

Found in practice: Strategies for naming methods

Jeroen Peeters
The University of Amsterdam

March 17, 2014

Die Grenzen meiner Sprache bedeuten die Grenzen meiner Welt
— Ludwig Wittgenstein —

Master Thesis
Software Engineering

Supervisor: dr. Alexander Serebrenik

— Table of Contents —

	Page
1 Introduction	4
2 Verification model	5
3 Acknowledgement	7
References	8
A Frequency of Java Dependencies	9
B Logbook	12

1 Introduction

1.1 What's in a name

Naming pieces of code is one of the most important tasks of a software developer. Giving a piece of code the correct name is important for understandability. But what is the right name? Providing good quality names is generally believed to be of high importance. Though there may be some guidelines, no exact or mechanical process exists to name pieces of code. This thesis seeks to discover the strategies people use to create names for existing pieces of code written in the Java programming language.

1.2 Characteristics in code

Zoals besproken op Skype:

De huidige opzet van mijn scriptie/onderzoek blijkt lastig te voltooien. Jij stelde het volgende voor; ik kom een week naar de UvA om studenten en misschien docenten te interviewen over de vraag hoe mensen komen tot de naamgeving van een methode, gegeven de implementatie en context. Deze kennis kan interessant zijn om code automatisch te labelen bij automatische refactoring. Jij merkte op dat het wellicht interessant is om uit te zoeken welke karakteristieken in de code nu echt belangrijk zijn voor het geven van een naam. Het idee is om code met veel verschillende karakteristieken te geven, een complex vraagstuk. Volgens Kahneman zal dit leiden tot een substituuat vraag, welke karakteristieken zijn nu nog belangrijk? Altijd dezelfde? Of altijd de eerste N die men tegenkomt bij het lezen van de code. Maakt de context überhaupt uit? of is alleen de implementatie genoeg? of is misschien juist alleen de context genoeg voor het geven van een naam? Worden andere namen gegeven bij afwezigheid van context of implementatie. Verder kan nog worden gekeken naar de strategie die mensen kiezen: beantwoorden ze elke vraag op eenzelfde manier waardoor de gegeven namen ook vergelijkbaar zijn opgebouwd. Mogelijkheid is ook dat mensen geen naam kunnen geven.

Jij noemde ook inconsistentie als mogelijkheid om naar te kijken en je noemde als voorbeeld het gebruik van naamgevingsconventies; zou je nog kort kunnen toelichten wat je hiermee bedoelde?, ik kon dit uit m'n aantekeningen niet direct meer opmaken :S

Het idee is, volgens mij, om het volgende te organiseren: over 3-4 weken een week op de Uva voor het afnemen van interviews. Ik bereid vragen voor op papier met random vragen uit een codebase van github van meestgebruikte open-source software. Interview is interactief volgens 'thinking aloud protocol'.

Deze week + weekend zal ik dit plan verder uitwerken + voorbeeld vragen voorbereiden die wij maandag 9 november kunnen bespreken. Ik zal dan in de middag, na de lunch (rond 13:00u?), staan bij C302.4.

Hans, hartelijk dank voor je tijd zover en je aanbieding om studenten te interviewen.

Groet, Jeroen

2 Verification model

2.1 Survey: Naming Java Methods

Naming pieces of code is one of the most important tasks of a software developer. Giving a piece of code the correct name is important for understandability. But what is a correct name? Providing good quality names is generally believed to be of high importance. Though there may be some guidelines, no exact or mechanical process exists to name pieces of code. In order to understand how developers create method names, I've executed a survey among a group of peers.

2.1.1 Survey setup

The survey is executed as an online questionnaire which is filled out individually without direct supervision. Participants are asked to provide names for shown nameless Java methods. Together with the method implementation limited contextual information is given, such as the name of the containing class and examples of how the method is used. Because the number of methods that can be named by each participant can vary greatly there's no maximum number of questions. Instead, there's a time limit of thirty minutes per participant. After this time, the questionnaire will stop automatically. At any moment, the participant can pause and resume the questionnaire at a later time.

2.1.2 Method corpus

In order to make sure that the methods included in the survey are a representative sample I've employed the following strategy. The methods used in the survey are taken from open-source Java projects with a high usage frequency. In other words, which projects are depended upon the most? See appendix A for detailed information on how the list was compiled.

- JUnit
- Log4J
- Commons IO
- Guava
- Commons-lang
- Mockito

In order to make sure that the methods taken from these projects are also representative I use the SIG maintainability model [1]. To make sure that any degree of small & large and simple & complex method are selected I use the following two properties:

- Complexity per unit
The complexity of source code units influences the systems changeability and its testability.
- Unit size
The size of units influences their analysability and testability and therefore of the system as a whole.

2.1.3 Results

2.2 Deducing the model

3 Acknowledgement

TBD...

References

- [1] Ilja Heitlager, Tobias Kuipers, and Joost Visser. A practical model for measuring maintainability. In *Quality of Information and Communications Technology, 2007. QUATIC 2007. 6th International Conference on the*, pages 30–39. IEEE, 2007.

A Frequency of Java Dependencies

This is an edited version of my original article as posted on my personal website: <http://www.jeroenpeeters.nl/articles/frequency-of-java-dependencies/>.

A.1 The Approach

Github exposes an API through which you can search for projects with a certain language, rating, etc. Furthermore it is possible to query a projects tree structure and obtain file data.

Because I needed a representative data set I choose to include mature and active projects only. For this purpose a project is considered active if it had at least one commit in the last year. Secondly a project is considered mature if it is older than at least one year.

The Java world has three major build and dependency management tools; Maven, Gradle and Ivy. I simply downloaded the according build files to obtain dependency related information.

For each Github project each dependency is only counted once. This means that if a project contains multiple modules each having its own dependency management, duplicated dependencies are counted as one occurrence. Furthermore, build tools specific dependencies (such as maven-compiler-plugin) are omitted from the results. This is because depending on these artifacts is a consequence of using the build tool. They would thus occur frequently and cloud the results.

A.2 The Results

From the years 2008 to 2012 I was able to retrieve 3.029 projects with dependency management files (of which 2502 (82%) Maven projects, 430 (14%) Gradle projects and 97 (3%) Ant+Ivy projects).

From these figures it is not difficult to conclude that the majority of Java projects use Maven as a build and dependency management tool.

These projects had a total of 26.235 unique dependencies. The following list details the top 5 projects on which others depend:

- junit 1883
- slf4j-api - 764
- oss-parent 700
- log4j - 671
- commons-io 543

The following graphic shows the top 25 most depended on projects. We observe that JUnit is by far the most depended on artifact, followed by slf4j-api and oss-parent.

We can see that a large portion of these top projects are related to testing (junit, mockito-all, spring-test, mockito-core) and logging (slf4j-api, log4j, slf4j-log4j12, common-logging, logback-classic).

A.3 The code

To obtain and analyze the Github data I had to implement two relatively small programs. Both of them are freely available under the GNU GPL from, of course, Github.

- <https://github.com/jeroenpeeters/github-dependency-analyzer>

A.4 The data

The full data set contains 25,243 projects. Table 1 lists the first 30 projects from the data set with their dependency frequency. The complete data set, in raw and analyzed form, can be downloaded from the above mentioned Github repository.

Project name	Dependency count
junit	1883
slf4j-api	764
oss-parent	700
log4j	671
commons-io	543
guava	520
servlet-api	489
slf4j-log4j12	482
commons-lang	444
commons-logging	436
mockito-all	385
commons-codec	335
lifecycle-mapping	333
spring-context	332
httpclient	290
spring-test	288
logback-classic	284
joda-time	283
jackson-mapper-asl	273
jcl-over-slf4j	264
testng	263
spring-core	263
mockito-core	261
android	244
spring-beans	235
spring-web	234
commons-collections	228
hsqldb	218
spring-webmvc	216
mysql-connector-java	215

Table 1: Dependency frequency: Thirty most depended-on projects

B Logbook

17 March 2014

- Created new draft of plan
- Setup of thesis document: titelpage, raw sketch of chapters
- Online survey software: measure user activity, intake questions

18 March 2014

- Processed new comments on plan
- Online survey software: saving answers
- Investigate open source projects from which methods could be used in the survey

15 March 2014

- Finalized plan
- Randomly select methods from corpus of selected projects.