pyne/blob/master/structure101/tajo\\_structure101.java.hsp](https://github.com/richardswesterhof/pyne/blob/master/structure101/tajo_structure101.java.hsp)

## 2.3 RUNNING THE DEPENDENCY CHECKER

To use the dependency checker, Java JDK 11+ and Maven is required. The installation itself is rather straightforward: Clone or download our fork of the Pyne repository. The dependency checker will be in there too.

```
1 git clone https://github.com/richardswesterhof/pyne.git
```

Install the dependency checker's dependencies and build an executable jar by running the following command inside the `dependency_checker/` folder:

```
1 mvn install
```

This will place a ready-to-execute jar file in the `target/` subfolder. Alternatively, a pre-built jar file is already located in the `dependency_checker/` folder.

### Executing

After the application has been built, one can use this application via CLI. This can be done by running the following command from the directory where the jar file is located:

```
1 java -jar dependency_checker.jar -s [PATH/TO/STRUCTURE101_FILE.csv] -p  
    ↪ [PATH/TO/PYNE_FILE.graphml] -d [CLASS|PACKAGE] [OPTIONS]
```

With `-d` you can specify whether you want to compare the dependencies between packages or classes, with `-p` you must specify the path to the `.graphml` output-file of pyne and with `-s` you must specify the path to the Structure101 `.csv` file. This file must either contain the relation of the classes or the packages depending on the specified option for `-d`. If the `.csv` file contains class relations, while `-d` is set to package, the dependency checker will not work. The other available options include:

The `-hr` option will create a human readable version of the output (with tabs and newlines, and more explicit definitions, so cross referencing IDs is not necessary). This is good for manually investigating the results.

The `-i` option will only indent the file (and place newlines), the `-hr` options implies this as well.

The `-c` option will generate a compact output, which aims to reduce duplicate data as much as possible. This is good for programmatically investigating the results.

Note that the options `-hr` and `-c` cannot be used together.

Given below is an example of how one could run the dependency checker:

```
1 java -jar dependency_checker.jar -s  
    ↪ ../graph-files/tajo_dependencies_by_structure101.csv -p  
    ↪ ../graph-files/tajo_dependencies_by_pyne.graphml -hr -d PACKAGE
```

# CHAPTER 3

## DESIGN

Here we will go into more detail about the design of Pyne itself. The project has been analysed using the existing documentation [1], the source code and the structural and dependency graphs created by Structure101 (figure 3.2).

### 3.1 REQUIREMENTS

According to the documentation [1] of the Pyne project, the project has the following requirements:

1. **The program should be able to create a graph from Java source files using git.**
2. The output graph should be the same as that of Arcan<sup>1</sup>, a Java software analysis tool.
3. The program should parse Java source files in such a way so it can find the different Java classes and packages.
4. The program needs to be able to use Git.
5. The program should provide a command line interface.

When inspecting the functionality of the program, one can see that all these requirements have been met. When inspecting the source code of the program and the documentation, one can even see how they implemented the requirements.

From the documentation it was made clear that most of the requirements were met using external dependencies. For example: Requirement 2 was met using Apache Tinkerpop<sup>2</sup> for graph creation, which is the same technology Arcan uses.

Requirement 3 was met using Spoon<sup>3</sup>, an extensive java source code parsing library. Lastly,

---

<sup>1</sup><https://essere.disco.unimib.it/wiki/arcan/>

<sup>2</sup><https://tinkerpop.apache.org/>

<sup>3</sup><http://spoon.gforge.inria.fr/>



Requirement 4 was met using Jgit<sup>4</sup>, a java git library.

In addition to that, using Structure101, we found that Requirement 5 is met using Apache Commons Cli<sup>5</sup> to provide a cli interface. Furthermore, Syncleus Ferma<sup>6</sup> was used as an extension to Tinkerpop. More information on this can be found in section 3.3.

## 3.2 IMPLEMENTATION

This application is divided into 3 different packages: pyne-demo, pyne-cli and pyne-api. The relations between these packages and their classes can be seen in figure 3.2. figure 3.1 represents the XS ("excess") diagram[3] of pyne. the XS diagram reflects the size of a method, package, class or other code item. It'll ensure that overly complex high level packages will have higher XS than a single method. This will reflect the likely negative impact on development. AS you can see from the diagram, there are two components used to measure the code items: fat and tangled. Fat measures the amount of basic complexity and tangled measures cyclic dependency. Tangled is mostly applied to higher level components like design and package.

Looking at figure 3.1, we can see that the fat of the methods in pyne is relatively small. This means that the methods in pyne are overall of a good size. However looking at the fat of the class and packages, we see that they are on the big side. To improve this, existing classes and packages could be split up into multiple classes/packages. looking at the design, we see that it is not really fat, but extremely tangled. This means that there are a lot of cyclic dependencies. Most of the time cyclic dependencies can have a negative impact on future development and maintenance of the application, so you would want to avoid them as much as possible.

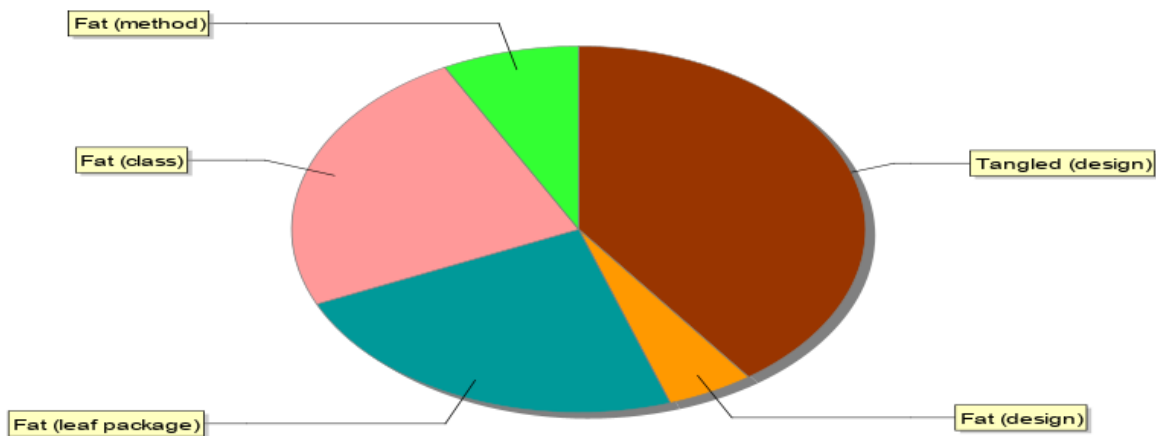


Figure 3.1: The XS diagram of the structure of Pyne, created by Structure101

<sup>4</sup><https://www.eclipse.org/jgit/>

<sup>5</sup><https://commons.apache.org/proper/commons-cli/>

<sup>6</sup><https://github.com/Syncleus/Ferma>

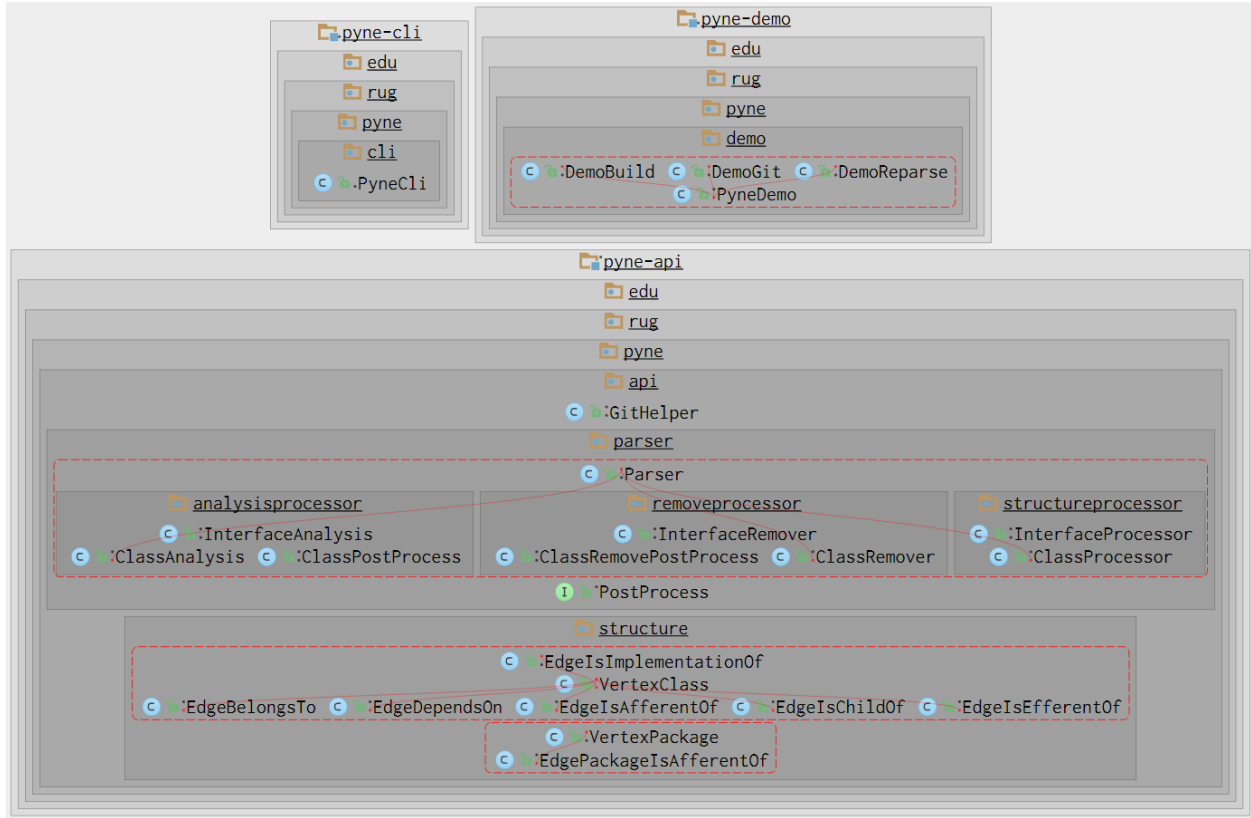


Figure 3.2: The full graph of the relations between the classes of Pyne, created by Structure101

Now that we have looked at the overall fat and tangles of pyne, we will briefly explain the 3 main packages within pyne. The pyne-api package will be explained in more detail, since this is the main component of the system.

### 3.2.1 PYNE-DEMO

This package acts as a testing ground for pyne. Most of the data in this package has been hard-coded to work only with a specific repository that was probably present on the creators local machine. Since we do not know which repository they used for testing and it is not present within the project, This package will not be able to work and we will consider this package deprecated.

### 3.2.2 PYNE-CLI

This package acts as a layer of communication between the user and the Pyne-api via a CLI (Command Line Interface). Given a git repository as a program argument, it will automatically start communicating with the pyne-api and create a dependency graph in `.graphml` format for each commit. A number of options can be given to the pyne-cli program. The first required argument is of course the repository itself, but one can also specify the

starting and end date to parse from, the interval between each commit and an input and output directory.

### 3.2.3 PYNE-API

This package is the heart of the program and is the one that is actually responsible for the dependency graph creation. This package is again subdivided into 2 packages: parser and structure, and contains one class: GitHelper, as can be seen in figure 3.3.

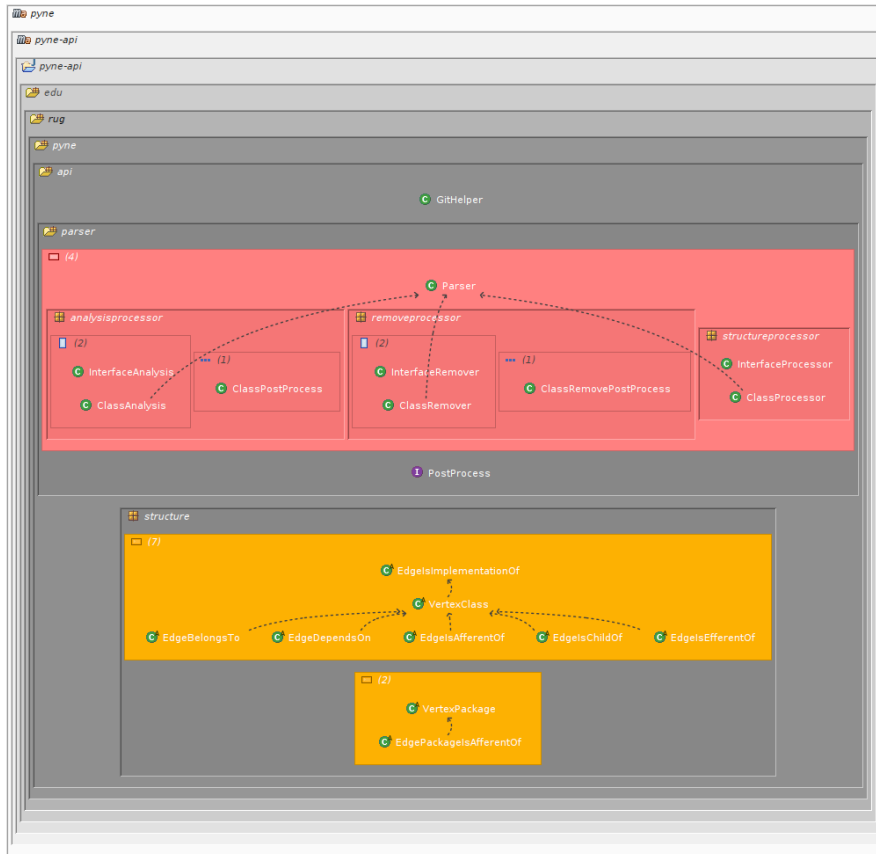


Figure 3.3: The graph of the relations between the classes of the pyne-api package created by Structure101

The GitHelper is a part of the API responsible for git related functionality like cloning a repository in a temporary location, parsing a commit or processing a commit. Most of the functionality of the GitHelper depends on the parser class to create the dependency graph.

The structure package does not contain any functionality, instead it provides abstract classes to build a dependency graph from. Special edges have been created for different kinds of relations: BelongsTo, DependsOn, AfferentOf, ChildOf, EfferentOf, ImplementationOf and

PackageisAfferentOF. Special vertexes have also been created to represent either a package or a class.

The parser package again consist of three sub-packages: analysisProcessor, removeProcessor, and structureProcessor. It also contains the Parser class and the PostProcess interface.

All communication with this package goes through the Parser class. This class is used to keep track of the files that have changed in the selected commit, and it holds lists of all types of processor classes, each of which are responsible for analysing the code in a different way.

The parser class will then in turn communicate with the analysis processor package, remove processor package and structure processor package to create a dependency graph. These packages provide the following functionality: The remove processor package will take care of removing nodes or edges from the graph if classes or relations have been removed or changed. The structure processor package takes care of adding nodes to the dependency graph and the analysis processor package analyses the different classes and packages and their relations to each other.

### 3.3 EXTERNAL DEPENDENCIES

The main external dependencies used by the project are:

- **Apache Tinkerpop**  
A flexible graphic framework. This library was chosen since Arcan, the software this project is suppose to mimic, also uses this framework.
- **Spoon**  
A java source code parsing library. This library was chosen so that a java parser did not need to be build from scratch[1].
- **Jgit**  
A java library that provides git functionality. This library was chosen so that the program is able to collect the source files using git.
- **Apache Commons CLI**  
A java library that provides CLI parsing and some basic CLI functionality. This library was used to not have to write a CLI parser from scratch.
- **log4j**  
A logging library. Used to log actions in the program.
- **Ferma**  
An extension to Apache Tinkerpop. Used to help building the graph.

Most of these dependencies were noted in the report [1], like Apache Tinkerpop, Spoon and Jgit. Others, like Apache Commons Cli, log4j and Ferma, were found using Structure101.

In figure 3.5 and figure 3.4 you will find the graphs depicting the relation between the packages of Pyne and the dependencies. Note that the used external dependencies have not been shown for pyne-demo. This is because pyne-demo does not provide any graph creation functionality, only a "demo" and it uses approximately 20 extra external dependencies that are not used in the rest of the project. It would therefore be unreadable and not really relevant to the workings of Pyne itself.

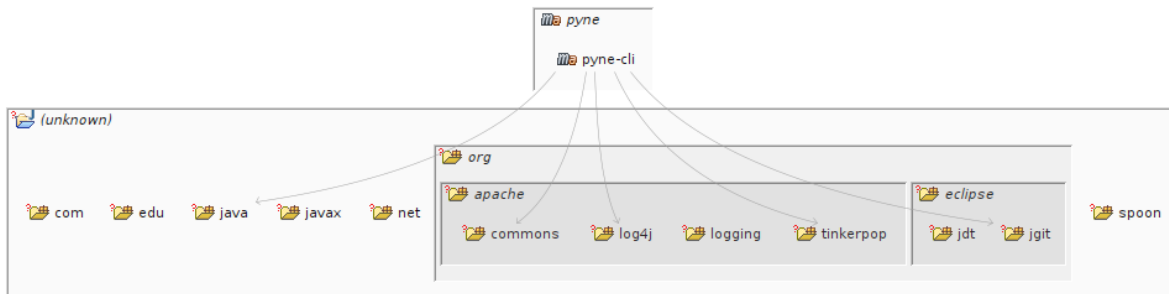


Figure 3.4: The graph of the external dependencies of pyne-cli, created by Structure101

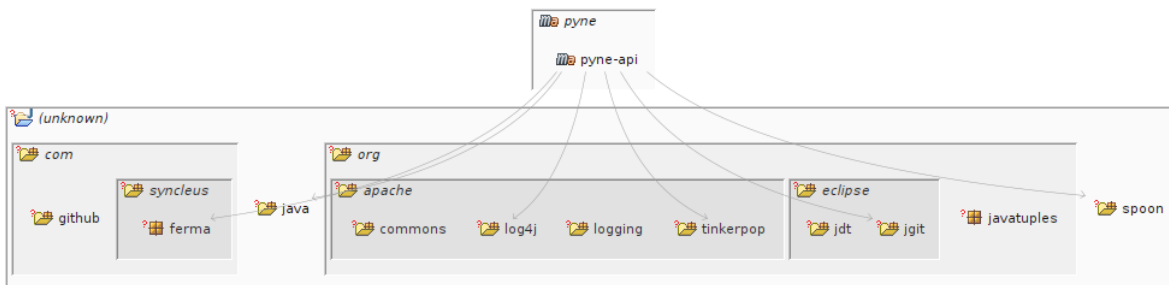


Figure 3.5: The graph of the external dependencies of pyne-api, created by Structure101

# CHAPTER 4

## PYNE VS STRUCTURE101

To compare Pyne to Structure101, we will test them both on a large existing system. For this we have selected Apache Tajo ([www.tajo.apache.org](http://www.tajo.apache.org)). The first section will describe the resulting graph for Tajo using Structure101. The following section will go into more detail about the design of Tajo using the graph generated by Structure101, then lastly we will run Pyne on Tajo and analyse the difference between Pyne and Structure101. Note that in the commands we show, we renamed the generated pyne-cli jar to simply `pyne-cli.jar` for brevity.

### 4.1 RUNNING STRUCTURE101 ON TAJO

Next we will run Structure101 on Apache Tajo. Since Structure101 uses a GUI instead of a CLI, we followed the steps in the GUI to create a Structure101 project. For ease of use, we have provided this Structure101 project in our GitHub repository under `structure101/tajo_structure101.java.hsp`. As mentioned before, the downside of Structure101 is the need to compile the target system before it will allow you to import it. This was difficult since we had first switched to java 11 to compile and run Pyne, but Tajo requires java 8. This was however not indicated clearly, so it took some time to figure out.

However, by being able to use a maven pom file, Structure101 does have the advantage of being able to find dependencies to external packages instead of only internal packages, like Pyne does.

Structure101 has more options for exporting as well, allowing users to export to `.png`, `.jpeg`, and "Graph as XML", which sounds a lot like `.graphml`, but they are not the same, and in fact the program we used to open the `.graphml` files did not display this file correctly.

Figure 4.1 shows what an exported image from Structure101 looks like.

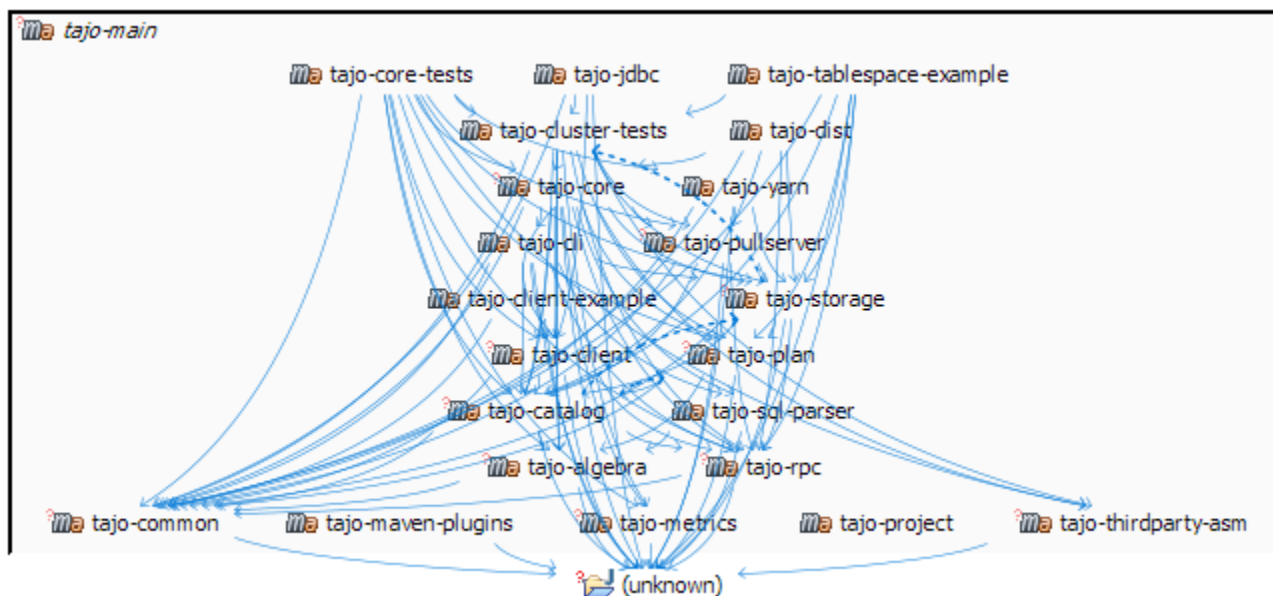


Figure 4.1: The graph of Tajo, created by Structure101

## 4.2 DESIGN OF TAJO

As you can see from figure 4.1, Tajo consists of several different packages working together. Firstly, We will briefly explain these components. Lastly we will go into more detail about the external dependencies tajo relies on.

### 4.2.1 COMPONENTS

Here we will list the components of Apache Tajo, describe what they do and how they interact with each other. Below you will find a list of all the main components according to a maven build file in the source directory and Structure101. The responsibilities of each component has been approximated using the build file and the source code.

- **tajo-core** is responsible for the main functionality of Tajo. This component uses a master worker pattern to perform its work. The master is responsible for query planning and coordinating the workers. It divides a query into sub-tasks and assigns these sub-tasks to individual workers. The workers are then responsible for executing these sub-tasks. The core structure can be seen in figure 4.2.
- **tajo-core-tests** and **tajo-cluster-tests** are responsible for testing the core and cluster respectively, where the focus is mainly on the correctness of query execution. Because they are testing libraries, they have only outgoing dependencies and are not really part of the core functionality. One exception is however the dependency of tajo-storage (more specifically the internal tajo-storage-kafka and tajo-storage-pgsql packages) on tajo-cluster-tests. This dependency was however caused by the fact that tajo-cluster-tests was added in as a dependency in the pom.xml file, while non of the

classes actually needed it. In other words, this dependency could've been removed from the pom.xml file.

- **tajo-plan** is responsible for the planning or scheduling within tajo. This package seems to be heavily used by the other components of the system. Tajo-core for example has around 4500 references to tajo-plan.
- In the dependency graph of 4.1 You can see the **tajo-project** component and that it has no dependencies. That is because it is not a java package. It is a parent POM for all Tajo Maven modules. All plugins and dependencies versions are defined here.
- **tajo-algebra** contains Algebraic expressions for database interaction. An example of a couple of expressions here are Join, DropTable, Sort and CreateDatabase. It only depends on the tajo-common package within the project and has several other components, among which the tajo-core component, depend on it.
- **tajo-catalog** contains a catalog of plugins for a server, client, drivers and common. According to Structure101, this package seems to be used by a lot of components, and seems to be depending on a few too. One could therefore say that it is an important component.
- **tajo-common** contains some common modules. This consists of some google protobuf[4] definitions and some helper functions. Since this is basically a util package, it does not depend on other components. Not surprisingly, all the other major components do seem to depend on it.
- **tajo-client** is responsible for the Client API and its implementation. The tajo-cli will communicate with the tajo-client, which will then in turn communicate with the server. Tajo-client has a few minor incoming and outgoing dependencies, but the most notable one is that of tajo-core with 460 references. It goes without saying that Tajo-client is an indispensable part of the system.
- **tajo-client-example** provides a Client API example and **tajo-tablespace-example** provides a table example. No components seem to depend on these parts. This proves that the do not provide any core functionality, they just act as examples.
- **tajo-cli** provides the command line interface for the user to communicate with tajo. It depends on a number of components in the system and even tajo-core seems to depend on it.
- **tajo-dist** This component is responsible for assembling the Tajo distribution. It depends on a few components, among which tajo-core, but no other component depends on it. Since this component should just assemble the Tajo distribution, no component should indeed depend on tajo-dist.
- **tajo-jdbc** represents the Tajo Java Database Connectivity driver. It is responsible for handling the communication with the database. Tajo-jdbc depends on tajo-common and tajo-client, but no other class seems to depend on it. It could be that this package is no longer used by the system.



- **tajo-maven-plugins** contains the maven plugins. From the dependency graph it can be seen that no component depends on tajo-maven-plugins. Therefore it is reasonable to assume that tajo-maven-plugins does not provide any functionality for tajo. It only contains some plugins
- **tajo-metrics** provides Metrics for quantitative assessment of tajo. just like tajo-algebra it just provides a few functions for tajo-core and the two test classes to apply a metric. Therefore only those classes depend on it and itself does not depend on anything else.
- **tajo-pullserver** is mainly responsible for fetching data from the server. The main components that use this service are tajo-core and tajo-yarn.
- **tajo-rpc** is responsible for the Remote procedure calls. This components seems to be used by most of the components in the system. It is therefore also a key component in the system.
- **tajo-sql-parser** is, as the name implies, responsible for the parsing of SQL commands. This component is only used by the cli and tajo-core, but still very important.
- **tajo-storage** is responsible for storage and the relevant plugins. This component seems to have a lot of dependencies both incoming and outgoing. Tajo-core even references it 530 times. It can therefore also be considered an important part of the system.
- **tajo-yarn** is presumable an extension of tajo in order to use it with yarn[5]. No main component seems to depend on yarn, so we can assume it is not an significantly important component.
- **thirdpart-asm** is a java code parser that provides functionality to parse compiled java code. It seems to be only used by tajo-core and the test components.

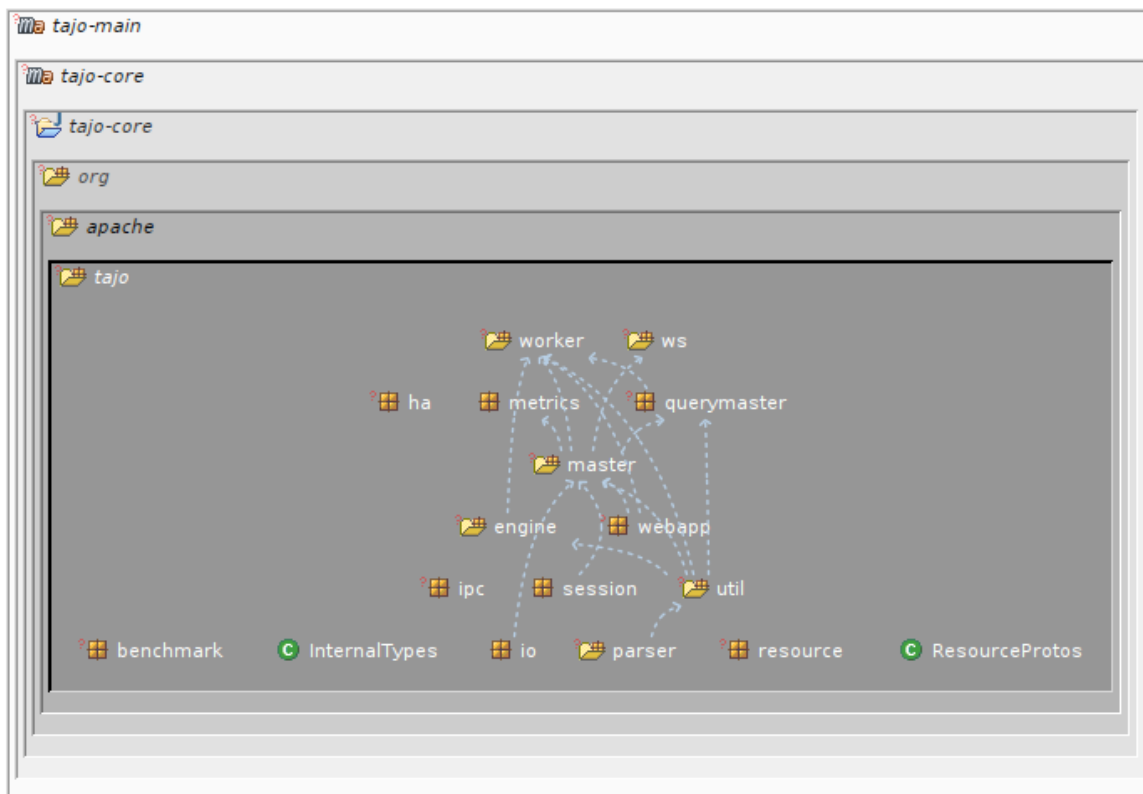


Figure 4.2: The graph of Tajo-core, created by Structure101

#### 4.2.2 EXTERNAL DEPENDENCIES

Here we will list a few of the main dependencies of tajo. This list will mainly include the external dependencies tajo-core uses, since that means those libraries are part of the core functionality of tajo. In addition to that, a few dependencies used by other components which can be considered important, will be listed as well.

One of the important external dependencies of Tajo is that to google protocol buffers[4]. These are Google’s language-neutral, platform-neutral, extensible mechanism for serializing structured data, much like JSON or XML, but smaller, faster and simpler. Tajo uses this to communicate with other services.

Another notable external dependency is that of amazonaws[6]. The features of this library are used by the tajo-storage component to enable authorisation or the enabling of retrieving credentials via the AWS service.

An important dependency used by tajo-core is codahale[7] and gmetric4j[8]. These two library provide metrics to give insight into what the code does in production and measure the behavior of critical components in the production environment.

Another dependency of tajo-core is minidev.json[9], which provides functionality of reading and parsing json-documents. The dependency jayway.jsonpath[10] is used in combination with minidev.json in a few minor cases to read the json documents.

Another remarkable dependency is the maxmind.geoip library[11]. This library is mainly used for identifying the location and other characteristics of an internet user. Looking at the source-code, it seems that tajo is only after the country-code and then only uses that information to find out the correct timezone you're in. It looks like tajo is not after your personalised location data for ad revenue.

two important dependencies for tajo-cli are the JLine library[12] and jansi[13], which provides functionality to handle and format console input.

antlr[14] is also used by both tajo-core and tajo-sql-parser to generate a parser for reading, processing, executing or translating structured text or binary files. Looking at the source code of tajo-core, it seems to be only used for SQL parsing.

A very important dependency of tajo-core is hadoop[15]. Not to mention that it is even required to have hadoop installed in order to use apache tajo. Hadoop provides a software framework for distributed storage and processing of big data.

Another important dependency used by tajo-core and tajo-rpc is glassfish.jersey[16]: An open source framework for developing RESTful Web Services in Java.

Snappy[17] is a dependency used by tajo-common for decompression/compression with the aim of highspeed and reasonable compression. This library is based of the original C++ version from google, but now written in Java. According to the git repository, it produces a byte-for-byte exact copy of the output created by the original C++ code.

Tajo-core uses mortbay jetty[18] as a java web server. More particularly Tajo-core uses this framework for their machine to machine communications.

Tajo-core also uses reflections[19]. Reflections scans your classpath, indexes the metadata, allows you to query it on runtime and may save and collect that information for many modules within your project.

another neat dependency is the SLF4j[20] or the simple logging facade for java. It allows the the end user to plug in the desired logging framework at deployment time.

## 4.3 RUNNING PYNE ON TAJO

We will start by running Pyne on Apache Tajo. The command we used to get a graph for the latest commit of Tajo is as follows:

```
1 java -jar pyne-cli.jar https://github.com/apache/tajo -s "2020-05-10" -e
   ↪ "2020-05-12" -p "DAY"
```

This creates a `.graphml` file<sup>1</sup>, but unfortunately this does have some issues. See Appendix A.

<sup>1</sup>available on our github repository, under `graphs/tajo_dependencies_pyne.graphml`

This is the only format Pyne can currently export, and opening a `.graphml` does require users to install a specialized program. We decided to use yEd Graph Editor<sup>2</sup> for this purpose. This has the advantage of being able to automatically rearrange the graph. Using this feature, we obtained figure 4.3.

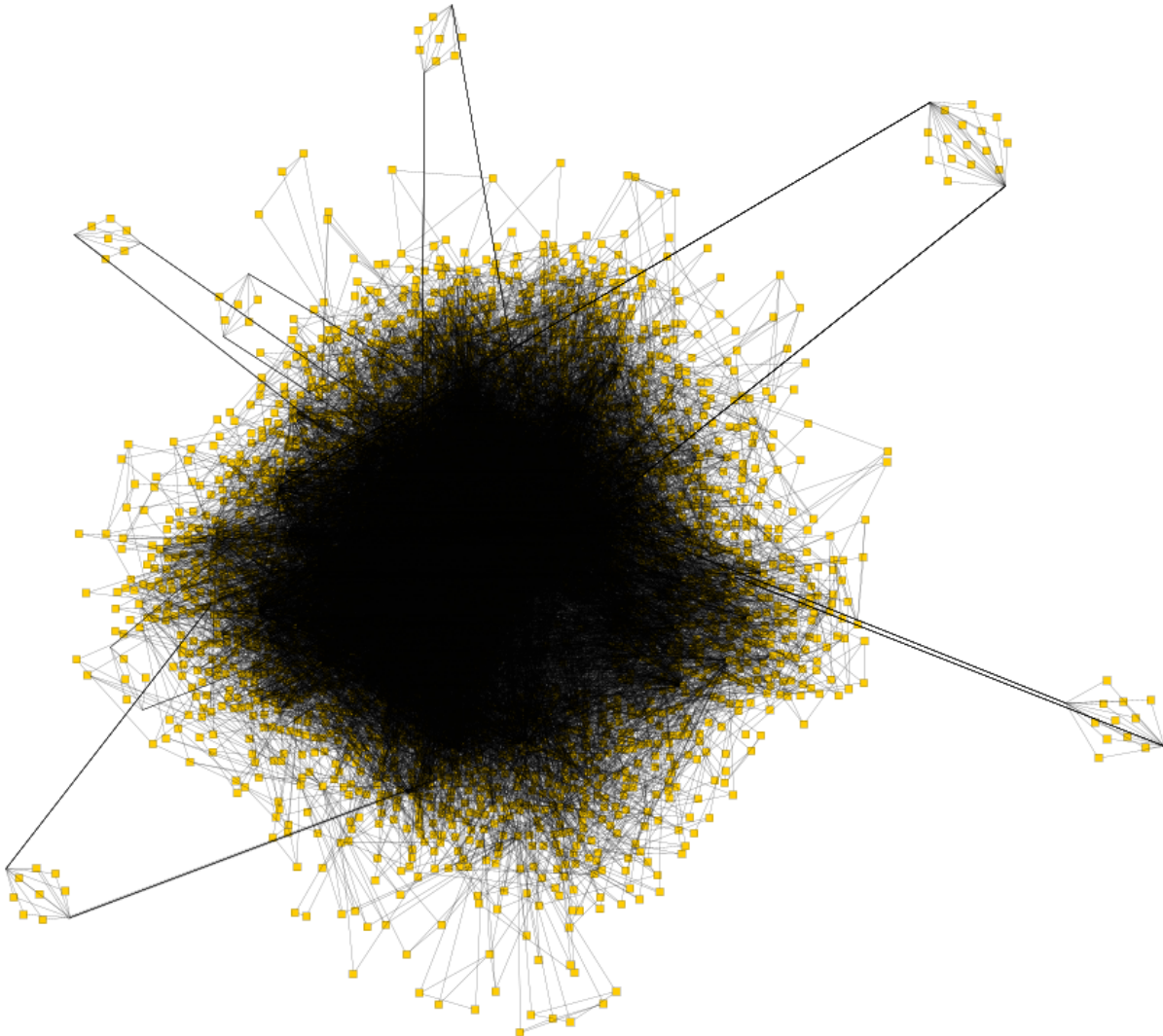


Figure 4.3: The automatically rearranged graph of Tajo, created by Pyne and rearranged by yEd

The nice thing about Pyne though, is that it can generate this graph without the need to compile the target system first, which in the case of Tajo, proved to be a difficult task on its own as we will see in the next section.

<sup>2</sup><https://www.yworks.com/products/yed>

## 4.4 COMPARING THE FEATURE SETS

Feature	Pyne	Structure101
GUI	NO	YES
IDE plugin support	NO	YES <sup>3</sup>
Needs compiled sources	NO	YES
Shows external dependencies	YES	YES
Shows fully qualified class names	YES	NO <sup>4</sup>
Open-source	YES	NO
Freely available	YES	NO
Multiple OS support	YES	YES
Allows direct analysis of remote repositories	YES	NO
.graphml export	YES	NO
.png export	NO	YES
.jpeg export	NO	YES
.csv export	NO	YES

Table 4.1: Comparing the features of Pyne and Structure101

## 4.5 COMPARING THE RESULTS OF PYNE TO STRUCTURE101

To compare the results of Pyne with Structure101, we needed to obtain a machine parsable format as the output of both tools. Pyne already does this, by exporting a `.graphml` file, however for Structure101 this was a little more difficult to find, since Structure101 is mainly focused on exporting images. However, Structure101 does have the option to export to `.csv`. After exploring the exported `.csv` from Structure101, we came to the conclusion that we could use this as our machine parsable format that we needed, though it is not an ideal file to work with.

Note that even though both of these programs are dependency graph creator engines, Pyne uses the source code to create the dependency graph, while Structure101 uses the `.jar` files to create the dependency graph. This will cause unavoidable differences in output for both programs, which are not necessarily wrong.

<sup>3</sup>though this is quite a bit more limited than using the full program

<sup>4</sup>for some reason, it only show the fully qualified names for packages

# CHAPTER 5

## COMPARING RESULTS PROGRAMMATICALLY

### 5.1 CREATING A DEPENDENCY CHECKER PROGRAM

Once we found a machine parsable format we could export from both tools, we wrote a Java program<sup>1</sup> which extracts the found dependencies from both of the generated files. It then checks if there any dependencies that were found by one tool, but not by the other, and provides some statistics on the results as well. A general flow diagram of this program can be found in Figure 5.1.

#### 5.1.1 DESIGN OF THE DEPENDENCY CHECKER

The dependency checker was made with the technologies from Pyne kept in mind. This means we used the same java version as Pyne does (version 11), and we export the results as an XML file. This way, not only does this mean that we don't have to awkwardly switch java version each time we want to run one or the other, but it also means we can easily migrate the code from being an external tool to being a feature in Pyne, if we ever want to do so. In the same spirit, the structure of the program was very much designed to be modular. This was achieved by creating what is known as a black box around the internal functionalities, and exposing a set of intuitive methods that execute all steps of the program.

A flow diagram of the dependency checker can be found in Figure 5.1.

---

<sup>1</sup>Available in our GitHub repository, under `dependency_checker/`

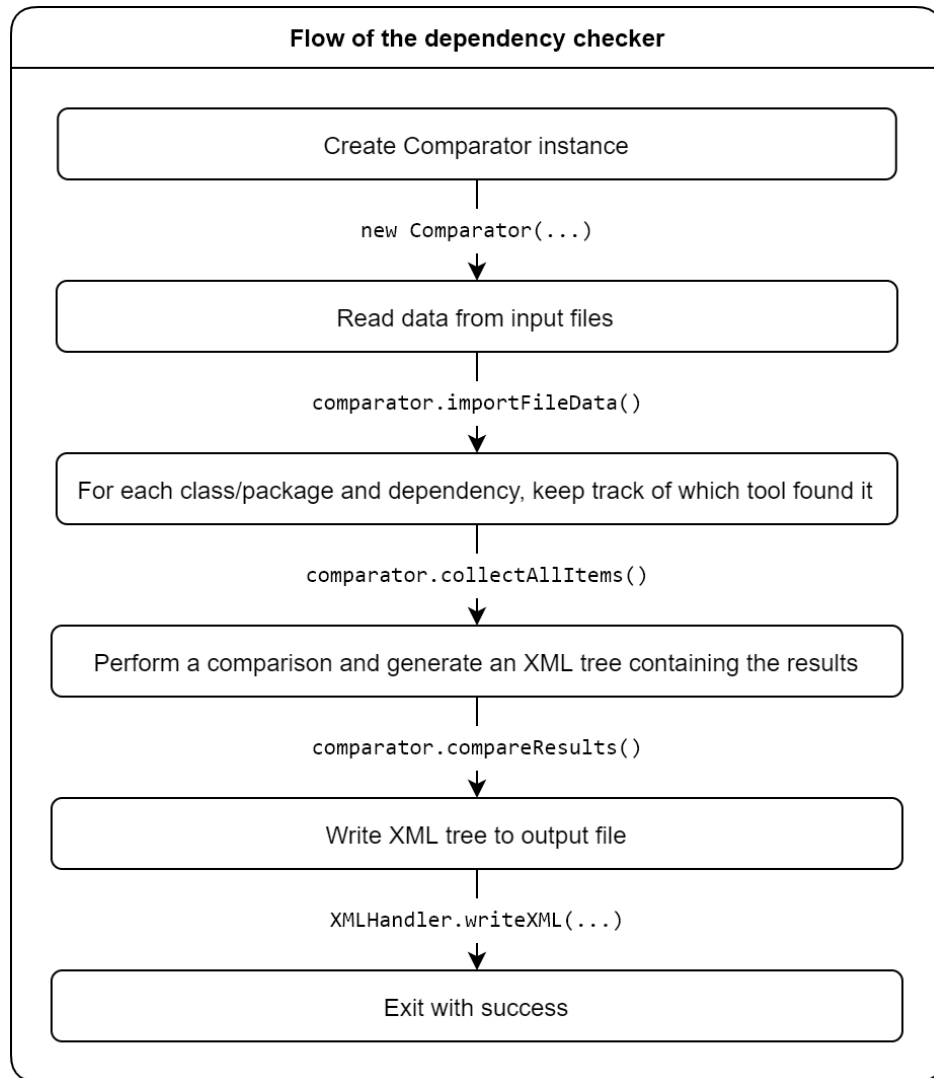


Figure 5.1: The flow diagram of the dependency checker program

The structure of the XML output of our dependency checker can be seen in Listing C.1, in appendix C.

The root element is called "results", which has 3 children: "allPackages"/"allClasses", "allDependencies" and "tools". The first contains the combined list of the packages or classes found by each individual tool. The second contains the combined list of the dependencies found by each individual tool. The last contains a list of all tools that we are comparing, where each tool contains:

- a list of packages or classes this tool did find
- a list of packages or classes this tool missed (by comparing the found packages or classes against all, mentioned earlier)
- a list of the dependencies this tool did find

- a list of dependencies this tool missed (by comparing the found dependencies against all dependencies, mentioned earlier)

Many elements also contain fields that provide some insights to the data without having to actually go through all items in the lists (which can get quite long for a big project, such as Apache Tajo). These fields include:

- "count", "countInternal", "countExternal", "countUnknown" tell the amount of items in a list, where "count" stands for the total amount of items, and the others tell the amount of items which are internal, external, or unknown, respectively.
- "internal" specifies whether a package is internal or not (i.e. external).
- "name" gives the name of the tool that the "tool" element represents.
- "percentageTotal", "percentageInternal", "percentageExternal", "percentageUnknown" specify the percentage that the amount of items in the list, which this field is on, is of all items of that category.

## 5.2 RUNNING OUR DEPENDENCY CHECKER PROGRAM

Using the fields on the XML elements of the output specified in subsection 5.1.1, we can obtain the statistics from comparing the unmodified version of Pyne to Structure101, using Apache Tajo as a benchmark again. These can be found in section C.2 in the Appendix.

From that output we can see that Pyne initially performed much worse than Structure101, only finding less than half of the package Structure101 did, and less than two thirds of the classes. We also see that Pyne found one package that Structure101 didn't find though. After investigating, it turns out that the dependency that Structure101 missed was `PrimitiveType`. Our first instinct was that Pyne could use this to label primitive types like `int` or `bool`, but this turned out not to be true. When we investigated further, we found that this package contains one class according to Pyne, namely `PrimitiveTypeNameConverter`. Searching for this term online returns a javadoc page from Apache Parquet<sup>2</sup>. The reason why it didn't get labelled by its full name (`org.apache.parquet.schema.PrimitiveType.PrimitiveTypeNameConverter`) is unknown.

Furthermore, Pyne missed 100% of the unknown classes. This makes sense, since Structure101 very often doesn't mention for classes whether they are internal or external (it only does this consistently for packages). Pyne on the other hand can always tell, so it makes sense it would not have any classes marked as unknown.

---

<sup>2</sup><https://www.javadoc.io/static/org.apache.parquet/parquet-column/1.8.2/org/apache/parquet/schema/PrimitiveType.html>



# CHAPTER 6

## IMPROVING PYNE

In this chapter we will go over the changes we made to Pyne in order for it to include all dependencies. We will be using the results from the dependency checker, which we discussed in chapter 5, and the source code of Pyne to find out what Pyne is missing and fix it. First of all, in section 6.1, we will list all the elements or "mistakes" in pyne that make it output different or less dependencies than structure101. Secondly, in section 6.2 we will list all these faults again and explain how we fixed them. Lastly, in section 6.3 we will compare the output of the dependency checker of the original version of pyne with the improved version of pyne (The version that has our implemented fixes). We will then evaluate to what extend our fixes have improved pyne's ability to find dependencies.

### 6.1 FAULTS IN PYNE

Here we will enumerate some of the "mistakes" in Pyne which makes its output different from Structure101. Note that the word mistake is a bit harsh here, some choices made in the implementation in Pyne are not necessarily wrong, but still give a different output. For example, Pyne does not include any folders with "test" or "example" in it, which is not necessarily wrong since those folders are not part of the main functionality of the program anyway. The first few points will discuss differences between Pyne and Structure101, which are not faults per se. The latter points will discuss the actual missing dependencies in Pyne and how it misses them. Most of these mistakes were made in the `ClassAnalysis` class, with most of them originating from the `getClassReferences()` function. The complete original code of pyne can be seen in the respective github repository[21]. We will explain how we fixed these issues in section 6.2.

#### 1. Files that are not in the source code are not parsed

Maybe an obvious fault of Pyne, if you can even call it a fault, but it should be mentioned anyway. Since Pyne uses the source files to create the dependency graph and Structure101 uses the .jar files, there are going to be unavoidable differences in the resulting dependency diagram. A lot of classes, mostly classes generated by Google Protobuf during compile time, are not in the source code yet and are only created

when the program is compiled. It also goes the other way around. Some packages, like tajo-docs and tajo-project are not included in the jar file, while they are in the source code.

## 2. Packages with no dependencies are skipped

This is a feature of Pyne explicitly implemented by the creator. Pyne has a whole post-process analysis where it traverses the graph and removes any nodes without edges. Structure101 does include packages without dependencies, so this feature will show a difference in output.

## 3. Pyne does not parse any folder with "test" or "example" in its name or packages that do not have the folder "src" or "main" in it.

Pyne explicitly filters out any folders without a folder named "src" or "main" in it. It also removes any folders that contain the word "test" and "example". It is not wrong to consider example and test files to not be part of the project and use the "src" or "main" keyword to find project folders, but it does result in a different output than Structure101.

## 4. Pyne does not check the constructor

Moving on to the actual faults in Pyne, or the missing dependencies, the first one we found was the missing check for the constructor in `getClassReferences()`. This function is responsible for returning the names (more specifically `CtTypeReferences`) of all classes this class depends on. It did go over each method in the class, but forgot about the constructor.

## 5. Pyne uses the return type of an invocation for its dependency, not the declaration for the method itself.

This was a mistake in the inner class `ExecutableConsumer`. The function `getClassReferences()` gave each invocation in the class, aka a method call like `dependencyClass.function()`, to the executable consumer which in turn should retrieve the name of the dependency. Instead of retrieving the class where the method was defined, it got the return type of the method. As can be seen in this code snippet:

```

1 public void accept(CtElement element) {
2     if (!(element instanceof CtExecutableReference <?>)) {
3         return;
4     }
5     CtExecutableReference executable = (CtExecutableReference)
    ↪ element;
6     try {
7         CtTypeReference executableType = executable.getType();
8         if (executableType != null &&
    ↪ executableType.getDeclaration() != null) {
9             dependencies.add(executableType.getTypeDeclaration());
10        }
11    } catch (NullPointerException e) {
12        // Ignore Spoon errors
13    }
14 }

```

If we ignore the empty catch statement made by the creator of Pyne (which is still painful to see) and focus our eyes on line 5 and 7, we can see where it goes wrong. On line 5, the element is cast to an executable, which is the invocation that is being checked. Then on line 7, the type from this invocation is retrieved with `getType()`. This will however not return the declaring type, but the return type of the invocation. The function that should've been used is `getDeclaringType()`.

6. **Pyne had to retrieve the type declaration (so the entire class + body) in order to make the dependency, but then later needed only the referenced-type (aka the full name).**

This caused Spoon a lot of trouble when looking up external libraries. Getting the declaration was not necessary, only the name would suffice. So Spoon had to go back and forth for finding the reference name. For an example, see the code snippet below:

```

1 private void processClassReferences(CtType clazz, VertexClass
  ↪ vertexClass) {
2     for (CtType referencedClass : getClassReferences(clazz)) {
3         if (referencedClass == null ||
  ↪ referencedClass.getReference() == null) {
4             continue;
5         }
6         VertexClass referencedClassVertex
7             =
  ↪ getOrCreateVertexClass(referencedClass.getReference());
8             vertexClass.addDependOnClass(referencedClassVertex);
9     }
10 }

1 private List<CtType> getClassReferences(CtType clazz) {...}

```

The `processClassReferences()` function takes a list of types (which is a definition for a class or interface). This list will then contain the class definition, methods and fields for each class this specific class depends on. But when creating the vertex on line 8, it only needs the reference, which is essentially just the package+name of the class. Since the type is retrieved from the reference in `getClassReferences()`, this step is totally unnecessary. `getClassReference` should return a `CtTypeReference`, not a `CtType`.

7. **Pyne did not look at the type of the parameters and return value of each method in the class.**

There is no retrieval of parameters or return value for each method in `getClassReferences()`. Class A does depend on Class B if one of it's methods uses Class B as a parameter or return value, so it should be considered.

8. **Pyne did not look at the type of fields in each class.**

Pyne does check the annotations of the fields, but not the type of the field. So fields in a class are completely skipped.

9. **Spoon has some issues getting declarations of nested executable references.**

This is actually not an issue with Pyne, but an issue with Spoon. To elaborate more

on what we mean with "nested executable references", it is something like this: "dependencyClass.function().anotherFuction()". Spoon has trouble finding out what the Class is that defines anotherFuction(). This is especially a problem for external packages, where Spoon does not have the source code for either the class that defines function() or the class that defines anotherFuction(). Since this is a Spoon issue and not an issue with Pyne, we would need to replace the entire parser and rewrite Pyne almost entirely from scratch.

## 6.2 CHANGES TO PYNE

Here we will go over each fault in Pyne and describe how we chose to fix it. Almost all of these fixes were performed in the `ClassAnalysis` class, of which the new source code can be found in our GitHub repository [22].

### 1. Files that are not in the source code are not parsed

For obvious reasons, this issue could not be fixed. However to filter out all of the files that Structure101 finds but Pyne does not, we manually inspect all the missing internal dependencies of Pyne. From the resulting output we see that the only missing packages are not there because they are either not in the source code or do not have dependencies at all.

### 2. Packages with no dependencies are skipped

Since this is also a core feature of Pyne, this issue could not be fixed. However we applied the same technique as before to inspect the missing internal dependencies of Pyne. As mentioned, From the resulting output we see that the only missing packages are not there because they are either not in the source code or do not have dependencies at all.

### 3. Pyne does not parse any folder with "test" or "example" in its name or packages that do not have the folder "src" or "main" in it.

We choose to include test and example folders in our project as well, since Structure101 does this too. This fix was as easy as removing the filters from the parser class. We did choose to keep the filter that specified that the folder must contain an "main" or "src" folder since Spoon would return errors if we removed this. The obvious problem is that Spoon cannot parse a directory that is not a source code directory.

### 4. Pyne does not check the constructor

As you can see from the original code from Pyne [21], Pyne previously used this line of code to retrieve all of its methods.

```
1 for (CtMethod<?> ctMethod : (Set<CtMethod<?>>) clazz.getMethods())
```

We added the constructor to this for loop with the following code:

```
1 ArrayList<CtExecutable<?>> executables = new ArrayList<>();
2 executables.addAll((Set<CtExecutable<?>>) clazz.getMethods());
3 if (clazz instanceof CtClass) {
```

```

4     executables.addAll((Set<CtExecutable<?>>) ((CtClass)
    ↪ clazz).getConstructors());
5 }
6 for (CtExecutable<?> ctExecutable : executables)

```

Essentially what we do is we add the list of methods and the list of constructors together and then loop over that instead. No real change needed to be made to the body of the for loop, since most of the functions used there came from `CtExecutable` anyway, not from the child class `CtMethod`.

#### 5. Pyne uses the return type of an invocation for its dependency, not the declaration for the method itself.

On first sight, this looked like an easy fix. Just replace `getType()` with `getDeclaringType()`. However the `executableConsumer` also handled constructor calls (like `new dependencyClass()`) for which `getType()` was the correct method to use. That is why we decided to remove the `executableConsumer` class and let `getClassReferences()` handle finding the references for constructor calls and invocations separately using the code below:

```

1 // Get all constructors in the method
2 List<CtConstructorCall<?>> constructorElements = body
3     .getElements(new TypeFilter<>(CtConstructorCall.class));
4
5 //add all references for the constructor calls in the method
6 for(CtConstructorCall<?> c : constructorElements){
7     references.add(c.getType());
8 }

```

```

1 // Get all invocations in the method
2 List<CtInvocation<?>> invocationElements = body
3     .getElements(new TypeFilter<>(CtInvocation.class));
4
5 // Retrieve the dependencies of all invocations
6 for(CtInvocation<?> c : invocationElements){
7     if(c.getExecutable().getDeclaringType() == null){
8         LOGGER.warn("Spoon cannot find the declaration of " + c);
9     }else {
10
11     ↪ if(!(c.getExecutable().getDeclaringType().getTypeDeclaration()
12     ↪ == null))
13         references.add(c.getExecutable().getDeclaringType());
14     else
15         LOGGER.warn("Spoon cannot find the declaring type of " +
16     ↪ c);
17     }
18 }

```

#### 6. Pyne had to retrieve the type declaration (so the entire class + body) in order to make the dependency, but then later needed only the referenced-type (aka the full name).

This was a simple fix. The only thing that needed to be done was rewrite statements like `c.getType().getTypeDeclaration()` to `c.getType()` and change the return value of `getClassReferences()`.

#### 7. Pyne did not look at the type of the parameters and return value of each method in the class.

To fix this a few extra checks needed to be added. This can be seen in the code below. This code snippet added the type of each method to the list of references.

```
1 references.add(ctExecutable.getType());
```

And this code snippet added all the annotations of the methods parameters and parameters types to the list of references.

```
1 for (CtParameter<?> parameter : ctExecutable.getParameters()) {
2     parameter.getAnnotations().forEach(annotationConsumer);
3     references.add(parameter.getType());
4     tempCheck(parameter.getType(), clazz, parameter.toString());
5 }
```

#### 8. Pyne did not look at the type of fields in each class.

To fix this a few extra checks needed to be added. This can be seen in the code below. This code snippet added all the annotations of the fields and field types to the list of references.

```
1 for (CtField<?> field : (List<CtField<?>>) clazz.getFields()) {
2     field.getAnnotations().forEach(annotationConsumer);
3     references.add(field.getType());
4     tempCheck(field.getType(), clazz, field.toString());
5 }
```

#### 9. Spoon has some issues getting declarations of nested executable references.

This unfortunately could not be fixed. This is an issue that does not lie in Pyne, but in Spoon instead. The only way to fix this is use an entirely new parser and rewrite Pyne entirely. Even then it would not ensure that this problem will be fixed.

## 6.3 RESULT

### 6.3.1 APPROACH

Here we will compare the original version of Pyne to our improved version with the changes mentioned in section 6.2. Before we compared the two versions we first took the following steps:

1. we ran the original version of pyne on Apache Tajo[2].
2. We compared the resulting graphml file to the csv file from structure 101 using the dependency checker. The resulting output can be found in section C.2 in the Appendix.

3. We ran the improved version of pyne on Apache Tajo[2].
4. We compared the resulting graphml file to the csv file from structure 101 using the dependency checker. The resulting output can be found in section C.3 in the Appendix.
5. We put both outputs of the dependency checker in tables 6.1 and 6.2.
6. We calculated the Recall, Precision and F-measure for each table using the following formulas:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F = \frac{2 * Precision * Recall}{Precision + Recall}$$

where

$TP$  = (found package dependencies by pyne—missed package dependencies by structure101)

$FP$  = missed package dependencies by structure101

$FN$  = missed package dependencies pyne

Using both table 6.1 and 6.2 we will try to compare the original version of Pyne to our improved version. This table contains information on the found and missed packages, found and missed package dependencies, found and missed classes and found and missed class dependencies of each program (pyne and structure101). In addition to that the Precision, recall and F-measure are given. These measures were calculated using the missed and found package dependencies field.

The "found -" rows in the table represent the code elements (aka packages, classes and dependencies) this particular program found. Note that in pyne's case, this does not mean that these are the true code elements. Pyne might misidentify or misname a code elements, resulting in a False positive. The "missed -" rows in the table represent the code elements this particular program was not able to find, but the other was. So in the case of pyne, these are all the true code elements pyne did not find. In the case of structure101, these are all the code elements pyne mislabeled. So the correctly identified code elements can be found by subtracting the missed packages/classes/dependencies of structure101 from the found packages/classes/dependencies of pyne.

Note that: the found/missed class (dependencies) were not used for the calculation of the precision, recall and f-measure. This is because structure101 does not provide a package name to the class when extracting the graph to a csv file, making the comparison to pyne rather difficult. This data can be used as an indication, but the package data is more reliable.

### 6.3.2 ANALYSES

Here we will analyse the table and actually compare the original version of Pyne to the improved version, while using the output of structure101 as ground truth. We will first note the difference in results of both table 6.1 and 6.2 and explain them. Secondly, we will compare the precision recall and f-measure. Lastly, we will put everything together and conclude the analysis.

	pyne-original	structure101
Found packages	193	429
Missed packages	237	1
Found package dependencies	1020	3327
Missed package dependencies	2538	231
Found classes	2348	3259
Missed classes	1622	711
Found class dependencies	7809	36319
Missed class dependencies	33457	4947
Precision	77.35%	
Recall	23.72%	
F-measure	36.30%	

Table 6.1: Results of the dependency checker for the original version of Pyne



	pyne-improved	structure101
Found packages	393	429
Missed packages	104	68
Found package dependencies	2431	3327
Missed package dependencies	1560	664
Found classes	3837	3259
Missed classes	957	1535
Found class dependencies	28889	36319
Missed class dependencies	29667	22237
Precision	72.69%	
Recall	53.11%	
F-measure	61.38%	

Table 6.2: Results of the dependency checker for the improved version of Pyne

The first thing you might be able to notice is the increase in missed packages, classes and dependencies of structure101 in table 6.2. This means that the improved version of Pyne finds packages, classes and dependencies that structure101 did not find. This issue can be explained however. Spoon sometimes has trouble retrieving the complete package name of (mostly external) dependencies, which makes it give only a subset of the package name. For example if the actual package name was `package.external.project.classFile`, spoon might only give `project.classFile`. In addition to that, most of the extra packages are protoclasses generated by google protobuf[4], which are generated on compile time and can produce issues with Spoon. However, pyne is able to find these dependencies, but because the dependency checker only sees an incomplete package name, it'll mark it as missing. So the number of missed package dependencies by structure101 is in reality much lower.

Looking at the missed package dependencies of pyne, we see that there is an overall decrease in the improved version. This is a very good thing, since this means that the improved version of pyne finds more of the classes found by structure101 than the original version. This alone proves that our changes improved Pyne drastically.

The remaining missed dependencies, after manual inspection of around 30 of them, can be explained as follows. All of the inspected internal missing dependencies and a large portion of the external ones are caused by the fact that the classes are not in the source-code, but are present in the jar file. Which makes it virtually impossible for spoon to find them. A lot of these classes are Google Protobuf[4] classes, which are presumably created on compile time.

One last issue explaining the rest of the inspected missing dependencies, is the fact that Spoon has trouble finding the declaration of nested dependencies. Since this is a Spoon specific issue, nothing can be changed to Pyne to fix it. Using another parser might be an option, but a full rewrite of Pyne will be needed and using a different parser might introduce new and different problems.

Looking at both table 6.1 and 6.2. We see that total amount of correctly identified packages

by Pyne, calculated by subtracting the missed packages of structure101 from the found packages of Pyne, increased from 192 to 325. This is an increase of 69.27%.

Furthermore, the amount of correctly identified package dependencies by Pyne, calculated by subtracting the missed package dependencies of structure101 from the found package dependencies of Pyne, increased from 789 to 1767. This is more than double the original dependencies.

Finally, we will be looking at the precision, recall and the f-measure of both the improved and original version of pyne. Looking at table 6.1 and 6.2, the precision goes from **77.35%** to **72.69%**, which is largely due to the fact that Spoon is not always able to find complete package names for external packages.

The Recall however has drastically increased with the improvements, from **23.72%** to **53.11%**. Which is due to the decrease in missed package dependencies.

Finally, looking at the F-measure we see an increase from **36.30%** to **61.38%**. With these increase in scores we can conclude that our improvements did increase the number of dependencies found by Pyne.

# CHAPTER 7

## DISCUSSION

While we were successful in our pursuit to improve the performance of Pyne, it still doesn't perform quite as well as Structure101 does, as we saw in section 4.5. This is where we already ran into the first issue of our project, namely that while we were able to find the differences in found packages, classes, and dependencies between those, it turned out to be quite difficult to diagnose the reason for these differences.

This is partially due to Structure101 giving very limited information in the files it exports. This does make some sense because it is not the primary use case of Structure101, however it is a shame that this feature was implemented in such a basic manner while the U.I. of Structure101 seems much more feature rich. Still, in the end this only means that Structure101 might not have been the ideal tool to use for our use case.

Another issue we had while trying to diagnose the reason for the differences between Structure101 and Pyne is that Tajo is a very large project. This means that Tajo consists of lots of classes, packages, and dependencies. This made it difficult for us to verify the results of our dependency checker program during its development, and to manually find any sort of pattern in the items missing from the results of one tool or the other. It was therefore easier to find issues with Pyne in the source code of Pyne itself rather than its output.

Some possible ways to mitigate these issues would be to use a smaller benchmark project instead of Tajo, maybe even one that is specifically meant to benchmark dependency graph creator engine programs, however finding one such project might be difficult in itself. Either way, a different benchmark project might have made it easier for us to manually identify patterns in the results.

Finally, there are some issues with our dependency checker program. Like mentioned in section 6.3, the program currently has issues recognizing two packages as the same when one of the two doesn't specify the full name of the package, but only a subset. For example, consider a package `a.b.c.d.e`. If one tool lists it by its full name, `a.b.c.d.e`, and another lists it as `c.d.e`, they will currently not be detected as the same package. This could be fixed, but it may lead to false positives, so special care would have to be taken to ensure this detection would be as accurate as possible.

## CHAPTER 8

# CONCLUSION AND FUTURE WORK

Looking at section 6.3 we saw an increase in found packages by the improved version of Pyne of nearly 70% compared to the original version. The amount of package dependencies even doubled compared to the original version.

In addition to that, looking at the calculated precision, recall and F-measure from section 6.3, we saw very promising improvements as well when comparing the improved version with the original version of Pyne. Although the precision slightly decreased from 77% to roughly 73%, recall increased from 24% to 53%. Finally, the F-measure shows the most promise with an increase from 36% to 61%. So with these results we concluded that our changes to Pyne indeed improved the program quite a bit.

There is still room for improvement however. For instance, Spoon seems to have a lot of trouble with finding the dependencies of nested invocations, as mentioned in section 6.3. Future work could include changing the Java source code parser used in Pyne from Spoon to something more mainstream like the Java parser library[23].

Another thing that might be a point of possible future improvement is method parsing. As already mentioned by the creators of Pyne [1], there is currently no functionality in Pyne that can add method to method dependencies to the graph. Since Structure101 does have this feature, it could be considered as a valuable option for Pyne as well.

As we already mentioned in Appendix A, the default CLI values of Pyne are quite confusing and illogical. The problem with this is that if the GitHub repository Pyne is analysing does not have any commits between now and 5 days ago, no graph is created and no error is thrown. A better solution for in the future would be to set the default period to the latest commit, and only analyse that. This would solve a lot of confusion for first time users.

To be able to find other problems in Pyne, future research could include looking at other programs than Tajo. Tajo was quite a big program and we might have missed some "missing" dependencies in Pyne. Applying Pyne to something smaller like Apache Tika [24], might show new missing dependencies or point to other existing problems withing Pyne.

# REFERENCES

- [1] Patrick Beuks. *Building a dependency graph from Java source code files*. 2019.
- [2] Apache Software Foundation. *Tajo*. Version 0.12.0-SNAPSHOT. Dec. 14, 2020. URL: <https://tajo.apache.org/>.
- [3] HeadwaySoftware. *XS – A Measure of Structural Over-Complexity*. 2006. URL: <https://structure101.com/static-content/pages/resources/documents/XS-MeasurementFramework.pdf>.
- [4] Google. *com.google.protobuf*. Dec. 14, 2020. URL: <https://developers.google.com/protocol-buffers>.
- [5] *yarn*. Dec. 14, 2020. URL: <https://yarnpkg.com/>.
- [6] Amazon. *com.amazonaws*. Dec. 14, 2020. URL: <https://aws.amazon.com/>.
- [7] Dropwizard. *io.dropwizard.metrics*. Dec. 14, 2020. URL: <https://metrics.dropwizard.io>.
- [8] ganglia. *info.ganglia.gmetric4j*. Version 1.0.3. Dec. 14, 2020. URL: <https://github.com/ganglia/gmetric4j>.
- [9] Minidev. *net.minidev*. Dec. 14, 2020. URL: <https://mvnrepository.com/artifact/net.minidev/json-smart>.
- [10] Jayway. *com.jayway.jsonpath*. Dec. 14, 2020. URL: <https://github.com/json-path/JsonPath>.
- [11] Maxmind. *com.maxmind.geoip*. Version 1.2.15. Dec. 14, 2020. URL: <https://www.maxmind.com/en/geoip2-services-and-databases>.
- [12] Maxmind. *jline*. Dec. 14, 2020. URL: <https://github.com/jline/jline3>.
- [13] Fusesource. *org.fusesource.leveldbjni*. Version 1.8. Dec. 14, 2020. URL: <https://github.com/fusesource/jansi>.
- [14] Terence Parr. *antlr*. Dec. 14, 2020. URL: <https://www.antlr.org/>.
- [15] Apache Software Foundation. *Hadoop*. Version 2.3.0+. Dec. 14, 2020. URL: <https://hadoop.apache.org>.
- [16] glassfish. *glassfish.jersey*. Dec. 14, 2020. URL: <https://eclipse-ee4j.github.io/jersey/>.
- [17] dain. *org.iq80.snappy*. Dec. 14, 2020. URL: <https://github.com/dain/snappy>.
- [18] Mortbay. *org.mortbay.jetty*. Dec. 14, 2020. URL: <https://mvnrepository.com/artifact/org.mortbay.jetty>.
- [19] ronmamo. *org.reflections*. Dec. 14, 2020. URL: <https://github.com/ronmamo/reflections>.
- [20] qos-ch. *slf4j*. Dec. 14, 2020. URL: <http://www.slf4j.org/>.

- [21] darius-sas. *Pyne*. Version 1.1-SNAPSHOT. Jan. 9, 2021. URL: <https://github.com/darius-sas/pyne>.
- [22] Job Heersink, Richard Westerhof. *Our Github Repository*. Jan. 9, 2021. URL: <https://github.com/richardswesterhof/pyne>.
- [23] Danny van Bruggen. *javaparser*. Version 3.18.0. Jan. 10, 2021. URL: <https://javaparser.org/>.
- [24] Apache. *Tika*. Version 1.24.1. Jan. 10, 2021. URL: <https://tika.apache.org/>.

# APPENDIX A

## TROUBLESHOOTING

Here we present several problems encountered and how we tried to solve them.

1. **Bug in graph creator:**

When I open a graph in a graph editor like draw.io or yEd graph editor, only one square appears. But the file is 2.25mb. Turns out that all the squares are stacked upon each-other, but the relation between the squares is still there (which is visible from the neighbours tab in yEd graph editor).

It seems that the squares itself do not contain any data on the classes they represent. They do not have any value assigned to them. Furthermore, opening any graph will give a warning about the 'linesOfCode' property being of type long, which apparently isn't supported. I sincerely hope that an integer would be enough to count the amount of lines of code anyway, so we'll change that.

2. **Default cli values:**

The default options for cli are not the best choice. Namely the period option is set to days by default, which will create a separate graph for each day. The start date is set to 5 periods (so 5 days on default) from the end date (which is the current date by default). It is relatively rare for a git repository to have commits from the past 5 days. We shall therefore change the default values.

3. **Error handling:**

In a few cases, where something goes wrong or no graph is created, no error is thrown and the user does not know why no graph has been made. For that reason more error logging statements need to be added.

4. **Java Version:**

Pyne requires java version 11, however to compile Tajo version 1.8 is needed. This leads to confusing errors if you try to compile Tajo for yourself after upgrading to java 11 to compile Pyne.

5. **First manual inspection:**

On the first manual inspection of the dependencies pyne is not able to see using

the source code and our created dependency checker, the following mistakes were recognized:

Pyne is not able to inspect the constructor. (adding this fixed two packages)

Pyne is not able to inspect arguments of methods and constructors.

Pyne is not able to inspect return values of methods.

#### 6. scanned directories:

After manually analysing the directories from the git repository added to pyne by the githelper, we already found some inconsistencies. `tajo-core-test`, `Tajo-tablespace-example`, `tajo-cluster-tests`, `tajo-client-example`, `tajo-project` are all missing from the list of added directories. In addition to that `tajo-docs` is present in pyne, while not present in `structure101`. The latter is probably due to the fact that `tajo-docs` is not present in the jar-file, but only in the git repository. Looking further into the source code of pyne, it seems that in the `findSourceDirectories()` method, pyne specifically filters out any directories with "test" or "example" in it and any directory that does not contain a `src/java` or `src/main` in their directory. Since there is no real reason why tests should not be tested on dependencies, this filter should be removed.

#### 7. missing packages in graph:

The missing packages in the graph are probably mainly due to the fact that pyne filters out packages without incoming or outgoing dependencies.

#### 8. inconsistencies structure101 and source code:

Since `structure101` analyses the .jar files and pyne the source files, there are going to be some obvious inconsistencies. When comparing the source code structure with the `structure101` diagram, you might notice that a few minor packages within the components are missing from the source code or missing from the diagram. We will have to keep this in mind.

#### 9. Compiling Tajo:

To let `Structure101` analyse Tajo, we had to compile it first. While the requirements for compiling Tajo<sup>1</sup> say that java 8 or above can be used, we weren't able to compile it with java 11, so we used java 8 instead. On top of that, it seems like there is one failing test case during the build stage, which prevents the build from finishing. To get around this, we compiled with the `-DskipTests` flag in Maven.

---

<sup>1</sup><https://github.com/apache/tajo#requirements>



# APPENDIX B

## CHANGE LOG

ID	Author	Student Number	Email Address
JH	Job Heersink	s3364321	j.g.heersink@student.rug.nl
RW	Richard Westerhof	s3479692	r.s.westerhof.2@student.rug.nl

Table B.1: The authors that contributed to this document

Ver.	ID	Date	Revision
0.1	JH	21-11-2020	Create User guide.
	JH	21-11-2020	Create Design section.
	JH	21-11-2020	Create problem (troubleshooting) section.
	RW	22-11-2020	Added classes and dependency graph to design
	RW	22-11-2020	Revised User guide
	RW	22-11-2020	Revised problem section

Ver.	ID	Date	Revision
0.2	JH	01-12-2020	Reformatted the document.
	JH	01-12-2020	Added frontpage.
	JH	01-12-2020	Added introduction.
	JH	01-12-2020	Rewrote design section.
	JH	01-12-2020	Moved and renamed troubleshooting section to appendix.
	RW	01-12-2020	Create chapter 4.
	RW	01-12-2020	Add "java version" item to Troubleshooting appendix.

Ver.	ID	Date	Revision
0.3	RW	03-12-2020	Incorporate TA feedback in document.
	JH	04-12-2020	Added external dependencies section.
	JH	05-12-2020	Added requirements section.
	RW	07-12-2020	Improve internal structure of document and improve consistency.
	RW	08-12-2020	Add XS diagram and graphs of Tajo created by both Structure101 and Pyne.

Ver.	ID	Date	Revision
0.4	JH	12-12-2020	Incorporate TA feedback in document.
	JH	12-12-2020	Add design of Tajo.
	JH	13-12-2020	Add external dependencies of Tajo.
	JH	14-12-2020	Added all software to reference list.
	RW	14-12-2020	Write code for automatic comparison, perform analysis using results from this comparison.
	RW	14-12-2020	Write section 4.5 (reporting on code and analysis).

Ver.	ID	Date	Revision
0.5	RW	21-12-20	Expanding section Creating a dependency checker program.
	RW	24-12-20	Write more details in Creating a dependency checker program.
	JH	29-12-20	Created chapter 6 and added faults of pyne.
	JH	30-12-20	Added improvements made to Pyne in chapter 6
	RW	02-01-21	Moved dependency checker sections to new chapter.
	RW	02-01-21	Added flow diagram of dependency checker.
	JH	02-01-21	Added comparison between original and improved Pyne in chapter 6.
	JH	03-01-21	Added explanation to figure 3.1.
	JH	03-01-21	Added section on running the dependency checker.
	RW	04-01-21	Added Tajo compilation to issues section.
	RW	04-01-21	Wrote user guide section for Structure101.
	RW	04-01-21	Expanded user guide section for dependency checker.
	JH	05-01-21	Updated introduction section.
	RW	05-01-21	Fixed spelling in the document.
	RW	05-01-21	Added results for class level analyses.

Ver.	ID	Date	Revision
0.6	RW	08-01-2021	Created discussion and conclusion chapters.
	JH	09-01-2021	Updated chapter 6.
	JH	09-01-2021	Updated Appendix C.
	JH	10-01-2021	Updated Chapter 8.
	RW	12-01-2021	Some spelling and grammar fixes.

# APPENDIX C

## DEPENDENCY CHECKER OUTPUT

This chapter contains all the output of the dependency checker program. Most of these outputs, except for the template, will contain information on the comparison of the dependencies found by structure101 to the dependencies found by some version of Pyne. Note that in the original xml output files contain the individual missing packages, classes and dependencies as well. Those have been left out to avoid cluttering the document.

### C.1 OUTPUT TEMPLATE

```
1 <?xml version="1.0" encoding="UTF-8" standalone="no"?>
2 <results>
3   <allPackages count="Integer" countExternal="Integer"
4     ↳ countInternal="Integer" countUnknown="Integer">
5     <pkg id="Integer" internal="Boolean">name.of.package</dependency>
6     etc.
7   </allPackages>
8   <allDependencies count="Integer">
9     <dep fromID="Integer" toID="Integer">
10      <!-- the following elements are only present when the
11      ↳ "human-readable" option is used -->
12      <fromIsInternal>Boolean</fromIsInternal>
13      <fromName>name.of.package</fromName>
14      <toIsInternal>Boolean</toIsInternal>
15      <toName>name.of.other.package</toName>
16    </dep>
17    etc.
18  </allDependencies>
19  <tools count="Integer">
20    <tool name="String">
21      <foundPackages count="Integer" countExternal="Integer"
22      ↳ countInternal="Integer" countUnknown="Integer"
23      ↳ percentageExternal="Float" percentageInternal="Float"
24      ↳ percentageTotal="Float" percentageUnknown="Float">
25        <!-- the "internal" attribute and the package name are not
26        ↳ present when the "compact" option is used -->
```

```

21     <pkg id="Integer" internal="Boolean">name.of.package</pkg>
22     etc.
23   </foundPackages>
24   <missedPackages count="Integer" countExternal="Integer"
→ countInternal="Integer" countUnknown="Integer"
→ percentageExternal="Float" percentageInternal="Float"
→ percentageTotal="Float" percentageUnknown="Float">
25     <!-- the "internal" attribute and the package name are not
→ present when the "compact" option is used -->
26     <pkg id="Integer" internal="Boolean">name.of.package</pkg>
27     etc.
28   </missedPackages>
29   <foundDependencies count="Integer" percentageTotal="Float">
30     <dep fromID="Integer" toID="Integer">
31       <!-- the following elements are only present when the
→ "human-readable" option is used -->
32       <fromIsInternal>Boolean</fromIsInternal>
33       <fromName>name.of.package</fromName>
34       <toIsInternal>Boolean</toIsInternal>
35       <toName>name.of.other.package</toName>
36     </dep>
37     etc.
38   </foundDependencies>
39   <missedDependencies count="Integer" percentageTotal="Float">
40     <dep fromID="Integer" toID="Integer">
41       <!-- the following elements are only present when the
→ "human-readable" option is used -->
42       <fromIsInternal>Boolean</fromIsInternal>
43       <fromName>name.of.package</fromName>
44       <toIsInternal>Boolean</toIsInternal>
45       <toName>name.of.other.package</toName>
46     </dep>
47     etc.
48   </missedDependencies>
49 </tool>
50 etc.
51 </tools>
52 </results>

```

Listing C.1: The structure of the output of our dependency checker

## C.2 ORIGINAL VERSION OF PYNE

### C.2.1 CLASSES

```

1 <results>
2   <allClasses count="3970" countExternal="237" countInternal="2184"
  ↳ countUnknown="1549"></allClasses>
3   <allDependencies count="41266"></allDependencies>
4   <tools count="2">
5     <tool name="STRUCTURE101">
6       <foundClasses count="3259" countExternal="223"
  ↳ countInternal="1487" countUnknown="1549"
  ↳ percentageExternal="94.09283" percentageInternal="68.08608"
  ↳ percentageTotal="82.09068" percentageUnknown="100.0"></foundClasses>
7       <missedClasses count="711" countExternal="14"
  ↳ countInternal="697" countUnknown="0" percentageExternal="5.907173"
  ↳ percentageInternal="31.913918" percentageTotal="17.909319"
  ↳ percentageUnknown="0"></missedClasses>
8       <foundDependencies count="36319"
  ↳ percentageTotal="88.01192"></foundDependencies>
9       <missedDependencies count="4947"
  ↳ percentageTotal="11.988077"></missedDependencies>
10      </tool>
11      <tool name="PYNE">
12        <foundClasses count="2348" countExternal="164"
  ↳ countInternal="2184" countUnknown="0" percentageExternal="69.19831"
  ↳ percentageInternal="100.0" percentageTotal="59.143578"
  ↳ percentageUnknown="0"></foundClasses>
13        <missedClasses count="1622" countExternal="73"
  ↳ countInternal="0" countUnknown="1549" percentageExternal="30.801687"
  ↳ percentageInternal="0" percentageTotal="40.856422"
  ↳ percentageUnknown="100.0"></missedClasses>
14        <foundDependencies count="7809"
  ↳ percentageTotal="18.923569"></foundDependencies>
15        <missedDependencies count="33457"
  ↳ percentageTotal="81.07643"></missedDependencies>
16      </tool>
17    </tools>
18 </results>

```

Listing C.2: A comparison of results on a class level before modifying Pyne

## C.2.2 PACKAGES

```

1 <results>
2   <allPackages count="430" countExternal="242" countInternal="188"
  ↪ countUnknown="0"></allPackages>
3   <allDependencies count="3558"></allDependencies>
4   <tools count="2">
5     <tool name="STRUCTURE101">
6       <foundPackages count="429" countExternal="241"
  ↪ countInternal="188" countUnknown="0" percentageExternal="99.58678"
  ↪ percentageInternal="100.0" percentageTotal="99.76744"
  ↪ percentageUnknown="0"></foundPackages>
7       <missedPackages count="1" countExternal="1" countInternal="0"
  ↪ countUnknown="0" percentageExternal="0.41322312"
  ↪ percentageInternal="0" percentageTotal="0.23255813"
  ↪ percentageUnknown="0"></missedPackages>
8       <foundDependencies count="3327"
  ↪ percentageTotal="93.50759"></foundDependencies>
9       <missedDependencies count="231"
  ↪ percentageTotal="6.492411"></missedDependencies>
10    </tool>
11    <tool name="PYNE">
12      <foundPackages count="193" countExternal="47"
  ↪ countInternal="146" countUnknown="0" percentageExternal="19.421488"
  ↪ percentageInternal="77.65958" percentageTotal="44.88372"
  ↪ percentageUnknown="0"></foundPackages>
13      <missedPackages count="237" countExternal="195"
  ↪ countInternal="42" countUnknown="0" percentageExternal="80.578514"
  ↪ percentageInternal="22.340425" percentageTotal="55.11628"
  ↪ percentageUnknown="0"></missedPackages>
14      <foundDependencies count="1020"
  ↪ percentageTotal="28.66779"></foundDependencies>
15      <missedDependencies count="2538"
  ↪ percentageTotal="71.33221"></missedDependencies>
16    </tool>
17  </tools>
18 </results>

```

Listing C.3: A comparison of results on a package level before modifying Pyne

## C.3 IMPROVED VERSION OF PYNE

### C.3.1 CLASSES

```

1 <results>
2   <allClasses count="4794" countExternal="1620" countInternal="2277"
   ↪ countUnknown="897">
3   <allDependencies count="58556">
4   <tools count="2">
5     <tool name="STRUCTURE101">
6       <foundClasses count="3259" countExternal="859"
   ↪ countInternal="1503" countUnknown="897"
   ↪ percentageExternal="53.024693" percentageInternal="66.007904"
   ↪ percentageTotal="67.980804" percentageUnknown="100.0">
7       <missedClasses count="1535" countExternal="761"
   ↪ countInternal="774" countUnknown="0" percentageExternal="46.975307"
   ↪ percentageInternal="33.992092" percentageTotal="32.01919"
   ↪ percentageUnknown="0">
8       <foundDependencies count="36319" percentageTotal="62.024384">
9       <missedDependencies count="22237" percentageTotal="37.975613">
10    </tool>
11    <tool name="PYNE">
12      <foundClasses count="3837" countExternal="1560"
   ↪ countInternal="2277" countUnknown="0" percentageExternal="96.296295"
   ↪ percentageInternal="100.0" percentageTotal="80.037544"
   ↪ percentageUnknown="0">
13      <missedClasses count="957" countExternal="60"
   ↪ countInternal="0" countUnknown="897" percentageExternal="3.7037036"
   ↪ percentageInternal="0" percentageTotal="19.962452"
   ↪ percentageUnknown="100.0">
14      <foundDependencies count="28889" percentageTotal="49.33568">
15      <missedDependencies count="29667" percentageTotal="50.66432">
16    </tool>
17  </tools>
18 </results>

```

Listing C.4: A comparison of results on a class level after modifying Pyne

## C.3.2 PACKAGES

```

1 <?xml version="1.0" encoding="UTF-8" standalone="no"?>
2 <results>
3   <allPackages count="497" countExternal="305" countInternal="192"
4   ↪ countUnknown="0">
5     <allDependencies count="3991">
6       <tools count="2">
7         <tool name="STRUCTURE101">
8           <foundPackages count="429" countExternal="237"
9           ↪ countInternal="192" countUnknown="0" percentageExternal="77.70492"
10          ↪ percentageInternal="100.0" percentageTotal="86.31791"
11          ↪ percentageUnknown="0">
12             <missedPackages count="68" countExternal="68"
13             ↪ countInternal="0" countUnknown="0" percentageExternal="22.295082"
14             ↪ percentageInternal="0" percentageTotal="13.682093"
15             ↪ percentageUnknown="0">
16               <foundDependencies count="3327" percentageTotal="83.362564">
17               <missedDependencies count="664" percentageTotal="16.637434">
18             </tool>
19             <tool name="PYNE">
20               <foundPackages count="393" countExternal="236"
21               ↪ countInternal="157" countUnknown="0" percentageExternal="77.37705"
22               ↪ percentageInternal="81.77083" percentageTotal="79.07445"
23               ↪ percentageUnknown="0">
24                 <missedPackages count="104" countExternal="69"
25                 ↪ countInternal="35" countUnknown="0" percentageExternal="22.62295"
26                 ↪ percentageInternal="18.229168" percentageTotal="20.925552"
27                 ↪ percentageUnknown="0">
28                   <foundDependencies count="2431" percentageTotal="60.912056">
29                   <missedDependencies count="1560" percentageTotal="39.087948">
30                 </tool>
31               </tools>
32             </results>

```

Listing C.5: A comparison of results on a package level after modifying Pyne