

Journal Pre-proof

An Empirical Study of Web Visual Flakiness: Characterisation and Fix Strategies

Yu Pei, Jeongju Sohn, Mike Papadakis

PII: S0164-1212(26)00060-9
DOI: <https://doi.org/10.1016/j.jss.2026.112826>
Reference: JSS 112826



To appear in: *The Journal of Systems & Software*

Received date: 1 November 2025
Revised date: 15 January 2026
Accepted date: 19 February 2026

Please cite this article as: Yu Pei, Jeongju Sohn, Mike Papadakis, An Empirical Study of Web Visual Flakiness: Characterisation and Fix Strategies, *The Journal of Systems & Software* (2025), doi: <https://doi.org/10.1016/j.jss.2026.112826>

This is a PDF of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability. This version will undergo additional copyediting, typesetting and review before it is published in its final form. As such, this version is no longer the Accepted Manuscript, but it is not yet the definitive Version of Record; we are providing this early version to give early visibility of the article. Please note that Elsevier's sharing policy for the Published Journal Article applies to this version, see: <https://www.elsevier.com/about/policies-and-standards/sharing#4-published-journal-article>. Please also note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2026 Elsevier Inc. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

An Empirical Study of Web Visual Flakiness: Characterisation and Fix Strategies

Yu Pei

yu.pei@uni.lu University of Luxembourg, Luxembourg

Jeongju Sohn¹

jeongju.sohn@knu.ac.kr Kyungpook National University, Korea

Mike Papadakis

michail.papadakis@uni.lu University of Luxembourg, Luxembourg

Abstract

Unpredictable behaviour in web front-ends undermines test reliability and user experience, often manifesting as inconsistent or unstable visual output. Prior work on web front-end flakiness has mainly focused on asynchronous waits, platform-specific behaviour, and logic errors, many of which stem from client-side execution, back-end dependencies, or timing variations across browsers. However, little attention has been given to pure visual flakiness—instabilities in layout, rendering, and styling that cause the same web content to appear differently across runs or environments.

In this work, we present a focused empirical study dedicated to web visual flakiness, defined as flakiness that manifests as structural or stylistic inconsistencies in web front-ends. We collected 262 cases of web visual flakiness: 144 from 31 open-source web projects and 118 from the Chromium project. Our findings show that web visual flakiness constitutes a notable portion of overall flaky cases and that visual changes are more prone to flakiness than non-visual ones, underscoring the need for focused investigation. We categorised the collected instances into two primary dimensions—structure-related (59.9%) and style-related (40.1%)—and further refined them into five concrete sub-categories by manifestation (three structure-related and two style-related). Finally, our analysis of developer fixes shows that while many address specific structural or stylistic inconsistencies, others are applied across multiple categories, revealing shared repair strategies.

18 *Keywords:*

19 Flaky Behaviour, Web Visuals, Web Testing, Empirical Study

20 1. Introduction

21 As web applications grow increasingly complex, ensuring their reliability
22 has become a critical challenge. Among many obstacles, flaky behaviours,
23 such as test outcomes that show non-deterministic results across repeated
24 runs, pose a constant threat to automated testing Luo et al. (2014); Parry
25 et al. (2021); Lam et al. (2019a). This issue is particularly pronounced in the
26 visual layer, where minor inconsistencies, such as browser-specific differences,
27 asynchronous content loading, and ordering of DOM updates, can lead to
28 unpredictable rendering and test failures Eck et al. (2019); Parry et al. (2022);
29 Habchi et al. (2022); Gruber and Fraser (2022); Haben et al. (2024).

30 The visual layer encompasses the code and frameworks that define and
31 render visual presentation in web rendering environments, including markup
32 (e.g., HTML), stylesheets (e.g., CSS), and client-side logic (e.g., JavaScript).
33 These components collectively control the structure and styling of the visual
34 layer, making them central to both user experience and test stability. How-
35 ever, even subtle changes to these components can lead to visual inconsisten-
36 cies that are difficult to detect and reproduce due to their non-deterministic
37 nature (i.e., flakiness). Responsive design and dynamic content further am-
38 plify this instability across browsers and environments.

39 While previous research has investigated flaky tests in web and mobile
40 environments, such as interaction-based GUI flakiness (e.g., Romano et al.
41 (2021)) and deterministic UI rendering defects (e.g., Mahajan et al. (2016);
42 Alameer and Halfond (2016); Althomali et al. (2021)), the focus has primar-
43 ily been on functional-level non-determinism (e.g., async waits, improper
44 test script logics, or state mismatches) or deterministic visual inconsistencies
45 (e.g., cross-browser layout bugs). Meanwhile, studies on web flaky tests have
46 primarily investigated non-determinism in functional or interaction-driven
47 contexts, addressing causes such as asynchronous waits, platform-specific
48 behaviour, or state management Parry et al. (2021); Lam et al. (2019a).

49 Taken together, prior studies have examined interaction-based GUI flaki-
50 ness and other forms of functional non-determinism in UI behaviour, as well
51 as deterministic visual faults, such as cross-browser layout bugs. However,
52 they have not explicitly isolated non-deterministic visual flakiness—often

53 manifesting as unstable layout, rendering, or styling—as a distinct empirical
54 category. In this work, we investigate such visual flakiness as instabilities
55 rooted in the visual layer that are inconsistently observable across execu-
56 tions and environments, bridging the gap between web flakiness research and
57 studies of visual correctness.

58 At the same time, web-visual testing is widely practised in industry, for
59 instance, through screenshot-based or layout-diff validation Yeh et al. (2009);
60 Choudhary et al. (2010, 2011); Stocco et al. (2018), yet its flaky behaviour
61 remains largely underexplored. Such flakiness is commonly observed in mod-
62 ern industrial web projects, highlighting the need for a dedicated empirical
63 investigation Rwemalika et al. (2019); Zou et al. (2014); Nass et al. (2023).

64 We present this study as a focused empirical investigation that system-
65 atically characterises non-deterministic visual flakiness in web applications
66 as a distinct analytical focus. To our knowledge, this study provides the
67 first systematic empirical analysis that explicitly focuses on web visual flak-
68 iness, examining its prevalence, characteristics, and developer responses. By
69 analysing 262 flaky cases from 31 open-source web projects and Chromium,
70 we identify how, where, and when visual flakiness manifests in real-world de-
71 velopment. Based on recurring patterns, we categorise them into structure-
72 related and style-related types. We also analyse how developers address
73 visual flakiness in practice, inspecting how these fixes relate to the defined
74 categories, and identifying common and cross-cutting fix strategies. In the
75 end, we address the following research questions:

- 76 • **RQ1:** What is the prevalence of flaky behaviours in web visuals?
- 77 • **RQ2:** What are the most common types of web visual flakiness, and
78 how do they manifest in practice?
- 79 • **RQ3:** What fix patterns can be employed to mitigate the flaky be-
80 haviour in web visuals?

81 Our findings, combined, highlight the importance of understanding visual
82 flakiness to improve the reliability of web testing. Below summarises the key
83 contributions:

- 84 • We conduct a focused empirical study of web visual flakiness as a dis-
85 tinct class of flaky behaviour, analysing 262 real-world cases across 31
86 open-source web projects and the Chromium codebase, and release the
87 resulting curated dataset to support future research.
- 88 • We provide a taxonomy of web visual flakiness, distinguishing two high-
89 level types, i.e., structure-related and style-related, based on recurring

90 manifestation patterns, and further refining each type into five concrete
91 subcategories.

- 92 • We investigate how developers handle web visual flakiness in practice,
93 identifying common, cross-cutting repair strategies and analysing their
94 relationships with different flakiness categories.

95 2. Methodology

96 2.1. *Motivating Example*

97 This study is motivated by the growing complexity of modern web front-
98 ends, where frequent modifications to DOM structure and visual styling can
99 introduce instability in the visual behaviour observed during execution. The
100 front-end of the web application is the layer where users directly interact with
101 the system, making visual instability particularly impactful. Even minor
102 inconsistencies, such as unexpected element positioning or missing visual
103 focus, can disrupt user actions and degrade user.

104 Consider the example in Figure 1a², where clicking the *Edit* menu item is
105 intended to automatically focus the text box, enabling users to begin editing
106 immediately (Figure 1b). However, as shown in Figure 1c, the text box some-
107 times fails to gain focus, requiring users to manually click it again. While
108 the underlying logic remains unchanged, this execution-time visual instabil-
109 ity, manifesting as inconsistent focus behaviour (i.e., flakiness), disrupts user
110 workflow and causes frustration, particularly in applications where users ex-
111 pect quick and fluid interactions. Such visual flakiness is often difficult to
112 detect and diagnose. Understanding the nature and recurring patterns of
113 such flakiness requires reasoning about instability in the visual layer, which
114 is inherently subtle, brittle, and sensitive to small variations Mahajan et al.
115 (2016); Pei et al. (2025).

116 2.2. *Study Design*

117 We conducted an empirical study to characterise flaky behaviours in web
118 visuals in real-world projects. To this end, we first collected commits and,
119 when available, their respective pull requests or issue reports related to web
120 visual flakiness from diverse web projects. We then investigated the observed
121 types of visual flakiness derived from recurring patterns and how developers

²<https://github.com/angular/components/issues/18750>

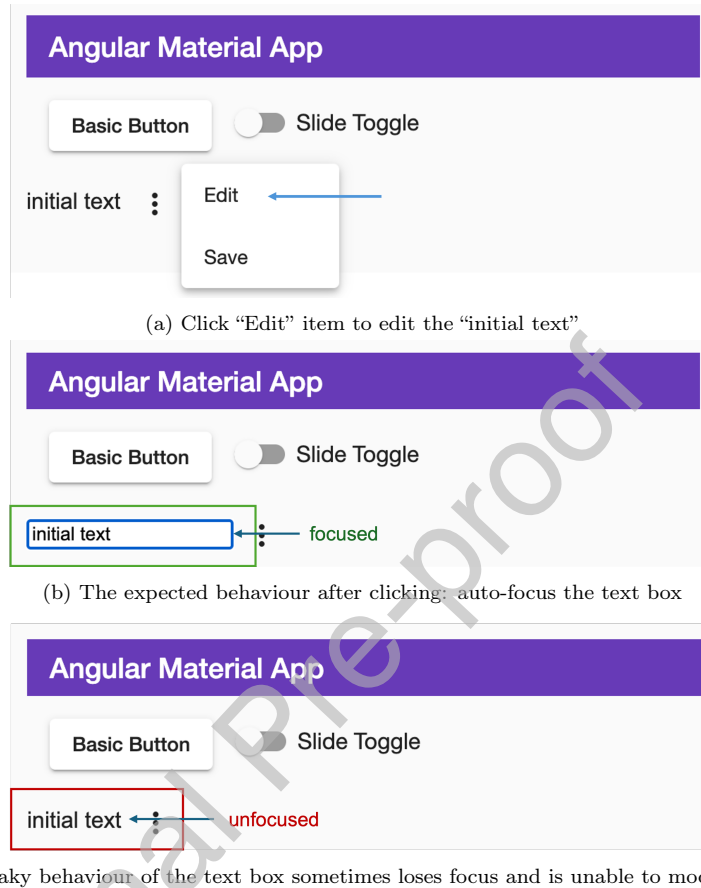


Figure 1: Example: flaky behaviour in the text box autofocus

122 address these issues in practice, including common fix strategies. To support
 123 this investigation, we performed a combination of manual inspection and
 124 LLM-assisted annotation to categorise both visual flakiness types and their
 125 corresponding fix strategies.

126 2.2.1. Terminology

127 To clarify the focus of our study, we introduce two key terms that define
 128 what we consider visual-related flakiness in web projects. These definitions
 129 distinguish our scope from general categories of web or UI flakiness that
 130 include functional or interaction-driven behaviours: *visual commits* and *web*
 131 *visual flakiness*.

- 132 • **Visual commits** are commits that include changes that influence the
133 visual behaviour or appearance of a web project, such as those related
134 to DOM structure, styling, or rendering logic.
- 135 • **Web visual flakiness** refers to non-deterministic, intermittent incon-
136 sistencies that manifest in the visual layer of web projects. Such in-
137 consistencies may appear, for instance, as layout shifts, rendering dif-
138 ferences, or style-related visual inconsistencies.

139 Based on the above terminology, we define a web visual flakiness-related
140 commit as any visual commit that exhibits flaky behaviour in the visual layer.
141 In contrast, non-visual commits are those whose changes do not influence the
142 rendering or presentation of web content, such as modifications unrelated to
143 the DOM structure or styling.

144 2.2.2. Dataset Collection

145 To collect web visual flakiness cases³, we followed a systematic procedure
146 similar to prior flakiness studies Luo et al. (2014); Romano et al. (2021);
147 Hashemi et al. (2022). We focus on open-source web application projects
148 on GitHub, where commit-level information is available and enables detailed
149 analysis of code changes and their context.

150 We first identified repositories primarily using web development languages
151 such as JavaScript and HTML and selected those with over 100 stars to
152 ensure their representativeness. In total, we collected 31 projects comprising
153 489,926 commits (Table 1, third column).

154 Next, we identified flaky commits by searching commit messages and,
155 when available, associated issues or pull requests for flakiness-related key-
156 words. We use “*flak**”, “*intermit**”, and “*unstable*” as our flakiness-related
157 keyword set. This yielded 1,377 candidate commits related to flakiness
158 (*#Flaky* in Table 1). Our primary interest lies in flakiness arising from
159 changes whose effects manifest in the visual behaviour of web applications.
160 Hence, in parallel with the flaky-commit search, we conducted a keyword
161 search for visual commits using visual-related keywords, which contains “*e2e*”,
162 “*UI*”, “*GUI*”, “*style*”, “*CSS*”, “*visual*”, “*layout*”, and “*display*”. If a commit
163 message, issue, or pull request contained any of these visual keywords, we

³We use the term flakiness case to refer to an instance of flaky behaviour, regardless of whether it originates from test code or application code.

Table 1: Subjects. $\#Commits$, $\#Commits^{Visual}$, $\#Flaky$, $\#FlakyVisual$ are the number of total commits, visual-related commits, all flaky-related cases, and visual-related flakiness cases

Projects	$\#Stars$	$\#Commits$	$\#Commits^{Visual}$	$\#Flaky$	$\#FlakyVisual$
Chromium	18.9k	-	-	7780	118
vets-website	240	29,633	3506	235	19
ionic-framework	51k	13,970	1736	62	17
web-platform-tests	4.9k	62,667	17,486	621	15
angular-components	24.3k	12,162	2339	40	12
matrix-react-sdk	1.1k	48,466	6749	78	10
storybook	84.1k	63,182	2858	64	7
material-components	17.1k	7,453	1674	31	7
fluentui	18.4k	18,590	2768	18	6
ant-design	92.1k	28,580	4315	14	5
eui	6.1k	6,187	1,127	18	5
openverse	243	11,200	621	44	5
material-ui	93.5k	25,656	3635	21	4
chakra-ui	37.6k	9,938	1865	8	4
tailwindcss	82.3k	5,568	806	7	3
bootstrap	170k	22,900	2378	5	2
vue-tify	39.7k	16,186	1831	5	2
quasar	25.8k	14,427	1889	5	2
svelte	78.7k	9,344	830	9	2
ghost	47.1k	39,504	4054	50	2
element	54.1k	4,612	527	4	2
blueprint	20.7k	3,581	447	4	2
balm-ui	506	3,131	129	2	2
bulma	49.2k	1,881	120	1	1
nuxt	54.3k	6,754	268	3	1
semi-design	8.4k	3,785	379	2	1
floating-ui	29.8k	2,221	246	3	1
vue.draggable	20.1k	521	22	1	1
vee-validate	10.8k	4,777	136	2	1
vue-material	9.9k	1,132	217	1	1
vite	67.7k	7,051	481	16	1
xyflow	25k	4,867	187	3	1
Open-source	240-170k	489,926	65,626	1377	144
TOTAL	-	489,926	65,626	9157	262

164 considered it likely to affect the visual parts of the application.⁴ In total,
 165 65,626 visual commits were identified ($\#_{commits}^{Visual}$, Table 1).

166 Lastly, to pinpoint web visual flakiness-related commits, we intersected
 167 the sets of flaky and visual commits. Two authors manually reviewed each
 168 candidate, examining the code changes⁵ and contextual information, such as
 169 commit messages and, when available, related issues, pull requests, or devel-
 170 oper comments, to ensure that each case genuinely involved web UI flakiness.
 171 This conservative strategy—combining the intersection of strict flakiness and
 172 visual keywords with manual validation—was chosen to prioritise data pre-
 173 cision over exhaustive coverage to support the confidence in our findings
 174 from subsequent qualitative analysis. Section 5.2 further discusses this data
 175 sampling strategy. The last column ($\#_{FlakyVisual}$) in the "Open-source"
 176 row of Table 1 reports 144 identified web visual flakiness cases across the 31
 177 open-source projects.

178 **Chromium.** To enhance the representativeness of our study, we also in-
 179 cluded Chromium, a major industrial browser project. For Chromium, we
 180 adapted the keyword-based matching procedure to its issue tracker⁶, where
 181 flaky behaviours are reported and managed.

182 Specifically, we searched issue reports for flakiness-related keywords and
 183 filtered them with visual-related terms, using the same keyword sets as for
 184 the open-source projects. Although Chromium maintains a mirrored GitHub
 185 repository, the correspondence between issue reports and individual commits
 186 is often incomplete, indirect, or missing. As a result, commit-level infor-
 187 mation (e.g., code diffs and commit statistics) is not consistently available
 188 for Chromium and, thus, commit-level analyses are not performed for this
 189 dataset. The collection process yielded 118 web visual flakiness cases from
 190 Chromium (the second row in Table 1), resulting in a total of 262 web visual

⁴While some of the collected commits labelled as "visual" based on the keyword match-
 ing may not directly interact with visual parts, this does not affect the conclusion of
 our analysis in Section 3.1, as the more precise identification decreases the denominator
 ($\#_{commits}^{Visual}$) in the frequency analysis, further supporting the findings.

⁵We mainly considered changes to files written in HTML, JavaScript, TypeScript, CSS,
 or SCSS, as these languages commonly define and influence the visual structure and style
 of web applications.

⁶<https://issues.chromium.org/issues>. Chromium issues hosted on this platform include
 both current records and historical issues migrated from the legacy Monorail system, which
 we analyse in this study.

191 flakiness cases across all studied subjects (the last row in Table 1).

192 The two primary data sources employed in this investigation exhibit a no-
193 table structural distinction. The open-source web projects dataset is derived
194 from GitHub commits, which enables the identification of web visual flakiness
195 and corresponding fixes through commit messages and code diffs. In contrast,
196 the Chromium dataset is derived from issue tracker records, where flaky be-
197 haviours are reported and discussed, often without a reliable correspondence
198 to individual commits. Consequently, analyses that rely on commit-level
199 granularity—such as LLM-assisted categorisation (Section 2.2.4) and flaki-
200 ness frequency per commit (RQ1)—are only conducted for the open-source
201 projects. For Chromium, the same taxonomic framework and validation pro-
202 cedure are applied to manually analyse issue descriptions, comments, and
203 linked patches, ensuring methodological consistency despite differences in
204 data availability.

205 2.2.3. Web Visual Flakiness Taxonomy

206 To better understand and analyse web visual flakiness, we manually in-
207 spected the collected instances in our dataset to search for recurring patterns.
208 This analysis revealed two broad families of issues: some involved unstable
209 DOM structures, where elements were not reliably present or positioned (e.g.,
210 elements missing, misaligned, or late), while others involved changes in ap-
211 pearance due to inconsistencies in style interpretation or resource loading
212 (e.g., CSS conflicts or delayed fonts). **Based on these recurring patterns, we**
213 **derive a taxonomy of web visual flakiness consisting of two high-level families:**
214 **structure-related and style-related.**

- 215 • **Structure-related flakiness:** occurs when visual inconsistencies arise
216 due to variations in how elements are rendered, positioned, or updated.
217 These inconsistencies may stem from changes to the DOM structure
218 (e.g., layout shifts due to dynamic element insertion or removal), asyn-
219 chronous or delayed updates to DOM elements, or failures in loading
220 content that result in missing or incomplete visual components; this
221 type of flakiness occurs regardless of whether the DOM structure is
222 explicitly modified, as even stable structures can exhibit inconsistent
223 rendering under dynamic or environment-dependent conditions. Hence,
224 structure-related flakiness involves all inconsistencies that relate to the
225 state of the DOM structure, including element presence, positioning,
226 and readiness, where positioning refers to whether an element is cor-

227 rectly realised in the layout rather than how its visual style is rendered.

- 228 • ***Style-related flakiness***: occurs when visual inconsistencies arise due
229 to variations in how style properties (e.g., CSS rules, font rendering,
230 or visual appearance attributes) are applied, interpreted, or available
231 across different environments, browsers, or rendering engines at run-
232 time. These inconsistencies affect the visual appearance of elements
233 without involving changes to the DOM structure [or the spatial pres-
234 ence and placement of elements](#). Style-related flakiness primarily affects
235 visual properties such as colour, size, and font.

236 In summary, [structure-related flakiness](#) concerns which visual elements
237 are present and where they appear on the page, while [style-related flakiness](#)
238 concerns how visual elements look. These two represent high-level categories
239 in our taxonomy. Each category includes multiple concrete subcategories that
240 capture distinct recurring manifestations observed in our dataset, which can
241 help developers diagnose and address flaky behaviours. For [structure-related
242 web visual flakiness](#), we identify three primary subcategories:

- 243 • **DOM Structure Issues**: flakiness where DOM elements (i.e., struc-
244 tural components) are not reliably or fully available in the DOM. This
245 includes elements being delayed, missing, removed, or incompletely ren-
246 dered due to asynchronous behaviour or race conditions, as well as
247 explicit structural changes (e.g., adding, removing, or reordering ele-
248 ments). These are existential problems independent of event execution
249 – *Whether and how elements are rendered or present*.
- 250 • **DOM Layout Instability**: flakiness where the spatial arrangement of
251 elements is inconsistent, leading to overlapping, misaligned, or shifted
252 components. Such effects may be triggered by structural changes such
253 as added, removed, or repositioned nodes. The key distinction from
254 DOM Structure Issues is that the inconsistency appears in layout po-
255 sition rather than in element presence or rendering – *How structural
256 elements are spatially positioned*.
- 257 • **Timing-sensitive Visual Triggers**: flakiness that occurs when visual
258 behaviours or event handling are triggered too early or too late relative
259 to when structural elements become ready in the DOM. Examples in-
260 clude situations where expected responses, such as focus, visibility, or
261 transitions, fail because the element was inserted, updated, or removed
262 at the wrong time relative to the event execution. In other words, the

263 failure arises from a synchronisation problem between event execution
264 and element readiness, and not about the presence or position within
265 the DOM. – *When those elements become available relative to an event*
266 *or trigger.*

267 For *style-related web visual flakiness*, we identified two primary subcate-
268 gories:

- 269 • **Style Interpretation Issues:** flakiness, where the visual appearance
270 of elements is inconsistent because style rules are interpreted or
271 applied differently at runtime. Examples include variations across envi-
272 ronments (e.g., differences in CSS properties, such as flexbox and font),
273 as well as inconsistencies within the same environment (e.g., cascade or-
274 der, specificity conflicts, or dynamic overrides via inline style injection)
275 – *How styles are interpreted and applied.*
- 276 • **Style Resource Latency:** flakiness where the visual appearance changes
277 inconsistently due to delays in loading style-dependent resources. Ex-
278 amples include elements appearing in a fallback or unstyled form due
279 to yet-to-be-loaded style resources, such as fonts or stylesheets – *When*
280 *style-dependent resources are applied.*

281 2.2.4. LLM-assisted categorisation

282 To support consistent and scalable annotation, we employed a large lan-
283 guage model (gpt-4o-mini) to assist in labelling verified web visual flakiness
284 instances from open-source web projects. This LLM-assisted procedure was
285 applied across multiple analysis dimensions defined in this study, including
286 the categorisation of visual flakiness types based on the taxonomy introduced
287 in Section 2.2.3, as well as the categorisation of corresponding developer fix
288 strategies analysed in the later research question (RQ3).

289 For each instance, a structured description including commit metadata
290 and messages, code diffs, and, if available, pull request details or linked issue
291 reports, was processed using a predefined prompt to generate a preliminary
292 label. After preliminary labels were generated through LLM, two authors
293 manually reviewed all label assignments. The LLM-generated labels served
294 as an initial reference during the manual review process, helping structure the
295 inspection and support consistent interpretation of the category definitions
296 across the dataset.

Prompt Template

```

// Definitions (provided to the model)
Flakiness Type Definitions. Descriptions of web visual flakiness types
(three structure-related and two style-related), e.g.,
1. DOM structure issues: ...
2. ...

Fix Strategy Definitions. Descriptions of fix strategies, e.g.,
1. Manage element changes: ...
2. ...

// Task
You are categorising the visual flakiness type and the corresponding fix strategy
of a web visual flakiness instance.

// Supporting Evidence
Consider code diffs and issue or pull request descriptions together with commit
metadata (e.g., author, date, statistics), commit messages, PR state/labels/-
commits/files, and issue labels/state as supporting evidence.

// Instructions
Decide the most likely single flakiness type and single fix strategy from the
allowed lists below. Use EXACT strings from the allowed sets. If the evidence
is sparse, select the closest plausible category and assign a lower confidence.

Allowed Flakiness Types: {...}
Allowed Fix Strategies: {...}

// Output Format
Return STRICT JSON with keys:
  flakiness_type (string),
  fix_strategy (string),
  confidence (float 0–1),
  rationale (string),
  evidence (string).
Use only allowed values.

```

Figure 2: System prompt used for LLM-assisted categorisation. The system prompt provides type definitions of visual flakiness and fix strategy. It restricts label(type) selection to a predefined set (*Allowed*), and specifies the required structured output format. Statements prefixed with "//" are comments added for readability and are not part of the prompt text itself.

Prompt Template

```
// Instance-specific context (user prompt)
```

```
Context: {A structured description of the flakiness instance, including commit metadata, commit messages, code diffs, and, when available, titles and descriptions from related pull requests or linked issue reports.}
```

Figure 3: User prompt for flakiness instance-specific context. The user prompt provides structured contextual information for each flakiness instance, including commit metadata, code diffs, and, when available, related issue or pull request descriptions.

297 The LLM’s confidence scores associated with each label assignment were
 298 used to prioritise manual review effort. For label assignments with confi-
 299 dence scores of 0.7 or higher, two authors performed a confirmatory check to
 300 verify that the suggested category—whether a visual flakiness subcategory
 301 or a fix strategy—aligned with the observed evidence in the code diffs and
 302 commit messages associated with resolving the reported flaky behaviour. If
 303 the confidence score was below 0.7, the authors jointly re-reviewed all rele-
 304 vant artefacts (e.g., code diffs, commit messages, and related pull requests),
 305 performing a more in-depth manual inspection and resolving any ambiguities
 306 through discussion until consensus was reached.

307 The final categorisation of each instance, across both visual flakiness types
 308 and fix strategies, was determined based on this combined LLM-assisted and
 309 manual review process. No additional visual flakiness subcategories emerged
 310 in the final categorisation results beyond those derived from the initial pat-
 311 tern analysis described in Section 2.2.3 during the categorisation; the result-
 312 ing fix strategy assignments are used solely to support the analysis in RQ3
 313 (Section 3.3), where the fix strategy categories are introduced and discussed.

314 Figures 2 and 3 illustrate the prompts used for the LLM-assisted categori-
 315 sation. The system prompt provides natural-language definitions of visual
 316 flakiness types (sub-categories) and fix strategies, specifies the categorisation
 317 task and the types of supporting evidence to consider, lists the allowed labels
 318 for each dimension, and defines the required structured output format and
 319 labelling rules. The user prompt provides instance-specific context for each
 320 case, such as commit metadata, code diffs, and, when available, related issue
 321 or pull request descriptions.

322 As discussed in Section 2.2.2, our dataset also includes web visual flak-
 323 iness cases from Chromium, for which commit-level artefacts required for

LLM-assisted categorisation (e.g., code diffs and commit statistics) are not consistently available. Therefore, the LLM-assisted categorisation procedure above is not applied to the Chromium cases. Instead, Chromium cases are manually labelled by two authors through collaborative review. Each case is assigned to one of the existing web visual flakiness subcategories defined in our taxonomy and derived from the open-source web projects, based on issue descriptions and discussion comments recorded in the issue tracker and, when explicitly available, linked patches or code review references. The same two authors jointly review each case and resolve disagreements through discussion until consensus is reached. No new flakiness subcategories were introduced beyond those identified from the open-source web projects.

2.2.5. Research Questions

This study aims to deepen the understanding of flaky behaviours in web visuals by exploring the following research questions:

RQ1: *What is the prevalence of flaky behaviours in web visuals?*

To assess the prevalence and impact of visual-related flakiness in web applications and their tests, we conduct three investigations from different perspectives. First, we examine the prevalence of visual commits to establish a baseline for how frequently visual components are modified in our target projects. This provides initial context for evaluating the role of visual changes in flaky behaviour. Next, we analyse the ratio of flaky commits to total commits within visual changes and within non-visual ones (e.g., backend or API). This step assesses whether visual commits are more likely to introduce flaky behaviour than other types of code changes. Lastly, we calculate the proportion of visual flaky commits among all flaky commits across projects. While the previous step compares the relative likelihoods of flakiness between visual flakiness and non-visual flakiness, this step evaluates the overall contribution of visual flakiness to the unreliability of web projects. Note that Chromium is excluded from this analysis, as its flakiness is tracked through issues, not commits (Section 3.1).

RQ2: *What are the most common types of web visual flakiness, and how do they manifest in practice?*

To answer this RQ, we analyse the distribution of visual flakiness across two primary categories—structure-related and style-related—and their respective concrete subtypes (three structure and two style). This two-level

359 categorisation allows hierarchical reasoning about the sources of visual insta-
 360 bility, enabling both an overall understanding and a fine-grained examination
 361 of individual patterns. In the following analysis, we further illustrate repre-
 362 sentative cases for each subtype to show how these flaky behaviours manifest
 363 in real-world projects, together with a quantitative analysis of their distribu-
 364 tion and trends.

365 **RQ3:** *What fix patterns can be employed to mitigate the flaky behaviour in*
 366 *web visuals?*

367 Flaky behaviours degrade test stability, complicate debugging, and thereby
 368 diminish user experience. While prior work has proposed general solutions
 369 for flakiness, web visual flakiness may require more targeted strategies. By
 370 analysing how developers address web visual flakiness through the changes
 371 applied to address it, we aim to provide practical insights into fix strategies
 372 for different flakiness types, helping developers mitigate them effectively.

373 3. Results

374 3.1. RQ1: Prevalence of Web Visual Flakiness

375 Figure 4 presents the prevalence of visual commits across the 31 open-
 376 source projects, measured as the ratio of visual commits to total commits
 377 in each project. Chromium is excluded from this analysis due to the un-
 378 availability of reliable commit-level data (Section 2.2.2). On average, visual
 379 commits constitute 11.7% of all commits in these studied projects. This
 380 finding suggests that changes to visual components account for a substantial
 381 portion of ongoing web development activities and may therefore contribute
 382 significantly to overall project instability, warranting dedicated inspection.

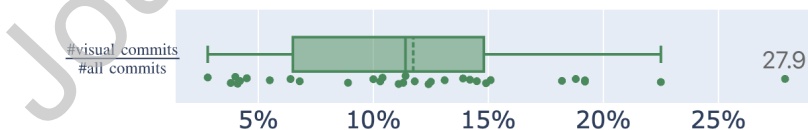
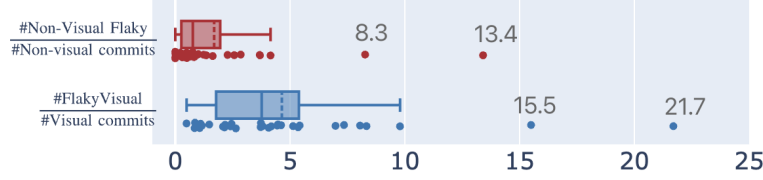
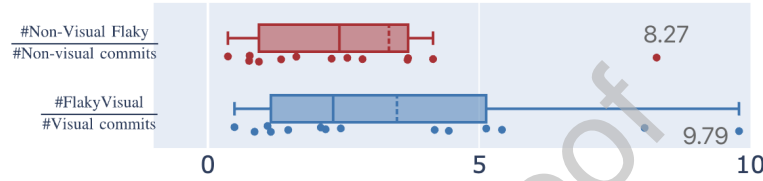


Figure 4: $\#$ visual commits to $\#$ all commits per project (%). The mean is 11.7%, and the median is 11.4%.

383 Next, we compare the frequency of flaky commits—i.e., the likelihood of
 384 commits exhibiting flakiness—between visual and non-visual commits. Fig-
 385 ure 5 presents the distribution of flaky commit ratios for each project (each
 386 dot represents one project) within the visual and the non-visual commits.



(a) $\#Flaky > 0$ (no threshold). For $flakyfreq_{visual}$, mean = 4.6, median = 3.7; for $flakyfreq_{non-visual}$, mean=1.7, median=0.8 (‰)



(b) $\#Flaky > 10$. For $flakyfreq_{visual}$, mean = 3.5, median = 2.3; for $flakyfreq_{non-visual}$, mean = 3.3, median = 2.4 (‰)

Figure 5: Flaky frequency. $flakyfreq_{visual} = \frac{\#FlakyVisual}{\#Visual\ commits}$ vs. $flakyfreq_{non-visual} = \frac{\#Non-Visual\ Flaky}{\#Non-visual\ commits}$. (‰) denotes commits per thousand.

387 In this study, we employed a keyword-based search to collect relevant
 388 commits and issue reports, following prior flakiness studies Romano et al.
 389 (2021). Despite combining multiple metadata sources (commit messages, is-
 390 sue reports, and pull requests), this approach inevitably misses some cases
 391 due to incomplete or inconsistent textual descriptions—for instance, when
 392 developers use ambiguous language or omit flakiness-(or visual) related key-
 393 words. Consequently, projects often contained too few identified flaky com-
 394 mits (Table 1), which has also been observed in prior flakiness studies Ro-
 395 mano et al. (2021); Hashemi et al. (2022). As this could bias the comparison,
 396 we excluded projects with fewer than ten identified flaky commits from the
 397 thresholded analysis to ensure reliability, resulting in 14 out of 31 projects.
 398 Accordingly, Figure 3 presents two groups of boxplots: one that includes all
 399 projects (Figure 5b) and another that applies the threshold (Figure 5a).

400 Overall, visual commits exhibit a higher frequency of flakiness than non-
 401 visual commits. The frequencies shown are normalised per thousand commits
 402 (‰). In both analyses, visual commits are consistently more likely to involve
 403 flaky behaviour. The median flakiness frequency for visual commits is 3.7‰
 404 (2.3‰ under the threshold), compared to 0.8‰ (2.4‰ under the thresh-
 405 old) for non-visual commits. The corresponding mean frequencies are 4.6‰
 406 (3.5‰) and 1.7‰ (3.3‰), respectively. These results indicate that the vi-

407 visual layers of web projects are often more prone to flaky behaviour than
 408 non-visual parts.

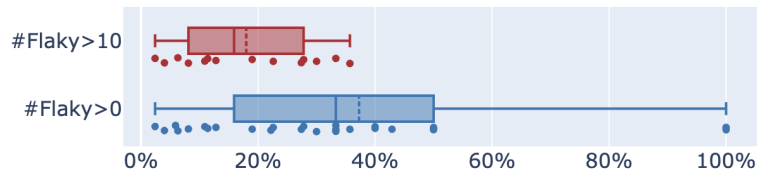


Figure 6: $\#FlakyVisual$ to $\#Flaky$ per project (%). $\#Flaky > 0$ and $\#Flaky > 10$ denote including only the projects with more than 0 (no threshold) and more than 10 flakiness instances. For $\#Flaky > 0$, mean = 37.2%, median = 33.3%; for $\#Flaky > 10$, mean = 17.9% and median=15.9%.

409 Having established that visual commits are both prevalent and more likely
 410 to exhibit flaky behaviour than non-visual ones, we lastly examined the over-
 411 all contribution of visual flakiness to all flaky cases across projects. Figure 6
 412 shows the distribution of the proportion of visual flaky commits among all
 413 flaky commits, with and without thresholding (i.e., including only projects
 414 with at least ten flaky commits). As shown in the lower boxplot of Figure 4,
 415 visual flakiness accounts for approximately 37.2% of all flaky cases, while in
 416 the thresholded subset (upper boxplot), 17.9% of all flaky commits are visu-
 417 ally related across the remaining 14 projects. Combined with the previously
 418 established prevalence and relatively high likelihood of visual commits being
 419 flaky, these observations highlight the need for a deeper investigation of the
 420 underlying characteristics of web visual flakiness.

Answer to RQ1. Visual commits, which account for an average of 11.7% of total commits, are often more prone to flakiness than non-visual ones. With and without the threshold of at least 10 flaky commits per project to be considered, visual flaky commits constitute a notable portion of all flaky cases (37.2% and 17.9%). These results combined underscore the inherently unstable nature of the visual layer of web projects and motivate a detailed examination of its underlying characteristics.

422 3.2. RQ2: Common Types of Web Visual Flakiness

423 Table 2 summarises the distribution of web visual flakiness in our dataset.
 424 Among the 262 analysed cases, 59.9% (157) are structure-related and 40.1%
 425 (105) are style-related. Similar trends appear in both Chromium (47.5%

Table 2: Distributions of Web Visual Flakiness. The *Webs* refer to the open-source web projects. The values within the parentheses denote the percentage of flakiness instances in each type. The last column *Total* shows the combined statistics of *Webs* and *Chromium*.

Type	Subtype of Flaky Behaviour	# Flaky Cases (%)		
		Webs	Chromium	Total
Structure -related	DOM Structure Issues	34 (23.6%)	25 (21.2%)	59 (22.5%)
	DOM Layout Instability	17 (11.8%)	19 (16.1%)	36 (13.7%)
	Timing-sensitive Visual Triggers	50 (34.7%)	12 (10.2%)	62 (23.7%)
	All structured-related	101 (70.1%)	56 (47.5%)	157 (59.9%)
Style -related	Style Interpretation Issues	35 (24.3%)	45 (38.1%)	80 (30.5%)
	Style Resource Latency	8 (5.6%)	17 (14.4%)	25 (9.5%)
	All style-related	43 (29.9%)	62 (52.5%)	105 (40.1%)
All	-	144 (100%)	118 (100%)	262 (100%)

426 vs. 52.5%) and the 31 open-source web projects (70.1% vs. 29.9%). In the
 427 open-source projects, although style-related flakiness forms a smaller portion
 428 than in Chromium, it still accounts for nearly one-third of all visual flaki-
 429 ness, showing its practical importance. Overall, these results indicate that
 430 both structure- and style-related flakiness are major sources of web visual
 431 instability, with neither overwhelmingly dominating the other.

432 The following section presents the detailed distribution of structure- and
 433 style-related flakiness and provides real-world examples illustrating how each
 434 type of web visual flakiness manifests in practice.

435 3.2.1. *Structure-related Flakiness*

436 Within structure-related flakiness (the last column in Table 2), the distri-
 437 bution is largely consistent across Chromium and the Web projects, except for
 438 Timing-sensitive Visual Triggers, which show the most notable variation—
 439 34.7% in the Web set but only 10.2% in Chromium. This gap may reflect
 440 the differences in execution environments. Web applications often rely on
 441 asynchronous updates or client-side scripts that are more directly affected by
 442 subtle variations in timing, whereas the rendering pipeline of Chromium is
 443 more centralised and internally synchronised. In contrast, DOM Structural
 444 Issues and DOM Layout Instability exhibit similar proportions across both
 445 sets, indicating that these types of structural flakiness (i.e., the presence, ap-
 446 pearance, and position of elements) occur at comparable rates in Chromium

447 and general web projects. The following examples illustrate how each type
 448 of structure-related flakiness manifests in practice.

449 **DOM Structure Issues (22.5%)** DOM structure issues refer to cases
 450 where the expected DOM hierarchy or elements are inconsistently constructed
 451 across executions. This can happen due to structural disruptions (e.g., el-
 452 ement insertions or removals) or timing-related factors (e.g., asynchronous
 453 rendering or race conditions) that affect *whether* they are properly built.

```

1 <div class="mat-mdc-snack-bar-label" #label>
2   /* fix: added aria-hidden, aria-live */
3 -   <ng-template cdkPortal></ng-template>
4 +   <div aria-hidden="true">
5 +     <ng-template cdkPortal></ng-template>
454 6 +   </div>
7 +   <div [attr.aria-live]="_live"></div>
8 </div>

```

Listing 1: Displaying and using elements dynamically

455 Listing 1⁷ presents an example where rendering flakiness occurred in dy-
 456 namically inserted content (i.e., a snack bar). The issue arose because the
 457 DOM structure lacked a properly defined `aria-live` region, causing unsta-
 458 ble update behaviour during content rendering. Lines 4-7 add a `div` with
 459 `aria-hidden="true"` (Lines 4-6) and another `div` marked with `aria-live`
 460 (Line 7), ensuring a stable and complete DOM structure that prevents pre-
 461 mature/inconsistent updates during dynamic rendering.

462 **DOM Layout Instability (13.7%)** DOM layout instability concerns flak-
 463 iness affecting *how* existing DOM elements are spatially positioned, rather
 464 than *whether* they are present. Such flakiness manifests as inconsistent spa-
 465 tial arrangements, leading to overlaps, misalignment, or displacement of ele-
 466 ments.

⁷<https://github.com/angular/components/commit/cf8e8eee>

```

1  /* fix: narrow viewport width flaky errors */
2  - <div #table style="margin: 16px">
3  + <div #table style="margin: 16px; max-width:
    90vw; max-height: 90vh;">
4      <mat-table [dataSource]="dataSource">
5          <ng-container matColumnDef="before">
6              </ng-container>
7          </mat-table>
8  </div>

```

467

Listing 2: Container size makes the page structure unstable

468 Listing 2⁸ shows a case of DOM layout instability due to uncontrolled
 469 expansion of a container `<div>`. The absence of maximum width and height
 470 constraints allowed the table to grow unpredictably, causing structural shifts
 471 such as overflow or misalignment. The fix directly modifies the DOM by
 472 adding bounding attributes, which constrain spatial behaviour and stabilise
 473 rendering. Though expressed via style-like properties, the impact is struc-
 474 tural, preventing disordered layouts rather than adjusting visual appearance.

475 **Timing-sensitive Visual Triggers (23.7%)** Timing-sensitive visual trig-
 476 gers refer to flakiness that manifests due to variation in *when* the visual state
 477 or DOM becomes observable or stable, rather than *whether* it appears at all.
 478 Such flakiness arise from timing mismatches between test actions and the
 479 completion of UI rendering or state updates.

⁸<https://github.com/angular/components/commit/d2b499da>

```

1 <input #textInput type="text" [value]="text"
2   (focus)="onAcceptEdit()"
3   (blur)="onAcceptEdit()"
4 /input>
5 export class TextEditorComponent {
6   @ViewChild("textInput")
480   ngAfterViewInit() {
8     /* fix: add delay to ensure the state */
9 +     window.setTimeout(() => {
10      this.textInput.nativeElement.focus();
11 +     }, 5);
12   }}

```

Listing 3: Unexpected element focus behaviour

481 In Listing 3⁹, DOM Structure flakiness occurs due to timing mismatches
482 in when the input field becomes available for interaction. The input field is
483 expected to receive focus when triggered, as indicated by the focus operation
484 ("onAcceptEdit") (Line 2). However, since rendering is not always synchro-
485 nised with when elements become fully available, focus may be applied too
486 early. This can cause the input field to fail to receive focus, preventing im-
487 mediate text input intermittently.

488 3.2.2. *Style-related Flakiness*

489 For style-related flakiness, Style Interpretation Issues are the most preva-
490 lent type in both Chromium (38.1%) and the Web projects (24.3%) (Ta-
491 ble 2). Style Resource Latency appears less frequently, though the magni-
492 tude of difference varies: it accounts for 14.4% of style-related flakiness in
493 Chromium but only 5.6% in the Web projects. This difference likely stems
494 from Chromium’s heavy reliance on diverse and externally hosted style re-
495 sources, which increases the likelihood of delayed or inconsistent loading,
496 potentially leading to more visual instability.

497 **Style Interpretation Issues (30.5%)** This type of flakiness often arises
498 when CSS properties or style interactions are handled inconsistently across
499 rendering environments. Even minor differences in CSS or interactions be-
500 tween style rules—such as differences in the evaluation of properties like
501 display or text-indent between platforms or browsers—can result in visually

⁹<https://github.com/angular/components/commit/4ae63b39>

502 different outcomes across executions. Such inconsistencies make it challeng-
 503 ing for developers to ensure a stable and uniform visual presentation in web
 504 applications.

```

1  <style>
2    input::-webkit-input-placeholder, textarea
3    ::-webkit-input-placeholder {
4      color: #aaaaaa;
5      text-indent: 0;}
505  /* fix: add rule improve compatibility */
6  +  textarea::-webkit-input-placeholder {
7  +    color: #aaaaaa;
8  +    text-indent: -3px;}
9  </style>

```

Listing 4: *text-indent* properties cause flakiness behaviour

506 Listing 4¹⁰ (i.e., style codes) presents an example of style-related flaki-
 507 ness caused by CSS attribute compatibility issues. The placeholder text and
 508 cursor become misaligned when the text area is focused due to variations
 509 in how browsers interpret the *text-indent* property for placeholders. This
 510 inconsistency arises from differences in CSS rule enforcement across render-
 511 ing engines. To improve compatibility, the fix explicitly defines *text-indent*
 512 for `textarea::-webkit-input-placeholder`, ensuring more consistent be-
 513 haviour across browsers.

514 **Style Resource Latency (9.5%)** Style Resource Latency occurs when es-
 515 sential assets — such as fonts, images, or external style sheets — are not
 516 fully loaded before the visual rendering is complete. Such delays can lead to
 517 temporary or inconsistent appearance, as the web page may be incorrectly
 518 styled or displayed without its intended formatting.

¹⁰<https://github.com/ionic-team/ionic-framework/commit/76ca4757>

```

1   main_resource.Write(HTML(
2     <!doctype html>
3     <head>
4       /* fix: cache aware font Loading */
5     <link rel="preload" as="font" type="font/
519   woff2"
6       href="https://example.com/font.woff
      ">))

```

Listing 5: Font loading contributes to flakiness

Listing 5¹¹ illustrates a case where flakiness arises due to delayed font loading (Line 6). If custom fonts take too long to load or fail to load, browsers fall back to default fonts, which may differ in size or spacing, causing misaligned or invisible elements. This mismatch between the expected and the actual leads to visual inconsistencies. Developers fixed this issue using preload hints to ensure timely font availability and prevent such flakiness.

Overall, Chromium shows a more even distribution of structure- and style-related flakiness compared to the Web projects. We suspect that the key contributor is the smaller proportion of Timing-sensitive Visual Triggers flakiness, possibly reflecting Chromium's more synchronised event scheduling. At the same time, its complex, multi-threaded rendering pipeline may introduce additional nondeterminism during layout computation and resource loading, contributing to the slightly higher proportions of DOM Layout Instability and Style Resource Latency observed. DOM construction remains largely deterministic within a single engine build. Hence, the modest increase in Style Interpretation Issues may instead reflect the broader redistribution of flaky cases from the decrease in timing-sensitive flakiness, rather than inherent nondeterminism in the style engine itself.

538 Flaky behaviours across different front-end frameworks.

539 Modern web development increasingly relies on front-end frameworks. In
540 this empirical study, we additionally examined whether the choice of frame-
541 work influences the distribution of visual flakiness types. Among 31 web
542 projects, 87 of 144 cases (60.6%) involved frameworks, such as React, Vue,
543 and Angular¹². As shown in Figure 7, similar patterns emerge across React,

¹¹<https://issues.chromium.org/issues/40114104>

¹²The remaining 39.4% of projects do not use front-end frameworks; some were developed before frameworks became widely adopted, while others are lightweight applications

544 Vue, and Angular, suggesting that the underlying framework does not have
 545 much influence on the manifestation of web visual flakiness.

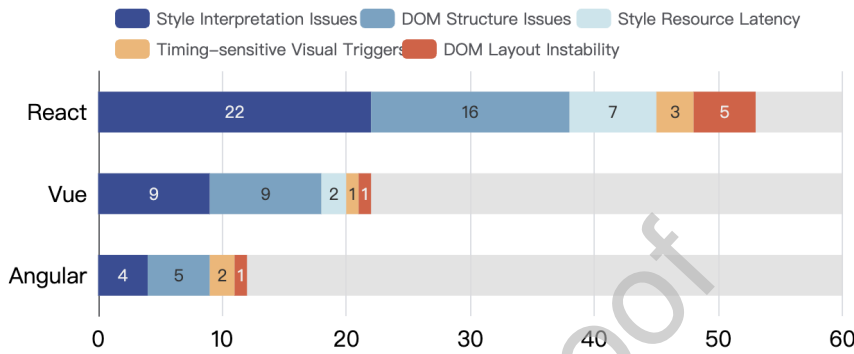


Figure 7: The distribution of flakiness across UI frameworks.

Answer to RQ2. The analysis shows that structure-related and style-related types together account for nearly equal portions of web visual flakiness. Among the structure-related types, Timing-sensitive visual triggers are the most frequent, whereas Style Interpretation Issues are the most common among style-related ones, followed by Style Resource Latency. Overall, both Chromium and the Web projects exhibit similar distributions of visual flakiness types. However, Chromium exhibits a lower proportion of Timing-sensitive visual triggers and a comparatively higher share of style-related flakiness. The smaller ratio of timing-sensitive triggers is likely attributed to more synchronised and centralised rendering and event pipeline of Chromium, reducing timing variability across executions. In contrast, the relatively higher ratio of style-related flakiness might partly result from fewer structure-related cases overall, rather than an inherent increase in style-level instability.

546

547 3.3. RQ3: Common Flakiness Fix Patterns of Developers

548 Based on the categorisation process described in Section 2.2.4, we identi-
 549 fied seven categories of recurring fix strategies developers use to address web
 550 visual flakiness.

or static websites that avoid the extra complexity that frameworks introduce.

- 551 • **Manage element changes:** insert, remove, or update DOM elements
552 to prevent inconsistencies in element presence or structure.
- 553 • **Delay:** introduce a conservative timing buffer to handle variability
554 in dynamic rendering or state transitions, ensuring that subsequent
555 operations are not processed prematurely.
- 556 • **Explicit size position:** define layout dimensions (e.g., width, height,
557 margins) to avoid instability caused by flexible or implicit sizing.
- 558 • **Use grid/flexbox:** employ structured layout systems such as CSS
559 Grid or Flexbox to improve spatial consistency and reduce layout mis-
560 alignments.
- 561 • **CSS compatibility:** modify CSS rules to ensure consistent behaviour
562 across different browsers and rendering engines by applying baseline
563 styles that override browser-specific defaults to standardise the initial
564 rendering behaviour across environments
- 565 • **Element-event coordination:** align target elements and triggering
566 events in timing and state to prevent missed or inconsistent interac-
567 tions.
- 568 • **Preload & Cache:** ensure key resources (e.g., fonts, images) are avail-
569 able when needed by preloading them and leveraging caching to prevent
570 delays or rendering failures.

571 We further analysed how these fix strategies are distributed across differ-
572 ent types of web visual flakiness to understand which strategies developers
573 often employ to handle structure- and style-related flakiness issues, as shown
574 in Table 3. Overall, developers tend to employ fix strategies that closely cor-
575 respond to how each type of web visual flakiness manifests. For instance, in
576 structure-related flakiness, such as DOM Structure issues, developers most
577 frequently applied Element–event coordination, used in 23 of 34 cases, to
578 stabilise interactions between elements and events. In Timing-sensitive Vi-
579 sual Triggers, Delay and Element–Event Coordination were the two dominant
580 strategies, both addressing asynchronous rendering and event ordering issues.
581 For DOM Layout Instability, developers often constrain layout flexibility to
582 prevent unintended spatial shifts by adopting the Explicit size position strat-
583 egy. For the Style Interpretation, nearly all fixes (33 of 35) were about CSS
584 Compatibility, addressing inconsistent rendering semantics through CSS ad-
585 justments. Lastly, for Style Resource Latency, timing- and loading-related
586 strategies, such as Delay and Preload & Cache, were most frequently applied

587 to ensure stable resource availability and rendering.

588 At the same time, we observed that specific fix strategies are not strictly
 589 confined to a single type of web visual flakiness. For instance, CSS Compati-
 590 bility was applied across all five categories of flakiness, and Delay was utilised
 591 in DOM Layout Instability, Timing-sensitive Visual Triggers, Style Interpre-
 592 tation Issues, and Style Resource Latency. Element-event coordination, the
 593 most adopted strategy in DOM Structure Issues, was also the most common
 594 strategy among all fixes for structure-related flakiness, appearing in 38 of 144
 595 cases (26.4%). We conjecture that this broad applicability stems from the
 596 general role these strategies play in stabilising cross-cutting sources of nonde-
 597 terminism. For instance, modifying CSS rules for CSS Compatibility could
 598 mitigate rendering inconsistencies that manifest across structure and style
 599 components. Similarly, the Delay strategy reduces premature or mismatched
 600 interactions by introducing timing buffers, while Element–event coordination
 601 can ensure the synchronisation of elements and their triggering events, thus
 602 generally preventing state mismatches in timing-sensitive scenarios.

603 **Developer priorities for different types.**

604 We further examined how developers prioritise addressing different types
 605 of web visual flakiness using the Chromium dataset, where each issue is as-
 606 signed an explicit priority level. Figure 8 shows the distribution of these
 607 priorities, P1 (high), P2 (medium), and P3 (low), across structure-related
 608 and style-related categories of web visual flakiness. Specifically, the first
 609 and third subfigures in Figure 5 illustrate priority distributions within style-
 610 related and structure-related flakiness, respectively. Overall, while a substan-
 611 tial proportion of both structure-related (44.6%) and style-related (59.7%)
 612 issues were assigned medium priority (P2), the structure-related issues were
 613 more frequently labelled as low priority (P3) than the style-related (41.1%
 614 vs. 24.2%). This implies that developers may allocate relatively more at-
 615 tention to style-related flakiness in Chromium. A likely explanation is that
 616 style-related flakiness often produces visually broader or more noticeable dis-
 617 ruptions, such as broken layouts or inconsistent rendering, compared to the
 618 structure-related ones that may remain localised to specific elements. In fact,
 619 our previous analysis on the distribution of web visual flakiness shows that
 620 style-related flakiness represents a larger overall share than in the open-source
 621 web projects.

622 At a finer granularity, DOM Layout Instability exhibited the highest pro-
 623 portion of high-priority (P1) assignments (21.2%), followed by Style Resource

Table 3: Repair strategies for 144 cases of web projects’ visual flakiness. Overall, developers typically apply fix strategies aligned with specific visual flakiness types (e.g., *Element–Event Coordination for DOM Structure Issues*, and *Delays for Timing-sensitive Visual Triggers*), while some strategies—such as *CSS compatibility*—are used across both structure- and style-related flakiness.

Type	Cause	Fix method	# Flaky Cases (%) Webs
Structure -related	DOM Structure Issues	CSS compatibility	3 (2.1%)
		Manage element changes	8 (5.6%)
		Element-event coordination	23 (16.0%)
	DOM Layout Instability	CSS compatibility	2 (1.4%)
		Delay	2 (1.4%)
		Manage element changes	1 (0.7%)
		Element-event coordination	1 (0.7%)
		Explicit size position	5 (3.5%)
		Use grid/flexbox	6 (4.2%)
	Timing-sensitive Visual Triggers	CSS compatibility	3 (2.1%)
		Delay	26 (18.1%)
		Element–event coordination	14 (9.7%)
Manage element changes		7 (4.9%)	
Style -related	Style Interpretation Issues	CSS compatibility	33 (22.9%)
		Delay	1 (0.7%)
		Manage element changes	1 (0.7%)
	Style Resource Latency	CSS compatibility	1 (0.7%)
		Delay	2 (1.4%)
		Preload & Cache	4 (2.8%)
		Manage element changes	1 (0.7%)

624 Latency (17.6%) and Style Interpretation Issues (15.6%). These results fur-
625 ther support the above conjectures, suggesting that developers often priori-
626 tise unstable behaviours with broad visual impact, especially those that affect
627 page layout or rendering consistency. Timing-Sensitive Visual Triggers were
628 mostly given medium priority (83.3%), indicating a medium level of urgency.
629 DOM Structure Issues had the highest proportion of low-priority (P3) as-
630 signments (56%), suggesting that element-specific inconsistencies are often

631 viewed as less critical to the overall user experience. This further explains
 632 the higher proportion of low priority in structure-related flakiness compared
 633 to style-related flakiness.

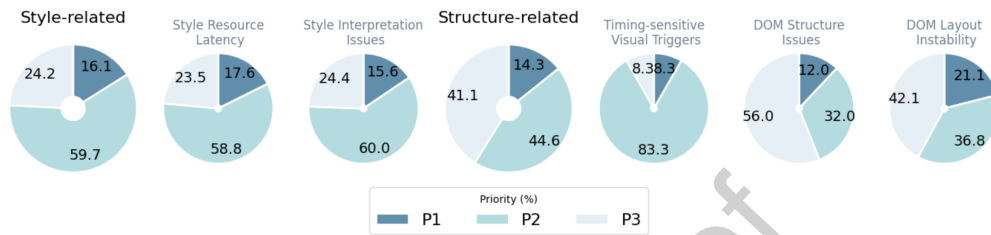


Figure 8: The distribution of developers addressing priorities (P1 (high), P2 (medium), and P3 (low)) under different causes, the proportions are in percentages (%).

Answer to RQ3. While developers tend to choose fix strategies that align closely with the specific types of web visual flakiness, some fix strategies appear across multiple types, reflecting their broad applicability. Additional analysis of issue priorities shows that developers tend to prioritise flakiness with broader visual impact, such as layout or style-related problems, over those with narrower or isolated effects.

634

635 4. Discussion

636 4.1. Positioning web visual flakiness against general UI issues and previous 637 Web flakiness studies

638 Unlike general web visual issues discussed in prior studies Althomali et al.
 639 (2021); Alameer and Halfond (2016); Mahajan et al. (2016); Romano et al.
 640 (2021), which are typically deterministic or consistently reproducible, web
 641 visual flakiness is inherently intermittent and non-deterministic. It often
 642 stems from asynchronous rendering, transient DOM states, or timing in-
 643 consistencies, and manifests across the structure and style dimensions that
 644 we categorised in this study. Such characteristics of these issues complicate
 645 their detection and resolution, as they depend on subtle runtime interactions
 646 rather than static or deterministic failures. Prior studies on Web flakiness
 647 have primarily examined broad or functional causes (e.g., network latency,
 648 environment configuration) or deterministic rendering faults Romano et al.

649 (2021). In contrast, our work isolates and empirically characterises flaky be-
650 haviours that specifically originate in the visual layer, i.e., those directly tied
651 to the presentation logic of web structures and styles, thereby addressing a
652 class of issues that has often been overlooked in prior work.

653 **5. Threats to Validity**

654 *5.1. Internal validity*

655 While our categorisations of flaky behaviours were performed through it-
656 erative discussion and cross-validation among the two authors, along with
657 LLM assistance, the process is inherently affected by subjective judgment.
658 Although we ensured consistency through shared criteria and multiple passes
659 over the data, potential bias or misclassification may remain. Flakiness
660 in web visuals can stem from subtle timing issues, rendering inconsisten-
661 cies, style misinterpretations, etc. While we applied multiple-round manual
662 checks, some edge cases may still have been misplaced. In addition, our cat-
663 egorisation may not capture all manifestations in diverse contexts. Future
664 work could refine and expand these criteria for a more precise.

665 *5.2. Construct validity*

666 Construct validity threats primarily stem from the use of keyword-based
667 filtering of commit messages, issue reports, and pull requests to identify vi-
668 sual commits and web visual flakiness. Keyword matching may introduce
669 misclassification, including false negatives when visual or flaky behaviour is
670 only implicitly described, as well as false positives when keywords do not
671 correspond to actual visual or flaky behaviour.

672 To mitigate this, we adopted a conservative sampling strategy focused on
673 precision. Specifically, we focused only on the interaction of flaky and visual
674 keywords, manually validating each candidate case. Two authors jointly in-
675 spected code changes and artifacts to confirm the presence of visual flakiness,
676 and thereby further eliminated false positives, increasing confidence in the
677 final dataset.

678 As this conservative choice likely results in false negatives, the reported
679 prevalence should be interpreted as a lower bound. However, since our cate-
680 gorisation and fix-pattern analyses focus on recurring behavioural character-
681 istics observed across consistently validated cases, rather than on complete
682 enumeration, this limitation is less likely to affect the qualitative findings.

683 5.3. *External validity*

684 A potential limitation of our study is the generalisability of the results
685 beyond the specific projects and test suites analysed. Although our find-
686 ings are rooted in web projects and Chromium, they may also be relevant to
687 other contexts, such as mobile applications and hybrid frameworks. Mobile
688 browsers, which exhibit a variety in hardware setups and screen sizes, are
689 also prone to similar structural and style inconsistencies, especially when it
690 comes to responsive development. Despite differences in technical implemen-
691 tation, manifestations of flakiness, like layout shifts and inconsistent style
692 application, remain conceptually aligned. This suggests that our categorisa-
693 tions and observed fix patterns may help inform broader efforts to address
694 visual flakiness across diverse platforms.

695 6. Related work

696 Flaky tests have been the subject of extensive research in the software
697 testing community, particularly in programming languages, such as Java,
698 Python, and JavaScript Luo et al. (2014); Lam et al. (2020a,b); Gruber et al.
699 (2021); Hashemi et al. (2022); Leong et al. (2019). Numerous studies have
700 examined the prevalence, causes, and mitigation strategies for flaky tests
701 in code-based systems Lam et al. (2020a); Pinto et al. (2020); Dutta et al.
702 (2020). For example, Luo et al. conducted a pioneering study on flaky
703 tests, analysing large-scale systems to identify the underlying causes of test
704 flakiness Luo et al. (2014). While these findings are valuable for traditional
705 software, they do not directly address the challenges posed by style-related
706 issues, where visual rendering discrepancies could also lead to flakiness.

707 For research on flaky tests in web environments, Hashemi et al. analysed
708 flaky behaviour in JavaScript frameworks, identifying browser-related issues
709 as significant contributors to flakiness testing Hashemi et al. (2022). Romano
710 et al. analysed flaky UI tests from 62 web and Android projects and identi-
711 fied four categories of causes for flaky UI tests Romano et al. (2021). Moran
712 et al. focused on identifying the root cause of flakiness in web applications
713 by analysing the execution of tests under diverse combinations of environ-
714 mental factors that may contribute to flakiness Morán et al. (2020). Pei et
715 al. investigated the async await flakiness in web tests, focusing on the DOM
716 and time issues Pei et al. (2024). However, these studies rarely isolate style-
717 related code as a distinct source of flaky behaviour, although CSS is highly
718 susceptible to the inconsistencies that lead to flakiness. Our study aims to

719 fill this gap by specifically investigating flaky behaviours in CSS, providing
720 a more nuanced understanding of the problem in web development.

721 Besides empirical studies on flaky tests, there has been considerable work
722 on mitigation strategies Ziftci and Cavalcanti (2020); Camara et al. (2021);
723 Haben et al. (2021); Rahman and Shi (2024); Fallahzadeh and Rigby (2022).
724 Bell et al. leveraged code evolution and code coverage to determine whether
725 new test failures between two commits are due to flaky tests Bell et al.
726 (2018). There have also been several tools that have been targeted to detect
727 or repair certain types of flaky tests, such as iDFlakies Lam et al. (2019b),
728 iFixFlakies Shi et al. (2019), FLASH Dutta et al. (2020), Flakify Fatima
729 et al. (2022), Flakeflagger Alshammari et al. (2021), and so on Dutta et al.
730 (2021); Cordy et al. (2022); Cai et al. (2024).

731 While there is a growing of research on flaky tests in software systems,
732 the focus on web visual-related flakiness is minimal. Our research builds on
733 previous studies of flaky tests by extending the analysis to the web visuals.
734 By focusing specifically on flaky behaviour in style and visual layers, this
735 study aims to bridge the gap in current literature and provide actionable
736 insights for developers dealing with flaky behaviours.

737 7. Conclusion

738 This paper presents a focused empirical study of web visual flakiness
739 as a distinct and explicitly characterised class of flakiness, based on 262
740 cases collected from 31 open-source projects (144 cases) and the Chromium
741 project (118 cases), revealing its recurring presence across diverse mod-
742 ern web projects. We found that visual commits, comprising 11.4% of all
743 commits, exhibit a higher frequency of flakiness than non-visual commits.
744 We categorised web visual flakiness into two primary dimensions—structure-
745 related and style-related—, which we further refined into three and two con-
746 crete subtypes, respectively. Through an inductive analysis of developers’
747 fixes, we found seven recurring fix strategies, which address common insta-
748 bility patterns, such as element-event coordination, timing, layout stabili-
749 sation, and style consistency. Although some are closely associated with
750 specific flakiness types, others span across multiple types, showing the inter-
751 twined nature of structure and style aspects in rendering. Together, these
752 findings highlight the need to treat web visual flakiness as a distinct subject
753 in web development and testing. To facilitate further research and repro-
754 ducibility, we make our dataset and replication package publicly available at:

755 <https://doi.org/10.5281/zenodo.17324581>

756 8. Acknowledgments

757 This work is supported by the Luxembourg National Research Funds
758 (FNR) through the CORE project grant C20/IS/14761415/TestFlakes and
759 partly supported by the National Research Foundation of Korea (NRF) grant
760 funded by the Korea government (MSIT) (No.2021R1A5A1021944).

761 References

- 762 Alameer, A., Halfond, W. G., 2016. An empirical study of international-
763 ization failures in the web. In: 2016 IEEE International Conference on
764 Software Maintenance and Evolution (ICSME). IEEE, pp. 88–98.
- 765 Alshammari, A., Morris, C., Hilton, M., Bell, J., 2021. Flakeflagger: Predict-
766 ing flakiness without rerunning tests. In: 2021 IEEE/ACM 43rd Interna-
767 tional Conference on Software Engineering (ICSE). IEEE, pp. 1572–1584.
- 768 Althomali, I., Kapfhammer, G. M., McMinn, P., 2021. Automated visual
769 classification of dom-based presentation failure reports for responsive web
770 pages. *Software Testing, Verification and Reliability* 31 (4), e1756.
- 771 Bell, J., Legunsen, O., Hilton, M., Eloussi, L., Yung, T., Marinov, D., 2018.
772 Deflaker: Automatically detecting flaky tests. In: Proceedings of the 40th
773 international conference on software engineering (ICSE). pp. 433–444.
- 774 Cai, X., Dong, Z., Wang, Y., Tiwari, A., Peng, X., 2024. Reproducing timing-
775 dependent gui flaky tests in android apps via a single event delay. In: Pro-
776 ceedings of the 33rd ACM SIGSOFT International Symposium on Software
777 Testing and Analysis (ISSTA). pp. 1504–1515.
- 778 Camara, B., Silva, M., Endo, A., Vergilio, S., 2021. On the use of test smells
779 for prediction of flaky tests. In: Proceedings of the 6th Brazilian Sympo-
780 sium on Systematic and Automated Software Testing (SAST). pp. 46–54.
- 781 Choudhary, S. R., Versee, H., Orso, A., 2010. Webdiff: Automated identifi-
782 cation of cross-browser issues in web applications. In: 2010 IEEE Interna-
783 tional Conference on Software Maintenance. IEEE, pp. 1–10.

- 784 Choudhary, S. R., Zhao, D., Versee, H., Orso, A., 2011. Water: Web appli-
785 cation test repair. In: Proceedings of the First International Workshop on
786 End-to-End Test Script Engineering. pp. 24–29.
- 787 Cordy, M., Rwemalika, R., Franci, A., Papadakis, M., Harman, M., 2022.
788 Flakime: Laboratory-controlled test flakiness impact assessment. In:
789 2021 IEEE/ACM 44rd International Conference on Software Engineering
790 (ICSE). IEEE.
- 791 Dutta, S., Shi, A., Choudhary, R., Zhang, Z., Jain, A., Misailovic, S., 2020.
792 Detecting flaky tests in probabilistic and machine learning applications.
793 In: Proceedings of the 29th ACM SIGSOFT international symposium on
794 software testing and analysis (ISSTA). pp. 211–224.
- 795 Dutta, S., Shi, A., Misailovic, S., 2021. Flex: fixing flaky tests in machine
796 learning projects by updating assertion bounds. In: Proceedings of the
797 29th ACM Joint Meeting on European Software Engineering Conference
798 and Symposium on the Foundations of Software Engineering (ESEC/FSE).
799 pp. 603–614.
- 800 Eck, M., Palomba, F., Castelluccio, M., Bacchelli, A., 2019. Understanding
801 flaky tests: The developer’s perspective. In: Proceedings of the 2019 27th
802 ACM Joint Meeting on European Software Engineering Conference and
803 Symposium on the Foundations of Software Engineering (ESEC/FSE). pp.
804 830–840.
- 805 Fallahzadeh, E., Rigby, P. C., 2022. The impact of flaky tests on historical
806 test prioritization on chrome. In: Proceedings of the 44th International
807 Conference on Software Engineering: Software Engineering in Practice
808 (ICSE-SEIP). pp. 273–282.
- 809 Fatima, S., Ghaleb, T. A., Briand, L., 2022. Flakify: A black-box, language
810 model-based predictor for flaky tests. IEEE Transactions on Software En-
811 gineering (TSE) 49 (4), 1912–1927.
- 812 Gruber, M., Fraser, G., 2022. A survey on how test flakiness affects developers
813 and what support they need to address it. In: 2022 IEEE Conference on
814 Software Testing, Verification and Validation (ICST). IEEE, pp. 82–92.

- 815 Gruber, M., Lukasczyk, S., Kroiß, F., Fraser, G., 2021. An empirical study of
816 flaky tests in python. In: 2021 14th IEEE Conference on Software Testing,
817 Verification and Validation (ICST). IEEE, pp. 148–158.
- 818 Habchi, S., Haben, G., Papadakis, M., Cordy, M., Le Traon, Y., 2022. A
819 qualitative study on the sources, impacts, and mitigation strategies of
820 flaky tests. In: 2022 IEEE Conference on Software Testing, Verification
821 and Validation (ICST). IEEE, pp. 244–255.
- 822 Haben, G., Habchi, S., Micco, J., Harman, M., Papadakis, M., Cordy, M.,
823 Le Traon, Y., 2024. The importance of accounting for execution failures
824 when predicting test flakiness. In: Proceedings of the 39th IEEE/ACM
825 International Conference on Automated Software Engineering (ASE). pp.
826 1979–1989.
- 827 Haben, G., Habchi, S., Papadakis, M., Cordy, M., Le Traon, Y., 2021. A
828 replication study on the usability of code vocabulary in predicting flaky
829 tests. In: 2021 IEEE/ACM 18th International Conference on Mining Soft-
830 ware Repositories (MSR). IEEE, pp. 219–229.
- 831 Hashemi, N., Tahir, A., Rasheed, S., 2022. An empirical study of flaky tests
832 in javascript. In: 2022 IEEE International Conference on Software Main-
833 tenance and Evolution (ICSME). IEEE, pp. 24–34.
- 834 Lam, W., Godefroid, P., Nath, S., Santhiar, A., Thummalapenta, S., 2019a.
835 Root Causing Flaky Tests in a Large-Scale Industrial Setting. In: Pro-
836 ceedings of the 28th ACM SIGSOFT International Symposium on Software
837 Testing and Analysis (ISSTA). pp. 101–111.
- 838 Lam, W., Muşlu, K., Sajnani, H., Thummalapenta, S., 2020a. A study on
839 the lifecycle of flaky tests. In: Proceedings of the ACM/IEEE 42nd Inter-
840 national Conference on Software Engineering (ICSE). pp. 1471–1482.
- 841 Lam, W., Oei, R., Shi, A., Marinov, D., Xie, T., 2019b. IDFlakies: A frame-
842 work for detecting and partially classifying flaky tests. Proceedings - 2019
843 IEEE 12th International Conference on Software Testing, Verification and
844 Validation (ICST), 312–322.
- 845 Lam, W., Winter, S., Astorga, A., Stodden, V., Marinov, D., 2020b. Un-
846 derstanding reproducibility and characteristics of flaky tests through test

- 847 reruns in java projects. In: 2020 IEEE 31st International Symposium on
848 Software Reliability Engineering (ISSRE). IEEE, pp. 403–413.
- 849 Leong, C., Singh, A., Papadakis, M., Le Traon, Y., Micco, J., 2019. Assess-
850 ing transition-based test selection algorithms at google. In: Proceedings
851 of the 41st International Conference on Software Engineering: Software
852 Engineering in Practice (ICSE-SEIP). IEEE, pp. 101–110.
- 853 Luo, Q., Hariri, F., Eloussi, L., Marinov, D., 2014. An empirical analysis
854 of flaky tests. In: Proceedings of the 22nd ACM SIGSOFT international
855 symposium on foundations of software engineering (FSE). pp. 643–653.
- 856 Mahajan, S., Gadde, K. B., Pasala, A., Halfond, W. G., 2016. Detecting and
857 localizing visual inconsistencies in web applications. In: 2016 23rd Asia-
858 Pacific Software Engineering Conference (APSEC). IEEE, pp. 361–364.
- 859 Morán, J., Augusto, C., Bertolino, A., De La Riva, C., Tuya, J., 2020. Flaky-
860 loc: flakiness localization for reliable test suites in web applications. *Journal of Web Engineering* 19 (2), 267–296.
- 862 Nass, M., Alégroth, E., Feldt, R., Coppola, R., 2023. Robust web element
863 identification for evolving applications by considering visual overlaps. In:
864 IEEE Conference on Software Testing, Verification and Validation, ICST
865 2023, Dublin, Ireland, April 16-20, 2023. IEEE, pp. 258–268.
866 URL <https://doi.org/10.1109/ICST57152.2023.00032>
- 867 Parry, O., Kapfhammer, G. M., Hilton, M., McMinn, P., 2021. A survey of
868 flaky tests. *ACM Transactions on Software Engineering and Methodology*
869 (TOSEM) 31 (1), 1–74.
- 870 Parry, O., Kapfhammer, G. M., Hilton, M., McMinn, P., 2022. Surveying
871 the developer experience of flaky tests. In: Proceedings of the 44th In-
872 ternational Conference on Software Engineering: Software Engineering in
873 Practice (ICSE-SEIP). pp. 253–262.
- 874 Pei, Y., Sohn, J., Habchi, S., Papadakis, M., 2024. Non-flaky and nearly-
875 optimal time-based treatment of asynchronous wait web tests. *ACM Trans-*
876 *actions on Software Engineering and Methodology (TOSEM)*.
- 877 Pei, Y., Sohn, J., Papadakis, M., 2025. An empirical study of web flaky tests:
878 Understanding and unveiling dom event interaction challenges. In: 2025

- 879 IEEE Conference on Software Testing, Verification and Validation (ICST).
880 IEEE, pp. 92–102.
- 881 Pinto, G., Miranda, B., Dissanayake, S., d’Amorim, M., Treude, C.,
882 Bertolino, A., 2020. What is the vocabulary of flaky tests? In: Proceed-
883 ings of the 17th International Conference on Mining Software Repositories
884 (MSR). pp. 492–502.
- 885 Rahman, S., Shi, A., 2024. Flakesync: Automatically repairing async flaky
886 tests. In: Proceedings of the IEEE/ACM 46th International Conference on
887 Software Engineering (ICSE). pp. 1–12.
- 888 Romano, A., Song, Z., Grandhi, S., Yang, W., Wang, W., 2021. An empirical
889 analysis of ui-based flaky tests. In: 2021 IEEE/ACM 43rd International
890 Conference on Software Engineering (ICSE). IEEE, pp. 1585–1597.
- 891 Rwemalika, R., Kintis, M., Papadakis, M., Traon, Y. L., Lorrach, P., 2019.
892 On the evolution of keyword-driven test suites. In: 12th IEEE Conference
893 on Software Testing, Validation and Verification, ICST 2019, Xi’an, China,
894 April 22-27, 2019. IEEE, pp. 335–345.
895 URL <https://doi.org/10.1109/ICST.2019.00040>
- 896 Shi, A., Lam, W., Oei, R., Xie, T., Marinov, D., 2019. ifixflakies: A frame-
897 work for automatically fixing order-dependent flaky tests. In: Proceedings
898 of the 2019 27th ACM Joint Meeting on European Software Engineering
899 Conference and Symposium on the Foundations of Software Engineering
900 (ESEC/FSE). pp. 545–555.
- 901 Stocco, A., Yandrapally, R., Mesbah, A., 2018. Visual web test repair. In:
902 Proceedings of the 2018 26th ACM Joint Meeting on European Software
903 Engineering Conference and Symposium on the Foundations of Software
904 Engineering. pp. 503–514.
- 905 Yeh, T., Chang, T.-H., Miller, R. C., 2009. Sikuli: using gui screenshots for
906 search and automation. In: Proceedings of the 22nd annual ACM symposi-
907 um on User interface software and technology. pp. 183–192.
- 908 Ziftci, C., Cavalcanti, D., 2020. De-flake your tests: Automatically locating
909 root causes of flaky tests in code at google. In: 2020 IEEE International
910 Conference on Software Maintenance and Evolution (ICSME). IEEE, pp.
911 736–745.

912 Zou, Y., Chen, Z., Zheng, Y., Zhang, X., Gao, Z., 2014. Virtual DOM cover-
913 age for effective testing of dynamic web applications. In: Pasareanu, C. S.,
914 Marinov, D. (Eds.), International Symposium on Software Testing and
915 Analysis, ISSTA '14, San Jose, CA, USA - July 21 - 26, 2014. ACM, pp.
916 60–70.
917 URL <https://doi.org/10.1145/2610384.2610399>

918 **Declaration of Interest Statement**

919 The authors declare that they have no known competing financial inter-
920 ests or personal relationships that could have appeared to influence the work
921 reported in this paper.