

The Dynamics of Innovation in Open Source Software Ecosystems

Gábor Mészáros* and Johannes Wachs^{†‡§}

*Center for Collective Learning, Corvinus University of Budapest, Hungary

[†]Department of Network Science, Corvinus University of Budapest, Hungary

[‡]HUN-REN Centre for Economics and Regional Studies, Budapest, Hungary

[§]Complexity Science Hub, Vienna, Austria

johannes.wachs@uni-corvinus.hu

Abstract—Software libraries are the elementary building blocks of open source software ecosystems, extending the capabilities of programming languages beyond their standard libraries. Although ecosystem health is often quantified using data on libraries and their interdependencies, we know little about the rate at which new libraries are developed and used. Here we study imports of libraries in 12 different programming language ecosystems within millions of Stack Overflow posts over a 15 year period. New libraries emerge at a remarkably predictable sub-linear rate within ecosystems per post. As a consequence, the distribution of the frequency of use of libraries in all ecosystems is highly concentrated: the most widely used libraries are used many times more often than the average. Although new libraries come out more slowly over time, novel combinations of libraries appear at an approximately linear rate, suggesting that recombination is a key innovation process in software. Newer users are more likely to use new libraries and new combinations, and we find significant variation in the rates of innovation between countries. Our work links the evolution of OSS ecosystems to the literature on the dynamics of innovation, revealing how ecosystems grow and highlighting implications for sustainability.

Keywords: Innovation, novelty, OSS ecosystems, maintenance, sustainability

I. INTRODUCTION

Software plays a key role in the modern digital economy [1], reflected in the explosive growth of software-related patents [2] and the increasing importance of software in traditional sectors like the automobile industry [3]. But how is the software industry itself evolving? Where do innovations in software come from? Empirically quantifying innovation in software is especially challenging [4], in part because the usual proxies used in the innovation literature like patents [5] or trademarks [6] are ill-suited to the software context. For instance patents tend to focus on end products and are broadly categorized; it also takes years from invention to award. They also ignore, by definition, open source software.

This is an important blindspot because of the crucial role that OSS plays in software development [7], and by extension on the digital economy; OSS ecosystems are referred to as digital infrastructure in analogy to roads and bridges [8], [9]. In particular, OSS libraries are like building blocks that

developers use and combine to build projects [10]. These blocks are codified, packaged, and made re-useable through social coding platforms like GitHub [11], creating large and growing ecosystems of interdependent code [12], [13]. The importance of these OSS building blocks in the broader economy will only grow as more products integrate digital components [14]–[18]. At the same time, researchers have raised concerns about the sustainability of these ecosystems because of their interdependencies and significant reliance on voluntary activity [19]–[22].

Indeed, these ecosystems thrive on the collaborative development and sharing of software libraries that developers can incorporate into their projects [23]. However, our understanding of the dynamics governing the introduction of new libraries and their novel combinations is surprisingly limited. This gap in knowledge hinders our ability to foster innovation effectively, allocate resources wisely, and sustain the health of OSS ecosystems. Despite recent work [10] on how developers combine libraries within projects, less is known about the overall rate of innovation at the ecosystem level. Understanding these rates is vital, as they shape the growth of dependency structures and reveal how libraries become central or peripheral to the broader ecosystem [24].

In this paper we examine the dynamics of innovation in OSS ecosystems through a quantitative study of OSS library use in Stack Overflow posts. By extracting libraries used in millions of posts written in 12 different programming languages across 15 years of data, we can observe the rates which new libraries emerge in different ecosystems. We borrow concepts from the Schumpeterian school, which frames innovation as an evolutionary process which can be meaningfully studied by observing the dynamics of elementary ingredients used in a system [25], [26]. In our framing, OSS libraries are these ingredients. A key insight from this literature, already recognized in the software engineering research community [10], is the key role that combinations of elementary ingredients play in the generation of impactful innovation [26]–[30].

We first show that across all ecosystems new libraries appear at slowing rate, suggesting that innovation within ecosystems may slow down over time. A consequence is that over time,

the distribution of library use becomes highly concentrated: for example, the 10% most frequently imported Python libraries account for about 80% of imports. On the other hand, the rate of novel pairs of libraries imported in posts grows at a linear rate in all ecosystems, suggesting that combinatorial innovation is the driver of growing ecosystems. This finding is similar to a study of over 200 years of novel codes and pairs of codes appearing on US patents [26], which the authors interpret as demonstrating a constant rate of exploration versus exploitation in innovation over the long run.

Finally we present results on the origin of innovation in software ecosystems. First we show that new users are significantly more likely to use new libraries and combinations, across all ecosystems studied. Second, after geolocating posting users, we report results on which countries are most likely to import new libraries or combinations. The geography of innovation in OSS ecosystems is surprisingly diverse.

Our results suggest that there are characteristic patterns of innovation in OSS ecosystems, and that these patterns can inform our understanding of ecosystem sustainability. New libraries emerge more slowly as ecosystems grow, concentrating use and hence risk in relatively few libraries. This also that maintenance of libraries can be prioritized. However, the steady growth of novel combinations of libraries used in code suggests that the need for maintenance to address issues emerging between ecosystem libraries (for instance through dependency links [31]–[34] or co-use relationships) will only grow. Finally, the observation that new and geographically diverse contributors are drivers of innovation in OSS ecosystems suggests the importance of dismantling barriers to OSS participation [35], [36] and facilitating onboarding [37] for the sake of long run ecosystem health.

II. THEORETICAL FRAMEWORK

We first review prior work on the dynamics of innovation, both generally and within software. We then examine how different issues in OSS maintenance relate to growing systems and their usage.

Dynamics of Innovation

Innovation systems are often analyzed by examining their elementary ingredients—fundamental units that, when combined, generate ideas, products, or technologies. For example, recent work in economic development draws an analogy between a country’s economic capabilities and the game of Scrabble, where the more letters (capabilities) a player gathers, the more high-value words (products) they can create [38].

A substantial body of research thus studies the rate at which new ingredients are introduced. Across a variety of domains, new ingredients appear at ever more slowly (i.e. sub-linear) rates. In linguistics the rate at which new words appear in a text decreases as the text grows longer [39]. Within individual software programs, program vocabulary grows sub-linearly with program length [40]. This pattern is known as Heaps’ Law [41], and has been observed in the emergence of online

social annotations (i.e. tagging systems) [42], Wikipedia pages [43], and cooking recipes [44].

One important consequence of Heaps’ Law is that the slowing rate of new ingredients in a system will naturally lead to a significant concentration of usage around specific ingredients [45]. Such distributions follow power laws or Zipf’s law [46]. In such distributions, a small number of elements are highly prevalent while the majority are infrequently used, reflecting a significant concentration of activity around certain key components. In other words, in growing systems with a slowly growing set of ingredients, some ingredients will be used many times more often than others. Such power laws appear at various levels of abstraction in software: in links between Java classes, dependencies between Perl libraries, and FreeBSD Ports, among others [47], [48].

A second, more abstract consequence of Heaps’ Law is that it suggests that either all such systems will tend to slow down over time, or that innovation is about more than just the emergence of new ingredients. Indeed, the literature since Schumpeter [25] has emphasized the importance of combinatorial innovation—the novel recombination of existing elements to create new functionalities or concepts. Indeed, if one considers new ways how ingredients can be combined, possibilities for innovation grow exponentially even if new ingredients emerge slowly [28]. The striking difference in the rates of introduction between single new components and novel combinations over time was confirmed through an analysis of 200 years of US patent data, which observed that while the rate of introducing new technological classes (individual novelties) slows down, patents registered with novel pairs of classes continue to appear at a linear rate [26].

Empirical studies across various fields support the importance of combinatorial innovation. In technology and patent research, firms innovate by recombining their existing knowledge and capabilities, resulting in new products and processes [27]. Analyses of patent data demonstrate that inventions stemming from novel combinations of existing technologies are more likely to be impactful, although they come with higher uncertainty [29]. In scientific research, high-impact papers often arise from atypical combinations of established ideas, suggesting that blending familiar concepts in novel ways can lead to significant breakthroughs [30].

What about software? Software is widely thought to be a highly innovative field [2]. Researchers interested in the sustainability of software systems have increasingly focused on open source libraries as the building blocks or infrastructure of software [8]. This suggests that OSS libraries are the ingredients of software—like the words that makes up novels, technology classes assigned to patents, or the literal ingredients used to cook a recipe—thus the right lens through which to study innovation in software. Indeed, open-source software itself in particular is considered an intensely creative activity [49], and reuse is a core concept in software engineering [23], [50].

In fact, previous work has explored innovation within individual OSS projects by examining how the adoption of novel

combinations of libraries influences project success [10]. The authors suggest an innovation-maintenance trade-off: using novel combinations of libraries enhances a project’s functionality and appeal but such combinations are more likely to be realized by smaller teams, which may be less able to maintain the project in the long run. Innovation may also be the burden of maintenance due to added complexity and dependency management. Unlike patents, where maintenance requirements are minimal after the invention is secured, software demands continuous updates and support to remain functional and secure. Despite this emerging tension between innovation and sustainability in software, we do not know much about the dynamics of innovation at the ecosystem level.

Sustainability of growing OSS ecosystems

OSS ecosystems are growing, both in terms of the number of libraries and the dependencies between them [24]. Implications for maintainability and sustainability are well-studied: issues arising because of vulnerabilities [33], [51], library abandonment [22], or update incompatibilities [31] spread beyond individual libraries through dependencies to affect larger systems.

At the same time, a better understanding of the dynamics of innovation in OSS ecosystems can inform our understanding of and barriers to their sustainability. For example, if OSS libraries appear at a sub-linear rate, it suggests that, over time, use of specific libraries will follow a highly skewed distribution. This implies that maintenance on a few specific libraries will be essential for sustainability. On the other hand, if growth in combinatorial innovation is linear, it suggests that cross-functionality of libraries is an important part of ecosystem health. Libraries used together are not necessarily near each other in the dependency tree, thus present a new perspective on links within an ecosystem.

On the other hand, innovation itself is likely an important, if overlooked aspect of ecosystem health [13]. Knowing about the rates of innovation in an OSS ecosystem, whether in terms of libraries or their combinations, can tell us if they go through different stages of maturity. If innovation is indeed important for ecosystem health, then it makes sense to know who is responsible for innovation. For instance, young people or people new to a scientific field are thought to be less beholden to established patterns, and thus be more able to innovate [52]; confirmation that new users, or that users from different geographic locations, are more likely to generate innovations in the OSS context would provide a new motivation for supporting developer onboarding and support for new contributors [35], [37], [53].

III. METHODS

Data Source

Our primary data source for this study is Stack Overflow, the leading Q&A platform for programmers. We accessed the complete Stack Overflow data dump, which includes all posts from the platform’s inception in 2008 up to 2024. This dataset

is widely utilized in research to explore various aspects of software development and developer behavior [54]–[57]. A wide variety of people, from beginner to experts are active on the platform [58]. Indeed, Stack Overflow plays a significant role in the software development community by influencing coding practices [55], [59] and even acting as a labor market signal [57]. Unlike traditional mailing lists or forums, Stack Overflow offers features such as collective curation and gamification, which enhance the quality and relevance of its content [56].

Stack Overflow is well-suited for studying innovation dynamics in OSS ecosystems for several reasons. First, its posts typically represent solutions or attempts at solutions to specific programming problems, offering practical insights into real-world challenges faced by developers. Code snippets are direct attempts to solve these problems. Unlike public OSS repos like GitHub, data mined from Stack Overflow partially reflects what developers working on closed-sourced projects are doing. Second, Stack Overflow is highly curated: duplicate and off-topic posts are systematically removed, bot activity is minimized, and community editors actively refine content. Users are encouraged to include detailed explanations, code snippets, and examples, enhancing the clarity and utility of shared knowledge. This rigorous curation ensures that the data reflects genuine developer activity. Note that we do not claim that the Stack Overflow posts that first mention a specific library are the actual moment of innovation. Rather these events are signals of innovative activity within an ecosystem.

The data dump contains over 50 million posts, including more than 20 million questions and 30 million answers, contributed by approximately 17 million registered users. For our analysis, we focus on posts related to the 12 programming languages with the highest number of posts on the platform during our study period. These languages are: Python, R, JavaScript, Java, C++, PHP, Ruby, Perl, Rust, Swift, Objective-C, and C#. This selection includes newer languages that were launched within our dataset’s timeline, such as Rust and Swift, as well as a language like Objective-C, which has reached planned obsolescence.

Data Extraction

We assign posts to languages using tags, with answers inheriting the tags of their questions. To extract code snippets from these posts, we adapt code by Baltes et al.[60]. As a post may be tagged with multiple programming languages, it is possible that an individual post is scanned for language-specific libraries multiple times.

Given a set of code snippets within a post, we extract libraries using language-specific regular expressions. If a post contains code in another language (for example if a post is tagged with multiple languages), it is unlikely that the regular expression would extract any libraries. We only retain posts that have at least one import statement. We show an example question from Stack Overflow, tagged with Python, and including a code snippet with import statements.

Our regular expressions handle multiple ways that users can import libraries in the different programming languages.

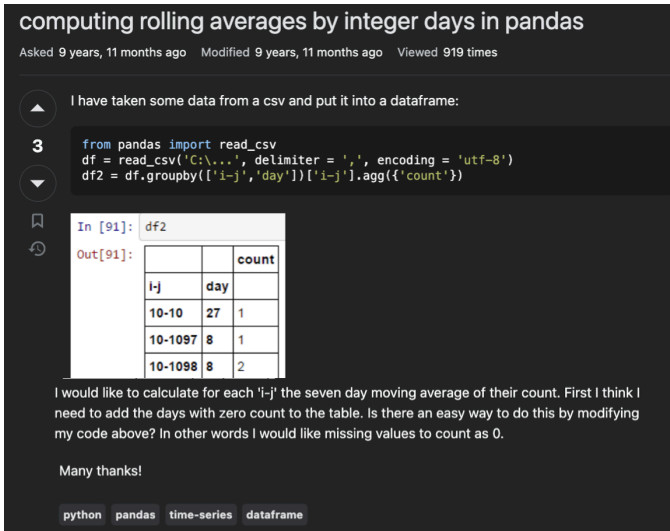


Fig. 1: Stack Overflow post tagged with the Python programming language and containing a code snippet with a library import.

For example, in the Python programming language there are multiple valid ways to import libraries:

```

1 import pandas as pd
2 from numpy import array
3 import glob
4 import requests, math

```

Checking validity

Stack Overflow code snippets are generally not valid programs, and there is no guarantee that import statements we detect refer to actual libraries. For example, a post may refer to libraries “foo”, “bar”, and “foobar” as part of an illustrative example. Our subsequent analyses rely on a relatively accurate identification of real imports of libraries. Thus, we checked how many library imports we observed in Python posts correspond to actual libraries recorded in the Python Package Index (PyPI) and the list of Python standard library packages (extracted for Python 2.7.18 and Python 3.12.6). We find that 85% of the Python libraries imported refer to libraries in the canonical list. We assume that the rate of error is similar across other ecosystems.

Quantifying novelty

We consider a library import a novelty in the respective programming language the first time it appears in a post. We consider two libraries a novel pair the first time they co-occur in a post. Note that, in order to form a novel pair, neither of the imported libraries has to be a novelty. For example, in the hypothetical scenario below, the posts contain 2, 1, and 0 novel libraries, respectively, assuming the posts appear in left-to-right order and no prior posts import any of the mentioned libraries. All of the listed posts introduce a novel pair.

In our analyses of novelty, we filter out libraries that were imported less than 10 times in subsequent posts. Our

```

1 import os
2 import sys

```

Post 1

```

1 import random
2 import sys

```

Post 2

```

1 import os
2 import random

```

Post 3

Fig. 2: Toy example of three sequential posts with import statements in the Python programming language. The first post contains two novel imports, the second contains one, and the third post contains none. The first post contains one novel combination of imports (os, sys), the second contains one novel combination (random, sys), and the third also contains one novel combination (os, random).

decision is driven by the goal of avoiding the analysis of small personal projects that did not achieve widespread adoption as innovations. Similarly, we only consider pairs of libraries a novel pair in our analyses if both libraries of the pair were imported at least 10 times throughout our dataset.

User characteristics

We consider two characteristics of posting users that may relate to their likelihood of posting novelties: their experience and geographic location. We define a user’s experience at the time of a post by the number of posts they had made previously in the programming language of the post. To infer a posting user’s country, we geocode the location they provide in their profile pages using the Microsoft Bing Maps API, an accurate tool which handles multiple languages and representations of locations [61]. We infer a country for 38% of posting users.

IV. RESULTS

In this section we present our findings. First we show how simple and combinatorial novelties emerge per post in each of the twelve ecosystems. We then report the concentration of library use across posts. Finally, we report on the relationships between posting user characteristics (experience and location) and novelties.

Novelties

In Figure 3 we report data on novel library imports and combinations of library imports per post. Novelties in all of the studied programming languages followed a sub-linear growth pattern. In other words, as these systems grow, fewer new libraries are introduced for the same level of activity. At the same time, the rate of combinatorial novelties grows at a remarkably steady rate in all ecosystems: linear fits are highly accurate. These two patterns replicate earlier findings on the introduction of new patent classes in a dataset of over 200 years of patents from the US [26]. There too, new classes emerged at an ever-slower rate, while novel pairs of classes emerged at a constant rate per patent across the entire dataset.

This finding has a few interesting implications. The first is that innovation in mature OSS ecosystems happens by combining libraries. Second, within an ecosystem, the rate

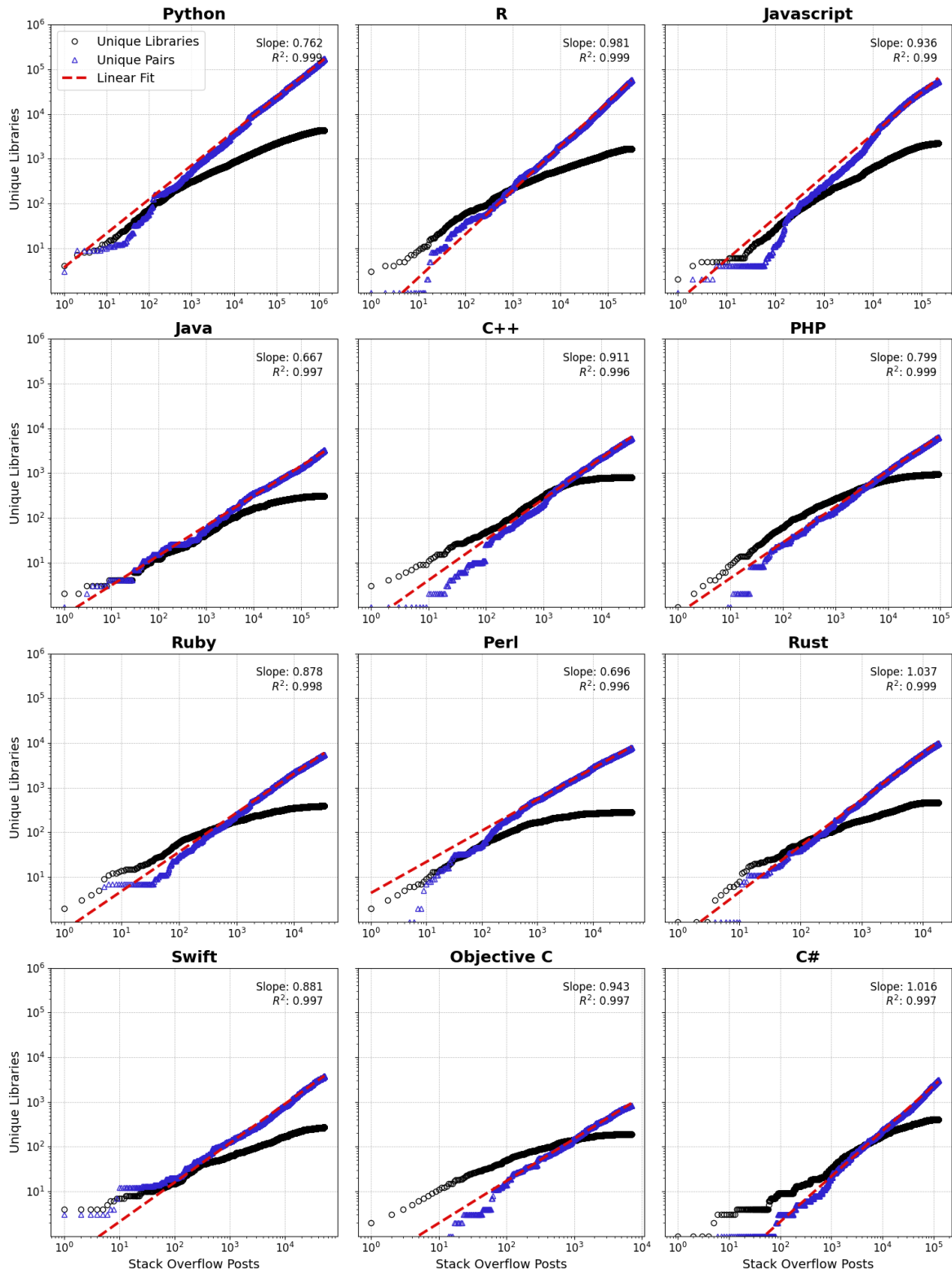


Fig. 3: Rates of simple and combinatorial novelties of library imports in Stack Overflow posts across programming languages.

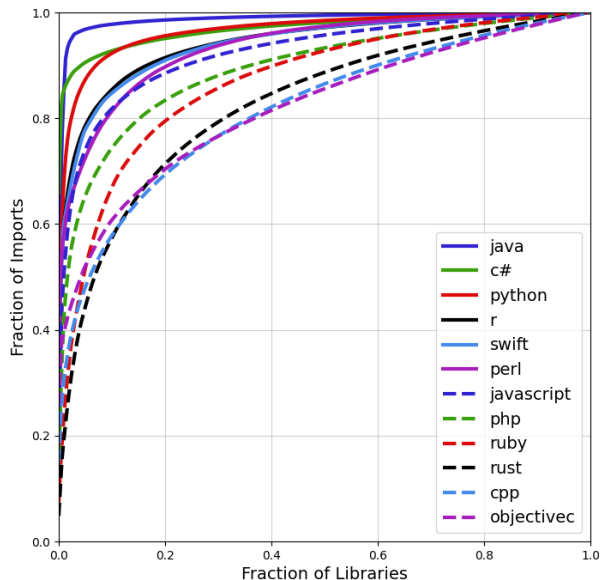


Fig. 4: The fraction of total imports in specific programming languages made of the fraction of the most commonly imported libraries. For example, imports of the 7% most frequently imported Python libraries account for 90% of all imports.

at which users explore new combinations of libraries or exploit existing combinations is roughly constant over time. The third implication, which we explore in greater detail in the subsequent subsection, is that, relative to the number of posts made, the number of libraries in an ecosystem is relatively small. In other words, we expect use of libraries to be significantly concentrated in a few key libraries.

Concentration of use

Given that the rate of new libraries introduced per post universally declines across all programming languages we study, we suspect that the use of libraries across all posts is highly concentrated [45]. To examine this hypothesized concentration, we draw Pareto curves [46] in Figure 4. These curves are often used to visualize the distribution of wealth in an economy, and are related to the “80-20” rule.

In our context, they report the cumulative share of all imports made in an ecosystem by individual libraries, recorded in descending order by number of imports. If each library were imported an equal number of times, we would expect a diagonal ($y = x$) line. If all imports were concentrated in a single library, we would observe a vertical line from (0,0) to (0,1), then a horizontal line from (0,1) to (1,1). Generally, the closer the observed curve is to the top left corner of the plot, the more concentrated the imports are among a few libraries.

Figure 4 confirms that imports are intensely concentrated among the top libraries in most ecosystems. For example, 90%

Language	Libraries	Imports	50%	80%	90%
Python	4,366	2,469,502	0.4	2.4	7.0
R	1,669	528,278	0.4	5.8	16.6
Objective-C	192	8,767	4.2	37.0	62.5
Javascript	2,205	413,073	1.0	8.4	23.8
Java	310	508,724	0.0	1.0	1.3
C++	801	44,341	5.9	36.0	59.9
PHP	936	106,829	1.7	15.3	36.1
Ruby	385	48,770	4.7	20.8	40.8
Perl	279	119,065	0.4	9.3	20.8
Rust	465	30,995	6.9	31.2	54.0
Swift	268	83,549	0.7	6.3	18.3
C#	409	169,781	0.0	0.5	4.4

TABLE I: Library import frequencies per programming language ecosystem. We report the unique number of libraries, frequency of their imports, and the share of top libraries responsible for 50%, 80% and 90% of imports. For example, the top 0.4% of libraries account for 50% of all imports in Python posts to Stack Overflow.

of all imports in Python are of the 7% most frequently imported libraries. To facilitate interpretation, we report several similar statistics for all language ecosystems in Table I.

These results suggest that a few libraries are used to solve most problems in any given ecosystem. This has consequences for maintenance of these ecosystems. On the one hand, attention and resources can effectively be concentrated on the most used libraries. On the other hand, errors and defects in these key libraries can have a very broad impact on end users. This analysis of library use complements previous work on the distribution of technical dependencies between libraries, which similarly suggests that a few key libraries are highly depended on [24].

Users

In our final set of analyses, we focus on user attributes which correlate with innovation. We first investigate how prior user experience in a language ecosystem correlates with the likelihood that a user will post a novel library or pair of libraries. On the one hand, new users may bring new perspectives, problems, and solutions to an ecosystem. On the other hand, experienced users may be the most knowledgeable about libraries available in an ecosystem, and may be dealing with more complex problems that require highly specialized libraries or combinations. If the former effect dominates, we would expect new users to be more likely to import novel libraries. If the latter is true, we would expect to see novelties appear more frequently on posts made by expert users.

To evaluate these competing predictions, for each post we quantify user experience as the prior number of posts made by that user in that ecosystem. We log-bin users into five categories of experience and plot the likelihood that a post by that user contains a novel library or pair in Figure 5. This plot shows estimates pooled across all language ecosystem; we find similar results in all individual ecosystems, reported in the Appendix.

For both kinds of novelty we observe a striking pattern: new users are significantly more likely to import novel libraries and

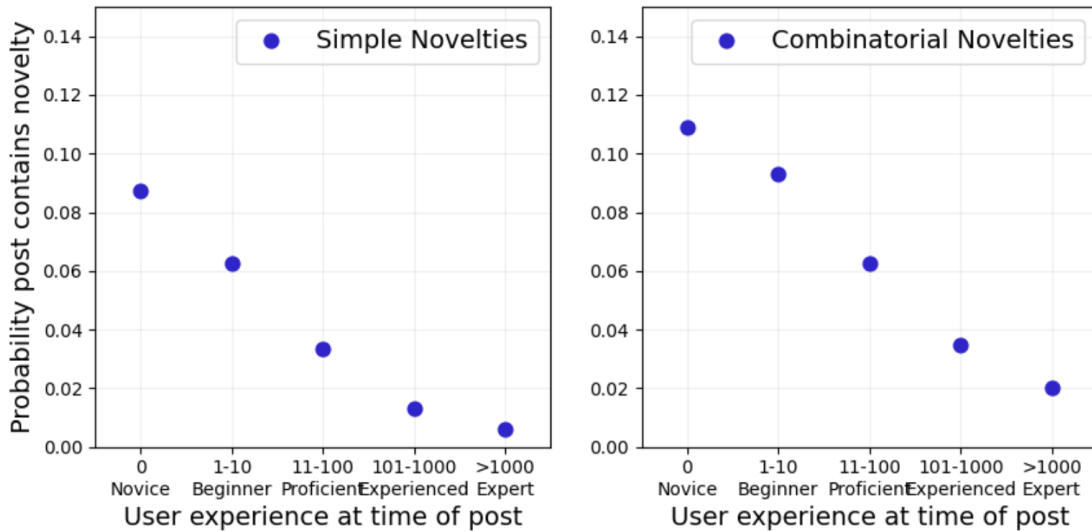


Fig. 5: The likelihood that a post contains a novel library or pair of libraries, respectively, as a function of the user’s previous posting experience at the time of the post. Note that we define user experience and novelty of posts at the programming language level, reporting pooled estimates here.

combinations of novel libraries. Across all languages pooled together, a beginner (having made 1-10 posts previously) is about four times more likely to make a post with a new library than an experienced user (101-1000 prior posts), and three times more likely to make a post with a combinatorial novelty. In other words, new users are introducing new ideas to OSS ecosystems. This finding also has implications for ecosystem sustainability, suggesting that the onboarding and inclusion of new contributors and users pays off in terms of innovation outcomes.

We check the robustness of this finding in a variety of ways. We first dropped novel libraries and pairs of libraries if those libraries or pairs were not used frequently in subsequent posts. Indeed new users may be more likely to import unsophisticated, less useful, or incoherent libraries or combinations of libraries. However, we found no substantive difference when keeping only novelties that were more broadly adopted (used in posts at least 1,000 times). We also found a similar pattern focusing only on posts made in specific years, alleviating the concern that our results are an artifact of the concentration of beginner status users and novelties in posts at the beginning of our dataset.

Finally, the second source of user-level heterogeneity in innovation we investigated was geography. In Figure 6 we report the likelihood that posts made by users from different countries contain novelties. Note that we report countries which we can associate to at least 1000 posts. We observe heterogeneities between countries: the rate of both simple and combinatorial novelty is several times greater in the leading countries versus the lagging countries. At the same time, we observe significant geographic diversity in signals of rates of innovation: no small, coherent group of countries dominates the others. This suggests that despite the previously

documented spatial concentration of OSS activity [61], new ideas can come from anywhere.

V. DISCUSSION

Our study investigates the dynamics of innovation within OSS ecosystems, revealing several key findings. We observed that the introduction of new libraries follows a sub-linear growth pattern, indicating a slowdown in creating entirely new libraries as ecosystems mature. Conversely, combinatorial innovation—where existing libraries are reused and combined in novel ways—exhibits steady linear growth. This trend suggests that mature OSS ecosystems tend to innovate by combining existing components rather than developing new foundational ones. Our analysis highlights a significant concentration in library usage, with a few key libraries becoming highly integrated across various projects. We found that innovation was more likely to be found in posts by newer users. Finally, the geography of innovative users was found to be quite diverse.

Here we discuss implications of our findings, focus on two aspects: OSS ecosystem sustainability and the dynamics of innovation in OSS compared with other systems. We then elaborate on the limitations of our work and suggest directions for future research.

a) Implications for OSS Ecosystem Sustainability: The concentrated use of certain libraries has both positive and negative implications for the sustainability of OSS ecosystems [22], [24]. On the positive side, heavily utilized core libraries can enhance stability, benefiting from robust community support and well-maintained codebases [19]. This shared foundation fosters consistent practices among developers. However, the reliance on a few key libraries also poses risks. If a widely used library becomes deprecated or encounters significant issues, it could threaten the entire ecosystem [33]. Another challenge our findings present is that combinations of libraries

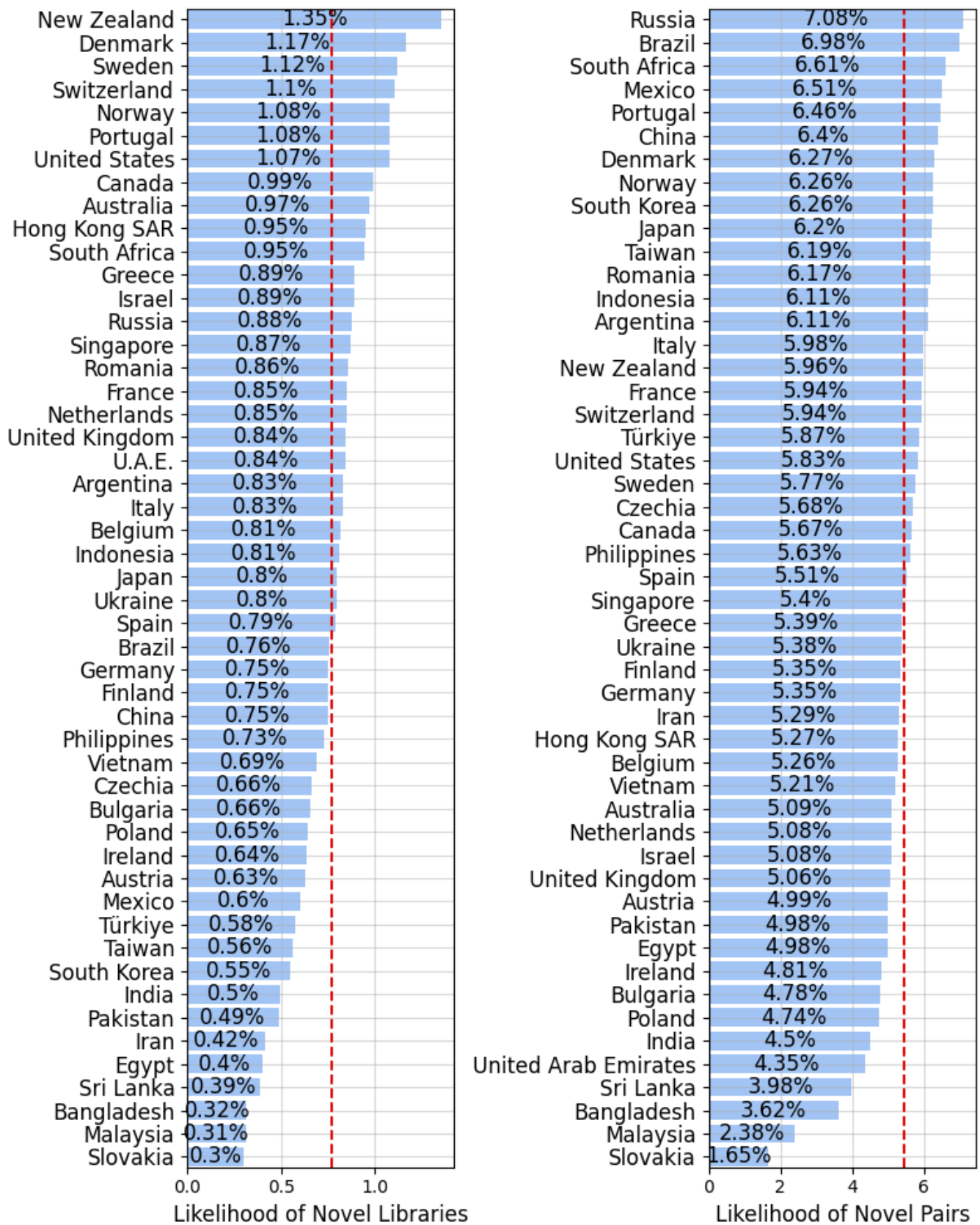


Fig. 6: Rates of novelty and combinatorial novelty in posts by users geolocated to different countries. Red vertical lines indicate global averages.

are used at a steadily growing rate in all ecosystems. Hence, ensuring that combinations of libraries, which may not depend on each other in a formal sense, work together is an especially daunting task in light of this growth pattern.

At the user level, we showed that new contributors play a pivotal role in driving innovation within OSS ecosystems [37], [53]. They often introduce novel combinations of libraries and experiment with approaches that may not occur to seasoned contributors [10]. This underscores the importance of community practices that lower barriers to entry and encourage diverse participation [35], [36]. Improving documentation, streamlining onboarding processes, and offering mentorship can help harness the innovative potential of newcomers [37]. Creating an inclusive and supportive environment where new ideas are valued is crucial for sustaining their contributions over time [36], [53].

Our findings also reveal substantial geographical diversity among OSS library contributors. Signals of new library and library-combination use come from all over the world, indicating that innovation in use is not confined to traditional software hubs [7]. This diversity brings a range of perspectives and problem-solving approaches, enhancing the resilience and adaptability of OSS ecosystems. However, our findings refer to the rates of innovation as a function of activity, and geographic disparities in activity certainly exist [35]. This suggests that increasing participation in and use of OSS can spur innovation in the overall system. Addressing this gap requires targeted outreach and support initiatives to empower contributors from emerging markets and underrepresented areas [36], [37].

In sum, our work suggests a few practical ideas for improving OSS ecosystem sustainability. First, prioritizing support for widely used libraries is crucial to ensure they remain well-maintained and resilient; this can be achieved through funding and community-driven initiatives [8], [19]; moreover, maintenance should also be understood as supporting the combinations of libraries. Second, fostering the involvement of new contributors is essential for ongoing innovation [37], [53]. Implementing structured onboarding programs, providing mentorship, and offering resources to guide newcomers can facilitate their integration and contribution [35], [37].

Innovation in OSS and other systems: We have shown that the patterns of innovation observed in OSS ecosystems are similar to those in other fields like technological patents and scientific research, where significant advances often result from novel recombinations of existing elements rather than entirely new creations [26], [28]. This parallel suggests that insights from evolutionary theories of the economy may be applied to understand OSS ecosystems [62]. In that tradition, firms adapt and imitate each other, driven by processes of variation, selection, and retention that influence their competitive advantage and sustainability. Similarly, OSS ecosystems evolve through cycles of new library introductions, novel code combinations, and the community’s selection of effective tools and practices. While software engineering literature has explored evolutionary concepts and the development of software ecosystems, it often emphasizes technical growth,

change propagation, and dependency management [63] or within-library evolution [64].

For instance, in the context of patents we know that social and collaboration network positions of individuals influence their ability to combine ingredients in a productive way [27], and that more complex and unorthodox combinations lead to higher variations in outcomes [29]. Similar studies of academic research papers suggests that actually realized combinations of ideas tend to be repetitive and imitative, and that the most successful papers tend to mix in just a few surprising ingredients [30]. Indeed recent work in software engineering references this tradition showing combinatorial innovation of libraries within software projects predicts success outcomes [10]. Likewise, collaboration network topology within software ecosystems, like in patenting, play a key role in the spread and adoption of new ingredients (libraries) [65].

Limitations and Future Work

While our study provides insights into the dynamics of innovation within OSS ecosystems, it is important to acknowledge its limitations, particularly regarding our reliance on Stack Overflow as the primary data source [60]. Stack Overflow offers a rich repository of developer interactions, which we utilize as signals to infer rates of innovation and the underlying dynamics of library usage [59]. However, we do not consider the specific posts or discussions on Stack Overflow as direct instances of innovation themselves. Instead, they serve as proxies indicating how developers engage with and combine existing libraries.

This reliance on a single platform introduces potential biases. Stack Overflow tends to highlight issues and topics that generate widespread interest or challenges among developers, possibly overlooking smaller-scale innovations or niche libraries that do not elicit as much public discussion [59]. Consequently, our analysis may underrepresent innovative activities occurring in less-discussed areas of the OSS ecosystem. Additionally, not all developers use Stack Overflow uniformly; some may prefer other forums, private communications, or may not seek external help at all, leading to gaps in the data captured.

To address these limitations, future research should incorporate a more diverse set of data sources. Direct analysis of code repositories, commit histories, package managers, and issue trackers could provide a more comprehensive view of library usage and innovation patterns. These sources can capture actual coding activities and library integrations that may not be discussed publicly. Combining multiple data sources would help mitigate the biases inherent in relying solely on Stack Overflow.

Furthermore, our study focuses on observable signals of innovation and may not fully capture the qualitative aspects of how developers innovate. Understanding the motivations, decision-making processes, and collaborative efforts behind library usage requires qualitative methods such as surveys or interviews. Incorporating these approaches could enrich the

analysis and provide deeper insights into the factors driving innovation within OSS ecosystems.

Lastly, the evolving nature of software development practices, including the rise of artificial intelligence and automated code generation tools, may influence innovation dynamics in ways not accounted for in our study [66]. Future investigations should explore how these emerging technologies impact the rates and patterns of innovation, potentially accelerating growth or introducing new forms of collaboration that reshape the OSS ecosystem.

VI. CONCLUSION

Software, especially OSS, is often thought of as a highly innovative and creative sector [49]. Here we show that the dynamics of innovation in OSS ecosystems are quite similar to the dynamics of innovation in patenting and scientific research. The implications, both for software maintenance and sustainability, and for our general understanding of OSS ecosystem evolution are numerous.

VII. DATA AND CODE AVAILABILITY

The Stack Overflow data dump is available at the Internet Archive: <https://archive.org/details/stackexchange>. Processed data (on post-level library use) is available at: <https://zenodo.org/uploads/14186439>. Code to replicate our analyses is available here: https://github.com/MeszarosGabor/SO_Post_Analyzer.

VIII. ACKNOWLEDGMENTS

We thank participants of the Danube Workshop on Software and the Digital Economy for helpful comments and suggestions. JW acknowledges support from the Hungarian National Scientific Fund (OTKA FK 145960). GM is funded by the European Union under Horizon EU project LearnData (101086712). The authors acknowledge use of the HUN-REN Cloud [67] in the “Geographies of Creation, Learning and Use in Software” project.

REFERENCES

- [1] M. Andreessen, “Why software is eating the world,” *The Wall Street Journal*, vol. 8, p. 20, 2011.
- [2] B. Chattergoon and W. R. Kerr, “Winner takes all? tech clusters, population centers, and the spatial transformation of us invention,” *Research Policy*, vol. 51, no. 2, p. 104418, 2022.
- [3] R. N. Charette, “How software is eating the car,” *IEEE Spectrum*, 2021.
- [4] N. bin Ali, H. Edison, and R. Torkar, “The impact of a proposal for innovation measurement in the software industry,” in *Proceedings of the 14th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*, 2020, pp. 1–6.
- [5] E. Mansfield, “Patents and innovation: An empirical study,” *Management science*, vol. 32, no. 2, pp. 173–181, 1986.
- [6] C. Castaldi, “To trademark or not to trademark: The case of the creative and cultural industries,” *Research Policy*, vol. 47, no. 3, pp. 606–616, 2018.
- [7] J. Musseau, J. S. Meyers, G. P. Sieniawski, C. A. Thompson, and D. German, “Is open source eating the world’s software? measuring the proportion of open source in proprietary software using java binaries,” in *Proceedings of the 19th International Conference on Mining Software Repositories*, 2022, pp. 561–565.
- [8] N. Eghbal, *Roads and bridges: The unseen labor behind our digital infrastructure*. Ford Foundation, 2016.
- [9] N. Eghbal, *Working in public: the making and maintenance of open source software*. Stripe Press, 2020.
- [10] H. Fang, J. Herbsleb, and B. Vasilescu, “Novelty begets popularity, but curbs participation—a macroscopic view of the python open-source ecosystem,” in *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, 2024, pp. 1–11.
- [11] L. Dabbish, C. Stuart, J. Tsay, and J. Herbsleb, “Social coding in github: Transparency and collaboration in an open software repository,” in *Proceedings of the ACM 2012 conference on computer supported cooperative work*, 2012, pp. 1277–1286.
- [12] T. Mens, M. Claes, P. Grosjean, and A. Serebrenik, “Studying evolving software ecosystems based on ecological models,” *Evolving software systems*, pp. 297–326, 2014.
- [13] T. Mens, C. De Roover, and A. Cleve, *Software Ecosystems: Tooling and Analytics*. Springer Nature, 2023.
- [14] M. Hoffmann, F. Nagle, and Y. Zhou, “The value of open source software,” *Harvard Business School Strategy Unit Working Paper*, no. 24-038, 2024.
- [15] G. Korkmaz, J. B. S. Calderón, B. L. Kramer, L. Guci, and C. A. Robbins, “From github to gdp: A framework for measuring open source software innovation,” *Research Policy*, vol. 53, no. 3, p. 104954, 2024.
- [16] J. Gortmaker, “Open source software policy in industry equilibrium,” Working paper, Tech. Rep., 2024.
- [17] V. Stojkoski, P. Koch, E. Coll, and C. A. Hidalgo, “Estimating digital product trade through corporate revenue data,” *Nature Communications*, vol. 15, no. 1, p. 5262, 2024.
- [18] S. Juhász, J. Wachs, J. Kaminski, and C. A. Hidalgo, “The software complexity of nations,” *arXiv preprint: 2407.13880*, 2024.
- [19] G. Avelino, E. Constantinou, M. T. Valente, and A. Serebrenik, “On the abandonment and survival of open source projects: An empirical investigation,” in *2019 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*, IEEE, 2019.
- [20] R. S. Geiger, D. Howard, and L. Irani, “The labor of maintaining and scaling free and open-source software projects,” *Proceedings of the ACM on human-computer interaction*, vol. 5, no. CSCW1, pp. 1–28, 2021.
- [21] J. Hejderup, “On the use of tests for software supply chain threats,” in *Proceedings of the 2022 ACM Workshop on Software Supply Chain Offensive Research and Ecosystem Defenses*, 2022, pp. 47–49.
- [22] C. Miller, C. Kästner, and B. Vasilescu, ““we feel like we’re winging it:” a study on navigating open-source dependency abandonment,” in *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2023, pp. 1281–1293.
- [23] V. R. Basili, L. C. Briand, and W. L. Melo, “How reuse influences productivity in object-oriented systems,” *Communications of the ACM*, vol. 39, no. 10, pp. 104–116, 1996.
- [24] A. Decan, T. Mens, and P. Grosjean, “An empirical comparison of dependency network evolution in seven software packaging ecosystems,” *Empirical Software Engineering*, vol. 24, no. 1, pp. 381–416, 2019.
- [25] J. A. Schumpeter, “Business cycles: A theoretical, historical and statistical analysis of the capitalist process,” 1939.
- [26] H. Youn, D. Strumsky, L. M. Bettencourt, and J. Lobo, “Invention as a combinatorial process: Evidence from us patents,” *Journal of the Royal Society interface*, vol. 12, no. 106, p. 20150272, 2015.
- [27] B. Kogut and U. Zander, “Knowledge of the firm, combinative capabilities, and the replication of technology,” *Organization Science*, 1992.
- [28] M. L. Weitzman, “Recombinant growth,” *The Quarterly Journal of Economics*, vol. 113, no. 2, pp. 331–360, 1998.
- [29] L. Fleming, “Recombinant uncertainty in technological search,” *Management science*, vol. 47, no. 1, pp. 117–132, 2001.
- [30] B. Uzzi, S. Mukherjee, M. Stringer, and B. Jones, “Atypical combinations and scientific impact,” *Science*, 2013.
- [31] A. Decan, T. Mens, and M. Claes, “An empirical comparison of dependency issues in oss packaging ecosystems,” in *2017 IEEE 24th international conference on software analysis, evolution and reengineering (SANER)*, IEEE, 2017, pp. 2–12.
- [32] I. Pashchenko, H. Plate, S. E. Ponta, A. Sabetta, and F. Massacci, “Vulnerable open source dependencies: Counting those that matter,” in *Proceedings of the 12th ACM/IEEE international symposium on empirical software engineering and measurement*, 2018, pp. 1–10.
- [33] M. Zimmermann, C.-A. Staicu, C. Tenny, and M. Pradel, “Small world with high risks: A study of security threats in the npm ecosystem,” in *28th USENIX Security symposium (USENIX security 19)*, 2019, pp. 995–1010.

- [34] W. Schueller and J. Wachs, "Modeling interconnected social and technical risks in open source software ecosystems," *Collective Intelligence*, vol. 3, no. 1, p. 26 339 137 241 231 912, 2024.
- [35] I. Steinmacher, T. Conte, M. A. Gerosa, and D. Redmiles, "Social barriers faced by newcomers placing their first contribution in open source software projects," in *Proceedings of the 18th ACM conference on Computer supported cooperative work & social computing*, 2015, pp. 1379–1392.
- [36] C. Mendez, H. S. Padala, Z. Steine-Hanson, et al., "Open source barriers to entry, revisited: A sociotechnical perspective," in *Proceedings of the 40th International conference on software engineering*, 2018, pp. 1004–1015.
- [37] B. Trinkenreich, M. Guizani, I. Wiese, A. Sarma, and I. Steinmacher, "Hidden figures: Roles and pathways of successful OSS contributors," *Proceedings of the ACM on human-computer interaction*, vol. 4, no. CSCW2, pp. 1–22, 2020.
- [38] R. Hausmann, C. A. Hidalgo, S. Bustos, M. Coscia, and A. Simoes, *The atlas of economic complexity: Mapping paths to prosperity*. MIT Press, 2014.
- [39] R. H. Baayen, *Word frequency distributions*. Springer Science & Business Media, 2012, vol. 18.
- [40] H. Zhang, "Discovering power laws in computer programs," *Information processing & management*, vol. 45, no. 4, pp. 477–483, 2009.
- [41] H. Heaps, *Information retrieval: Computational and theoretical aspects*, 1978.
- [42] C. Cattuto, A. Barrat, A. Baldassarri, G. Schehr, and V. Loreto, "Collective dynamics of social annotation," *Proceedings of the National Academy of Sciences*, vol. 106, no. 26, pp. 10 511–10 515, 2009.
- [43] F. Tria, V. Loreto, V. D. P. Servedio, and S. H. Strogatz, "The dynamics of correlated novelties," *Scientific Reports*, vol. 4, no. 1, p. 5890, 2014.
- [44] Y.-X. Zhu, J. Huang, Z.-K. Zhang, Q.-M. Zhang, T. Zhou, and Y.-Y. Ahn, "Geography and similarity of regional cuisines in china," *PLoS one*, vol. 8, no. 11, e79161, 2013.
- [45] F. Tria, V. Loreto, and V. D. Servedio, "Zipf's, heaps' and taylor's laws are determined by the expansion into the adjacent possible," *Entropy*, vol. 20, no. 10, p. 752, 2018.
- [46] M. E. Newman, "Power laws, pareto distributions and zipf's law," *Contemporary physics*, vol. 46, no. 5, pp. 323–351, 2005.
- [47] S. Valverde, R. F. Cancho, and R. V. Sole, "Scale-free networks from optimal design," *Europhysics Letters*, vol. 60, no. 4, p. 512, 2002.
- [48] P. Louridas, D. Spinellis, and V. Vlachos, "Power laws in software," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 18, no. 1, pp. 1–26, 2008.
- [49] K. R. Lakhani and R. G. Wolf, *Why hackers do what they do: Understanding motivation and effort in free/open source software projects*. MIT Press, 2005.
- [50] W. B. Frakes and K. Kang, "Software reuse research: Status and future," *IEEE Transactions on Software Engineering*, vol. 31, no. 7, pp. 529–536, 2005.
- [51] A. Decan, T. Mens, and E. Constantinou, "On the impact of security vulnerabilities in the npm package dependency network," in *Proceedings of the 15th international conference on mining software repositories*, 2018, pp. 181–191.
- [52] A. Dietrich and N. Srinivasan, "The optimal age to start a revolution," *The Journal of Creative Behavior*, vol. 41, no. 1, pp. 54–74, 2007.
- [53] M. Gerosa, I. Wiese, B. Trinkenreich, et al., "The shifting sands of motivation: Revisiting what drives contributors in open source," in *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, IEEE, 2021, pp. 1046–1058.
- [54] A. Anderson, D. Huttenlocher, J. Kleinberg, and J. Leskovec, "Discovering value from community activity on focused question answering sites: a case study of Stack Overflow," in *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2012.
- [55] B. Vasilescu, V. Filkov, and A. Serebrenik, "Stackoverflow and github: Associations between software development and crowdsourced knowledge," in *2013 International conference on social computing*, IEEE, 2013, pp. 188–195.
- [56] B. Vasilescu, A. Serebrenik, P. Devanbu, and V. Filkov, "How social q&a sites are changing knowledge sharing in open source software communities," in *Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing*, 2014, pp. 342–354.
- [57] L. Xu, T. Nian, and L. Cabral, "What makes geeks tick? a study of stack overflow careers," *Management Science*, vol. 66, no. 2, pp. 587–604, 2020.
- [58] S. L. Vadlamani and O. Baysal, "Studying software developer expertise and contributions in stack overflow and github," in *2020 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, IEEE, 2020, pp. 312–323.
- [59] Y. Wu, S. Wang, C.-P. Bezemer, and K. Inoue, "How do developers utilize source code from stack overflow?" *Empirical Software Engineering*, vol. 24, pp. 637–673, 2019.
- [60] S. Baltes, C. Treude, and S. Diehl, "Sotorrent: Studying the origin, evolution, and usage of stack overflow code snippets," in *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*, IEEE, 2019, pp. 191–194.
- [61] J. Wachs, M. Nitecki, W. Schueller, and A. Polleres, "The geography of open source software: Evidence from github," *Technological Forecasting and Social Change*, vol. 176, p. 121 478, 2022.
- [62] R. R. Nelson and S. G. Winter, *An evolutionary theory of economic change*. Harvard University Press, 1985.
- [63] T. Mens, A. Serebrenik, A. Cleve, et al., *Evolving Software Systems*. Springer, 2014, vol. 190.
- [64] C. F. Kemerer and S. Slaughter, "An empirical approach to studying software evolution," *IEEE transactions on software engineering*, vol. 25, no. 4, pp. 493–509, 1999.
- [65] H. Fang, P. Park, J. Evans, J. Herbsleb, and B. Vasilescu, *Weak ties explain open source innovation*, 2024. arXiv: 2411.05646 [cs.SE]. [Online]. Available: <https://arxiv.org/abs/2411.05646>.
- [66] R. M. del Rio-Chanona, N. Laurentsyeva, and J. Wachs, "Large language models reduce public knowledge sharing on online q&a platforms," *PNAS Nexus*, vol. 3, no. 9, p. 400, 2024.
- [67] M. Héder, E. Rigó, D. Medgyesi, et al., "The past, present and future of the elkh cloud," *INFORMÁCIÓS TÁRSADALOM: TÁRSADALOM-TUDOMÁNYI FOLYÓIRAT*, vol. 22, no. 2, pp. 128–137, 2022.

APPENDIX

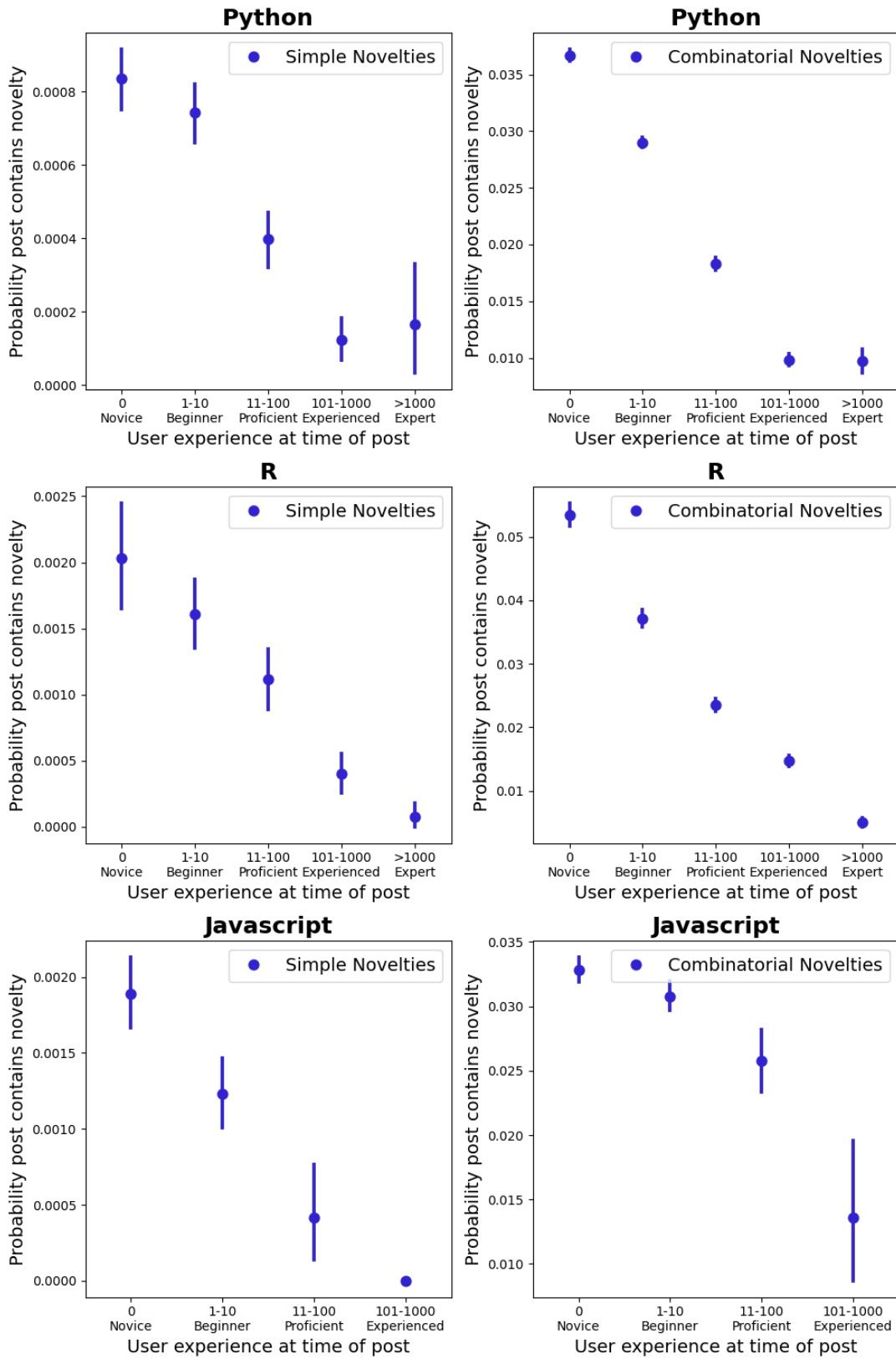


Fig. 7: User Experience Distribution for Simple and Combinatorial Novelties per Programming Languages I.

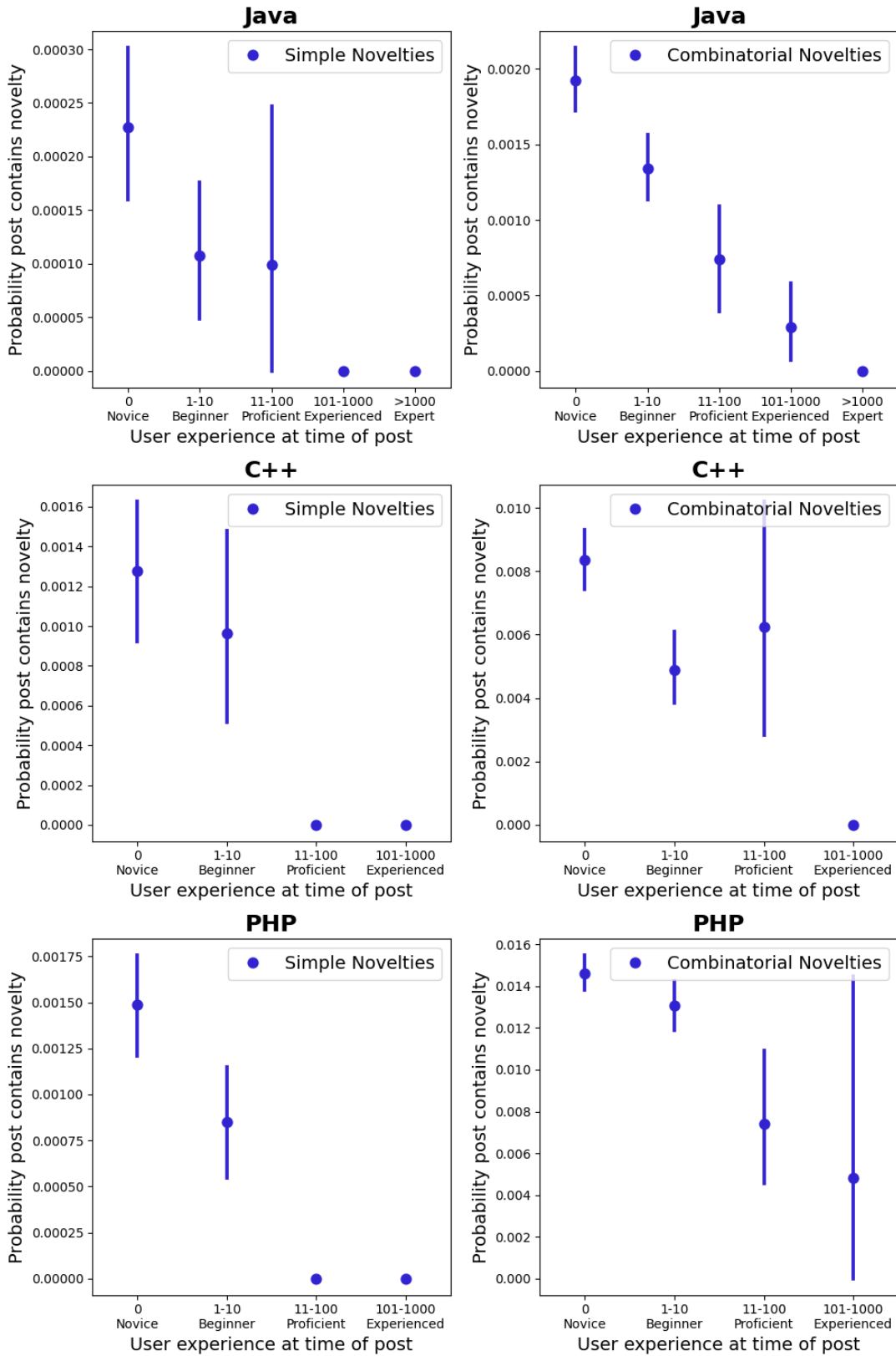


Fig. 8: User Experience Distribution for Simple and Combinatorial Novelties per Programming Languages II.

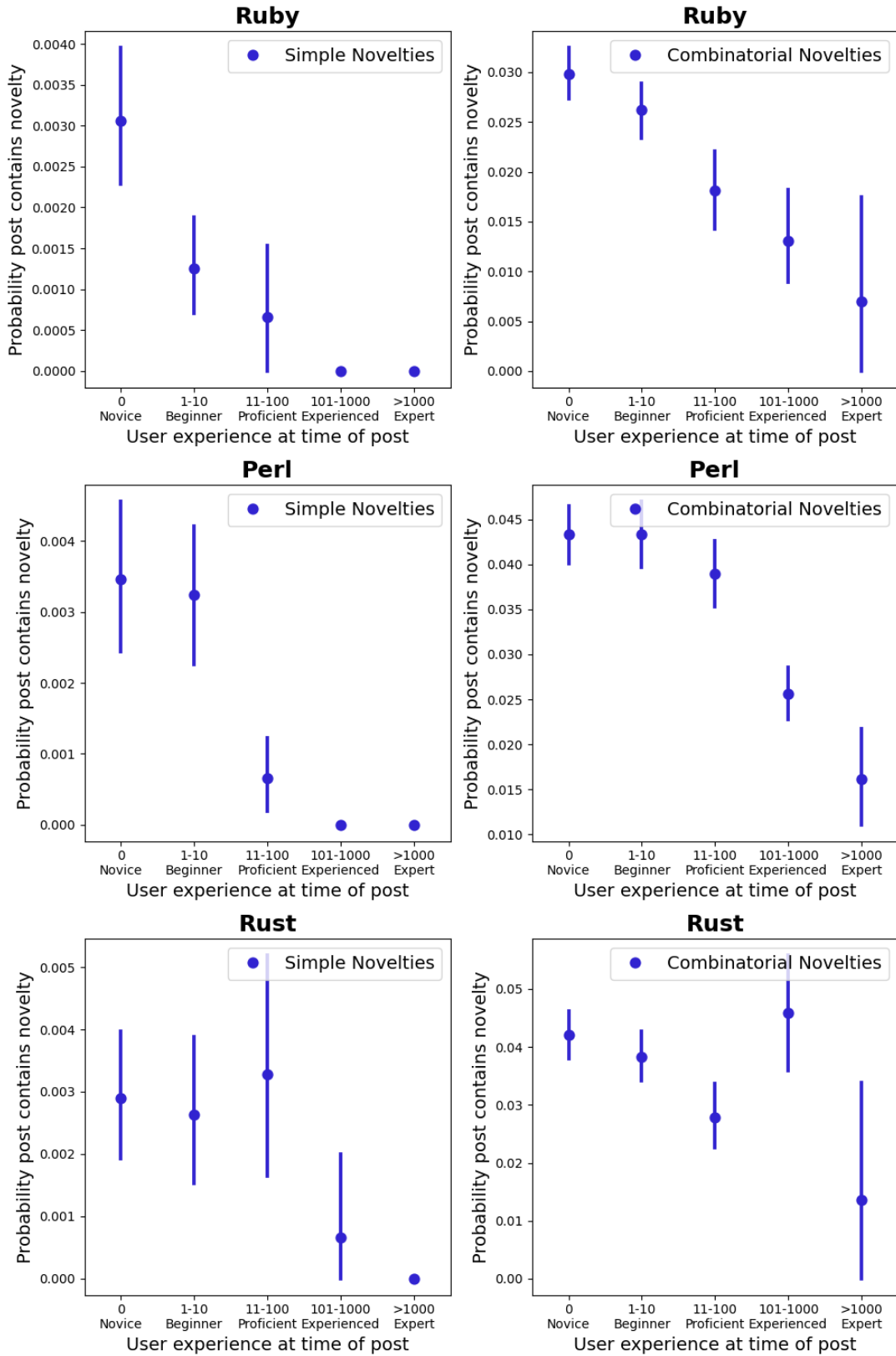


Fig. 9: User Experience Distribution for Simple and Combinatorial Novelties per Programming Languages III.

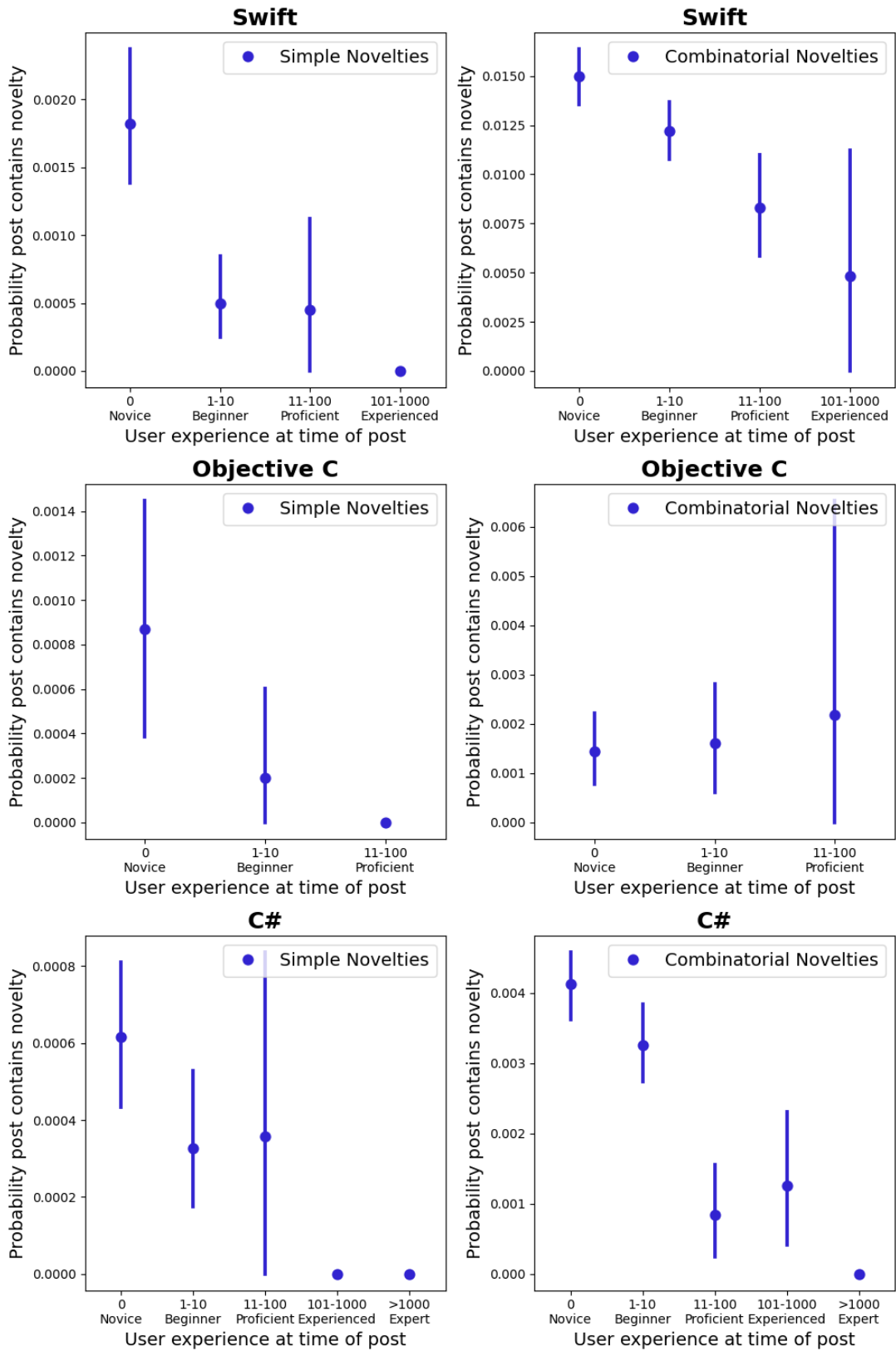


Fig. 10: User Experience Distribution for Simple and Combinatorial Novelties per Programming Languages IV.

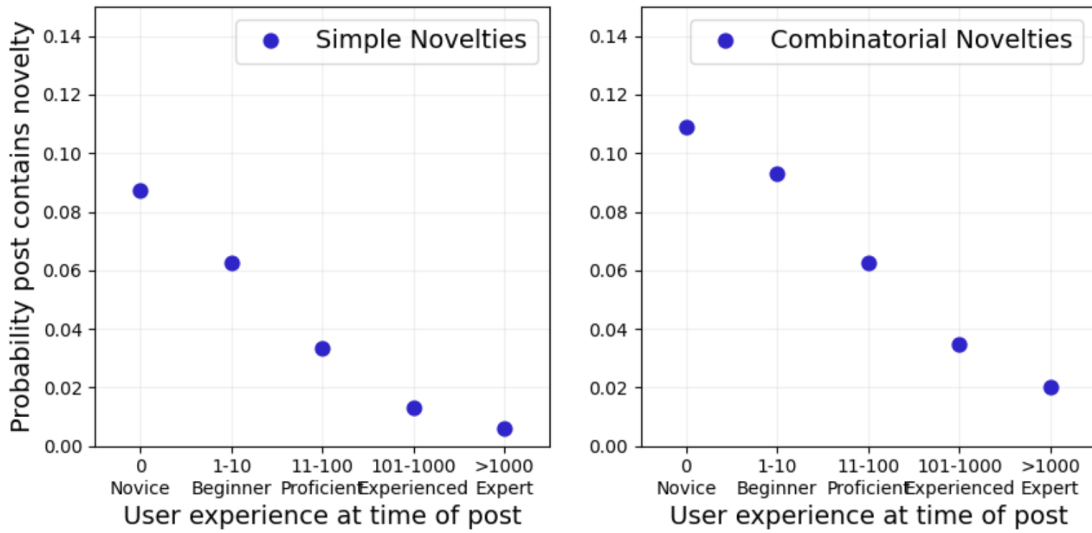


Fig. 11: User Experience Distribution for Simple and Combinatorial Novelties among libraries appearing at least 1000 times.

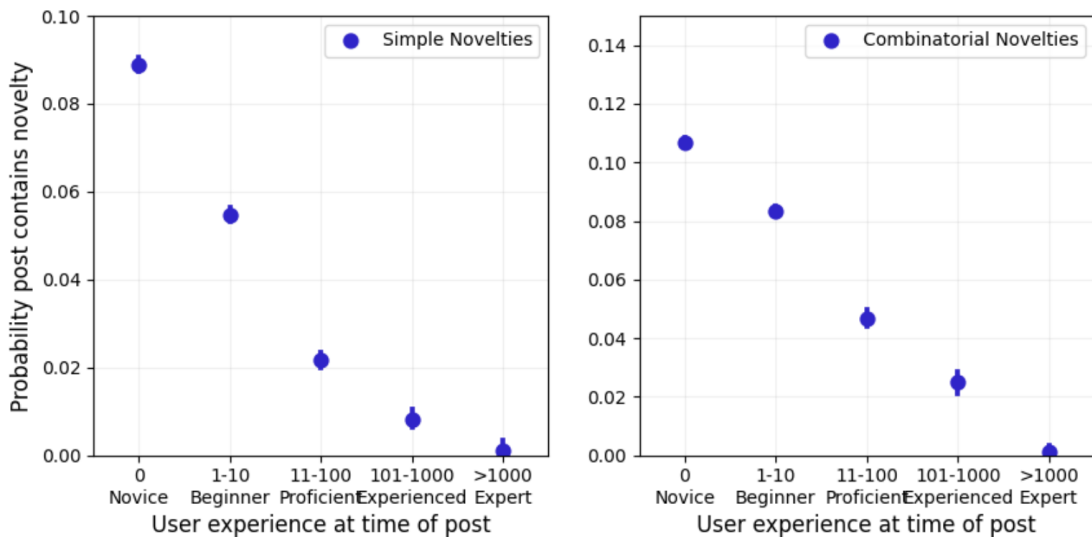


Fig. 12: User Experience Distribution for Simple and Combinatorial Novelties among libraries appearing at least 100 times - restricted view to 2016.

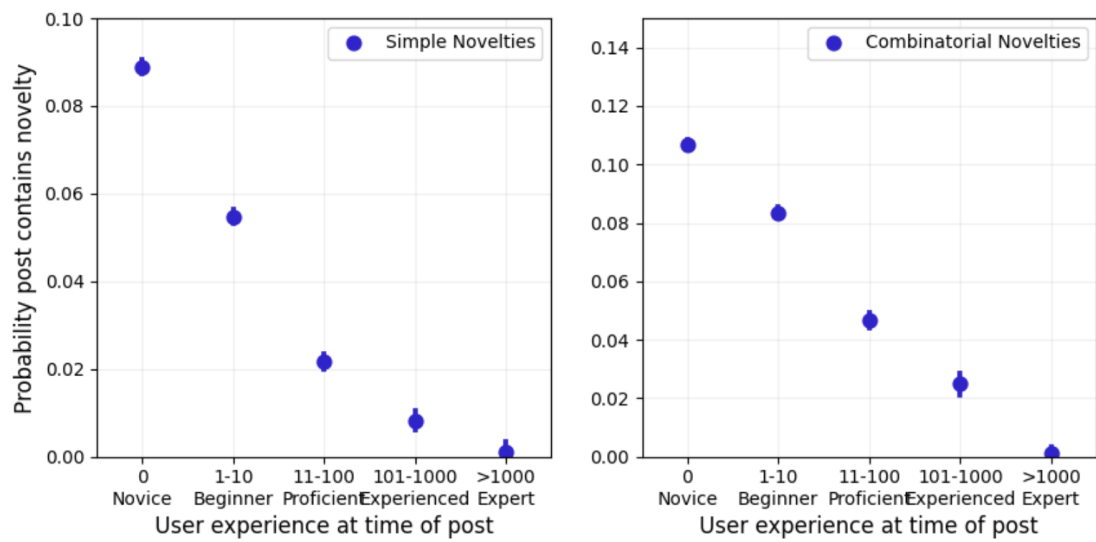


Fig. 13: User Experience Distribution for Simple and Combinatorial Novelties among libraries appearing at least 1000 times - restricted view to 2016.