

The Fault in Our Stars

Designing Reproducible Large-scale Code Analysis Experiments

Petr Maj

Czech Technical University

Stefanie Muroya

Institute of Science and Technology

Konrad Siek

Czech Technical University

Luca Di Grazia

Università della Svizzera Italiana

Jan Vitek

Charles University and Northeastern University

Abstract

Large-scale software repositories are a source of insights for software engineering. They offer an unmatched window into the software development process at scale. Their sheer number and size holds the promise of broadly applicable results. At the same time, that very size presents practical challenges for scaling tools and algorithms to millions of projects. A reasonable approach is to limit studies to representative samples of the population of interest. Broadly applicable conclusions can then be obtained by generalizing to the entire population. The contribution of this paper is a standardized experimental design methodology for choosing the inputs of studies working with large-scale repositories. We advocate for a methodology that clearly lays out what the population of interest is, how to sample it, and that fosters reproducibility. Along the way, we discourage researchers from using extrinsic attributes of projects such as stars, that measure some unclear notion of popularity.

2012 ACM Subject Classification Software and its engineering → Ultra-large-scale systems;

Keywords and phrases software, mining code repositories, source code analysis

Digital Object Identifier 10.4230/LIPIcs.ECOOP.2023.23

1 Introduction

And so it begins...

count the number of stars associated with each repository. The number of stars relate to how many people are interested in that project. Thus, we assume that stars indicate the popularity of a project. We select the top 50 projects in each language

Sentences like these appear in the methodology sections of software engineering papers, with sometimes, little more in terms of experimental design. This paper aims to convince readers that using extrinsic features of projects, such as popularity, may limit applicability of results of the studies relying on them. Instead one should select projects based on their intrinsic features and spell out expectations as well as threats to validity.

Empirical software engineering studies are experiments performed on a corpus of software to validate some hypotheses. For instance, one could take projects written in various languages and attempt to show that some language feature has an impact on the quality of the code written using it. The value of large-scale corpus study often does not lie in what we learn about the projects that were analyzed, but rather in what these can teach us about the larger



© Anon;

licensed under Creative Commons License CC-BY 4.0

37th European Conference on Object Oriented Programming (ECOOP 2023).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:23

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

43 population. There is limited value in, say, finding out that a new Java language feature is
 44 beneficial in a handful of cases if we cannot generalize that result to a broader portion of the
 45 Java ecosystem.

46 Yet, many papers in the field do not articulate how broadly applicable their results are
 47 expected to be. Even the simple question of how the projects that were analyzed were
 48 selected is not clear. While large-scale code repositories, such as GitHub, are a boon to
 49 the software engineering community, their sheer size requires care. We argue that better
 50 experimental design will strengthen research done in the field.

51 Consider Table 1 which has a meta-study of three years of the *Mining Software Repositories*
 52 conference. Forty-one papers relied on subsets of GitHub. Out of those, five papers lacked
 53 sufficient information about their dataset to determine how they selected their inputs, twenty-
 54 one used GitHub stars to obtain a subset of projects, ten used simple combinations of
 55 attribute thresholds and only five relied on random sampling over the entire population.

papers	class	description	# of projects
5	Unknown	Unknown or proprietary	1–35K
21	Stars	Filter projects using stars	5–2M
10	Other	Other filter for projects	7–290K
5	Random	Filter and sample randomly	6–51K

■ **Table 1** Experimental design in MSR 2019, 2020 and 2021

56 What is the right choice here? None of the papers analyzed the entire ecosystem as that
 57 would mean tens to hundreds of millions of projects. The question thus becomes how to
 58 sample the population of projects hosted on GitHub. In this paper, we criticize the use of
 59 GitHub stars as they are appealing and popular, yet also dangerous. But, really, our point
 60 generalizes to any extrinsic attributes of a project. So, again, what is the right choice? The
 61 answer is, of course, that it depends. The rest of this paper attempts to shed some light on
 62 how to make a reasoned choice of inputs.

63 First, let us return to stars and ponder why they play such an important role in our
 64 experimental methodology? We believe expectations and pragmatics are the explanation.
 65 Community standards are largely set by the papers we publish. The literature codifies
 66 expectations for authors of the next batch of papers. These expectations slowly evolve in
 67 response to reviewer attitudes. So, we use stars because our peers do. And just as importantly
 68 for the pragmatic reason that GitHub does not provide an index of projects, nor does it
 69 allow to query over intrinsic attributes of code. Finding inputs is thus hobbled by limitations
 70 of our tools. Stars play a double role. They are a queryable index of projects as GitHub
 71 does provide an interface to obtain them. They also come with an expectation that starred
 72 projects enjoy some notion of quality [8]. This paper will show that stars do not necessarily
 73 correlate with quality and that they introduce reproduction barriers.

74 We propose a methodology for designing reproducible software engineering experiments
 75 over large-scale repositories with the explicit goal of improving the generalizability of our
 76 results. The methodology is in line with evolving community standards [23] but specific to
 77 large-scale code analysis. We propose to follow the following protocol when designing a new
 78 experiment:

- 79 **1. Population Hypothesis:** Give a brief description of the population of interest the
 80 research should generalize to; it may be narrow such as “programs written by students
 81 learning JavaScript” or as broad as “commercial code”.

- 82 **2. Frame Oracle:** Give a procedure for deciding if a project belongs to the population of
 83 interest. The procedure should be efficiently computable over intrinsic attributes of a
 84 project. An oracle could, say, return projects with a single JavaScript file created by a
 85 user with no previous commits.
- 86 **3. Sampling Strategy:** Describe a strategy for selecting a subset of the entire population.
 87 Ideally, specified algorithmically. An example is random sampling without replacement
 88 from a known seed.
- 89 **4. Validity:** Give an argument as to how the oracle and the sampling strategy are valid
 90 means to obtain representative samples. This can be a discussion of how to check result
 91 quality, such as manual inspection of samples written by beginners, and threats to validity.
- 92 **5. Reproduction Artifacts:** Publish an artifact that reproduces exactly the reported
 93 results, and supports changes to either the input or the details of the experiment.

94 Reproducibility has nuances. Our emphasis is on providing support for the following three use
 95 cases: *Repetitions* which run the reproduction artifact to obtain bit-for-bit equal results. This
 96 is the most stringent use case and often requires a reproduction artifact that bundles code and
 97 inputs. *Reanalysis* alters either the method or its input, it requires an executable artifact and
 98 a method for acquiring new inputs. Finally, *reproductions* are independent implementations
 99 that require the paper to have an unambiguous description of all experimental details.¹
 100 Supporting reproducibility can be greatly simplified with appropriate tooling. Our work
 101 builds on the CodeDJ infrastructure (codedj-prg.github.io). Our contributions are:

- 102 **1. A dataset** of 2Mio+ projects with intrinsic attributes precomputed.
- 103 **2. A characterization of stars** as a means to select inputs for code analysis experiments.
- 104 **3. A methodology** that can be readily adopted to improve reproducibility.
- 105 **4. A reproduction** that highlights challenges to generalization due to project selection.

106 Our community has been moving towards broader adoption of the practice of artifact
 107 evaluation [11]. While artifacts are clearly helpful as they make papers providing them
 108 easier to reproduce, the selection of inputs is often hardwired and not considered part of
 109 the reproduction. The impact of our proposal, if adopted, would be to encourage authors of
 110 large-scale code studies to consider the collection of the inputs to their work to be part of
 111 the experiment and thus make it easy to change the way inputs are selected.

112 **Road map** The structure of the paper is as follows.

- 113 ■ Section 2 begins with a short overview of the state with respect to methodologies for
 114 project selection and tooling to support it.
- 115 ■ Section 3 takes four practical examples, papers published at the Mining Software
 116 Repositories conference, and attempts to couch their experimental design in the terms
 117 introduced above. These example suggest that authors are not always clear about their
 118 intent and strategy. While looking at these papers we found a number of practical
 119 impediments to reproducibility.
- 120 ■ Section 4 describes characteristics of the projects hosted on Github and argues that stars
 121 cannot yield a representative sample of developed projects.
- 122 ■ Section 5 outlines our proposal for how to design large-scale program analysis experiments.
- 123 ■ Section 6 follows our guidelines and attempts to repeat the studies of Section 3 while
 124 perturbing the experimental inputs.

¹ The terminology comes from [24] and was used by SIGPLAN artifact evaluation committees.

23:4 Designing Reproducible Large-scale Code Analysis Experiments

- 125 ■ Section 7 is responsive to reviewers of this paper and their request to reproduce an
126 experiment from a paper with a verified artifact.
- 127 ■ Section 8 concludes and gives some parting thoughts. This paper improves on the state
128 of the art in that it argues for a structured experimental design that relies on tooling for
129 input selection.

130 **2** Related work

131 We review relevant advice, warnings and the state of tooling.

132 **2.1** Community standards

133 A push towards reproducibility is underway. The standards framework of Ralph et al. [23]
134 includes a section on experimental design and specifically on sampling. This is further explored
135 by Baltes and Ralph [1]. They argue that software engineering faces a generalizability crisis.
136 In their meta-analysis of 120 papers, they find that convenience sampling² is widely used
137 to select projects to analyze from a large population. Convenience sampling rarely leads to
138 representative samples, and – without a careful study of potential sources of bias – can lead
139 to conclusions that do not generalize. They explain this state of affairs by a fundamental
140 challenge: the lack of appropriate sampling frames to access elements of the population of
141 interest. Earlier work by Nagappan et al. [19] already attempted to address this problem
142 by defining the notion of sample coverage as a way to assess the quality of the data used as
143 input to an experiment. Even closer to our paper is the study by Cosentino et al. [4] which
144 reported that out of 93 large corpus papers, 63 papers failed to provide replication datasets.
145 Most papers did not use random samples and omitted mentions of any limitations.

146 **2.2** Mining repositories

147 GitHub is a popular data source. Warnings about its perils go back to the work of
148 Kalliamvakou et al. [10] which highlighted “noise” among hosted projects. In particular, they
149 point out, tiny and inactive projects dominate the platform. Lopes et al. [13] poured oil on
150 that fire, showing that up to 95% of the files containing code in some language ecosystems
151 were copies of one another and filtering by stars reduces the proportion of duplicates without
152 eliminating them. One way researchers have strived to find signal in GitHub’s sea of noise is
153 to use stars. But what do stars mean? We would like them to be correlated with quality
154 code, code worth analyzing. Borges and Valente [2] conducted a user survey that found
155 the most common reasons for starring a project was to show appreciation (e.g. *starred this*
156 *repository because it looks nice*) and bookmark it (e.g. *starred it because I wanted to try*
157 *it later*). They also warn against promotional campaigns to drive up ratings. Popularity
158 of projects was studied by Han et al. [8], they found that while users believe stars are a
159 measure of a project’s popularity, intrinsic attributes such as branches, open issues and
160 contributors are better predictors. Expanding on that result, Munaiah et al. [18] propose
161 classifier for *engineered projects*, which they define as projects that leverage sound software
162 engineering principles. They show that the classifier outperforms stars. Pickerill et al. [22]
163 further improved classification with an approach based on time-series clustering.

² The Wikipedia definition of convenience sampling is a type of non-probability sampling that involves the sample being drawn from that part of the population that is close to hand.

164 2.3 Tools for miners

165 A number of infrastructures have been developed to assist researchers in the field. The most
166 ubiquitous was, the now defunct, GHTorrent [6]. The project offered a continuously updated
167 database of metadata about public projects that was a valuable building block for other
168 tools. Boa is complementary as it lets users write sophisticated queries over source code [5].
169 CodeDJ is a newer infrastructure that supports queries over both meta-data and file contents
170 and is language agnostic [15]. Recent works address performance issues of querying at scale
171 [14, 17]. Of these, only CodeDJ ensures reproducible queries.

172 3 State of practice

173 How do people design experiments for large-scale code studies? This section gives some
174 examples that we believe to be representative which we will revisit later when we attempt to
175 reproduce the results with different inputs. For each paper, we provide a brief summary of
176 the scientific claims made by the authors. Then, we attempt, with our best understanding of
177 the work, to reverse engineer a version of the protocol laid out in the introduction. We, thus,
178 give an account of each paper’s population hypothesis, a description of the frame oracle,
179 sampling strategy, validation and reproduction artifacts. We conclude the section with some
180 observations general reproduction issues that show up in these papers.

181 3.1 MSR 2020: What is software

182 “Software” has an intuitive definition, namely code, but there is more. The paper by
183 Pfeiffer [21] classifies the content of repositories in categories such as code, data and
184 documentation. They, then, observe that software is more than just code. Documentation
185 is an integral constituent of software, and software without data is often correlated with
186 libraries, and finally that software without code is rare, but exists. The paper answers
187 the question “*what are the constituents of software and how are they distributed?*” The
188 paper argues that existing definitions of the term are non-descriptive, inconclusive and even
189 contradictory.

190 Population Hypothesis: Implicitly, the population is all inclusive.

191 Frame Oracle: Given the lack of details, we assume all projects on GitHub belong.

192 Sampling Strategy: the authors carry out convenience sampling by choosing popular
193 repositories. Stating “*by popularity we mean the starred criteria with which GitHub users*
194 *express liking similar to likes in social networks.*” Most-starred projects in 25 languages were
195 acquired by executing one query by language, saying that “*without language qualifier, the*
196 *API returns only 1,020 repositories in total, which we decided is not enough for our study.*”

197 Validity: No discussion of relevant issues or threats.

198 Reproducibility Artifacts: A listing of files and repositories is provided along with the code
199 of the classifier and a notebook. Repository contents were not preserved.

200 3.2 MSR 2020: Method chaining

201 In an object-oriented language, a *method chain* occurs when the result of a method invocation
202 is the receiver of a subsequent invocation. In Java, chains manifest as sequences of calls
203 connected by dots. Nakamura et al. [20] analyze trends in usage of method chains and
204 conclude that they increase over a period of eight years.

205 Population Hypothesis: Java projects developed “*by real-world programmers.*” The authors
 206 state that they “*did not apply any filter to the collected repositories. This supports the*
 207 *generalizability of our results.*” The authors also consider generalization beyond Java, saying
 208 “*our results are more likely to be applied to a language that does not provide such a construct*
 209 *(e.g. PHP and JavaScript). The empirical study of this hypothesis is future work.*”

210 Frame Oracle: Implicitly, all Java projects hosted on GitHub.

211 Sampling Strategy: The authors use convenience sampling, taking 2,814 projects that
 212 appeared at least once in the list of the 1K most-starred projects in November 2019. Projects
 213 were deduplicated and filtered for syntactically invalid files.

214 Validity: No discussion of relevant issues.

215 Reproducibility Artifacts: Project metadata and computed chain lengths are available.
 216 Communication with the authors reveals that their reproduction package is not available.

217 3.3 MSR 2019: Style analyzer

218 Each software project seems to develop its own formatting conventions. Markotsev et al. [16]
 219 demonstrate that an unsupervised learning algorithm can automate project-specific code
 220 formatting. They reproduce styles with a high degree of precision for a set of repositories.

221 Population Hypothesis: The authors speak of “*real projects*” and their artifact support
 222 JavaScript, so we assume an expectation that the projects “developed” in a sense similar
 223 to [18].

224 Frame Oracle: All developed JavaScript projects hosted on GitHub.

225 Sampling Strategy: Convenience sampling yielded 19 JavaScript projects with high star
 226 counts.

227 Validity: Authors manually inspected projects in the selection.

228 Reproducibility Artifacts: A GitHub repository containing the tool and a file with project
 229 URLs along with their head and base commits is provided. Contents of repositories are not
 230 included. Run scripts did not run out of the box.

231 3.4 MSR 2020: Code smells

232 Code smells are programming idioms correlated with correctness or maintenance issues.
 233 Jebnoun et al. [9] contrast code smells in projects related to deep learning and general purpose
 234 software. Their claim is that for large and small projects there is a statistical difference in
 235 the occurrence of code smells, whereas medium sized projects are indistinguishable.

236 Population Hypothesis: The paper focuses on two populations: projects related to deep
 237 learning, and general purpose software. For pragmatic reasons, they focus on Python as it is
 238 popular for machine learning.

239 Frame Oracle: Python projects with keywords indicating machine learning, discarding
 240 tutorials. Furthermore, the authors “*also carefully select popular and mature DL projects*
 241 *from them by employing maturity and popularity metrics (e.g., issue count, commit count,*
 242 *contributor count, fork count, stars).*”

243 Sampling Strategy: A staged strategy was employed. The authors relied on judgment
 244 sampling to manually select 59 deep learning projects. For general purpose projects, they
 245 used a top-starred list of 106 Python projects from [3] and randomly sampled 59 projects.
 246 Projects were further clustered into small ($\leq 4,000$), medium, and large ($\geq 15,000$) based
 247 on size.

248 Validity: No issues were discussed.

249 Reproducibility Artifacts: A listing of the 59 deep learning projects is provided.

250 3.5 Summary and discussion

251 The papers we have reviewed do not explicitly talk about any of the four points in our
252 protocol, in all cases we had to reverse engineer (or guess) some of them. This suggests that
253 our proposal would improve the generalizability of the research.

254 While the mentioned research projects were done with care, there were challenges
255 reproducing them out of the box. Common sources of reproducibility failures that occur in
256 the papers we have reviewed are:

- 257 ■ *Missing descriptions*: Failure to specify either one of: population hypothesis, frame oracle
258 or sampling strategy. Reproduction is fraught with perils and an apple-to-apple comparison
259 between papers is difficult. This affects [21, 20, 16, 9] as their descriptions are open to
260 interpretation.
- 261 ■ *Missing projects*: Even with a list of URLs, the corresponding projects may vanish at any
262 time (e.g., deleted or made private). Reproductions are partial at best, we have seen a
263 project disappear while being downloaded. This affects [21, 20, 16, 9].
- 264 ■ *Fading stars*: Stars are volatile. [20] observed close to 3,000 projects in the top 1K
265 during a period of two months. Without a history of star attribution and a timestamp,
266 reconstructing the star listings is not possible. Stars volatility also caused problems for [9].
- 267 ■ *Shifting contents*: The contents of a project change with new commits. To reconstruct the
268 data, ids of the last observed commit must be specified. Even that is not foolproof as Git
269 histories can be updated destructively. This affects [21, 20, 9].
- 270 ■ *Language attribution*: Projects contain code in many languages. For reproduction
271 attribution must be specified. While delegating to, e.g. GitHub, is reasonable, one
272 should be aware that GitHub has changed their attribution algorithm several times.
273 Double counting a project is sometimes valid. This affects [21, 20, 16].
- 274 ■ *Deterministic replay*: Non-determinism must be limited. Random sampling seeds should
275 be specified. This affects [16].

276 4 Mapping the GitHub landscape

277 The meta-study of Table 1 highlights the dominant position of GitHub as a data source in
278 large-scale code analysis studies. The size of GitHub is such that it is necessary to resort to
279 sampling to yield manageable datasets. As shown in the previous section, authors often look
280 for some notion of “developed” projects, that is, they want projects that contains code of
281 some quality.

282 We claimed that convenience sampling using stars as a proxy for various other characteristics
283 of “real-world” software is flawed. While this may sound plausible to some readers, it should
284 be backed up with data. Given the size of GitHub, this section uses sampling to answer
285 the following questions: *Are starred projects a representative sample of all projects?* and
286 *Are starred projects a representative sample of developed projects?* where what it means for
287 a project to be developed is purposefully left open as there is no agreement on a precise
288 definition of the term.

289 Since the later parts of this paper require Java, Python and JavaScript, we acquire
290 samples of these three ecosystems. We use CodeDJ to do this. It is an open source project

291 that allows users to create a dedicated input project database and ensure reproducibility of
 292 queries.

293 We used random sampling over the entire GHTorrent dataset to select which projects to
 294 acquire in each of the languages of interest. The number of downloaded projects is somewhat
 295 arbitrary as it is based on available hardware during the acquisition phase. The datastore
 296 has 1,111,950 Java projects, 216,602 Python projects and 1,259,856 JavaScript projects. To
 297 give an idea of the scale, our Java dataset accounts for 20% of all non-forked GitHub Java
 298 projects. To get a manageable size, we down-sample further, randomly selecting 1Mio Java
 299 and JavaScript projects, and 200K Python projects.

300 4.1 Attributes

301 With CodeDJ, it is easy to write queries that compute project attributes. For this paper, we
 302 calculate five attributes that highlight the differences between projects:

- 303 ■ **C-index:** A developer handle has a C-index of n if that developer was party to at least
 304 n commits to n projects (i.e. n^2 commits). The C-index of a project is the highest such
 305 number across developers. This measures developer expertise.
- 306 ■ **Age:** The age of a project is the number of days separating the first commit and the most
 307 recent commit. This correlates with the maturity of a project.
- 308 ■ **Devs:** The count of unique developer handles in the git logs; includes both the author of
 309 a code change and the committer of that change. Devs approximates the size of a team,
 310 as some individuals may have more than one handle this is an upper bound.
- 311 ■ **Locs:** The total number of lines in files that are recognized as code, in any language, and
 312 appear in the head of the default branch.
- 313 ■ **Versions:** A version is implicitly created for each commit touching a file, be that for
 314 creation, deletion or update. This counts versions in the entire project's history including
 315 branches. Versions measure the activity in a project.

316 While we make no claims that these attributes suffice to fully describe a software project, we
 317 have found them to be an effective summary in many interesting dimensions.

318 4.2 Stars v. All

319 What do these attributes tell us about the overall population and about starred projects?
 320 If starred projects were representative of the entire population, they would share the same
 321 statistical distribution.

322 Fig. 1 is a histogram of each attribute; the x-axis is log scaled values in the unit of each
 323 attribute, the y-axis is the proportion of projects normalized for the maximum height. Grey
 324 denotes the whole population, red and blue denote the 1K most starred projects written
 325 in Java and Python respectively. The black, red and blue bars denote the median of their
 326 respective populations.

327 Consider the grey bars for the whole population, when comparing Java and Python, we
 328 see the same general shapes. The C-index is low, with a median of 2 for Java and Python.
 329 This means that half of the projects hosted on GitHub, have developers who have made at
 330 most two commits. The median age of Java projects is 7 days, while Python projects trend
 331 slightly older, 46 days. The median number of developers for both languages is 2. As for
 332 median lines of code, Java project are slightly larger than Python, 655 compared to 448. The
 333 median number of commits (versions) is 16 for both languages. Overall, this confirms that
 334 most projects are small, short-lived and created by relative newcomers.

335 The top 1K starred projects have a very different make up. Visually it is clear from the
 336 fact that every distribution is shifted right. Starred projects are larger, older, with more
 337 experienced developers. While there are slight differences between languages, the overall
 338 picture is consistent.

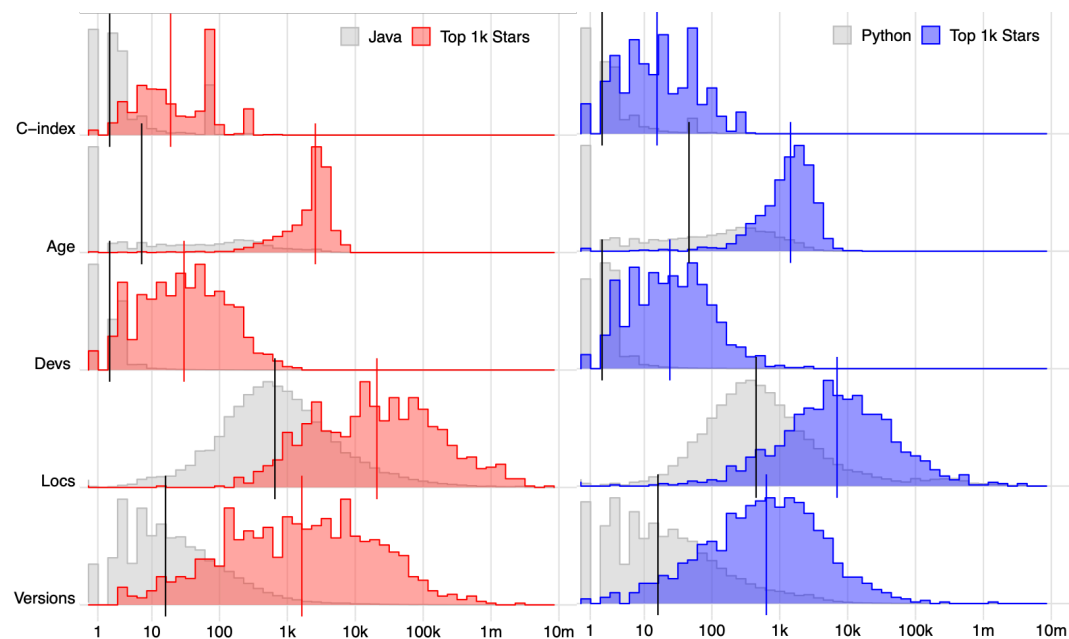
339 Consider for instance, the C-index and age attributes. While many starred projects are
 340 team efforts, a significant number of projects have few contributors. Their C-index is high,
 341 with median of 19 for Java and 15.5 for Python, suggesting that experienced developers tend
 342 to contribute to popular repositories. The median age projects is more surprising with 2,581
 343 and 1,440 days. Manual inspections suggest that many starred projects are indeed long lived
 344 but also have been inactive for years. Projects rarely “lose” stars, so if a project gets to the
 345 top there is a chance it will stay there long past its useful lifetime.

346 The answer to our first question is clearly negative. Starred projects are not representative
 347 of the overall population. This is not necessarily a bad thing, as folklore suggests that most
 348 of GitHub is uninteresting. Perhaps it is the case that starred projects are more “interesting.”

349 4.3 Stars v. Developed

350 Researchers often look for *engineered* [18, 22] or *developed* projects – informally, projects
 351 created with some care – alas there is little agreement on a precise definition.

352 Slightly easier, perhaps, is to settle on what we do not want, the projects that are
 353 uninteresting, one that are clearly of little value for any reasonable research question.
 354 Moreover, one could hope that the complement of uninteresting projects are the projects
 355 researchers look for. Let us define a project as *uninteresting* if it has less than 100 lines of
 356 code, fewer than 7 days old, and fewer than 10 commits. When this definition is used to
 357 filter projects, this rather low bar manages to eliminate 71% of Java and 55% of Python
 358 projects. For the purpose of this discussion we term the remaining projects are developed or
 359 interesting. It would be nifty if stars were a proxy for filtering out uninteresting projects.



■ **Figure 1** Comparing datasets

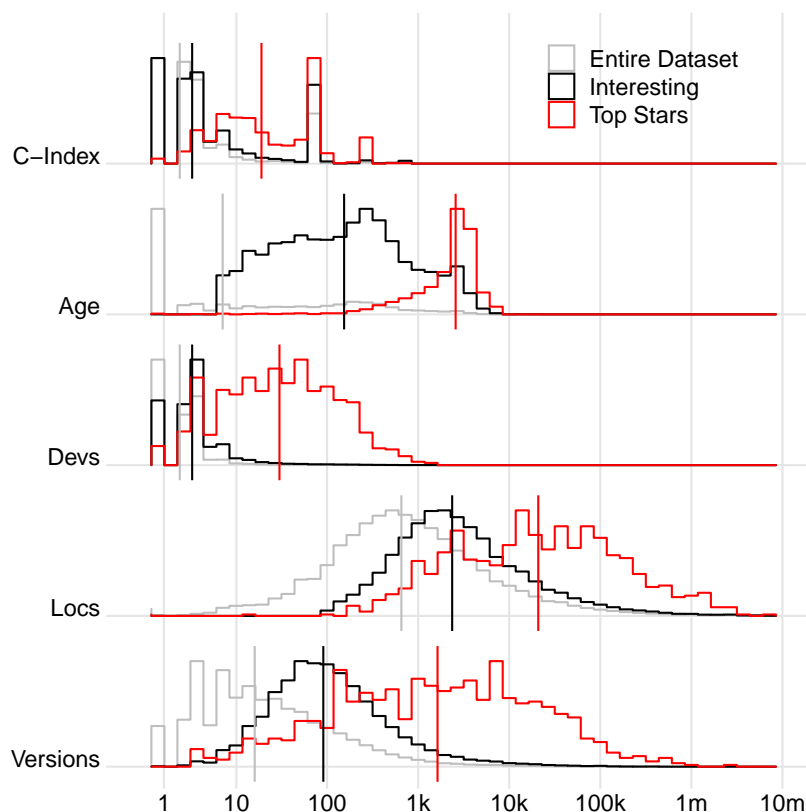
23:10 Designing Reproducible Large-scale Code Analysis Experiments

360 Fig. 2 shows the distribution of the whole population in Grey (in the same way as in
361 Fig. 1), the interesting projects in Black and the top 1K starred projects in Red. Clearly, the
362 shape of the Black and Red distributions do not align suggesting that stars do not represent
363 interesting projects.

364 Manual inspection of the starred project highlights their main issue – stars are extrinsic
365 properties without a direct connection to any attributes of a project. Unlike our computed
366 attributes, stars grow monotonically. Their meaning is unclear as users award them for
367 various reasons including humor and shock value. Some projects earned stars because of a joke
368 not fit for this audience (e.g. `github.com/dickrnn/dickrnn.github.io`), or have dared
369 users to star junk (`github.com/gaopu/java`). Stars do not measure quality or usefulness of
370 repositories.

371 To further illustrate the limitation of stars as a filter, we take, for each attribute, the 20
372 projects with the lowest score for that attribute. Table 2 shows a manual classification of
373 these projects. None of these projects is useful: externals lack histories, widgets are small
374 and biased by their application domain, babies are too small to yield much insights, and the
375 remaining ones only have code snippets.

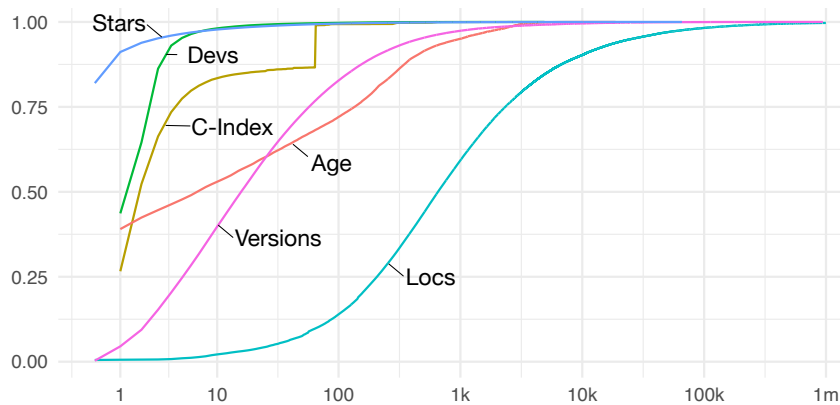
376 The answer to our second question is also negative. Starred projects are not representative
377 of the interesting ones. To summarize what we learned about stars, they capture extrinsic
378 characteristics of GitHub projects and are at best an indirect and noisy proxy for a robust
379 frame oracle.



■ **Figure 2** Comparing developed and starred projects

Category	Java	Python	Description
Externals	9%	5%	Infrequent synchronization with another repository.
Widgets	43%	0%	Tiny projects with little activity, popular UI widgets or plugins.
Docs	4%	15%	Interview questions, course materials, games, knitting patterns.
Tutorials	17%	9%	Educational materials, tutorials and example applications.
Babies	16%	32%	Valid but extremely small projects with little activity.
Artifacts	0%	21%	Research artifacts developed elsewhere and deposited for sharing.
Deprecates	1%	5%	Deprecated projects, no code on the main branch.

■ **Table 2** Categorizing 200 starred projects



■ **Figure 3** Cumulative Density Functions

380 4.4 How to select projects?

381 What to use for project selection if not stars? We argue that selection must be based on
 382 intrinsic features – measurable attributes of a project’s contents. While one may use machine
 383 learning [18, 22] to build classifiers, in this paper we will use our five computed attributes
 384 (we leave machine learning as an interesting area of future work).

385 Fig. 3 is the cumulative density function of the various attributes for Java (the shapes of
 386 the curves for Python are similar). The interpretation of each line is what percentage of the
 387 dataset is filtered for a particular attribute value. Project selection can be performed by a
 388 combination of attributes with cutoffs. We do not argue for a particular formula; researchers
 389 must make their own choices in this respect.

390 For instance, if one were to use 10 days of age as a cutoff, then 52% of the dataset would
 391 be filtered out. Whereas picking a 10 star cutoff, filters out 98% of projects.³

392 4.5 Validity of our dataset

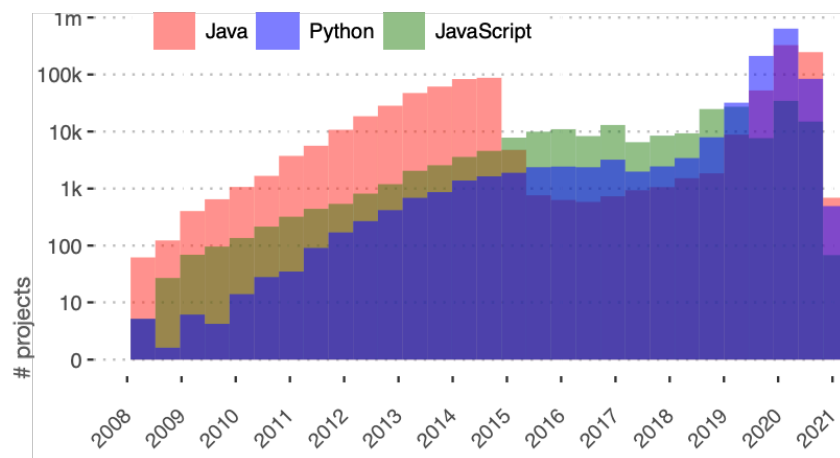
393 We noticed an oddity around project ages in our dataset. Experience with GitHub trained
 394 us to expect the unexpected. Our investigation started with a plot of creation dates.

395 Fig. 4 shows the log scaled counts of new projects over time. While there is a steady

³ The discontinuity of C-Index at 65 is odd. After manual investigation, we found that there is a single developer with that C-index, it turns out that the “developer” is a bot doing automated updates.

23:12 Designing Reproducible Large-scale Code Analysis Experiments

396 progression in the count of projects created each year, we see a significant drop in 2015 and
397 a plateau until 2019.



■ **Figure 4** Creation date

398 We reviewed our pipeline to no avail. We use GHTorrent to acquire all available URLs.
399 Then, we randomly sample projects from that list. We validated both acquisition and
400 sampling. This leaves us with two hypotheses. First is a consistent flaw in the CodeDJ
401 downloader causing some projects to fail to download. 17% URLs obtained from GHTorrent
402 point to dead projects, but there is no apparent bias. Second some projects could be missing
403 in GHTorrent.

404 Another issue showed up on inspection, JavaScript project ages are significantly higher
405 than those of other languages. We found that GitHub timestamps are frequently inconsistent.
406 Why should JavaScript be more affected? Until an explanation can be found, we removed
407 JavaScript from the overall comparison and use JavaScript projects in the reproduction with
408 extreme care.

409 **5** Reproducible large-scale analysis experiment design

410 This paper proposes that researchers conducting experiments over large-scale software
411 repositories follow a specific experimental design methodology to ensure their work can
412 be reproduced and increase chances that their results generalize as expected. While the
413 mechanics of reproducibility of the actual experiment itself vary, the setup of the experiment
414 is a common problem. The proposed methodology has five steps, we encourage researchers
415 to document each of these steps explicitly.

416 **5.1** Population hypothesis

417 Formulating a population hypothesis lets researchers stake a claim about the applicability
418 of their work. This represents the population to which the result of an experiment should
419 generalize to. The statement of that hypothesis can be brief and appeal to intuition, the
420 other parts of the description flesh out the details.

421 Ideally, we would like our results to be as broadly applicable as possible, but pragmatically
422 designing experiments that back up overly broad claims is difficult. Some populations of
423 interest are difficult to sample, for instance “commercial software” is a relatively simple

424 and unambiguous description but one that we typically cannot sample from as most of the
425 commercial software is not in the public domain. Other populations can be difficult to
426 identify. Imagine a study of the challenges linked to retraining imperative programmers to
427 use functional idioms. Finding code written by such developers can be done manually but
428 is difficult to automate. It is often easier to describe a population by intrinsic features of
429 projects such as the language used to write the code or some estimate of the size of the
430 project.

431 5.2 Frame oracle

432 A frame oracle is a, possibly noisy, deterministic algorithm for deciding if a project belongs
433 to the population of interest. The oracle is our best approximation of the population of
434 interest. An executable and reproducible oracle allows to compare different papers with the
435 same selection. The description of the oracle should specify the data source along with any
436 information required to acquire projects. The procedure for evaluating a project should be
437 clear and based on intrinsic attributes. A paper should at least have a short description of
438 the oracle, full details should be given in the reproduction artifact.

439 5.3 Sampling strategy

440 The literature has an abundant advice on sampling (see e.g. [12]). Briefly, a sampling
441 strategy picks the type of sampling (probabilistic or non-probabilistic) and describes the
442 steps used to obtain a sample. The sampling implementation is expected to be found in the
443 reproduction artifact.

444 Many works use convenience sampling as it is simpler, cheaper and less time consuming.
445 A better alternative is some form of probabilistic sampling as it is more likely to yield a
446 representative sample. Probabilistic sampling can be staged if the structure of the population
447 is more complex. The simplest approach is random sampling where each element has the same
448 chance of being picked. We often have to resort to stratified sampling when the population
449 is divided into subgroups of different sizes. Typically we sample without replacement as we
450 do not want to pick the same project multiple times.

451 5.4 Validity

452 The validity section argues, when there are reasons for doubt, why using the frame oracle
453 and the sample strategy results in representative samples of the population of interest. This
454 section should address potential sources of bias and attempts by the authors to control for
455 them. This section also should address any foreseen challenges to reproducibility and offer
456 means to mitigate them.

457 5.5 Reproducibility artifacts

458 Finally, we advocate to link the paper to a reproduction artifact that contains code and data
459 to support experimental repeatability and reanalysis.

460 Section 3 listed issues with reproducibility. In some cases, the authors did not give a
461 precise description of the steps needed for reproduction. Following our proposed methodology
462 along with a reproduction artifact will greatly help.

463 The second category of issues are more pragmatic, it is difficult to repeat the analysis
464 of a paper because some aspect of the data used is not available. We suggest that research
465 infrastructures should support this task.

23:14 Designing Reproducible Large-scale Code Analysis Experiments

466 An example of an infrastructure is CodeDJ which is both a continuously updated datastore
467 and a database that can be queried by a DSL. We adopted it for our work and illustrate how
468 it helps with reproducibility. The implementation of a frame oracle and sampling strategy
469 can be combined into a single expression. Fig. 5 shows a query which starts by filtering out
470 projects containing fewer than 80% JavaScript code, then it uses pre-computed attributes
471 Locs, Age and Devs to filter further. The last stage of filtering involves computing an
472 attribute on the fly, here we sum up the commits in the project, before performing random
473 sampling.

```
database.projects().filter(|p| {  
    p.language_composition().map_or(false, |langs| {  
        langs.into_iter().any(|(lang, rate)| { lang == JavaScript && rate >= 80 })  
    })  
})  
    .filter_by(AtLeast(Locs, 5000))  
    .filter_by(AtLeast(Age, Duration::from_months(12)))  
    .filter_by(AtLeast(Devs, 2))  
    .filter_by(AtLeast(Count(Commits), 100))  
    .sample(Random(30, SEED))
```

■ **Figure 5** Project selection with CodeDJ

474 CodeDJ is split between a persistent datastore in which every data item is timestamped,
475 and an ephemeral database used to service queries. A reproducible query is a Rust crate
476 archived in a git repository associated to the datastore. Running the query produces a *receipt*
477 which is the hash of a commit automatically added to the archive repository. The receipt
478 can be used to share the query (exactly as executed) and its results (exactly as produced). It
479 can be used to retrieve the Rust crate and re-execute the code. Code re-execution is helped
480 by the fact that queries are deterministic and the crate contains a list of all dependencies,
481 a timestamp, and all random seeds. When a query with an embedded date is executed,
482 CodeDJ accesses the exact state of the datastore at the specified date. Since CodeDJ stores
483 the contents of files, entire experiments can be fully reproduced.

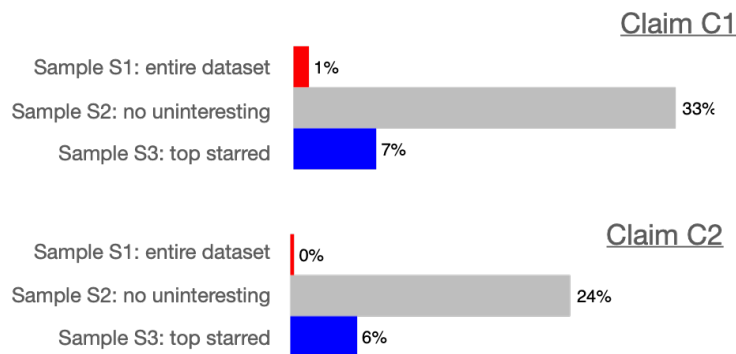
484 **6 Reproductions**

485 We illustrate the use of the proposed methodology by revisiting the papers we discussed
486 in section 3. For each paper, we attempt reproduction where we vary the input. The
487 original papers used stars in their selection, we will explore different inputs based on intrinsic
488 attributes.

489 If the results of the reproduction match the original results, then it suggests that stars
490 were an appropriate filter. If the reproduction departs, this suggests that there may be need
491 to conduct further experiments to be confident in the results.

492 For each paper, we picked a subset of the scientific claims to fit the reproduction in the
493 available space. We use our proposed methodology to describe the details of the reproduction.

494 One may wonder how we selected the paper to reproduce. Our criteria were (1) papers
495 that used automated techniques to analyze properties of the code contained in Github
496 projects, (2) their population of interest was a large subset of Github, (3) a working artifact
497 could be located, and (4) the input could be changed with ease. We did not cherry-pick as
498 even positive results would be interesting. Our choice was limited by the fact that many
499 papers either did not have artifacts or that they were not working anymore. Furthermore,
500 some papers had hardwired their selection of projects by, for example, preprocessing the



■ **Figure 6** Proportion of projects without code, data or documentation

501 input data and omitting to include the tooling to repeat that processing. Given more time,
 502 more works could be reproduced.

503 6.1 Reproduction: What is software

504 This reproduction aims to validate two findings of [21]: (C1) 4% of repositories do not contain
 505 code, data and documentation; (C2) 2% of repositories do not contain documentation.

506 *Population Hypothesis:* The universe of software projects.

507 *Frame Oracle:* To understand the impact of project selection we consider two oracles. O1
 508 accepts any project hosted on GitHub. O2 removes uninteresting projects (as defined above).

509 *Sampling Strategy:* We report on three samples. S0 is a convenience sample of starred
 510 following [21]. S1 and S2 are random samples without replacement from O1 and O2
 511 respectively, stratified by language and deduplicated.

512 *Validity:* Our reproduction differs in the number of languages (3 v. 25) and by categorizing
 513 files based on the file path alone. We tested stability of our results with multiple samples of
 514 varying sizes and manually inspected the produced labels.

515 *Reproduction Artifact:* Our artifact contains a CodeDJ receipt for this query.

516 6.1.1 Results

517 Fig. 6 shows results for claims C1 and C2. Compare the percentages between S0 (original)
 518 and S1 (target population). Statistical analysis is not required to see that the difference is
 519 significant. The sample S2 (without uninteresting projects) is there to illustrate the impact
 520 of slightly more developed population, but even these are still quite different.

521 Would the results agree if we included more languages? The three languages we
 522 downloaded account for most of GitHub, it is conceivable that other languages could affect
 523 results, but that would just push the generalizability issue somewhere else as the claims
 524 would become language-specific.

525 6.2 Reproduction: Method chaining

526 Nakamura et al. [20] claim that 50% of projects in 2018 had method chains longer than 7
 527 while in 2010 that number was 42%. They state that “chains of length 8 are unlikely to
 528 be composed by programmers who tend to avoid method chaining, this result is another
 529 supportive evidence for the widespread use of method chaining.”

23:16 Designing Reproducible Large-scale Code Analysis Experiments

530 Population Hypothesis: The universe of real-world Java programs.

531 Frame Oracle: We accept any Java project hosted on GitHub and delegate to Github for
532 language attribution.

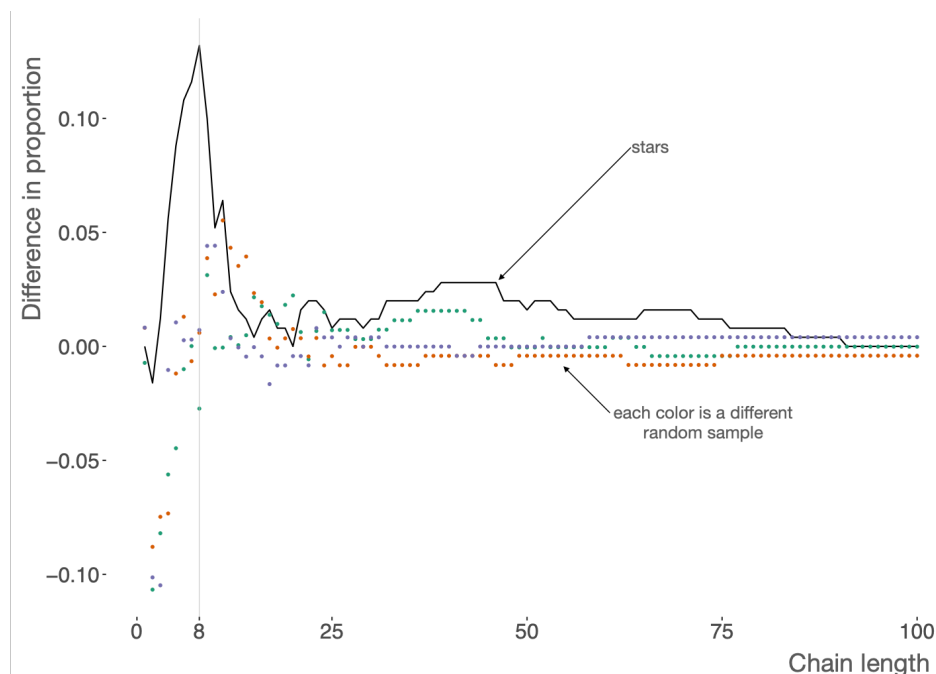
533 Sampling Strategy: Stratified sampling to randomly select projects with commits in 2010
534 and 2018.

535 Validity: To reproduce the original results, we performed stratified sampling to get top
536 starred projects active in the target years. The authors used a different sample of top stars.
537 The original paper had different sample sizes for each year, but those are not specified. We
538 fix the sample size to 250. The authors could not locate the code of their chain detector, so
539 we use our own implementation.

540 Reproduction Artifact: Our artifact contains a CodeDJ receipt for this query.

541 6.2.1 Results

542 Fig. 7 shows the difference in proportion of projects at various chain lengths. The solid
543 line uses stars, colors represent different random samples. For instance, if we pick chains of
544 length 8, the number used by [20], the difference is a 13% increase in the number of projects
545 between 2011 and 2018. The differences for our random samples are -2%, 0.6% and 0.7%.
546 In other words, the samples from this particular population do not seem to show the effect
547 expected by the authors. We surmise that some notion of developed project may show more
548 favorable results, but without more guidance in the population hypothesis it is hard to guess
549 which to pick.



■ **Figure 7** Difference in chain lengths

6.3 Reproduction: Style analyzer

Markovtsev et al. [16] builds a model of the style of a repository and apply this model on a held-out part of that repository to produce corrections. Their experiment uses 19 top-starred JavaScript project to gauge the precision with which the tool flags formatting discrepancies and the relationship between this precision and the size of the project. They report a precision of 94% (average, weighed by project size) and better overall performance for large projects and projects with better style guidelines.

Population Hypothesis: Developed JavaScript projects.

Frame Oracle: JavaScript projects with at least 80% JavaScript code, $\text{Loc} \geq 5000$, $\text{Age} \geq 12 * 31$ and $\text{Devs} \geq 2$.

Sampling Strategy: We randomly select 10 sets of 30 projects. This is more projects than the original sample to account for errors in processing. After processing is finished, following the original paper, we randomly select 19 projects out of the pool of successfully processed projects.

Validity: Given that the author’s artifact lacks a configuration, we used the default one. This increases project size, as compared to the published numbers, by 38% per project (up to a maximum of 154%) and causes precision to diverge by 2.2% on average, and up to 7.9%.

The tool failed to process 4 projects: `freecodecamp` and `atom` due to errors in unicode processing, `express` due to a programming bug, and `30-seconds-of-code` due to bad file identification. Three of the missing projects were located close to the median in terms of precision, prediction rate, and project size in the original paper, while `axios` was in the lower quartile for sample count.

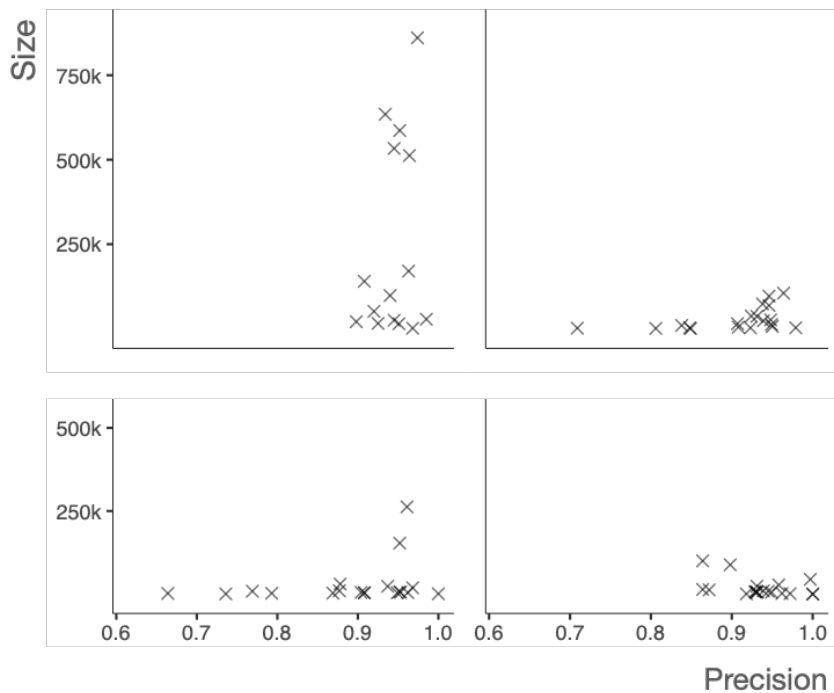
Style analyzer analyzes each project at two points in its history specified by a base commit and a head commit. The base commit is a point in the past which the tool checks out to learn the project’s formatting style. The head commit is a more recent point used to evaluate the model and calculate precision. The original paper provides head and base commits for each project in their experiment, but does not specify the method of selecting these commits. We pick the current head of the default branch as the head commit. For base commit we pick one that lies at an offset equal to 10% of the number of all commits in the default branch from the head commit. This retrieves different commits than the original paper, which causes a 3% median change in precision (up to 17%—`telescope`) and a median project size increase of 76%, and up to 311% (`reveal.js`).

Reproduction Artifact: Datasets, receipts from submitted queries, style analyzer’s reports and scripts for the entire experimental pipeline are included in the artifact.

6.3.1 Results

We recreate a plot of the effect of the number of items in the training set on precision from the original paper in Fig. 8. The training set consists of snippets created around tokens/AST nodes relevant to formatting (whitespace, indentation, quotes, zero-length gaps). We plot the selection from the original paper along three selections from our interesting project frames. In addition, we plot the distributions of precision in each selection in

In Fig. 9, we compare the precision scores in each sample with the selection used in the original paper using a Mann-Whitney U test to show which samples performed statistically differently from the original. The scatter plots show a different grouping of results from the original paper. The groupings in the scatter plot visibly differ between selections. The distribution comparison shows that our selections generate significantly smaller training



■ **Figure 8** Relationship between label groups and precision

sets in all cases and yield lower precision. In addition, 3 out of the 10 interesting project selections produced significantly lower precision, with the remainder producing a statistically equivalent distribution.

Overall, we see our selections yielding precision between 0.9 and 0.95 (the paper sets a precision of 0.95 as a benchmark for success). We also do not see a clear relationship between the number of label groups and precision, such as the one the authors note in the original paper.

6.4 Reproduction: Code smells

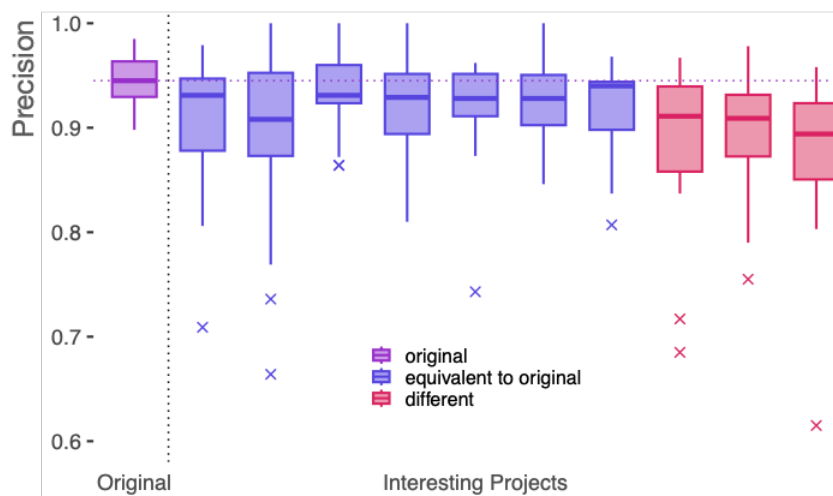
We seek to validate the claim of [9] that for large and small projects there is a statistical difference in the occurrence of code smells between machine learning and most popular Python repositories, whereas medium sized projects are indistinguishable.

Population Hypothesis: Mature Python projects in all application domains including machine learning.

Frame Oracle: Projects with C-Index ≥ 5 , or Age ≥ 180 , or Locs $\geq 10,000$, or Versions ≥ 100 .

Sampling Strategy: The deep learning projects were provided by the authors. Out of 59 projects, 57 were still accessible on August 2nd 2021. At download time there were 6 small, 13 medium, and 38 large deep learning projects. For the reproduction of the original results, we used a staged strategy, first convenience sampling the top starred Python projects and amongst those used stratified sampling to select 57 projects with a similar distribution of sizes. To generalize the results we used quota sampling to match the size distribution.

Validity: Our reproduction uses the Locs reported by CodeDJ. The date the authors



■ **Figure 9** Comparing label group count and precision

617 downloaded the repositories is unknown. We use the content of the main branch of each
 618 repository as of April 1st, 2020. The authors say “*each of repositories is pre-processed*
 619 *and prepared for code smell detection*”, however details are missing. We used the default
 620 thresholds of their tool.

621 Reproduction Artifact: A CodeDJ receipt is included in our reproduction package along
 622 with code to run the experiment.

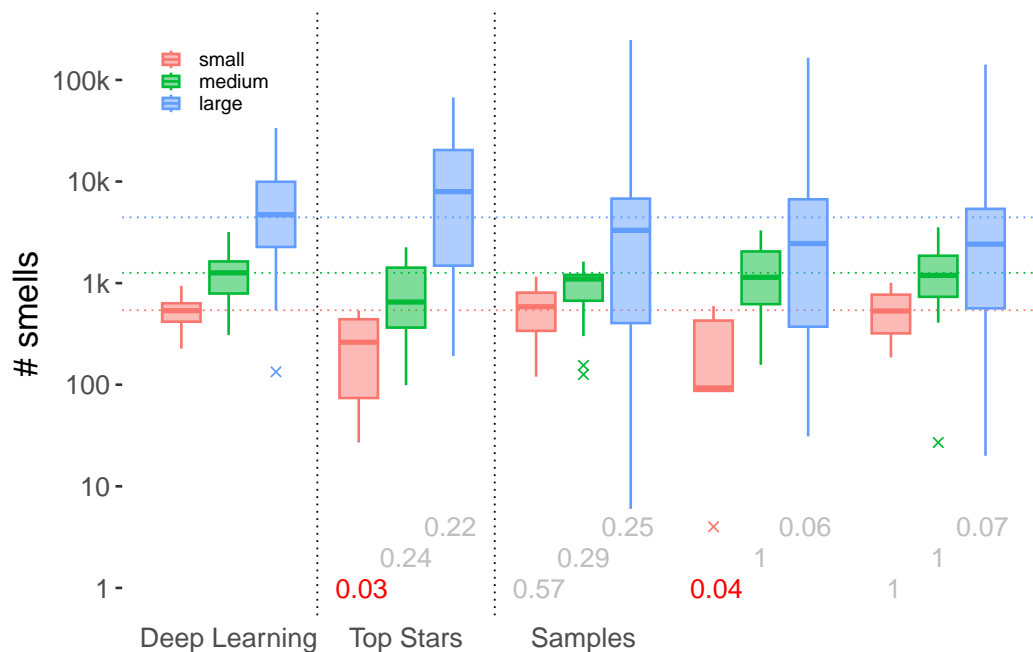
623 6.4.1 Results

624 Fig. 10 contrasts the distribution of code smells for deep learning projects, top starred
 625 projects, and three random samples. Computing the p-values with the non-parametric
 626 Mann-Whitney Wilcoxon shows that while we were able to reproduce the statistically significant
 627 results for the small projects, we disagree on the large most popular projects with the original.
 628 The disagreement is even greater in the random samples where no large projects and only
 629 one small project is statistically different. Generalizability of the results is thus questionable.

630 7 Collaborative Reproduction

631 We perform one last reproduction in which we obtain the assistance of the study’s authors
 632 to validate whether stars are a good input selection strategy. For this reproduction we
 633 selected a distinguished paper from the Foundations of Software Engineering conference
 634 (2020), “The Evolution of Type Annotations in Python” [7]. The paper has a reusable artifact
 635 for repetition of the original results. Our reproduction only required to change the list of
 636 GitHub URLs used as input in the analysis. The authors helpfully allowed us to run the
 637 rather computationally intensive workload on their machine.

638 The original study reported the following protocol for input selection. “*We group projects*
 639 *by creation date, considering projects created in the years 2010 to 2019, into ten groups. We*
 640 *sort each group by number of stars and select the top-1000 per group, which yields a total of*
 641 *10,000 projects. The rationale for first grouping and then sampling is to avoid biasing our*
 642 *study toward projects created in a particular time frame, e.g., mostly old projects. Removing*
 643 *projects that we could not clone, e.g., because they became unavailable since the beginning of*



■ **Figure 10** Comparing Smells. Numbers are p-Values indicating a significance of the difference from the deep learning projects. Statistically significant different p-Values (cutoff at 0.05) are shown in color, insignificant ones are in gray.

644 *our study, the total number of analyzed repositories is 9,655.*

645 The first research questions was related to the evolution in the number of type annotations
 646 in projects. The main insight from the work was that *“better developer training and automated
 647 techniques for adding type annotations are needed, as most code still remains unannotated,
 648 and they call for a better integration of gradual type checking into the development process.”*

649 For this reproduction, we discussed input selection criteria with the authors and arrived
 650 at the following formulation.

651 *Population Hypothesis:* Python projects in all application domains with earliest commit
 652 date in 2015⁴ as this was the year when early adoption of type annotations began.

653 *Frame Oracle:* Python projects whose life span is longer or equal to 7 days and that have
 654 over 100 lines of Python code. The study authors intended their work to be representative of
 655 most of the Python ecosystem, but closer inspection of some of the small projects suggested
 656 that they would introduce noise. The particular cutoffs were chosen heuristically.

657 *Sampling Strategy:* Projects were grouped by year active and, for each year, a random
 658 sample of 1200 projects was selected. The goal was to get close to a thousand usable projects
 659 for each year. As some of the projects in our database are no longer available, the sample
 660 size is increased heuristically.

661 Fig. 11 illustrates the reproduction results (and mirrors Fig. 2 in [7]). The plot on the
 662 left shows the number of type annotations found per thousand lines of code in the projects
 663 being looked at. The red line has the new data, the blue one is for the original study. The

⁴ 2015 is when type annotations were added to Python.

664 y-axis is logarithmic. The two lines start at zero in 2015. Both data sets tell a similar story:
 665 type annotations are gradually added to projects. The plot on the right shows the total
 666 number of annotations found each year in all projects. The y-axis is in thousands. In 2021,
 667 the original data had close to 800,000 annotation while the new data is under 250,000.

668 Both data sets are large. The original one contains 1,123,393 commits and the new data set
 669 set 1,535,824 – suggesting that projects are slightly larger in the randomly selected data than
 670 in the most starred projects. In both cases, only a fraction of the repositories have types. In
 671 the old data 668 repositories are type-annotated, whereas in the new data 1,040 projects
 672 have at least one type. The fraction of commits that change a type annotation is small in
 673 both cases 5.5% in the original data and 2.1% in the new data.

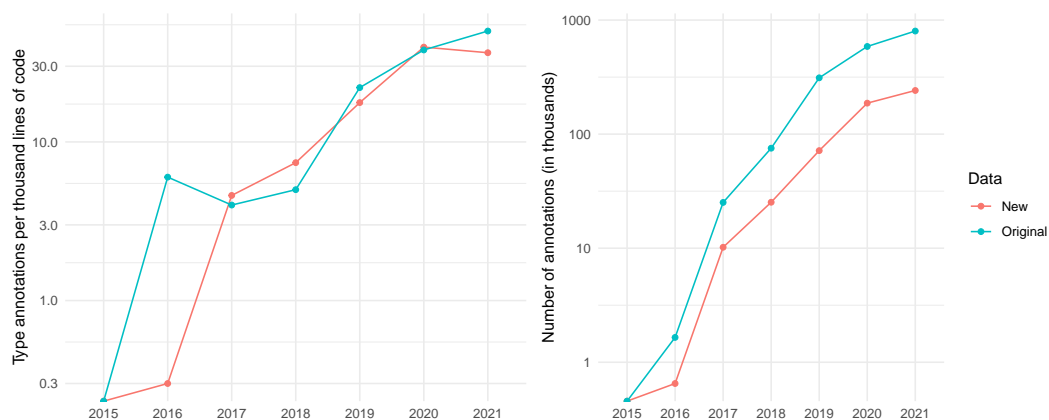
674 Overall, the reproduction verifies and, even, strengthens the conclusion of the original
 675 paper. Five years after introduction of type annotations, their use remains rather limited.
 676 Having said this, it is true that actual values reported are different enough to be noticeable.

677 8 Conclusions

678 Sometimes, doing it wrong is so much easier than the alternative, that we convince ourselves
 679 that a little wrong can be right enough. Our paper is unusual. While it purports to contain a
 680 call to arms for better experimental practices, it is just as much a record of our own journey
 681 to that goal. What reads as criticism is really written in self-reflection. So, what can a
 682 researcher take away from this paper? There are three ideas we would like to leave the reader
 683 with.

684 *Generalizability.* The value of an experiment often lies as much in what it generalizes to,
 685 as in the experiment's outcome. We found that many researchers rely on GitHub stars to
 686 pick representative samples of software projects, yet starred projects tend to be larger in
 687 most dimensions than typical ones, also that they are more likely to be inactive, and that
 688 their ranking is not a measure of intrinsic qualities of the code. Hopefully, this paper is the
 689 last nail in that coffin. More generally, we advocate for the use of probabilistic sampling over
 690 populations defined by intrinsic attributes of software, and also for clear and standardized
 691 documentation of experimental design.

692 *Reproducibility.* The value of a scientific experiment also lies in our ability to reproduce
 693 it. Carrying out reproducible experiments over large-scale software repositories is hard.
 694 Especially when aiming to support the three reproduction modalities: repetition, as practiced



677 ■ **Figure 11** Types in Python

695 in artifact evaluation, where an artifact is re-executed to obtain identical results; reanalysis,
 696 where the artifact or its input are modified; and independent reproduction, where the entire
 697 experiment is re-implemented from scratch. The first modality requires faithful replay and
 698 is best served if all data used is included with the artifact. The second, requires support
 699 for automatically acquiring new representative samples. The third needs an unambiguous
 700 description of all experimental steps. We advocate for reproductions artifacts that supports
 701 the first two modes, and a detailed description of the experiment for the last.

702 *Tooling.* Generalizability and reproducibility, while worthy goals, represent much work,
 703 and they are work that is orthogonal to the scientific goals of researchers. The only
 704 reasonable answer is to provide tooling that automates acquisition of representative samples
 705 and generation of reproduction artifacts. In this paper, we used CodeDJ and found it helpful
 706 as it let us specify queries over attributes of the code for many projects, while also supporting
 707 experimental repetition and reanalysis through historical queries. It has its limitations, we
 708 found execution times to be somewhat long and doubt it will scale to the whole of GitHub.

709 Our vision for a brighter future is one where the community agrees on standard tools and
 710 techniques for this kind of experiment, tools which automate the acquisition and packaging
 711 of input datasets and the re-execution of entire experiments.

712 Acknowledgments

713 We would like to thank Digital Ocean for their involuntary contribution of computational
 714 resources during the early data gathering phase of our research. This work was supported by
 715 the Czech Ministry of Education, Youth and Sports under program ERC-CZ, grant agreement
 716 LL2325, NSF grants CCF-1910850, CNS-1925644, and CCF-2139612, as well as the GACR
 717 EXPRO grant 23-07580X. We acknowledge the reviewers of ICSE'22, and thank the reviewers
 718 of ECOOP'23 for their encouragements and for sticking around until 2024.

719 References

- 720 1 S Baltés and P Ralph. Sampling in software engineering research: a critical review and
 721 guidelines. *Empir. Softw. Eng.*, 27(4):94, 2022. doi:10.1007/s10664-021-10072-8.
- 722 2 H Borges and M Tulio Valente. What's in a github star? understanding repository starring
 723 practices in a social coding platform. *Journal of Systems and Software*, 2018. doi:10.1016/j.
 724 jss.2018.09.016.
- 725 3 Z Chen et al. Understanding metric-based detectable smells in python software. *Information*
 726 *and Software Technology*, 2018. doi:10.1016/j.infsof.2017.09.011.
- 727 4 V Cosentino, J Izquierdo, and J Cabot. Findings from GitHub: Methods, datasets and
 728 limitations. In *Mining Software Repositories (MSR)*, 2016. doi:10.1145/2901739.2901776.
- 729 5 R Dyer, H Nguyen, H Rajan, and T Nguyen. Boa: A language and infrastructure for analyzing
 730 ultra-large-scale software repositories. In *Int. Conf. on Software Engineering (ICSE)*, 2013.
 731 doi:10.5555/2486788.2486844.
- 732 6 G Gousios and D Spinellis. GHTorrent: GitHub's data from a firehose. In *Mining Software*
 733 *Repositories (MSR)*, 2012. doi:10.1109/MSR.2012.6224294.
- 734 7 L Di Grazia and M Pradel. The evolution of type annotations in python: An empirical study.
 735 In *European Software Engineering Conference and Symposium on the Foundations of Software*
 736 *Engineering (ESEC/FSE)*, 2022. doi:10.1145/3540250.3549114.
- 737 8 J Han et al. Characterization and prediction of popular projects on GitHub. In *Computer*
 738 *Software and Applications Conf. (COMPSAC)*, 2019. doi:10.1109/COMPSAC.2019.00013.
- 739 9 H Jebnoun et al. The scent of deep learning code. In *Mining Software Repositories (MSR)*,
 740 2020. doi:10.1145/3379597.3387479.

- 741 10 E Kalliamvakou et al. The promises and perils of mining GitHub. In *Mining Software*
742 *Repositories (MSR)*, 2014. doi:10.1145/2597073.2597074.
- 743 11 S Krishnamurthi and J Vitek. The real software crisis: repeatability as a core value. *Commun.*
744 *ACM*, 58(3), 2015.
- 745 12 S Lohr. *Sampling: Design and Analysis*. 2010.
- 746 13 C Lopes et al. Déjà Vu: A map of code duplicates on GitHub. *Proc. ACM Program. Lang.*,
747 (OOPSLA), 2017. doi:10.1145/3133908.
- 748 14 Y Ma et al. World of code: enabling a research workflow for mining and analyzing the universe
749 of open source vcs data. *Empirical Softw. Eng.*, 2021. doi:10.1007/s10664-020-09905-9.
- 750 15 P Maj et al. CodeDJ: Reproducible queries over large-scale software repositories. In *European*
751 *Conf. on Object-Oriented Programming (ECOOP)*, 2021. doi:10.1145/2658987.
- 752 16 V Markovtsev et al. Style-analyzer: fixing code style inconsistencies with interpretable
753 unsupervised algorithms. In *Mining Software Repositories (MSR)*, 2019. doi:10.1109/MSR.
754 2019.00073.
- 755 17 T Mattis, P Rein, and R Hirschfeld. Three trillion lines: Infrastructure for mining github
756 in the classroom. In *Conf. on Art, Science & Eng. of Programming <Programming>*, 2020.
757 doi:10.1145/3397537.3397551.
- 758 18 N Munaiah et al. Curating github for engineered software projects. *Empirical Software*
759 *Engineering*, 2017. doi:10.1007/s10664-017-9512-6.
- 760 19 M Nagappan, T Zimmermann, and C Bird. Diversity in software engineering research. In
761 *Foundations of Software Engineering (FSE)*, 2013. doi:10.1145/2491411.2491415.
- 762 20 T Nakamaru et al. An empirical study of method chaining in Java. In *Mining Software*
763 *Repositories (MSR)*, 2020. doi:10.1145/3379597.3387441.
- 764 21 R Pfeiffer. What constitutes software? In *Mining Software Repositories (MSR)*, 2020.
765 doi:10.1145/3379597.3387442.
- 766 22 P Pickerill et al. Phantom: curating github for engineered software projects using time-series
767 clustering. *Empir Software Eng*, 2020. doi:10.1007/s10664-020-09825-8.
- 768 23 P Ralph. SIGSOFT empirical standards released. *Softw. Eng. Notes*, 46(1):19, 2021. doi:
769 10.1145/3437479.3437483.
- 770 24 J Vitek and T Kalibera. R3: Repeatability, reproducibility and rigor. *SIGPLAN Not.*, 2012.
771 doi:10.1145/2442776.2442781.