

A peer-reviewed version of this preprint was published in PeerJ on 2 March 2016.

[View the peer-reviewed version](https://peerj.com/articles/cs-49) (peerj.com/articles/cs-49), which is the preferred citable publication unless you specifically need to cite this preprint.

Wagner S, Abdulkhaleq A, Bogicevic I, Ostberg J, Ramadani J. 2016. How are functionally similar code clones syntactically different? An empirical study and a benchmark. PeerJ Computer Science 2:e49
<https://doi.org/10.7717/peerj-cs.49>

How are functionally similar code clones different?

Stefan Wagner, Asim Abdulkhaleq, Ivan Bogicevic, Jan-Peter Ostberg, Jasmin Ramadani

Background. Today, redundancy in source code, so-called “clones”, caused by copy&paste can be found reliably using clone detection tools. Redundancy can arise also independently, however, caused not by copy&paste. At present, it is not clear how only *functionally similar clones* (FSC) differ from clones created by copy&paste. Our aim is to understand and categorise the differences in FSCs that distinguish them from copy&paste clones in a way that helps clone detection research. **Methods.** We conducted an experiment using known functionally similar programs in Java and C from coding contests. We analysed syntactic similarity with traditional detection tools and explored whether concolic clone detection can go beyond syntax. We ran all tools on 2,800 programs and manually categorised the differences in a random sample of 70 program pairs. **Results.** We found no FSCs where complete files were syntactically similar. We could detect a syntactic similarity in a part of the files in < 16 % of the program pairs. Concolic detection found 1 of the FSCs. The differences between program pairs were in the categories algorithm, data structure, OO design, I/O and libraries. We selected 58 pairs for an openly accessible benchmark representing these categories. **Discussion.** The majority of differences between functionally similar clones are beyond the capabilities of current clone detection approaches. Yet, our benchmark can help to drive further clone detection research.

How Are Functionally Similar Code Clones Different?

Stefan Wagner*, Asim Abdulkhaleq, Ivan Bogicevic, Jan-Peter
Ostberg, and Jasmin Ramadani

Institute of Software Technology, University of Stuttgart, Germany

Abstract

Background. Today, redundancy in source code, so-called “clones”, caused by copy&paste can be found reliably using clone detection tools. Redundancy can arise also independently, however, caused not by copy&paste. At present, it is not clear how only functionally similar clones (FSC) differ from clones created by copy&paste. Our aim is to understand and categorise the differences in FSCs that distinguish them from copy&paste clones in a way that helps clone detection research.

Methods. We conducted an experiment using known functionally similar programs in Java and C from coding contests. We analysed syntactic similarity with traditional detection tools and explored whether concolic clone detection can go beyond syntax. We ran all tools on 2,800 programs and manually categorised the differences in a random sample of 70 program pairs.

Results. We found no FSCs where complete files were syntactically similar. We could detect a syntactic similarity in a part of the files in < 16 % of the program pairs. Concolic detection found 1 % of the FSCs. The differences between program pairs were in the categories *algorithm*, *data structure*, *OO design*, *I/O* and *libraries*. We selected 58 pairs for an openly accessible benchmark representing these categories.

Discussion. The majority of differences between functionally similar clones are beyond the capabilities of current clone detection approaches. Yet, our benchmark can help to drive further clone detection research.

*Corr. author: Stefan Wagner, Universitätsstr. 38, 70569 Stuttgart, Germany, phone +49 711 685 88455, stefan.wagner@informatik.uni-stuttgart.de

29 1 Introduction

30 As software is now a key ingredient of current complex systems, the size of
31 software systems is continuously increasing. While software with a code size of
32 several thousand lines has been considered large in the seventies and eighties,
33 we now reach code sizes of hundreds of millions of lines of code. This has strong
34 effects on the complexity and manageability of these systems and, as a result,
35 on the cost of maintaining them.

36 By abstraction and code generation, modern programming languages and
37 development techniques help to reduce the amount of code we have to under-
38 stand. Nevertheless, it still tends to be overwhelming. A factor that aggravates
39 the situation is that there is unnecessary code in these huge code bases: un-
40 reachable code, never executed code and redundant code. The latter has been
41 of increasing interest in the software engineering research community under the
42 term “cloning”. Especially clones that resulted from copy&paste can now be
43 detected reliably. In our own research with companies, we often found rates of
44 redundant code caused by copy&paste in the range of 20 % – 30 % [Wagner,
45 2013].

46 More challenging is the detection of functionally similar source code. We
47 will refer to it as *functionally similar clones* (FSCs). FSCs were not created by
48 copy&paste but developers independently needed and implemented certain func-
49 tionalities in their code base. We deliberately go beyond *type-4 clones* [Koschke,
50 2007] or *simions* [Deissenboeck et al., 2012] which only include functional equiv-
51 alence. In a refactoring session to reduce the size of a code base, a developer
52 would still be interested in mostly similar and not only exactly equivalent func-
53 tionality. Although this problem is in general undecidable, there have been
54 several heuristic efforts [Marcus and Maletic, 2001, Jiang and Su, 2009, Deis-
55 senboeck et al., 2012, Kim et al., 2011].

56 Juergens, Deissenboeck and Hummel [Juergens et al., 2010b] showed that
57 traditional clone detection approaches and tools are hardly able to detect func-
58 tionally equivalent clones because they rely on syntactic similarities.

59 1.1 Problem Statement

60 So far, the work by Juergens, Deissenboeck and Hummel [Juergens et al., 2010b]
61 is the only study investigating the differences in functionally similar clones.
62 Furthermore, their study is limited: they use only programs implementing a
63 single specification in Java. Therefore, we have no clear understanding of what
64 differences make a functionally similar clone really different from copy&paste
65 clones. Hence, a realistic, open benchmark for comparing and improving such
66 approaches is also lacking although it is necessary for faster progress in the
67 field [Lakhotia et al., 2003].

68 1.2 Research Objectives

69 The objective of this study is to better understand the differences that make
70 up functionally similar clones to support future research on their detection. In
71 particular, we want to classify and rate differences and build a representative
72 benchmark.

73 1.3 Contribution

74 We contribute a large-scale quantitative study combined with a qualitative anal-
75 ysis of the differences. We selected 2,800 Java and C programs which are solu-
76 tions to the Google Code Jam programming contest and are therefore function-
77 ally similar. We identified copy&paste clones by using two clone detection tools
78 (ConQAT and Deckard) to quantify syntactic similarities. We explored how a
79 type-4 detection tool (CCCD) using cocolic detection performs in detecting the
80 not syntactically similar FSCs. We created a categorisation of differences be-
81 tween undetected clones and quantified these categories. Finally, we derived a
82 benchmark based on real FSCs covering the categories and degrees of differences
83 which can drive the improvement of clone detection tools.

84 As there is a large diversity in how the terms around FSCs are used, we
85 provide definitions for the clone types we investigate in this paper. Moreover,
86 we define terms for granularities of the software programs under analysis in
87 Tab. 1.

88 The structure of the remainder of the paper follows the guidelines in [Jedl-
89 itschka and Pfahl, 2005].

90 2 Related Work

91 A *code clone* consists of at least two pieces of code that are similar according
92 to a definition of similarity. Most commonly, **clone detection** approaches look
93 for exact clones (also called *type-1*) and clones with simple changes such as re-
94 naming (also called *type-2*). These types of clones are detectable today in an
95 efficient and effective way. Even clones with additional changes (*inconsistent*,
96 *near-miss* or *type-3* clones) can be detected by several detection approaches
97 and tools [Kamiya et al., 2002, Deissenboeck et al., 2008, Jiang et al., 2007a].
98 There are also two surveys [Koschke, 2007, Roy and Cordy, 2007] and a system-
99 atic literature review [Rattan et al., 2013] on this topic. Tiarks, Koschke und
100 Falke [Tiarks et al., 2011] investigated in particular type-3 clones and also their
101 differences. They concentrated, however, on differences in code metrics (e.g.
102 fragment size), low level edits (e.g. variable) and abstracted them only slightly
103 (e.g. to type substitution).

104 Juergens, Deissenboeck and Hummel [Juergens et al., 2010b] report on an
105 experiment to investigate the differences between syntactical/representational
106 and semantic/behavioural similarities of code and the detectability of these sim-
107 ilarities. They use a simple student assignment called *email address validator*
108 and also inspect the open-source software *JabRef*. Both of them are in Java.

Table 1: Terminology

Type-1 clone	Similar code fragments except for variation in whitespace, layout and comments [Bellon et al., 2007]
Type-2 clone	Similar code fragments except for variation in identifiers, literals, types, whitespaces layouts and comments [Bellon et al., 2007]
Type-3 clone	Similar code fragments except that some statements may be added or deleted in addition to variation in identifiers, literals, types, whitespaces, layouts or comments [Bellon et al., 2007]
Type-4 clone	Code fragments that perform the same function but are implemented quite differently [Bellon et al., 2007]
Functionally similar clone (FSC)	Code fragments that perform a similar function but are implemented quite differently
Solution file	A single program in one file implementing the solution to a programming problem
Solution set	A set of solution files all solving the the same programming problem
Clone pair	Two solution files from the same solution set which we assume to be functionally similar

109 To detect the clones of types 1–3, they use the clone detection tools ConQAT
 110 and Deckard. They review the open-source system manually to identify if be-
 111 haviourally similar code that does not result from copy&paste can be detected
 112 and occurs in real-world software. The results indicate that behaviourally sim-
 113 ilar code of independent origin is highly unlikely to be syntactically similar.
 114 They also report that the existing clone detection approaches cannot identify
 115 more than 1 % of such redundancy. We build our work on their study but
 116 concentrate on understanding the differences in more detail based on a diverse
 117 sample with a larger sample size and different programming languages.

118 Several researchers have proposed to move away from the concrete syntax
 119 to detect what they call **semantic clones**. Marcus and Malefic [Marcus and
 120 Maletic, 2001] used information retrieval techniques on source code to detect se-
 121 mantic similarities. Krinke [Krinke, 2001] proposed to use program dependence
 122 graphs (PDG) for abstracting source code. Komondoor and Horwitz [Komon-
 123 door and Horwitz, 2001] also use PDGs for clone detection and see the possibility
 124 to find non-contiguous clones as a main benefit. Gabel, Jiang and Su [Gabel
 125 et al., 2008] combine the analysis of dependence graphs with abstract syntax
 126 trees in the tool Deckard to better scale the approach.

127 A very different approach to detecting semantic clones comes from Kim et

128 al. [Kim et al., 2011] who use static analysis to extract the memory states
129 for each procedure exit point. They can show that they find more semantically
130 similar procedures as clones than previous clone detectors including PDG-based
131 detectors. Nevertheless, the used approach as well as the examples of found
132 semantic clones suggest that the syntactic representation still plays a role and
133 that the clones have been created by copy&paste.

134 These semantic clone detection techniques cannot guarantee that they also
135 find all functionally similar clones as a completely different structure and mem-
136 ory states can generate similar functionality.

137 Jiang and Su [Jiang and Su, 2009] were the first to comprehensively detect
138 **functionally similar code** by using random tests and comparing the output.
139 Hence, they were also the first who were able to detect clones without any
140 syntactic similarity. They claim they are able to detect “*functionally equivalent*
141 code fragments, where functional equivalence is a particular case of semantic
142 equivalence that concerns the input/output behavior of a piece of code.” They
143 were able to detect a high number of functionally equivalent clones in a sorting
144 benchmark and the Linux kernel. Several of the detected clones are dubious,
145 however, as it is not clear how useful they are. They state: “Assuming the
146 input and output variables identified by EQMINER for these code fragments are
147 appropriate, such code fragments are indeed functionally equivalent according to
148 our definition. However, whether it is really useful to consider them functionally
149 equivalent is still a question worth of future investigation.”

150 Deissenboeck et al. [Deissenboeck et al., 2012] followed an analogous ap-
151 proach to Jiang and Su [Jiang and Su, 2009] to detect functionally similar code
152 fragments in Java systems based on the fundamental heuristic that two func-
153 tionally similar code fragments will produce the same output for the same ran-
154 domly generated input. They implemented a prototype based on their toolkit
155 ConQAT. The evaluation of the approach involved 5 open-source systems and
156 an artificial system with independent implementations of the same specification
157 in Java. They experienced low detection results due to the limited capability of
158 the random testing approach. Furthermore, they mention that the similarities
159 are missed due to chunking, i.e. if the code fragments perform a similar com-
160 putation but use different data structures at their interfaces. They emphasise
161 that further research is required to understand these issues.

162 CCCD [Krutz and Shihab, 2013] also claims to detect functionally similar
163 code for C programs based on concolic analysis. Its creators evaluated their
164 implementation of the approach on the benchmarks mentioned below and found
165 a 92 % recall even in the type-4 clones in those benchmarks. As the tool is freely
166 available in a virtual machine, we were able to include it in our experiment.

167 A clear comparison and measurement of the improvement in clone detection
168 research would require a comprehensive **benchmark**. There have been few
169 approaches [Lakhotia et al., 2003, Roy et al., 2009a, Tempero, 2013] trying to
170 establish a benchmark but they are either small and artificial or do not con-
171 tain (known) FSCs. The only exception is the recent *BigCloneBench* [Svajlenko
172 et al., 2014] which has a huge number of clones mined from source code reposi-
173 tories. Yet, they do not classify the types of differences and also state “there is

174 no consensus on the minimum similarity of a Type-3 clone, so it is difficult to
175 separate the Type-3 and Type-4 clones”.

176 3 Experimental Design

177 To reach our research objectives, we developed a study design based on the idea
178 that we investigate sets of programs which we knew to be functionally similar:
179 accepted submissions to programming contests. We formulated four research
180 questions which we all answer by analysing these programs and corresponding
181 detection results. All instrumentation, analysis scripts and results are freely
182 available in [Wagner et al., 2014].

183 3.1 Research Questions

184 As we have independently developed but functionally similar programs, we first
185 wanted to establish how much syntactic similarity is in these programs. We can
186 investigate this by quantifying the share of type-1–3 clones

187 **RQ 1: What share of independently developed similar programs are**
188 **type-1–3 clones?**

189 Then we wanted to understand what is different in clones not of type-1–
190 3. This should result in a categorisation and rating of the differences between
191 FSCs.

192 **RQ 2: What are the differences between FSC that go beyond type-1–3**
193 **clones?**

194 Although we could not fully evaluate type-4 detectors, we wanted at least
195 to explore what a modern clone detection approach can achieve on our FSCs.
196 This should give us an indication how much more research is needed on those
197 detection approaches.

198 **RQ 3: What share of FSC can be detected by a type-4 clone detector?**

199 Finally, to make our results an operational help for clone detection research,
200 we wanted to create a representative benchmark from the non-detected clones.

201 **RQ 4: What should a benchmark contain that represents the differ-**
202 **ences between FSC?**

203 3.2 Hypotheses

204 We define two hypotheses regarding RQ 1. As we investigate the share of
205 detectable Type-1–3 clones, we wanted to understand if there are differences
206 between the used tools and analysed languages because this might have an
207 influence on the generalisability of our results. We formulated the two null
208 hypotheses:

209 **H1: There is no difference in the share of detected Type-1–3 clones**
210 **between programming languages.**

211 **H2: There is no difference in the share of detected Type-1–3 clones**
212 **between clone detection tools.**

213 Moreover, in RQ 2, we wanted to understand the characteristics of non-
 214 detected clone pairs and, therefore, categorised them. In this categorisation,
 215 we also rated the degree of difference in each category. An ideal categorisa-
 216 tion would have fully orthogonal categories and, hence, categories would not be
 217 correlated in the degree of difference:

218 **H3: There is no correlation between the degrees of difference between**
 219 **categories.**

220 Furthermore, we could imagine that different programming languages might
 221 cause disparately strong differences in certain categories. As this again has an
 222 impact on the generalisability of our results, we formulated this null hypotheses:

223 **H4: There is no difference in the degree of difference between pro-**
 224 **gramming languages.**

225 3.3 Design

226 The overall study design is a combination of quantitative and qualitative anal-
 227 ysis. For the quantitative part of our study we used a factorial design with
 228 two factors (programming language and clone detection tool). As applying the
 229 treatments of both factors was mostly automated we could apply almost all
 230 factor levels to all study object programs (which we call *solutions*). Only if a
 231 detection tool did not support a certain programming language, we would not
 232 apply it. We tried to minimise that but to include a contemporary tool, we
 233 accepted an unbalanced design. Table 2 shows the factors in our experiment.

Table 2: The factorial design used in this experiment

		Programming language	
		Java	C
Clone	CCCD	–	✓
detection	ConQAT	✓	✓
tool	Deckard	✓	✓

234 We will describe the programming languages, clone detection tools and cor-
 235 responding programs under analysis in more detail in the next subsection.

236 3.4 Objects

237 The general idea of this experiment was that we analyse accepted solutions to
 238 programming contests because we know that for a given problem, the solutions
 239 must be functionally similar. Therefore, our selection of study objects needed
 240 to include clone detection tools we could access and execute as well as solutions
 241 in programming languages supported by most of the detection tools.

242 3.4.1 Clone Detection Tools

243 Primarily, we needed clone detection tools for detecting type-1–3 clones to in-
244 vestigate with RQ 1 the syntactic similarity of FSCs. We did a literature and
245 web search for available tools.

246 Many research prototypes were not available or could not be brought to execute
247 correctly. Commercial tools were not exact enough in what they detect. Several
248 tools were not included in the study due their lower performance and scalabil-
249 ity or their lack of support for some clone types. CloneDR and CPMiner have
250 lower performance and scalability compared to Deckard [Jiang et al., 2007a].
251 CCFinder has also lower performance than Deckard and does not support type-
252 3 clones [Svajlenko and Roy, 2014].

253 In the end, we chose two clone detection tools that both can analyse Java and
254 C programs: **ConQAT** [Deissenboeck et al., 2008] and **Deckard** [Jiang et al.,
255 2007a]. They have been described as most up-to-date implementations of to-
256 ken based and AST based clone detection algorithms [Juergens et al., 2010b]
257 Choosing tools based on different techniques, we allow distinct approaches in
258 finding code clones [Roy et al., 2009b].

259 **ConQAT** is a stable open-source dashboard toolkit also used in industry.
260 It is a general-purpose tool for various kinds of code measurement and analysis.
261 For our experiment, ConQAT offers several specific clone detection configu-
262 rations for various programming languages including Java, C/C++, C# and
263 Cobol. It has separate detection algorithms for type-1/2 clones and type-3
264 clones. We employed the latter algorithm. ConQAT has been used in various
265 studies on clone detection [Juergens et al., 2009, Juergens et al., 2010a] including
266 the study we build on [Juergens et al., 2010b].

267 The language-independent clone detection tool **Deckard** works on code in
268 any programming language that has a context-free grammar. Deckard uses an
269 efficient algorithm for identifying similar subtrees and applies it to tree represen-
270 tations of source code. It automatically generates a parse tree builder to build
271 parse trees required by its algorithm. By a similarity parameter it is possible
272 to control whether only type-1/2 clones or type-3 clones are detected. Deckard
273 is a stable tool used in other studies [Gabel et al., 2008, Jiang et al., 2007b]
274 including the study we build on.

275 To explore the state of type-4 clone detection tools, we also searched for
276 such tools. Most existing tools, however, could not be used. For example,
277 EqMiner [Jiang and Su, 2009] was too tightly coupled with the Linux kernel
278 and MeCC [Kim et al., 2011] could not detect clones across files. Finally, we
279 were able to only include a single type-4 detector.

280 **CCCD** [Krutz and Shihab, 2013] is a novel clone detection tool that uses
281 concolic analysis as its primary approach to detect code clones. Concolic anal-
282 ysis combines symbolic execution and testing. CCCD detects only clones in
283 programs implemented in C. The concolic analysis allows CCCD to focus on
284 the functionality of a program rather than the syntactic properties. Yet, it has
285 the restriction that it only detects function-level clones.

286 3.4.2 Solution Sets and Solutions

287 We looked at several programming contests and the availability of the submitted
 288 solutions. We found that *Google Code Jam*¹ provided us with the broadest
 289 selection of programming languages and the highest numbers of submissions.
 290 Google Code Jam is an annual coding contest organised by Google. Several
 291 tens of thousands of people participate each year. In seven competition rounds,
 292 the programmers have to solve small algorithmic problems within a defined time
 293 frame. Although over one hundred different programming languages are used,
 294 the majority of the solutions are in C, C++, Java and Python. Most solutions
 295 of the participants are freely available on the web.²

296 We define a *solution* as a single code file delivered by one participant during
 297 the contest. We define a *solution set* as a set of solutions all solving the same
 298 problem. A typical solution set consists of several hundred to several thousand
 299 solutions. We can be sure that all solutions in a solution set should be FSCs be-
 300 cause they passed the judgement of the programming contest. Even if there are
 301 differences in the programs, e.g. in the result representation, these are instances
 302 of similarity instead of equivalence.

303 We selected 14 out of 27 problem statements of the Google Code Jam 2014.
 304 For every problem we randomly chose 100 solutions in Java and 100 solutions
 305 in C from sets of several hundreds to several thousands of solutions. Table 3
 306 shows a summary of the size of the chosen solution sets. Hence, on average a C
 307 solution has a length of 46 LOC and a Java solution of 94 LOC.

Table 3: Summary of the Solution Sets

	#No. Sets	#Files/Set	Size	LOC
C	14	100	6 MB	64,826
Java	14	100	7 MB	131,398

308 In Table 4, we detail the size of the selected Java solution sets and in Table 5
 309 of the C solution sets. The solution sets differ in size but the means all lie
 310 between 33 and 133 LOC per solution.

311 3.5 Data Collection Procedure

312 3.5.1 Preparation of Programs Under Analysis

313 We implemented an instrumentation which automatically downloaded the solu-
 314 tions from the website, sampled the solution sets and solutions and normalised
 315 the file names. The instrumentation is freely available as Java programs in our
 316 GitHub project. Every downloaded solution consisted of a single source code
 317 file.

¹<https://code.google.com/codejam/>

²<http://www.go-hero.net/jam/14/>

Table 4: Information on the Java Solution Sets

Set	#Files	LOC	#Functions
1	100	11,366	823
2	100	7,825	523
3	100	10,624	575
4	100	6,766	473
5	100	7,986	585
6	100	10,137	611
7	100	13,300	869
8	100	8,568	614
9	100	8,580	717
10	100	9,092	459
11	100	8,536	584
12	100	11,412	648
13	100	9,436	465
14	100	7,770	357

318 3.5.2 Configuration of Clone Detection Tools

319 We installed ConQAT, Deckard and CCCD and configured the tools with a
 320 common set of parameters. As far as the parameters between the tools were
 321 related to each other, we tried to set the same values based on the configuration
 322 in [Juergens et al., 2010b]. We set the parameters conservatively so that the
 323 tools find potentially more clones as we would normally consider valid clones.
 324 This ensured that we do not reject our null hypotheses because of configura-
 325 tions. For example, we set the minimal clone length in ConQAT to 6 statements.
 326 All the detailed configurations are available on GitHub. For CCCD, we only in-
 327 cluded clone pairs which have a Levenshtein similarity score below 35 as advised
 328 in [Krutz and Shihab, 2013] where the score is calculated for the concolic output
 329 for each function. The detection failed for 4 of the solutions of our sample set.
 330 We had to exclude them from further analysis.

331 3.5.3 Executing Clone Detection Tools

332 We manually executed the clone detection tools for every solution set. Con-
 333 QAT generated an XML file for every solution set containing a list of found
 334 clone classes and clones. Deckard and CCCD generate similar CVS files. Our
 335 instrumentation tool parsed all these result files and generated reports in a uni-
 336 fied format. The reports are tables in which both rows and columns represent
 337 the solutions. The content of the table shows the lowest detected clone type
 338 between two files. Additionally, our tool calculated all shares of syntactic sim-
 339 ilarity as described in the next section and wrote the values into several CSV
 340 files for further statistical analysis. We also wrote all the detected clones into

Table 5: Information on the C Solution Sets

Set	#Files	LOC	#Functions
1	100	3,917	233
2	100	3,706	167
3	100	4,750	265
4	100	3,928	219
5	100	4,067	187
6	100	6,840	166
7	100	4,701	263
8	100	4,679	176
9	100	6,831	227
10	100	4,063	159
11	100	4,624	266
12	100	3,574	163
13	100	3,335	168
14	100	5,811	249

341 several large CSV files. Altogether, the tools reported more than 9,300 clones
 342 within the Java solutions and more than 22,400 clones within the C solutions.

343 3.6 Analysis Procedure

344 3.6.1 Share of Syntactic Similarity (RQ 1)

345 All solutions in a solution set solve the same programming problem and were
 346 accepted by Google Code Jam. Hence, their functionality can only differ slightly
 347 and, therefore, they are functionally similar. To understand how much of this
 348 similarity is expressed in syntactic similarity, we calculate the share of FSCs
 349 which are also type-1-2 or type-1-3 clones.

350 Inspired by [Juergens et al., 2010b], we distinguish partial and full syntactic
 351 similarity. The *share of full syntactic similarity* is the ratio of clone pairs where
 352 all but an defined looseness of the statements of the solutions of the pair were
 353 detected as a clone in relation to all clone pairs. We set the threshold of this
 354 looseness to a maximum of 16 lines of code difference within a clone pair, which
 355 leads to ratios of 5% to 33 % of difference based on the functions lines of code.

$$\text{Share of full synt. similarity} = \frac{|\text{Found full clone pairs}|}{|\text{All clone pairs}|} \quad (1)$$

356 Because we expected the share of full syntactic similarity to be low, we
 357 wanted to check whether there are at least some parts with syntactic similarity.
 358 It would give traditional clone detection tools a chance to hint at the FSC.
 359 Furthermore, it allowed us to inspect more closely later on what was not detected
 360 as a clone. We called the ratio *share of partial syntactic similarity*.

$$\text{Share of partial synt. similarity} = \frac{|\text{Found partial clone pairs}|}{|\text{All clone pairs}|} \quad (2)$$

361 For a more differentiated analysis, we calculated two different shares each
362 representing certain types of clones. We first computed the share for type-1-
363 2 clones. This means we only need to accept exact copies, reformatting and
364 renaming. Then, we determined the shares for type-1-3 clones which includes
365 type-1-2 and adds the additional capability to tolerate smaller changes.

366 In ConQAT and Deckard, we can differentiate between type-1/2 clones and
367 type-3 clones by configuration or result, respectively. In ConQAT, clones with
368 a gap of 0 are type-1/2 clones. In Deckard, analysis results with a similarity
369 of 1 are type-1/2 clones. The others are type-3 clones. The instrument tooling
370 described in Sec. 3.5 directly calculated the various numbers. We computed
371 means per clone type and programming language.

372 For a further statistical understanding and to answer the hypotheses H1-
373 H4, we did statistical hypotheses tests. For answering H1 and H2, we performed
374 an analysis of variance (ANOVA) on the recall data with the two factors *pro-*
375 *gramming language* and *detection tool*. We tested the hypotheses at the 0.05
376 level. All analyses implemented in R together with the data are available in our
377 GitHub project.

378 The combined descriptive statistics and hypothesis testing results answered
379 RQ 1.

380 3.6.2 Classifying Differences (RQ 2)

381 For the categorisation of the differences of FSCs that were not syntactically
382 similar, we took a random sample of these clone pairs. As we had overall 69,300
383 clone pairs for Java and C, we needed to restrict the sample for a manual anal-
384 ysis. We found in an initial classification (see also Sec. 3.7) that a sample of
385 0.5 ‰ per language and fully/partially different clone pairs is sufficient for find-
386 ing repeating categories and getting a quantitative impression of the numbers
387 of clone pairs in each category. With larger samples, the categories just kept
388 repeating. Therefore, we took a sample of 2 ‰ of the syntactically different
389 clone pairs: 70 pairs each of the fully and partially different clone pairs (35 C
390 and 35 Java).

391 The set of *fully syntactically different clone pairs* is the set of all pairs in
392 all solution sets minus any pair detected by any of the type-1-3 detection. We
393 apply random sampling to get pairs for further analysis: First, we randomly
394 selected one of the solution sets in a language. Second, we randomly selected
395 a solution file in the solution set and checked if it was detected by Deckard or
396 ConQAT. If it was detected, we would discard it and select a new one. Third,
397 we randomly picked a second solution file, checked again if it was detected and
398 discard it if it was.

399 The set of *partially syntactically different clone pairs* is then the superset of
400 all partially different clone pairs minus the superset of all fully different clone
401 pairs. From that set, we randomly selected clone pairs from all partially different

402 pairs of a programming language and checked if it was fully different. If that
403 was the case, we would discard it and take a new random pair. We found their
404 analysis to be useful to understand also smaller syntactic differences.

405 We then employed qualitative analysis. We manually classified the charac-
406 teristics in the clone pairs that differed and, thereby, led to being not detected as
407 type-1–3 clone. This classification work was done in pairs of researchers in three
408 day-long workshops in the same room. It helped us to discuss the categories
409 and keep them consistent. The result is a set of categories of characteristics that
410 describe the differences. We added quantitative analysis to it by also counting
411 how many of the sampled clone pairs have characteristics of the found types.

412 After the creation of the categories we also assessed the degree of difference
413 (high, medium, or low) per category. From the discussion of the categories, we
414 discovered that this gave us a comprehensive yet precise way to assign clone
415 pairs to the categories. Furthermore, it gave us additional possibilities for a
416 quantified analysis. First, we wanted to understand better how we categorised
417 and assessed the degrees of difference as well as answer H3. Therefore, we per-
418 formed correlation analysis on them. We chose Kendall's tau as the correlation
419 coefficient and tested all correlations on the 0.05 level.

420 For answering H4, we performed a multivariate analysis of variance (MANOVA)
421 which allows more than one dependent variable to be used. Here, our depen-
422 dent variables are the degrees of difference and the independent variable is the
423 programming language. In this analysis, we have a balanced design because we
424 ignored the category *OO design* which was only applicable to Java programs.
425 We use the Pillar-Bartlett statistic for evaluating statistical significance. We
426 checked H4 also on the 0.05 level.

427 These categories with frequencies as well as the results of the hypothesis
428 tests answered RQ 2.

429 3.6.3 Running a Type-4 Detector (RQ 3)

430 As this part of the study is only for exploratory purposes, we focused on the
431 recall of CCCD in the FSCs. As all solutions contain a *main* function, we
432 expected it to find each *main*-pair as clone. We calculate the recall as the
433 number of detected clone pairs by the sum of all clone pairs. A perfect clone
434 detection tool would detect all solutions from a solution set as clones.

435 3.6.4 Creating a Benchmark (RQ 4)

436 After the categorisation to answer RQ 2, we had a clear picture of the various
437 differences. Therefore, we could select representative examples of each difference
438 for each programming language and put them into our new benchmark. To check
439 that the clone pairs cannot be detected by the tools, we run the tools again on
440 the benchmark. If one of the tools still detected a clone, we would replace the
441 clone pair by another representative example until no clones are detected.

442 We created the benchmark by choosing clone pairs consisting of two source
443 code files out of the same solution set. The two files therefore solve the same

444 problem. We selected three pairs where the difference between the files belong
445 to that category for each of the categories we created by answering RQ 2. We
446 chose three pairs for all of the three levels of difference. The other categories of
447 the pairs are very low, ideally zero. Additionally, we added one extra clone pair
448 with extreme differences in all categories.

449 Preferably, we would provide the source code of the chosen solutions directly
450 all in one place. Yet, the copyright of these solutions remains with their authors.
451 Therefore, we provide source files following the same structure as the original
452 files but not violating the copyright.

453 A final set of clone pairs that are not detected as full clones by any of the
454 tools constitutes the benchmark and answered RQ 4.

455 **3.7 Validity Procedure**

456 To avoid selection bias, we performed random sampling where possible. We
457 randomly selected the solution sets and solutions that we use as study objects.
458 In addition, before we manually analysed the category of syntactically different
459 clone pairs, we chose random samples of clone pairs.

460 To avoid errors in our results, we manually checked for false positives and
461 clone types with samples of clones in the solution sets. Furthermore, by working
462 in pairs during all manual work, we controlled each other and detected problems
463 quickly. Overall, the manual inspection of 70 clone pairs for RQ 2 also was a
464 means to detect problems in the detection tools or our instrumentation.

465 For the manual categorisation, we started by categorising 30 syntactically
466 different clone pairs to freely create the categories of undetected clone pairs.
467 Afterwards, we discussed the results among all researchers to come to a unified
468 and agreed categorisation. The actual categorisation of clone pairs was then
469 performed on a fresh sample. Additionally, we performed an independent cate-
470 gorisation of a sample of 10 categorised clone pairs and calculated the inter-rater
471 agreement using Cohen's kappa.

472 **4 Analysis and Results**

473 We structure the analysis and results along our research questions. All quanti-
474 tative and qualitative results are also available in [Wagner et al., 2014].

475 **4.1 Share of Syntactic Similarity (RQ 1)**

476 We summarised the results of the calculated shares for fully and partially syn-
477 tactically similar clone pairs in Tab. 6. We divided the results by programming
478 languages, detection tools and detected clone types. The results differ quite
479 strongly from tool to tool but only slightly between the programming languages.
480 The average syntactic similarities and the standard deviations (SD) are all very
481 low. ConQAT detects more full and partial clones in clone pairs.

Table 6: Full and partial syntactic similarity (in %)

Lang.	Tool	Partially similar				Fully similar			
		Type 1-2		Type 1-3		Type 1-2		Type 1-3	
		Mean	SD	Mean	SD	Mean	SD	Mean	SD
Java	ConQAT	6.36	0.05	11.53	0.07	0.00	0.00	0.00	0.00
	Deckard	0.33	0.00	0.87	0.01	0.00	0.00	0.00	0.00
	Mean	3.35	0.03	6.11	0.04	0.00	0.00	0.00	0.00
C	ConQAT	5.24	0.09	11.48	0.13	1.30	0.00	1.73	0.00
	Deckard	0.28	0.00	1.44	0.01	0.01	0.00	0.01	0.00
	Mean	1.82	0.00	4.32	0.06	0.47	0.00	0.58	0.00
Grand mean		2.45	0.04	5.07	0.04	0.26	0.00	0.35	0.00

482 Table 7 shows the ANOVA results which we need for answering hypotheses
 483 H1 and H2. As our experiment is unbalanced, we use the Type II sum of squares.
 484 This is possible because we found no significant interactions between the factors
 485 in any of the ANOVA results.

486 The results give us no single evaluation of the hypotheses H1 and H2. We
 487 have to differentiate between partial and full syntactic similarity. For the par-
 488 tial similarity, we consistently see a significant difference in the variation in the
 489 detection tools but not in the programming languages. Hence, for partial clone
 490 similarity, we corroborate H1 that there is no difference in recall between pro-
 491 gramming languages. Yet, we reject H2 in favour of the alternative hypothesis
 492 that there is a difference in the similarity share between different tools. For
 493 full similarity, we reject H1 in favour of the alternative hypothesis that there is
 494 a difference between the programming languages. Instead, we accept H2 that
 495 there is no difference between the detection tools.

496 How can we interpret these results? *The overall interpretation is that share*
 497 *of syntactic similarity in FSCs is very small.* There seem to be many possibili-
 498 ties to implement a solution for the same problem with very different syntactic
 499 characteristics. When we only look at the full syntactic similarity, the results
 500 are negligible. Both tools detect none in Java and only few clone pairs for C.
 501 Hence, the difference between the tools is marginal. The difference is significant
 502 between C and Java, however, because we found no full clone pairs in Java. As
 503 we saw in manual inspection, the full detection is easier in C if the developers
 504 implement the whole solution in one main function.

505 For partial syntactic similarity, we get higher results but still stay below
 506 12 %. Hence, for almost 90 % of the clone pairs, we do not even detect smaller
 507 similarities. We have no significant difference between the languages but the
 508 tools. ConQAT has far higher results than Deckard in the type-1-3 clones. The
 509 distinct detection algorithms seem to make a difference here. For the further

Table 7: ANOVA results for variation in recalls (Type II sum of squares, * denotes a significant result)

Partial type 1-2	Sum of Squares	F value	Pr(>F)	
Language	0.0005	0.2352	0.6294	
Tool	0.0491	12.2603	$3 \cdot 10^{-5}$	*

Partial type 1-3	Sum of Squares	F value	Pr(>F)	
Language	0.0010	0.0210	0.8853	
Tool	0.1884	20.5846	$1 \cdot 10^{-7}$	*

Full type 1-2	Sum of Squares	F value	Pr(>F)	
Language	$1 \cdot 10^{-7}$	7.8185	0.0072	*
Tool	$2 \cdot 10^{-8}$	1.1566	0.2871	

Full type 1-3	Sum of Squares	F value	Pr(>F)	
Language	$2 \cdot 10^{-7}$	7.7757	0.0074	*
Tool	$5 \cdot 10^{-8}$	1.9439	0.1692	

510 analysis, we accept an FSC as syntactically similar if one of the tool detected
 511 it.

512 4.2 Categories of Differences (RQ 2)

513 Initially, we created 18 detailed categories. In our qualitative analysis and dis-
 514 cussions, we finally reduced them to five main categories of characteristics de-
 515 scribing the differences between the solutions in a clone pair. The five categories
 516 are *algorithm*, *data structure*, *object-oriented design*, *input/output* and *libraries*.
 517 We could assign each of the 18 initial categories there and realised that we can
 518 assign them to different *degrees of difference*. Therefore, we ended up with a
 519 categorisation including an ordinal quantification of the degree of difference with
 520 the levels *low*, *medium* and *high*. The overall categorisation is shown in Fig. 1.
 521 The centre of the dimensions would be a type-1 clone. The further out we go
 522 on each dimension, the larger the difference.

523 To make the categories and degrees of difference clearer, we give examples of
 524 characteristics in clone pairs in Tab. 8 that led us to classify them in the specific
 525 degrees of difference. The guiding principle was how much effort (in terms of
 526 edit operations) it would be to get from one solution to the other.

527 The two main aspects in any program are its *algorithms* and its *data struc-*
 528 *tures*. This is reflected in our two main categories. Our corresponding degrees
 529 of difference reflect that there might be algorithms that are almost identical
 530 with e.g. only a *switch* instead of nested *if* statements up to completely dif-
 531 ferent solutions, e.g. iterative vs. recursive. Similarly, in data structures, we

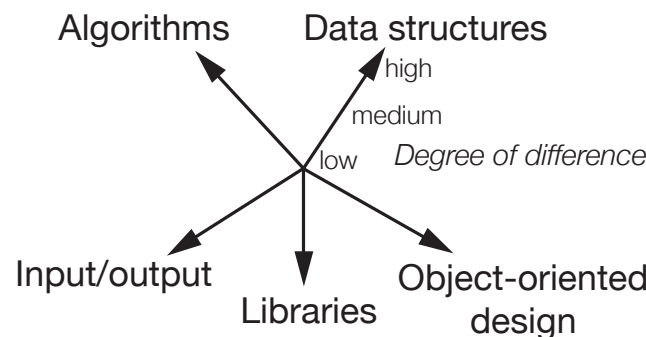


Figure 1: The categories of characteristics of differences between clone pairs

532 can have very simple type substitutions which change the behaviour but are still
 533 functionally very similar (e.g. from *int* to *long*) but also completely user-defined
 534 data types with strong differences.

535 Related to data structures is the category *OO design*. We made this a
 536 separate category because it only applies to OO languages and it had a particular
 537 kind of occurrence in the programs we inspected. Some developers tended to
 538 write Java programs like there were no object-oriented features while others
 539 created several classes and used their objects.

540 As our programming environments and languages are more and more defined
 541 by available *libraries*, this was also reflected in the differences between solutions.
 542 If one developer of a solution knew about a library with existing functionality
 543 needed, and the other developer implemented it herself, this created code that
 544 looks strongly different but can have similar functionality.

545 Finally, maybe a category that arose because the programming contest did
 546 not specify if the input and output should come from a console or a file was
 547 the usage of *I/O*. Nevertheless, we think that this might also be transferable to
 548 other FSCs and contexts because we might be interested in functionally similar
 549 code even if one program writes the output on a network socket while the other
 550 writes into a file.

551 Table 9 shows descriptive statistics for the categories in our sample of unde-
 552 tected clone pairs. The column *Share* shows the ratio of clone pairs with a degree
 553 of difference higher than 0 in relation to all clone pairs in that language. The
 554 median and median absolute deviation (MAD) give the central tendency and
 555 dispersion of the degrees in that category. For that, we encoded *no difference*
 556 = 0, *low* = 1, *medium* = 2 and *high* = 3.

557 All categories occur in the majority of clone pairs. The categories *algorithm*
 558 and *libraries* even occur in about three quarters of the clone pairs. The occur-
 559 rence of categories is consistently smaller in C than in Java. The medians are
 560 mostly low but with a rather large deviation. Only *input/output* in C has a median
 561 of 0. This is consistent with our observation during the manual inspection

Table 8: Examples for the levels in the degree of difference per category

Algorithm	<i>low</i>	Only syntactic variations
	<i>medium</i>	Similarity in the control structure but different method structure
	<i>high</i>	No similarity
Data structure	<i>low</i>	Different data types, e.g. int – long
	<i>medium</i>	Related data types with different interface, e.g. array vs. List
	<i>high</i>	Standard data types vs. own data classes or structs
OO design	<i>low</i>	Only one/few static methods vs. object creation
	<i>medium</i>	Only one/few static methods vs. data classes or several methods
	<i>high</i>	Only one/few static methods vs. several classes with methods
Library	<i>low</i>	Different imported/included but not used libraries
	<i>medium</i>	Few different libraries or static vs. non-static import
	<i>high</i>	Many different or strongly different libraries
I/O	<i>low</i>	Writing to file vs. console with similar library
	<i>medium</i>	Strongly different library, e.g. Scanner vs. FileReader
	<i>high</i>	Strongly different library and writing to file vs. console

562 that I/O is done similarly in the C programs.

563 For evaluating H3, we calculated Kendall's correlation coefficients for all
 564 combinations of categories. The results are shown in Tab. 10. The statistical
 565 tests for these correlations showed significant results for all the coefficients.
 566 Therefore, we need to reject H3 in favour of the alternative hypothesis that there
 567 are correlations between the degrees of difference between different categories.

568 Finally, for evaluating H4, we show the results of the MANOVA in Tab. 11.
 569 We can reject H4 in favour of the alternative hypothesis that there is a difference
 570 between the degrees of difference between the programming languages. This is
 571 consistent with the impression from the descriptive statistics in Tab. 9.

572 In summary, we interpret these results such that *there are differences in*
 573 *FSC pairs in their algorithms, data structures, input/output and used libraries.*
 574 *In Java, there are also differences in the object-oriented design.* On average,
 575 *these differences are mostly small but the variance is high.* Hence, we believe
 576 that with advances in clone detectors for tolerating the smaller differences, there
 577 could be large progress in the detection of FSCs. Yet, there will still be many
 578 medium to large differences. We also saw that the programming languages vary
 579 in the characteristics of undetected difference. Therefore, it might be easier to
 580 overcome those differences in non-object-oriented languages, such as C, than in

Table 9: Descriptive statistics of degrees of difference over categories and programming languages

Lang.	Category	Share	Median	MAD
Java	Algorithm	96 %	3	0.0
	Libraries	86 %	1	1.5
	I/O	83 %	2	1.5
	Data structure	72 %	1	1.5
	OO design	71 %	1	1.5
C	Algorithm	76 %	2	1.5
	Libraries	73 %	1	1.5
	Data structure	66 %	1	1.5
	I/O	38 %	0	0.0
Total	Algorithm	86 %	2	1.5
	Libraries	79 %	1	1.5
	OO design	71 %	1	1.5
	Data structure	69 %	1	1.5
	I/O	60 %	1	1.5

581 object-oriented languages which offer even more possibilities to express solutions
 582 for the same problem. Yet, we were impressed by the variety in implementing
 583 solutions in both languages during our manual inspections.

584 Our categories are significantly correlated with each other. This can mean
 585 that there might be other, independent categories with less correlation. Never-
 586 theless, we believe the categories are useful because they describe major code
 587 aspects in a way that is intuitively understandable to most programmers. It
 588 would be difficult to avoid correlations altogether. For example, a vastly differ-
 589 ent data structure will always lead to a very different algorithm.

590 4.3 Type-4 Detection (RQ 3)

591 Table 12 shows the recall of fully and partially detected clone pairs in our sample.
 592 CCCD has a considerable recall for partial clones in the clone pairs of about
 593 16 %. It does, however, detect almost none of the clone pairs a full clones.
 594 The overlap with ConQAT and Deckard, and therefore type-1–3 clones, is tiny
 595 (0.05 % of the recall).

596 We interpret this result such that also contemporary type-4 detection tools
 597 have still problems detecting real-world FSCs and to handle the differences we
 598 identified in RQ 2.

Table 10: Correlation matrix with Kendall's correlation coefficient for the category degrees (all are significant)

	Algo.	Data struct.	OO design	I/O	Libraries
Algorithm	1.00	0.38	0.44	0.15	0.31
Data struct.	0.38	1.00	0.26	0.25	0.21
OO design	0.44	0.26	1.00	0.29	0.39
I/O	0.15	0.25	0.29	1.00	0.27
Libraries	0.31	0.21	0.39	0.27	1.00

Table 11: MANOVA results for variation in degree of differences (Type I sum of squares, * denotes a significant result)

	Pillai-Bartlett	approx. F	Pr(>F)	
Language	0.1513	6.0196	0.0002	*

599 4.4 Benchmark (RQ 4)

600 The number of study objects used in our analysis is quite high. As described
 601 above, we examined 1,400 Java files and 1,400 C files. For many demonstrations
 602 and clone detection tool analyses a much smaller file set is sufficient. We call
 603 this smaller set of files *benchmark*.

604 The first half of the benchmark we provide consists of 29 clone pairs. For
 605 Java, we include 16 clone pairs. The set of clone pairs we provide for C is
 606 structured in exactly the same way as the Java samples except that we do not
 607 have the three clone pairs that differ only in object-oriented design. Therefore,
 608 we do not have 16 samples here but 13 which make the 29 clone pairs for both
 609 languages.

610 Figure 2 shows a rating of an example clone pair in the benchmark set where
 611 the two files only differ significantly in the kind of input/output, but not in the
 612 other categories.

613 We provide this distribution of clone pairs for both partial clones and full

Table 12: Full and partial clone recall means over solution sets for CCCD (in %)

	Mean	SD
Partial	16.03	0.07
Full	0.10	0.00

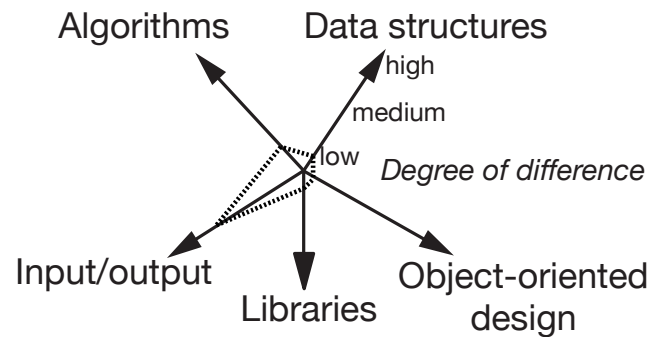


Figure 2: Example category rating of a clone pair in the benchmark set

```

public static void main(String [] args) {
    reader = new BufferedReader(
        new FileReader("A-large.in"));
    writer = new PrintWriter("a.out");
}

public static void main(String [] args) {
    File file = new File(System.in);
    try (Scanner scanner = new Scanner(
        new FileReader(file))) {
        File out = new File(System.out);
        try (PrintWriter writer = new ...

```

Figure 3: Example of a high difference in the category Input/Output

614 clones. Hence, the total number of clone pairs within the benchmark is 58.
 615 Figure 4 shows an overview of the structure of the whole benchmark set. This
 616 structure enables developers of a clone detection tool to test their tool easily as
 617 well analyse the nature of the clones found and not found by a tool.

618 Our benchmark provides several advantages to the research community.
 619 First, *developers of a clone detection tool can easily test their tool with the*
 620 *source files as input.* They can see whether their tool detects the clones or they
 621 can analyse why it did not. Second, the clones in the benchmark are easily
 622 understandable examples for the categories we created. Third, *the clones in*
 623 *the benchmark were not built artificially; the solutions were implemented inde-*
 624 *pendently by at least two persons during the Code Jam contest.* Despite our
 625 modifications to avoid copyright problems, neither changing structure nor algo-
 626 rithm, the code clones are more realistic than fully artificial copies where one
 627 file is modified as part of a study.

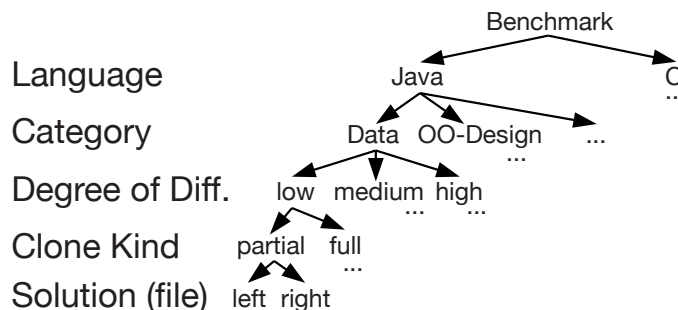


Figure 4: Structure of the benchmark set (overview)

628 5 Threats to Validity

629 We analyse the validity threats for this study following common guidelines for
 630 empirical studies [Yin, 2003, Wohlin et al., 2012].

631 5.1 Conclusion Validity

632 As most of our measurements and calculations were performed automatically,
 633 the threats to conclusion validity are low. For the corresponding hypothesis
 634 tests, we checked all necessary assumptions. Only the classification and rating
 635 of the degree of difference is done manually and, hence, could be unreliable.
 636 We worked in pairs to reduce this threat. Furthermore, one of the researchers
 637 performed an independent classification of a random sample of 10 clone pairs
 638 to compare the results. We calculated Cohen's kappa for the categories of
 639 differences between clone pairs as presented in Table 13.

640 We interpret the kappa results according to the classification by Landis and
 641 Koch [Landis and Koch, 1977]. Hence, our results are a moderate agreement

Table 13: Kappa values for difference categories

Category	Kappa
Data structures	0.41
OO design	0.35
Algorithms	0.47
Libraries	0.36
Input/Output	0.47

642 between the categories: data structures, algorithms and input/output. For the
643 categories object-oriented design and libraries we have a fair agreement. We
644 consider this to be reliable enough for our investigations.

645 5.2 Internal Validity

646 There is the threat that the implementation of our instrumentation tooling may
647 contain faults and, therefore, compute incorrect results for the detected clones
648 and recalls. We reduced this threat inherently by the manual inspections done
649 to answer RQ 2 and independently to investigate the type-4 clones.

650 A further threat to internal validity is that we took our solution sets from
651 Google Code Jam. We cannot be sure that all the published solutions of the
652 Code Jam within a solution set are actually functionally similar. We rely on
653 the fact that the organisers of the Code Jam must have checked the solutions to
654 rank them. Furthermore, we assume to have noticed in the manual inspections
655 if there were solutions in a solution set with highly differing functionality.

656 5.3 Construct Validity

657 To fully understand the effectiveness of a clone detection approach, we need to
658 measure precision as well as recall. In our study, we could not measure precision
659 directly because of the large sample size. We checked for false positives during
660 the manual inspections and noted only few rather short clones. Our minimal
661 clone length is below recommended thresholds. This is a conservative approach
662 to the problem. By that we will find more clones than in an industrial approach.
663 We decided to use this threshold to be sure that we cover all the interesting clone
664 pairs that would be lost due to variation in the precision of the tools.

665 There is a threat because we count each clone pair only once. In partial
666 clones, one clone pair might contain a type-2 as well as a type-3 partial clone.
667 In those cases, we decided that the lower type – the easier detection – should be
668 recorded. Hence, the assignment to the types might be imprecise. We accept
669 this threat as it has no major implication for the conclusions of the experiment.

670 5.4 External Validity

671 There is also a threat to external validity in our usage of solutions from Google
672 Code Jam. The submitted programs might not represent industrial software
673 very well. Participants had a time limit set for turning in their solutions. Fur-
674 thermore, the programming problems contained mostly reading data, perform-
675 ing some calculations on it and writing data. This might impact the method
676 structure within the solutions. This threat reduces the generalisability of our
677 results. Yet, we expect that other, more complex software will introduce new
678 kinds of difference categories (e.g. differences in GUI code) and only extend but
679 not contradict our results.

680 For the study, we chose three well-known and stable clone detection tools.
681 Two of them analyse Java and C programs detecting type-1 to type-3 clones.
682 The third one detects type 4 clones and supports only programs written in C
683 and only finds clones in complete functions. Overall, we are confident that these
684 tools represent the available detection tools well.

685 6 Conclusions and Future Work

686 In this paper, we investigated the characteristics of clones not created by copy&paste.
687 We base our study on [Juergens et al., 2010b], but this is the first study with pro-
688 grams implementing different specifications in diverse programming languages
689 including CCCD as concolic clone detector for type-4 clones. We found that a
690 full syntactic similarity was detected in less than 1 % of clone pairs. Even partial
691 syntactic similarity was only visible in less than 12 %. The concolic approach
692 of CCCD can detect FSCs without syntactic similarity as type-4 clones. Yet, a
693 full detection was only possible in 0.1 % of clone pairs.

694 Our categorisation of the differences of clone pairs not syntactically similar
695 showed that usually several characteristics make up these differences. On av-
696 erage, however, the differences were mostly small. Hence, we believe there is a
697 huge opportunity to get a large improvement in detection capabilities of type-4
698 detectors even with small improvements in tolerating additional differences. We
699 provide a carefully selected benchmark with programs representing real FSCs.
700 We hope it will help the research community to make these improvements.

701 6.1 Relation to Existing Evidence

702 We can most directly relate our results to Juergens, Deissenboeck and Hum-
703 mel [Juergens et al., 2010b]. We support their findings that using type-1-3
704 detectors, below 1 % is fully and below 10 % is partially detected. We can
705 add that with the type-4 detection of CCCD, the partial clone recall can reach
706 16 %. They introduce categories which were derived from other sources but not
707 created them with a systematic qualitative analysis. Yet, there are similarities
708 in the categories.

709 Their category *syntactic variation* covers “if different concrete syntax con-
710 structs are used to express equivalent abstract syntax”. We categorised this

711 as small algorithm difference. Their category *organisational variation* “occurs
712 if the same algorithm is realized using different partitionings or hierarchies of
713 statements or variables”. We categorise these differences as a medium algorithm
714 difference. Their category *delocalisation* “occurs since the order of statements
715 that are independent of each other can vary arbitrarily between code fragments”
716 is covered as difference in algorithm in our categorisation. Their category *gen-
717 eralisation* “comprises differences in the level of generalization” which we would
718 cover under *object-oriented design*. They also introduce *unnecessary code* as
719 category with the example of a debug statement. We did not come across such
720 code in our sample but could see it as potential addition.

721 Finally, they clump together *different data structure and algorithm* which
722 we categorised into separate categories. We would categorise these variations
723 as either data structure or algorithm differences with probably a high degree
724 of difference. They found that 93 % of their clone pairs had a variation in
725 the category *different data structure or algorithm*. We cannot directly support
726 this value but the tendency. We found that 91 % of the inspected clone pairs
727 had a difference at least in either *algorithm* or *data structure* and especially for
728 algorithm the difference was on average large.

729 Tiarks, Koschke und Falke [Tiarks et al., 2011] created a categorisation for
730 differences in type-3 clones. Therefore, their focus was on classifying syntactic
731 differences that probably hail from independent evolution of initially copied
732 code. Yet, the larger the differences, the more their categories are similar to
733 ours. For example, they abstract edit operations to *type substitution* or *different
734 algorithms*. We believe, however, that our categorisation is more useful for FSCs
735 and to improve clone detection tools along its lines.

736 6.2 Impact

737 Independently developed FSCs have very little syntactic similarity. Therefore,
738 type-1–3 clone detectors will not be able to find them. Newer approaches,
739 such as CCCD, can find FSCs but their effectiveness still seems limited. Hence
740 more research in approaches more independent of syntactic representations is
741 necessary. We will need to find ways to transfer the positive results of Jiang
742 and Su [Jiang and Su, 2009] with the Linux kernel to other languages and
743 environments while overcoming the challenges in such dynamic detections as
744 discussed, for example, in Deissenboeck et al. [Deissenboeck et al., 2012]. We
745 hope our benchmark will contribute to this.

746 6.3 Limitations

747 The major limitation of our study is that we did not use a wide variety of
748 types of programs that exist in practice. The programs from Google Code Jam
749 all solve structurally similar problems, for example, without any GUI code.
750 We expect, however, that such further differences would rather decrease the
751 syntactic similarity even more. The categories might have to be extended to
752 cover these further differences. Nevertheless, the investigated programs were all
753 developed by different programmers and are not artificial.

754 Furthermore, we had to concentrate on three clone detectors and two pro-
755 gramming languages. Other tools and languages might change our results but
756 we are confident that our selection is representative of a large class of detectors
757 and programming languages.

758 **6.4 Future Work**

759 We plan to investigate the differences between the tools and the detected clone
760 pairs of different types in more detail. In particular, we would like to work
761 with researchers who have built type-4 detectors to test them against our clone
762 database and to inspect the found and not found clones.

763 **Acknowledgment**

764 The authors would like to thank Benjamin Hummel, Lingxiao Jiang and Daniel
765 Krutz for their help in getting their tools to work and Kornelia Kuhle for feed-
766 back on the text.

767 **References**

- 768 [Bellon et al., 2007] Bellon, S., Koschke, R., Antoniol, G., Krinke, J., and
769 Merlo, E. (2007). Comparison and evaluation of clone detection tools. *IEEE*
770 *Transactions on Software Engineering*, 33(9):577–591.
- 771 [Deissenboeck et al., 2012] Deissenboeck, F., Heinemann, L., Hummel, B., and
772 Wagner, S. (2012). Challenges of the dynamic detection of functionally similar
773 code fragments. In *Proc. 16th European Conference on Software Maintenance*
774 *and Reengineering (CSMR)*, pages 299–308. IEEE.
- 775 [Deissenboeck et al., 2008] Deissenboeck, F., Juergens, E., Hummel, B., Wag-
776 ner, S., y Parareda, B. M., and Pizka, M. (2008). Tool support for continuous
777 quality control. *IEEE Software*, 25(5):60–67.
- 778 [Gabel et al., 2008] Gabel, M., Jiang, L., and Su, Z. (2008). Scalable detec-
779 tion of semantic clones. In *Proc. 30th International Conference on Software*
780 *Engineering (ICSE '08)*, pages 321–330. ACM.
- 781 [Jedlitschka and Pfahl, 2005] Jedlitschka, A. and Pfahl, D. (2005). Reporting
782 guidelines for controlled experiments in software engineering. In *Proc. 4th*
783 *International Symposium on Empirical Software Engineering (ISESE)*. IEEE.
- 784 [Jiang et al., 2007a] Jiang, L., Mishergchi, G., Su, Z., and Glondu, S. (2007a).
785 Deckard: Scalable and accurate tree-based detection of code clones. In *Proc.*
786 *29th International Conference on Software Engineering (ICSE)*, pages 96–
787 105. IEEE.
- 788 [Jiang and Su, 2009] Jiang, L. and Su, Z. (2009). Automatic mining of func-
789 tionally equivalent code fragments via random testing. In *Proc. Eighteenth In-*
790 *ternational Symposium on Software Testing and Analysis (ISSTA '09)*, pages
791 81–92. ACM.
- 792 [Jiang et al., 2007b] Jiang, L., Su, Z., and Chiu, E. (2007b). Context-based
793 detection of clone-related bugs. In *Proc. 6th Joint Meeting of the European*
794 *Software Engineering Conference and the ACM SIGSOFT Symposium on The*
795 *Foundations of Software Engineering (ESEC/FSE)*, pages 55–64. ACM.
- 796 [Juergens et al., 2010a] Juergens, E., Deissenboeck, F., Feilkas, M., Hummel,
797 B., Schaetz, B., Wagner, S., Domann, C., and Streit, J. (2010a). Can clone
798 detection support quality assessments of requirements specifications? In
799 *Proc. ACM/IEEE 32nd International Conference on Software Engineering*
800 *(ICSE'10)*, pages 79–88. ACM.
- 801 [Juergens et al., 2010b] Juergens, E., Deissenboeck, F., and Hummel, B.
802 (2010b). Code similarities beyond copy & paste. In *Proc. 14th European*
803 *Conference on Software Maintenance and Reengineering (CSMR)*, pages 78–
804 87. IEEE.

- 805 [Juergens et al., 2009] Juergens, E., Deissenboeck, F., Hummel, B., and Wag-
806 ner, S. (2009). Do code clones matter? In *Proc. 31st International Conference*
807 *on Software Engineering (ICSE'09)*, pages 485–495. IEEE.
- 808 [Kamiya et al., 2002] Kamiya, T., Kusumoto, S., and Inoue, K. (2002).
809 CCFinder: A multilinguistic token-based code clone detection system for large
810 scale source code. *IEEE Transactions on Software Engineering*, 28(7):654–
811 670.
- 812 [Kim et al., 2011] Kim, H., Jung, Y., Kim, S., and Yi, K. (2011). MeCC: Mem-
813 ory comparison-based clone detector. In *Proc. 33rd International Conference*
814 *on Software Engineering (ICSE '11)*, pages 301–310. ACM.
- 815 [Komondoor and Horwitz, 2001] Komondoor, R. and Horwitz, S. (2001). Us-
816 ing slicing to identify duplication in source code. In *Proc. 8th International*
817 *Symposium on Static Analysis (SAS'01)*, pages 40–56. Springer.
- 818 [Koschke, 2007] Koschke, R. (2007). Survey of research on software clones. In
819 *Dagstuhl Seminar Proc. Duplication, Redundancy, and Similarity in Software*.
- 820 [Krinke, 2001] Krinke, J. (2001). Identifying similar code with program depen-
821 dence graphs. In *Proc. Eighth Working Conference on Reverse Engineering*
822 *(WCRE)*, pages 301–309. IEEE.
- 823 [Krutz and Shihab, 2013] Krutz, D. E. and Shihab, E. (2013). CCCD: Con-
824 colic code clone detection. In *Proc. 20th Working Conference on Reverse*
825 *Engineering (WCRE'13)*. IEEE.
- 826 [Lakhota et al., 2003] Lakhota, A., Li, J., Walenstein, A., and Yang, Y.
827 (2003). Towards a clone detection benchmark suite and results archive. *11th*
828 *IEEE International Workshop on Program Comprehension*, pages 285–286.
- 829 [Landis and Koch, 1977] Landis, R. J. and Koch, G. G. (1977). The measure-
830 ment of observer agreement for categorical data. *Biometrics*, 33(1):159–74.
- 831 [Marcus and Maletic, 2001] Marcus, A. and Maletic, J. I. (2001). Identification
832 of high-level concept clones in source code. In *Proc. 16th Annual International*
833 *Conference on Automated Software Engineering (ASE 2001)*, pages 107–114.
834 IEEE.
- 835 [Rattan et al., 2013] Rattan, D., Bhatia, R., and Singh, M. (2013). Software
836 clone detection: A systematic review. *Information and Software Technology*,
837 55(7):1165–1199.
- 838 [Roy and Cordy, 2007] Roy, C. K. and Cordy, J. R. (2007). A survey on soft-
839 ware clone detection research. Technical Report 2007-541, Queen’s University,
840 Kingston, Canada.
- 841 [Roy et al., 2009a] Roy, C. K., Cordy, J. R., and Koschke, R. (2009a). Compari-
842 son and evaluation of code clone detection techniques and tools: A qualitative
843 approach. *Science of Computer Programming*, 74(7):470–495.

- 844 [Roy et al., 2009b] Roy, C. K., Cordy, J. R., and Koschke, R. (2009b). Compari-
845 son and evaluation of code clone detection techniques and tools: A qualitative
846 approach. *Sci. Comput. Program.*, 74(7):470–495.
- 847 [Svajlenko et al., 2014] Svajlenko, J., Islam, J. F., Keivanloo, I., Roy, C. K., and
848 Mia, M. M. (2014). Towards a big data curated benchmark of inter-project
849 code clones. In *Proc. International Conference on Software Maintenance and*
850 *Evolution (ICSME'14)*, pages 476–480. IEEE.
- 851 [Svajlenko and Roy, 2014] Svajlenko, J. and Roy, C. K. (2014). Evaluating
852 Modern Clone Detection Tools. In *Proceedings of the 30th International Con-*
853 *ference on Software Maintenance and Evolution*, pages 321–330. IEEE.
- 854 [Tempero, 2013] Tempero, E. (2013). Towards a curated collection of code
855 clones. In *Proc. 7th International Workshop on Software Clones (IWSC'13)*,
856 pages 53–59. IEEE.
- 857 [Tiarks et al., 2011] Tiarks, R., Koschke, R., and Falke, R. (2011). An extended
858 assessment of type-3 clones as detected by state-of-the-art tools. *Software*
859 *Quality Journal*, 19(2):295–331.
- 860 [Wagner, 2013] Wagner, S. (2013). *Software Product Quality Control*. Springer.
- 861 [Wagner et al., 2014] Wagner, S., Abdulkhaleq, A., Bogicevic, I., Ostberg, J.-
862 P., and Ramadani, J. (2014). Detection of functionally similar code clones:
863 Data, analysis software, benchmark. DOI 10.5281/zenodo.12646.
- 864 [Wohlin et al., 2012] Wohlin, C., Runeson, P., Höst, M., Ohlsson, M. C., Reg-
865 nell, B., and Wesslén, A. (2012). *Experimentation in Software Engineering*.
866 Springer.
- 867 [Yin, 2003] Yin, R. K. (2003). *Case Study Research: Design and Methods*.
868 Applied Social Research Methods. SAGE Publications.