

Control flow in the wild : a first look at 13K Java projects (Extended abstract)

Citation for published version (APA):

Landman, D., Serebrenik, A., & Vinju, J. J. (2013). Control flow in the wild : a first look at 13K Java projects (Extended abstract). In T. Mens, M. Claes, M. Goeminne, & S. Drobisz (Eds.), *12th Belgian-Netherlands Software Evolution Seminar (BENEVOL'13), December 16-17, 2013, Mons, Belgium* (pp. 33-36). Universit  de Mons.

Document status and date:

Published: 01/01/2013

Document Version:

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

Control Flow in the Wild

A first look at 13K Java projects

Davy Landman*, Alexander Serebrenik*[†], Jurgen Vinju*

* Centrum Wiskunde & Informatica, Amsterdam, The Netherlands
{Davy.Landman, Jurgen.Vinju}@cwi.nl

[†] Eindhoven University of Technology, Eindhoven, The Netherlands
a.serebrenik@tue.nl

I. INTRODUCTION

We are interested in the understandability of software. Maintainability models such as the SIG model use cyclomatic complexity to measure understandability. However, doubts have been raised about the relation between cyclomatic complexity and understanding of the code. In a grounded theory approach we first observe control flow in a large corpus. Which in the long term will enable us to find categories and create well-founded metrics or indicators for understandability.

We present our early observations of Control Flow Patterns (CFPs) [1] in the Sourcerer Corpus [2], a set of 13 thousand Java projects. We observe saturations when CFPs belonging to two or more systems are considered, but no saturation when all patterns are considered. Most observed patterns are unique, only present in one system, moreover they are small, less than 20 statements. We conclude with questions for future research.

II. EXPERIMENT

We took the Sourcerer Corpus which contains 18K (13K non empty) Java projects. Using a Java grammar and RASCAL [3] we parsed all Java files. All methods were transformed [1] into CFPs.

A CFP is an AST created by removing all statements not related to control flow. Table I contains a list of Java’s language constructs kept. The last step is to change all expressions inside the arguments of the constructs into an empty string.

TABLE I
JAVA LANGUAGE CONSTRUCTS USED IN A CFP.

if	if else	switch	case	labeled	continue	break
for	while	do while	return	throw	synchronized	try

Table II describes how large the Sourcerer corpus is, and how many CFPs we extracted and how many of those CFPs were unique to one system.

TABLE II
SIZE OF SOURCERER CORPUS AND EXTRACTED CFPs

Size	Files	LOC [†]	Methods	CFPs	CFPs _{unique} [‡]
19GB	2M	477M	23M	678K	516K

[†] measured using `wc -l`

[‡] CFPs only observed in one system.

III. OBSERVATIONS

Figure 1 shows the amount of CFPs observed, where we see that almost every time when we add a new systems, we observe new patterns. Narrowing our definition of a pattern, only considering patterns present in 2 or more systems, we observe a saturation. Even more so for patterns shared by 3 or more and 4 or more. Figure 2 shows these narrowed definitions in more detail.

Unique CFPs are patterns only occurring in exactly one system. The almost linear growth in Figure 1 raises the question what causes it. Figure 3 shows that this is not primarily caused by large patterns, that most unique patterns are actually smaller then 20 control flow statements.

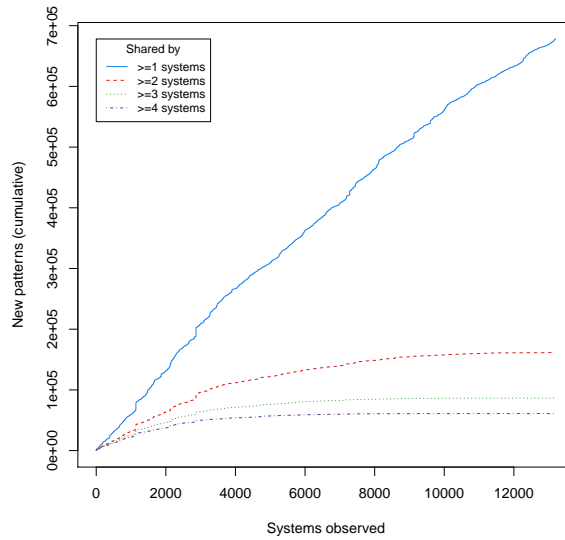


Fig. 1. Saturation of the patterns in the Sourcerer corpus. The four lines represent the saturation of patterns appearing in x or more systems.

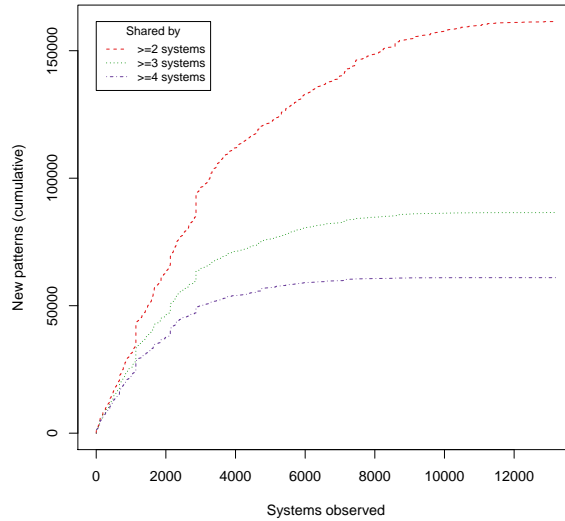


Fig. 2. Zoomed in on the patterns shared in more than one system.

The theoretical reason for so many small *unique* CFPs, is the exponential growth in possible patterns. For size 4 there are already 2,474,634 possible CFPs. Figure 4 shows how many different CFPs per size were observed and it shows the theoretical maximum.

Figure 5 shows the distribution of the size of a CFP and in how many systems it occurs. Here we can see that even the larger CFPs are shared. Eye-balling these larger shared CFPs revealed code clones and code generated by the same generator. We also observed code clones where the full library was embedded.

IV. OPEN QUESTIONS

As future work we have the following questions:

- 1) Are systems with a lot of CFPs not using OO constructs?
- 2) Can we find categories of CFPs?
- 3) Are CFPs abstract enough?
- 4) Can we find a relation between the naming of a method and its CFP?
- 5) If we observe more systems, would the saturation change?
- 6) If we analyse non Java systems, would we find similar patterns and saturations?
- 7) What would be the impact of removing clones on the amount of shared patterns?

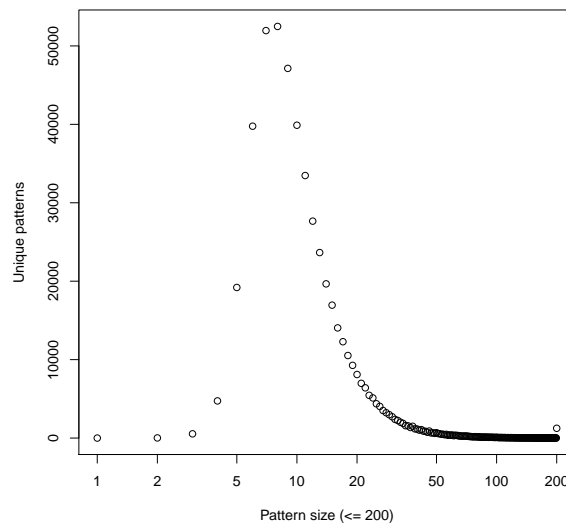


Fig. 3. How many unique patterns of a certain size are found. Patterns equal or larger than 200 are grouped to show that the long-tail does not contribute that much to the continuous growth of patterns observed.

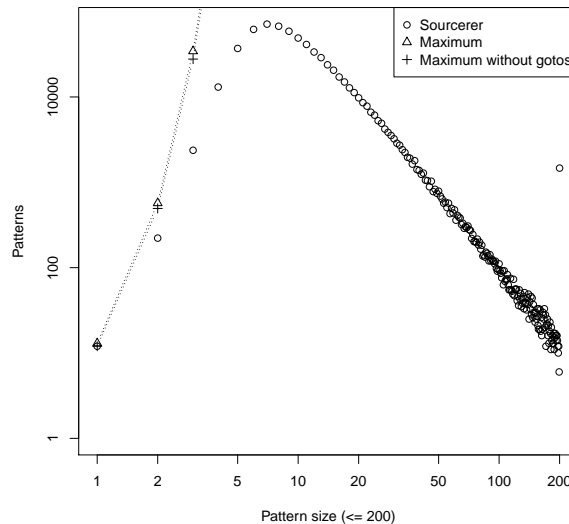


Fig. 4. How many different patterns of a certain size were observed and the maximum possible (only shown for 1 – 3 due to the exponential growth). The “maximum without GOTOs” does not take into account the “structured” GOTOs in Java.

- 8) Could CFPs be used to fingerprint systems?
- 9) Why is there so much control flow in an OO language?

REFERENCES

- [1] J. J. Vinju and M. W. Godfrey, “What does control flow really look like? eyeballing the cyclomatic complexity metric,” in *Ninth IEEE International Working Conference on Source Code Analysis and Manipulation (SCAM)*. IEEE Computer Society, 2012.
- [2] E. Linstead, S. K. Bajracharya, T. C. Ngo, P. Rigor, C. V. Lopes, and P. Baldi, “Sourcerer: mining and searching internet-scale software repositories,” *Data Min. Knowl. Discov.*, vol. 18, no. 2, pp. 300–336, 2009.
- [3] P. Klint, T. van der Storm, and J. J. Vinju, “RASCAL: A Domain Specific Language for Source Code Analysis and Manipulation,” in *Proc. 9th IEEE International Working Conference on Source Code Analysis and Manipulation (SCAM)*. IEEE, September 2009, pp. 168–177.

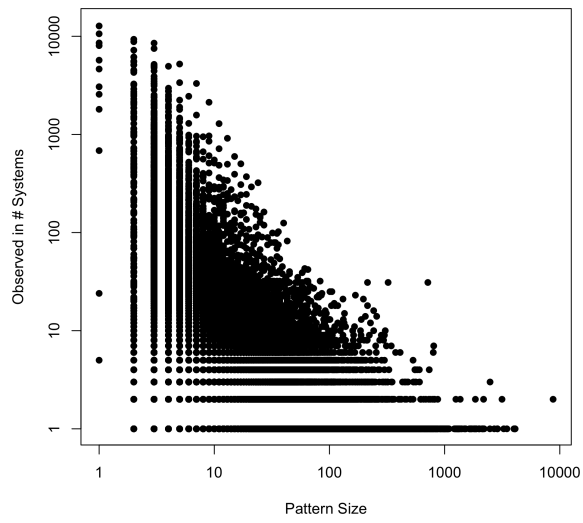


Fig. 5. A scatterplot of the pattern's sizes and in how many projects they occurred.