
MobileDev-Bench: A Benchmark for Issue Resolution in Mobile Application Development

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) have shown strong performance on automated soft-
2 ware engineering tasks, yet existing benchmarks focus primarily on library-style
3 repositories, leaving mobile application development largely unexplored despite its
4 framework-specific build systems, heterogeneous artifact types, and coordinated
5 multi-file fix requirements. We introduce MobileDev-Bench, a benchmark com-
6 prising 415 real-world issue-resolution tasks collected from 19 production mobile
7 applications spanning Android Native (Java/Kotlin), React Native (TypeScript),
8 and Flutter (Dart). Each task pairs a verified developer-reported issue with exe-
9 cutable test patches, enabling fully automated validation of model-generated fixes
10 within mobile build environments. The benchmark exhibits substantially greater
11 patch complexity than prior benchmarks: fixes modify 13.5 files and 344.2 lines
12 on average, and 42.2% of instances require coordinated changes across multiple
13 artifact types, such as source, build configuration, and resource files. Evaluation of
14 four frontier LLMs (Claude Sonnet 4.5, Qwen3-Coder, GPT-5.2, and Gemini 2.5
15 Flash) yields end-to-end resolution rates of only 3.90%–4.28% under automated
16 fault localization and at most 6.07% under oracle localization, well below resolu-
17 tion rates reported on existing benchmarks. We release MobileDev-Bench with
18 task instances, an evaluation harness, and containerized environments to support
19 reproducible research on AI-assisted mobile application development.

20 1 Introduction

21 Large language models (LLMs) such as GPT-4 [OpenAI, 2023], Claude [Anthropic, 2023], and Code
22 Llama [Roziere et al., 2023] have demonstrated strong progress in automated software engineering,
23 particularly in repository-level issue resolution. Benchmarks such as SWE-bench [Jimenez et al.,
24 2024] establish a realistic evaluation paradigm by requiring models to resolve real-world GitHub
25 issues through patch generation, where correctness is validated by executing the relevant test suite.
26 Subsequent extensions [Yang et al., 2024b, Zhang et al., 2024a, Kabir, 2025, Mhatre et al., 2025,
27 Zan et al., 2024, 2025] broaden coverage across programming languages and evaluation settings,
28 revealing persistent performance gaps even for frontier code-capable LLMs.

29 Despite these advances, existing benchmarks primarily target issue resolution in general-purpose
30 libraries or web applications. However, mobile application (app) development remains largely
31 underexplored in benchmark design, despite introducing distinct engineering constraints. Mobile
32 apps are built within platform-defined architectures with lifecycle-managed execution [Android,
33 2020]. They frequently employ declarative and event-driven programming [Flutter, 2026] and rely
34 on asynchronous state management across user interface and service components. Compilation and
35 testing are mediated by framework-specific build systems [Android, 2025], and correctness often
36 depends on coordinated changes across heterogeneous artifacts, including source files, manifests,
37 resource definitions, and configuration metadata [Android, 2024, Flutter, 2025]. These characteristics

Benchmarks	Domain	#Tasks	#Repos	Verified	Difficulty level	Multi-type artifact fixes	Patch # Files	Complex. # Lines
SWE-bench	Python libraries	2294	12	✗	✗	✗	1.7	32.8
SWE-bench Verified	Python libraries	500	12	✓	✗	✗	1.2	14.3
SWE-bench Multimodal	JavaScript libraries	617	17	✓	✗	✓	3.9	131.6
SWE-PolyBench	Multi-lang. libraries	2110	21	✗	✗	✗	2.6	51.2
SWE-bench Multilingual	Multi-lang. libraries	300	42	✗	✗	✗	1.7	47.8
OmniGIRL	Multi-lang. repos	959	15	✗	✗	✗	1.2	46.3
SWE-bench-java	Java libraries	91	6	✓	✗	✗	2.3	38.7
SWE-Sharp-bench	C# libraries	150	17	✓	✗	✗	4.9	131.2
Rust-SWE-bench	Rust repositories	500	34	✓	✗	✗	9.8	139.9
Multi-SWE-bench	Multi-lang. repos	1632	39	✓	✓	✗	4.9	163.3
MobileDev-Bench	Mobile Apps	415	19	✓	✓	✓	13.5	344.2

Table 1: Comparison of issue-resolution benchmarks. Unlike prior work that primarily evaluates source-code fixes in library-style repositories, MobileDev-Bench targets mobile apps, incorporates human verification and difficulty stratification, evaluates multi-artifact fixes (e.g., source and resources artifacts such as manifest files), and restores pre-fix buildability to enable execution-based evaluation.

38 make mobile app repositories structurally different from the library-centric settings that dominate
39 existing issue-resolution benchmarks.

40 We introduce *MobileDev-Bench*, a benchmark of 415 manually verified issue-resolution tasks col-
41 lected from 19 production mobile app repositories spanning Android Native (Java/Kotlin), React
42 Native (TypeScript), and Flutter (Dart). Each instance pairs an authentic developer-reported issue with
43 its merged pull request and executable test patches, enabling automated validation of model-generated
44 fixes. To support reliable evaluation, MobileDev-Bench executes tests under three configurations: the
45 base commit, the test patch only, and the full fix, capturing transitions such as Fail-to-Pass, None-to-
46 Pass, Pass-to-Pass, and Pass-to-Fail. Unlike prior benchmarks that exclude instances with compilation
47 failures, we explicitly model transitions from non-existent or non-executable test states and restore
48 projects so that they compile when necessary. This compilation-aware design is particularly important
49 for *statically typed languages* common in mobile app development, such as Java and Kotlin, where
50 successful compilation is required before tests can execute.

51 Table 1 summarizes how MobileDev-Bench differs from prior issue-resolution benchmarks along four
52 key dimensions. First, it targets production mobile apps drawn from mature and actively maintained
53 open-source repositories, each with at least 400 GitHub stars. Second, it evaluates fixes that span
54 multiple artifact types, including source code, resources, and configuration metadata. Third, it
55 incorporates *manual verification* and difficulty stratification to ensure benchmark quality, with tasks
56 distributed across Easy (139), Medium (140), and Hard (136) tiers. Finally, it exhibits substantially
57 higher patch complexity: fixes modify 13.5 files and 344.2 lines on average, and 42.2% of instances
58 require coordinated changes across multiple artifact types. These properties capture the multi-file
59 and multi-artifact dependencies common in mobile app development that are largely absent from
60 existing library-centric benchmarks. A detailed comparison of MobileDev-Bench against existing
61 issue-resolution benchmarks, including patch complexity analysis and per-repository statistics, is
62 provided in Appendix C.

63 We extend Agentless [Xia et al., 2024] to support Java, Kotlin, Dart, and TypeScript and evaluate four
64 state-of-the-art LLMs on MobileDev-Bench. Despite strong performance on prior issue-resolution
65 benchmarks, models achieve low resolution rates of 3.90%–4.28% (Agentless) and 1.94%–6.07%
66 (Oracle) on MobileDev-Bench. Our analysis reveals recurring failure modes specific to mobile devel-
67 opment, including misaligned resource updates, improper dependency configuration, and incomplete
68 cross-artifact coordination. Across models, fault localization emerges as the dominant bottleneck,
69 with file-level recall ranging from 13.5% to 18.3%.

70 This work makes the following contributions:

- 71 • We introduce MobileDev-Bench, a benchmark of 415 verified mobile app issue-resolution
- 72 tasks drawn from 19 production repositories across Android, React Native, and Flutter.
- 73 • We develop a compilation-aware, execution-based evaluation framework for issue resolution
- 74 across Java, Kotlin, Dart, and TypeScript mobile app projects and release MobileDev-
- 75 Bench with task instances, an evaluation harness, and containerized environments at <https://huggingface.co/datasets/MobileDev-Bench/mobiledev-bench>.
- 76

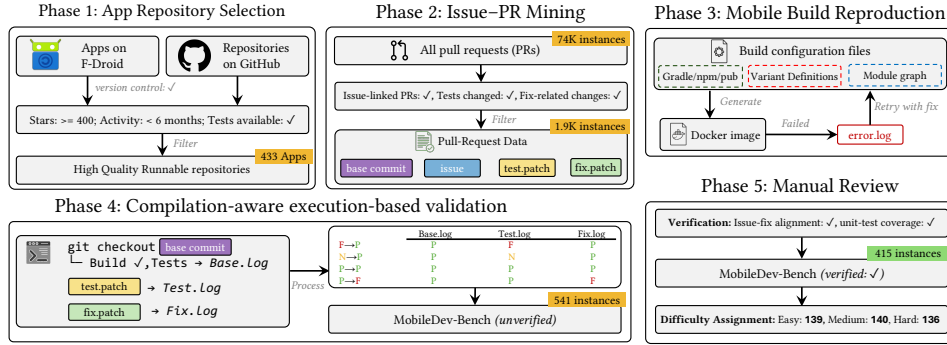


Figure 1: Overview of MobileDev-Bench construction pipeline.

- 77 • We provide a comprehensive evaluation of four frontier code-capable LLMs, showing that
 78 resolution rates drop to 3.90%–6.07% and that fault localization emerges as the dominant
 79 bottleneck for mobile app issue resolution.

80 2 Related Work

81 Early benchmarks such as HumanEval [Chen et al., 2021] and MBPP [Austin et al., 2021] established
 82 a foundation for assessing the programming capabilities of LLMs. These benchmarks primarily
 83 focused on small, self-contained problems: 95.8% target Python and over 83% evaluate function- or
 84 statement-level tasks [Cao et al., 2024], overlooking the real-world repository-level complexity.

85 SWE-bench [Jimenez et al., 2024] introduced repository-scale evaluation with 2,294 instances derived
 86 from real GitHub issues in Python projects, requiring models to generate patches validated through test
 87 execution. Subsequent work extended this direction to additional languages. Multi-SWE-bench [Zan
 88 et al., 2025] spans seven programming languages with 1,632 instances, while SWE-bench Java [Zan
 89 et al., 2024] provides 91 Java-specific tasks. SWE-PolyBench [Rashid et al., 2025] further shows
 90 that realistic repair tasks often require edits across multiple files. Despite this progress, existing
 91 benchmarks primarily focus on library-style repositories [Guo et al., 2025, Mhatre et al., 2025]
 92 and do not capture platform-specific app development environments. Alongside these benchmarks,
 93 LLM-based repair systems including SWE-agent [Yang et al., 2024a], AutoCodeRover [Zhang et al.,
 94 2024b], MAGIS [Tao et al., 2024], OpenDevin [Wang et al., 2024], CodeR [Chen et al., 2024],
 95 MASAI [Arora et al., 2024], and Agentless [Xia et al., 2024] have demonstrated strong results on
 96 these benchmarks, yet none target mobile application repositories.

97 Mobile app development *introduces additional challenges* not represented in current benchmarks,
 98 including framework-specific architectures, complex build systems, and coordination across hetero-
 99 geneous artifacts such as source code, resources, and configuration files. Mobile app development
 100 spans Android, iOS, and cross-platform frameworks such as React Native [Facebook, 2015] and
 101 Flutter [Google, 2018]. Prior research in mobile apps has focused on static analysis and testing tools
 102 such as FlowDroid [Arzt et al., 2014], AutoComply [Fakorede and Farooq, 2025], and DroidBot [Li
 103 et al., 2017]. More recently, LLM-based approaches have targeted mobile UI task automation [Wen
 104 et al., 2024] and mobile agent evaluation [Deng et al., 2024], while benchmarks for evaluating
 105 AI-assisted *code issue resolution* in mobile repositories remain *largely unexplored*.

106 MobileDev-Bench addresses this gap by targeting production mobile app repositories with heteroge-
 107 neous build environments and multi-artifact fixes, complementing existing library-centric benchmarks
 108 and broadening evaluation for AI-assisted software engineering systems.

109 3 MobileDev-Bench Construction

110 MobileDev-Bench is constructed from real issue-pull request (PR) pairs in open-source mobile app
 111 repositories. The construction process follows a five-phase pipeline as shown in Figure 1, designed
 112 to ensure that each instance corresponds to a real software engineering task and supports reliable
 113 execution-based validation.

114 3.1 Phase 1: App Selection

115 We begin by constructing a pool of candidate mobile app repositories spanning three major frame-
116 works: Android Native (Java/Kotlin), React-Native (TypeScript/JavaScript), and Flutter (Dart).
117 Repositories are collected from F-Droid [F-Droid, 2024], a catalog of free and open-source Android
118 apps, and from GitHub searches filtered by mobile framework tags. We prioritize production apps
119 distributed through app stores, as they typically represent mature codebases with active development
120 and testing practices. This process yields an initial pool of 433 candidate repositories.

121 To improve quality and benchmark viability, we filter repositories based on community adoption
122 (≥ 400 GitHub stars), recent maintenance, an open-source license, a GitHub-hosted issue tracker, and
123 verified buildability in a clean environment. The detailed filter criteria are provided in Appendix A.

124 3.2 Phase 2: PR Collection and Filtering

125 From the selected repositories, we retrieve merged pull requests (PR) using the GitHub API, collecting
126 metadata such as issue descriptions, pull request text, commit history, and patch information. To
127 ensure stable repository states, we only consider PRs merged into the repository’s default branch.

128 To construct valid issue-resolution tasks, we apply the following filtering steps: (i) We retain only
129 PRs that are explicitly linked to a GitHub issue using standard closure patterns (e.g., “fixes #123”,
130 “resolves #45”). This ensures that each task has a clear problem description derived from the issue
131 report. (ii) PRs must modify or introduce test files so that the correctness of the fix can be verified
132 through automated execution. (iii) For each retained PR, we extract two patches: a test patch
133 containing the test changes and a fix patch containing the changes. These patches are anchored to
134 the PR’s base commit to enable controlled replay of the development history. Thus, each candidate
135 instance contains four components: (a) the base commit, (b) the issue description, (c) the test patch,
136 and (d) the fix patch. These filters yield a set of candidate issue-PR pairs for subsequent validation.

137 3.3 Phase 3: Environment Runtime Setup

138 Mobile app builds require precise environment configurations spanning language runtimes, framework
139 SDKs, build tools, and dependencies. To support reproducibility, we automatically construct build
140 environments by analyzing repository metadata and generating framework-specific Dockerfiles. For
141 each candidate instance, we inspect the repository at the base commit and parse framework configura-
142 tion files (`build.gradle/.kts` and `gradle/libs.versions.toml` for Android; `package.json`
143 for React Native; `pubspec.yaml` for Flutter) to infer toolchains and dependencies. When configura-
144 tion files are incomplete, we supplement them with README documentation and GitHub Actions
145 workflow files. We then generate a per-repository Dockerfile. The environment remediation details
146 are provided in Appendix A.

147 3.4 Phase 4: Execution-Based Validation

148 We validate candidate instances through execution-based filtering by running tests under three
149 configurations within isolated Docker environments: (i) the base commit, (ii) the base commit with
150 the test patch applied (`test.patch`), and (iii) the base commit with both test and fix patches applied
151 (`fix.patch`). Rather than generic framework commands, we run instance-specific test commands
152 derived from the PR’s CI configuration, targeting only the tests relevant to each instance. Android
153 instrumentation tests that require emulators are excluded due to reproducibility concerns.

154 For each execution, we record four states: PASS, FAIL, NONE (test absent in that configuration), and
155 SKIP (see Appendix A for more). We derive test-state transitions and retain instances satisfying three
156 criteria: (i) at least one test exhibits a *Fail-to-Pass (F2P)* or *None-to-Pass (N2P)* transition indicating
157 the fix resolves the issue, (ii) no test exhibits a *Pass-to-Fail (P2F)* regression, and (iii) both the
158 base and fix configurations build and execute tests successfully (we provide formal definitions in
159 Appendix A). This execution-based filtering reduces the candidate pool from 1,939 to 541 instances.

160 3.5 Phase 5: Human Verification

161 Following SWE-bench Verified [OpenAI, 2024], we perform a manual review to ensure benchmark
162 quality. Two annotators independently examine each of the 541 candidate instances that passed
163 execution-based filtering and resolve disagreements through discussion. Review focuses on four

Table 2: Per-repository statistics for MobileDev-Bench. #Files and #LoC report codebase size at the base commit. # is the instance count. #Tok is the average issue description length in tokens. Fix patch columns (#L: changed lines, #H: hunks, #F: files) and test transition columns (F2P: Fail-to-Pass, N2P: None-to-Pass, P2P: Pass-to-Pass) report per-instance averages.

Repository	Codebase		Instances		Fix Patch (Avg.)			Tests (Avg.)		
	#Files	#LoC	#	#Tok	#L	#H	#F	F2P	N2P	P2P
Flutter (Dart)										
talawa	3,264	236k	12	280	256.6	18.0	5.9	1.0	4.2	3.8
zulip-flutter	514	141k	51	325	254.9	18.4	7.9	1.2	53.7	49.9
Android Native (Java/Kotlin)										
antennapod	1,143	115k	9	1,567	348.1	35.7	12.9	0.0	15.7	0.0
apps-android-commons	1,439	150k	10	548	1356.4	74.3	27.7	0.7	68.5	0.8
element-x-android	6,504	341k	72	154	387.1	41.8	18.7	1.9	31.5	22.8
medtimer	556	41k	9	162	472.0	32.1	16.6	0.4	4.7	0.6
geto	451	23k	1	156	172.0	33.0	5.0	0.0	52.0	0.0
jerboa	337	122k	4	218	12.8	1.8	1.2	4.0	11.8	1.0
neostumbler	627	39k	2	282	12.0	2.0	1.0	1.0	0.0	0.5
openhab-android	1,173	52k	8	228	584.1	13.5	8.1	0.6	10.9	5.6
voice	649	37k	4	259	273.0	10.0	7.5	0.3	4.8	0.5
streetcomplete	5,534	629k	18	411	625.4	52.4	27.1	0.6	68.3	7.1
thunderbird-android	6,389	365k	101	408	341.1	27.9	16.8	0.8	30.2	3.1
tusky	1,593	169k	13	464	234.8	11.2	5.4	1.3	17.9	7.2
wordpress-android	6,342	776k	84	159	174.9	17.4	7.2	0.3	28.9	0.8
React-Native (TypeScript)										
artsy/eigen	2,962	261k	1	287	2311.0	116.0	18.0	3.0	0.0	0.0
expensify/app	7,923	2,652k	2	547	429.0	42.0	17.5	1.5	1.0	2.0
NMF-earth/nmf-app	900	21k	12	143	439.8	40.3	17.3	3.9	0.2	0.4
rocket.chat/reactnative	1,601	134k	2	442	24.5	9.0	5.0	3.5	3.5	0.5

164 aspects: (i) Problem Statement Clarity to ensure that the GitHub issue provides sufficient context
 165 to understand the task; (ii) Test Coverage Appropriateness by verifying that the tests correctly
 166 correspond to the issue being resolved; (iii) Task Difficulty by labeling instances as Easy, Medium, or
 167 Hard based on conceptual complexity, number of files modified, and required framework knowledge;
 168 and (iv) Task Category by classifying each instance into one of the task types: Bug Fix, New Feature,
 169 or Feature Optimization.

170 Instances with severely underspecified issues or misaligned tests are excluded. In total, 126 instances
 171 (23.3%) fail manual verification, yielding a final benchmark of 415 instances. Difficulty and category
 172 labels are retained as metadata for analyzing model performance across task complexity levels and
 173 task types. More details are provided in Appendix B.

174 4 Benchmark Characteristics

175 4.1 Benchmark Statistics

176 MobileDev-Bench comprises 415 manually verified instances drawn from 19 production mobile
 177 repositories spanning three ecosystems: Flutter (Dart), Android Native (Java/Kotlin), and React-
 178 Native (TypeScript). Per-repository statistics are reported in Table 2. Repository sizes range from 337
 179 to 7,923 source files (332K lines on average), reflecting the diversity of mature open-source codebases
 180 included in the benchmark. Instance counts vary from 1 to 101 per repository, with the four largest
 181 contributors being thunderbird-android (101), wordpress-android (84), element-x-android (72), and
 182 zulip-flutter (51), together accounting for 74.2% of the dataset.

183 Table 3 reports aggregate statistics across the dataset. Issue descriptions average 310 tokens (median
 184 173), with a long tail reaching 11,501 tokens for the most detailed reports, reflecting the context-rich
 185 nature of app bug reports that typically include reproduction steps, platform versions, and device-
 186 specific conditions. Fix patches are larger than those in comparable benchmarks as shown in Figure
 187 3: the mean patch modifies 344 lines across 13.5 files and 28.4 hunks, while the median is 97 lines,
 188 5 files, and 12 hunks. The gap between mean and median reflects a right-skewed distribution, with
 189 outlier fixes in repositories such as apps-android-commons (avg. 1,356 lines) and openhab-android
 190 (avg. 584 lines) substantially inflating the mean. None-to-Pass is the dominant validation signal: 355

Table 3: Summary statistics for MobileDev-Bench.

	Issue	Codebase		Fix Patch			Tests		
	Tokens	#Files	#Lines (K)	#Lines	#Files	#Hunks	F2P	N2P	P2P
Mean	310	4,601	386	344	13.5	28.4	1.1	31.5	11.8
Median	173	6,342	341	97	5	12	0	11	0
Max	11,501	7,923	2,652	11,915	237	592	30	375	267

of 415 instances (85.5%) yield at least one N2P transition, with a mean of 31.5 and a median of 11. N2P transitions arise from two sources: tests that did not exist at the base commit and are introduced by the fix, and tests that were present but could not be executed due to compilation failures at the base commit, a common occurrence in statically typed languages such as Java and Kotlin where the fix itself is required to restore compilability. Additionally, 148 instances (35.7%) yield at least one Fail-to-Pass transition (mean 1.1), indicating cases where pre-existing tests directly captured the regression. The high N2P mean is driven by repositories with large test suites (e.g., zulip-flutter avg. 53.7 N2P per instance), where the fix enables previously non-executable test modules or introduces newly added tests.

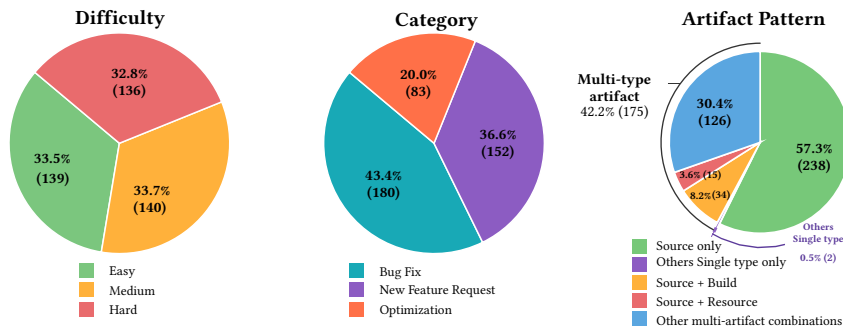


Figure 2: Distribution of MobileDev-Bench tasks by difficulty, category, and artifact pattern. Tasks are balanced across difficulty tiers, with bug fixes as the largest category. Most tasks modify source files only, while 42.2% require coordinated changes across multiple artifact types.

Figure 2 shows the distribution of instances in MobileDev-Bench across difficulty and task category. Difficulty is broadly balanced across Easy (under 15 min, 33.5%), Medium (15 min to 1 hr, 33.7%), and Hard (over 1 hr, 32.8%), ensuring the benchmark exercises a range of problem complexities. By task category, 43.4% of instances are Bug Fixes, 36.6% are New Features, and 20.0% are Feature Optimizations.

4.2 Benchmark Features

MobileDev-Bench captures characteristics of mobile app development that are absent from existing repository-level benchmarks. In particular, several fixes span multiple artifact types and may involve more than one language within a single patch.

Multi-Artifact Fix Patterns. Mobile app fixes often require coordinated changes across heterogeneous artifact types. We classify modified files in each patch into eight categories: (1) Source (.kt, .java, .dart, .ts, .tsx, .js); (2) Resource (res/ subtree, image and media assets); (3) i18n (values-*/strings.xml, Flutter .arb); (4) Build (build.gradle, libs.versions.toml, *.kts); (5) Manifest (AndroidManifest.xml, pubspec.yaml); (6) Config (CI/CD workflows, linter config files); (7) Docs (Markdown, changelog fragments, license files); and (8) Other (files not matched by the above). Source files appear in nearly all instances (99.5%). However, many fixes also require non-source changes: 21.2% modify resources, 16.9% modify localization files, 19.3% modify build files, and 3.6% modify manifest files. As summarized in Figure 2, 42.2% of instances modify two or more artifact types and 21.9% modify three or more, demonstrating that cross-artifact coordination is a common property of mobile app fixes. Full combination-level counts with per-language breakdowns are reported in Table 9 in the appendix.

Intra-Patch Language Diversity. Fix patches also frequently span multiple languages. We identify cross-language instances by examining the set of source-language file extensions modified within each fix patch. In total, 20.7% of instances modify source files written in more than one language. The most common pattern is Kotlin–Java co-modification, representing 19.8% of all instances and

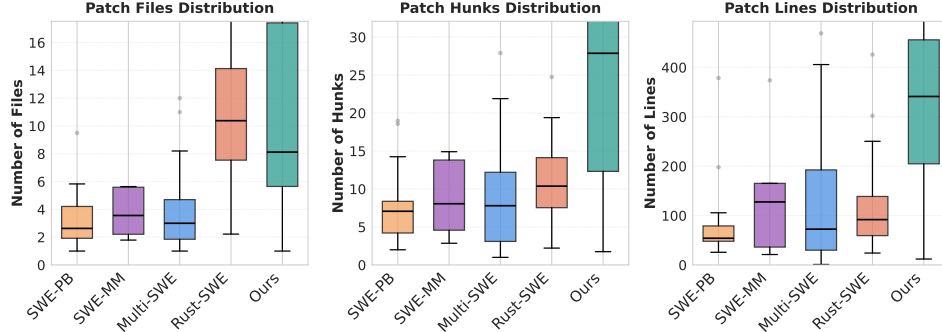


Figure 3: Distribution of patch complexity metrics (files modified, hunks, changed lines) across five benchmarks (SWE-PB: SWE-PolyBench; SWE-MM: SWE-bench Multimodal; Multi-SWE: Multi-SWE-bench; Rust-SWE: Rust-SWE-bench; Ours: MobileDev-Bench). Each data point is the per-repository average. Y-axes are capped at the 90th percentile for visual clarity.

225 25.2% of Kotlin-based fixes. This reflects the ongoing migration of Android projects from Java to
 226 Kotlin, with co-modification rates varying widely across repositories depending on their stage of
 227 migration: projects retaining residual Java code exhibit rates above 40%, while repositories that
 228 have fully transitioned to Kotlin show none. A second cross-language pattern occurs in React-
 229 Native repositories, where TypeScript–JavaScript co-modification appears in 18.8% of React-Native
 230 instances, reflecting gradual TypeScript adoption in existing JavaScript codebases.

Table 4: Resolution rates (%) by difficulty level under Agentless and Oracle settings.

Difficulty	Tasks	Agentless				Oracle			
		Claude	Qwen	GPT	Gemini	Claude	Qwen	GPT	Gemini
Easy	139	8.0	8.5	8.0	8.0	12.9	6.7	11.6	4.3
Medium	140	4.4	1.5	1.5	1.5	3.6	0.7	2.2	1.4
Hard	136	0.0	3.0	2.2	3.0	1.5	0.0	1.5	0.0
Overall	415	4.15	4.28	3.90	4.16	6.07	2.52	5.19	1.94

231 5 Experimental Setup

232 We evaluate models under two complementary settings designed to isolate the two stages of automated
 233 repair [Xia and Zhang, 2023]: fault localization and patch generation.

234 **Agentless.** We adopt Agentless [Xia et al., 2024] as our primary evaluation framework, selecting
 235 it as a competitive baseline whose modular architecture supports extension to new programming
 236 languages. We extend it to support Java, Kotlin, TypeScript, and Dart by replacing its Python-specific
 237 AST parser with tree-sitter [Brunsfeld et al., 2018] (see Appendix E.2 for details).

238 **Oracle.** We additionally evaluate an *oracle localization* setting [Jimenez et al., 2024], in which the
 239 ground-truth set of modified files is provided directly to the model, simulating a scenario where fault
 240 localization is perfect. The model receives the issue description and the exact files it must modify, and
 241 is asked solely to generate the patch. This setting serves two purposes. First, it establishes an upper
 242 bound on repair performance achievable under ideal localization conditions. Second, by comparing
 243 Oracle and Agentless resolution rates, we can attribute the performance gap to localization failure
 244 versus repair failure, isolating which component is the primary bottleneck.

245 In both settings, generated patches are applied to the base commit and validated by executing the
 246 repository test suite within our Docker-based [Merkel, 2014] containerized environment.

247 **Models.** We evaluate four frontier code-capable models spanning proprietary and open-weight
 248 systems: Claude Sonnet 4.5 [Anthropic, 2025], GPT-5.2 [OpenAI, 2025], Gemini 2.5 Flash [Comanici
 249 et al., 2025], and Qwen-3-Coder [Qwen Team, 2025]. To ensure deterministic patch generation
 250 and fair comparison, we disable extended reasoning modes where available, following the original
 251 Agentless evaluation protocol.

252 **Metrics.** Following prior work [Jimenez et al., 2024, Xia et al., 2024], we report two complementary
 253 metrics evaluated under both settings. (i) *Resolution Rate*: the percentage of tasks where the generated

Table 5: Resolution rate breakdown by artifact diversity. Resolved column shows the share of uniquely resolved instances in each category; Ratio = Resolved% / Benchmark%.

Category	Bench.	Agentless		Oracle	
		Resolved	Ratio	Resolved	Ratio
Single-type artifact	57.8%	84.4%	1.46x	86.8%	1.50x
Multi-type artifact	42.2%	15.6%	0.37x	13.2%	0.31x

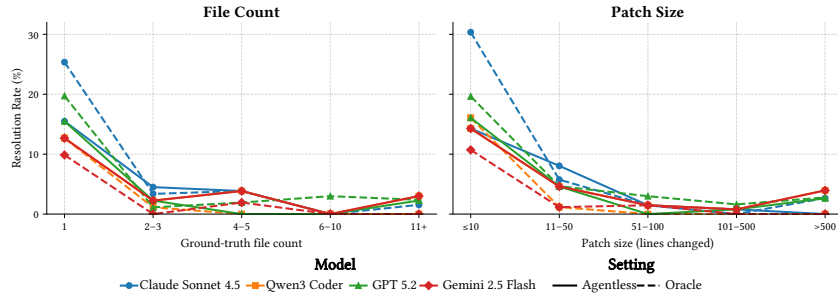


Figure 4: Resolution rates by ground-truth file count (left) and patch size (right). Single-file tasks achieve 13–16% (Agentless) and 10–25% (Oracle); rates collapse for multi-file and large-patch tasks.

254 patch passes the full test suite, measuring end-to-end repair success. (ii) *Localization Metrics*: file-
 255 level Precision, Recall, and F1 (macro-averaged across tasks), measuring how accurately the model’s
 256 patch covers the ground-truth modified files. In the Agentless setting these metrics reflect the quality
 257 of automated fault localization; in the Oracle setting they reflect how completely the model generates
 258 edits across all files it was given. Comparing the two settings reveals whether the dominant failure
 259 mode is localization (finding the right files) or repair (generating correct changes once files are
 260 known). Formal metric definitions are in Appendix E.

261 6 Results

262 6.1 Repair Performance

263 Table 4 reports resolution rates across difficulty levels under both settings. Performance declines
 264 with task difficulty in both cases. Under Agentless, Easy tasks reach 8.0–8.5%, Medium tasks drop
 265 to 1.5–4.4%, and Hard tasks are rarely solved (0.0–3.0%). Under Oracle, Easy task rates rise to
 266 4.3–12.9%, confirming that localization is a key bottleneck on simpler tasks, while Hard tasks remain
 267 difficult (0.0–1.5%) even with ground-truth files. More details are provided in Appendix F.

268 6.2 Resolution Characteristics

269 Table 5 analyzes how resolved instances differ from the overall benchmark distribution across artifact
 270 diversity. The comparison highlights structural factors that influence repair success.

271 **Artifact Type Diversity.** Resolution success is strongly associated with artifact diversity. While
 272 single-artifact tasks constitute 57.8% of the benchmark, they account for 84.4% of resolved Agentless
 273 instances, representing a 1.46x over-representation. In contrast, multi-artifact tasks require coordi-
 274 nated modifications across heterogeneous files and account for 42.2% of the benchmark but only
 275 15.6% of resolved instances (0.37x). This gap indicates that current models struggle to coordinate
 276 fixes across multiple artifact types, such as source code, build configuration, and resource files.

277 **File Count Effects.** File count shows the strongest relationship with repair success. Single-file tasks,
 278 which represent only 17.1% of the benchmark, account for 65.6% of all resolved Agentless instances,
 279 a 3.84x overrepresentation. As the number of required file modifications increases, resolution success
 280 drops sharply. Tasks requiring 4–10 files are particularly challenging, accounting for 29.4% of
 281 benchmark tasks but only 6.3% of resolved instances. Tasks with 11+ files remain difficult but exhibit
 282 a modest recovery in resolution share (12.5%), a pattern further analyzed in Appendix F.3.

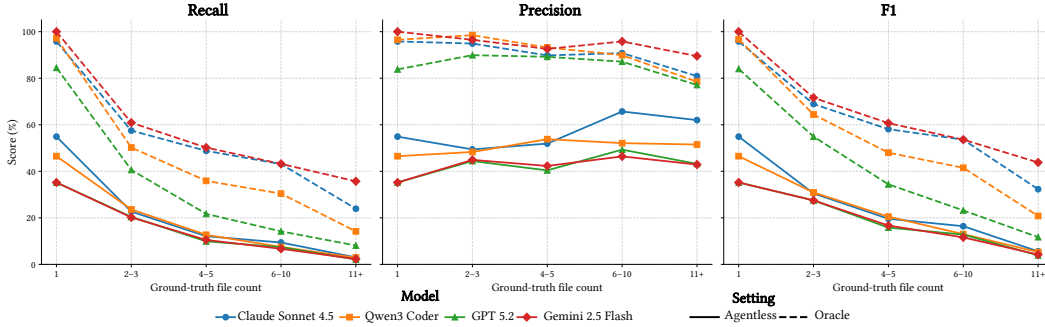


Figure 5: Recall decreases as the number of modified files increases.

283 6.3 Fault Localization Performance

284 Figure 5 shows file-level recall as a function of the number of ground-truth files under the Agentless
 285 setting. Across models, overall recall ranges from 13.5% to 18.3%, indicating that models frequently
 286 miss files that must be modified for a correct fix. In contrast, precision remains higher (42–58%,
 287 Appendix Table 12), suggesting that when models choose to edit a file, it is often relevant, but many
 288 required files are not identified.

289 Recall decreases as task complexity increases. Single-file tasks achieve the highest recall (35–55%),
 290 but recall drops to 20–24% for tasks with 2–3 files, 10–13% for tasks with 4–5 files, and 7–9% for
 291 tasks with 6–10 files. For the most complex tasks requiring 11+ files, recall falls to ~ 2 –3% across
 292 all models. This trend reflects the structural characteristics of mobile apps. Issues frequently require
 293 coordinated changes across multiple layers – source code, resources, and configuration artifacts.
 294 Current models rarely capture these cross-file dependencies, leading to missed modification sites.

295 **Fault localization is the dominant bottleneck.** Fault localization is widely recognized as a key
 296 determinant of automated repair success [Wong et al., 2016, Lou et al., 2020]. The low recall in the
 297 Agentless setting imposes a hard ceiling on resolution: even perfect patch generation could achieve at
 298 most 13–18% resolution rate (matching recall). Observed rates (3.90%–4.28% under Agentless) fall
 299 well below this ceiling, confirming that unresolved tasks fail first at localization.

300 The Oracle setting (Table 12) isolates the repair stage by providing ground-truth files directly. Even
 301 under these ideal conditions, recall reaches only 30.9–55.2% overall—models receive the exact files
 302 to modify but still generate patches that omit changes to many of them. Recall in the Oracle setting
 303 degrades similarly with file count: from 85–100% on single-file tasks to 8–36% on tasks with 11+
 304 files. This repair-side failure is compounded on top of localization failure in the Agentless setting,
 305 where models must both discover the right files and modify all of them correctly.

306 Comparing the two settings quantifies each failure mode. The large gap between Oracle and Agentless
 307 recall (roughly 17–42 percentage points overall) confirms that fault localization is the dominant
 308 bottleneck. However, the Oracle recall falling well below 100% on multi-file tasks reveals a secondary
 309 challenge: even given perfect file knowledge, models struggle to generate coherent edits across many
 310 files simultaneously. Improving MobileDev-Bench resolution rates therefore requires advances on
 311 both fronts: dependency-aware multi-file retrieval for localization, and long-range patch coherence
 312 for repair. Comprehensive resolution rate breakdowns, per-model analysis, and additional resolution
 313 characteristics are provided in Appendix F.

314 7 Conclusion

315 We introduced MobileDev-Bench, a benchmark for evaluating LLM-based program repair on mobile
 316 app development tasks. It contains 415 verified instances from 19 production repositories across
 317 Android Native, React Native, and Flutter, capturing several properties of mobile development such
 318 as multi-artifact coordination and framework-driven execution. Evaluation of four frontier LLMs
 319 shows low resolution rates, with fault localization as the primary bottleneck. File-level recall remains
 320 limited and declines further for complex tasks, particularly multi-file and multi-artifact fixes. These
 321 results highlight the need for a better understanding of mobile app frameworks to support multi-file
 322 dependency reasoning and artifact-aware localization. MobileDev-Bench provides a reproducible
 323 benchmark and baseline results for future research in AI-assisted mobile app development.

324 References

- 325 Android. The activity lifecycle. 2020. URL [https://developer.android.com/guide/](https://developer.android.com/guide/components/activities/activity-lifecycle)
326 [components/activities/activity-lifecycle](https://developer.android.com/guide/components/activities/activity-lifecycle). Accessed 2026-02-12.
- 327 Android. App resources overview. 2024. URL [https://developer.android.com/guide/](https://developer.android.com/guide/topics/resources/providing-resources)
328 [topics/resources/providing-resources](https://developer.android.com/guide/topics/resources/providing-resources). Accessed 2026-02-12.
- 329 Android. Configure your build, 2025. URL <https://developer.android.com/build>. Accessed
330 2026-02-12.
- 331 Anthropic. Claude: A large language model by anthropic. 2023. URL [https://www.anthropic.](https://www.anthropic.com)
332 [com](https://www.anthropic.com).
- 333 Anthropic. Claude sonnet 4.5. <https://www.anthropic.com/news/claude-sonnet-4-5>, 2025.
334 Accessed: 2026-02-17.
- 335 Apple Inc. Xcode and Apple SDKs Agreement, 2024. URL [https://www.apple.com/legal/](https://www.apple.com/legal/sla/docs/xcode.pdf)
336 [sla/docs/xcode.pdf](https://www.apple.com/legal/sla/docs/xcode.pdf). Accessed 2026-02-12.
- 337 D. Arora, A. Sonwane, N. Wadhwa, A. Mehrotra, S. Utpala, R. Baire, A. Kanade, and N. Natarajan.
338 Masai: Modular architecture for software-engineering ai agents. *arXiv preprint arXiv:2406.11638*,
339 2024.
- 340 S. Arzt, S. Rasthofer, C. Fritz, E. Bodden, A. Bartel, J. Klein, Y. Le Traon, D. Octeau, and P. McDaniel.
341 Flowdroid: Precise context, flow, field, object-sensitive and lifecycle-aware taint analysis for
342 android apps. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language*
343 *Design and Implementation*, pages 259–269, 2014.
- 344 J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry,
345 Q. Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*,
346 2021.
- 347 M. Brunsfeld et al. tree-sitter: An incremental parsing system for programming tools. <https://tree-sitter.github.io/tree-sitter/>, 2018. Accessed: 2026-02-17.
- 349 J. Cao, Z. Chen, J. Wu, S.-C. Cheung, and C. Xu. Javabench: A benchmark of object-oriented code
350 generation for evaluating large language models. *arXiv preprint arXiv:2406.12902*, 2024.
- 351 D. Chen, S. Lin, M. Zeng, D. Zan, J.-G. Wang, A. Cheshkov, J. Sun, H. Yu, G. Dong, A. Aliev, et al.
352 Coder: Issue resolving with multi-agent and task graphs. *arXiv preprint arXiv:2406.01304*, 2024.
- 353 M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda,
354 N. Joseph, G. Brockman, et al. Evaluating large language models trained on code. *arXiv preprint*
355 *arXiv:2107.03374*, 2021.
- 356 G. Comanici, E. Bieber, M. Schaekermann, I. Pasupat, et al. Gemini 2.5: Pushing the frontier with
357 advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025.
358 URL <https://arxiv.org/abs/2507.06261>.
- 359 S. Deng, W. Xu, H. Sun, W. Liu, T. Tan, J. Liu, A. Li, J. Luan, B. Wang, R. Yan, and S. Shang.
360 Mobile-bench: An evaluation benchmark for llm-based mobile agents, 2024. URL <https://arxiv.org/abs/2407.00993>.
361 <https://arxiv.org/abs/2407.00993>.
- 362 F-Droid. F-Droid. <https://f-droid.org>, 2024. Free and Open Source Android App Repository.
- 363 Facebook. React native, 2015. URL <https://reactnative.dev>.
- 364 M. Fakorede and U. Farooq. Understanding and detecting platform-specific violations in android
365 auto apps, 2025. URL <https://arxiv.org/abs/2503.04003>.
- 366 Flutter. Start thinking declaratively, 2025. URL [https://docs.flutter.dev/](https://docs.flutter.dev/data-and-backend/state-mgmt/declarative)
367 [data-and-backend/state-mgmt/declarative](https://docs.flutter.dev/data-and-backend/state-mgmt/declarative). Accessed 2026-02-12.
- 368 Flutter. Introduction to declarative ui, 2026. URL [https://docs.flutter.dev/flutter-for/](https://docs.flutter.dev/flutter-for/declarative)
369 [declarative](https://docs.flutter.dev/flutter-for/declarative). Accessed 2026-02-12.

- 370 Google. Flutter, 2018. URL <https://flutter.dev>.
- 371 L. Guo, W. Tao, R. Jiang, Y. Wang, J. Chen, X. Liu, Y. Ma, M. Mao, H. Zhang, and Z. Zheng.
372 Omnigirl: A multilingual and multimodal benchmark for github issue resolution, 2025. URL
373 <https://arxiv.org/abs/2505.04606>.
- 374 C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. Narasimhan. Swe-bench: Can
375 language models resolve real-world github issues? 2024. URL [https://arxiv.org/abs/2310.](https://arxiv.org/abs/2310.06770)
376 06770.
- 377 K. Kabir. Swe-bench multilingual, 2025. URL <https://kabirk.com/multilingual>.
- 378 Y. Li, Z. Yang, Y. Guo, and X. Chen. Droidbot: A lightweight ui-guided test input generator for
379 android. In *2017 IEEE/ACM 39th International Conference on Software Engineering Companion*
380 (*ICSE-C*), pages 23–26. IEEE, 2017.
- 381 Y. Lou, A. Ghanbari, X. Li, L. Zhang, H. Zhang, D. Hao, and L. Zhang. Can automated program
382 repair refine fault localization? a unified debugging approach. In *Proceedings of the 29th ACM*
383 *SIGSOFT International Symposium on Software Testing and Analysis*, pages 75–87, 2020.
- 384 D. Merkel. Docker: Lightweight linux containers for consistent development and deployment. *Linux*
385 *journal*, 2014(239):2, 2014.
- 386 S. Mhatre, Y. Bajpai, S. Gulwani, E. Murphy-Hill, and G. Soares. Swe-sharp-bench: A reproducible
387 benchmark for c# software engineering tasks, 2025. URL [https://arxiv.org/abs/2511.](https://arxiv.org/abs/2511.02352)
388 02352.
- 389 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 390 OpenAI. Swe-bench verified. 2024. URL <https://openai.com/swe-bench-verified>.
- 391 OpenAI. Introducing GPT-5.2. <https://openai.com/index/introducing-gpt-5-2/>, 2025.
392 Accessed: 2026-02-17.
- 393 Qwen Team. Qwen3-coder: A new era for code intelligence. [https://qwen.ai/blog?id=](https://qwen.ai/blog?id=qwen3-coder)
394 [qwen3-coder](https://qwen.ai/blog?id=qwen3-coder), 2025. Accessed: 2026-02-17.
- 395 M. S. Rashid, C. Bock, Y. Zhuang, A. Buchholz, T. Esler, S. Valentin, L. Franceschi, M. Wistuba, P. T.
396 Sivaprasad, W. J. Kim, A. Deoras, G. Zappella, and L. Callot. Swe-polybench: A multi-language
397 benchmark for repository level evaluation of coding agents, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2504.08703)
398 [abs/2504.08703](https://arxiv.org/abs/2504.08703).
- 399 B. Roziere, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, T. Remez, J. Rapin,
400 et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- 401 W. Tao, Y. Zhou, Y. Wang, W. Zhang, H. Zhang, and Y. Cheng. MAGIS: LLM-based multi-agent
402 framework for github issue resolution. *arXiv preprint arXiv:2403.17927*, 2024.
- 403 X. Wang, B. Li, Y. Song, F. F. Xu, X. Tang, M. Zhuge, J. Pan, Y. Song, B. Li, J. Singh, et al.
404 Opendevin: An open platform for ai software developers as generalist agents. *arXiv preprint*
405 *arXiv:2407.16741*, 2024.
- 406 H. Wen, Y. Li, G. Liu, S. Zhao, T. Yu, T. J.-J. Li, S. Jiang, Y. Liu, Y. Zhang, and Y. Liu. Autodroid:
407 Llm-powered task automation in android. In *Proceedings of the 30th Annual International*
408 *Conference on Mobile Computing and Networking*, 2024.
- 409 W. E. Wong, R. Gao, Y. Li, R. Abreu, and F. Wotawa. A survey on software fault localization. *IEEE*
410 *Transactions on Software Engineering*, 42(8):707–740, 2016.
- 411 C. S. Xia and L. Zhang. Automated program repair in the era of large pre-trained language models.
412 *arXiv preprint arXiv:2210.14179*, 2023.
- 413 C. S. Xia, Y. Deng, S. Dunn, and L. Zhang. Agentless: Demystifying llm-based software engineering
414 agents, 2024. URL <https://arxiv.org/abs/2407.01489>.

- 415 J. Xiang, W. He, X. Wang, H. Tian, and Y. Zhang. Evaluating and improving automated repository-
416 level rust issue resolution with llm-based agents, 2026. URL [https://arxiv.org/abs/2602.](https://arxiv.org/abs/2602.22764)
417 [22764](https://arxiv.org/abs/2602.22764). ICSE 2026.
- 418 J. Yang, C. E. Jimenez, A. Wettig, K. Lieret, S. Yao, K. Narasimhan, and O. Press. Swe-agent:
419 Agent-computer interfaces enable automated software engineering, 2024a. URL [https://arxiv.](https://arxiv.org/abs/2405.15793)
420 [org/abs/2405.15793](https://arxiv.org/abs/2405.15793).
- 421 J. Yang, C. E. Jimenez, A. L. Zhang, K. Lieret, J. Yang, X. Wu, O. Press, N. Muennighoff, G. Syn-
422 naeve, K. R. Narasimhan, D. Yang, S. I. Wang, and O. Press. SWE-bench Multimodal: Do ai
423 systems generalize to visual software domains? *arXiv preprint arXiv:2410.03859*, 2024b.
- 424 D. Zan, Z. Huang, A. Yu, S. Lin, Y. Shi, W. Liu, D. Chen, Z. Qi, H. Yu, L. Yu, D. Ran, M. Zeng,
425 B. Shen, P. Bian, G. Liang, B. Guan, P. Huang, T. Xie, Y. Wang, and Q. Wang. Swe-bench-java: A
426 github issue resolving benchmark for java, 2024. URL <https://arxiv.org/abs/2408.14354>.
- 427 D. Zan, Z. Huang, W. Liu, H. Chen, L. Zhang, S. Xin, L. Chen, Q. Liu, X. Zhong, A. Li, S. Liu,
428 Y. Xiao, L. Chen, Y. Zhang, J. Su, T. Liu, R. Long, K. Shen, and L. Xiang. Multi-swe-bench: A
429 multilingual benchmark for issue resolving, 2025. URL <https://arxiv.org/abs/2504.02605>.
- 430 L. Zhang, D. Zan, Q. Yang, Z. Huang, D. Chen, B. Shen, T. Liu, Y. Gong, P. Huang, X. Lu,
431 G. Liang, L. Cui, and Q. Wang. Codev: Issue resolving with visual data, 2024a. URL [https:](https://arxiv.org/abs/2412.17315)
432 [//arxiv.org/abs/2412.17315](https://arxiv.org/abs/2412.17315).
- 433 Y. Zhang, H. Ruan, Z. Fan, and A. Roychoudhury. Autocoderover: Autonomous program improve-
434 ment. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing*
435 *and Analysis*, pages 1592–1604, 2024b.

Table 6: Summary and licenses for all GitHub repositories represented in MobileDev-Bench.

Repository	Summary	License
palisadoesfoundation/talawa	Community organization management	GNU GPLv3
zulip/zulip-flutter	Zulip mobile apps for Android and iOS	Apache 2.0
antennapod/antennapod	Podcast manager for Android	GNU GPLv3
commons-app/apps-android-commons	Wikimedia Commons Android uploader	Apache 2.0
element-hq/element-x-android	Android Matrix messenger (Jetpack Compose)	GNU AGPLv3
futsch1/medtimer	Medication reminder app	MIT
jackeblan/geto	Per-app device settings manager	GNU GPLv3
lemmy.net/jerboa	Native Android client for Lemmy	GNU AGPLv3
mjaakko/neostumbler	OpenStreetMap geolocation stumbler	MIT
openhq/openhab-android	openHAB client for Android	EPL 2.0
paulwoitaschek/voice	Minimalist audiobook player	GNU GPLv3
streetcomplete/streetcomplete	Easy-to-use OpenStreetMap editor	GNU GPLv3
thunderbird/thunderbird-android	Thunderbird for Android (fka K-9 Mail)	Apache 2.0
tuskyapp/tusky	Android client for Mastodon	GNU GPLv3
wordpress-mobile/wordpress-android	WordPress for Android	GNU GPLv2
artsy/eigen	Artsy’s iOS/Android art marketplace	MIT
expensify/app	Financial collaboration platform	MIT
NMF-earth/nmf-app	Carbon footprint tracker for daily life	GNU GPLv3
rocketchat/rocket.chat.reactnative	Secure communications platform	MIT

436 A Additional Construction Details

437 This appendix provides supplementary details for the construction pipeline described in Section 3.

438 A.1 Repository Selection Filters

439 The following five filters are applied during Phase 1 (App Selection) to determine benchmark viability:

- 440 (i) A minimum of 400 GitHub stars to ensure community adoption and sustained development.
- 441 (ii) Recent maintenance within the past six months, verified through commit activity, merged
442 pull requests, and responsive issue management.
- 443 (iii) A permissive open-source license (e.g., Apache 2.0, MIT, or GPL variants) allowing research
444 use.
- 445 (iv) An issue tracker hosted on GitHub, enabling reliable retrieval of issue–PR mappings and
446 contextual information.
- 447 (v) The latest commit must build and execute its test suite in a clean environment. For Android
448 Native apps, we install the required JDK version (11/17/21) and run Gradle build and test
449 tasks. For React Native apps, we verify Node.js compatibility and run tests via npm or yarn.
450 For Flutter apps, we confirm Flutter SDK compatibility and execute the test framework.
451 Repositories that fail this check are excluded.

452 A.2 Test File Identification Patterns

453 During Phase 2 (PR Collection and Filtering), test files are identified using the following framework-
454 specific naming conventions:

- 455 • **Android Native:** `*Test.kt`, `*Test.java`, or files under `androidTest/` and `test/` source
456 sets.
- 457 • **React Native:** `*.test.ts`, `*.test.tsx`, `*.test.js`, or `*.spec.ts` files.
- 458 • **Flutter:** `*_test.dart` files or any file under the `test/` directory.

459 A.3 Repository Composition

460 MobileDev-Bench draws from 19 open-source Android and cross-platform mobile repositories
 461 spanning utility apps, media players, productivity tools, and developer tooling. Table 6 lists all 19
 462 repositories with a brief description of their purpose and their open-source license. The selection
 463 prioritises projects with active CI pipelines, high test coverage, and sustained PR activity to ensure a
 464 sufficient supply of execution-validated instances.

465 A.4 Pipeline Conversion Statistics

466 Table 7 reports per-repository counts at each stage of the conversion pipeline: PRs crawled, attribute-
 467 filtered candidates, execution-validated instances, and manually verified benchmark instances. Over-
 468 all, the pipeline reduces 74,195 crawled PRs to 415 verified instances (0.6% retention rate). The
 469 primary bottleneck is attribute filtering (Phase 1), which discards PRs that lack an explicit issue
 470 link, do not modify test files, or touch more files than the complexity threshold permits. Execution
 471 validation (Phase 4) constitutes the second largest reduction, as a substantial fraction of candidates
 472 fail to build or produce non-deterministic test outcomes.

Table 7: Benchmark repository instance conversion pipeline showing progression from collected pull requests through attribute-based filtering, execution-based validation, and manual verification.

	PRs Crawled	Filtered	Validated	Verified
palisadoesfoundation/talawa	1,743	212 _{↓1531}	13 _{↓199}	12 _{↓1}
zulip/zulip-flutter	1,051	96 _{↓955}	51 _{↓45}	51 ₀
antennapod/antennapod	2,031	89 _{↓1942}	11 _{↓78}	9 _{↓2}
commons-app/apps-android-commons	1,477	200 _{↓1277}	17 _{↓183}	10 _{↓7}
element-hq/element-x-android	3,899	303 _{↓3596}	77 _{↓226}	72 _{↓5}
futsch1/medtimer	643	25 _{↓618}	12 _{↓13}	9 _{↓3}
jackeban/geto	195	20 _{↓175}	6 _{↓14}	1 _{↓5}
lemmyynet/jerboa	1,048	15 _{↓1033}	5 _{↓10}	4 _{↓1}
mjaakko/neostumbler	681	20 _{↓661}	5 _{↓15}	2 _{↓3}
openhab/openhab-android	1,607	27 _{↓1580}	12 _{↓15}	8 _{↓4}
paulwoitaschek/voice	1,364	10 _{↓1354}	6 _{↓4}	4 _{↓2}
streetcomplete/streetcomplete	1,167	96 _{↓1071}	27 _{↓69}	18 _{↓9}
thunderbird/thunderbird-android	2,728	245 _{↓2483}	110 _{↓135}	101 _{↓9}
tuskyapp/tusky	2,110	28 _{↓2082}	19 _{↓9}	13 _{↓6}
wordpress-mobile/wordpress-android	4,199	438 _{↓3761}	140 _{↓298}	84 _{↓56}
artsy/eigen	10,020	3 _{↓10017}	2 _{↓1}	1 _{↓1}
expensify/app	34,301	30 _{↓34271}	8 _{↓22}	2 _{↓6}
NMF-earth/nmf-app	181	23 _{↓158}	15 _{↓8}	12 _{↓3}
rocketchat/rocket.chat.reactnative	3,750	59 _{↓3691}	5 _{↓54}	2 _{↓3}
Total	74,195	1,939 _{↓72256}	541 _{↓1398}	415 _{↓126}

473 A.5 Environment Dockerfile Generation and Remediation

474 After generating a candidate Dockerfile from parsed configuration (as described in Phase 3, Section 3),
 475 we run a full build-and-test cycle to validate correctness. Failed builds are manually inspected;
 476 common fixes include adding missing system packages (e.g., libssl-dev for React Native native
 477 modules), pinning a breaking transitive dependency, or correcting a misidentified JDK version.
 478 Repositories requiring more than three remediation rounds are excluded.

479 For Android projects, GitHub Actions workflow files (.github/workflows/*.yml) often pro-
 480 vide the most precise environment specification, pinning exact Gradle wrapper versions, JDK
 481 versions, and Android SDK command-line tool versions that are absent from build.gradle or
 482 gradle.properties.

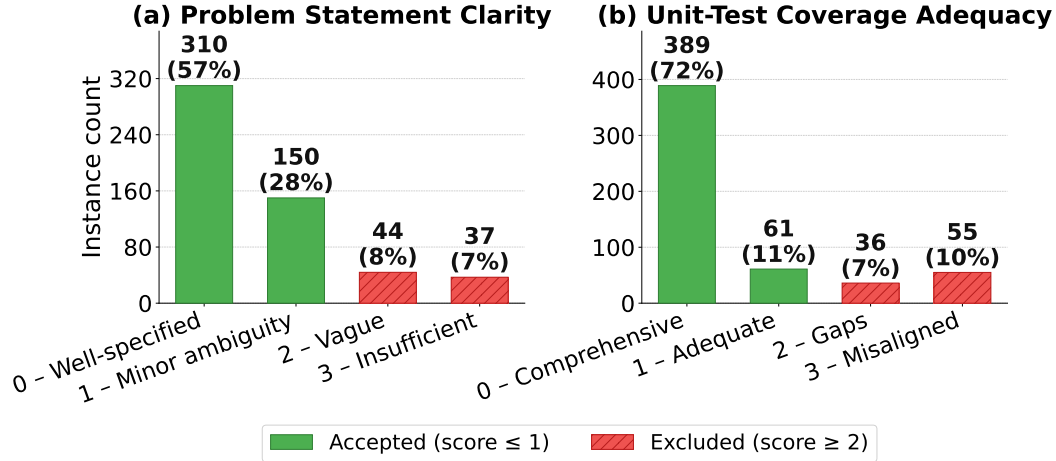


Figure 6: Distribution of manual annotation scores for (a) problem-statement clarity and (b) unit-test coverage adequacy, across all 541 annotated candidate instances. Green bars (scores 0–1) indicate instances accepted into the benchmark; hatched red bars (scores 2–3) indicate excluded instances. The majority of exclusions arise from vague issue descriptions (score 2/3 on dimension a) or tests that impose undue implementation constraints (score 2/3 on dimension b).

483 A.6 NONE State and None-to-Pass Semantics

484 The NONE test state arises when a test class references types or methods introduced only by the fix
 485 patch; the class is absent from the compiled base commit, so the test runner cannot locate or execute
 486 it. This is distinct from FAIL, where the test compiles and runs but produces an incorrect result. A
 487 NONE→PASS (N2P) transition therefore indicates a newly introduced test that becomes executable
 488 only once the fix is applied. Repositories that follow a test-driven development workflow may exhibit
 489 exclusively N2P transitions with zero F2P counts, because their test patches introduce test classes
 490 that are entirely absent from the base commit.

491 A.7 Multi-Module Gradle Test Targeting

492 For multi-module Android repositories, we execute module-specific Gradle tasks (e.g.,
 493 `:feature:login:testDebugUnitTest`) rather than the top-level test task, targeting only mod-
 494 ules containing files touched by each PR’s test patch. Relevant modules are identified by mapping
 495 modified test files to their enclosing Gradle module via `settings.gradle` declarations. We default
 496 to the Debug build variant, which is the standard configuration used by the CI pipelines of the Android
 497 repositories in the benchmark.

498 B Manual Validation Details

499 Following SWE-bench-Verified [OpenAI, 2024], we perform human review of each instance that
 500 passes execution-based filtering. The review criteria are described in Section 3 (Phase 5). Figures 6
 501 and 7 visualise the score distributions and dataset composition, respectively.

502 B.1 Annotation Results

503 Among the 541 annotated candidates, 310 instances (57.3%) have unambiguously well-specified
 504 problem statements (score 0) and 150 (27.7%) have minor ambiguity but a clear interpretation
 505 (score 1). The remaining 81 (15.0%) are excluded by problem-statement score: 44 receive score 2
 506 (vague) and 37 score 3 (insufficient information).

507 Test-coverage adequacy shows a higher baseline: 389 instances (71.9%) receive score 0 and 61
 508 (11.3%) score 1. A total of 91 instances (16.8%) are excluded for inadequate test coverage, of which

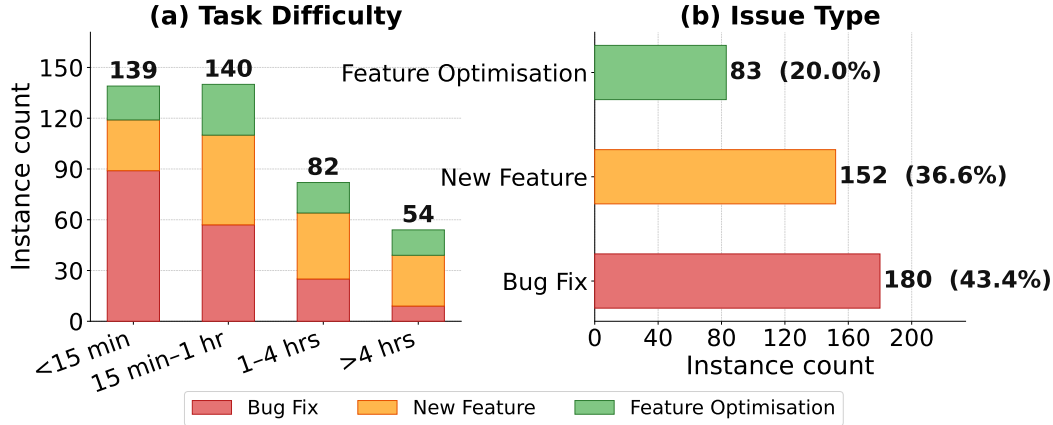


Figure 7: Task characteristics of the 415 verified benchmark instances. (a) Difficulty distribution stacked by issue type, using the four time-based difficulty labels assigned during annotation. (b) Overall issue-type breakdown. The benchmark spans the full difficulty range and covers all three major categories of mobile development activity.

509 36 score 2 and 55 score 3. Fifty-two instances (9.6%) are excluded for *both* reasons simultaneously,
 510 yielding 126 total exclusions (23.3%) and 415 verified benchmark instances (76.7%).

511 B.2 Task Characteristics

512 Figure 7 summarises the composition of the 415 verified instances by estimated resolution time and
 513 issue type. Difficulty is broadly distributed: 139 instances (33.5%) are estimated to take under 15
 514 minutes, 140 (33.7%) 15 minutes to one hour, 82 (19.8%) one to four hours, and 54 (13.0%) more
 515 than four hours. By issue type, 180 instances (43.4%) are Bug Fixes, 152 (36.6%) are New Features,
 516 and 83 (20.0%) are Feature Optimisations.

517 C Benchmark Comparison

518 C.1 Comparison with Existing Benchmarks

519 Table 1 situates MobileDev-Bench among existing issue-resolution benchmarks across four dimen-
 520 sions: domain coverage, verification, difficulty stratification, and patch complexity. Most prior
 521 benchmarks focus on library-style repositories, often in a single language or a small set of languages.
 522 Their gold patches are typically small, averaging between 1.2 and 4.9 modified files and between 14
 523 and 163 changed lines. In contrast, MobileDev-Bench targets production mobile apps and exhibits
 524 substantially larger fixes, averaging 13.5 modified files and 344.2 changed lines.

525 **Patch Complexity.** Figure 3 compares per-repository median patch size across five benchmarks
 526 using three metrics: files modified, hunks, and changed lines. SWE-PolyBench [Rashid et al., 2025],
 527 SWE-bench Multimodal [Yang et al., 2024b], and Multi-SWE-bench [Zan et al., 2025] occupy a
 528 compact low-to-moderate range, with median file counts of 2.6, 3.6, and 3.0, and median changed
 529 lines of 54, 128, and 73, respectively. Rust-SWE-bench [Xiang et al., 2026] is an outlier on the files
 530 axis (median 10.4 files) owing to the broad cross-module scope of typical Rust fixes, yet its median
 531 changed lines (92) remain well below those of MobileDev-Bench.

532 MobileDev-Bench presents the largest median hunk count (27.9) and changed line count (341.2)
 533 across all five benchmarks. Compared with SWE-PolyBench, the closest match on repository count,
 534 MobileDev-Bench fixes are $3.1\times$ larger in files, $3.9\times$ in hunks, and $6.3\times$ in changed lines. Relative to
 535 Multi-SWE-bench, the ratios are $2.7\times$, $3.6\times$, and $4.7\times$, whereas relative to SWE-bench Multimodal,
 536 they are $2.3\times$, $3.5\times$, and $2.7\times$. While Rust-SWE-bench has a slightly higher median file count
 537 than MobileDev-Bench (10.4 vs. 8.1), MobileDev-Bench exceeds it by $2.7\times$ in hunks and $3.7\times$ in
 538 changed lines, indicating that mobile app fixes involve more densely modified changes rather than

539 touching more files. This pattern reflects the coordinated nature of mobile app fixes, where a single
540 logical change often requires updates across source code, resources, and build configuration files.

541 **Design Dimensions and Scale.** MobileDev-Bench contains 415 manually verified instances from 19
542 repositories across four languages, comparable in scale to SWE-PolyBench Verified (382 instances,
543 20 repositories) and roughly four-fifths of SWE-bench Verified (500 instances, 12 repositories),
544 while targeting a distinct application domain. Beyond scale, it differs from prior benchmarks in
545 three design aspects. First, it covers production mobile apps spanning three frameworks (Android
546 Native, React Native, and Flutter) across four programming languages, a domain absent from all
547 existing benchmarks. Second, it combines human verification with difficulty stratification, enabling
548 controlled evaluation across task complexity levels. Third, it explicitly captures multi-artifact fixes,
549 where resolving an issue requires coordinated changes across heterogeneous project components
550 rather than isolated source edits. Issue descriptions are also longer on average (310 tokens), reflecting
551 the richer contextual information typical of mobile bug reports, such as reproduction steps, device
552 configurations, and platform versions.

553 C.2 Conceptual Differences from Backend Benchmarks

554 Existing benchmarks, such as SWE-bench [Jimenez et al., 2024] and its multimodal [Yang et al.,
555 2024b] and multilingual variants [Zan et al., 2025, Rashid et al., 2025, Kabir, 2025], Existing
556 benchmarks primarily target issue resolution in library-style repositories with relatively direct control
557 flow. These benchmarks often evaluate library-centric or backend repositories where control flow is
558 direct and execution environments are relatively straightforward (e.g., standard Python or JavaScript
559 runtimes).

560 **Execution and Lifecycle:** Mobile application frameworks change this control flow significantly. The
561 operating system dictates the lifecycle of an application, instantiating UI components (such as Activi-
562 ties in Android or Widgets in Flutter) and maintaining complex asynchronous states. MobileDev-
563 Bench tests an LLM’s ability to reason over these implicit callbacks and event-driven architectures
564 that are entirely absent in CLI or library benchmarks.

565 **Artifact Diversity:** Furthermore, a core distinction of MobileDev-Bench is its emphasis on multi-
566 artifact coordination. While conventional benchmarks evaluate patches applied strictly to source
567 code files, mobile fixes routinely demand coordinated changes across declarative user interface
568 definitions (e.g., XML layouts), localized string resources, manifestation configurations, and heavy
569 build scripts (Gradle). Models that perform well on algorithmic tasks may face additional challenges
570 when synchronization across artifacts is required between a deeply nested source logic and external
571 framework resources.

572 **Scale of Modifications:** As detailed in the comparison in Figure 3, the sheer scale of patches
573 in MobileDev-Bench far exceeds those typical of SWE-bench paradigms. The necessity to touch
574 multiple files reflects the cross-cutting operational reality of mobile deployments, reflecting a higher
575 degree of integration complexity than prior settings.

576 D Dataset Analysis

577 We provide additional quantitative analyses of the MobileDev-Bench dataset that complement the
578 summary statistics in Section 4.

579 D.1 Framework-wise Distribution and Code Change Characteristics

580 Table 8 reports per-framework instance counts and patch statistics. Android Native dominates
581 the dataset (335 instances, 80.7%), reflecting the larger number of mature open-source Android
582 repositories. Flutter contributes 63 instances (15.2%) from two high-activity repositories, and React
583 Native contributes 17 instances (4.1%) from four repositories. The fix patch statistics reveal that
584 Android Native instances involve substantially more code changes on average (14.5 modified files,
585 29.6 hunks) than Flutter (7.6 files, 18.3 hunks), reflecting the greater modularisation typical of large
586 Android codebases. Notably, Flutter test patches add considerably more lines on average (234.3) than
587 their fix counterparts (204.6), indicating that Flutter contributors tend to write comprehensive test
588 suites alongside their fixes.

Table 8: Framework-wise instance distribution and code change characteristics in MobileDev-Bench

Metric	Android Native	Flutter	React-Native
<i>Instance distribution</i>			
Repositories	13	2	4
Instances	335	63	17
<i>Fix patch statistics</i>			
Avg. # files	14.5	7.6	15.9
Avg. # hunks	29.6	18.3	41.3
Lines added	234.5	204.6	312.5
Lines removed	118.6	50.6	187.2
<i>Test patch statistics</i>			
Avg. # files	6.6	4.0	3.9
Avg. # hunks	16.2	13.9	15.1
Lines added	143.5	234.3	98.4
Lines removed	50.1	43.5	67.4

589 D.2 Distribution of Artifact Type Combinations

590 Table 9 reports the distribution of artifact type combinations present in fix patches. The majority
 591 of instances (57.3%) modify only source files, and 57.8% are single-artifact overall, indicating that
 592 the primary challenge is code comprehension and logic repair. However, 42.2% of instances require
 593 coordinated changes across multiple artifact types, testing a model’s ability to reason about the full
 594 mobile-development toolchain. The most common multi-artifact pattern is Source + Build (8.2%),
 595 where fixes also modify Gradle build scripts, dependency version catalogues, or build configuration
 596 files. Android Native accounts for most multi-artifact instances: patterns combining source with
 597 resources, internationalisation string files, or build scripts are exclusive to or heavily concentrated in
 598 Android projects, whereas Flutter instances are predominantly source-only (74.6%).

599 D.3 Background Knowledge Required by Framework

600 Figure 8 characterises the technical knowledge domains implicitly required by each framework’s
 601 instances, inferred from import statements, annotations, and identifiers in the added lines of fix patches.
 602 The six domains span the main layers of modern mobile engineering: asynchronous programming,
 603 declarative UI frameworks, architecture components, dependency injection, platform and native APIs,
 604 and networking.

605 For Android Native, Architecture Components are the most frequently required domain (36.4%),
 606 reflecting widespread use of `ViewModel`, `Room`, `WorkManager`, and `Navigation` across the benchmark
 607 repositories. Declarative UI knowledge (Jetpack Compose) is needed in 33.7% of Android Native
 608 instances, Asynchronous/Concurrency patterns (Kotlin coroutines, `Flow`) in 30.4%, Platform SDK
 609 APIs in 17.3%, and Dependency Injection via `Hilt` or `Dagger` in 17.3%. Together, these figures
 610 confirm that Android Native instances stress a broad stack of modern Android engineering beyond
 611 core Kotlin syntax.

612 Flutter instances primarily demand Widget and UI framework knowledge (39.7%) alongside asyn-
 613 chronous Dart patterns (31.7%), consistent with Flutter’s widget-centric programming model. Archi-
 614 tecture and state-management patterns (BLoC, `Provider`) are needed in only 6.3% of Flutter instances,
 615 reflecting the dataset’s focus on high-activity bug-fix PRs rather than large architectural changes.

616 React Native percentages are indicative only ($n = 17$): UI framework patterns appear in 41.2%
 617 of instances and Async/Concurrency in 23.5%, with Architecture Components in 17.6% and Plat-
 618 form/Native APIs in 11.8%; Dependency Injection and Networking domains register zero, plausibly
 619 reflecting these repositories’ reliance on abstracted, library-level patterns not captured by source-level
 620 identifiers alone.

Table 9: Distribution of artifact type combinations in MobileDev-Bench. 42.17% of instances require coordinated changes across multiple artifact types. Framework columns show the % of instances within each framework exhibiting that combination (AN = Android Native, FI = Flutter, RN = React Native).

Combination	Overall		% by Framework		
	#	%	AN	FI	RN
<i>Single Artifact Type</i>					
Src	238	57.35	54.62	74.60	47.06
Res	1	0.24	0.30	0.00	0.00
Build	1	0.24	0.30	0.00	0.00
<i>Subtotal</i>	<i>240</i>	<i>57.83</i>	<i>55.22</i>	<i>74.60</i>	<i>47.06</i>
<i>Two Artifact Types</i>					
Src + Build	34	8.19	10.15	0.00	0.00
Src + Res	15	3.62	4.48	0.00	0.00
Src + i18n	15	3.62	0.00	20.63	11.77
Src + Docs	14	3.37	4.18	0.00	0.00
Src + Other	4	0.96	0.89	0.00	5.88
Other combinations	2	0.48	0.30	1.59	0.00
<i>Subtotal</i>	<i>84</i>	<i>20.24</i>	<i>20.00</i>	<i>22.22</i>	<i>17.65</i>
<i>Three or More Artifact Types</i>					
Src + Res + i18n	27	6.51	8.06	0.00	0.00
Src + Res + i18n + Build	11	2.65	3.29	0.00	0.00
Src + Build + Docs	7	1.69	1.79	1.59	0.00
Src + Manifest + Config	5	1.21	0.00	1.59	23.53
Src + Res + i18n + Docs	5	1.20	1.49	0.00	0.00
Other combinations	36	8.67	10.15	0.00	11.76
<i>Subtotal</i>	<i>91</i>	<i>21.93</i>	<i>24.78</i>	<i>3.18</i>	<i>35.29</i>
Total (Single)	240	57.83	55.22	74.60	47.06
Total (Multi)	175	42.17	44.78	25.40	52.94

621 E Experimental Setup

622 This section provides detailed information about the experimental framework, model selection, and
623 evaluation metrics used to assess automated program repair performance on MobileDev-Bench.

624 E.1 Evaluation Framework

625 We evaluate MobileDev-Bench using the localization and repair pipeline from Agentless [Xia et al.,
626 2024], a state-of-the-art autonomous code repair framework that has demonstrated strong performance
627 on Python-centric benchmarks. Agentless decomposes the repair task into two sequential phases:

- 628 1. **Fault Localization:** Given an issue description and repository context, the model identifies
629 potentially buggy files and functions. This phase employs retrieval-augmented generation to
630 narrow the search space from the entire codebase to a manageable set of candidate locations.
- 631 2. **Patch Generation:** The model generates a patch as a unified diff format targeting the
632 localized code elements. The generated patch must be syntactically valid and applicable to
633 the base commit.

634 This two-phase decomposition enables us to separately measure localization capability (can the model
635 identify where to fix?) and repair capability (can the model generate the correct fix?), providing
636 diagnostic insights into failure modes.

637 We validate generated patches using a Docker-based test harness. For each task instance, we: (1) clone
638 the repository at the base commit, (2) apply the model-generated patch, (3) execute the repository
639 test suite in an isolated Docker container, and (4) compare test outcomes against the ground truth. A
640 patch is considered successful if it passes all tests that pass with the developer’s ground truth patch,
641 including previously failing tests that should now pass.

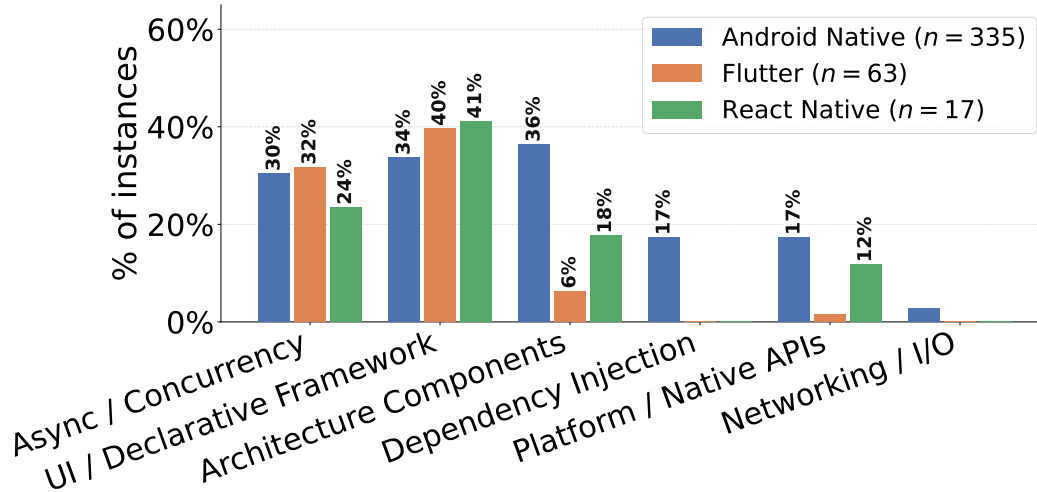


Figure 8: Percentage of instances per framework whose fix patches involve each technical knowledge domain, inferred from import statements, annotations, and identifiers in the added lines of fix patches. An instance may span multiple domains.

642 This evaluation approach mirrors the workflow used in SWE-bench [Jimenez et al., 2024] and Multi-
 643 SWE-bench [Zan et al., 2025], enabling direct comparison with existing benchmarks while avoiding
 644 the compounding errors and non-determinism of full agentic systems that include web browsing,
 645 terminal access, and multi-turn interactions.

646 E.2 Adaptation for Mobile Development Languages

647 The original Agentless implementation relies on Python’s Abstract Syntax Tree (AST) module for
 648 parsing and extracting code structures. To support the diverse programming languages in MobileDev-
 649 Bench (Java, Kotlin, TypeScript, and Dart), we incorporate tree-sitter¹, a parsing library that provides
 650 unified syntax tree generation across multiple languages.

651 Similar to the MagentLess² adaptation for Multi-SWE-bench, we replace Python-specific parsing
 652 calls with tree-sitter-based parsing while maintaining the same localization and patch generation
 653 algorithms. This ensures that:

- 654 • Code structure extraction (classes, methods, functions) works consistently across Android
 655 (Java/Kotlin), Flutter (Dart), and React Native (TypeScript/JavaScript) codebases.
- 656 • Function-level localization operates with the same granularity across all four languages.
- 657 • The core Agentless algorithm remains unchanged, ensuring fair comparison with results on
 658 Python-based benchmarks.

659 For non-code artifacts (XML layouts, JSON i18n files, Gradle build scripts), which lack traditional
 660 function boundaries, we perform file-level localization only.

661 E.3 Model Selection

662 We evaluate four frontier language models: Claude Sonnet 4.5 [Anthropic, 2025] (Anthropic’s
 663 proprietary model with 1M token context window), GPT-5.2 [OpenAI, 2025] (OpenAI’s frontier
 664 model with enhanced reasoning capabilities), Gemini 2.5 Flash [Comanici et al., 2025] (Google’s
 665 fast-inference model with native multimodal support), and Qwen-3-Coder [Qwen Team, 2025] (an
 666 open-weight model from Alibaba Cloud, fine-tuned for code tasks).

667 All models are accessed via OpenRouter³.

¹<https://pypi.org/project/tree-sitter-language-pack/>

²<https://github.com/multi-swe-bench/MagentLess>

³<https://openrouter.ai/>

Table 10: Resolution rates by ground truth file count under the Agentless and Oracle evaluation settings. Rate is resolved / evaluated instances (empty-patch predictions excluded from the denominator in the Oracle setting; instances without a submitted prediction excluded in Agentless). Column abbreviations: C=Claude Sonnet 4.5, Q=Qwen3 Coder, G=GPT 5.2, Ge=Gemini 2.5 Flash.

#Files	Tasks	Resolved				Rate (%)			
		C	Q	G	Ge	C	Q	G	Ge
<i>Agentless</i>									
1 file	71	11	9	11	9	15.5	12.7	15.5	12.7
2-3 files	89	4	2	2	2	4.5	2.2	2.2	2.2
4-5 files	52	2	2	0	2	3.8	3.8	0.0	3.8
6-10 files	70	0	0	0	0	0.0	0.0	0.0	0.0
11+ files	133	0	4	3	4	0.0	3.0	2.3	3.0
Overall	415	17	17	16	17	4.15	4.28	3.90	4.16
<i>Oracle</i>									
1 file	71	18	9	14	7	25.4	12.7	19.7	9.9
2-3 files	89	3	1	1	0	3.4	1.1	1.1	0.0
4-5 files	52	2	0	1	1	3.8	0.0	1.9	1.9
6-10 files	70	0	0	2	0	0.0	0.0	3.0	0.0
11+ files	133	2	0	3	0	1.5	0.0	2.4	0.0
Overall	415	25	10	21	8	6.07	2.52	5.19	1.94

668 **Agentless setting.** We use temperature 0.0 for fault localization. For patch generation, we generate one greedy sample (temperature 0.0) and multiple diverse samples (temperature 0.8), then rank and select the best patch. To ensure fair comparison, we disable extended reasoning modes: reasoning_effort=low for GPT-5.2, no extended thinking for Claude Sonnet 4.5, and default settings for Gemini 2.5 Flash and Qwen-3-Coder.

673 **Oracle setting.** The fault localization step is bypassed; models are given the exact set of files modified in the ground truth patch. Patch generation uses a single greedy sample (temperature 0.0). No reasoning mode overrides are applied.

676 E.4 Evaluation Metrics

677 Following prior work on automated program repair [Xia et al., 2024, Rashid et al., 2025], we report two complementary families of metrics that separately assess end-to-end repair capability and fault localization accuracy.

680 E.4.1 Resolution Rate

681 Let T denote the total number of tasks in the benchmark. For each task t , we define an indicator function $\mathbf{1}(\text{Pass}_t)$ that equals 1 if the model-generated patch passes all tests when applied to the base commit, and 0 otherwise. We compute Resolution Rate as:

$$\text{Resolution Rate} = \frac{1}{T} \sum_{t=1}^T \mathbf{1}(\text{Pass}_t)$$

684 This metric measures end-to-end repair capability at the task level. A high Resolution Rate indicates that the model successfully localizes the bug, generates a syntactically correct patch, and produces semantically correct fixes that satisfy all test requirements.

687 Resolution Rate is the primary metric for comparing model performance, as it directly measures the practical utility of automated repair systems. However, it conflates localization and generation capabilities, providing limited diagnostic insight when rates are low.

690 **E.4.2 Retrieval Metrics: Precision, Recall, and F1**

691 To isolate fault localization performance from patch generation quality, we compute file-level retrieval
692 metrics. For each task t , let F_t^{GT} denote the set of ground truth modified files (extracted from the
693 developer’s patch) and let F_t^{Pred} denote the set of files modified by the model-generated patch.

694 We compute per-task retrieval metrics as:

$$\begin{aligned} \text{Recall}_t &= \frac{|F_t^{GT} \cap F_t^{Pred}|}{|F_t^{GT}|} \\ \text{Precision}_t &= \frac{|F_t^{GT} \cap F_t^{Pred}|}{|F_t^{Pred}|} \\ \text{F1}_t &= \frac{2 \cdot \text{Precision}_t \cdot \text{Recall}_t}{\text{Precision}_t + \text{Recall}_t} \end{aligned}$$

695 We then report macro-averaged metrics across all T tasks:

$$\begin{aligned} \text{Recall} &= \frac{1}{T} \sum_{t=1}^T \text{Recall}_t \\ \text{Precision} &= \frac{1}{T} \sum_{t=1}^T \text{Precision}_t \\ \text{F1} &= \frac{1}{T} \sum_{t=1}^T \text{F1}_t \end{aligned}$$

696 These metrics measure fault localization performance independently of patch correctness:

- 697 • **Recall** captures the fraction of ground truth files that the model correctly identifies as
698 requiring modification. Low recall indicates the model misses many necessary modification
699 sites.
- 700 • **Precision** measures the fraction of predicted files that are actually relevant (i.e., appear in
701 the ground truth). High precision indicates the model avoids editing irrelevant files, while
702 low precision suggests over-editing or incorrect localization.
- 703 • **F1** provides a harmonic mean balancing both aspects, penalizing models that achieve high
704 recall through indiscriminate file editing or high precision through overly conservative
705 predictions.

706 Together with Resolution Rate, these metrics distinguish localization capability from repair capability.
707 High recall with low Resolution Rate suggests that the model identifies the correct files but fails to
708 generate correct fixes. Conversely, low recall limits achievable Resolution Rate, since a model cannot
709 repair files it fails to localize. High precision indicates the model avoids over-editing irrelevant files,
710 which is important for minimizing unintended side effects in production deployments.

711 **E.5 Evaluation Infrastructure**

712 Each task instance is executed in an isolated Docker container with repository-specific dependencies
713 installed from the original project’s build configuration (Gradle for Android, pub for Flutter, npm/yarn
714 for React Native).

715 Test execution timeouts are set to 30 minutes per task to accommodate large test suites in production
716 repositories. Tasks that exceed this limit are marked as timeouts and counted as failures.

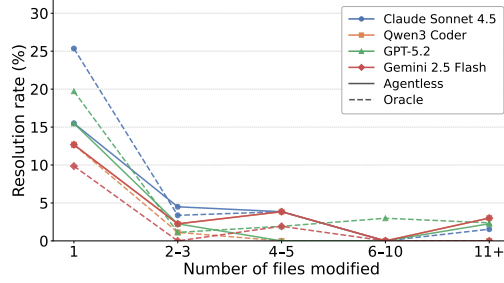


Figure 9: Resolution rates by number of files modified, for Agentless (solid lines) and Oracle (dashed lines) settings. Single-file tasks achieve 12.7–15.5% under Agentless and up to 25.4% under Oracle, dropping to near zero for 6+ files.

717 F Evaluation Results

718 This section presents comprehensive evaluation results for four frontier language models on
 719 MobileDev-Bench. We analyze resolution rates, fault localization performance, and the relationship
 720 between task characteristics and repair success.

721 F.1 Resolution Rate Analysis

722 Table 10 shows resolution rates for both evaluation settings on 415 verified tasks (rate = resolved /
 723 evaluated, excluding empty-patch predictions from the denominator). Under the Agentless setting,
 724 rates range from 3.90% to 4.28% (16–17 resolved per model, 397–410 evaluated). Under the Oracle
 725 setting, rates range from 1.94% (Gemini 2.5 Flash, 8/412) to 6.07% (Claude Sonnet 4.5, 25/412).
 726 Tasks in MobileDev-Bench require modifying an average of 13.5 files (median: 5), reflecting the
 727 multi-file nature of mobile development.

728 Multi-File Modifications Show Lower Resolution Rates.

729 Figure 9 reveals a dramatic degradation in resolution rates as the number of required file modifications
 730 increases. Under the Agentless setting, single-file tasks (71 tasks) achieve 12.7–15.5% resolution
 731 rates, substantially higher than the 3.9–4.3% overall rates. Under the Oracle setting, single-file rates
 732 rise further to 9.9–25.4%, reflecting that localization is the dominant bottleneck on simpler tasks.
 733 However, success rapidly diminishes with task complexity under both settings.

734 Tasks requiring 2–3 file modifications show a steep drop to 2.2–4.5% under Agentless and 0.0–
 735 3.4% under Oracle. Mid-complexity tasks requiring 4–10 files prove nearly unsolvable (0.0–3.8%
 736 Agentless; 0.0–3.8% Oracle). This difficulty cliff highlights models’ fundamental limitations in
 737 coordinating changes across multiple files.

738 Tasks requiring 11+ files show modest recovery under Agentless (0.0–3.0%), with Qwen3 Coder and
 739 Gemini 2.5 Flash each resolving 4 of 133 such instances. However, detailed analysis (Section F.3)
 740 reveals this recovery stems from repository-specific architectural patterns rather than improved multi-
 741 file reasoning capabilities. Specifically, the majority of resolved 11+ file instances come from a
 742 single repository (element-hq/element-x-android), suggesting that certain codebases facilitate
 743 large-scale automated changes through consistent structure and comprehensive test coverage.

744 Across all models, only 32 unique instances were successfully resolved under Agentless, accumulating
 745 to 67 total resolutions due to overlap. The remaining 383 tasks (92.3%) were not resolved by
 746 any Agentless model, underscoring the fundamental challenge of multi-file reasoning in mobile
 747 development contexts.

748 Larger Patches Correlate with Lower Success Rates.

749 Patch complexity, measured by lines changed in the ground truth fix, strongly predicts resolution
 750 success (Figure 10). Under Agentless, tiny patches (≤ 10 lines) achieve 14.3–16.1% resolution rates,
 751 while very large patches (> 500 lines) drop to 0.0–3.9%. Under Oracle, rates for ≤ 10 -line patches

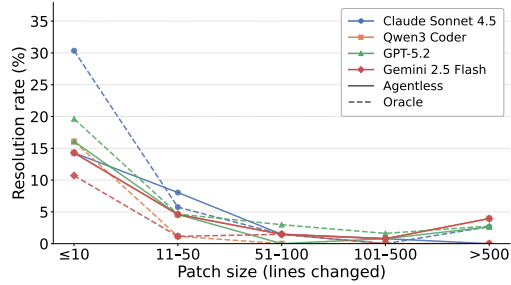


Figure 10: Resolution rates by patch size, for Agentless (solid lines) and Oracle (dashed lines) settings. Tiny patches (≤ 10 lines) achieve 14.3–16.1% under Agentless and 10.7–30.4% under Oracle, dropping to near zero for patches > 100 lines.

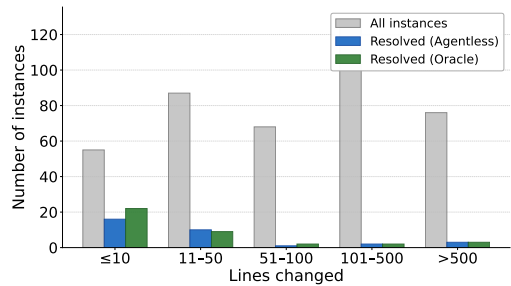


Figure 11: Patch size distribution for all instances versus resolved instances. Resolved instances skew heavily toward patches ≤ 50 lines (81% vs. 34%).

752 reach 10.7–30.4%, and > 500 -line patches remain low at 0.0–2.7%. This inverse relationship reveals
 753 that models struggle not only with multi-file coordination but also with generating extensive code
 754 changes within identified locations.

755 Resolved instances exhibit dramatically smaller patches than unresolved tasks. The median patch
 756 size for resolved instances is 9.5 lines versus 109.0 lines for unresolved instances (11.5x difference).
 757 Mean patch sizes show a similar pattern: 102.6 lines for resolved versus 364.4 lines for unresolved
 758 (0.28x ratio). This 72% size reduction indicates that successful repairs target localized, focused
 759 changes rather than broad refactorings.

760 Patch size distributions further emphasize this disparity (Figure 11). Small patches (≤ 50 lines)
 761 comprise 81% of resolutions but only 34% of the benchmark, a 2.4x overrepresentation. Conversely,
 762 patches > 50 lines represent 66% of the benchmark but only 19% of resolutions. Tiny patches
 763 (≤ 10 lines) show the most extreme skew: 53.1% of resolutions versus 13.5% of benchmark (3.9x
 764 overrepresentation). Resolution rates collapse at medium patch sizes (51–100 lines: 0.0–1.5%) and
 765 remain near zero for larger patches (101–500 lines: 0.8–0.8%; > 500 lines: 0.0–3.9%).

766 This finding complements the file count analysis, revealing two orthogonal dimensions of repair
 767 difficulty: spatial distribution (how many files) and change magnitude (how many lines). While
 768 file count measures coordination complexity, patch size measures generation complexity. Large
 769 patches require models to reason about extensive code interactions, maintain consistency across many
 770 modified lines, and ensure all changes compile and pass tests. Current models excel at small, surgical
 771 fixes but fail when extensive code generation is required.

772 F.2 Characteristics of Resolved Instances

773 Table 11 analyzes how resolved instances differ from the overall benchmark distribution across three
 774 dimensions: artifact type diversity, file count distribution, and framework. This analysis reveals
 775 structural factors that influence repair success.

Table 11: Resolution rate breakdown by task characteristics (415 tasks). Resolved column shows the share of uniquely resolved instances (32 Agentless, 38 Oracle) in each category vs. the overall benchmark share. Ratio = Resolved% / Benchmark%.

Category	Bench.	Agentless		Oracle	
		Resolved	Ratio	Resolved	Ratio
<i>Artifact Type Diversity</i>					
Single-artifact	57.8%	84.4%	1.46x	86.8%	1.50x
Multi-artifact	42.2%	15.6%	0.37x	13.2%	0.31x
<i>File Count Distribution</i>					
1 file	17.1%	65.6%	3.84x	68.4%	4.00x
2–3 files	21.4%	15.6%	0.73x	10.5%	0.49x
4–5 files	12.5%	6.3%	0.50x	5.3%	0.42x
6–10 files	16.9%	0.0%	0.00x	5.3%	0.31x
11+ files	32.1%	12.5%	0.39x	10.5%	0.33x
<i>Framework Distribution</i>					
Android Native	80.7%	71.9%	0.89x	81.6%	1.01x
Flutter	15.2%	28.1%	1.85x	13.1%	0.86x
React Native	4.1%	0.0%	0.00x	5.3%	1.29x

776 **Artifact Type Diversity.** We classify tasks as single-artifact (all files of one type) or multi-artifact
777 (spanning source code, resources, build files, manifests). While 57.8% of benchmark tasks are single-
778 artifact, 84.4% of Agentless resolutions come from single-artifact tasks (1.46x overrepresentation).
779 Multi-artifact tasks show 0.37x underrepresentation (15.6% of resolutions vs. 42.2% of benchmark),
780 highlighting difficulty coordinating changes across heterogeneous file types.

781 **File Count Distribution.** Single-file tasks dominate Agentless resolutions (65.6% of resolutions vs.
782 17.1% of benchmark, 3.84x overrepresentation), while multi-file tasks show severe underrepresenta-
783 tion. Tasks requiring 6–10 files exhibit 0.0% of Agentless resolutions versus 16.9% of benchmark
784 (0.00x ratio), the steepest underrepresentation. The 11+ file category shows modest recovery (0.39x
785 ratio), driven by repository-specific architectural regularity (Section F.3).

786 **Framework Distribution.** Framework distribution among resolved tasks broadly follows benchmark
787 composition but with notable deviations (Table 11). Android Native tasks are near-proportionally
788 represented (71.9% of Agentless resolutions vs. 80.7% of benchmark, 0.89x), while Flutter tasks are
789 overrepresented at 28.1% of resolutions despite comprising only 15.2% of the benchmark (1.85x),
790 likely reflecting the structural characteristics of Flutter repositories, which tend toward single-artifact
791 changes with strong semantic signals in issue descriptions. React Native tasks see complete failure
792 (0.0% vs. 4.1%), mirroring patterns seen across complex multi-layer architectures where cross-file
793 dependency tracing breaks down. These differences reflect repository architecture and framework
794 conventions rather than underlying language difficulty.

795 F.3 Analysis of 11+ File Resolved Instances

796 The modest recovery in resolution rates for 11+ file tasks (0.0–3.0%) stems from the concentration of
797 resolved instances in a single repository rather than general model capability. All 4 unique resolved
798 instances with 11+ files under the Agentless setting come from `element-hq/element-x-android`,
799 whose modular architecture, consistent naming conventions, and comprehensive test coverage facili-
800 tate large-scale automated repair.

801 These instances share structural characteristics: averaging 33.5 files across 21.5 directories, predomi-
802 nantly multi-artifact changes (3/4), and 71% source code modifications. This finding indicates that
803 repository architecture, not task scale alone, determines repair difficulty.

804 F.4 Fault Localization Performance

805 Table 12 shows that fault localization is the primary bottleneck. Overall retrieval performance
806 ranges from 13.5–18.3% recall, 42.5–57.5% precision, and 17.3–23.0% F1, with Claude Sonnet
807 4.5 achieving the highest scores. The large precision-recall gap (3x across all models) indicates
808 conservative file selection: models identify relevant files but miss many necessary modifications.

Table 12: Retrieval metrics by ground truth file count for Agentless and Oracle settings. All metrics are percentages, averaged across tasks in each file count bin. Column abbreviations: C=Claude Sonnet 4.5, Q=Qwen3 Coder, G=GPT 5.2, Ge=Gemini 2.5 Flash.

#Files	Tasks	Recall (%)				Precision (%)				F1 (%)			
		C	Q	G	Ge	C	Q	G	Ge	C	Q	G	Ge
<i>Agentless</i>													
1 file	71	54.9	46.5	35.2	35.2	54.9	46.5	35.2	35.2	54.9	46.5	35.2	35.2
2–3 files	89	22.7	23.6	20.4	20.2	49.4	48.3	44.4	44.9	30.5	30.9	27.5	27.5
4–5 files	52	12.1	12.7	9.9	10.5	51.9	53.8	40.4	42.3	19.6	20.5	15.8	16.7
6–10 files	70	9.4	7.5	7.4	6.7	65.7	52.1	49.3	46.4	16.4	13.0	12.8	11.6
11+ files	133	3.0	2.9	2.1	2.3	62.0	51.5	43.2	42.9	5.6	5.3	3.9	4.2
Overall	415	18.3	16.8	13.6	13.5	57.5	50.4	42.8	42.5	23.0	21.0	17.3	17.3
<i>Oracle</i>													
1 file	71	95.8	97.2	84.5	100.0	95.8	96.5	83.8	100.0	95.8	96.7	84.0	100.0
2–3 files	89	57.5	50.2	40.6	60.9	94.9	98.5	89.9	96.5	68.9	64.4	54.9	71.6
4–5 files	52	48.8	35.9	21.7	50.2	89.8	93.2	89.2	92.6	58.1	48.0	34.4	60.7
6–10 files	70	43.1	30.4	14.2	43.2	90.8	89.9	87.1	95.8	53.6	41.5	23.2	53.6
11+ files	133	23.9	14.2	8.1	35.7	80.9	78.5	77.1	89.5	32.3	20.8	11.7	43.8
Overall	415	49.8	41.6	30.9	55.2	89.2	89.6	84.2	94.3	57.8	50.0	38.1	63.1

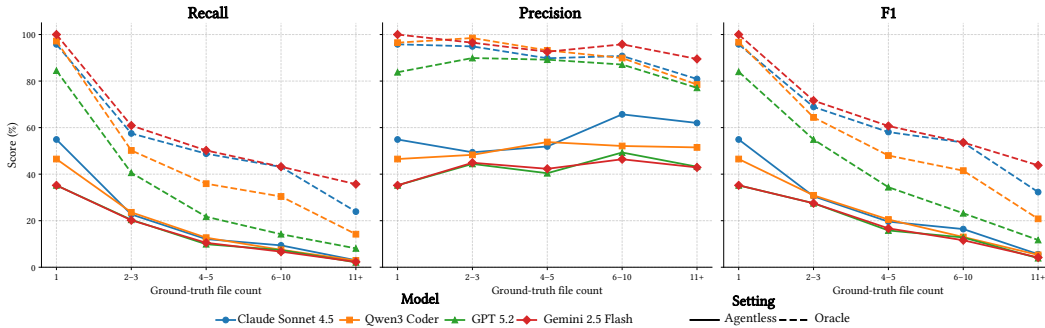


Figure 12: Retrieval metrics (recall, precision, F1) by ground-truth file count across all four models. Recall collapses from 35–55% on single-file tasks to 2–3% on 11+ file tasks, while precision remains comparatively stable (35–66%), widening the precision-recall gap with complexity. F1 degrades sharply, driven entirely by recall failure.

809 **Recall Degradation by File Count.** Recall degrades severely with file count: from 35–55% on
 810 single-file tasks to 2–3% on 11+ file tasks, representing a 33–53 percentage point decline. For tasks
 811 requiring 6–10 files, models identify fewer than 10% of necessary files. This catastrophic degradation
 812 reflects models’ inability to trace cross-layer dependencies in mobile architectures, where changes
 813 cascade across network, data, presentation, view, and test layers.

814 **Precision Trends.** Precision increases with file count (35–55% for single-file tasks to 46–66% for
 815 6–10 file tasks), indicating conservative behavior on complex changes. Localization errors stem
 816 primarily from missed files rather than false positives: models identify a core subset correctly but fail
 817 to expand to all necessary modifications.

818 **Localization Patterns.** Three distinct patterns emerge:

- 819 • **Perfect localization** (5–8% of tasks): Almost exclusively single-file changes with strong
 820 semantic signals in issue descriptions.
- 821 • **Zero localization** (46–59% of tasks): Models edit entirely incorrect files, common when
 822 ground truth modifies non-code artifacts or peripheral modules. 162 tasks (39.0%) show
 823 zero recall across all models.

824 • **Partial localization** (36–47% of tasks): Models find some relevant files but miss cascading
825 modifications across architectural layers.

826 **F.4.1 Implications**

827 Low recall imposes an upper bound on resolution rates: even perfect patch generation would achieve
828 only 13–18% resolution (matching the recall ceiling). Observed rates (3.90%–4.28% under Agentless)
829 are lower still, indicating localization dominates repair difficulty. This inverts traditional assumptions
830 in automated program repair, where patch generation is typically the primary challenge.

831 Improving performance requires advances in multi-file fault localization, including dependency-aware
832 retrieval, architectural priors for mobile app patterns (MVVM, repository pattern), and specialized
833 strategies for non-code artifacts underrepresented in embedding spaces.

834 **G Discussion**

835 Our evaluation reveals overall resolution rates of 3.90%–4.28% (Agentless) and 1.94%–6.07%
836 (Oracle) on MobileDev-Bench. Three architectural properties of mobile development explain this low
837 performance. First, *multi-artifact coordination*: models achieve 1.46x success on single-artifact tasks
838 but only 0.37x on multi-artifact tasks requiring synchronization across source code, resources, build
839 configurations, and manifests. Second, *framework-mediated execution*: platform lifecycles obscure
840 control flow from static analysis. Third, *cross-module dependencies*: recall drops from 35–55% on
841 single-file tasks to 7–9% on multi-file tasks (6–10 files), and further to 2–3% on large-scale tasks
842 (11+ files).

843 Localization is the primary bottleneck. File-level recall of 13.5–18.3% imposes a hard ceiling on
844 resolution rates. For mobile development, the challenge lies not in generating correct patches but in
845 identifying which files need modification. This finding suggests research should prioritize multi-file
846 dependency tracing, artifact-aware retrieval, and framework-specific reasoning over patch generation
847 improvements.

848 The non-monotonic relationship between file count and resolution rate provides insight: single-file
849 tasks achieve 13 to 16% resolution, mid-complexity tasks (4 to 10 files) drop to 0 to 3.8%, while
850 11+ file tasks show modest recovery (0.0 to 3.0%). This pattern suggests models struggle with
851 *architectural complexity* (interconnected changes across heterogeneous components) rather than
852 raw file count. Repository-specific phenomena hint that certain architectural patterns may facilitate
853 automated repair.

854 **H Future Work**

855 Our findings suggest three high-priority research directions. First, given that localization is the primary
856 bottleneck, techniques combining semantic search with static analysis for multi-file dependency
857 tracing warrant immediate investigation. Second, repository-specific success patterns hint that
858 certain architectural characteristics aid automated repair; systematically analyzing what makes
859 codebases "LLM-friendly" could inform both tool design and best practices. Third, mobile issues
860 often include screenshots and crash logs; evaluating multi-modal models' ability to ground repairs in
861 visual artifacts may improve UI-related bug localization. Extending to iOS, developing framework-
862 aware representations encoding lifecycle semantics, and testing iterative agentic workflows provide
863 additional promising directions.

864 **I Limitations**

865 MobileDev-Bench focuses on Android Native and cross-platform apps, including React-Native
866 and Flutter repositories, and does not include iOS-native projects written in Swift or Objective-C.
867 This exclusion is a deliberate reproducibility constraint rather than a coverage gap. iOS builds
868 are tightly coupled to Xcode and macOS: the Xcode and Apple SDKs Agreement [Apple Inc.,
869 2024] explicitly prohibits execution on non-Apple-branded hardware, and macOS itself cannot be
870 legally run in standard Linux container environments. Consequently, the containerized build-and-test
871 execution that is central to our evaluation design is not achievable for iOS repositories without

872 proprietary infrastructure. Extending MobileDev-Bench to iOS would require macOS-based CI
873 runners or hardware-in-the-loop setups, both of which introduce environment heterogeneity that
874 undermines the reproducibility guarantees the benchmark is designed to provide. Evaluation relies
875 on execution-based validation using repository test suites, which may not fully capture issues with
876 insufficient or incomplete test coverage. Additionally, the benchmark is constructed from issue-pull
877 request pairs in open-source repositories, which may not represent all mobile development activities,
878 such as UI design iterations or refactoring tasks without associated issues. Finally, our evaluation
879 uses a single-pass repair pipeline, whereas app development often involves iterative debugging and
880 interactive workflows.

881 **J Potential Societal Impact**

882 MobileDev-Bench provides an execution-validated benchmark for mobile-specific program repair,
883 exposing failure modes invisible to library-centric evaluations, guiding practitioners on where LLM
884 assistants are and are not reliable in mobile CI pipelines, and highlighting the need for human over-
885 sight when deploying model-generated patches in accessibility-, localisation-, or security-sensitive
886 components.

887 **AI Assistance.** AI-based writing tools were used only for minor grammatical and wording improve-
888 ments. All technical content and analyses were produced and verified by the authors.

889 **NeurIPS Paper Checklist**

890 **1. Claims**

891 Question: Do the main claims made in the abstract and introduction accurately reflect the
892 paper’s contributions and scope?

893 Answer: [Yes]

894 Justification: The abstract and introduction state three claims that are directly supported by
895 the paper: (1) MobileDev-Bench is a benchmark of 415 verified mobile app issue-resolution
896 tasks (Section 3, Table 2); (2) the evaluation framework is compilation-aware and execution-
897 based across four languages (Section 3.4, Appendix E); and (3) four frontier LLMs achieve
898 low resolution rates of 3.90–6.07%, with fault localization as the dominant bottleneck
899 (Section 6, Table 4, Figure 5). No aspirational claims are presented as achieved results.

900 Guidelines:

- 901 • The answer [N/A] means that the abstract and introduction do not include the claims
902 made in the paper.
- 903 • The abstract and/or introduction should clearly state the claims made, including the
904 contributions made in the paper and important assumptions and limitations. A [No] or
905 [N/A] answer to this question will not be perceived well by the reviewers.
- 906 • The claims made should match theoretical and experimental results, and reflect how
907 much the results can be expected to generalize to other settings.
- 908 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
909 are not attained by the paper.

910 **2. Limitations**

911 Question: Does the paper discuss the limitations of the work performed by the authors?

912 Answer: [Yes]

913 Justification: Appendix I discusses limitations explicitly, including the exclusion of iOS-
914 native repositories (Swift/Objective-C) due to Apple Inc. licensing constraints that prevent
915 containerized execution, reliance on repository test suites that may not fully cover fix
916 correctness, the restriction to open-source issue–PR pairs which may not represent propri-
917 etary workflows, and the use of a single-pass repair pipeline that does not model iterative
918 debugging.

919 Guidelines:

- 920 • The answer [N/A] means that the paper has no limitation while the answer [No] means
921 that the paper has limitations, but those are not discussed in the paper.
- 922 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 923 • The paper should point out any strong assumptions and how robust the results are to
924 violations of these assumptions (e.g., independence assumptions, noiseless settings,
925 model well-specification, asymptotic approximations only holding locally). The authors
926 should reflect on how these assumptions might be violated in practice and what the
927 implications would be.
- 928 • The authors should reflect on the scope of the claims made, e.g., if the approach was
929 only tested on a few datasets or with a few runs. In general, empirical results often
930 depend on implicit assumptions, which should be articulated.
- 931 • The authors should reflect on the factors that influence the performance of the approach.
932 For example, a facial recognition algorithm may perform poorly when image resolution
933 is low or images are taken in low lighting. Or a speech-to-text system might not be
934 used reliably to provide closed captions for online lectures because it fails to handle
935 technical jargon.
- 936 • The authors should discuss the computational efficiency of the proposed algorithms
937 and how they scale with dataset size.
- 938 • If applicable, the authors should discuss possible limitations of their approach to
939 address problems of privacy and fairness.

940 • While the authors might fear that complete honesty about limitations might be used by
941 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
942 limitations that aren't acknowledged in the paper. The authors should use their best
943 judgment and recognize that individual actions in favor of transparency play an impor-
944 tant role in developing norms that preserve the integrity of the community. Reviewers
945 will be specifically instructed to not penalize honesty concerning limitations.

946 3. Theory assumptions and proofs

947 Question: For each theoretical result, does the paper provide the full set of assumptions and
948 a complete (and correct) proof?

949 Answer: [No]

950 Justification: The paper makes no theoretical claims and presents no theorems, lemmas, or
951 proofs. All contributions are empirical.

952 Guidelines:

- 953 • The answer [N/A] means that the paper does not include theoretical results.
- 954 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
955 referenced.
- 956 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 957 • The proofs can either appear in the main paper or the supplemental material, but if
958 they appear in the supplemental material, the authors are encouraged to provide a short
959 proof sketch to provide intuition.
- 960 • Inversely, any informal proof provided in the core of the paper should be complemented
961 by formal proofs provided in appendix or supplemental material.
- 962 • Theorems and Lemmas that the proof relies upon should be properly referenced.

963 4. Experimental result reproducibility

964 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
965 perimental results of the paper to the extent that it affects the main claims and/or conclusions
966 of the paper (regardless of whether the code and data are provided or not)?

967 Answer: [Yes]

968 Justification: The benchmark instances, evaluation harness, and containerized Docker en-
969 vironments are released at [https://huggingface.co/datasets/MobileDev-Bench/
970 mobiledev-bench](https://huggingface.co/datasets/MobileDev-Bench/mobiledev-bench). Appendix E details the Agentless adaptation (tree-sitter parsing, lan-
971 guage extensions), model identifiers and versions, temperature settings, reasoning-mode
972 configurations, and the Docker-based test harness. Instance-specific test commands are
973 included in each benchmark record.

974 Guidelines:

- 975 • The answer [N/A] means that the paper does not include experiments.
- 976 • If the paper includes experiments, a [No] answer to this question will not be perceived
977 well by the reviewers: Making the paper reproducible is important, regardless of
978 whether the code and data are provided or not.
- 979 • If the contribution is a dataset and/or model, the authors should describe the steps taken
980 to make their results reproducible or verifiable.
- 981 • Depending on the contribution, reproducibility can be accomplished in various ways.
982 For example, if the contribution is a novel architecture, describing the architecture fully
983 might suffice, or if the contribution is a specific model and empirical evaluation, it may
984 be necessary to either make it possible for others to replicate the model with the same
985 dataset, or provide access to the model. In general, releasing code and data is often
986 one good way to accomplish this, but reproducibility can also be provided via detailed
987 instructions for how to replicate the results, access to a hosted model (e.g., in the case
988 of a large language model), releasing of a model checkpoint, or other means that are
989 appropriate to the research performed.
- 990 • While NeurIPS does not require releasing code, the conference does require all submis-
991 sions to provide some reasonable avenue for reproducibility, which may depend on the
992 nature of the contribution. For example

- 993 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
994 to reproduce that algorithm.
995 (b) If the contribution is primarily a new model architecture, the paper should describe
996 the architecture clearly and fully.
997 (c) If the contribution is a new model (e.g., a large language model), then there should
998 either be a way to access this model for reproducing the results or a way to reproduce
999 the model (e.g., with an open-source dataset or instructions for how to construct
1000 the dataset).
1001 (d) We recognize that reproducibility may be tricky in some cases, in which case
1002 authors are welcome to describe the particular way they provide for reproducibility.
1003 In the case of closed-source models, it may be that access to the model is limited in
1004 some way (e.g., to registered users), but it should be possible for other researchers
1005 to have some path to reproducing or verifying the results.

1006 5. Open access to data and code

1007 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1008 tions to faithfully reproduce the main experimental results, as described in supplemental
1009 material?

1010 Answer: [Yes]

1011 Justification: Task instances (problem statements, base commits, test patches, fix patches,
1012 test commands, and transition labels) are publicly released at <https://huggingface.co/datasets/MobileDev-Bench/mobiledev-bench> under a CC-BY 4.0 license. The
1013 evaluation harness and per-repository Dockerfiles are released alongside the dataset.
1014

1015 Guidelines:

- 1016 • The answer [N/A] means that paper does not include experiments requiring code.
- 1017 • Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1018 • While we encourage the release of code and data, we understand that this might not
1019 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
1020 including code, unless this is central to the contribution (e.g., for a new open-source
1021 benchmark).
- 1022 • The instructions should contain the exact command and environment needed to run to
1023 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1024 • The authors should provide instructions on data access and preparation, including how
1025 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1026 • The authors should provide scripts to reproduce all experimental results for the new
1027 proposed method and baselines. If only a subset of experiments are reproducible, they
1028 should state which ones are omitted from the script and why.
- 1029 • At submission time, to preserve anonymity, the authors should release anonymized
1030 versions (if applicable).
- 1031 • Providing as much information as possible in supplemental material (appended to the
1032 paper) is recommended, but including URLs to data and code is permitted.
- 1033 •
- 1034 •

1035 6. Experimental setting/details

1036 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
1037 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1038 Answer: [Yes]

1039 Justification: Section 5 and Appendix E cover all evaluation details: the two evaluation
1040 settings (Agentless and Oracle), the four models evaluated with their versions and access
1041 method (OpenRouter), temperature settings (0.0 for localization, 0.0 greedy and 0.8 diverse
1042 for Agentless patch generation, 0.0 greedy for Oracle patch generation), reasoning-mode
1043 disabling for Agentless (not applied under Oracle), the tree-sitter adaptation, file-level
1044 localization for non-code artifacts, and the 30-minute per-task timeout. No training is
1045 performed; the benchmark is evaluation-only.

1046 Guidelines:

- 1047 • The answer [N/A] means that the paper does not include experiments.
- 1048 • The experimental setting should be presented in the core of the paper to a level of detail
- 1049 that is necessary to appreciate the results and make sense of them.
- 1050 • The full details can be provided either with the code, in appendix, or as supplemental
- 1051 material.

1052 7. Experiment statistical significance

1053 Question: Does the paper report error bars suitably and correctly defined or other appropriate

1054 information about the statistical significance of the experiments?

1055 Answer: [No]

1056 Justification: Error bars are not reported. Following prior work on issue-resolution bench-

1057 marks (SWE-bench, Multi-SWE-bench), we report deterministic resolution rates using

1058 temperature 0.0 for localization and a single greedy sample for the primary Agentless

1059 result. Rerunning the full evaluation across 415 tasks \times 4 models \times 2 settings would be

1060 prohibitively expensive.

1061 Guidelines:

- 1062 • The answer [N/A] means that the paper does not include experiments.
- 1063 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
- 1064 intervals, or statistical significance tests, at least for the experiments that support the
- 1065 main claims of the paper.
- 1066 • The factors of variability that the error bars are capturing should be clearly stated (for
- 1067 example, train/test split, initialization, random drawing of some parameter, or overall
- 1068 run with given experimental conditions).
- 1069 • The method for calculating the error bars should be explained (closed form formula,
- 1070 call to a library function, bootstrap, etc.)
- 1071 • The assumptions made should be given (e.g., Normally distributed errors).
- 1072 • It should be clear whether the error bar is the standard deviation or the standard error
- 1073 of the mean.
- 1074 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 1075 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 1076 of Normality of errors is not verified.
- 1077 • For asymmetric distributions, the authors should be careful not to show in tables or
- 1078 figures symmetric error bars that would yield results that are out of range (e.g., negative
- 1079 error rates).
- 1080 • If error bars are reported in tables or plots, the authors should explain in the text how
- 1081 they were calculated and reference the corresponding figures or tables in the text.

1082 8. Experiments compute resources

1083 Question: For each experiment, does the paper provide sufficient information on the com-

1084 puter resources (type of compute workers, memory, time of execution) needed to reproduce

1085 the experiments?

1086 Answer: [Yes]

1087 Justification: Appendix E.5 specifies that each task instance is executed in an isolated Docker

1088 container with a 30-minute timeout, with repository-specific dependencies installed from

1089 the project's build configuration (Gradle for Android, pub for Flutter, npm/yarn for React

1090 Native). Model inference is accessed via the OpenRouter API, requiring no local GPU

1091 resources. The total number of evaluated instances (415 tasks \times 4 models \times 2 settings) and

1092 the per-task timeout bound the compute envelope.

1093 Guidelines:

- 1094 • The answer [N/A] means that the paper does not include experiments.
- 1095 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 1096 or cloud provider, including relevant memory and storage.
- 1097 • The paper should provide the amount of compute required for each of the individual
- 1098 experimental runs as well as estimate the total compute.

- 1099 • The paper should disclose whether the full research project required more compute
1100 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1101 didn't make it into the paper).

1102 9. Code of ethics

1103 Question: Does the research conducted in the paper conform, in every respect, with the
1104 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1105 Answer: [Yes]

1106 Justification: The benchmark is constructed entirely from public open-source GitHub
1107 repositories under permissive licenses (Apache 2.0, MIT, GPL variants; listed in Table 6).
1108 No personal data, private repositories, or proprietary code is included. Human annotators
1109 reviewed benchmark instances; annotations were performed by the authors as part of their
1110 research duties.

1111 Guidelines:

- 1112 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
1113 Ethics.
- 1114 • If the authors answer [No], they should explain the special circumstances that require a
1115 deviation from the Code of Ethics.
- 1116 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1117 eration due to laws or regulations in their jurisdiction).

1118 10. Broader impacts

1119 Question: Does the paper discuss both potential positive societal impacts and negative
1120 societal impacts of the work performed?

1121 Answer: [Yes]

1122 Justification: Appendix J discusses positive impacts (exposing failure modes of LLM-based
1123 repair in mobile CI pipelines, guiding practitioners on where AI-generated patches can and
1124 cannot be trusted) and raises the need for human oversight when deploying model-generated
1125 patches in accessibility-, localisation-, or security-sensitive components. No direct negative
1126 societal impacts are identified, as the benchmark does not release a model or enable new
1127 harmful capabilities.

1128 Guidelines:

- 1129 • The answer [N/A] means that there is no societal impact of the work performed.
- 1130 • If the authors answer [N/A] or [No], they should explain why their work has no societal
1131 impact or why the paper does not address societal impact.
- 1132 • Examples of negative societal impacts include potential malicious or unintended uses
1133 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1134 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1135 groups), privacy considerations, and security considerations.
- 1136 • The conference expects that many papers will be foundational research and not tied
1137 to particular applications, let alone deployments. However, if there is a direct path to
1138 any negative applications, the authors should point it out. For example, it is legitimate
1139 to point out that an improvement in the quality of generative models could be used to
1140 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
1141 that a generic algorithm for optimizing neural networks could enable people to train
1142 models that generate Deepfakes faster.
- 1143 • The authors should consider possible harms that could arise when the technology is
1144 being used as intended and functioning correctly, harms that could arise when the
1145 technology is being used as intended but gives incorrect results, and harms following
1146 from (intentional or unintentional) misuse of the technology.
- 1147 • If there are negative societal impacts, the authors could also discuss possible mitigation
1148 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1149 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1150 feedback over time, improving the efficiency and accessibility of ML).

1151 11. Safeguards

1152 Question: Does the paper describe safeguards that have been put in place for responsible
1153 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1154 image generators, or scraped datasets)?

1155 Answer: [N/A]

1156 Justification: MobileDev-Bench is a benchmark dataset derived from public GitHub pull
1157 requests and issues. It does not release a pre-trained model, image generator, or scraped
1158 dataset containing sensitive content. The constituent repositories are already publicly avail-
1159 able under open-source licenses, and no new risk surface is introduced by their aggregation
1160 into benchmark instances.

1161 Guidelines:

- 1162 • The answer [N/A] means that the paper poses no such risks.
- 1163 • Released models that have a high risk for misuse or dual-use should be released with
1164 necessary safeguards to allow for controlled use of the model, for example, by requiring
1165 that users adhere to usage guidelines or restrictions to access the model or implementing
1166 safety filters.
- 1167 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1168 should describe how they avoided releasing unsafe images.
- 1169 • We recognize that providing effective safeguards is challenging, and many papers do
1170 not require this, but we encourage authors to take this into account and make a best
1171 faith effort.

1172 12. Licenses for existing assets

1173 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1174 the paper, properly credited and are the license and terms of use explicitly mentioned and
1175 properly respected?

1176 Answer: [Yes]

1177 Justification: Table 6 lists all 19 source repositories with their licenses (GNU GPLv3,
1178 Apache 2.0, GNU AGPLv3, MIT, EPL 2.0, GNU GPLv2). The Agentless framework [Xia
1179 et al., 2024] and tree-sitter [Brunsfield et al., 2018] are cited. All repositories are used in
1180 accordance with their open-source licenses, which permit research use.

1181 Guidelines:

- 1182 • The answer [N/A] means that the paper does not use existing assets.
- 1183 • The authors should cite the original paper that produced the code package or dataset.
- 1184 • The authors should state which version of the asset is used and, if possible, include a
1185 URL.
- 1186 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1187 • For scraped data from a particular source (e.g., website), the copyright and terms of
1188 service of that source should be provided.
- 1189 • If assets are released, the license, copyright information, and terms of use in the
1190 package should be provided. For popular datasets, `paperswithcode.com/datasets`
1191 has curated licenses for some datasets. Their licensing guide can help determine the
1192 license of a dataset.
- 1193 • For existing datasets that are re-packaged, both the original license and the license of
1194 the derived asset (if it has changed) should be provided.
- 1195 • If this information is not available online, the authors are encouraged to reach out to
1196 the asset’s creators.

1197 13. New assets

1198 Question: Are new assets introduced in the paper well documented and is the documentation
1199 provided alongside the assets?

1200 Answer: [Yes]

1201 Justification: MobileDev-Bench is released at [https://huggingface.co/datasets/
1202 MobileDev-Bench/mobiledev-bench](https://huggingface.co/datasets/MobileDev-Bench/mobiledev-bench) with a dataset card describing fields, construction
1203 methodology, and intended use. Each instance includes structured metadata (instance
1204 ID, framework, base commit, difficulty, category, artifact types, test transitions, and test
1205 command). The evaluation harness and Dockerfiles are released alongside the dataset.

1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [Yes]

Justification: Manual verification (Phase 5, Section 3.5 and Appendix B) was performed by two annotators who are co-authors of the paper, not external crowdworkers. Appendix B describes the four review criteria (problem statement clarity, test coverage appropriateness, task difficulty, task category), the scoring rubric (0–3 per dimension), and the inter-annotator disagreement resolution process. Figure 6 shows the resulting score distributions. No external compensation was involved.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [N/A]

Justification: The annotation work described in the paper was performed by the authors themselves and does not constitute research with human subjects under standard IRB definitions. No external participants were recruited, and no personal data was collected.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

1257 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1258 non-standard component of the core methods in this research? Note that if the LLM is used
1259 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
1260 scientific rigor, or originality of the research, declaration is not required.

1261 Answer: [Yes]

1262 Justification: LLMs (Claude Sonnet 4.5, GPT-5.2, Gemini 2.5 Flash, Qwen-3-Coder) are the
1263 evaluated systems and are central to the paper's experiments; their usage is fully described in
1264 Section 5 and Appendix E.3. AI-based writing tools were used only for minor grammatical
1265 and wording improvements and did not affect methodology or results, as noted in Appendix J.

1266 Guidelines:

- 1267 • The answer [N/A] means that the core method development in this research does not
1268 involve LLMs as any important, original, or non-standard components.
- 1269 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
1270 be described.