

SWE-EFFICIENCY: CAN LANGUAGE MODELS OPTIMIZE REAL-WORLD REPOSITORIES ON REAL WORKLOADS?

Anonymous authors

Paper under double-blind review

ABSTRACT

Optimizing the performance of large-scale software repositories demands expertise in code reasoning and software engineering (SWE) to reduce runtime while preserving program correctness. However, most benchmarks emphasize what to fix rather than how to fix code. We introduce SWE-EFFICIENCY, a benchmark for evaluating repository-level performance optimization on real workloads. Our suite contains 498 tasks across nine widely used data-science, machine-learning, and HPC repositories (e.g., `numpy`, `pandas`, `scipy`): given a complete codebase and a slow workload, an agent must investigate code semantics, localize bottlenecks and relevant tests, and produce a patch that matches or exceeds expert speedup while passing the same unit tests. To enable this how-to-fix evaluation, our automated pipeline scrapes GitHub pull requests for performance-improving edits, combining keyword filtering, static analysis, coverage tooling, and execution validation to both confirm expert speedup baselines and identify relevant repository unit tests. Empirical evaluation of state-of-the-art agents reveals significant underperformance. On average, agents achieve less than $0.15\times$ the expert speedup: agents struggle in localizing optimization opportunities, reasoning about execution across functions, and maintaining correctness in proposed edits. We release the benchmark and accompanying data pipeline to facilitate research on automated performance engineering and long-horizon software reasoning.

1 INTRODUCTION

Language models (LMs) are becoming an increasingly substantial part of software engineering, from LM-powered auto-complete to autonomous software-engineering agents that plan, implement, and verify changes in large repositories. Recent agentic systems show that LMs can fix functional bugs and implement small features (Jimenez et al., 2024; Jain et al., 2024b; Yang et al., 2024; Wang et al., 2025). However, most benchmarks for these systems focus on what gets fixed or resolved, not the properties of code implementations—overlooking runtime performance, memory efficiency, style, and other software-engineering concerns. As we reach the limits of hardware, software optimizations become critical and have tremendous impact: Jain et al. (2024a) show that pure software changes can reduce high-utilization workload throughput by 10% on Google’s datacenter compute, saving estimated millions of dollars. Recent benchmarks begin to probe code performance (e.g., KernelBench, Ouyang et al. (2025); PIE, Shypula et al. (2024); EffiBench, Huang et al. (2024)), but they avoid real-world, end-to-end workloads on real repositories. We therefore ask: *to what extent can LM agents optimize the runtime of real-world repositories on real-world workloads?*

Recent work has begun to evaluate whether LMs can improve repo-level software runtime, most notably, GSO (Shetty et al., 2025) and SWE-Perf (Fan et al., 2025): both benchmarks tackle the problem of repo-level code optimization. GSO provides each task with an oracle script verifying functional equivalence and runs hidden performance tests at evaluation. SWE-Perf adapts existing repo unit-tests for both correctness and performance measurement, with task instructions pointing agents to optimize specific functions. However, software repositories commonly separate correctness and performance tests (ISO/IEC, 2011), and immediate correctness oracles are usually unavailable. Thus, these setups still insufficiently assess a core part of performance engineering: investigating an unfamiliar repository to recover code semantics and correctness from the codebase alone.

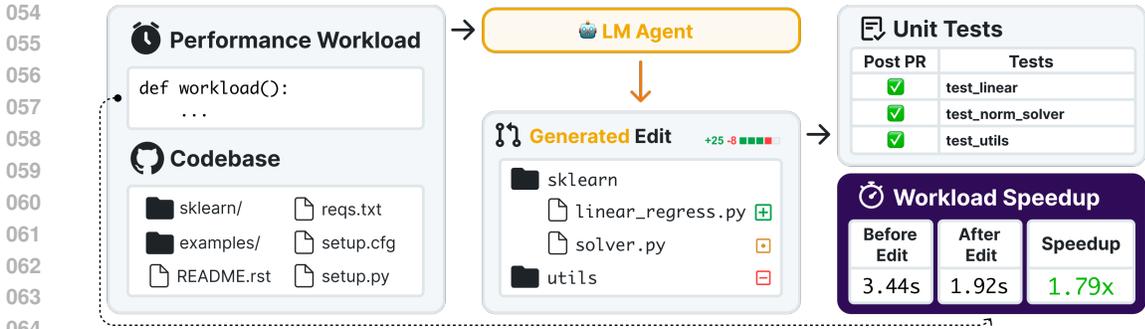


Figure 1: SWE-FFICIENCY evaluates the investigative, pass-to-pass workflow of performance engineering: given an existing codebase state and a performance workload of interest, agents must edit the codebase to speed up that workload while keeping relevant repo unit tests green.

Performance engineers *characterize* a workload (which can be slow for any myriad of reasons); *localize* where to intervene; and, just as importantly, *localize tests*—identifying and executing existing unit-tests to be confident that an optimization does not introduce new functionality. We design our benchmark to target this challenging and open-ended investigative workflow.

To address these gaps, we propose SWE-FFICIENCY (pronounced *swee-FISH-uhn-see*), a new benchmark to evaluate how well LMs can improve the performance of real-world workloads through modifying software repositories (Figure 1). To build SWE-FFICIENCY, we propose a novel, systematic data collection pipeline for optimization task instances, which uses attribute filtering, static analysis, code coverage, and execution validation. This anchors the realism of the benchmark and usefulness of the optimization tasks—a large and diverse set of 498 tasks across 9 codebases across data science, machine learning, and high performance computing. We score LM systems using *speedup ratio* (*SR*), which evaluates how well models match or improve upon expert edits and motivates long-term progress on our benchmark.

We also conduct a holistic evaluation of LMs on SWE-FFICIENCY to better understand strengths and limitations. We reveal systemic gaps: on average, LMs achieve less than $0.15\times$ expert speedup and often introduce correctness bugs via proposed edits. Models struggle to localize the same expert optimization opportunities and prefer superficial speedups than more principled expert algorithmic rewrites. Thus, while LMs exhibit promise in other SWE tasks, substantial advances in repo-level reasoning, systems optimization, and long-horizon planning are needed to close this expert gap.

Our contributions. (1) A scalable, oracle-free benchmark of 498 tasks across 9 repos, requiring deep codebase investigation and test localization. (2) A systematic pipeline for extracting realistic and reproducible performance engineering tasks from GitHub repos. (3) An evaluation metric, speedup ratio, that measures parity with experts and encourages long term benchmark progress. (4) Empirical and qualitative analysis revealing large gaps between LMs and experts in edit localization and principled optimizations. (5) Open-sourced dataset, benchmark harness, and pipeline to accelerate research on automated performance engineering and long-horizon software reasoning.

2 SWE-FFICIENCY OVERVIEW

SWE-FFICIENCY is a benchmark containing real performance-optimization GitHub pull requests from popular repositories. The task is to generate a pull request that modifies the codebase to make a given workload faster while preserving the correctness of existing repo tests.

2.1 DATA COLLECTION PROCEDURE

We scrape nine popular Python GitHub repos, including `astropy`, `dask`, `matplotlib`, `numpy`, `pandas`, `scikit-learn`, `scipy`, `sympy`, and `xarray`. Figure 2 shows how our method extends the scraping recipe from SWE-bench, modifying the attribute and execution filtering (Stages 2 and 5) to select previously-excluded performance edits and introducing test coverage filtering and workload annotation (Stages 3 and 4) to identify reproducible, verifiable optimization tasks.

Table 1: SWE-FFICIENCY jointly (i) evaluates the runtime of performance workloads, (ii) verifies correctness using a repository’s own tests, and (iii) uses separate correctness and performance workloads. For more details on related benchmarks, see Section 5.

Benchmark	Evaluates Runtime	Repo Level	Correctness Eval: Using Repo’s Own Tests	Performance Eval: Separate end-to-end system test	# of Optimization Tasks
SWE-BENCH	✗	✓	✓	✗	0
EFFIBENCH	✓	✗	✗	✗	1000
MERCURY	✓	✗	✗	✗	1889
PIE	✓	✗	✗	✗	978
KERNELBENCH	✓	✗	✗	✗	250
ALGOTUNE	✓	✗	✗	✗	154
GSO	✓	✓	✗	✓	102
SWE-PERF	✓	✓	✓	✗	140
SWE-FFICIENCY (OURS)	✓	✓	✓	✓	498

Stage I: Repo selection and instance scraping. We target GitHub pull requests (PRs) from popular data science, machine learning, and high-performance computing repositories—these domains are performance-sensitive and contain PRs where authors explicitly optimize runtime. Widely-used libraries surface optimizations that are actually useful in real-world software.

Stage II: Performance regression attribute filtering. We prune away PRs that clearly are not performance-related or that introduce new behavior: in contrast, issue-resolution benchmarks like SWE-bench filter for the opposite by choosing tasks that add new tests. We select PRs only when (i) metadata includes performance keywords (i.e. `perf`, `speedup`, `benchmark`); (ii) PRs do not modify tests, to avoid behavior-changing edits misrepresenting as optimizations; and (iii) edits meaningfully modify the file’s abstract syntax tree (AST), excluding no-op or docs-only diffs.

Stage III: Identifying covering correctness tests. To enforce our benchmark’s invariant specification (i.e. all relevant tests must continue to pass after an edit), we require that at least one existing unit test exercises the modified code. Per instance, we build a Docker image with pinned dependencies, run the repository’s test suite, and use line coverage to confirm the edit is exercised.

Stage IV: Annotating performance workloads. Unit tests often do not capture runtime behavior, and software generally separates correctness from performance tests (ISO/IEC, 2011). Using PR descriptions and discussion as context, we manually annotate each task—writing a workload script that, when run before and after the PR’s edit, shows a measurable performance improvement. Although PR info often includes ad-hoc demo scripts, these are not reliably auto-extractable; likewise, LM-based workload generation from patches fails to consistently elicit claimed gains (see Sec. 4.2).

Stage V: Execution-based filtering. To curate a final set of verifiable and consistently reproducible optimization tasks, we run each instance’s unit tests and annotated workload in a controlled environment (containerization, resource pinning) to ensure no interference with speedup measurements.

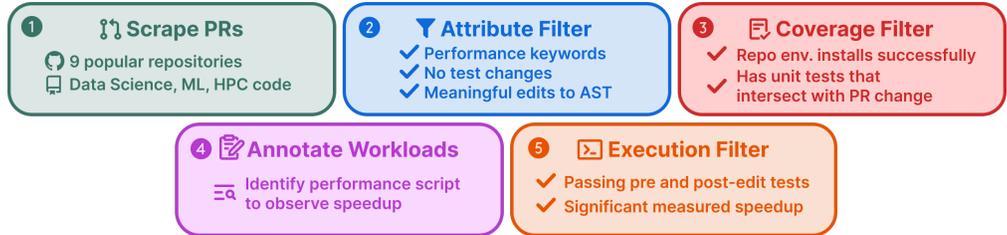


Figure 2: SWE-FFICIENCY collects tasks through a multi-stage scraping pipeline: each stage prunes candidate tasks that introduce new behavior, are unlikely to be performance related, or unsuitable for reproducible benchmarking. This yields a set of tasks, each of which have an accompanying expert or *gold* patch. See Appendix C for stage-specific details.

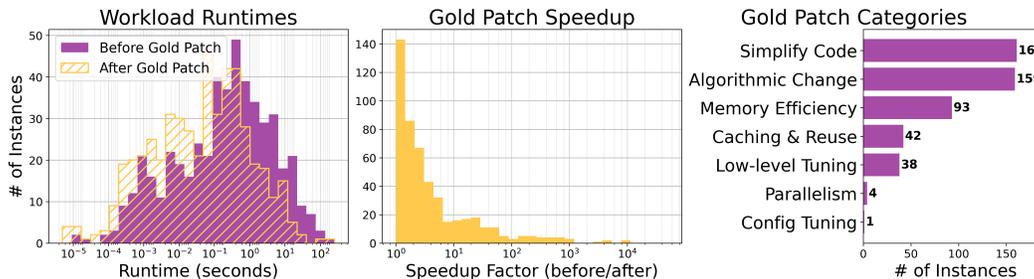


Figure 3: SWE-FFICIENCY contains a diverse distribution over performance workload runtime (left); over gold patch speedup (speedup achieved from expert PR edit); and over types of optimizations made by the expert (right). We use an LM to categorize the gold patch for each instance (for high-level analysis only) and manually verify a randomly chosen subset: see Appendix B.

We retain only instances that demonstrate significant speedups (runtime improvement greater than $2\times$ measurement std. dev.) and log test statuses for benchmark correctness checks.

2.2 SWE-FFICIENCY DATASET DISTRIBUTION AND UNIQUE BENCHMARK FEATURES

Open-ended but precise evaluation criteria. Providing agents with just an executable code snippet (workload) and codebase makes the task very open-ended: agents can choose any approach, including changes different from the expert’s *gold* patch. This mirrors the flexibility of real-world performance engineering—there rarely is a single prescriptive path to faster implementations. Our evaluation is made precise by grounding correctness in a set of unit tests and measuring LM speedup against gold patch speedup (i.e. how much speedup the expert achieved). Our benchmark encourages creativity in edit strategies while guaranteeing unambiguous criteria for strong optimizations.

Clear distinction between performance and correctness tests. Repositories generally require unit tests to run quickly (unlike more substantial performance workloads), and software standards encourage performance benchmarks to be clearly separated from correctness tests (ISO/IEC, 2011). Thus, defining a performance workload using unit-tests or a combined correctness-performance oracle is not fully reflective of actual performance engineering. Instead, SWE-FFICIENCY clearly separates performance evaluation workloads from repo correctness tests.

Preserving existing correctness during optimization. Unlike SWE-bench’s issue-resolution setting (evaluating bug-fixes that flip failing tests to passing), our benchmark targets *pass-to-pass* optimization—speeding up already-correct code without introducing new behavior. Specifically, we choose PRs that do not introduce new behavior: edits that introduce new features (and new tests) may have unintended performance effects on other workloads, and confound our specific evaluation of code optimization abilities. Evaluating how agents perform edits in this constrained task provides more confidence they can be deployed in real codebases without disrupting existing behavior.

2.3 TASK FORMULATION

Model input. An agent is given a complete codebase and a performance workload exercising codebase functionality. We task the agent with modifying the codebase so that workload runtime improves while expected repository unit tests still pass. Expert performance engineers only require a reported slow workload and codebase to start optimizing: first characterizing the workload’s bottlenecks, modifying files, verifying speedup against the workload and identifying relevant unit tests to check for no regressions. For examples of performance workloads, see Appendix B.2.

Evaluation metrics. Our evaluation metric is *speedup ratio (SR)*, which answers the question: *normalized to the expert edit, how well does the LM’s generated edit perform?*. We apply an LM’s submitted patch to the codebase and run repository test files associated with each instance. If the patch applies successfully and all tests pass, we compute the instance *speedup ratio* as $SR = Speedup_{LM}/Speedup_{gold}$ where gold speedup is $Speedup_{gold} = T_{pre}/T_{post-gold-patch}$ and LM

speedup is $Speedup_{LM} = T_{pre}/T_{post-LM-patch}$. For example, if the expert achieves a gold speedup of $5\times$ and the LM achieves $1.2\times$ on the same instance, $SR = 1.2/5 = 0.24$. To aggregate across instances, we take the *harmonic mean* of each SR: if a system submits an empty patch or a patch that fails unit tests, the instance’s speedup ratio is $SR = 1/Speedup_{gold}$. We use harmonic mean since it is most appropriate for averaging speedup ratios (Smith, 1988; Eeckhout, 2024).

Why a factor-based evaluation metric (not % solved)? We adopt speedup ratio because it provides a long runway for progress and explicitly rewards going beyond human parity. Percentage-style metrics collapse to two regimes—near 0% today and nearer 100% once tests are routinely passed—leaving little room to compare systems once the benchmark begins to saturate. It motivates *continued progress*: anchoring the scale at expert parity ($1\times$) turns super-human performance into a first-class goal and keeps the leaderboard competitive after models reach expert performance. This means our benchmark stays meaningful for the community both now (when systems only reach $0.15\times$ of expert performance) and later (when models might consistently score above $1\times$).

3 EVALUATION SETUP

Machine Configuration. We containerize each task environment: instance Docker images are built and uploaded to a registry for reproducibility and easy integration with agent harnesses. All evaluations are run on a single Google Cloud `n2-standard-64` VM (64 vCPUs, 256GB Memory). To parallelize the benchmark for faster evaluation without interference from parallel workers, we pin each worker to an exclusive set of *physical* CPU cores (4 vCPUs), the CPU’s corresponding memory node, and assigning a memory limit (16GB) per worker. For more details, see Appendix D.

Agent Scaffold. We provide baseline performance on two open-source agent harnesses, `OPENHANDS` (`CodeActAgent-v0.51.1`) (Wang et al., 2025) and `SWE-AGENT` (`v1.1.0`) (Yang et al., 2024). Both scaffolds provide file-editing tools and a bash terminal interface for LM agents to easily edit code and execute commands. We configure agents with a 3-hour time limit per task, a 30-minute timeout per step, a maximum action count of 100, and provide the same number of vCPUs and memory as the evaluation setting. The agent is provided a task prompt and repo-specific commands for rebuilding (to support possible C/C++/Cython edits) and executing arbitrary test files. In the `SWE-AGENT` setting, we configure the underlying LM with a \$1 token-spend max per task to observe performance under limited inference cost. We provide evaluations on both harnesses specifically for `CLAUDE-3.7-SONNET`, `GPT-5 MINI`, and `GEMINI 2.5 FLASH` models. See Appendix E for agent prompts and details.

Models. We evaluate several frontier models from OpenAI, Anthropic, Google, Z.ai, Moonshot AI, and DeepSeek: `GPT-5` (OpenAI, 2025a), `GPT-5 MINI` (OpenAI, 2025b), `CLAUDE 4.1 OPUS` (Anthropic, 2025a), `CLAUDE 4.5 SONNET` (Anthropic, 2025c), `CLAUDE-3.7-SONNET` (Anthropic, 2025b), `GEMINI-2.5 PRO` (Google, 2025b), `GEMINI 2.5 FLASH` (Google, 2025a), `GLM-4.6` (Z.ai, 2025), `KIMI K2-9005` (Moonshot AI, 2025), and `DEEPSEEK V3.1` (DeepSeek, 2025). For each instance, we sample a single trajectory and report the aggregated speedup ratio. We focus on $pass@1$ because it best matches both agent capabilities and realistic human workflows: (i) expert pull requests are effectively $pass@1$ (infeasible to review multiple PR submissions) and (ii) agentic LMs can still explore multiple edits, execute workloads repeatedly, and iterate over alternatives within a single trajectory. This also follows common practice in prior benchmarks, including HumanEval and SWE-bench, where $pass@1$ is the primary metric (Chen et al., 2021).

4 EXPERIMENTS AND RESULTS

On `SWE-EFFICIENCY`, leading LM agents trail experts and often introduce correctness bugs. Our benchmark enables key quantitative and qualitative observations about LM agent behavior, namely how models solve easier cases, falter on harder ones, and exhibit convenience bias—small, input-specific, harder-to-maintain edits—underscoring the gap to expert-level performance engineering.

Table 2: SWE-FFICIENCY results across several frontier models (higher is better; human-expert speedup ratio (SR) is 1.0×). SR is *pass@1*: each system submits a single patch per instance to be evaluated. SR is calculated by normalizing the speedup from the LM-generated edit to the speedup from the gold (human-written) patch, aggregated across all tasks via harmonic mean. All experiments below used OPENHANDS: SWE-AGENT achieves similar results, see Appendix F.

System	Speedup Ratio (↑)
GPT-5	0.150×
CLAUDE 4.1 OPUS	0.098×
QWEN3 CODER PLUS	0.064×
CLAUDE 3.7 SONNET	0.047×
CLAUDE 4.5 SONNET	0.041×
GLM-4.6	0.026×
GPT-5 MINI	0.019×
KIMI K2-0905	0.008×
GEMINI 2.5 FLASH	0.008×
DEEPSEEK V3.1	0.007×
GEMINI 2.5 PRO	0.007×

Table 3: Further breakdown of patch outcomes by system. Pre-edit denotes codebase before any edits. “Passes correctness tests” refers to functional correctness only (not necessarily perf. optimal). All experiments below used OPENHANDS.

System	Fails Tests (↓)	Passes Correctness Tests		
		Slower than Pre-edit (↓)	Faster than Pre-edit (↑)	Faster than Expert (↑)
COMPOSER-1 (CURSOR CLI)	16%	13%	48%	23%
GPT-5	18%	4%	32%	46%
CLAUDE 4.1 OPUS	15%	4%	43%	38%
QWEN3 CODER PLUS	22%	11%	42%	24%
CLAUDE 3.7 SONNET	35%	12%	32%	20%
CLAUDE 4.5 SONNET	19%	5%	44%	33%
GLM-4.6	33%	13%	38%	16%
GPT-5 MINI	45%	15%	27%	13%
KIMI K2-0905	26%	11%	44%	19%
GEMINI 2.5 FLASH	39%	14%	34%	12%
DEEPSEEK V3.1	19%	18%	44%	18%
GEMINI 2.5 PRO	40%	18%	34%	8%

4.1 OVERALL PERFORMANCE

Leading agents struggle on SWE-FFICIENCY. Across all agents, we observe that LM agents struggle to achieve more than 0.15× of expert level performance. Table 2 summarizes speedup ratio performance of leading software-engineering agents on SWE-FFICIENCY. We see a substantial capability transfer gap: GPT-5 MINI (OPENHANDS) achieved 0.019× of expert speedup, while the same system scored 62.6% on SWE-bench Verified. This indicates that current agents, while successful on issue-resolution and bug-fix tasks, currently do not immediately transfer to efficiency-oriented program changes, showing substantial headroom for improvement.

Agents often introduce bugs during optimization. LM agents often propose edits that cause repository unit tests to newly fail, invalidating any optimizations made. Table 3 unpacks agent performance across different unit test and performance outcomes: even when patches are functionally correct, the majority of edits are still slower than the expert. Strikingly, with the exceptions of GPT-5, CLAUDE 4.1 OPUS, and CLAUDE 4.5 SONNET, fewer than a quarter of solutions are both correct and outperform expert-level speedups.

Strong on easy wins, weak on harder speedups. We identify three measures of task difficulty: (1) *pre-edit workload runtime* (longer duration workloads likely require more algorithmic insight); (2) *gold patch length* (harder instances require editing more lines); and (3) *the speedup factor* that

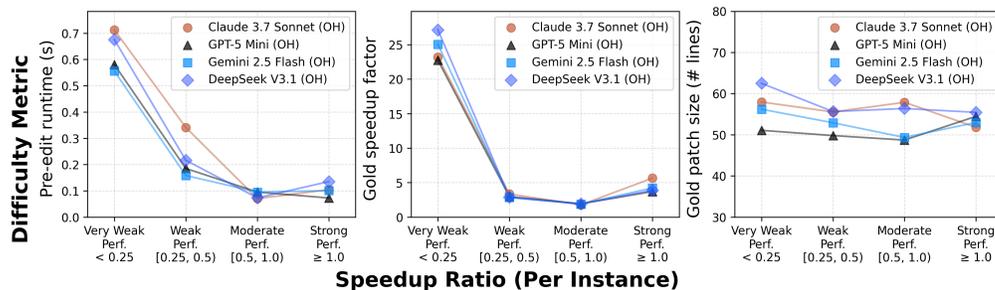


Figure 4: LMs achieve strong performance on easier problems but struggle on tasks with longer workload runtime duration and larger baseline expert speedups. We bucket LM submissions by per-instance speedup ratio and compute the geometric mean per-bucket of (i) pre-edit workload runtime, (ii) the gold (expert) patch speedup, and (iii) the number of lines in the gold patch.

the expert edit achieves (instance is harder if expert speedup is larger). Figure 4 shows a breakdown of benchmark performance in relation to these task difficulty measures. Across all three measures of task complexity, LMs are able to match expert performance on lower-complexity tasks. However, LMs struggle to solve tasks with longer duration workloads or larger feasible speedup opportunities.

Function-level mislocalization severely limits LM performance. Much of LM underperformance appears to stem from *failing to optimize the same functions as the expert*. If we view expert (gold) speedup as “mass” distributed over edited files and functions, Figure 6 shows that over 68% of expert gains occur in functions the LM never edits. Although LM and expert modify the same files over 55% of the time, they miss the functions carrying most of the expert’s speedup. Likewise, Figure 7 visualizes a function-level breakdown where the LM makes an attempted optimization in a deeper function and fails to match the expert’s improvement. For more details, see Appendix G.

4.2 QUALITATIVE ANALYSIS

LMs make satisficing optimizations, giving up before expert parity. Figure 5 shows that the shortest sequences of agent actions (i.e. file-editing, running scripts) happen when LMs achieve speedup ratios exceeding $1\times$ —expert-level wins are found early. When LMs underperform experts, median trajectory action counts sit at less than mid length (30–50 turns), well below the 100 action cap. This pattern fits a *satisficing* story: once the model secures a measurable speedup, it tends to stop instead of pushing any closer to expert parity. Future agents can employ “don’t-stop-early” triggers when code heuristics show larger possible speedups.

Shortcut bias and caching as a crutch vs. systemic cost reduction. LMs preferentially add localized shortcuts—identity checks, ad-hoc early exits, and memoization—such as self-equality fast paths or persistent caches. Experts instead restructure code to reduce per-element cost. Figure 8

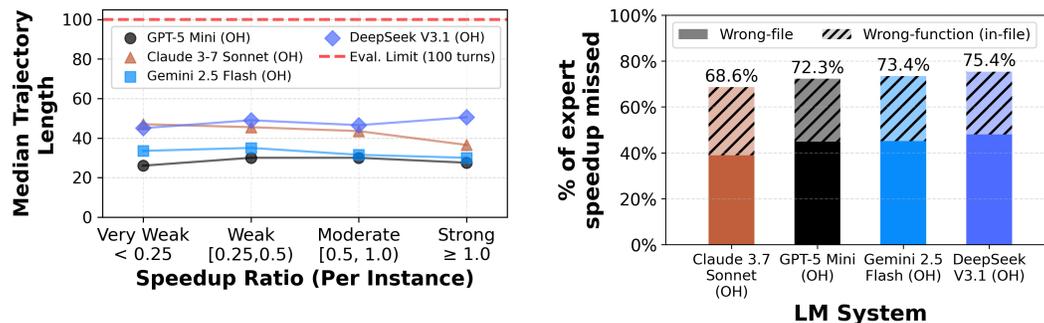
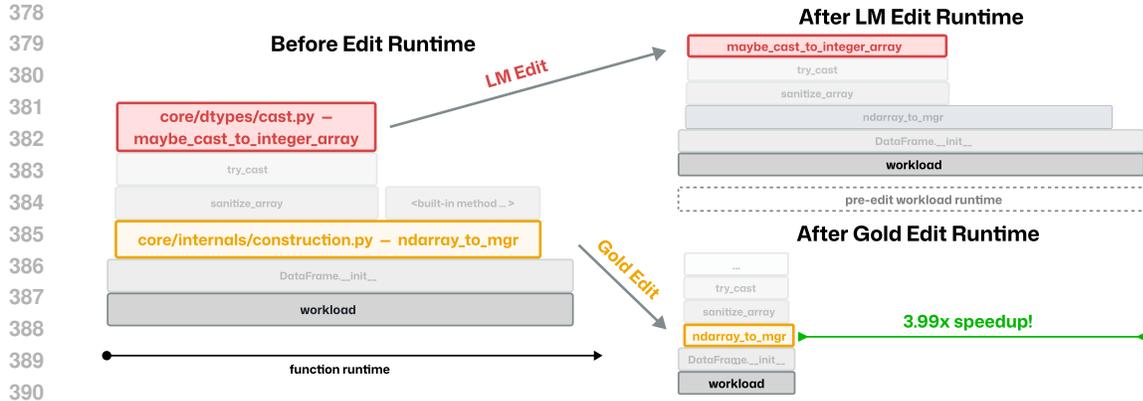


Figure 5: LMs find expert-level wins earlier in action trajectories. When they underperform experts, LMs submit satisficing optimizations rather than trying on for expert parity.

Figure 6: LMs leave a significant portion of expert-achievable speedup on the table due to wrong file/function selection and localization.



391 Figure 7: LMs prefer to edit different functions than the gold patch, missing out on major speedups.
392 For a workload flamegraph for task `pandas-dev__pandas-52054`, CLAUDE 3.7 SONNET
393 (SWE-AGENT) (red) chooses a different function (and file) than the expert (gold): it does not
394 achieve the expert’s overall workload speedup, since the expert’s speedup is at a shallower scope.
395

396
397
398
399
400
401
402
403
404
405
406
407

```

--- a/pandas/core/arrays/arrow/array.py
+++ b/pandas/core/arrays/arrow/array.py
@@ -406,8 +406,14 @@ def _cmp_method(self, other, op):
-     result = result.to_numpy()
-     return BooleanArray._from_sequence(result)
+     if result.null_count > 0:
+         values = pc.fill_null(result, False).
+         to_numpy()
+         mask = result.is_null().to_numpy()
+     else:
+         values = result.to_numpy()
+         mask = np.zeros(len(values), dtype=np.
+         bool_)
+     return BooleanArray(values, mask)

def _evaluate_op_method(self, other, op,
arrow_funcs):
    pc_func = arrow_funcs[op.__name_]
  
```

408
409
410
411
412

407 Figure 8: **Left:** Expert’s edit (gold patch) on instance `pandas-dev__pandas-50524` optimiz-
408 ing a workload via avoiding a conversion to object dtype (20.5× speedup). **Right:** CLAUDE 3.7
409 SONNET (OPENHANDS) instead identifies a different fast path optimization when no null elements
410 are present, but only achieves a 2.3× speedup (scoring a speedup ratio of 0.113×).
411
412

413 shows a flamegraph both after an LM versus an expert edit, where the expert optimizes by keeping
414 work in fast Arrow kernels and producing a `BooleanArray` from a values/mask pair without
415 materializing slow object-dtypes. Experts also use faster backends (Cython/Pythran/BLAS) to
416 reduce Python overhead or remove Python-level work entirely—vectorizing, moving loops to com-
417 piled code, or dispatching to type-aware fast paths. LMs yield strong speedups only when these
418 shortcut conditions hold, whereas systemic reductions are more broadly robust.
419

420 **Workload overfitting and semantic drift.** Another LM pattern is to bake benchmark properties
421 into patches, producing impressive but brittle wins. This sometimes crosses into correctness drift—
422 e.g., returning the original `DataFrame` from `groupby.apply` or monkey-patching `np.arange`.
423 Experts instead target generalizable structure (i.e. multi-index skipping, per-dimensional slice reuse)
424 while preserving functional behavior. With an preliminary version of our evaluation harness, some
425 agents exploited function stackframe info to detect when code is being run in our evaluation environ-
426 ment: we consequently added robust checks for this in the harness (see Appendix J.1).
427

428 **Maintainability of generated edits.** LM edits are frequently invasive—global monkey-patching,
429 module-level mutable caches, or fast paths tied to dynamic object attributes (an example shown in
430 Fig. 9). Expert patches are localized and composable—adding a function call with precomputed
431 constants, a Cython helper mirroring existing logic, or reusing shallow copies of constructor argu-
ments. Expert edits have a lower blast radius of code edits and are more maintainable long term.

```

432                                     --- a/pandas/core/series.py
433                                     +++ b/pandas/core/series.py
434                                     @@ -1816,9 +1816,18 @@ class Series(base.
                                        IndexOpsMixin, NDFrame): # type: ignore[misc]
                                        else:
435                                     -         return into_c(k, v) for k, v in self.
436                                     +         items()
437                                     +         values = getattr(self, "_values", None)
438                                     +         if values is None:
439                                     +             return into_c((k, v) for k, v in
440                                     +                 self.items())
441                                     +             try:
442                                     +                 list_vals = values.tolist()
443                                     +                 except Exception:
444                                     +                     # fallback to generic iteration
445                                     +                     list_vals = [v for v in values]
446                                     +                 return into_c(zip(self.index, list_vals)
447                                     +                     ))
448
449                                     --- a/pandas/core/series.py
450                                     +++ b/pandas/core/series.py
451                                     @@ -1818,7 +1818,7 @@ def to_dict(self, into: type[
452                                     dict] = dict) -> dict:
453                                     else:
454                                     -         return into_c(k, v) for k, v in self.
455                                     +         items()
456                                     +         return into_c(self.items())
457
458                                     -         return into_c(k, v) for k, v in self.
459                                     +         items()
460                                     +         return into_c(self.items())
461
462                                     -         return into_c(k, v) for k, v in self.
463                                     +         items()
464                                     +         return into_c(self.items())
465
466                                     -         return into_c(k, v) for k, v in self.
467                                     +         items()
468                                     +         return into_c(self.items())
469
470                                     -         return into_c(k, v) for k, v in self.
471                                     +         items()
472                                     +         return into_c(self.items())
473
474                                     -         return into_c(k, v) for k, v in self.
475                                     +         items()
476                                     +         return into_c(self.items())
477
478                                     -         return into_c(k, v) for k, v in self.
479                                     +         items()
480                                     +         return into_c(self.items())
481
482                                     -         return into_c(k, v) for k, v in self.
483                                     +         items()
484                                     +         return into_c(self.items())
485
486                                     -         return into_c(k, v) for k, v in self.
487                                     +         items()
488                                     +         return into_c(self.items())
489
490                                     -         return into_c(k, v) for k, v in self.
491                                     +         items()
492                                     +         return into_c(self.items())
493
494                                     -         return into_c(k, v) for k, v in self.
495                                     +         items()
496                                     +         return into_c(self.items())
497
498                                     -         return into_c(k, v) for k, v in self.
499                                     +         items()
500                                     +         return into_c(self.items())
501
502                                     -         return into_c(k, v) for k, v in self.
503                                     +         items()
504                                     +         return into_c(self.items())
505
506                                     -         return into_c(k, v) for k, v in self.
507                                     +         items()
508                                     +         return into_c(self.items())
509
510                                     -         return into_c(k, v) for k, v in self.
511                                     +         items()
512                                     +         return into_c(self.items())
513
514                                     -         return into_c(k, v) for k, v in self.
515                                     +         items()
516                                     +         return into_c(self.items())
517
518                                     -         return into_c(k, v) for k, v in self.
519                                     +         items()
520                                     +         return into_c(self.items())
521
522                                     -         return into_c(k, v) for k, v in self.
523                                     +         items()
524                                     +         return into_c(self.items())
525
526                                     -         return into_c(k, v) for k, v in self.
527                                     +         items()
528                                     +         return into_c(self.items())
529
530                                     -         return into_c(k, v) for k, v in self.
531                                     +         items()
532                                     +         return into_c(self.items())
533
534                                     -         return into_c(k, v) for k, v in self.
535                                     +         items()
536                                     +         return into_c(self.items())
537
538                                     -         return into_c(k, v) for k, v in self.
539                                     +         items()
540                                     +         return into_c(self.items())
541
542                                     -         return into_c(k, v) for k, v in self.
543                                     +         items()
544                                     +         return into_c(self.items())
545
546                                     -         return into_c(k, v) for k, v in self.
547                                     +         items()
548                                     +         return into_c(self.items())
549
550                                     -         return into_c(k, v) for k, v in self.
551                                     +         items()
552                                     +         return into_c(self.items())
553
554                                     -         return into_c(k, v) for k, v in self.
555                                     +         items()
556                                     +         return into_c(self.items())
557
558                                     -         return into_c(k, v) for k, v in self.
559                                     +         items()
560                                     +         return into_c(self.items())
561
562                                     -         return into_c(k, v) for k, v in self.
563                                     +         items()
564                                     +         return into_c(self.items())
565
566                                     -         return into_c(k, v) for k, v in self.
567                                     +         items()
568                                     +         return into_c(self.items())
569
570                                     -         return into_c(k, v) for k, v in self.
571                                     +         items()
572                                     +         return into_c(self.items())
573
574                                     -         return into_c(k, v) for k, v in self.
575                                     +         items()
576                                     +         return into_c(self.items())
577
578                                     -         return into_c(k, v) for k, v in self.
579                                     +         items()
580                                     +         return into_c(self.items())
581
582                                     -         return into_c(k, v) for k, v in self.
583                                     +         items()
584                                     +         return into_c(self.items())
585
586                                     -         return into_c(k, v) for k, v in self.
587                                     +         items()
588                                     +         return into_c(self.items())
589
590                                     -         return into_c(k, v) for k, v in self.
591                                     +         items()
592                                     +         return into_c(self.items())
593
594                                     -         return into_c(k, v) for k, v in self.
595                                     +         items()
596                                     +         return into_c(self.items())
597
598                                     -         return into_c(k, v) for k, v in self.
599                                     +         items()
600                                     +         return into_c(self.items())
601
602                                     -         return into_c(k, v) for k, v in self.
603                                     +         items()
604                                     +         return into_c(self.items())
605
606                                     -         return into_c(k, v) for k, v in self.
607                                     +         items()
608                                     +         return into_c(self.items())
609
610                                     -         return into_c(k, v) for k, v in self.
611                                     +         items()
612                                     +         return into_c(self.items())
613
614                                     -         return into_c(k, v) for k, v in self.
615                                     +         items()
616                                     +         return into_c(self.items())
617
618                                     -         return into_c(k, v) for k, v in self.
619                                     +         items()
620                                     +         return into_c(self.items())
621
622                                     -         return into_c(k, v) for k, v in self.
623                                     +         items()
624                                     +         return into_c(self.items())
625
626                                     -         return into_c(k, v) for k, v in self.
627                                     +         items()
628                                     +         return into_c(self.items())
629
630                                     -         return into_c(k, v) for k, v in self.
631                                     +         items()
632                                     +         return into_c(self.items())
633
634                                     -         return into_c(k, v) for k, v in self.
635                                     +         items()
636                                     +         return into_c(self.items())
637
638                                     -         return into_c(k, v) for k, v in self.
639                                     +         items()
640                                     +         return into_c(self.items())
641
642                                     -         return into_c(k, v) for k, v in self.
643                                     +         items()
644                                     +         return into_c(self.items())
645
646                                     -         return into_c(k, v) for k, v in self.
647                                     +         items()
648                                     +         return into_c(self.items())
649
650                                     -         return into_c(k, v) for k, v in self.
651                                     +         items()
652                                     +         return into_c(self.items())
653
654                                     -         return into_c(k, v) for k, v in self.
655                                     +         items()
656                                     +         return into_c(self.items())
657
658                                     -         return into_c(k, v) for k, v in self.
659                                     +         items()
660                                     +         return into_c(self.items())
661
662                                     -         return into_c(k, v) for k, v in self.
663                                     +         items()
664                                     +         return into_c(self.items())
665
666                                     -         return into_c(k, v) for k, v in self.
667                                     +         items()
668                                     +         return into_c(self.items())
669
670                                     -         return into_c(k, v) for k, v in self.
671                                     +         items()
672                                     +         return into_c(self.items())
673
674                                     -         return into_c(k, v) for k, v in self.
675                                     +         items()
676                                     +         return into_c(self.items())
677
678                                     -         return into_c(k, v) for k, v in self.
679                                     +         items()
680                                     +         return into_c(self.items())
681
682                                     -         return into_c(k, v) for k, v in self.
683                                     +         items()
684                                     +         return into_c(self.items())
685
686                                     -         return into_c(k, v) for k, v in self.
687                                     +         items()
688                                     +         return into_c(self.items())
689
690                                     -         return into_c(k, v) for k, v in self.
691                                     +         items()
692                                     +         return into_c(self.items())
693
694                                     -         return into_c(k, v) for k, v in self.
695                                     +         items()
696                                     +         return into_c(self.items())
697
698                                     -         return into_c(k, v) for k, v in self.
699                                     +         items()
700                                     +         return into_c(self.items())
701
702                                     -         return into_c(k, v) for k, v in self.
703                                     +         items()
704                                     +         return into_c(self.items())
705
706                                     -         return into_c(k, v) for k, v in self.
707                                     +         items()
708                                     +         return into_c(self.items())
709
710                                     -         return into_c(k, v) for k, v in self.
711                                     +         items()
712                                     +         return into_c(self.items())
713
714                                     -         return into_c(k, v) for k, v in self.
715                                     +         items()
716                                     +         return into_c(self.items())
717
718                                     -         return into_c(k, v) for k, v in self.
719                                     +         items()
720                                     +         return into_c(self.items())
721
722                                     -         return into_c(k, v) for k, v in self.
723                                     +         items()
724                                     +         return into_c(self.items())
725
726                                     -         return into_c(k, v) for k, v in self.
727                                     +         items()
728                                     +         return into_c(self.items())
729
730                                     -         return into_c(k, v) for k, v in self.
731                                     +         items()
732                                     +         return into_c(self.items())
733
734                                     -         return into_c(k, v) for k, v in self.
735                                     +         items()
736                                     +         return into_c(self.items())
737
738                                     -         return into_c(k, v) for k, v in self.
739                                     +         items()
740                                     +         return into_c(self.items())
741
742                                     -         return into_c(k, v) for k, v in self.
743                                     +         items()
744                                     +         return into_c(self.items())
745
746                                     -         return into_c(k, v) for k, v in self.
747                                     +         items()
748                                     +         return into_c(self.items())
749
750                                     -         return into_c(k, v) for k, v in self.
751                                     +         items()
752                                     +         return into_c(self.items())
753
754                                     -         return into_c(k, v) for k, v in self.
755                                     +         items()
756                                     +         return into_c(self.items())
757
758                                     -         return into_c(k, v) for k, v in self.
759                                     +         items()
760                                     +         return into_c(self.items())
761
762                                     -         return into_c(k, v) for k, v in self.
763                                     +         items()
764                                     +         return into_c(self.items())
765
766                                     -         return into_c(k, v) for k, v in self.
767                                     +         items()
768                                     +         return into_c(self.items())
769
770                                     -         return into_c(k, v) for k, v in self.
771                                     +         items()
772                                     +         return into_c(self.items())
773
774                                     -         return into_c(k, v) for k, v in self.
775                                     +         items()
776                                     +         return into_c(self.items())
777
778                                     -         return into_c(k, v) for k, v in self.
779                                     +         items()
780                                     +         return into_c(self.items())
781
782                                     -         return into_c(k, v) for k, v in self.
783                                     +         items()
784                                     +         return into_c(self.items())
785
786                                     -         return into_c(k, v) for k, v in self.
787                                     +         items()
788                                     +         return into_c(self.items())
789
790                                     -         return into_c(k, v) for k, v in self.
791                                     +         items()
792                                     +         return into_c(self.items())
793
794                                     -         return into_c(k, v) for k, v in self.
795                                     +         items()
796                                     +         return into_c(self.items())
797
798                                     -         return into_c(k, v) for k, v in self.
799                                     +         items()
800                                     +         return into_c(self.items())
801
802                                     -         return into_c(k, v) for k, v in self.
803                                     +         items()
804                                     +         return into_c(self.items())
805
806                                     -         return into_c(k, v) for k, v in self.
807                                     +         items()
808                                     +         return into_c(self.items())
809
810                                     -         return into_c(k, v) for k, v in self.
811                                     +         items()
812                                     +         return into_c(self.items())
813
814                                     -         return into_c(k, v) for k, v in self.
815                                     +         items()
816                                     +         return into_c(self.items())
817
818                                     -         return into_c(k, v) for k, v in self.
819                                     +         items()
820                                     +         return into_c(self.items())
821
822                                     -         return into_c(k, v) for k, v in self.
823                                     +         items()
824                                     +         return into_c(self.items())
825
826                                     -         return into_c(k, v) for k, v in self.
827                                     +         items()
828                                     +         return into_c(self.items())
829
830                                     -         return into_c(k, v) for k, v in self.
831                                     +         items()
832                                     +         return into_c(self.items())
833
834                                     -         return into_c(k, v) for k, v in self.
835                                     +         items()
836                                     +         return into_c(self.items())
837
838                                     -         return into_c(k, v) for k, v in self.
839                                     +         items()
840                                     +         return into_c(self.items())
841
842                                     -         return into_c(k, v) for k, v in self.
843                                     +         items()
844                                     +         return into_c(self.items())
845
846                                     -         return into_c(k, v) for k, v in self.
847                                     +         items()
848                                     +         return into_c(self.items())
849
850                                     -         return into_c(k, v) for k, v in self.
851                                     +         items()
852                                     +         return into_c(self.items())
853
854                                     -         return into_c(k, v) for k, v in self.
855                                     +         items()
856                                     +         return into_c(self.items())
857
858                                     -         return into_c(k, v) for k, v in self.
859                                     +         items()
860                                     +         return into_c(self.items())
861
862                                     -         return into_c(k, v) for k, v in self.
863                                     +         items()
864                                     +         return into_c(self.items())
865
866                                     -         return into_c(k, v) for k, v in self.
867                                     +         items()
868                                     +         return into_c(self.items())
869
870                                     -         return into_c(k, v) for k, v in self.
871                                     +         items()
872                                     +         return into_c(self.items())
873
874                                     -         return into_c(k, v) for k, v in self.
875                                     +         items()
876                                     +         return into_c(self.items())
877
878                                     -         return into_c(k, v) for k, v in self.
879                                     +         items()
880                                     +         return into_c(self.items())
881
882                                     -         return into_c(k, v) for k, v in self.
883                                     +         items()
884                                     +         return into_c(self.items())
885
886                                     -         return into_c(k, v) for k, v in self.
887                                     +         items()
888                                     +         return into_c(self.items())
889
890                                     -         return into_c(k, v) for k, v in self.
891                                     +         items()
892                                     +         return into_c(self.items())
893
894                                     -         return into_c(k, v) for k, v in self.
895                                     +         items()
896                                     +         return into_c(self.items())
897
898                                     -         return into_c(k, v) for k, v in self.
899                                     +         items()
900                                     +         return into_c(self.items())
901
902                                     -         return into_c(k, v) for k, v in self.
903                                     +         items()
904                                     +         return into_c(self.items())
905
906                                     -         return into_c(k, v) for k, v in self.
907                                     +         items()
908                                     +         return into_c(self.items())
909
910                                     -         return into_c(k, v) for k, v in self.
911                                     +         items()
912                                     +         return into_c(self.items())
913
914                                     -         return into_c(k, v) for k, v in self.
915                                     +         items()
916                                     +         return into_c(self.items())
917
918                                     -         return into_c(k, v) for k, v in self.
919                                     +         items()
920                                     +         return into_c(self.items())
921
922                                     -         return into_c(k, v) for k, v in self.
923                                     +         items()
924                                     +         return into_c(self.items())
925
926                                     -         return into_c(k, v) for k, v in self.
927                                     +         items()
928                                     +         return into_c(self.items())
929
930                                     -         return into_c(k, v) for k, v in self.
931                                     +         items()
932                                     +         return into_c(self.items())
933
934                                     -         return into_c(k, v) for k, v in self.
935                                     +         items()
936                                     +         return into_c(self.items())
937
938                                     -         return into_c(k, v) for k, v in self.
939                                     +         items()
940                                     +         return into_c(self.items())
941
942                                     -         return into_c(k, v) for k, v in self.
943                                     +         items()
944                                     +         return into_c(self.items())
945
946                                     -         return into_c(k, v) for k, v in self.
947                                     +         items()
948                                     +         return into_c(self.items())
949
950                                     -         return into_c(k, v) for k, v in self.
951                                     +         items()
952                                     +         return into_c(self.items())
953
954                                     -         return into_c(k, v) for k, v in self.
955                                     +         items()
956                                     +         return into_c(self.items())
957
958                                     -         return into_c(k, v) for k, v in self.
959                                     +         items()
960                                     +         return into_c(self.items())
961
962                                     -         return into_c(k, v) for k, v in self.
963                                     +         items()
964                                     +         return into_c(self.items())
965
966                                     -         return into_c(k, v) for k, v in self.
967                                     +         items()
968                                     +         return into_c(self.items())
969
970                                     -         return into_c(k, v) for k, v in self.
971                                     +         items()
972                                     +         return into_c(self.items())
973
974                                     -         return into_c(k, v) for k, v in self.
975                                     +         items()
976                                     +         return into_c(self.items())
977
978                                     -         return into_c(k, v) for k, v in self.
979                                     +         items()
980                                     +         return into_c(self.items())
981
982                                     -         return into_c(k, v) for k, v in self.
983                                     +         items()
984                                     +         return into_c(self.items())
985
986                                     -         return into_c(k, v) for k, v in self.
987                                     +         items()
988                                     +         return into_c(self.items())
989
990                                     -         return into_c(k, v) for k, v in self.
991                                     +         items()
992                                     +         return into_c(self.items())
993
994                                     -         return into_c(k, v) for k, v in self.
995                                     +         items()
996                                     +         return into_c(self.items())
997
998                                     -         return into_c(k, v) for k, v in self.
999                                     +         items()
1000                                    +         return into_c(self.items())

```

Figure 9: **Left:** Expert’s edit on pandas-dev__pandas-50089, optimizing Series.to_dict by replacing a key-value pair generator with a items() view, eliminating per-element tuple allocation ($1.38\times$ speedup). **Right:** GPT-5 MINI (OPENHANDS) converts the underlying array to a Python list, zipping with the index to reduce Python-level boxing when iterating ($1.98\times$ speedup).

Manually annotated workloads outperform LM generation. Using our evaluation harness, we also study how well LMs can generate performance workloads. We compare the runtime improvement of each expert patch under two workloads: (i) an LM-generated (GEMINI 2.5 FLASH) workload produced from the gold patch and relevant files, and (ii) SWE-FFICIENCY’s manually annotated workload (Stage 4, Fig. 2). Our annotations show stronger performance deltas 76% of the time, with 47% of LM workloads showing no significant speedups. Since performance engineering involves both bottleneck workload identification and code optimization, we show how SWE-FFICIENCY can be further used to probe performance understanding in LMs. For more details, see Appendix I.

5 RELATED WORK

Foundational optimization and synthesis. Classic superoptimization approaches examined code-to-code transformations (Massalin, 1987; Bansal & Aiken, 2006; Schkufza et al., 2013; Solar-Lezama, 2008; Torlak & Bodik, 2014). Profile-guided methods (e.g., Graham et al., 1982; Pettis & Hansen, 1990) and RL-for-performance (e.g., AlphaDev (Mankowitz et al., 2023)) added steerability into code edits. However, these lines of work prioritize transformation quality and search, not the evaluation scaffolding needed for repo-scale, regression-free, workload improvement.

Function-level efficiency benchmarks. MERCURY (Du et al., 2024), EFFIBENCH (Huang et al., 2024), and PIE (Shypula et al., 2024) quantify how often model-generated functions are slower than human references and study feedback- and goal-conditioned improvement. ECCO (Waghjale et al., 2024) emphasizes the necessity of correctness-preserving edits. Domain-focused work such as KERNELBENCH (GPU kernels; Ouyang et al. (2025)) and ALGOTUNE (algorithmic redesign; Press et al. (2025)) further stress wall-clock runtime as the metric. These function level efficiency benchmarks hold utilized libraries fixed and modify user-level workloads to be more efficient. We note that SWE-fficiency instead focuses on a “library maintainer” view of performance optimization, keeping the user workload fixed and editing the library to improve runtime.

Repository-scale SWE benchmarks. Benchmarks like SWE-BENCH (Jimenez et al., 2024; Yang et al., 2025), COMMIT0 (Zhao et al., 2024), and SWT-BENCH (Mündler et al., 2024) established that long-horizon reasoning over code repos is substantially harder than snippet tasks, but they mostly target fixing bugs, developing features and writing tests, rather than performance. Agentic systems (e.g. SWE-agent (Yang et al., 2024); OpenHands, (Wang et al., 2025)) supply the tooling to navigate, edit, run, and profile codebases, improving long-horizon outcomes.

Repository-level performance datasets. Closer to our setting, GSO (Shetty et al., 2025) and SWE-PERF (Fan et al., 2025) curate tasks from GitHub commits and evaluate repo-level runtime. SWE-PERF reuses repository unit tests for both correctness and performance and instructs agents to optimize specified functions; GSO employs both LM-generated correctness tests and performance

workloads, providing a correctness oracle but not exposing performance workloads to agents. While suitable for measuring speedups, these designs insufficiently assess the *localization* skills central to performance engineering—*characterizing* a workload, *localizing* bottlenecks and edits, and *localizing tests* by discovering in-repo unit tests—capabilities our benchmark targets directly.

6 DISCUSSION

Limitations. SWE-FFICIENCY is primarily Python/Cython changes across nine widely-used libraries; extending to lower-level stacks (C/C++/Rust) requires generalizing our coverage test selection and adding language specific build awareness. However, prior Python-only benchmarks (e.g. HumanEval, SWE-bench) have lead to accelerated progress in their respective research directions, and we believe that SWE-FFICIENCY can similarly motivate the research community. Finally, while we focus on a controlled CPU-only setup, scaling to longer-running workloads and heterogeneous hardware would further evaluate agent planning and measurement. Our curation and measurement methodology, prebuilt containers, performance isolation provide a solid foundation to build upon.

Conclusion. We present SWE-FFICIENCY, a repo-level benchmark of 498 optimization tasks across nine widely used repositories. Each task combines a performance workload, an expert patch with a significant speedup, and correctness tests covering the expert diff, enabling evaluation of *pass-to-pass* optimization. Our pipeline rigorously combines regression and AST filters, coverage-guided test selection, manual workload annotation, and reproducibility checks. We present our metric, speedup ratio, for expert parity comparison and find that current agents remain well below expert performance. Our task containerization integrates with open agent frameworks, and we expose qualitative gaps with LM agents like mislocalization and shortcut bias. Our benchmark motivates long-term progress towards autonomous performance engineering and agentic-first codebases.

7 ETHICS STATEMENT

SWE-FFICIENCY is collected entirely from public repositories with licenses that permit software usage that our contributions are in accordance with. Details of the licenses are included in Table 4. We do not collect information about GitHub pull request authors during data collection or evaluation, and SWE-FFICIENCY does not use GitHub data beyond what is available via public API. Our work did not involve any human subject participation: we did not crowdsource or recruit human task workers for any part of SWE-FFICIENCY. For instance environment setup and annotation, the authors conducted all manual and semi-manual tasks. For the benchmark release, we plan to open source the SWE-FFICIENCY task instances, task collection and evaluation infrastructure, and the experimental results and model trajectories from the paper. We will also clearly document each component according to best practices and include channels for communication to engage the community to improve SWE-FFICIENCY.

8 REPRODUCIBILITY STATEMENT

We provide our codebase and all configuration details for the evaluation environment, including container images, CPU pinning, and memory limits; see Appendix D for reproducibility techniques and Appendix E for agent prompts and harness specifics. Our dataset construction steps, filters, and thresholds are documented in 2.1 and Appendix C; example workloads are shown in Appendix B.2. All instances are containerized with pinned and pre-installed dependencies; we prebuild images to avoid run-to-run variance. We report exact evaluation settings in Section 3 and Appendix E. The benchmark artifacts (containers, instance manifests, and harness package) are provided in the supplemental materials for review, and we plan to release a PyPI package and leaderboard website for easy community usage.

REFERENCES

Anthropic. Claude 4.1 Opus, 2025a. URL <https://www.anthropic.com/news/claude-sonnet-4-5>.

- 540 Anthropic. Claude 3.7 Sonnet, 2025b. URL <https://www.anthropic.com/news/claude-3-7-sonnet>.
- 541
- 542
- 543 Anthropic. Claude 4.5 Sonnet, 2025c. URL <https://www.anthropic.com/news/claude-sonnet-4-5>.
- 544
- 545 Sorav Bansal and Alex Aiken. Automatic generation of peephole superoptimizers. In *ASPLOS*, pp. 394–403, 2006. URL <https://theory.stanford.edu/~sbansal/pubs/asplos06.pdf>.
- 546
- 547
- 548
- 549 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- 550
- 551
- 552
- 553
- 554
- 555
- 556
- 557
- 558
- 559
- 560 DeepSeek. DeepSeek V3.1, 2025. URL <https://api-docs.deepseek.com/news/news250821>.
- 561
- 562
- 563 Mingzhe Du, Anh Tuan Luu, Bin Ji, Qian Liu, and See-Kiong Ng. Mercury: A code efficiency benchmark for code large language models. In *NeurIPS Datasets and Benchmarks*, 2024. URL <https://arxiv.org/abs/2402.07844>.
- 564
- 565
- 566 Lieven Eeckhout. R.i.p. geomean speedup use equal-work (or equal-time) harmonic mean speedup instead. *IEEE Computer Architecture Letters*, 23(1):78–82, 2024. doi: 10.1109/LCA.2024.3361925.
- 567
- 568
- 569
- 570 Zhijie Fan, Yiming Huang, Zejian Yuan, Zejun Ma, Qian Liu, et al. SWE-Perf: Can language models optimize code performance on real-world repositories? *arXiv*, 2025. URL <https://arxiv.org/abs/2507.12415>.
- 571
- 572
- 573 Google. Gemini 2.5 Flash, 2025a. URL <https://deepmind.google/models/gemini/flash/>.
- 574
- 575
- 576 Google. Gemini 2.5 Pro, 2025b. URL <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/>.
- 577
- 578
- 579 Susan L. Graham, Peter B. Kessler, and Marshall K. McKusick. gprof: A call graph execution profiler. In *SIGPLAN Symposium on Compiler Construction*, pp. 120–126, 1982. URL <https://web.eecs.umich.edu/~weimerw/2012-4610/reading/graham-gprof.pdf>.
- 580
- 581
- 582 Dong Huang, Yuhao Qing, Weiyi Shang, Heming Cui, and Jie M. Zhang. Effibench: Benchmarking the efficiency of automatically generated code. In *NeurIPS Datasets and Benchmarks*, 2024. URL <https://arxiv.org/abs/2402.02037>.
- 583
- 584
- 585
- 586 ISO/IEC. Systems and software engineering – Systems and software quality requirements and evaluation (SQuaRE) – System and software quality models. ISO/IEC 25010, International Organization for Standardization, Geneva, Switzerland, 2011.
- 587
- 588
- 589 Akanksha Jain, Hannah Lin, Carlos Villavieja, Baris Kasikci, Chris Kennelly, Milad Hashemi, and Parthasarathy Ranganathan. Limoncello: Prefetchers for scale. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, ASPLOS '24, pp. 577–590, New York, NY, USA, 2024a. Association for Computing Machinery. ISBN 9798400703867. doi: 10.1145/3620666.3651373. URL <https://doi.org/10.1145/3620666.3651373>.
- 590
- 591
- 592
- 593

- 594 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
595 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
596 evaluation of large language models for code, 2024b. URL [https://arxiv.org/abs/
597 2403.07974](https://arxiv.org/abs/2403.07974).
- 598 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R
599 Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth
600 International Conference on Learning Representations*, 2024. URL [https://openreview.
601 net/forum?id=VTF8yNQM66](https://openreview.net/forum?id=VTF8yNQM66).
- 602 Daniel J. Mankowitz, A. Michi, A. Zhernov, et al. Faster sorting algorithms discovered using deep
603 reinforcement learning. *Nature*, 618:257–263, 2023. doi: 10.1038/s41586-023-06004-9. URL
604 <https://www.nature.com/articles/s41586-023-06004-9.pdf>.
- 605 Henry Massalin. Superoptimizer: A look at the smallest program. In *Proceedings of the
606 2nd Intl. Conf. on Architectural Support for Programming Languages and Operating Systems
607 (ASPLOS-II)*, pp. 122–126, 1987. doi: 10.1145/36206.36194. URL [https://courses.cs.
608 washington.edu/courses/cse501/15sp/papers/massalin.pdf](https://courses.cs.washington.edu/courses/cse501/15sp/papers/massalin.pdf).
- 609 Moonshot AI. Kimi K2, 2025. URL <https://moonshotai.github.io/Kimi-K2/>.
- 610 Niels Müндler, Mark Niklas Mueller, Jingxuan He, and Martin Vechev. SWT-bench: Testing and
611 validating real-world bug-fixes with code agents. In *The Thirty-eighth Annual Conference on
612 Neural Information Processing Systems*, 2024. URL [https://openreview.net/forum?
613 id=9Y8zUO11EQ](https://openreview.net/forum?id=9Y8zUO11EQ).
- 614 OpenAI. GPT-5, 2025a. URL <https://openai.com/gpt-5/>.
- 615 OpenAI. GPT-5 Mini, 2025b. URL <https://openai.com/gpt-5/>.
- 616 Anne Ouyang, Simon Guo, Simran Arora, Alex L. Zhang, William Hu, Christopher Ré, and Azalia
617 Mirhoseini. Kernelbench: Can llms write efficient GPU kernels? *arXiv*, 2025. URL [https://arxiv.org/abs/
618 2502.10517](https://arxiv.org/abs/2502.10517).
- 619 Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe
620 Zhang. Training software engineering agents and verifiers with swe-gym, 2025. URL [https://arxiv.org/abs/
621 2412.21139](https://arxiv.org/abs/2412.21139).
- 622 Karl Pettis and Robert C. Hansen. Profile guided code positioning. In *PLDI*, pp. 16–27, 1990. URL
623 <https://dblp.org/rec/conf/pldi/PettisH90>.
- 624 Ori Press, Brandon Amos, Haoyu Zhao, Yikai Wu, Samuel K. Ainsworth, Dominik Krupke, Patrick
625 Kidger, Touqir Sajed, Bartolomeo Stellato, Jisun Park, Nathanael Bosch, Eli Meril, Albert Steppi,
626 Arman Zharmagambetov, Fangzhao Zhang, David Perez-Pineiro, Alberto Mercurio, Ni Zhan,
627 Talor Abramovich, Kilian Lieret, Hanlin Zhang, Shirley Huang, Matthias Bethge, and Ofir Press.
628 Algotune: Can language models speed up general-purpose numerical programs?, 2025. URL
629 <https://arxiv.org/abs/2507.15887>.
- 630 Eric Schkufza, Rahul Sharma, and Alex Aiken. Stochastic superoptimization. In *ASPLOS*, pp.
631 305–316, 2013. URL [https://theory.stanford.edu/~aiken/publications/
632 papers/asplos13.pdf](https://theory.stanford.edu/~aiken/publications/papers/asplos13.pdf).
- 633 Manish Shetty, Naman Jain, Jinjian Liu, Vijay Kethanaboyina, Koushik Sen, and Ion Stoica.
634 Gso: Challenging software optimization tasks for evaluating swe-agents, 2025. URL [https://arxiv.org/abs/
635 2505.23671](https://arxiv.org/abs/2505.23671).
- 636 Alexander Shypula, Aman Madaan, Yimeng Zeng, Uri Alon, Jacob Gardner, Milad Hashemi, Gra-
637 ham Neubig, Parthasarathy Ranganathan, Osbert Bastani, and Amir Yazdanbakhsh. Learning
638 performance-improving code edits. In *ICLR*, 2024. URL [https://arxiv.org/pdf/2302.
639 07867](https://arxiv.org/pdf/2302.07867).
- 640 J. E. Smith. Characterizing computer performance with a single number. *Commun. ACM*, 31(10):
641 1202–1206, October 1988. ISSN 0001-0782. doi: 10.1145/63039.63043. URL [https://doi.
642 org/10.1145/63039.63043](https://doi.org/10.1145/63039.63043).

- 648 Armando Solar-Lezama. *Program Synthesis by Sketching*. PhD thesis, EECS Department, Univer-
649 sity of California, Berkeley, 2008. URL [https://www2.eecs.berkeley.edu/Pubs/
650 TechRpts/2008/EECS-2008-177.html](https://www2.eecs.berkeley.edu/Pubs/TechRpts/2008/EECS-2008-177.html).
- 651 Emina Torlak and Rastislav Bodik. A lightweight symbolic virtual machine for solver-aided host
652 languages. In *PLDI*, pp. 530–541, 2014. URL [https://homes.cs.washington.edu/
653 ~bodik/ucb/Files/2014/rosette-pldi2014.pdf](https://homes.cs.washington.edu/~bodik/ucb/Files/2014/rosette-pldi2014.pdf).
- 655 Siddhant Waghjale, Vishruth Veerendranath, Zora Zhiruo Wang, and Daniel Fried. ECCO: Can we
656 improve model-generated code efficiency without sacrificing functional correctness? *arXiv*, 2024.
657 URL <https://arxiv.org/abs/2407.14044>.
- 658
- 659 Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan,
660 Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng,
661 Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert
662 Brennan, Hao Peng, Heng Ji, and Graham Neubig. Openhands: An open platform for AI soft-
663 ware developers as generalist agents. In *The Thirteenth International Conference on Learning
664 Representations*, 2025. URL <https://openreview.net/forum?id=OJd3ayDDoF>.
- 665 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan,
666 and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering.
667 In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL
668 <https://arxiv.org/abs/2405.15793>.
- 669 John Yang, Carlos E Jimenez, Alex L Zhang, Kilian Lieret, Joyce Yang, Xindi Wu, Ori Press,
670 Niklas Muennighoff, Gabriel Synnaeve, Karthik R Narasimhan, Diyi Yang, Sida Wang, and Ofir
671 Press. SWE-bench multimodal: Do AI systems generalize to visual software domains? In
672 *The Thirteenth International Conference on Learning Representations*, 2025. URL [https:
673 //openreview.net/forum?id=riTiq3i21b](https://openreview.net/forum?id=riTiq3i21b).
- 674
- 675 Z.ai. GLM-4.6, 2025. URL <https://docs.z.ai/guides/llm/glm-4.6>.
- 676
- 677 Wenting Zhao, Nan Jiang, Celine Lee, Justin T Chiu, Claire Cardie, Matthias Gallé, and Alexan-
678 der M Rush. Commit0: Library generation from scratch, 2024. URL [https://arxiv.org/
679 abs/2412.01769](https://arxiv.org/abs/2412.01769).

681 A LLM USAGE

682

683 Language models were used to polish writing, help with grammatical errors and typos, and to help
684 check with compliance against ICLR’s author guide. Beyond the LM usage in our benchmark eval-
685 uations and experiments, they were not used in any other part of writing this work.

687 B ADDITIONAL DATASET INFORMATION

688

689 In this section, we provide more details on the dataset summary and distribution of SWE-
690 EFFICIENCY. We verify that all repositories used have permissive licenses, allowing for data mining
691 and inclusion into the SWE-EFFICIENCY benchmark, as shown in Table 4. Figure 10 and 11 show
692 additional details on the distribution of task instances, repositories, and corresponding information.
693 We observe that compared to SWE-bench, SWE-EFFICIENCY gold patches are larger on average and
694 have significantly larger numbers of related tests (as checking for correctness regression is stricter
695 than SWE-bench’s pass criteria of issue resolution).

696

697 We also note that SWE-bench does not contain any code optimization tasks as noted in Table 1.
698 Firstly, as we select for changes that do not introduce test file changes, this disqualifies any of SWE-
699 bench instances from passing Stage 2 of Figure 2 (attribute filtering). Furthermore, we also run
700 SWE-bench instances through our performance keyword pipeline and randomly sample 50 of the
701 581 resulting instances for human review: all of these instances had issue statements or hints text
that clearly show the change introduces new behavior (which is tested by the SWE-bench instance’s
test_patch).

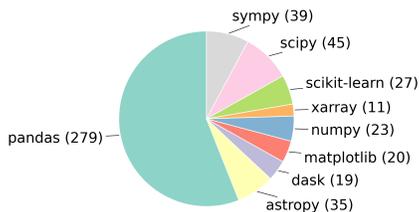


Figure 10: Repository distribution of task instances in the SWE-FFICIENCY dataset.

Category	Metric	Mean	Max
Codebase	# Instances	498	
	# Repos	9	
Workload	# of Lines	25.47	180
	Runtime (s)	4.47	257.09
Gold Patch	# Lines edited	49.1	2445
	# Files edited	2.2	12
	Speedup	2.64×	249k×
Tests	# Pass to Pass	54k	222k

Figure 11: Additional summary statistics for SWE-FFICIENCY dataset. Arithmetic mean is used in all cases, except for speedup (harmonic mean).

B.1 REPOSITORY DESCRIPTIONS AND PERMISSIVE LICENSES

All SWE-FFICIENCY task instances are scraped from public GitHub repos with permissive licenses via the public GitHub API as stated in our Ethics Statement. Repository specific licenses are shown in Table 4, where links to custom licenses are provided inline.

Table 4: SWE-FFICIENCY GitHub repositories, package description, and their permissive licenses.

Repository	Description	License
astropy/astropy	Astronomy and astrophysics core library	BSD-3-Clause
dask/dask	Parallel computing library for analytics	BSD-3-Clause
matplotlib/matplotlib	Plotting and graphics library	Custom
numpy/numpy	Core scientific computing library	Custom
pandas-dev/pandas	Core data analysis library	BSD-3-Clause
pydata/xarray	Multi-dimensional array library	Apache 2.0
scikit-learn/scikit-learn	Machine learning in Python	BSD-3-Clause
scipy/scipy	Package for math, science, and engineering	BSD-3-Clause
sympy/sympy	Computer algebra system written in Python	Custom

B.2 EXAMPLE PERFORMANCE WORKLOADS

We provide some examples of the performance workloads associated with each task instance. To recap, each performance workload consists of (i) necessary imports, (ii) an optional `setup` function, which sets up work that should not be runtime benchmarked, (iii) a `workload` function, which runs some repository functionality and runtime to be measured, and (iv) performance measurement harness code. Each problem runs the workload and setup multiple times to generate a distribution of runtimes, from which we compute the mean and standard deviation. During our dataset curation pipeline, we reject task instances and workloads that fail to show statistically significant improvements in the execution validation stage (Appendix C.5).

Performance Workload for `dask__dask-6293`

```
import timeit
import statistics

import dask
import dask.array as da
import numpy as np

def setup():
    global stacked, sub_arrays
```

```

756
757     sub_arrays = [
758         da.from_delayed(
759             dask.delayed(np.zeros)((100000,)), dtype="int64"),
760             shape=(100000,),
761             dtype="int64",
762             name=idx
763         )
764     ]
765     stacked = da.stack(sub_arrays)
766
767     def workload():
768         global stacked, sub_arrays
769         for i in range(len(sub_arrays)):
770             stacked[i]
771
772     runtimes = timeit.repeat(workload, number=1, repeat=5, setup=setup)
773
774     # Print runtime mean and std deviation.
775     print("Mean:", statistics.mean(runtimes))
776     print("Std Dev:", statistics.stdev(runtimes))
777

```

Performance Workload for `numpy__numpy`-11720

```

778     import timeit
779     import statistics
780
781     import numpy as np
782
783     b = np.random.random((5, 2))
784     t = np.random.random((5, 5, 2))
785     p = np.random.random((2, 5))
786
787     def workload():
788         out = np.einsum('ij,ixy,ji->xy', b, t, p)
789
790     runtimes = timeit.repeat(workload, number=100, repeat=10**4)
791
792     # Print runtime mean and std deviation.
793     print("Mean:", statistics.mean(runtimes))
794     print("Std Dev:", statistics.stdev(runtimes))
795

```

Performance Workload for `scipy__scipy`-12001

```

796     import timeit
797     import statistics
798     import numpy as np
799     from scipy.stats import maxwell
800
801     def setup():
802         global data
803         data = maxwell.rvs(loc=2.0, scale=3.0, size=100000,
804             ↪ random_state=42)
805
806     def workload():
807         global data
808         _ = maxwell.fit(data)
809
810     runtimes = timeit.repeat(workload, number=1, repeat=200, setup=setup)
811

```

```

810
811 print("Mean:", statistics.mean(runtimes[-100:]))
812 print("Std Dev:", statistics.stdev(runtimes[-100:]))

```

B.3 LM CLASSIFICATION OF GOLD PATCH EDIT TYPES

To examine the diversity of gold (expert) patch optimizations that SWE-fficiency submissions are graded against (rightmost plot in Figure 3), we prompt GEMINI 2.5 FLASH with the following task prompt to extract optimization categories. We then randomly sampled 50 LM classifications for manual review and confirmed the high level categorization and explanation for each. We emphasize that during SWE-FFIENCY evaluation, LM agents are free to make any optimization desired to improve the performance workload, and we use this categorization strictly to show the diversity of our collected dataset.

Performance Diff Classification Task Prompt

```

825 You are an excellent performance engineer. Given a code diff and an
826 ↪ affected performance workload that shows a speedup as a result of
827 ↪ this diff, output a single high-level performance bucket, the
828 ↪ concrete signals from the diff that justify that bucket, a
829 ↪ mechanism-level explanation of why the specific code edits improve
830 ↪ performance, and a confidence score. Prefer software-side
831 ↪ mechanisms; ignore hardware/microarchitecture unless explicitly
832 ↪ cited.
833
834 ## Inputs
835
836 * Performance workload and code diffs.
837
838 ## Buckets (choose exactly one `classification`)
839
840 1. Algorithmic / Data Structure Improvements -- Better asymptotic
841 ↪ complexity or more suitable data structures; removes redundant
842 ↪ passes.
843 2. Memory Efficiency & Management -- Fewer allocations/copies;
844 ↪ pooling/reuse; layout/locality changes; alignment; `reserve()`;
845 ↪ SoA/AoS.
846 3. Concurrency & Parallelism -- Threading/async; work
847 ↪ partitioning; lock scope/structure; atomics; SIMD/vectorization.
848 4. I/O and Storage Efficiency -- Fewer syscalls;
849 ↪ buffering/batching; async I/O; (de)serialization changes; payload
850 ↪ trimming.
851 5. Code Simplification / Dead-Code Elimination -- Early exits;
852 ↪ pruning logging/asserts on hot paths; removing unnecessary
853 ↪ work/branches.
854 6. Compiler / Build / Low-level Tuning -- Flags (LTO/PGO),
855 ↪ inlining hints, intrinsics, UB fixes enabling optimization,
856 ↪ branch hints.
857 7. Configuration / Parameter Tuning --
858 ↪ Constants/thresholds/buffer sizes, thread-pool/GC settings,
859 ↪ feature flags altering performance behavior.
860 8. Caching & Reuse -- Memoization, caches, reuse of precomputed
861 ↪ artifacts/results, avoiding repeated expensive calls.
862 9. Unknown / Not Enough Information -- Claimed speedup but
863 ↪ mechanism not inferable from available changes.
864
865 ## Secondary Tags (optional)
866
867 * Other relevant buckets or keywords, if any (e.g., "Memory
868 ↪ Efficiency & Management" and "Caching & Reuse" could both apply).
869 * These can be more specific, e.g., "memoization", "lock-free",
870 ↪ "buffered I/O", but try to use standard terms where possible.

```

```

864
865
866   ### Disambiguation rules
867
868   * Algorithmic vs Caching: If an algorithm was fundamentally
869   ↪ changed and a cache was added as a helper, choose
870   ↪ Algorithmic; add `memoization` in `mechanism_signals`.
871   * Concurrency vs Algorithmic: If parallelism is added without
872   ↪ changing the algorithm, choose Concurrency & Parallelism.
873   * I/O vs Memory: If copies were removed primarily to cut syscalls
874   ↪ or shrink payloads, choose I/O; if focused on
875   ↪ allocation/locality/pressure, choose Memory.
876   * Compiler/Build: Choose Compiler / Build when source logic
877   ↪ is the same but flags/hints/toolchain changed.
878   * Benchmark-only: Choose Workload/Benchmark-Only when only
879   ↪ harness/warmup/affinity/timers changed.
880
881   ## What to extract as "signals"
882
883   Short, concrete phrases tied to the diff, e.g.:
884
885   * containers/algos: "`vector`→`unordered_map`", "removed nested
886   ↪ loop", "added binary search", "streaming parse"
887   * memory: "added `reserve()`", "object pool", "moved to stack
888   ↪ allocation", "SoA layout"
889   * concurrency: "introduced thread pool", "reduced lock scope",
890   ↪ "lock-free queue", "SIMD intrinsics"
891   * I/O: "batched writes", "buffered reader", "protobuf→flat
892   ↪ serialization", "compression level tuned", "fewer syscalls"
893   * simplification: "early return before parse", "pruned logging on hot
894   ↪ path", "deleted dead branch"
895   * compiler/build: "enabled LTO/PGO", "added
896   ↪ `inline`/`cold`/`likely`", "UB fix unblocking vectorization"
897   * config: "increased read buffer to 1MB", "thread pool size = cores"
898   * caching: "added LRU cache", "memoized function result"
899   * meta: "only bench harness changed"
900
901   ## Output Requirements (STRICT)
902
903   * Output only the JSON object -- no prose, no Markdown, no code
904   ↪ fences.
905   * Keep `explanation` ≤ 6 sentences and tie it to specific
906   ↪ lines/files/patterns from the diff.
907   * If evidence is weak or ambiguous, use `classification: "Unknown /
908   ↪ Not Enough Information"` and lower `confidence`.
909
910   ### JSON Schema
911
912   ```json
913   {
914     "classification": "<one of the 10 buckets>",
915     "secondary_tags": ["<optional: other relevant buckets or keywords,
916     ↪ if any>"],
917     "mechanism_signals": ["<short phrases pulled from the diff that
918     ↪ justify the classification>"],
919     "affected_components": ["<files/modules/functions inferred from
920     ↪ paths/symbols>"],
921     "explanation": "<mechanism-level rationale grounded in the diff:
922     ↪ what changed, how it reduces
923     ↪ work/contention/latency/allocs/syscalls, and why that maps to
924     ↪ the chosen bucket>",
925     "confidence": "<high|medium|low>"
926   }
927   ...

```

```

918
919
920 ## Final sanity check (do this before emitting JSON)
921
922 1. Have I picked exactly one bucket that best explains the
923   ↪ performance mechanism?
924 2. Do my `mechanism_signals` cite concrete code changes from the diff
925   ↪ that motivate that bucket?
926 3. Is the explanation mechanism-centric and grounded in the edits
927   ↪ (not benchmarks)?
928
929 ---
930
931 ### Mini-examples
932
933 **A. Algorithmic / Data Structure Improvements**
934
935 ```json
936 {
937   "classification": "Algorithmic / Data Structure Improvements",
938   "secondary_tags": ["asymptotic complexity"],
939   "mechanism_signals": ["removed nested O(n^2) scan", "introduced
940     ↪ unordered_set", "added reserve() to avoid rehash"],
941   "affected_components": ["src/import/dedupe.cpp",
942     ↪ "Importer::dedupeRecords"],
943   "explanation": "The patch replaces a quadratic duplicate search
944     ↪ with hash-based membership checks and preallocates the table to
945     ↪ avoid rehash, removing repeated comparisons across the hot
946     ↪ loop.",
947   "confidence": "high"
948 }
949 ```
950
951 **B. Concurrency & Parallelism**
952
953 ```json
954 {
955   "classification": "Concurrency & Parallelism",
956   "secondary_tags": ["parallelism", "lock contention"],
957   "mechanism_signals": ["introduced thread pool", "tile-based work
958     ↪ partitioning", "narrowed mutex scope"],
959   "affected_components": ["decoder/pipeline.cc", "decoder/tiler.cc"],
960   "explanation": "Work is partitioned by tile across a shared pool
961     ↪ and critical sections are reduced to short regions, lowering
962     ↪ contention and enabling parallel execution of the same
963     ↪ algorithm.",
964   "confidence": "high"
965 }
966 ```
967
968 **C. Unknown**
969
970 ```json
971 {
972   "classification": "Unknown / Not Enough Information",
973   "secondary_tags": [],
974   "mechanism_signals": ["broad refactor", "no evident hot-path
975     ↪ edits"],
976   "affected_components": ["loader/*"],
977   "explanation": "Large refactor touches many files without showing
978     ↪ hot-path changes or recognizable performance mechanisms, so the
979     ↪ cause of any improvement cannot be determined from the diff
980     ↪ alone.",
981   "confidence": "low"
982 }
983 ```

```

```
}
...
```

B.4 DATASET SCHEMA

For clarity, we describe the schema and description of our dataset in Table 5. We upload our dataset to HuggingFace ([swefficiency-anon/swefficiency](https://huggingface.co/swefficiency-anon/swefficiency)) for easy community download, and refer readers to our code supplementary material to see how the dataset is directly used during evaluation.

Table 5: SWE-FFICIENCY dataset columns and description.

Column Name	Description
<code>repo</code>	(str) Repository identifier for the task (e.g., <code>owner/repo</code> on GitHub).
<code>instance_id</code>	(str) Unique ID for this dataset instance.
<code>base_commit</code>	(str) Git commit SHA to check out before applying any patches; defines the baseline state under evaluation.
<code>patch</code>	(str) Expert git patch with source-code changes which solves the task and shows a performance optimization. For evaluation purposes, this should never be provided to the agent.
<code>created_at</code>	(str) ISO-8601 timestamp indicating when this instance was created.
<code>version</code>	(str) Repository version string for this instance (used for repository environment building).
<code>environment_setup_commit</code>	Commit SHA (or ref) that pins environment setup artifacts (e.g., dependency files) for reproducible evaluation.
<code>workload</code>	(str) Python workload script used for performance measurement (script/benchmark/entrypoint) See Appendix B.2 for examples..
<code>test_cmd</code>	Shell command prefix used to run the test suite for this repo (e.g., <code>pytest -q</code>).
<code>rebuild_cmd</code>	Shell command to (re)build/reinstall the project between runs for agent usage (e.g., <code>pip install -e .</code>)
<code>image_name</code>	Container image tag/name providing the canonical evaluation environment (OS, toolchain, deps).
<code>covering_tests</code>	(list of str) List of tests paths that exercise the changed code regions (e.g., from coverage or curated mappings). For evaluation purposes, this should never be provided to the agent.
<code>single_thread_tests</code>	List of tests that must run serially (to avoid flakiness or resource contention) during evaluation. For evaluation purposes, this should never be provided to the agent.
<code>PASS_TO_PASS</code>	(list of str) List of test identifiers that pass after the expert edit and are expected to pass after an LM generated edit (regression guard). For evaluation purposes, this should not be provided to the agent.

C ADDITIONAL DETAILS ON DATA COLLECTION PROCEDURE

In this section, we provide more concrete details on the dataset collection procedure explained in Section 2.1 and in Figure 2. To recap, our dataset pipeline is comprised of the following stages (with numbers on the yield after each stage and direct references to appendix subsections):

1. Scrape pull requests from nine widely used open-source Python repositories in data science, machine learning, and high-performance computing—chosen for mature test suites and stringent performance requirements—yielding **96457 PRs** (Appendix C.1).
2. Filter for candidate performance-regression instances by requiring performance-related keywords, excluding PRs that modify tests (and introduce new behavior), and retain only edits that meaningfully change the abstract syntax tree (AST) — criteria that notably targets instances intentionally excluded by SWE-bench due to its test-change filter. After attribute filtering and prior to checking for meaning full changes to the AST, we retain **9257 PRs** (-90.4%) at this stage (Appendix C.2).
3. Construct an executable Docker environment per instance with manually curated, version-pinned dependencies (tests are often version-sensitive), run repository unit tests, and use line coverage to keep only instances with at least one “guarding” test whose executed lines intersect the edited code. We retain **1041 PRs** (-88.8%) up until this stage (Appendix C.3).
4. Annotate each surviving instance with a minimal workload script that reliably exposes the pre/post performance delta (Appendix C.4).
5. Run correctness tests and the performance workload before/after applying the patch, retaining only instances with a statistically significant improvement, recording post-patch test outcomes, and ensuring reproducibility via Docker with 4 CPU cores and 16GB RAM. This yields our final **498 tasks** across 9 repos (Appendix C.5).

C.1 REPO SELECTION AND INSTANCE SCRAPING

We provide more details on Stage 1 from Figure 2 in how we selected repositories and scraping raw instance information. In this stage, using **March 12th, 2025** as a cutoff date, we scrape merged pull requests from nine (9) repositories (astropy, dask, matplotlib, numpy, pandas, scikit-learn, scipy, sympy, and xarray). Note that we also relax the SWE-bench requirement that a valid PR require a linked GitHub issue as a problem statement: from our observation, we see that *many performance optimization PRs are opportunistic and do not always have a linked issue created ahead of time*.

We note that some repositories overlap with SWE-bench and SWE-Gym (Pan et al., 2025), such as astropy, matplotlib, xarray, scikit-learn, sympy, dask, and pandas, while others are exclusive to SWE-EFFICIENCY like numpy and scipy. We examined all the repositories in SWE-bench and SWE-Gym to start and found that most of the repositories in those datasets that are not included in our benchmark contain very few performance related changes: This is due to some of those repositories focusing on general functionality rather than performance optimization specific changes. For example, packages like flask and django in SWE-bench focus on rapid iteration and high-level design rather than optimizing for high degrees of performance, which reflects in their list of merged pull requests (PRs) having very few occurrences of the word `perf`.

C.2 PERFORMANCE REGRESSION ATTRIBUTE FILTERING

As discussed in Section 2.1, we modify the attribute filtering from SWE-bench to prune away instances that are not regression-free performance optimization. Specifically, we keep an instance only if it satisfies the following:

1. *Does not contribute test changes*: We intentionally drop instances if they add test changes, since this indicates, with high likelihood, that new behavior is introduced by that PR. In contrast, SWE-bench selects for the opposite (intentionally selects for new tests, and masks out those tests for instance evaluation), meaning that *our dataset is exactly instance-wise disjoint* with both SWE-bench and SWE-gym. We attribute this to “test-driven development” in general: performance PRs generally do not change inputs or output behavior and thus do not need new tests to guard this new behavior.
2. *Contains performance related keywords or tags*: We check if the pull request meta-data includes any of the following keywords: performance, speedup, speeds up, speed-up, speed up, faster, memory, optimize, optimization, profiling, accelerate, fast, runtime, efficiency, benchmark,

latency, throughput, multithreading, parallel, concurrency, concurrent, profiling, CPU usage, memory usage, resource usage, cache, caching, `timeit`, and `asv`. We also check if pull request has been tagged with any repo specific performance tags and keep those pull requests as well.

3. *PR contains meaningful changes to AST*: We finally check that the PR edit has made meaningful changes to each changed file’s abstract syntax tree (AST) as parsed by `tree-sitter`. This helps us ignore no-op changes that are comment or doc-string only and select more specifically for substantial performance related changes, which are almost guaranteed to require a modification to a code file’s AST.

C.3 IDENTIFYING COVERING CORRECTNESS TESTS

In this third stage, we retain a task instance only if the repository installs with all required dependencies, if the tests execute properly, and if we can identify unit tests that intersect the PR diff via line coverage. Installation is usually the most brittle step: a repo may install successfully yet fail at test time due to mismatched or missing dependencies. In practice, this requires manual curation—pinning versions and resolving transitive constraints—to map each instance to a working environment, beyond the constants provided by SWE-bench and SWE-Gym.

Once tests run cleanly, we execute the full test suite with coverage enabled and record, for each test, the lines of each source file that are executed. For every test file, we align this dynamic coverage with the lines, functions, classes, and modules modified in the PR to determine whether the test intersects the change. We keep a task instance only if at least one unit test intersects the original PR edit. Recall that SWE-bench selected PRs that introduced tests: since we intentionally select tests without test changes, this coverage step is required for us to identify guarding tests in the code repo.

Because our correctness check targets performance edits intended to be semantics-preserving, any check violation must appear on executions that traverse the modified regions. We therefore restrict the necessary correctness tests to those whose coverage intersects the PR diff (aggregated across lines, functions, classes, and modules). This change-focused selection is a conservative form of test-impact analysis: tests that never execute the modified code cannot surface regressions, yet they would inflate wall-clock time and noise in a benchmark setting. Limiting evaluation to intersecting tests preserves detection power for performance regressions, reduces spurious failures from unrelated tests, and yields stable, low-cost runs—making the benchmark practical for repeated use and community adoption.

C.4 ANNOTATING PERFORMANCE WORKLOADS

Given performance-related candidate task instances for which we can easily check that edits maintain correctness of code, we need a way of also grading whether edits improve performance. We explored using an LM generated pipeline to generate workloads (see Appendix I), but found that manual annotation based on GitHub issue PR and issue metadata was a more effective strategy and yielded more realistic workloads (i.e. the same workloads that PR authors used as a baseline to implement their optimization edits). Thus, for each candidate task instance in this stage, we examine its linked GitHub pull request and issue info and generate a workload script which shows a performance delta: we double check these workloads in the next stage to verify the reproducibility and statistical significance of the performance improvements for benchmark inclusion.

Each workload consists of four items: (i) required imports (including `timeit` and `statistics` for computing runtime distributions); (ii) an optional `setup()` function for initializing any parts of the performance workload that should not be measured for runtime; (iii) a `workload()` function, which encapsulates the key functionality of interest to measure; and (iv) timing-specific code to run workloads multiple times to generate consistent runtime distributions for evaluation and analysis. See Appendix B.2 for examples of performance workloads and Appendix I for detailed workload creation instructions.

Notably, we find that automatically extracting workloads that show the performance delta from PR information is difficult. For example, `pandas` PRs #43274, #49596, #59608 each contain at least one (of many) codeblocks with a performance script from the PR author, showing the intended performance delta. However, note that each block uses a different format, timing mechanism, and

1134 method of executing programs, as shown in Figures 12, 13, 14. Unifying these consistently without
 1135 hugely reducing dataset yield is non-trivial, as shown in other works like GSO (Shetty et al., 2025).
 1136 We leave LM-pipeline based approaches to automate this extraction to future work.
 1137

1138 PR Performance Script Codeblock for pandas-dev__pandas-43274

```
1139 In [1]: from asv_bench.benchmarks.indexing import InsertColumns
1140
1141 In [2]: self = InsertColumns()
1142
1143 In [3]: self.setup()
1144
1145 In [4]: %timeit self.time_assign_with_setitem()
1146 27.6 ms ± 68.1 µs per loop (mean ± std. dev. of 7 runs, 10 loops
1147 ↪ each) #master
1148
1149 In [4]: %timeit self.time_assign_with_setitem()
1150 8.6 ms ± 6.43 µs per loop (mean ± std. dev. of 7 runs, 100 loops
1151 ↪ each) #PR
```

1152 Figure 12: This PR uses existing asv benchmarks in the repository and measures performance
 1153 improvement with the timeit command line entypoint.
 1154

1155 PR Performance Script Codeblock for pandas-dev__pandas-49596

```
1156 import pandas as pd
1157 import pandas._testing as tm
1158
1159 vals = pd.Series(tm.rands_array(10, 10**6), dtype="string")
1160 df = pd.DataFrame({"cat": vals.astype("category")})
1161
1162 %timeit df.groupby("cat").size()
1163
1164 1.21 s ± 4.71 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
1165 ↪ <- main
1166 15.1 ms ± 274 µs per loop (mean ± std. dev. of 7 runs, 100 loops
1167 ↪ each) <- PR
```

1168 Figure 13: This PR uses a bespoke workload (non asv benchmark and measures a specific func-
 1169 tionality, again using timeit to measure speedup.
 1170

1171 PR Performance Script Codeblock for pandas-dev__pandas-59608

```
1172 import pandas as pd
1173 import pyarrow as pa
1174 import pyarrow.csv as csv
1175 import time
1176
1177 NUM_ROWS = 10000000
1178 NUM_COLS = 20
1179
1180 # Example Multi-Index DataFrame
1181 df = pd.DataFrame(
1182     {
1183         f"col_{col_idx}": range(col_idx * NUM_ROWS, (col_idx + 1) *
1184 ↪ NUM_ROWS)
1185         for col_idx in range(NUM_COLS)
1186     }
1187
```

```

1188 )
1189 df = df.set_index(["col_0", "col_1"], drop=False)
1190
1191 # Timing Operation A
1192 start_time = time.time()
1193 df.to_csv("file_A.csv", index=False)
1194 end_time = time.time()
1195 print(f"Operation A time: {end_time - start_time} seconds")
1196
1197 # Timing Operation B
1198 start_time = time.time()
1199 df_reset = df.reset_index(drop=True)
1200 df_reset.to_csv("file_B.csv", index=False)
1201 end_time = time.time()
1202 print(f"Operation B time: {end_time - start_time} seconds")

```

Figure 14: This PR uses both a bespoke workload and non-timeit, Python timing functionality to measure speedup.

To ensure a consistent and measurement environment for every PR, we pre-built Docker images with dependencies pre-installed. Each image encapsulated the repository at the specific base commit (i.e., immediately before the PR’s changes were applied) along with all necessary dependencies to execute the project’s code and test suites. Prior to commencing the full annotation task, the two author annotators calibrated their methodology and interpretation. They jointly annotated a set of five (5) example workloads, discussing discrepancies and establishing a consistent standard for what constituted a valid and representative workload. This calibration ensured alignment on the annotation process, including how to interpret PR descriptions, locate relevant code, and structure the final workload script. During annotation, the authors communicated regularly to ensure they stayed calibrated throughout. In total, workload annotation took around 200 author-hours for 1041 instances, 498 of which made it through final execution validation.

Annotator Workload Template

```

1219 # Import statements and one-time global work here
1220 import timeit
1221 import statistics
1222 ...
1223
1224 def setup():
1225     # Any one time work that needs to be done before every new run.
1226     pass
1227
1228 def workload():
1229     # Workload to be measured.
1230     pass
1231
1232 # Fill in number of repeats to get tight error bars
1233 runtimes = timeit.repeat(workload, number=1, repeat=10, setup=setup)
1234
1235 print("Mean:", statistics.mean(runtimes))
1236 print("Std Dev:", statistics.stdev(runtimes))

```

Figure 15: Example template provided to annotators to fill in and adapt with PR specific content.

For each of the 1041 PRs, the annotation task involved pulling the instance Docker image locally, then a deep analysis of the PR’s description, code differential (`diff`), and any associated discussion threads, code blocks, or linked issues (including browsing GitHub and the state of the code-base at that commit). The goal was to identify or reconstruct the specific code path or usecase that the PR author intended to optimize, with high preference towards existing code blocks and

PR author performance scripts. Annotators then scripted this workload into a standardized Python `def workload()` template (as shown in 15), which can then be executed against the repository. A rigorous two-step verification process was applied to every annotated workload before its inclusion in the benchmark:

1. **Peer Verification:** First, all workloads were verified by both annotators via discussion and approval. This side-by-side review process ensured the workload’s logic was sound, it accurately captured the optimization described in the PR, and it was a faithful representation of the performance test (either explicitly provided by the PR author or inferred from the code changes).
2. **Execution Verification:** Second, each workload was programmatically executed to confirm its correctness and efficacy. We ran every workload script against two distinct versions of the code within its Docker container: (1) the `base` commit (pre-optimization) and (2) the post-patch state (post-optimization). A workload was only accepted if it executed successfully on both commits *and* demonstrated a statistically significant performance speedup on the `head` commit, thereby empirically validating the PR’s performance claim.

Only workloads that passed both peer and execution verification were included in the final set of tasks for our benchmark.

C.5 EXECUTION-BASED FILTERING

We finally verify task instances by (i) executing and collecting their after-gold-edit test statuses and (ii) verifying that performance optimizations are statistically significant. Each annotated workload script contains a `workload()` function and a measurement harness to run a repeated number of iterations, generating a distribution of runtimes with both a mean and standard deviation (pre-edit as μ_{pre} and σ_{pre} and μ_{post} and σ_{post}). We filter away any instances where $\mu_{pre} - \mu_{post} \leq 2\sigma_{post}$ (i.e. runtime speedup is larger than two post-edit runtime standard deviations). In this stage we also run correctness tests ten times to filter out any possible flaky tests, as those would cause our aggregated speedup ratio to be lower than the actual value.

D TECHNIQUES FOR IMPROVING PERFORMANCE REPRODUCIBILITY

This section describes two implementation choices we use in our benchmark to reduce incidental variability in measured runtime and throughput: (i) prebuilding Docker containers so that environment resolution and installation never occur on the critical path of an evaluation run, and (ii) CPU pinning that separates container execution from Docker management daemons and assigns containers to non-overlapping groups of logical cores.

D.1 PREBUILDING INSTANCE DOCKER IMAGES

We containerize each benchmark task and *prebuild* the corresponding Docker images prior to any timed evaluation (uploading it to a public Docker image registry). Thus, evaluation runs start from a fully built image; they do not perform package installation, environment resolution, or other setup work that would otherwise consume CPU cycles and introduce run-to-run variance. This design ensures that the CPU resources measured during evaluation are dedicated to the containerized program and harness rather than to container initialization. It also makes parallel execution more stable: because images are prepared ahead of time, concurrent workers do not contend for CPU due to on-the-fly dependency installation or environment setup. We provide these scripts in our code artifact release.

D.2 PINNING CONTAINERS AND DOCKER DAEMON TO CPU CORES

We implement a CPU-affinity policy that (a) assigns containers to disjoint groups of logical cores and (b) reserves a separate set of physical CPUs for Docker’s background services. The policy proceeds as follows.

Grouping logical cores. Since each instance is evaluate on 4 vCPUs and 16GB of RAM, we first identify the logical-core (vCPU) topology and partition the available vCPUs into groups of four (4), with the constraint that *no two vCPUs in the same group share a physical core*. This grouping helps reproducibility because cache lines are generally isolated per physical core (and thus isolated between groups of 4 vCPUs), so execution within one group is more insulated from core-level contention within that group.

Isolating Docker management. We pin the Docker daemon and `containerd` to a dedicated set of *physical* CPUs (and their corresponding logical cores) that is disjoint from $\bigcup_i G_i$. As a result, container- and image-management activity (e.g., image downloading and setup) is confined to these reserved CPUs and cannot steal cycles from the cores executing benchmark containers. This separation allows us to parallelize workers across multiple groups G_i without coupling their performance to background Docker activity.

Memory limits and NUMA node binding. In addition to CPU affinity, we restrict memory on a per-container basis, allowing each container to only consume 16GB and assign each container to use the NUMA (Non-Uniform Memory Access) memory node corresponding to the physical cores of the vCPUs that the container is assigned to. This bounds each worker’s memory footprint and makes memory allocation more predictable and reduces cross-container interference due to host-level memory pressure, complementing the CPU isolation described above.

E AGENT HARNESS PROMPTS AND DETAILS

This section describes the prompt and harness specific details used to generate our evaluation results in Section 4. More details can also be found in our attached code artifact.

E.1 CODE OPTIMIZATION TASK PROMPTS

The prompt provided to OPENHANDS and SWE-AGENT is provided below, asking an LM agent to optimize a specific workload given a repository, file utilities, and bash execution abilities in a containerized environment. Note that the agent is also given the commands for (1) rebuilding/reinstalling the repository and (2) the generic prefix command for running an arbitrary unit-test file.

Code Optimization Task Prompt

```
<uploaded_files>
{{working_dir}}
</uploaded_files>
```

```
I've uploaded a python code repository in the directory
↪ workspace_dir_name. Consider the following python workload
↪ showing a specific usage and measured performance of the
↪ repository:
<performance_workload>
{{workload}}
</performance_workload>
```

```
Can you help me implement the necessary changes to the repository so
↪ that the runtime of the `workload()` function is faster? Basic
↪ guidelines:
```

1. Your task is to make changes to non-test files in the `/workspace` directory to improve the performance of the code running in `workload()`. Please do not directly change the implementation of the `workload()` function to optimize things: I want you to focus on making the workload AS IS run faster by only editing the repository containing code that the `workload()` function calls.
2. Make changes while ensuring the repository is functionally equivalent to the original: your changes should not introduce new bugs or cause already-passing tests to begin failing after your changes. However, you do not need to worry about tests that already fail without any changes made. For relevant test files you find in the repository, you can run them via the bash command `{{test_cmd}} <test_file>` to check for correctness. Note that running all the tests may take a long time, so you need to determine which tests are relevant to your changes.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

3. Make sure the ``workload()`` function improves in performance after
 ↪ you make changes to the repository. The workload can potentially
 ↪ take some time to run, so please allow it to finish and be
 ↪ generous with setting your timeout parameter: for faster
 ↪ iteration, you should adjust the workload script to use fewer
 ↪ iterations. Before you complete your task, please make sure to
 ↪ check that the `**original performance workload**` and ``workload()``
 ↪ function runs successfully and the performance is improved.

4. You may need to reinstall/rebuild the repo for your changes to
 ↪ take effect before testing if you made non-Python changes.
 ↪ Reinstalling may take a long time to run, so please be patient
 ↪ with running it and allow it to complete if possible. You can
 ↪ reinstall the repository by running the bash command
 ↪ ``{{rebuild_cmd}}`` in the workspace directory.

5. All the dependencies required to run the ``workload()`` function are
 ↪ already installed in the environment. You should not install or
 ↪ upgrade any dependencies.

Follow these steps to improve performance:

1. As a first step, explore the repository structure.

2. Create a Python script to reproduce the performance workload,
 ↪ execute it with `python <workload_file>`, and examine the printed
 ↪ output metrics.

3. Edit the source code of the repository to improve performance.
 ↪ Please do not change the contents of the ``workload()`` function
 ↪ itself, but focus on optimizing the code in the repository that
 ↪ the original ``workload()`` function uses.

4. If non-Python changes were made, rebuild the repo to make sure the
 ↪ changes take effect.

5. Rerun your script to confirm that performance has improved.

6. If necessary, identify any relevant test files in the repository
 ↪ related to your changes and verify that test statuses did not
 ↪ change after your modifications.

7. After each attempted change, please reflect on the changes
 ↪ attempted and the performance impact observed. If the performance
 ↪ did not improve, consider alternative approaches or
 ↪ optimizations.

8. Once you are satisfied, please use the `finish` command to complete
 ↪ your task.

Please remember that you should not change the implementation of the
 ↪ ``workload()`` function. The performance improvement should solely
 ↪ come from editing the source files in the code repository.

E.2 DETAILS ON LANGUAGE MODEL SAMPLING PARAMETERS

We elaborate on the evaluation settings discussed in 3. For all models, we perform the recommended greedy sampling within the OpenHands and SWE-agent harnesses. For GPT-5 MINI, we sample at a temperature of $t = 1$ (as mandated by the API as of August 2025) and at $t = 0$ for all other models. In the SWE-agent setting, we enforce a token spending limit of \$1, meaning that, in addition to the 100-turn action limit and time-limits specified, evaluation runs per-instance are stopped when (API/token-spending) cost is exceeded, and their patches as of that last action are immediately submitted (as is common practice with cost-limited SWE-agent runs on SWE-bench).

Table 6: SWE-EFFICIENCY results across several frontier models (higher is better; human-expert speedup ratio (SR) is 1.0 \times). SR is *pass@1*: each system submits a single patch per instance to be evaluated. SR is calculated by computing the speedup from the LM-generated edit, normalized by the speedup from the gold (human-written) patch, and aggregated across all tasks via harmonic mean.

System	Speedup Ratio
Expert	1.0 \times
GPT-5 (OPENHANDS)	0.150 \times
CLAUDE 4.1 OPUS (OPENHANDS)	0.098 \times
QWEN3 CODER PLUS (OPENHANDS)	0.064 \times
CLAUDE 3.7 SONNET (OPENHANDS)	0.047 \times
CLAUDE 4.5 SONNET (OPENHANDS)	0.041 \times
GLM 4.6 (OPENHANDS)	0.026 \times
GPT-5 MINI (OPENHANDS)	0.019 \times
KIMI K2-0905 (OPENHANDS)	0.008 \times
GEMINI 2.5 FLASH (OPENHANDS)	0.008 \times
DEEPSEEK V3.1 (OPENHANDS)	0.007 \times
GEMINI 2.5 PRO (OPENHANDS)	0.007 \times
CLAUDE 3.7 SONNET (SWE-AGENT)	0.041 \times
GPT 5 MINI (SWE-AGENT)	0.026 \times
GEMINI 2.5 FLASH (SWE-AGENT)	0.006 \times

We believe this lower-resource, cost-constrained setting is important from an *efficiency* standpoint of eventually yielding systems that can solve optimization tasks at reasonable dollar costs. In the OpenHands setting, our results only have the action count and time limits enforced. We also provide links to our forks of those agent harnesses for evaluation, which will also be merged upstream with corresponding harness libraries for community reproducibility.

F ADDITIONAL DETAILS ON MAIN EVALUATION RESULTS

F.1 BENCHMARK PERFORMANCE FULL RESULTS

We provide full versions of evaluation results below for both OPENHANDS and SWE-AGENT below in Tables 6 and 7. Recall that for both harnesses, we set a 3 hour wall-clock time limit and a 100 turn interaction limit: in the SWE-AGENT case, we also limit LMs to a maximum cost of \$1.

Table 7: Distribution of patch outcomes by system. “Passes correctness tests” denotes functional correctness only (and not necessarily performance-optimal).

System	Fails Tests (\downarrow)	Passes Correctness Tests		
		Slower than Pre-edit (\downarrow)	Faster than Pre-edit (\uparrow)	Faster than Expert (\uparrow)
GPT-5	18.3%	4.4%	31.5%	45.8%
CLAUDE 4.1 OPUS	15.2%	4.0%	42.8%	38%
QWEN3 CODER PLUS (OPENHANDS)	22.5%	11.0%	42.0%	24.5%
CLAUDE 3.7 SONNET (OPENHANDS)	34.7%	12.7%	43.8%	32.9%
CLAUDE 4.5 SONNET (OPENHANDS)	18.7%	4.7%	32.1%	20.7%
GLM 4.6	32.7%	12.8%	38.1%	16.3%
GPT 5 MINI (OPENHANDS)	45.2%	14.7%	26.5%	13.9%
KIMI K2-0905 (OPENHANDS)	38.2%	6.4%	36.6%	19.3%
GEMINI 2.5 FLASH (OPENHANDS)	39.2%	14.5%	34.1%	12.7%
DEEKSEEK V3.1 (OPENHANDS)	19.5%	17.9%	44.2%	18.5%
GEMINI 2.5 PRO (OPENHANDS)	41.8%	17.7%	32.9%	7.6%
CLAUDE 3.7 SONNET (SWE-AGENT)	39.6%	8.6%	27.5%	24.7%
GPT 5 MINI (SWE-AGENT)	25.5%	11.5%	35.7%	27.7%
GEMINI 2.5 FLASH (SWE-AGENT)	44.2%	12.0%	35.7%	8.4%

F.2 HOW DOES LM PERFORMANCE SCALE WITH DATASET DIFFICULTY?

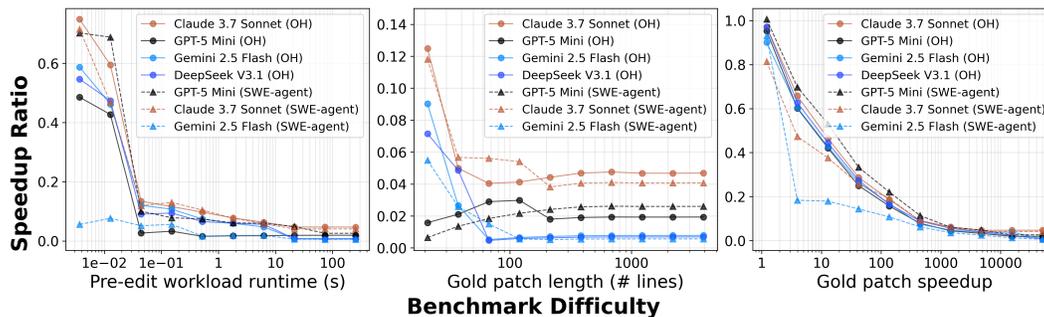


Figure 16: LMs achieve easier wins on lower difficulty problems, but struggle as higher difficulty tasks are included across multiple “definitions” of difficulty. For each difficulty measure and each measure upper bound τ , we restrict to instances with difficulty measure $\leq \tau$ and report the resulting aggregate speedup ratio, generating our curves shown.

In addition to the bucketed trends shown in Figure 4, Figure 16 shows curves demonstrating how models perform on our dataset as we increasingly include tasks across the same three dimensions of (i) pre-edit workload duration, (ii) number of lines modified in the gold (expert) patch, and (iii) speedup factor achieved by the gold patch. We notice that as increased difficulty tasks are included, speedup ratio across the dataset decreases, further indicating that LMs are achieving easier wins on lower difficulty problems but struggling at higher difficulties.

F.3 EXAMINING MORE EXPENSIVE REASONING MODELS: COMPARING GEMINI 2.5 PRO VS. FLASH

We initially (as of September 19th 2025) could not run full benchmark results on full reasoning models like GPT-5, OPUS 4.1, and GEMINI 2.5 PRO *due to budget and runtime limitations*: runs would cost a significant amount and also take much longer per inference call (even with parallel requests) to reasonably complete in time.

Instead, we share results from a selected subset of 100 SWE-EFFICIENCY problems, designated as SWE-EFFICIENCY LITE. For this subset, we sample to be representative with respect to pre-edit workload runtime, gold patch speedup, and number of lines in gold patch from the distributions in Figure 3. Specifically, we construct a small, distribution-matched “lite” split by log-spacing each difficulty metric into bins and assigning instances to bins. We allocate a per-metric quota for a target size $N = 100$ in proportion to each bin’s population and sample without replacement from those bins (using a fixed random seed). We take the union across metrics to cover diverse regions of each marginal distribution and top up any remaining slots by sampling from the unsampled pool with weights inversely proportional to the density of each instance’s 3-way bin signature, which promotes rare metric combinations and preserves joint structure. If the union overshoots $N = 100$, we trim instances uniformly at random until we reach the desired amount.

In Table 8, we see that GEMINI 2.5 PRO performs similarly to its medium-compute counterpart GEMINI 2.5 FLASH on SWE-EFFICIENCY LITE, while being more than $5\times$ as expensive dollar-wise and incurring an extra 2.74 total hours of inference latency. This suggests that more-expensive state-of-the-art reasoning models also still heavily struggle on SWE-EFFICIENCY and larger, agentic advances are needed to make models that reason more in-depth about repo-level performance and that can iterate in harnesses quickly.

G PROFILING-BASED ATTRIBUTION AND COVERAGE METRICS

In this section, we elaborate on how we computed the profiling and function localization results shared in Section 4 and Figure 6, which show that LMs often miss out on expert-level speedup due to function level mislocalization.

What is a function-level profiler? A function-level profiler instruments program execution to record, for every function invocation, (i) the *exclusive* or *self* time spent in the function body (excluding callees), often called *tottime*, (ii) the *inclusive* or *cumulative* time spent in the function and

Table 8: SWE-EFFICIENCY LITE results between GEMINI 2.5 PRO and GEMINI 2.5 FLASH (higher is better; human-expert parity is $1.0\times$). “Passes tests” indicates passing functional correctness tests only. LM cost is total token spend (including prompt-caching). Inference latency is sum of total request latency over all requests.

System	Speedup Ratio	Passes Tests	LM Cost	Inference Latency (hrs)
GEMINI 2.5 PRO (OPENHANDS)	$0.008\times$	60%	\$509.52	9.03
GEMINI 2.5 FLASH (OPENHANDS)	$0.007\times$	65%	\$98.83	6.29

its transitive callees, often called *cumtime*, (iii) call counts, and (iv) the caller–callee relationships that induce a directed call graph. In our setting we profile a benchmark `workload` entry point before and after a code edit, obtaining two traces per actor (expert vs. LLM). We compute all metrics *only* on instances that pass correctness checks (e.g., unit or regression tests), so any measured speedup does not come at the expense of functional correctness.

G.1 DATA AND NOTATION

Let \mathcal{G} denote the set of functions observed by the profiler. We uniquely identify a function $g \in \mathcal{G}$ by its source file $f(g)$, line number $\ell(g)$, and name $n(g)$. For an actor $A \in \{\text{Expert, LLM}\}$ and a profiling phase $p \in \{\text{pre, post}\}$, let

$$\tau_{A,p}(g) \in \mathbb{R}_{\geq 0} \quad \text{and} \quad T_{A,p}(g) \in \mathbb{R}_{\geq 0}$$

denote the exclusive (*tottime*) and inclusive (*cumtime*) runtime attributed to g , respectively. We define per-function improvements (positive means faster) as

$$\Delta_A^{\text{tot}}(g) := \tau_{A,\text{pre}}(g) - \tau_{A,\text{post}}(g), \quad \Delta_A^{\text{cum}}(g) := T_{A,\text{pre}}(g) - T_{A,\text{post}}(g).$$

Let W denote the top-level entry point (`workload`); its end-to-end improvement for actor A is

$$\delta_A^W := T_{A,\text{pre}}(W) - T_{A,\text{post}}(W).$$

We additionally report whole-trace speedups normalized by pre-edit workload time,

$$\text{Speedup}_A^W := \frac{\delta_A^W}{T_{A,\text{pre}}(W)}, \quad \text{Speedup}_A^{\text{tot}} := \frac{\sum_g \tau_{A,\text{pre}}(g) - \sum_g \tau_{A,\text{post}}(g)}{T_{A,\text{pre}}(W)}.$$

Call graph and depths. We first need to isolate the call graph (and function runtimes) that are strictly attributed to the `workload` function in each performance workload (and disregard function runtimes from any other source, like the `setup` function). From the *pre*-edit profile we build a directed call graph $\mathcal{C} = (\mathcal{G}, E)$ whose edges point from caller to callee. We define the *workload depth* $d(g)$ as the minimum caller distance from any node named `workload` to g in \mathcal{C} ; nodes not reachable from `workload` have undefined depth. Depths are used only for selection and the diagnostic depth metric in Appendix G.5.

Patch-based file filters. We parse unified diffs to extract modified files for each actor and restrict candidate functions to those files. This ensures we attribute improvements to edited regions and reduces noise from unrelated code.

G.2 SELECTING CORE FUNCTION-LEVEL IMPROVEMENTS WITHOUT DOUBLE COUNTING

Naively summing Δ_A^{cum} over functions double-counts speedups because a caller’s inclusive time subsumes callee improvements. We therefore select a *deepest non-overlapping* set of improved functions for each actor via a depth-aware greedy procedure.

Thresholding and candidates. We set a per-instance absolute threshold

$$\theta := \max(\theta_{\text{sec}}, \theta_{\text{frac}} \cdot \delta_{\text{Expert}}^W),$$

where θ_{sec} is an absolute time floor (seconds) and θ_{frac} is a fraction of the expert’s end-to-end improvement (default 0.02). Candidates functions are defined below where W is the node corresponding to the `workload` function entry point (see B.2 for a workload example):

$$\mathcal{C}_A := \left\{ \begin{array}{l} \Delta_A^{\text{cum}}(g) > 0, \Delta_A^{\text{cum}}(g) \geq \theta, \\ g \in \mathcal{G} : g \text{ is reachable from } W, \\ \text{and } f(g) \text{ was edited by } A \end{array} \right\}.$$

1566 Let $S_A^+ := \sum_{g \in \mathcal{C}_A} \Delta_A^{\text{cum}}(g)$ denote the total positive mass among candidates.
 1567

1568 **Greedy deepest-first selection.** We sort candidates by (i) larger depth $d(g)$ first, (ii) larger share
 1569 $\Delta_A^{\text{cum}}(g) / \max(T_{A,\text{pre}}(W), \epsilon)$, then (iii) larger $\Delta_A^{\text{cum}}(g)$ (ties broken arbitrarily), and greedily build
 1570 a set $E_A \subseteq \mathcal{C}_A$ such that no selected function is an *ancestor* (caller, transitively) of another selected
 1571 function in the pre-edit call graph. Intuitively, we select a set of functions closest to the scope of the
 1572 speedup (i.e. the function scope where the speedup has the most significant percentage improvement
 1573 over that time scope). For example, a speedup could occur in function C but is called by B , which is
 1574 then called by A : C would show the largest percentage speedup relative to the total amount of time
 1575 spent in that function, since B and A have other runtime overhead that was not optimized. We stop
 1576 when the accumulated mass reaches a configurable cap $\rho \in (0, 1]$:

$$1577 \sum_{g \in E_A} \Delta_A^{\text{cum}}(g) \geq \rho \cdot S_A^+ \quad (\text{default } \rho = 1).$$

1580 If the procedure would select nothing, we include the single best candidate. This yields our set of
 1581 functions E_{Expert} and E_{LLM} .

1582 G.3 EXPERT-RELATIVE COVERAGE (ERC) AND LOSS DECOMPOSITION

1583 We measure how well the LLM’s edits localize to the same *places of improvement* as the expert. Let
 1584 $\Phi(S) := \{f(g) : g \in S\}$ map a set of functions to its set of files. To define the expert’s *attribution*
 1585 *mass* we (i) keep only functions above threshold and (ii) restrict to files the expert actually optimized
 1586 (to guard against spurious activity in unrelated files):
 1587

$$1588 \mathcal{M}_{\text{Expert}} := \{g \in \mathcal{G} : \Delta_{\text{Expert}}^{\text{cum}}(g) \geq \theta, f(g) \in \Phi(E_{\text{Expert}})\}.$$

1589 Define per-function expert mass $s_{\text{exp}}(g) := \Delta_{\text{Expert}}^{\text{cum}}(g)$ for $g \in \mathcal{M}_{\text{Expert}}$ and 0 otherwise, and let
 1590 $S_{\text{exp}} := \sum_g s_{\text{exp}}(g)$.
 1591

1592 **Coverage at file and function granularity.** We compute ERC at two granularities:
 1593

$$1594 \text{ERC}_{\text{file}} := \frac{\sum_{f \in \Phi(E_{\text{LLM}})} \sum_{g \in \mathcal{M}_{\text{Expert}}: f(g)=f} s_{\text{exp}}(g)}{S_{\text{exp}}},$$

$$1596 \text{ERC}_{\text{func}} := \frac{\sum_{g \in E_{\text{LLM}}} s_{\text{exp}}(g)}{S_{\text{exp}}}.$$

1599 Intuitively, ERC_{file} asks “did the LLM edit the *right files*?” while ERC_{func} asks “did it optimize the
 1600 *right functions* within those files?”
 1601

1602 **Loss decomposition.** We decompose the portion of expert mass *not* captured at the function level
 1603 into two orthogonal failure modes:

$$1604 \text{WrongFileLoss} := 1 - \text{ERC}_{\text{file}},$$

$$1605 \text{InFileLoss} := \max\{0, \text{ERC}_{\text{file}} - \text{ERC}_{\text{func}}\}.$$

1607 This yields a tight partition of expert mass:
 1608

$$1609 \text{WrongFileLoss} + \text{InFileLoss} + \text{ERC}_{\text{func}} = 1,$$

1610 where *WrongFileLoss* captures file-selection mistakes and *InFileLoss* captures localization mistakes
 1611 *within* the right files (e.g., editing non-bottleneck functions).
 1612

1613 G.4 EDITED-FILE OVERLAP

1614 We also report the Jaccard similarity of edited files between actors:
 1615

$$1616 \text{Jaccard} := \frac{|\Phi(E_{\text{Expert}}) \cap \Phi(E_{\text{LLM}})|}{|\Phi(E_{\text{Expert}}) \cup \Phi(E_{\text{LLM}})|}.$$

1618 When the union is empty (neither actor meets the selection threshold), the quantity is undefined and
 1619 we omit it.

G.5 DEPTH OF OPTIMIZATION FROM THE WORKLOAD

As a diagnostic we compute a depth-of-optimization statistic with respect to the pre-edit call graph and the workload root. Let $[x]_+ := \max\{x, 0\}$ and define weights $w_A(g) := [\Delta_A^{\text{tot}}(g)]_+$ for $g \in E_A$. The *weighted average workload depth* and its coverage are

$$\bar{d}_A := \frac{\sum_{g \in E_A \cap \text{reach}(W)} w_A(g) d(g)}{\sum_{g \in E_A \cap \text{reach}(W)} w_A(g)}, \quad \text{ReachShare}_A := \frac{\sum_{g \in E_A \cap \text{reach}(W)} w_A(g)}{\sum_{g \in E_A} w_A(g)}.$$

\bar{d}_A reflects whether improvements concentrate near the entry point or deep in the call tree; ReachShare_A indicates how much of the selected mass is reachable from the workload (should be close to 1 in well-instrumented runs). We use $\theta_{\text{sec}} = 0$, $\theta_{\text{frac}} = 0.02$ (i.e., a per-function floor at 2% of the expert’s end-to-end gain δ_{Expert}^W to disregard speedups that are due to measurement noise), and cap $\rho = 1.0$, and we restrict candidate functions to edited files for each actor. The call graph used for depths and ancestry tests is always taken from the pre-edit trace to avoid post-edit structural confounds. Note that we only compute these statistics over instances that have passed functional correctness tests.

Why Δ^{cum} for selection and Δ^{tot} for depth weights? We select by Δ^{cum} to capture inclusive speedups (including callee effects) while the depth statistic weights by Δ^{tot} to avoid double-counting along a chain. The deepest-first greedy constraint further prevents attributing the same improvement to both a caller and its callee.

Interpretation. High ERC_{file} with low ERC_{func} indicates that the model navigated to the right files but failed to touch the expert-optimized functions (*within-file localization gap*). Low ERC_{file} indicates a file-selection gap. Because S_{exp} is defined over expert-selected files above threshold, the metrics focus on *where* expert improvements actually occurred, rather than on unrelated noisy regions.

G.6 FULL SPEEDUP ATTRIBUTION RESULTS

System	ERC_{file}	ERC_{func}	WrongFileLoss	InFileLoss	Jaccard (files)
CLAUDE 3.7 SONNET (SWE-AGENT)	0.630	0.298	0.370	0.332	0.636
CLAUDE 3.7 SONNET (OPENHANDS)	0.611	0.314	0.389	0.297	0.604
GPT-5 MINI (OPENHANDS)	0.551	0.278	0.449	0.274	0.559
GEMINI 2.5 FLASH (OPENHANDS)	0.549	0.265	0.451	0.283	0.556
DEEPSEEK V3.1 (OPENHANDS)	0.519	0.246	0.481	0.273	0.531

Table 9: Expert-Relative Coverage (ERC) and related losses. Means computed over instances that pass correctness and have speedup ≥ 1 . File-overlap Jaccard is averaged over instances where it is defined.

System	Number Correct Instances	Depth _{exp}	Depth _{llm}
CLAUDE 3.7 SONNET (SWE-AGENT)	196	4.85	4.20
CLAUDE 3.7 SONNET (OPENHANDS)	252	4.61	4.18
GPT-5 MINI (OPENHANDS)	193	4.59	4.13
GEMINI 2.5 FLASH (OPENHANDS)	211	4.51	3.83
DEEPSEEK V3.1 (OPENHANDS)	246	4.89	4.28

Table 10: Dataset size (n) and weighted average optimization depth from the workload entry point in the pre-edit call graph (expert vs. LLM).

We compute results only over LM patches that passed correctness and achieve a speedup from the pre-edit runtime (but not necessarily faster than the expert). Across systems, the ERC and loss metrics in Table 9 show a consistent pattern: *WrongFileLoss* is roughly 37%–45% across models (mean $\approx 41.5\%$), while *InFileLoss* is roughly 27%–33% (mean $\approx 29.7\%$). Taken together, this implies that models choose the wrong function (either by editing the wrong file or the wrong function within the right file) about $X+Y \approx 69\%$ – 73% of the time (mean $\approx 71.2\%$), consistent with $\text{ERC}_{\text{func}} \approx 0.26$ – 0.31 . For context on number of instances analyzed per system and how deep

optimizations occur in the call tree, Table 10 reports n per system and the weighted average depth from the workload root (experts typically operate at call-stack depth 4.5–4.9, LLMs at 3.8–4.2).

H COMPARISON OF LM GENERATED EDITS VERSUS EXPERTS

We provide the raw diffs with in-line comments from Figures 8 and 9 below in Figures 17 and 18. We also include an additional diff for 19 comparison. In the main text, we removed comments and surrounding lines to focus only on the lines changed between expert and LM generated diff.

I SYNTHETICALLY GENERATING PERFORMANCE WORKLOADS

We provide more details on our investigation on whether LMs can capably generate performance workloads, as discussed in the end of Section 4.2.

Rationale. We generate LLM-based workloads to mirror our human curation process and test a key hypothesis: when the *gold patch is held fixed*, workloads curated by expert annotators (our pipeline) expose larger, statistically reliable performance deltas than workloads produced by an LLM from the same evidence.

Inputs per instance. For each SWE-FICIENCY benchmark instance we use:

- The unified diff of the expert (gold) patch,
- The pre-edit source files corresponding to paths touched in the diff, repository and commit identifiers.

Prompt construction. We parse the diff headers to identify touched files and fetch their pre-edit contents at the base commit. The model is given:

1. The full patch (pre/post diff), and
2. The concatenated pre-edit files for those paths

Instruction parity with human annotators. We prompt the LLM (GEMINI 2.5 FLASH) with similar instructions as we used ourselves for workload annotation: produce a *self-contained Python workload script* with `setup()` (realistic inputs; seeded randomness) and `workload()` (representative, non-trivial call path), executed via `timeit.repeat(...)`, and printing exactly two lines—mean and standard deviation. This parity isolates the effect of workload design quality, not interface differences. We provide the instruction prompt below.

Synthetic Workload Generation Prompt

```
You are a performance testing expert. You will be provided a code
↪ edit as a git diff and the pre-edit source files. You need to
↪ generate a **self-contained Python performance workload script**
↪ that measures performance of code paths or APIs changed in the
↪ diff.
```

Guidelines for the workload script contents.

- Use a `setup()` function to prepare any realistic, non-trivial data
 - ↪ or environment needed for the test.
 - Data must be representative of real-world usage (avoid trivial
 - ↪ arrays or easily optimizable patterns).
 - Prefer real datasets or realistic synthetic data with
 - ↪ reproducibility (set a random seed).
 - All expensive or one-time setup (e.g., file download,
 - ↪ preprocessing) must be in `setup()`, not in the workload.
- Use a `workload()` function to run the actual operation(s) being
 - ↪ timed.
 - The workload should reflect a **representative and challenging
 - ↪ real-world use case** of the API or library under test.
 - Avoid corner cases that could be trivially optimized.

```

1728
1729 - Inputs should be varied enough to prevent caching or
1730 ↪ constant-folding from affecting results.
1731
1732 - Run the benchmark using `timeit.repeat(workload, number=...,
1733 ↪ repeat=..., setup=setup)`.
1734 - `number` should match a realistic single-run execution count (do
1735 ↪ not batch multiple runs for cumulative timing).
1736 - `repeat` should be high enough to gather stable statistics.
1737
1738 - Print the mean and standard deviation of the last set of runtimes
1739 ↪ using `statistics.mean()` and `statistics.stdev()`.
1740 - Output should be clear and ready for performance comparison.
1741
1742 - The output must be a complete Python script containing only:
1743 1. import statements
1744 2. `setup()` function
1745 3. `workload()` function
1746 4. the `timeit.repeat()` call
1747 5. mean/stddev printing
1748
1749 The script should only print two lines at the end: the mean of
1750 ↪ measured runtimes and the standard deviation of runtimes.
1751
1752 Example workload to follow (please strictly follow this format of
1753 ↪ imports, setup function, workload function, timeit call, and
1754 ↪ print statements). In particular, make sure the mean and standard
1755 ↪ deviation print statements are exactly as shown below.
1756
1757 ```python
1758 import timeit
1759 import statistics
1760 import numpy as np
1761
1762 def setup():
1763     global arr
1764     np.random.seed(42)
1765     arr = np.random.rand(5000, 5000)
1766
1767 def workload():
1768     global arr
1769     _ = arr @ arr.T
1770
1771 runtimes = timeit.repeat(workload, number=1, repeat=10, setup=setup)
1772
1773 print("Mean:", statistics.mean(runtimes))
1774 print("Std Dev:", statistics.stdev(runtimes))
1775 ```
1776
1777 Here's a commit and it's information that does some optimization in
1778 ↪ the {repo_name} repository that might be relevant to writing the
1779 ↪ test:
1780 ## Commit Diff:
1781 ...
1782 {commit_diff}
1783 ...
1784
1785 ## Pre-edit source files:
1786 {pre_edit_code}
1787

```

Evaluation (holding the patch fixed). For each instance i , we evaluate both workloads—LLM-generated and manually annotated—against the same code states:

1. **Base:** Repository at the pre-gold patch commit.

1782 2. **Patched:** Repository with the gold patch applied.

1783
1784 We then compare the magnitude of improvements (speedups) between the LM-generated (GEMINI
1785 2.5 FLASH) workload and the manual workload for the same instance and compute if values are
1786 statistically significant and if the LM generated workload outperforms the annotated one. This
1787 ablation directly measures whether an LLM, given the same diff and file context and the same
1788 instructions as experts, can generate workloads that surface the patch’s performance gains as reliably
1789 and strongly as manual curation. We see that, in 76% of cases, our manual annotations show a larger
1790 performance delta than the LM generated workloads on the expert patch, with 47% of workloads
1791 showing a non-significant performance delta at all.

1792 **J PREVENTING MODEL REWARD HACKING**

1793
1794 In this section, we outline two techniques implemented in our evaluation harness to prevent LM
1795 reward hacking behavior. The first, *stack-frame-based* reward hacking covers when models try to
1796 determine their caller (i.e. whether they are being called in a performance runtime environment),
1797 while the second covers *run-to-run caching*, which is when models try to cache computations across
1798 evaluation runs to improve performance (but technically cheating, as we’d like to evaluate specific
1799 workloads under non-cached conditions).

1800 **J.1 PREVENTING STACK-FRAME-BASED REWARD HACKING**

1801 **Overview.** We observed that some LLM-generated patches only speed up workloads when
1802 they detect they are being timed (i.e. under `timeit` or from being called from a function
1803 with name `workload`, as is done in our benchmark). The common mechanism is Python
1804 stack introspection (e.g., `inspect.currentframe()`, `traceback.extract_stack()`,
1805 `sys._getframe()`, or reaching frame objects via `f_back/tb_frame`). These edits can short-
1806 circuit code paths or memoize based on caller identity, inflating measured speedups without actually
1807 improving the underlying algorithm.

1808 In general, code changes (and in-particular performance improving edits) should *never* require
1809 caller identity information to improve performance. Other types of changes, such as tuning specifi-
1810 cally based on input attributes (like size or data-type), are actually quite common in perfor-
1811 mance optimization PRs: we intentionally permit these changes in SWE-FFICIENCY. We verify
1812 that none of the expert PRs and gold patches use stackframe information: the one exception is
1813 `pandas-dev__pandas-45247`, which optimizes `find_stack_level` in `pandas`, a utility
1814 to more readably show exception stackframes (and where the expert patch uses these utilities).

1815 **Goal.** We want to flag *newly introduced* stack-introspection logic in a submitted patch while
1816 tolerating any pre-existing usage in the codebase (many mature projects legitimately use
1817 `inspect/traceback` during import/configuration) as well as making sure that expert/gold
1818 patches are not falsely flagged.

1819 **Implementation.** Our checker takes a unified diff as input and analyzes only the *post-image* of
1820 files touched by that diff. It reports an error *only* if the diff *adds* lines that contain stack-introspection
1821 primitives. Existing occurrences are ignored by design.

- 1822
- 1823 1. **Patch-scope extraction.** We parse the unified diff to recover, for each touched file, the set
1824 of line numbers that are newly added on the “+” side. This yields a map `added_lines :`
1825 `path ↦ {new line numbers}`. We also track which files are brand-new in the patch and the
1826 set of all post-image paths seen in `+++` headers.
 - 1827 2. **Standalone-new filtering.** Brand-new files (introduced by the patch) that are *not imported*
1828 by any other file touched in the patch are treated as “standalone” and excluded from the
1829 check. This avoids flagging ad hoc scripts (e.g., local reproducer/benchmark drivers) that
1830 do not affect the library under test, as often LM systems might produce ad hoc scripts
1831 like these for debugging and introspection in a scratchpad form (which may use some
1832 stackframe inspection). This lets us ignore scratch-pad like files introduced by patches,
1833 while still checking newly-created files that are imported and used in the repository. We
1834 determine “referencedness” by building a lightweight import graph among touched files:

- 1835 (a) For each patch edited file, collect its imported module names via AST (both
 `import m` and `from m import x`).

- 1836 (b) For each brand-new file, derive candidate module names from its path (e.g.,
 1837 `foo/bar/baz.py` \rightarrow `{baz, bar.baz, foo.bar.baz}`).
- 1838 (c) Mark the new file as referenced if any other touched file imports one of its candidates
 1839 (exact or dotted-suffix match).
- 1840
- 1841 3. **Post-edit AST scan with alias resolution.** For each edited file in `added_lines`, we read
 1842 the *post-edit* source file. We parse each post-edit file source with `ast` and walk the tree
 1843 once, collecting “findings” whenever the code contains introspection-like constructs. This
 1844 scan is robust to renaming via an an AST-based import resolve, described below
- 1845 • **Imports:** record module aliases (e.g., `import inspect as ins`) and function
 1846 aliases (e.g., `from inspect import currentframe as cf`).
 - 1847 • **Direct calls:** resolve callee to (module, attribute) pair and match against a denylist,
 1848 which includes the following list:
 - 1849 - `inspect`.`{currentframe, stack, getouterframes,`
 1850 `getinnerframes, trace, getframeinfo, getsource,`
 1851 `getsourcefile}`
 - 1852 - `traceback`.`{extract_stack, format_stack, print_stack,`
 1853 `walk_stack};`
 - 1854 - `sys`.`{_getframe, settrace, setprofile}`
 - 1855 - `gc`.`{get_referrers, get_objects}`
 - 1856 • **Dynamic imports:** Detect dynamic imports such as `__import__('inspect')`
 1857 and `importlib.import_module('inspect')`.
 - 1858 • **Frame-object attributes:** flag attribute reads commonly used to climb or expose
 1859 frames (`.f_back`, `.tb_frame`, `.gi_frame`, `.cr_frame`, `.ag_frame`), re-
 1860 gardless of receiver type (conservative heuristic).
- 1861 4. **Added-line projection.** After AST scanning, we *project* identified usages (and line number
 1862 occurrences) onto the added-line set for that file and retain only those whose source line
 1863 number is in `added_lines[path]`. This makes the check purely *diff-relative*: modifications
 1864 that reuse pre-existing introspection do not fail the patch (as those usages are valid).
- 1865 5. **Reporting.** If any filtered findings remain (after standalone-new filtering and pragma sup-
 1866 pression), a patch fails the check and the LM generated patch is considered incorrect and
 1867 fails correctness in our evaluation.

1868 We use AST parsing instead of regex/simple grep since simple regexes miss aliased imports and
 1869 produce many false positives/negatives on strings or comments. Our AST pass is cheap, robust,
 1870 and semantically aware: it resolves `inspect` aliases, identifies function calls regardless of whites-
 1871 pace/nesting, and recognizes dynamic imports. We verify that this mechanism (i) fails our previously
 1872 identified LM generated patches that exploit stackframe info and (ii) passes our gold, expert edits
 1873 (i.e. does not flag any false positives).

1874 J.2 PREVENTING RUN-TO-RUN CACHE REWARD HACKING

1875 **Overview.** Repeated measurements within a single Python process allow module- and process-
 1876 local state to leak across runs (e.g. module-level dictionaries, `@lru_cache`, global arrays,
 1877 ad hoc memo tables). Such state can make later iterations appear faster without changing the true
 1878 cost of the underlying algorithm. To make per-iteration results robust, we isolate runs so that no
 1879 Python-level caches (or mutated globals) can persist between repetitions.

1880 **Implementation.** We transform our annotated `workload` scripts into an equivalent program
 1881 that executes each timing repetition in a *fresh child process*. Concretely, we (i) parse the origi-
 1882 nal script with `ast`, (ii) preserve its logic and output formatting, and (iii) replace the in-process
 1883 `timeit.repeat` loop with a small harness built on `multiprocessing` using the `spawn` start
 1884 method. The `fork` policy starts a brand-new copied process for every repetition, guaranteeing a
 1885 module namespace and empty caches equivalent to the original parent run. This implementation
 1886 allows us to provide simple scripts at inference time in problem statements to LM agents, while also
 1887 being able to use those scripts as inputs to yield memory-isolated runtime scripts.

- 1888 1. **Locate the timing site.** We walk the AST to find the assignment to
 1889 `runtimes = timeit.repeat(...)` (or an equivalent import form), then extract the
workload callable, optional *setup* callable, and the numeric `number/repeat` parameters.

- 1890 We also record if the script later slices the results (e.g., `runtimes[-10000:]`) so that
 1891 summary statistics are computed over the same view.
- 1892 2. **Preserve surrounding code.** The transformer keeps all top-level declarations and state-
 1893 ments *except* the original `timeit.repeat` assignment and the immediately following
 1894 summary prints. Any statements that originally lived between those two points are pre-
 1895 served and either (i) executed after the isolated timing (if they are harmless post-processing)
 1896 or (ii) moved into a guarded `finally` block if they look like teardown of temporary
 1897 files/directories (simple heuristic over `os/shutil` calls).
- 1898 3. **Fork-per-run harness.** We synthesize a minimal harness:
- 1899 • a child-side function that constructs a `timeit.Timer(workload,`
 1900 `setup=setup)` and calls `Timer.timeit(number)`,
 - 1901 • a top-level *picklable* target that runs the child once and returns the duration through a
 1902 `multiprocessing.Queue`.
 - 1903 • a driver `_run_isolated(number, repeat, start_method="fork")`
 1904 that loops `repeat` times: for each iteration it creates a new process/context, exe-
 1905 cutes the child, checks the exit code, and appends the reported duration.
- 1906 We deliberately use **fork** so the child interpreter starts from the same memory state as the
 1907 parent (with copy-on-write) such that any edits to parent memory objects, such as module
 1908 level Python caches (including `@lru_cache` and ad hoc dictionaries) cannot carry over
 1909 between repetitions.
- 1910 4. **Result and summary fidelity.** After the harness returns the list of durations, we recon-
 1911 struct the original slicing intent (if any) into a `runtimes_view` and compute `Mean` and
 1912 `Std Dev` exactly as in the input script. Any non-teardown statements that originally ran
 1913 between the timing and the summary are executed afterward to preserve observable side
 1914 effects.

1915 Writing the transformation allows us to guarantee each repetition runs in a brand-new interpreter
 1916 process; thus module-level state, Python memo tables, and global variables cannot influence sub-
 1917 sequent repetitions. Import-time effects reoccur per iteration, making “first-run vs. warmed-run”
 1918 behavior explicit in the measurement. Random number generators also begin from the child’s fresh
 1919 state unless the user seeds them in `setup`, in which case seeding is applied identically per run.
 1920 By enforcing *fork-per-run*, the rewritten benchmarks are robust to module-level caching and other
 1921 intra-process artifacts. The transformation preserves user-visible behavior (including result slicing
 1922 and post-processing) while ensuring that any speedups reflect genuine algorithmic improvements
 1923 rather than residual state from previous iterations.

Original Raw Workload for `pandas-dev__pandas-56508`

```

1926 import timeit
1927 import statistics
1928
1929 import pandas as pd
1930 import numpy as np
1931
1932 np.random.seed(0)
1933
1934 N = 100_000
1935 data = np.arange(N)
1936 arr = pd.array(np.where(np.random.rand(N) > 0.1, data, np.nan),
1937               ↪ dtype="Int32")
1938
1939 def workload():
1940     arr._hash_pandas_object(encoding='utf-8',
1941                             ↪ hash_key="1000000000000000", categorize=False)
1942
1943 runtimes = timeit.repeat(workload, number=5, repeat=1000)
1944
1945 # Print runtime mean and std deviation.
1946 print("Mean:", statistics.mean(runtimes))
  
```

```

1944
1945 print("Std Dev:", statistics.stdev(runtimes))
1946
1947
1948

```

Transformed and Memory Isolated Workload for pandas-dev__pandas-56508

```

1950
1951 import timeit
1952 import statistics
1953 import pandas as pd
1954 import numpy as np
1955 np.random.seed(0)
1956 N = 100000
1957 data = np.arange(N)
1958 arr = pd.array(np.where(np.random.rand(N) > 0.1, data, np.nan),
1959               ↪ dtype='Int32')
1960 def workload():
1961     arr._hash_pandas_object(encoding='utf-8',
1962                             ↪ hash_key='10000000000000000', categorize=False)
1963
1964 # ---- AUTO-GENERATED ISOLATION HARNESS (timeit-in-child, fork-safe)
1965 ↪ ----
1966 import statistics as _statistics
1967 import multiprocessing as _mp
1968 import timeit as _timeit
1969
1970 def _child_once(number: int):
1971     setup_fn = (lambda: None)
1972     t = _timeit.Timer(workload, setup=setup_fn)
1973     return t.timeit(number)
1974
1975 # TOP-LEVEL target: picklable under 'fork'
1976 def _child_target(q, number):
1977     try:
1978         dur = _child_once(number)
1979         q.put(dur)
1980     except BaseException:
1981         import traceback
1982         traceback.print_exc()
1983         q.put(None)
1984
1985 def _run_isolated(number: int, repeat: int, start_method: str =
1986 ↪ "fork"):
1987     ctx = _mp.get_context(start_method)
1988     results = []
1989     for _ in range(repeat):
1990         q = ctx.Queue()
1991         p = ctx.Process(target=_child_target, args=(q, number))
1992         p.start()
1993         dur = q.get()
1994         p.join()
1995         if dur is None or p.exitcode != 0:
1996             raise RuntimeError("Child run failed--see traceback
1997 ↪ above.")
1998         results.append(dur)
1999     return results
2000
2001 if __name__ == "__main__":
2002     _number = 5
2003     _repeat = 1000
2004     runtimes = _run_isolated(_number, _repeat, start_method="fork")
2005     runtimes_view = runtimes

```

1998
 1999
 2000
 2001
 2002
 2003
 2004
 2005
 2006
 2007
 2008
 2009
 2010
 2011
 2012
 2013
 2014
 2015
 2016
 2017
 2018
 2019
 2020
 2021
 2022
 2023
 2024
 2025
 2026
 2027
 2028
 2029
 2030
 2031
 2032
 2033
 2034
 2035
 2036
 2037
 2038
 2039
 2040
 2041
 2042
 2043
 2044
 2045
 2046
 2047
 2048
 2049
 2050
 2051

```
print("Mean:", _statistics.mean(runtimes_view))
print("Std Dev:", _statistics.stdev(runtimes_view) if
  ↪ len(runtimes_view) > 1 else 0.0)
```

K COMPARISON OF LM GENERATED EDITS VERSUS EXPERTS

We provide the raw diffs with in-line comments from Figures 8 and 9 below in Figures 17 and 18. We also include an additional diff for 19 comparison. In the main text, we removed comments and surrounding lines to focus only on the lines changed between expert and LM generated diff.

L SPEEDUP RATIO STABILITY ANALYSIS

We also conduct a small study analyzing the variability and stability in our aggregate metric, Speedup Ratio (SR), since individual workload performance measurement can have minor variability. However, we find that SR does not vary significantly between independent runs. While raw runtime may exhibit micro-variance, the expert-grounded speedup normalization used in SR acts as an effective noise filter.

To demonstrate this, we performed three independent evaluation runs of the patches generated by three selected LM agents as well as the expert edit, computing the SR for each run. As shown in Table 11, the variance in the final SR score is negligible ($< 0.5\%$), confirming that measurement noise does not affect the relative ranking or assessment of agent capabilities. This observation aligns with related work finding that conversion to normalized or percentage-based speedup metrics inherently reduces variability in final scores ?.

Table 11: Stability of Speedup Ratio (SR) across three independent runs. The low spread confirms that SR is robust to microbenchmark runtime variance.

System	SR Run 1	SR Run 2	SR Run 3	SR Spread
Expert	1.000×	0.9986×	1.0026×	0.0040×
Claude 3.7 Sonnet (OpenHands)	0.04761×	0.04796×	0.04779×	0.00035×
GPT-5 Mini (OpenHands)	0.02115×	0.02115×	0.02128×	0.00013×
Gemini 2.5 Flash (OpenHands)	0.007941×	0.007947×	0.007935×	0.00012×

```

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

```

<pre> --- a/pandas/core/arrays/arrow/array.py +++ b/pandas/core/arrays/arrow/array.py @@ -406,8 +406,14 @@ def _cmp_method(self, other, op): f"{op.__name__} not implemented for { type(other)}") - - result = result.to_numpy() - return BooleanArray._from_sequence(result) + + if result.null_count > 0: + # GH50524: avoid conversion to object for + # better perf + values = pc.fill_null(result, False). + to_numpy() + mask = result.is_null().to_numpy() + else: + values = result.to_numpy() + mask = np.zeros(len(values), dtype=np. + bool_) + return BooleanArray(values, mask) + + def _evaluate_op_method(self, other, op, + arrow_funcs): + pc_func = arrow_funcs[op.__name__] </pre>	<pre> --- a/pandas/core/arrays/arrow/array.py +++ b/pandas/core/arrays/arrow/array.py @@ -406,8 +406,16 @@ class ArrowExtensionArray(OpsMixin, ExtensionArray): f"{op.__name__} not implemented for { type(other)}") - - result = result.to_numpy() - return BooleanArray._from_sequence(result) + + # Fast path: if there are no nulls, we can + # avoid the expensive BooleanArray creation + if result.null_count == 0: + # Ensure we get a boolean numpy array + result_np = result.to_numpy().astype(bool +) + return BooleanArray(result_np, np.zeros(+ len(result), dtype=bool)) + + # Slow path: handle nulls + result_np = result.to_numpy().astype(bool) + mask = result.is_null().to_numpy() + return BooleanArray(result_np, mask) + + def _evaluate_op_method(self, other, op, + arrow_funcs): + pc_func = arrow_funcs[op.__name__] </pre>
--	---

Figure 17: **Left:** Expert patch on instance pandas-dev__pandas-50524 optimizing a workload via avoiding a conversion to object dtype (20.5× speedup). **Right:** CLAUDE 3.7 SONNET (OPENHANDS) instead identifies a different fast path optimization when no null elements are present, but only achieves a 2.3× speedup (scoring a speedup ratio of 0.113×).

```

2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

```

<pre> --- a/pandas/core/series.py +++ b/pandas/core/series.py @@ -1818,7 +1818,7 @@ def to_dict(self, into: type[dict] = dict) -> dict: else: # Not an object dtype => all types will # be the same so let the default # indexer return native python type - return into_c((k, v) for k, v in self. + items()) + return into_c(self.items()) + + def to_frame(self, name: Hashable = lib. + no_default) -> DataFrame: + """ </pre>	<pre> --- a/pandas/core/series.py +++ b/pandas/core/series.py @@ -1816,9 +1816,18 @@ class Series(base. IndexOpsMixin, NDFrame): # type: ignore[misc] if is_object_dtype(self): return into_c((k, maybe_box_native(v)) for k, v in self.items()) else: - - # Not an object dtype => all types will - # be the same so let the default - # indexer return native python type - return into_c((k, v) for k, v in self. + items()) + + # Not an object dtype => use numpy + # fastpath to produce native python types + # by converting the underlying values + # to a python list in C and zipping + # with the index. This reduces Python- + # level boxing overhead. + values = getattr(self, "_values", None) + if values is None: + return into_c((k, v) for k, v in + self.items()) + try: + list_vals = values.tolist() + except Exception: + # fallback to generic iteration + list_vals = [v for v in values] + return into_c(zip(self.index, list_vals +)) + + def to_frame(self, name: Hashable = lib. + no_default) -> DataFrame: + """ </pre>
--	---

Figure 18: **Left:** Expert edit on pandas-dev__pandas-50089, optimizing Series.to_dict by replacing a generator of (k, v) pairs with the view self.items(), eliminating per-element tuple allocation. **Right:** GPT-5 MINI (OPENHANDS) converts the underlying array to a Python list and zips with the index to reduce Python-level boxing during iteration (achieving a 1.98× speedup vs. the expert’s 1.38×).

```

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

diff --git a/lib/matplotlib/transforms.py b/lib/
matplotlib/transforms.py
index c3d1b79291..f85025efff 100644
--- a/lib/matplotlib/transforms.py
+++ b/lib/matplotlib/transforms.py
@@ -1992,6 +1992,11 @@ class Affine2D(Affine2DBase):
     self.invalidate()
     return self

+ # Cache for rotation matrices
+ _rotation_cache = {}
+ _last_theta = None
+ _last_rotate_mtx = None
+
def rotate(self, theta):
    """
    Add a rotation (in radians) to this transform
    in place.
@@ -2000,10 +2005,17 @@ class Affine2D(Affine2DBase):
    calls to :meth:`rotate`, :meth:`rotate_deg`,
    :meth:`translate`
    and :meth:`scale`.
    """
    a = math.cos(theta)
    b = math.sin(theta)
    rotate_mtx = np.array([[a, -b, 0.0], [b, a,
0.0], [0.0, 0.0, 1.0]],
float)
# Cache the rotation matrix for repeated
calls with the same angle
if theta == self._last_theta and self.
_last_rotate_mtx is not None:
    rotate_mtx = self._last_rotate_mtx
else:
    a = math.cos(theta)
    b = math.sin(theta)
    rotate_mtx = np.array([[a, -b, 0.0], [b,
a, 0.0], [0.0, 0.0, 1.0]],
float)
    self._last_theta = theta
    self._last_rotate_mtx = rotate_mtx

self._mtx = np.dot(rotate_mtx, self._mtx)
self.invalidate()
return self

diff --git a/lib/matplotlib/transforms.py b/lib/
matplotlib/transforms.py
index c3d1b7929128..b0456773f6a4 100644
--- a/lib/matplotlib/transforms.py
+++ b/lib/matplotlib/transforms.py
@@ -2002,9 +2002,16 @@ def rotate(self, theta):
    """
    a = math.cos(theta)
    b = math.sin(theta)
    rotate_mtx = np.array([[a, -b, 0.0], [b, a,
0.0], [0.0, 0.0, 1.0]],
float)
    self._mtx = np.dot(rotate_mtx, self._mtx)
    mtx = self._mtx
    # Operating and assigning one scalar at a
time is much faster.
    (xx, xy, x0), (yx, yy, y0), _ = mtx.tolist()
    # mtx = [[a -b 0], [b a 0], [0 0 1]] * mtx
    mtx[0, 0] = a * xx - b * yx
    mtx[0, 1] = a * xy - b * yy
    mtx[0, 2] = a * x0 - b * y0
    mtx[1, 0] = b * xx + a * yx
    mtx[1, 1] = b * xy + a * yy
    mtx[1, 2] = b * x0 + a * y0
    self.invalidate()
    return self

```

Figure 19: **Left:** Expert patch on instance matplotlib__matplotlib-22108 optimizing a rotation transform workload, avoiding numpy arithmetic overhead by operating and assigning one scalar at a time (1.9× speedup). **Right:** CLAUDE 3.7 SONNET (OPENHANDS) instead identifies a last rotation caching mechanism, and achieves a 2.4× speedup (scoring a speedup ratio of 1.292×).