How to Debug Inclusivity Bugs? A Debugging Process with Information Architecture

Mariam Guizani Oregon State University Corvallis, Oregon, USA guizanim@oregonstate.edu

Abrar Fallatah Oregon State University Corvallis, Oregon, USA fallataa@oregonstate.edu Igor Steinmacher Northern Arizona University Flagstaff, AZ, USA igor.steinmacher@nau.edu

Margaret Burnett Oregon State University Corvallis, Oregon, USA burnett@engr.orst.edu Jillian Emard Oregon State University Corvallis, Oregon, USA emardj@oregonstate.edu

Anita Sarma Oregon State University Corvallis, Oregon, USA anita.sarma@oregonstate.edu

ABSTRACT

Although some previous research has found ways to find inclusivity bugs (biases in software that introduce inequities), little attention has been paid to how to go about fixing such bugs. Without a process to move from finding to fixing, acting upon such findings is an ad-hoc activity, at the mercy of the skills of each individual developer. To address this gap, we created Why/Where/Fix, a systematic inclusivity debugging process whose inclusivity fault localization harnesses Information Architecture(IA)—the way user-facing information is organized, structured and labeled. We then conducted a multi-stage qualitative empirical evaluation of the effectiveness of Why/Where/Fix, using an Open Source Software (OSS) project's infrastructure as our setting. In our study, the OSS project team used the Why/Where/Fix process to find inclusivity bugs, localize the IA faults behind them, and then fix the IA to remove the inclusivity bugs they had found. Our results showed that using Why/Where/Fix reduced the number of inclusivity bugs that OSS newcomer participants experienced by 90%.

Lay Abstract: Diverse teams have been shown to be more productive as well as more innovative. One form of diversity, cognitive diversity — differences in cognitive styles — helps generate diversity of thoughts. However, cognitive diversity is often not supported in software tools. This means that these tools are not inclusive of individuals with different cognitive styles (e.g., those who like to learn through process vs. those who learn by tinkering), which burdens these individuals with a cognitive "tax" each time they use the tool. In this work, we present an approach that enables software developers to: (1) evaluate their tools, especially those that are information-heavy, to find "inclusivity bugs"—cases where diverse cognitive styles are unsupported, (2) find where in the tool these bugs lurk, and (3) fix these bugs. Our evaluation in an open source

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICSE-SEIS'22, May 21–29, 2022, Pittsburgh, PA, USA © 2022 Association for Computing Machinery. ACM ISBN 978-1-4503-9227-3/22/05...\$15.00 https://doi.org/10.1145/3510458.3513009

project shows that by following this approach developers were able to reduce inclusivity bugs in their projects by 90%.

KEYWORDS

Diversity, Information Architecture, Open Source, Inclusivity Bugs

ACM Reference Format:

Mariam Guizani, Igor Steinmacher, Jillian Emard, Abrar Fallatah, Margaret Burnett, and Anita Sarma. 2022. How to Debug Inclusivity Bugs? A Debugging Process with Information Architecture. In *Software Engineering in Society (ICSE-SEIS'22), May 21–29, 2022, Pittsburgh, PA, USA*. ACM, New York, NY, USA, 12 pages. https://doi.org/10.1145/3510458.3513009

1 INTRODUCTION

Although in recent times diversity initiatives have become common, sometimes we forget *why* diversity is important to so many organizations. Besides social justice reasons, what many organizations hope to gain from diverse backgrounds (cultural, ethnic, education, gender, etc.) is diversity of information and of thought [46]—i.e., *cognitive diversity*. Diversity's accompanying diversity of thought has been shown to have many positive effects on organizations, including better ability to innovate, better reputation as ethical corporate citizens, and a better "bottom line" for businesses [31, 44, 46]. However, efforts to support diversity rarely consider either cognitive diversity or inclusivity of technology environments.

In this paper, we consider these aspects together: how to *increase* support for *cognitive diversity* within *technology environments*, especially information-heavy ones. The setting for our investigation is an information-heavy environment that is particularly challenged in attracting diverse populations: Open Source Software (OSS) communities.

This study complements the existing literature: previous work has investigated OSS-specific challenges [20, 28, 32, 58] and the inclusivity issues affecting OSS [5, 12, 27, 35, 39, 42, 49, 51, 63], but has not focused on how to *debug* OSS projects' *technology* to support cognitive diversity.

1.1 Why/Where/Fix: An IA-based Inclusivity-Debugging Process

A debugging perspective suggests that OSS practitioners who want to improve inclusivity of their project's infrastructure will need three capabilities. (1) First, they need to find "inclusivity failures" (analogous to testing [1]). Since the failure is about inclusivity (not

about producing a wrong output), OSS practitioners will also need to be able to discern *why* the observed phenomenon is considered an inclusivity failure. (2) Second, the practitioners will need to tie an inclusivity failure to *where* the "inclusivity fault(s)" occur (analogous to fault localization [3]); so that (3) the inclusivity faults can be *fixed* to stop the associated inclusivity failure from occurring. In this paper, we term the inclusivity debugging capabilities as "Why/Where/Fix" (Fig. 1), and investigate its efficacy at debugging inclusivity bugs.

Debugging requires a definition of a bug. We derive our definition from the testing community's notion of a software failure. Ammann and Offutt define a "failure" as "...external, incorrect behavior with respect to the requirements or ... expected behavior" [1]. Our analogous requirement/expected behavior is inclusivity across diverse cognitive styles, so we define *inclusivity failures/bugs* as user-visible features or workflows that do not equitably support users with diverse cognitive styles. As with Ammann/Offutt's definition, an inclusivity bug is a barrier but not necessarily a "show-stopper". That is, if groups of users eventually complete their tasks but disproportionately experience barriers along the way (e.g., confusion, missteps, workarounds), these too are inclusivity bugs.

To find such inclusivity bugs and their "Why"s, we used Gender-Mag [11], an empirically validated method [9, 25, 43, 65] with a dual gender/cognitive focus. GenderMag integrates finding an inclusivity bug with its "Why", because using GenderMag includes identifying cognitive mismatches that pinpoint which users disproportionately run into barriers using a system. In this paper, an OSS team used GenderMag to find inclusivity bugs in their OSS project.

After finding a bug, the next step in debugging is to figure out what and where a bug's causes are, referred to as "faults" in SE literature. According to Avizienis et al. [3] a fault is the underlying cause of an error, a condition that may lead to a failure; and fault localization is the act of identifying the locations of faults. Building upon these definitions, we define an *inclusivity fault* as the userfacing components (e.g., UI elements, user-facing documentation, workflow) of the system that produced an inclusivity bug; and *inclusivity fault localization* as the process of identifying the locations of these faults in these user facing components.

Thus, for Why/Where/Fix's "Where", we devised a systematic inclusivity fault localization approach that harnesses Information Architecture (IA) [38]. IA is the "blueprint" for the structure, arrangement, labeling, and search affordances of information content, and is especially pertinent to information-rich environments [55]. Although substantial research exists on how Information Architectures can support usability, navigation, and understandability [17, 22, 29, 36, 52],



Figure 1: The Why-Where-Fix process. The *Why* is to produce the inclusivity bugs and the cognitive styles behind them; the *Where* is to localize the faults behind the bugs to the user-facing IA elements; and the *Fix* is to change the IA elements to expand the software's inclusivity. The grey roundtangle highlights Why/Where/Fix's new contributions.

research has not considered how different Information Architectures do or do not support populations with diverse cognitive styles, or how IA can be used for inclusivity fault localization.

To use IA to tie together the above "Why" and "Where" foundations to point to the fixes, we supplemented the GenderMag process for finding inclusivity bugs with a mechanism by which evaluators specified any IA elements (the faults) implicated in the inclusivity bugs found along the way. Thus, the Why/Where/Fix process in Figure 1, is: find the bugs using cognitive styles, which contribute the Why (using GenderMag), enumerate the implicated IA elements involved in the bug (Where), and change those IA elements (Fix).

1.2 Can IA Squash the Inclusivity Bugs?

We have pointed out that the Why capability (finding inclusivity bugs) is already possible using GenderMag. But debugging requires getting from finding to fixing, and this capability of Why/Where/Fix rests on IA.

Thus, to empirically investigate IA's effectiveness in the Why/ Where/Fix debugging process, we used a three-stage combination of field work (Stage One and Stage Two) and lab work (Stage Three), as follows:

In Stage One (Why \rightarrow Where), we worked in the field with an OSS team who used GenderMag to detect cognitive inclusivity bugs in their project's infrastructure, to investigate RQ1: Is IA implicated in inclusivity bugs? If so, how? In Stage Two (Where \rightarrow Fix), the OSS team worked alone to change the project infrastructure's IA using what they had learned in Stage One, which enabled us to investigate RQ2: Can practitioners use IA to fix inclusivity bugs? If so, how? In Stage Three (Lab), we brought OSS newcomers into the lab to investigate whether the team's IA-localized faults and fixes decreased the inclusivity bugs those newcomers experienced.

Our primary contributions in this paper are:

- (1) Presents and empirically investigates the first inclusivity debugging process, including systematic fault localization.
- (2) Empirically investigates whether Information Architecture can itself be the cause of inclusivity bugs.
- (3) Reveals ways OSS projects can improve their infrastructures' Information Architecture to improve their project's inclusivity.

2 BACKGROUND AND RELATED WORK

2.1 Information Architecture

The term "Information Architecture" was coined in the mid-70's as a way of "making the complex clear" [66]. This paper follows the definition of Morville and Rosenfeld [38], referred to as the "bible" of IA, that defines IA as a set of four component systems (Figure 2).

The first is the *Organization System (Org)*, analogous to the architectural arrangement of a building's "rooms", which has an organization scheme *OrgScheme* and an organization structure *OrgStruct*. The organization scheme is the way content is arranged or grouped (e.g., alphabetical or by task). An architect chooses the scheme according to the situations they want the Information Architecture to support, such as alphabetical (*OrgScheme-Alpha*) to support exact look-ups, or task-based (*OrgScheme-Task*) to facilitate high priority tasks. The organization structure defines the relationship between content groups (e.g., hierarchical (*OrgStruct-Hierarchy*)).

Second, the *Navigation System (Nav)*, analogous to adding doors and windows to a building, enables users to traverse the information groupings and structure. Some of the navigation system is embedded in the content (e.g., contextual links (*Nav-ContextualLink*)), while others are supplemental (e.g., site maps). Third, the *Labeling System* (Label) adds signposts (also known as "cues" in Information Foraging literature [47]) to the "doors", such as the labels on contextual links (*Label-ContextualLink*), headers (*Label-Header*), cues/keywords (*Label-IndexTerm*), etc. Fourth, the *Search System*, when provided, supplements the rest of the IA, to enable users to retrieve information using a particular term or phrase.

The majority of IA research has focused on the design and evaluation of websites, but some has explored other domains. For example, IA has been used in the design of usable security tools [14], as the basis of a semantic web structuring tool [7, 8, 15], to investigate the accessibility, use and reuse of information across multiple devices [40], to evaluate different information visualization tools [33], and for mobile applications screen-reader navigation [19, 67]. One body of research has compared IA to other attributes of information sites. For example, Aranyi et al.'s empirical evaluation of a news website showed that the content and its IA were the main problems [2]. Petri and Power's study likewise found prominent IA problems when evaluating six government websites, with IA accounting for about 9% of user-reported bugs [45].

Other IA research has evaluated the usability of different subsets (organizational vs. labeling schemes) of IA. For example, Gullikson et al. evaluated the IA of an academic website and reported that although participants were satisfied with the content of the site, they found its (IA) labeling to be confusing [23], and were especially dissatisfied with the IA's organization system. Resnick and Sanchez found that user-centric labels significantly improved user performance and satisfaction as compared to user-centric organization, which only improved performance if labels were of low quality [50]. Similarly, others have found that navigation success depends more on the quality of labels than the structure of a page [37, 56].

Of particular interest is IA research on supporting diverse populations. Lachner et al. used IA to promote cultural diversity and used Hofstede et al. power distance cultural dimension [26] to design and evaluate culturally-specific collaborative Q&A websites [30]. Accessibility and IA has been studied by others. Swierenga et al. showed that IA's organization and labeling system create barriers for visually impaired and low-vision individuals [62]. A multitude of

*Organization	*Navigation	*Labeling	Search	
"Building rooms"	"Adding doors and windows"	"Adding door signs"	"Asking for directions"	
*Scheme: Rationales behind content grouping (e.g., *topic, *task)	*Embedded: Inherent to the structure of a system (e.g.,*contextual link)	Labels of (e.g., *contextual link, *header, *index term)	Supplements the IA systems by allowing information retrieval using particular terms	
*Structure: Relationship between content groups (e.g., *hierarchy)	Supplemental: Auxiliary to the structure of a system (e.g., site map)		or phrases (e.g., search engine)	

Figure 2: IA's four component systems [38]. The organization and navigation systems have subsystems (underlined). *s mark IA (sub)systems and elements used in this paper.

research [4, 16, 53, 54, 64, 67] has investigated IA auditory systems for designing and evaluating accessible websites for low-vision users. Ghahari et al., for example, showed how topic- and list-based aural navigation strategies can enhance user's navigation effectiveness and efficiency [53]. However, we cannot locate any research on how IA can support cognitive diversity.

2.2 Diversity and the GenderMag Method

GenderMag, a method used to find and fix inclusivity bugs, provides a dual lens—gender- and cognitive-diversity—to evaluate workflows. It considers five dimensions ("facets" in GenderMag) of cognitive styles (Table 1), each backed by extensive foundational research [11, 61]. Each facet has a range of possible values. A few values within each facet's range are brought to life by the three GenderMag personas: "Abi", "Pat", and "Tim." Abi's facets are statistically more common among women than other people, Tim's are statistically more common among men, and Pat has a mix of Abi's/Tim's facets plus a few unique ones.

Each persona is a "multi-"persona [25]—their demographics can be customized to match those of the system's target audience. For example, any gender, any photo, any educational background, or any pronoun can be integrated (e.g., she/her, he/him, they, ze, etc.). Their cognitive facets, however, remain fixed. Figure 3 shows portions of the OSS team's customization of Abi, which they used in Stage One.)

Evaluation teams, such as the OSS team in this paper, use Gender-Mag to walk through a use-case in the project they are evaluating using Abi, Pat, or Tim. At each step of the walkthrough, the team writes down the answers to three questions: (1) whether <Persona> would have the subgoal the project owners hoped for and why, (2) whether <Persona> would take the action the project owners hoped for and why, and (3) if <Persona> did take the hoped-for action, would they know they did the right thing and were making progress toward their goal, and why. When the answer to any of these questions is negative, it identifies a potential bug; if the "why" relates to a particular cognitive style, this shows a disportionate effect on people who have that cognitive style—i.e., an *inclusivity bug*. Thus, a team's answers to these questions become their inclusivity bug report, which they can then process and prioritize in the same way they would do with any other type of bug report.

Abi (Abigail/Abishek)



Abi is a second-year engineering student... She is comfortable with the technologies she uses regularly... She is interested in branching out to the world of open source..., but their software systems are new to her... She likes Math...

Abi's facets are listed and described here

Figure 3: Portions of the OSS team's Abi persona. The photo(s) and blue text are customizable; the black text is not. Abi's facets (gray block) are as per Table 1. (The supplemental document [21] includes the full Abi persona used in Stage One.)

The method and its derivatives have been used in a variety of domains, such educational software, digital libraries, search engines, and software tools [11, 13, 18, 24, 35, 57, 65]. Particularly pertinent to this paper, in a study of OSS professionals, over 80% of the barriers they found in OSS projects were gender inclusivity bugs, which were later confirmed by OSS newcomers [43].

However, prior work has left largely to the practitioners' judgment how exactly to fix such inclusivity bugs (e.g., [65]). This paper aims to pave a path from finding to fixing with an IA-based process to systematically localize inclusivity faults.

3 METHODOLOGY

We conducted a multi-stage (in-the-field and in-the-lab) empirical investigation to analyze whether changing the IA of an OSS project infrastructure would help support newcomers across a range of diverse cognitive styles. ¹

Table 1: The GenderMag cognitive facet values for each persona. The research behind each facet is enumerated in [11].

Facet	Cognitive facet value for each persona
Motivations	Uses technology <i>Abi</i> : Only as needed for the task at hand. Prefers familiar and comfortable features to keep focused on the primary task. <i>Tim</i> : To learn what the newest features can help accomplish. <i>Pat</i> : Like Abi in some situations and like Tim in others.
Self- Efficacy	Abi: Lower self-efficacy than their peers about unfamiliar computing tasks. If tech problems arise, often blames self, and might give up as a result. Tim: Higher self-efficacy than their peers with technology. If tech problems arise, usually blames the technology. Sometimes tries numerous approaches before giving up. Pat: Medium self-efficacy with technology. If tech problems arise, keeps on trying for quite awhile.
Attitude Toward Risk	Abi and Pat: Risk-averse, little spare time; like familiar features because these are predictable about the benefits and costs of using them. Tim: Risk tolerant; ok with exploring new features, and sometimes enjoys it.
Information Processing	Abi and Pat: Gather and read everything comprehensively before acting on the information. Tim: Pursues the first relevant option, backtracking if needed.
Learning Style	Abi: Learns best through process-oriented learning; (e.g., processes/algorithms, not just individual features). Tim: Learns by tinkering (i.e., trying out new features), but sometimes tinkers addictively and gets distracted. Pat: Learns by trying out new features, but does so mindfully, reflecting on each step.

¹We did not *recruit* participants with any particular cognitive style as a criterion; rather, we *collected* cognitive style data as part of the investigation.

For the field aspect, we gathered in-the-field data from an OSS project team (Team F) that was interested in increasing diverse newcomers' participation in their project (Project F). The empirical investigation had three stages:

- Stage One (Why → Where), in the field: We worked with Team F
 to detect IA-based inclusivity bugs. Team F then worked alone
 to select which of these bugs to fix.
- Stage Two (Where → Fix), in the field: Team F worked alone to derive IA-based cognitive diversity-inspired fixes to Project F's Information Architecture.
- Stage Three, in the lab: We brought OSS newcomers into the lab to evaluate the inclusivity bugs they encountered with the original Project F vs. the new version of Project F.

3.1 Stage 1, Team F, RQ1 (in the field): Why \rightarrow Where

Stage One had two purposes. First, for ecological validity, we wanted to avoid artificially creating inclusivity bugs; thus, Stage One provided a way to harvest them from a real OSS project. For this purpose, we used the GenderMag method (Section 2.2). To facilitate IA-based fault localization, we then added the following IA-based Where question to the GenderMag question set: "What in the UI helped/confused <Persona> in this step?" Both the original and IA-supplemented GenderMag forms, and all our study materials, are provided in the supplemental document [21].

Note that Stage One's purpose was *not* to investigate whether an OSS team can use GenderMag to point out inclusivity bugs, because its validity with OSS project teams has already been validated [43]. The GenderMag method has also been empirically validated in other lab [11, 25, 65] and field [9] studies. As with other cognitive walk-through (CW) methods, its reliability (precision) is very high: CW methods tend to have false-positive rates of 5%-10%, and Gender-Mag's false-positive rates have been 5% or lower [11, 43, 65].

For Stage One's bug harnessing purpose, Team F worked with two researchers using the IA-supplemented GenderMag method to find inclusivity bugs in four use-cases (Table 2). Team F selected these use-cases for their importance for Project F newcomers. Analyzing these use-cases produced both a list of likely inclusivity bugs with the facets that caused them (Why), and IA-localized faults that may have produced these bugs (Where).

The second purpose of Stage One was the beginning of our RQ1 investigation into whether some IA elements are indeed implicated in such real-world inclusivity bugs. For this purpose, Team F worked alone, without our help.

Team F began by deciding which of the bugs to take forward into the next stage of the investigation. They selected these bugs using the criteria that the bug (1) had at least one cognitive facet that the Information Architecture did not support; and (2) was associated with the project itself and *not* the UI of the hosting platform (e.g., GitLab, GitHub). These criteria produced 6 bugs (Table 2).

Along the way, Team F had noticed some general usability bugs not related to any cognitive facet. To prevent these from influencing Stage Three, Team F fixed these bugs and brought the project up to GitHub's recommended content standards [41], resulting in the prototype we call the *Original* version.

3.2 Stage 2, Team F, RQ2 (in the field): Where \rightarrow Fix

Team F then worked alone to derive fixes for each of these 6 bugs by changing the IA elements they had identified as the probable causes of the bugs, so as to better support the previously unsupported cognitive facets without loss of support for the supported facets. We refer to the "fixed" version of Project F as the *DiversityEnhanced* version.

3.3 Stage 3, OSS Newcomers, RQ1+RQ2 (in the lab)

We then brought OSS newcomers into the lab to investigate: (1) whether OSS newcomers trying to use the Original version would run into the bugs Team F had found in the Original version, and (2) whether the IA fixes Team F had derived for the DiversityEnhanced version would actually improve support for cognitively diverse OSS newcomers.

We recruited the OSS newcomers from a large US university. Our recruiting criteria were people with no prior experience contributing to OSS projects. All 31 respondents came from a variety of science and engineering majors. Because the investigation focuses only on cognitive diversity (not on disabilities), we did not seek out participants with any particular cognitive style or with a disability. Because none of the experimental tasks required programming, we did not collect their programming experience.

Participants filled out a cognitive facet questionnaire [9, 18, 65] (provided in our supplemental document [21]) in which participants answered Likert-scale items about their cognitive styles. Using their responses and genders, we selected 18 respondents to gender-balance and to include a wide range of cognitive styles (Figure 4). Of the 18 selected participants, 8 identified as women, 9 identified as men, and one participant declined to specify their gender.

We assigned participants to the Original or DiversityEnhanced treatments, balancing the cognitive styles between the treatments based on the participants' cognitive facet questionnaire responses.

Table 2: The four use-cases and associated bugs. Team F provided these use-cases, which were important to their project.

Use-Case	Descriptions	Bugs
U1-Find	Finding an issue to work on	Bug 1 & 2
U2-Document	Contribute to the documentation	Bug 3
U3-FileIssue	File an issue	Bug 4
U4-Setup	Set up the environment	Bug 5 & 6

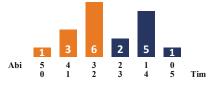


Figure 4: Number of participants with more Abi facets (left half, orange) or more Tim (right half, blue). For example: the first column says that 1 participant had 5 Abi facets and no Tim facets. Table 1 explains Abi, Tim, and their facets.

Because facet values are relative to one's peer group, the median response for each facet served to divide closer-to-Abi facet values from closer-to-Tim facet values. This produced identical facet distributions (Figure 5) for both groups. We audio-recorded each participant as they talked-aloud while working on the use-cases presented earlier in Table 2. We transcribed the recordings, and counted how often the participants encountered one of the 6 bugs that Team F had attempted to fix.

We qualitatively coded cognitive facets that participants verbalized when they encountered one of these bugs, which enabled us to compare participants' *in-situ* reactions to their cognitive facet questionnaire responses. For example, we coded P2-O's verbalization "...this leads me to a page with the bare minimum of instructions... I have no idea where to go from here" as "learning style: processoriented", which aligned with their questionnaire response. To ensure reliability of the coding, two researchers independently coded 20% of the data and calculated IRR using the Jaccard index. Jaccard, a measure of "consensus" interrater reliability [60], is useful when multiple codes per segment are used, as in our case. The consensus level was 90.2%. Given this level of consensus, the researchers split up coding the remainder of the data.

At session end, participants filled out a subset of the System Usability Scale (SUS) survey [6] (supplemental document [21]).

4 RESULTS

We begin with "whether" answers to both research questions—for RQ1, whether Information Architecture was implicated in the inclusivity bugs, and for RQ2, whether Team F's IA fixes increased inclusivity for OSS newcomers.

As Table 3 shows, both answers were yes. Regarding RQ1, with the Original version, OSS newcomers ran into inclusivity bugs in the Information Architecture 20 times. Regarding RQ2, Team F's inclusivity fixes to the IA reduced the number of inclusivity bug experiences in the DiversityEnhanced version to only 2. In total, Team F's IA fixes cut the number of bugs participants experienced by 90% (Table 3).

To answer the *how* aspects of our RQs, Table 4 summarizes, for each bug, Team F's *Why* analyses (first column) of the cognitive facets involved in the bug, their *Where* analyses to localize the faults to IA elements (second column), and how they implemented their IA *Fixes* (third column). The following sections discuss them in depth.

4.1 Bug 1 & 2 in Depth: Issues with the "issue list"

The first two rows in Table 4 show how Team F addressed Bug 1 & 2, the IA-based inclusivity bugs that Team F identified in Stage One





Figure 5: Number of participants with Abi (bottom, orange) vs. Tim (top, blue) facets who used the Original (columns 1-5) vs. DiversityEnhanced (columns 6-10) versions of the OSS project. (The two distributions are identical.)

in the context of use-case U1-Find (finding a task to work on). As Table 4 shows, for Bug 1, Team F predicted that Abi-like newcomers would face problems in understanding the *process* of finding an issue. Their *why* analysis (Table 4 row 1 col. 1) pointed out that the lack of information about finding an issue could be problematic to comprehensive information processors, risk averse, or process-oriented newcomers. As Stage Three Participant 1 using the Original version later put it:

P1-O: "I just feel like I wouldn't have enough to go on."

Team F localized the fault (*wheres*, Table 4's row 1 col. 2) to the IA's link labeling (*Label-ContextualLink*) and to the absence of keywords (*Label-IndexTerm*), which could lead newcomers to follow wrong link(s) and never obtain the kind of information they were seeking.

Once a newcomer was past Bug 1, Team F predicted that the Issue List provided too little information to enable some newcomers to select an issue appropriate to their skills (Bug 2). Team F's *why* analysis showed that this bug would be particularly pertinent to newcomers with a comprehensive information processing style, low self-efficacy, or risk aversion.

Team F localized the fault behind Bug 2 (IA *wheres*) to the issue list's nondescript titles, uninformative descriptions, and limited labeling. Team F realized that, with this IA, the Issue List gave little indication as to whether an issue would fit a newcomer's skill level (*Label-IndexTerm, Label-Header*). Stage Three proved Team F to be right: Bug 1 & 2 did affect several participants (Figure 6):

P1-O: "...I don't really know...I would say if I had to fix [an issue from the issue list], I'd probably just ask someone for help."

To fix Bug 1 (Table 4 rol 1 col. 3), Team F made several changes to the IA. They created better cues for the link to the contribution guidelines by changing its label (Label-ContextualLink) from the file name ("contributing.md") to "contributing guidelines" and including additional keywords about what to expect from the link. They also modified the IA of the "contributing.md" to point out specific task-oriented instructions for finding an issue (OrgScheme-Task) including a header (Label-Header)—"Find an issue" (Fig 7), a link to the "issue list"(Nav-ContextualLink, Label-ContextualLink), and additional keywords (Label-IndexTerm) to add support for processoriented and risk-averse newcomers.

Team F fixed Bug 2 (Table 4 row 2 col. 3) with improved issue headers and labels (*Label-Header, Label-IndexTerm*). The labels signaled attributes of the open issues in the project (Figure 8). Team F also rewrote some issue descriptions to support newcomers with a comprehensive information processing style.

In Stage Three, the OSS newcomer participants showed that Bug 1 & 2 were pervasive; *all* participants using the Original version faced problems related to Bug 1 and/or 2 (Figure 6). But were these

Table 3: The number of participants who ran into the bug(s), out of the 18 participants (9/ group).

Bug ID	Original	DiversityEnhanced
Bug 1 & 2	9/9	1/9
Bug 3	2/9	0/9
Bug 4	0/9	0/9
Bug 5 & 6	9/9	1/9
Total bugs encountered	20	2

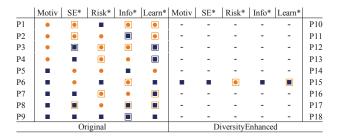


Figure 6: In Bug 1 & 2, all Original participants ran into bugs (left), but only 1 DiversityEnhanced participant (right). Participant ID numbering is from the most Abi-like to the most Tim-like.

*: facet the fix(es) targeted; circles | squares: the facet values from the participants facet questionnaire for Abi-like and Tim-like facet values respectively; square outline | square outline: Abi-like | Tim-like facet values participants expressed when they ran into a bug.



Figure 7: Bug 1 before the fix, the screen appeared as shown without the call-out, giving little guidance on how to find a suitable issue. The fix added the "Find an issue" process description.

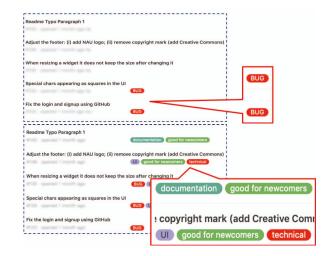


Figure 8: Top: Bug 2 before the fix had only one label ("Bug"). Bottom: The fix added multiple descriptive labels.

bugs *inclusivity* bugs, i.e., *disproportionately* affecting people with particular cognitive styles?

Table 4: For each use-case's bug(s), excerpts from Team F's Stage One analysis, the Bug's Why's (facets impacted), Where's (IA involved), and their Stage Two IA fixes.

Bug's Why: Facets	Bug's Where: IA involved	Bug's Fixes and IA elements changed	
	" may click [the wrong link]" IA: Label-ContextualLink, Label-IndexTerm	 In README.md: Label-IndexTerm: added cue/keyword to guide to "contribut finding an issue. Label-ContextualLink: changed a link label to clarify what it In contributing.md: Nav-ContextualLink, Label-ContextualLink: added a link to Label-IndexTerm: added cues/keywords to guide issue choic orgscheme-Task, Label-Header: added a header following a organization scheme. Other: added more information. 	leads to. the "issue list".
	"labels will help, but there aren't labels for every issuelike 'good for newcomer'. Headings are missing info should be a bit more detailed" IA: Label-IndexTerm, Label-Header		ve. See Figure 8
"[The instructions are] all about technical contributions, nothing about documentation changes [So] she may think that she needs to do all the technical setup before editing the README (which is a lot)" Facets: Motiv, Learn, SE, Risk	"README and contribute files may confuse her. The README is here but there is no clear indication [cue/keyword] of what she needs to do to change the file." IA: Label-IndexTerms	 In README.md: Label-IndexTerms: added cue/keyword to guide to "contrib documentation contributions. In contributing.md: 	ntation contribution.
an issue is part of contributions. No clear instruction about what	"doesn't say where to find the issue listMaybe adding an indication [cue/keyword] or a link would be helpful." IA: Label-IndexTerm, Nav-ContextualLink, Label-ContextualLink		sue.
" nothing that explicitly says set up the envShe would read through step 0 and think it's not for mac [OS]." Facets: Info, SE, Risk	"no hint [cue/keyword] about how to set up the environment in the readme More about Ubuntu and Linux and not about Windows and Macmaybe this file needs to be more high level." IA: Label-IndexTerm, OrgScheme, OrgStruct	 In README.md: Label-IndexTerm: added cue/keyword to "contributing guid up the environment. In contributing.md section "Help us with code":	er of abstraction.
" No explanation about the different things to install and where to install them". Facets: Info, Motiv, Learn, SE, Risk	"sees all this code and does not know where and how to run it. Maybe a hint about using the terminal [cue/ keyword] and copying and pasting the code would be helpful." IA: Label-IndexTerm	• In OS instruction sub-pages: - Label-IndexTerm: added cues/ keywords about where to exc - Other: added additional explanation about each command.	ecute commands. See Supp.Doc.

Figure 6 answers this question. Counting up the colored outlines, which show the facets Stage Three participants verbalized *when they ran into those bugs*, shows that Bug 1 & 2 disproportionately impacted Abi-like facet values: 74% (14/19) of the facets participants verbalized with Bug 1 & 2 were Abi-like facet values (orange square outlines in Figure 6, left).

Although Bug 1 & 2 disproportionately affected participants with Abi-like facet values, targeting these facets helped participants across the entire cognitive style spectrum, *both for Abi-like and Tim-like newcomers* (Figure 6, right). Further, only one participant of the

DiversityEnhanced treatment (P15-D, Figure 6) ran into these bugs—compared to all 9 participants in the Original treatment (Table 3).

Even when participants veered off track, the label fixes (*Label-IndexTerm*) (Figure 8) helped them find their way back. For example, P17-D initially chose an issue labeled "good for newcomers" and "technical", but soon found that they would have needed more coding experience. P17-D realized that issues that did not include the "technical" label would be a better fit.

P17-D: "...and in fear of not making the same mistake, I'm just going to go with a [issue], which only says good for newcomers and documentation."

4.2 Bug 3: "I would expect something linear"

When evaluating the documentation contribution use-case (U2-Document), Team F predicted that newcomers might think that they have to go through all the technical setup in order to make any contribution, even a documentation contribution (Bug 3). Team F's *why* analysis (Table 4's third row) pointed to four of Abi's cognitive styles: task-oriented motivations, process-oriented learning, relatively low self efficacy, and risk aversion. Team F localized Bug 3's fault in the IA (*wheres*) to point to the absence of keywords that could guide newcomers in contributing documentation.

In Stage Three, Team F's prediction was borne out: two lab participants did run into Bug 3 (Figure 9). For example:

P2-O (risk-averse as per facet questionnaire responses): "Should I be doing this? Like, should I be coding just to change an N to an M? Seems a little unnecessary?...I'm stuck."

The lack of a task-centric organization scheme for the instructions also impacted P2-O, a process-oriented learner according to their facet questionnaire responses:

P2-O: "I would expect something linear."

As Table 4 row 3 col. 3 summarizes, Team F fixed the IA by mentioning "contributions with documentation" in the README.md (Label-IndexTerm), and by organizing contributing.md information with a header (Label-Header) that followed a task-based organization scheme (OrgScheme-Task), to support people with Abilike motivations. Team F also added step-by-step instructions, keywords (Label-IndexTerm) and links to detailed information (Nav-ContextualLink), to support diverse learning and information processing styles.

The results of Stage 3 showed that the changes had positive effects. As Figure 9 shows, although two participants ran into Bug 3 with the Original version, nobody did using the DiversityEnhanced version.

4.3 Bug 4: Where to go to file an issue

For Bug 4 Team F decided that, in trying to file an issue (use-case U3-FileIssue), newcomers might not know where to go, especially those who are risk-averse, those with comprehensive information processing styles or relatively low self-efficacy (Table 4 row 4 col. 1). The elements of IA *where* the team found these problems were in *Nav-ContextualLink*, *Label-IndexTerm*, and *Label-ContextualLink* elements.

	Motiv*	SE*	Risk*	Info	Learn*	Motiv*	SE*	Risk*	Info	Learn*	
P1	-	-	-	-	-	-	-	-	-	-	P10
P2	•					-	-	-	-	-	P11
P3	-	-	-	-	-	-	-	-	-	-	P12
P4	-	-	-	-	-	-	-	-	-	-	P13
P5			•			-	-	-	-	-	P14
P6	-	-	-	-	-	-	-	-	-	-	P15
P7	-	-	-	-	-	-	-	-	-	-	P16
P8	-	-	-	-	-	-	-	-	-	-	P17
P9	-	-	-	-	-	-	-	-	-	-	P18
		Oı	iginal				D	iversityE	nhance	d	

Figure 9: Two Original treatment participants ran into Bug 3, but nobody using the DiversityEnhanced version did. *, circles, squares: see Figure 6.

However, Team F was wrong—in Stage Three, none of the Original version lab participants ran into Bug 4. The reason was a flaw in Team F's analysis of this use-case as it related to newcomers' prior experience. In the Stage Three task sequence, participants had already been to the "issue list" in context of an earlier use-case (U1-Find). Thus, as P5-O put it:

P5-O: "Since I already spent some time on that issue page [issue list]. That part [filing an issue] was not too hard."

Still, Stage Three had not yet occurred, and Team F made the IA fixes in Stage Two to fix the bug (Table 4 row 4 col. 3). The Stage Three participants who then used the DiversityEnhanced version experienced no problems. Thus, the question of whether newcomers *would have* run into these problems if they had not previously learned the features remains unanswered. However, the question of whether newcomers ran into problems in the changed version is answered: nobody ran into any problems in the DiversityEnhanced version (Table 3).

4.4 Bug 5 & 6: What, where, and how to set up

In use-case U4-Setup, Team F's analysis revealed Bug 5 (Table 4's fifth row), namely that newcomers with comprehensive information processing style, low self-efficacy, or risk aversion could run into problems finding the setup instructions for their particular operating system (OS). Team F identified the underlying faults to be the *Label-IndexTerm*, *OrgScheme* and *OrgStruct*, none of which were pointing out where different OSs' setup instructions might be.

Even if a newcomer overcame Bug 5 and found the right instructions, Team F realized that an OSS newcomer might not necessarily "just know" what each command in the instructions actually did or exactly where to run them (Bug 6: Table 4's sixth row). As the table shows, Team F's why analysis suggested that this inclusivity bug could particularly affect a newcomer with any of Abi's cognitive style values, due in part to the absence of hints with clarifying keywords (e.g., "command line terminal...") (Label-IndexTerm).

Stage Three's results confirmed Team F's predictions: all Original participants ran into one or both of these bugs (Figure 10). Also as per Team F's prediction, when participants ran into the bugs, they verbalized mostly Abi-like facet values: for Bug 5 & 6, 81% (17/21) were Abi-like facet values (orange square outlines left half Figure 10). For example:

P1-O (low-self-efficacy): "I feel like they [the OSS developers] put up barriers because they would want people that really knew what they were doing..."

P1-O (continues): "I'd probably just, like, not work on it."

The lab participants also pointed out mismatches for cognitive styles like process-oriented learning, comprehensive information processing, and risk-aversion to using commands they did not completely understand:

P1-O: "These instructions aren't working super good for me ... if there was explanations a little more."

P3-O: "I don't completely understand ... where to move it [a command] or where to put it."

To address Bug 5, Team F restructured the "Help us with code" section by adding a layer of hierarchy to structurally identify general information about code contributions (*OrgStruct-Hierarchy*). They also reorganized the section topically by OS type (*OrgScheme-Topic*) (Figure 11). Moreover, they added keywords (*Label-IndexTerm*)

in the README.md similar to Bug 3's fix, to more clearly guide newcomers to the right setup instructions for their OS. To fix Bug 6, Team F added explanations to each step in the instructions, in which they made explicit the reason for each step and the need to use a command line terminal for the commands (*Label-IndexTerm*).

Team F's IA fixes paid off: both Abi-like and Tim-like participants improved and the number of participants who ran into problems decreased from 9 to 1, an 89% improvement (Figure 10). Further, although *none* of the Original participants completed the task successfully, *all* participants using the DiversityEnhanced version were able to complete the task—even P14-D, who at first ran into a problem, but overcame it and eventually succeeded.

5 DISCUSSION

5.1 The IA Fixes: Equity and Inclusion

As the results sections have shown, the IA fixes that differentiated the DiversityEnhanced version from the Original version led to a 90% reduction in the bugs that Team F had found to be inclusivity bugs (Section 4's Table 3). However, this leaves unanswered whether these fixes actually contributed to the goals of making the project's

	Motiv*	SE*	Risk*	Info*	Learn*	Motiv*	SE*	Risk*	Info*	Learn*	
P1	•					-	-	-	-	-	P10
P2	•		•			-	-	-	-	-	P11
P3	•		•			-	-	-	-	-	P12
P4	•					-	-	-	-	-	P13
P5		•						•			P14
P6		•				-	-	-	-	-	P15
P7			•			-	-	-	-	-	P16
P8						-	-	-	-	-	P17
P9						-	-	-	-	-	P18
	Original						D	iversityE	nhance	d	

Figure 10: All Original participants but only 1 DiversityEnhanced participant ran into Bug 5 & 6. *, circles, squares: see Figure 6.



Figure 11: Top: Bug 5 before the fix: no scheme or cues/keywords to enable finding instructions for different OS's. Bottom: Bug 5's fix added topic-based scheme and linked to instructions for each OS.

infrastructure (1) more *equitable* and (2) more *inclusive*. For example, equitability could be achieved by helping one group at the expense of another, but that would not achieve inclusivity. Team F's goal was to do both.

First we consider equity. A dictionary definition of equity is "the quality of being fair and impartial" [48]. We measured equity analyzing the lab participants' data, because the participants covered an almost equal number of Abi and Tim facets (recall Figure 5: 22 Abi facet values and 23 Tim facet values in each treatment). Thus, if the lab participants' number of "Abi facets" affected by a bug was greater than the number of "Tim facets", or vice-versa, we conclude that the bug was inequitable in the ways it affected the participants.

By this measure, Bugs 1 & 2 in the Original version were inequitable: together they affected 14 of participants' Abi facets (orange outlines for Figure 6's Original version), compared to only 5 Tim facets (blue outlines). Applying the same measure to the DiversityEnhanced version shows that, although the DiversityEnhanced version was still slightly inequitable—two of participants' Abi facet inequities (2 orange outlines), and zero Tim facet inequities—it was less inequitable than the Original version. Applying the same measures to Bug 3 (Figure 9 - Original: 5 Abi/1 Tim; DiversityEnhanced: 0 Abi/0 Tim) and to Bugs 5 & 6 (Figure 10 - Original: 17 Abi/4 Tim; DiversityEnhanced 2 Abi/1 Tim) also show that the IA fixes likewise reduced the inequities. Thus, we can conclude that the IA fixes did make Project F's infrastructure more equitable.

Inclusion can be computed using a different measure on the same data. According to the dictionary, inclusion is "the action or state of including or of being included within a group or structure" [48]. Applying this definition to being included by a bug fix, we will conclude that the bug fix was inclusive if the number of lab participants' facets affected by a bug decreased from the Original version to the DiversityEnhanced version for participants' Abi facets *and* for participants' Tim facets.

Applying this measure to Bugs 1 & 2 (Figure 6) reveals that, after the fix, participants' Abi facets affected decreased by 12 (from 14 facets affected to 2). Likewise, participants' Tim facets affected decreased by 5 (from 5 facets affected to 0). Since the number of participants' facets affected decreased for participants' Abi facets and for participants' Tim facets, we conclude that the fixes improved inclusivity. Applying the same measures to Bug 3 (Figure 9 - Abi:

Table 5: Inclusivity summary: Team F's IA fixes' effects on the Abi-like facet values (top) and the Tim-like facet values (bottom) were all positive, showing that the IA fixes increased the inclusivity of the prototype across all cognitive styles.

+:More successes in Version DE; -:fewer (zero occurrences). Grayed out: nobody with these facets ran into this bug.

Bug ID	Motiv	SE	Risk	Info	Learn
Bug 1 & 2	+	+	+	+	+
Bug 3	+	+	+		+
Bug 4					
Bug 5 & 6	+	+	+	+	+
Bug 1 & 2	+	+	+	+	+
Bug 3	+			+	
Bug 4					
Bug 5 & 6	+	+	+	+	+

Table 6: Participants' SUS rating scores. (Maximum possible for the subset we used: 32.)

	Original	DiversityEnhanced
Men's Average	12 (6 Men)	19 (3 Men)
Women's Average	12 (3 Women)	22 (5 Women)
Gender-not-stated	N/A	32
Overall Average	12	22

5 Original/0 DiversityEnhanced, Tim: 1 Original/0 DivEnhanced) and Bugs 5 & 6 (Figure 10 - Abi: 17 Orig/2 DivEnhanced, Tim: 4 Orig/1 DivEnhanced) shows that they also improved inclusivity. As Table 5 shows, for every bug and every facet value, participants' Abi-facets and Tim-facets all ran into fewer barriers in the DiversityEnhanced version.

5.2 What about gender?

In some prior literature (e.g., [65]), analyses of these cognitive styles have revealed gender differences. That was also the case for our Stage Three participants' cognitive styles. The participants displayed a range of facet values, but as in other studies, women's facet values tended more "Abi-wards" than the other participants' (Figure 12). These results agree with previous literature that explain how these facets tend to cluster by gender [11]. These results also, when taken together with Figure 6, Figure 9, and Figure 10, show that most of the facets affected by the bugs were those of the women participants.

However, the SUS usability ratings did not differ much by gender. First, as Table 6 shows, the SUS scores of participants who used the Original project were equally low across gender, which may suggest that the Original had a long way to go from everyone's perspective. Second, the SUS scores for participants who used the DiversityEnhanced project were much higher across gender, adding to the body of evidence (e.g., [34, 65]) that designing for oftenoverlooked populations (here, Abi) can benefit everyone.

5.3 The Facet Questionnaire's Validity

As a few other researchers have also done [18, 24, 65], we used the cognitive facet questionnaire (Section 3.3) to collect the participants' facet values. However, we also collected facet values from a second source: participants' verbalization during their tasks. These two sources enabled us to consider the consistency of the questionnaire's responses with the facets that actually arose among the participants.

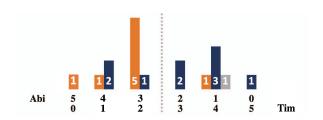


Figure 12: # of women (orange), men (black), and decline-tospecify (gray) with each combination of facets (from facet questionnaire), using the same x-axis scheme (from 5 Abi facets to 5 Tim facets) as Figure 4. Note that the right half of the graph contains only 1 of the 8 women participants.

The data comparing participants' facet questionnaire responses with their actual *in-situ* facet occurrences were detailed earlier in Figure 6, Figure 9, and Figure 10. Outline colors depict the *in-situ* facet occurrences that arose; the shape's fill color depicts the participant's questionnaire response for that facet. (No outline color simply means no evidence arose *in-situ* about that facet.) Thus, when an outline color matches the shape's fill color (questionnaire response), then the questionnaire captured that participant's facet value correctly for the situation.

Overall, 78% of participants' *in-situ* facet verbalizations aligned with their facet questionnaire responses which suggests that the facet questionnaire was a reasonable measure of participants' facet values.

6 THREATS TO VALIDITY

As with any empirical research, our investigation has threats to validity. In this section, we explain threats related to our investigation and ways we guarded against them.

During Stage One, Team F reported the issues found in their project from the perspective of one type of newcomer based on GenderMag's Abi persona. Past research has suggested using the Abi persona first [24], since Abi's facet values tend to be more undersupported in software than those of the other personas (e.g., [10]). However, fixing problems from only this persona's perspectives could leave non-Abi-like newcomers less supported than before. We mitigated this risk by empirically evaluating the fixes with both Abi-like and Tim-like newcomers. That said, some cognitive facets are not considered at all by GenderMag personas, such as memory or attention span, which could be particularly pertinent to people with even mild cognitive disorders. Our investigation did not account for those types of cognitive facets.

As with any investigation with a lab study component, we needed to choose a setting, and our setting (OSS Project F) may not generalize to other OSS projects. The relatively small number of participants (18 in total), which was necessary for tractability of qualitative analysis, also threatens generalizability. In addition, our Stage Three investigation could have uncontrolled differences between the two participant groups. To partially mitigate this threat, we used participants' facet questionnaire responses to assign them to treatments with identical facet distributions (recall Fig 5).

In Stage Three, the identical sequence of the tasks (use-cases), which reflects a workflow common for OSS contributions [59], may have created learning effects that could have influenced the results. Finally, our comparison of facet questionnaire results against verbalizations had only partial data available, since we coded facets from only participants' verbalizations when they encountered a bug, and P5-O's audio for Bug 1 & 2 were corrupted, so we only had observation notes for that participant.

Threats like these can be addressed only by additional studies across a spectrum of empirical methods that isolate particular variables and establish the generality of findings over different types of OSS projects, populations, and other information rich-environments.

7 CONCLUSION

This paper has presented Why/Where/Fix, a systematic inclusivity debugging process. Why/Where/Fix harnesses information architecture, so we also investigated how IA can create inclusivity bugs.

Our setting was an OSS project's technology infrastructure. The "whether" aspects of our RQ1 results revealed that IA can indeed cause inclusivity bugs. In our investigation, the OSS newcomer participants ran into IA-related inclusivity bugs 20 times (Table 3). Our RQ2 "whether" results also revealed that IA can be part of the solution. In our investigation, Team-F's IA fixes reduced the number of inclusivity bugs the participants experienced by 90% (Table 3).

Team F's *hows* of the above results lay in the fault localization capabilities IA brought to Why-Where-Fix:

- *IA and where's*: In Stage One, Team F localized the IA where's behind the inclusivity bugs (Section 4 and Table 4), all but one which the OSS newcomers verified.
- IA and fixes: In Stage Two, Team F fixed the faults, by changing the IA in the ways detailed in Section 4 and summarized in Table 4. The participants in Stage Three showed that Team F's IA fixes helped across the cognitive diversity range of the newcomers in our investigation (Tables 3 and 5).

Key to these results is that these inclusivity fixes lay not in supporting one population at the expense of another, and not in "compromising" to give each population a little less than they need. Rather, as Table 5 illustrated, the fixes produced positive effects across diverse cognitive styles. These results provide encouraging evidence that the Why-Where-Fix process may provide an effective way to increase the equity and inclusion of information-rich environments like OSS projects.

ACKNOWLEDGMENTS

We thank all the study participants for their time and insight. This work is partially supported by the National Science Foundation grants 1901031, 2042324, and 2008089; DARPA grant N66001-17-2-4030; USDA-NIFA/NSF grant 2021-67021-35344; and CNPq grant #313067/2020-1.

REFERENCES

- Paul Ammann and Jeff Offutt. 2016. Introduction to software testing. Cambridge University Press.
- [2] Gabor Aranyi, Paul Van Schaik, and Philip Barker. 2012. Using think-aloud and psychometrics to explore users' experience with a news Web site. *Interacting with Computers* 24, 2 (2012), 69–77.
- [3] Algirdas Avizienis, J-C Laprie, Brian Randell, and Carl Landwehr. 2004. Basic concepts and taxonomy of dependable and secure computing. *IEEE transactions* on dependable and secure computing 1, 1 (2004), 11–33.
- [4] Davide Bolchini, Sebastiano Colazzo, Paolo Paolini, and Daniele Vitali. 2006. Designing aural information architectures. In ACM international conference on Design of Communication. 51–58.
- [5] Amiangshu Bosu and Kazi Zakia Sultana. 2019. Diversity and inclusion in open source software (OSS) projects: Where do we stand?. In 2019 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM). IEEE, 1–11.
- [6] John Brooke. 1996. SUS Quick and Dirty Usability Scale. Usability Evaluation in Industry 189, 194 (1996), 4–7.
- [7] Josep Maria Brunetti. 2013. Design and evaluation of overview components for effective semantic data exploration. In *International Conference on Web Intelligence, Mining and Semantics*. 1–8.
- [8] Josep Maria Brunetti, Rosa Gil, Juan Manuel Gimeno, and Roberto García. 2012. Improved linked data interaction through an automatic information architecture. *International Journal of Software Engineering and Knowledge Engineering* 22, 03 (2012), 325–343.
- [9] Margaret Burnett, Robin Counts, Ronette Lawrence, and Hannah Hanson. 2017. Gender HCI and Microsoft: Highlights from a Longitudinal Study. In *IEEE Symposium on Visual Languages and Human-Centric Computing*. IEEE, 139–143.
- [10] Margaret Burnett, Anicia Peters, Charles Hill, and Noha Elarief. 2016. Finding Gender-inclusiveness Software Issues with GenderMag: A Field Investigation. In ACM Conference on Human Factors in Computing Systems (Santa Clara, California, USA) (CHI '16). ACM, 2586–2598.

- [11] Margaret Burnett, Simone Stumpf, Jamie Macbeth, Stephann Makri, Laura Beckwith, Irwin Kwan, Anicia Peters, and William Jernigan. 2016. GenderMag: A Method for Evaluating Software's Gender Inclusiveness. *Interacting with Computers* 28, 6 (2016), 760–787.
- [12] Gemma Catolino, Fabio Palomba, Damian A. Tamburri, Alexander Serebrenik, and Filomena Ferrucci. 2019. Gender Diversity and Women in Software Teams: How Do They Affect Community Smells?. In ACM/IEEE International Conference on Software Engineering: Software Engineering in Society (Montreal, Quebec, Canada). IEEE Press, 11–20.
- [13] Sally Jo Cunningham, Annika Hinze, and David M Nichols. 2016. Supporting gender-neutral digital library creation: A case study using the GenderMag Toolkit. In *International Conference on Asian Digital Libraries*. Springer, 45–50.
- [14] André de Lima Salgado, Felipe Silva Dias, João Pedro Rodrigues Mattos, Renata Pontin de Mattos Fortes, and Patrick CK Hung. 2019. Smart toys and children's privacy: usable privacy policy insights from a card sorting experiment. In ACM International Conference on the Design of Communication. 1–8.
- [15] Roberto García, Josep Maria Brunetti, Antonio López-Muzás, Juan Manuel Gimeno, and Rosa Gil. 2011. Publishing and interacting with linked data. In International Conference on Web Intelligence, Mining and Semantics. 1–12.
- [16] Chrysoula Gatsou, Anastasios Politis, and Dimitrios Zevgolis. 2012. Novice User involvement in information architecture for a mobile tablet application through card sorting. In IEEE Federated Conference on Computer Science and Information Systems (FedCSIS). IEEE, 711–718.
- [17] Asif Qumer Gill, Nathan Phennel, Dean Lane, and Vinh Loc Phung. 2016. IoT-Enabled Emergency Information Supply Chain Architecture for Elderly People: The Australian Context. *Information Systems* 58 (2016), 75–86.
- [18] Catarina Gralha, Miguel Goulao, and Joao Araujo. 2019. Analysing Gender Differences in Building Social Goal Models: A Quasi-experiment. In IEEE International Requirements Engineering Conference (RE 2019). 12 pages.
- [19] Mikaylah Gross, Joe Dara, Christopher Meyer, and Davide Bolchini. 2018. Exploring Aural Navigation by Screenless Access. In *Internet of Accessible Things*. 1–10.
- [20] Mariam Guizani, Amreeta Chatterjee, Bianca Trinkenreich, Mary Evelyn May, Geraldine J. Noa-Guevara, Liam James Russell, Griselda G. Cuevas Zambrano, Daniel Izquierdo-Cortazar, Igor Steinmacher, Marco A. Gerosa, and Anita Sarma. 2021. The Long Road Ahead: Ongoing Challenges in Contributing to Large OSS Organizations and What to Do. Proc. ACM Hum.-Comput. Interact. 5, CSCW2, Article 407 (Oct. 2021), 30 pages. https://doi.org/10.1145/3479551
- [21] Mariam Guizani, Igor Steinmacher, Jillian Emard, Abrar Fallatah, Margaret Burnett, and Anita Sarma. Feb, 2022. Supplemental Document for How to Debug Inclusivity Bugs? A Debugging Process with Information Architecture. Available at https://figshare.com/s/36e3d2ca390863402790.
- [22] Shelley Gullikson, Ruth Blades, Marc Bragdon, Shelley McKibbon, Marnie Sparling, and Elaine G. Toms. 1999. The Impact of Information Architecture on Academic Web Site Usability. *The Electronic Library* 17, 5 (1999), 293–304.
- [23] Shelley Gullikson, Ruth Blades, Marc Bragdon, Shelley McKibbon, Marnie Sparling, and Elaine G Toms. 1999. The impact of information architecture on academic web site usability. *The Electronic Library* (1999).
- [24] Claudia Hilderbrand, Christopher Perdriau, Lara Letaw, Jillian Emard, Zoe Steine-Hanson, Margaret Burnett, and Anita Sarma. 2020. Engineering Gender-Inclusivity into Software: Ten Teams' Tales from the Trenches. In ACM/IEEE International Conference on Software Engineering.
- [25] Charles G Hill, Maren Haag, Alannah Oleson, Chris Mendez, Nicola Marsden, Anita Sarma, and Margaret Burnett. 2017. Gender-Inclusiveness Personas vs. Stereotyping: Can We Have It Both Ways?. In ACM Conference on Human Factors in Computing Systems (CHI 17). ACM, 6658–6671.
- [26] Geert Hofstede. 2011. Dimensionalizing cultures: The Hofstede model in context. Online readings in psychology and culture 2, 1 (2011), 8.
- [27] Daniel Izquierdo, Nicole Huesman, Alexander Serebrenik, and Gregorio Robles. 2019. Openstack Gender Diversity Report. *IEEE Software* 36, 1 (Jan 2019), 28–33.
- [28] Carlos Jensen, Scott King, and Victor Kuechler. 2011. Joining Free/Open Source Software Communities: An Analysis of Newbies' First Interactions on Project Mailing Lists. In Proceedings of the 2011 44th Hawaii International Conference on System Sciences (Kauai, HI, USA – 4-7 January 2011) (HICSS '11). IEEE Computer Society, Washington, DC, USA, 1–10. https://doi.org/10.1109/HICSS. 2011.264
- [29] Flávia Lacerda, Mamede Lima-Marques, and Andrea Resmini. 2017. An Information Architecture Framework for the Internet of Things. *Philosophy & Technology* (2017), 1–18.
- [30] Florian Lachner, Mai-Anh Nguyen, and Andreas Butz. 2018. Culturally sensitive user interface design: a case study with German and Vietnamese users. In Second African Conference for Human Computer Interaction: Thriving Communities. ACM 1
- [31] Meredith B Larkin. 2020. Board gender diversity, corporate reputation and market performance. *International Journal of Banking and Finance* 9, 1 (2020), 1–26.
- [32] Amanda Lee, Jeffrey C Carver, and Amiangshu Bosu. 2017. Understanding the impressions, motivations, and barriers of one time code contributors to FLOSS

- projects: a survey. In 2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE). IEEE, 187–197.
- [33] Mingran Li, Ruimin Gao, Xinghe Hu, and Yingjie Chen. 2017. Comparing infovis designs with different information architecture for communicating complex information. *Communication Design Quarterly Review* 5, 1 (2017), 43–56.
- [34] Sara Ljungblad and Lars Erik Holmquist. 2007. Transfer scenarios: grounding innovation with marginal practices. In ACM Conference on Human Factors in Computing Systems. ACM, 737–746.
- [35] Christopher Mendez, Hema Susmita Padala, Zoe Steine-Hanson, Claudia Hilderbrand, Amber Horvath, Charles Hill, Logan Simpson, Nupoor Patil, Anita Sarma, and Margaret Burnett. 2018. Open source barriers to entry, revisited: A sociotechnical perspective. In Proceedings of the 40th International conference on software engineering. 1004–1015.
- [36] Craig S. Miller and Roger W. Remington. 2004. Modeling Information Navigation: Implications for Information Architecture. *Human–Computer Interaction* 19, 3 (2004), 225–271.
- [37] Craig S Miller and Roger W Remington. 2004. Modeling information navigation: Implications for information architecture. *Human-computer interaction* 19, 3 (2004), 225–271.
- [38] Peter Morville and Louis Rosenfeld. 2006. Information architecture for the World Wide Web: Designing large-scale web sites. O'Reilly Media, Inc.
- [39] Dawn Nafus. 2012. "Patches Don't Have Gender": What Is Not Open in Open Source Software. New Media & Society 14, 4 (2012), 669–683.
- [40] Gerard Oleksik, Hans-Christian Jetter, Jens Gerken, Natasa Milic-Frayling, and Rachel Jones. 2013. Towards an information architecture for flexible reuse of digital media. In *International Conference on Mobile and Ubiquitous Multimedia*. 1–10
- [41] Open Source Guides. 2019. Retrieved September 12, 2019 from https://opensource.guide/. Accessed on: Sept-3-2019.
- [42] Susmita Hema Padala, Christopher John Mendez, Luiz Felipe Dias, Igor Steinmacher, Zoe Steine Hanson, Claudia Hilderbrand, Amber Horvath, Charles Hill, Logan Dale Simpson, Margaret Burnett, et al. 2020. How gender-biased tools shape newcomer experiences in oss projects. IEEE Transactions on Software Engineering (2020).
- [43] Susmita Hema Padala, Christopher John Mendez, Luiz Felipe Dias, Igor Steinmacher, Zoe Steine Hanson, Claudia Hilderbrand, Amber Horvath, Charles Hill, Logan Dale Simpson, Margaret Burnett, et al. 2020. How Gender-Biased Tools Shape Newcomer Experiences in OSS Projects. IEEE Transactions on Software Engineering (2020).
- [44] Scott E Page. 2019. The diversity bonus: How great teams pay off in the knowledge economy. Princeton University Press.
- [45] Helen Petrie and Christopher Power. 2012. What do users really care about? A comparison of usability problems found by users and experts on highly interactive websites. In ACM Conference on Human Factors in Computing Systems. 2107– 2116.
- [46] Katherine W Phillips, Douglas Medin, Carol D Lee, Megan Bang, Steven Bishop, and DN Lee. 2014. How diversity works. *Scientific American* 311, 4 (2014), 42–47.
- [47] Peter Pirolli. 2007. Information Foraging Theory: Adaptive Interaction with Information. Oxford University Press.
- [48] Oxford University Press. 2019. Lexico US Dictionary. https://www.lexico.com/
- [49] Huilian Sophie Qiu, Alexander Nolte, Anita Brown, Alexander Serebrenik, and Bogdan Vasilescu. 2019. Going Farther Together: The Impact of Social Capital on Sustained Participation in Open Source. In ACM/IEEE International Conference on Software Engineering (Montreal, Quebec, Canada) (ICSE '19). IEEE Press, Piscataway, NJ, USA, 688–699.
- [50] Marc L Resnick and Julian Sanchez. 2004. Effects of organizational scheme and labeling on task performance in product-centered and user-centered retail web sites. *Human factors* 46, 1 (2004), 104–117.
- [51] Gregorio Robles, Laura Arjona Reina, Alexander Serebrenik, Bogdan Vasilescu, and Jesús M González-Barahona. 2014. Floss 2013: A Survey Dataset about Free Software Contributors: Challenges for Curating, Sharing, and Combining. In ACM 11th Working Conference on Mining Software Repositories (MSR 2014). ACM, 306–309
- [52] Álvaro Rocha and Jorge Freixo. 2015. Information Architecture for Quality Management Support in Hospitals. *Journal of Medical Systems* 39, 10 (2015), 125.
- [53] Romisa Rohani Ghahari, Mexhid Ferati, Tao Yang, and Davide Bolchini. 2012. Back navigation shortcuts for screen reader users. In ACM International Conference on Computers and Accessibility. 1–8.
- [54] Romisa Rohani Ghahari, Jennifer George-Palilonis, and Davide Bolchini. 2013. Mobile web browsing with aural flows: an exploratory study. *International Journal of Human-Computer Interaction* 29, 11 (2013), 717–742.
- [55] Louis Rosenfeld, Peter Morville, and Jorge Arango. 2015. Information Architecture: For the Web and Beyond. O'Reilly Media, Inc.
- [56] Paul Van Schaik, Raza Habib Muzahir, and Mike Lockyer. 2015. Automated computational cognitive-modeling: goal-specific analysis for large websites. ACM Transactions on Computer-Human Interaction (TOCHI) 22, 3 (2015), 1–29.

- [57] Arun Shekhar and Nicola Marsden. 2018. Cognitive Walkthrough of a learning management system with gendered personas. In 4th Conference on Gender & IT. 191–198.
- [58] Igor Steinmacher, Tayana Conte, Marco Aurélio Gerosa, and David Redmiles. 2015. 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. 1379–1392.
- [59] Igor Steinmacher, Tayana Uchoa Conte, Christoph Treude, and Marco Aurélio Gerosa. 2016. Overcoming Open Source Project Entry Barriers with a Portal for Newcomers. In ACM/IEEE International Conference on Software Engineering (ICSE'16). ACM, 273–284.
- [60] Steven E Stemler. 2004. A Comparison of Consensus, Consistency, and Measurement Approaches to Estimating Interrater Reliability. *Practical Assessment, Research & Evaluation* 9, 4 (2004), 1–19.
- [61] Simone Stumpf, Anicia Peters, Shaowen Bardzell, Margaret Burnett, Daniela Busse, Jessica Cauchard, and Elizabeth Churchill. 2020. Gender-Inclusive HCI Research and Design: A Conceptual Review. Foundations and Trends in Human-Computer Interaction 13, 1 (2020), 1–69.
- [62] Sarah J Swierenga, Jieun Sung, Graham L Pierce, and Dennis B Propst. 2011. Website design and usability assessment implications from a usability study with visually impaired users. In *International Conference on Universal Access in Human-Computer Interaction*. Springer, 382–389.
- [63] Bogdan Vasilescu, Daryl Posnett, Baishakhi Ray, Mark GJ van den Brand, Alexander Serebrenik, Premkumar Devanbu, and Vladimir Filkov. 2015. Gender and Tenure Diversity in Github Teams. In ACM Conference on Human Factors in Computing Systems (Seoul, Republic of Korea) (CHI '15). ACM, ACM, New York, NY, USA, 3789–3798.
- [64] Markel Vigo and Simon Harper. 2013. Challenging information foraging theory: screen reader users are not always driven by information scent. In ACM Conference on Hypertext and Social Media. 60–68.
- [65] Mihaela Vorvoreanu, Lingyi Zhang, Yun-Han Huang, Claudia Hilderbrand, Zoe Steine-Hanson, and Margaret Burnett. 2019. From Gender Biases to Gender-inclusive Design: An Empirical Investigation. In ACM Conference on Human Factors in Computing Systems (Glasgow, Scotland Uk) (CHI '19). ACM, New York, NY, USA, Article 53, 14 pages.
- [66] Richard Saul Wurman and Joel Katz. 1975. Beyond graphics: The architecture of information. AIA Journal 10 (1975), 40–45.
- [67] Tao Yang, Mexhid Ferati, Yikun Liu, Romisa Rohani Ghahari, and Davide Bolchini. 2012. Aural browsing on-the-go: listening-based back navigation in large web architectures. In ACM Conference on Human Factors in Computing Systems. 277–286.