

# Magistrate: A Repository-Context Aware Multi-Agent Framework for Automated Code Review

Anonymous ACL submission

## Abstract

As Large Language Models make code writing faster than ever, the bottleneck has shifted to the critical task of reviewing code before merge: ensuring correctness, security, and compliance with project guidelines. We introduce **Magistrate**, an automated code review framework designed to augment human reviewers by detecting complex, cross-file logic errors that single-pass analysis misses. Unlike existing tools that analyze isolated diffs, Magistrate employs a hierarchical multi-agent architecture: a *Delegator* partitions changed files into semantically coherent batches, while parallel *IssueDetector* agents combine static analysis (ast-grep) with LLM-based semantic reasoning over the full repository context. We also present **Magistrate-Bench**, a benchmark of 2,042 Pull Requests across 12 programming languages, with complete repository workspaces cached at evaluation time to enable realistic dependency tracing. Evaluated on a 108-PR subset across four frontier models (Gemini 3 Flash, Devstral, Minimax M2, and Grok 4.1 Fast), Magistrate consistently improves F1 scores by 2.2–5.5× over single-shot baselines, with the best model (Gemini 3 Flash) achieving an F1 of 0.214 and a hallucination rate of just 0.6%. Across all models, Magistrate identified **997 valid issues** that human reviewers had overlooked, demonstrating its utility as a complementary filter for objective correctness issues.

## 1 Introduction

Automated code review remains a formidable challenge for AI systems despite significant advances in code understanding. While Large Language Models (LLMs) have achieved state-of-the-art results in code *generation* (Jimenez et al., 2024), their application to *review* has been constrained by insufficient access to systemic context. Current benchmarks and tools often treat a Pull Request (PR) as an isolated text-processing task (feeding a

model a diff snippet and expecting a comprehensive critique). This “Diff-Only Myopia” (Zeng et al., 2025) prevents agents from reasoning about global implications of local changes, such as breaking downstream dependencies or violating project-wide invariants.

Recent benchmarks have begun addressing this limitation. SWR-Bench (Zeng et al., 2025) and CodeFuse-CR-Bench (Guo et al., 2025) establish that repository-level context is essential for realistic evaluation. ContextCRBench (Hu et al., 2025) further demonstrates the importance of enriched metadata such as issue descriptions and PR context. However, while these benchmarks provide *data* for repository-level review, they do not prescribe *architectural solutions* for agents to effectively consume massive context without suffering from information overload or context window limitations.

We present **Magistrate**, a multi-agent framework designed to orchestrate hybrid static and semantic analysis over the complete repository state. Unlike single-model approaches or simple tool loops, Magistrate employs a structured three-phase architecture: **(1) Delegation**, **(2) Detection**, and **(3) Aggregation**. A Delegator agent partitions changed files into semantically coherent batches based on module boundaries and test-implementation coupling. Specialized IssueDetector agents then analyze each batch in parallel, combining static pre-scanning via ast-grep with LLM-based reasoning augmented by code exploration tools (ripgrep, git log, cat). A final aggregation phase deduplicates findings and produces a structured review report.

Critically, we position Magistrate not as a replacement for human reviewers, but as a **complementary filter**. Our analysis reveals that human reviewers predominantly focus on subjective style and design decisions, often overlooking subtle logic errors. By contrast, Magistrate is architected to prioritize objective correctness. This distinct

084	focus explains the system’s divergent performance	2.3 Multi-Agent Systems	131
085	characteristics: while it may not reproduce every	Multi-agent architectures like RevAgent (Li et al.,	132
086	human stylistic comment (lower Recall on ground	2025) and CodeAgent (Ma et al., 2024) assign	133
087	truth), it excels at surfacing critical logic bugs that	specialized roles (Reviewer, Author) to agents.	134
088	humans miss.	However, conversational systems can suffer from	135
089	Our core contributions are:	"prompt drifting." Agentless (Xia et al., 2024)	136
090	1. <b>Magistrate Architecture:</b> A multi-agent	showed that hierarchical localization can outper-	137
091	framework that orchestrates delegation and	form complex agents on repair tasks. Magistrate	138
092	hybrid static/semantic analysis to detect com-	adopts this hierarchical philosophy (Delegation →	139
093	plex cross-file issues.	Detection) but retains the tool-use capabilities of	140
094	2. <b>Magistrate-Bench:</b> A dataset of 2,042 PRs	SWE-agent (Yang et al., 2024) for active verifica-	141
095	across 12 languages, featuring full repository	tion, avoiding drift through "Stateless Batching."	142
096	workspaces cached at evaluation time to en-	2.4 Evaluation Metrics	143
097	able realistic dependency tracing.	Traditional metrics like BLEU correlate poorly	144
098	3. <b>Evaluation Protocol:</b> A rigorous “Valid	with review quality. We adopt the philosophy of	145
099	Extras” methodology that uses workspace-	CRScore (Naik et al., 2025) and DeepCRCEval (Lu	146
100	backed verification to distinguish between	et al., 2025) by using a "Valid Extras" methodology	147
101	hallucinations and valid issues overlooked by	that validates discoveries against the workspace,	148
102	human reviewers.	distinguishing real bugs from hallucinations.	149
103	<b>2 Related Work</b>	<b>3 Magistrate Framework</b>	150
104	We situate our contributions within code review	Magistrate is a hierarchical multi-agent system de-	151
105	automation, benchmarks, and multi-agent systems.	signed around a central MagistrateOrchestrator	152
106	<b>2.1 Code Review Automation</b>	that coordinates the review process. When a pull	153
107	Code review is fundamentally a knowledge task	request is submitted for review, the orchestrator	154
108	requiring broad contextual awareness. Recent sur-	first dispatches the task to a Delegator agent,	155
109	veys (Tufano et al., 2024; Heumüller and Ortmeier,	which analyzes the changed files and creates	156
110	2025) document the limitations of LLM-based ap-	an intelligent execution plan. This plan is then	157
111	proaches in handling complex, context-dependent	executed by parallel IssueDetector agents, each	158
112	scenarios. A key emerging direction is integrat-	analyzing a batch of related files. Finally, the	159
113	ing static analysis: Jaoua et al. (Jaoua et al., 2025)	orchestrator aggregates all findings into a unified	160
114	and BitsAI-CR (Sun et al., 2025) demonstrate that	review report (Figure 1).	161
115	grounding LLMs in static findings reduces hallu-	<b>3.1 Delegation</b>	162
116	cinations. Magistrate extends this by combining	The Delegator agent performs two critical func-	163
117	deterministic static analysis (ast-grep) with active	tions. First, it analyzes the PR metadata to gen-	164
118	workspace exploration, prioritizing correctness	erate a high-level PRSummary, establishing the	165
119	over convention.	intent of the changes. Second, it generates a	166
120	<b>2.2 Benchmarks for Code Review</b>	DelegationPlan that groups files into semanti-	167
121	Benchmark design has evolved from function-level	cally coherent batches using heuristics such as	168
122	to repository-level evaluation. SWR-Bench (Zeng	module cohesion and test-implementation pairing.	169
123	et al., 2025) and CodeFuse-CR-Bench (Guo et al.,	<b>3.2 Analysis</b>	170
124	2025) established that diff-only evaluation is in-	In the analysis stage, each DelegationBatch is	171
125	sufficient. ContextCRBench (Hu et al., 2025)	processed by a dedicated IssueDetector agent	172
126	demonstrated the value of enriched textual context.	instance running in parallel. Each detector oper-	173
127	Table1 summarizes the landscape. Magistrate-	ates in a hybrid mode combining static and seman-	174
128	Bench extends these by caching complete reposi-	tic analysis:	175
129	tory workspaces at evaluation time, enabling realis-	<b>Static Pre-Scan.</b> The agent first executes	176
130	tic dependency tracing across 12 languages.	ast-grep with a suite of security and quality rules.	177

Feature	Magistrate-Bench	SWR-Bench	CodeFuse-CR	ContextCRBench	SWE-bench
Task Unit	Pull Request	Pull Request	Pull Request	Code Chunk	GitHub Issue
Scale	2,042 PRs	1,000 PRs	601 PRs	67,910 Chunks	2,294 Issues
Languages	12	Python	Python	9	Python
Code Context	Full Repository	Full Repository	Full Repository	Diff Only	Full Repository
Verification	Tool-Verified	Text-Match	Model Score	Text-Match	Test-Execution

Table 1: Comparison of code review benchmarks. Magistrate-Bench distinguishes itself by combining **Full Repository** context (unlike ContextCRBench) with active **Tool-Verified** evaluation (unlike SWR-Bench).

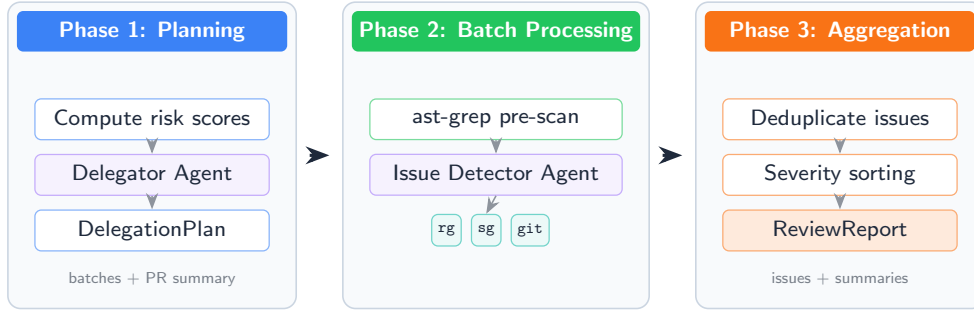


Figure 1: Magistrate execution flow. Phase 0 prepares the workspace from cached repo tarballs; Phase 1 Delegator pulls README/CONTRIBUTING and commit context, batches related files, and emits the PR summary. Phase 2 runs parallel IssueDetector instances with ast-grep prescan, tool-backed evidence gathering (batch\_execute\_tool), and post-LLM verification. Phase 3 deduplicates by file/line/title, orders by severity, and writes review.json.

These deterministic findings are injected into the LLM’s context as hints.

**Semantic Analysis.** The LLM validates the pre-scan findings while searching for deeper logic errors that static analysis cannot catch. It utilizes a guarded tool palette (detailed in Appendix C) to explore the codebase.

**Verification.** A programmatic post-processing layer validates that every reported issue cites a file path and line number that physically exists in the PR, filtering out hallucinated references.

### 3.3 Aggregation

Upon the completion of all detector agents, the MagistrateOrchestrator aggregates the validated issues from every batch, performing a deduplication pass to merge identical findings before producing the final ReviewReport.

### 3.4 End-to-End Execution Workflow

The complete review pipeline proceeds through the following stages:

- Task Submission:** A review request is queued, triggering a dedicated worker container for isolated execution.
- Workspace Preparation:** The worker restores a cached repository workspace corre-

sponding to the exact base and head commits, ensuring byte-for-byte reproducibility.

- Phase 0 (Static Prescan):** The IssueDetector executes a static rule-based prescan using ast-grep, producing grounded evidence injected into subsequent prompts.
- Phase 1 (Delegation):** The Delegator produces a PRSummary and DelegationPlan, grouping files into semantically coherent batches to mitigate context-window overflow.
- Phase 2 (Parallel Batch Analysis):** Each batch is processed by an IssueDetector instance. Detectors operate in parallel to improve throughput.
- Phase 2b (Verification and Filtering):** All issues generated by detectors undergo a structured post-processing stage. Invalid line numbers and nonexistent file paths are removed.
- Phase 3 (Aggregation):** The orchestrator deduplicates issues, normalizes severity labels, and emits the final review.json artifact.
- Return to Server:** The worker uploads results and terminates. No state persists, guaranteeing isolation.

227	This workflow enforces strict grounding, keeps	<b>5 Experimental Design</b>	273
228	LLM context under control, and ensures consistent	Evaluating automated code review presents a funda-	274
229	behavior through typed I/O schemas.	mental challenge: agent-reported issues not present	275
230		in ground truth may represent either hallucinations	276
231	<b>4 Magistrate-Bench Dataset</b>	<i>or</i> valid issues that human reviewers overlooked.	277
232	To rigorously test repository-aware agents, we	Traditional precision metrics conflate these cases,	278
233	constructed <b>Magistrate-Bench</b> , a new benchmark	penalizing systems for finding real bugs. We	279
234	designed to overcome the contextual myopia of	address this through a rigorous verification pipeline	280
235	existing datasets.	that distinguishes <i>Hallucinations</i> from <i>Valid Ex-</i>	281
236		<i>tras</i> .	282
237		<b>5.1 Research Questions</b>	283
238		We structure our evaluation around the following:	284
239	<b>4.1 Construction Pipeline</b>	• <b>RQ1 (Effectiveness):</b> Does repository-aware	285
240	Our collection process ensures high-quality, repro-	orchestration improve the detection of logic	286
241	ducible ground truth through a rigorous multi-stage	errors (Recall/F1) compared to single-shot	287
242	pipeline (Figure 2).	baselines?	288
243	<b>Repository &amp; PR Selection:</b> We targeted active	• <b>RQ2 (Reliability):</b> Does the workspace-	289
244	repositories (> 500 stars, > 1500 PRs, active	backed verification layer significantly reduce	290
245	within 150 days) across 12 languages, explicitly	hallucination rates?	291
246	excluding tutorials and algorithm lists. We selected	• <b>RQ3 (Efficiency):</b> What is the trade-off be-	292
247	only merged PRs with a minimum of 6 distinct	tween detection quality and computational	293
248	human review threads.	cost?	294
249	<b>Temporal State Reconstruction:</b> Review com-	<b>5.2 Experimental Setup</b>	295
250	ments typically address the state <i>before</i> fixes are	We evaluate on a stratified subset of 108 PRs from	296
251	pushed. We algorithmically identify the last com-	Magistrate-Bench (Python, JS/TS, Java, Rust). For	297
252	mit hash immediately preceding the first human	each PR, we compare (i) a single-shot <b>Raw</b> base-	298
253	review timestamp. We then fetch the cumulative	line (prompting the model with the diff) against	299
254	diff from the PR base to this specific commit,	(ii) the orchestrated <b>Magistrate</b> system, using the	300
255	ensuring the agent evaluates the exact state seen by	same underlying LLM backbone.	301
256	the human reviewer.	<b>5.3 Verification Methodology</b>	302
257	<b>Hybrid Extraction &amp; Verification:</b> We employ	To ensure the soundness of our "Valid Extras"	303
258	a hybrid strategy to minimize noise. We extract	(issues found by agents but missed by humans), we	304
259	exact file paths, lines, and comment bodies via the	employ a multi-stage verification process designed	305
260	GitHub API. An LLM agent (Temperature=0) acts	to minimize "LLM-as-a-judge" bias:	306
261	as a classifier, categorizing threads into issue types	<b>1. Structural Filtering:</b> We programmatically	307
262	(e.g., Logic, Security). Finally, a programmatic	verify that every reported file path and line number	308
263	verification step cross-references every extracted	exists in the PR diff. This step rigorously filters	309
264	issue against the actual diff, discarding findings	out hallucinations referencing non-existent files or	310
265	that reference unchanged files.	invalid line coordinates, ensuring that all candidate	311
266	<b>Diversity &amp; Standardization:</b> We enforced	issues are spatially grounded before semantic veri-	312
267	a strict cap of 10 PRs per repository to ensure	fication.	313
268	generalization. Extracted issues are mapped to a	<b>2. Claim-Blind Evidence Gathering:</b> To verify	314
269	standardized taxonomy to facilitate direct compari-	semantic claims (e.g., "This function call is missing	315
270	son with existing benchmarks. See Appendix D for	a required argument"), we use a strictly <b>Claim-</b>	316
271	detailed distribution statistics.	<b>Blind</b> protocol to eliminate confirmation bias. An	317
272	<b>4.2 Dataset Statistics</b>	"Evidence Agent" is given a neutral query derived	318
Magistrate-Bench exceeds the scale and diversity	of prior benchmarks (see Table 5 in Appendix). It	from the issue (e.g., "What are the arguments for	319
contains 2,042 PRs covering 12 languages, with an	average of 5.94 issues per PR.	function X?") <i>without access to the agent's claim</i>	320

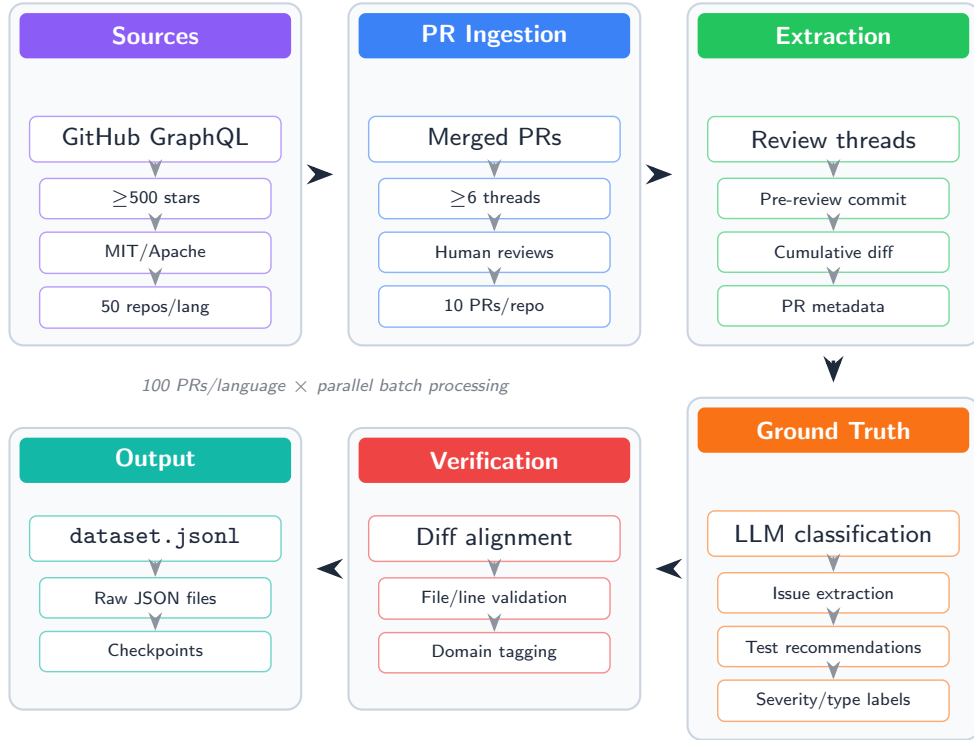


Figure 2: Magistrate-Bench construction. Filters select merged PRs with human reviews across diverse languages. We scrape metadata/diffs/threads and clone repos; build unified diffs and cache full-repo workspaces for evaluation; parse review threads into structured issues/tests with an LLM parser (T=0); remove bots/LGTM noise; publish JSONL/CSV.

321 *or verdict*. It must use repository tools (`grep`,  
 322 `cat`) to retrieve the ground truth code state. This  
 323 prevents the verifier from simply agreeing with the  
 324 detector’s plausible-sounding rhetoric.

325 **3. Objective Verdict:** A separate Verifier  
 326 compares the retrieved ground-truth evidence against  
 327 the agent’s original claim. This separation ensures  
 328 the system relies on codebase reality, not rhetorical  
 329 plausibility.

### 330 5.4 Enhanced Metrics

331 We report verification-aware metrics:

- 332 • *TP*: Matches ground truth.
- 333 • *V (Valid Extras)*: Unmatched but verified as  
 334 correct.
- 335 • *H (Hallucinations)*: Unmatched and verified  
 336 as incorrect.

337 **Enhanced Precision:**  $P = \frac{TP+V}{TP+V+H}$

338 **Recall:**  $R = \frac{TP}{|GT|}$ , where  $|GT|$  is the total number  
 339 of issues in the ground truth.

340 **F1 Score:**  $F1 = 2 \cdot \frac{P \cdot R}{P+R}$

341 **Hallucination Rate:**  $HR = \frac{H}{TP+V+H}$

342 These metrics reward agents for genuine discover-  
 343 ies (*V*) while penalizing incorrect hallucinations

(*H*), providing a fair assessment of real-world  
 344 utility. 345

## 346 6 Results

347 We report verified results on the 108-PR evalu-  
 348 ation subset described in Section 5, comparing  
 349 single-shot **Raw** baselines against the orchestrated  
 350 **Magistrate** system across four frontier models.

### 351 6.1 Effectiveness (RQ1)

352 Table 2 summarizes performance. Magistrate im-  
 353 proves F1 by 2.2–5.5× and increases verified *Valid*  
 354 *Extras*. Gemini 3 Flash achieves the strongest  
 355 results, while Devstral shows the largest relative  
 356 improvement over its baseline.

357 It is important to contextualize the absolute  
 358 Recall scores. Human-annotated ground truth is  
 359 heavily skewed toward subjective style and design  
 360 preferences (see Table 3), which Magistrate is often  
 361 prompted to ignore in favor of objective correctness.  
 362 Consequently, matching ground truth requires not  
 363 just finding bugs, but mimicking human stylistic  
 364 preferences. The "Valid Extras" metric is therefore  
 365 a more accurate proxy for the system’s utility in

Model	Agent	Enh. Prec.	Recall	F1	Valid Extras	Hallucinations	Hall. Rate
Gemini 3 Flash	Raw	0.804	0.052	0.097	94	32	19.6%
	Magistrate	<b>0.994</b>	<b>0.120</b>	<b>0.214</b>	<b>432</b>	3	<b>0.6%</b>
Devstral	Raw	0.600	0.007	0.014	13	12	40.0%
	Magistrate	<b>0.903</b>	<b>0.040</b>	<b>0.077</b>	<b>194</b>	24	<b>9.7%</b>
Minimax M2	Raw	0.737	0.038	0.072	46	26	26.3%
	Magistrate	<b>0.820</b>	<b>0.099</b>	<b>0.177</b>	<b>197</b>	59	<b>18.0%</b>
Grok 4.1 Fast	Raw	0.625	0.013	0.026	10	12	37.5%
	Magistrate	<b>0.787</b>	<b>0.038</b>	<b>0.073</b>	<b>174</b>	55	<b>21.3%</b>

Table 2: Overall performance comparison across models. Magistrate consistently improves over Raw baseline across all models, with F1 improvements of 2.2–5.5× and substantial reductions in hallucination rate.

detecting objective defects. Magistrate identifies **997** verified issues absent from human review, indicating strong complementary coverage.

## 6.2 Reliability (RQ2)

Magistrate consistently reduces hallucination rates (Table 2). Gemini 3 Flash drops from 19.6% to 0.6%, with other models showing similar trends. This validates the effectiveness of the workspace-backed verification layer.

## 6.3 Qualitative Alignment

Table 3 reveals a clear divergence: human reviewers focus on **Style/Design** (500 issues), while Magistrate excels at **Bug/Logic** (200 Valid Extras for Gemini). This confirms our hypothesis that AI agents function best as **complementary filters** for objective correctness, allowing humans to focus on high-level architectural decisions and subjective style.

Category	GT	Gemini VE	Devstral VE	Minimax VE
Style/Design	500	103	41	66
Bug/Logic	121	<b>200</b>	<b>99</b>	<b>87</b>
Security	7	8	3	3
Docs	85	50	14	20
Testing	0	22	32	15

Table 3: Condensed category analysis (Valid Extras).

## 6.4 Cost-Efficiency (RQ3)

Magistrate consumes significantly more tokens (10–60×) than single-shot baselines. However, Table 4 shows that when measured per *valid issue found*, the trade-off is favorable. Grok 4.1 Fast achieves near-parity cost efficiency while finding 10× more issues.

While the absolute token cost (\$0.50–\$1.00 per PR) is higher than a simple API call, it must be weighed against the cost of a software defect.

Model	Token OH	Valid OH	Efficiency
Gemini 3 Flash	62.6×	4.0×	15.8×
Devstral	53.2×	12.4×	4.3×
Minimax M2	42.1×	3.7×	11.5×
Grok 4.1 Fast	10.8×	11.2×	1.0×

Table 4: Token efficiency analysis. Token OH = total token overhead (Magistrate/Raw), Valid OH = valid issues ratio (Magistrate/Raw), Efficiency = tokens per valid issue ratio (Magistrate/Raw). Despite higher token usage, Magistrate achieves better issue detection density, with efficiency ratios ranging from 2.45× (Devstral on small PRs) to 22.9× (Gemini on small PRs).

Given that the industry average cost to fix a bug post-release is orders of magnitude higher than finding it during review, paying a marginal token cost to uncover critical logic errors (which humans missed) represents a high-ROI investment.

## 7 Conclusion and Future Work

We presented Magistrate, a repository-aware automated code review framework comprising: (1) a hierarchical multi-agent architecture that orchestrates delegation, detection, and aggregation; (2) Magistrate-Bench, a benchmark of 2,042 PRs across 12 languages with cached repository workspaces; and (3) an evaluation protocol that distinguishes hallucinations from valid discoveries through workspace-backed verification.

Our experiments on a 108-PR subset across four frontier models (Gemini 3 Flash, Devstral, Minimax M2, and Grok 4.1 Fast) demonstrate that orchestrated multi-agent review consistently improves F1 scores by 2.2–5.5× over single-shot baselines. The best-performing model (Gemini 3 Flash) achieves an F1 of 0.214 with an exceptionally low hallucination rate of 0.6%. Critically, across all models, Magistrate identified **997** verified issues that are absent from the human-review

ground truth, with individual models contributing 432 (Gemini), 194 (Devstral), 197 (Minimax), and 174 (Grok) Valid Extras respectively. This validates Magistrate’s utility as a complementary filter that excels at objective correctness issues rather than stylistic concerns.

The consistency of improvements across diverse model providers (Google, Mistral, Minimax, xAI) suggests that the Magistrate architecture provides robust value independent of the specific LLM backbone, making it a practical choice for production deployment.

**Future Work.** Three directions merit investigation: (1) *Dynamic Analysis Integration* (verifying suspected bugs through test execution or sandboxed evaluation to reduce false positives); (2) *Personalization* (adapting review agents to repository-specific conventions and maintainer preferences); and (3) *Heterogeneous Model Orchestration* (using high-capability models for delegation while routing granular detection to cost-effective alternatives).

## 8 Limitations

We identify several limitations in our current work concerning internal consistency and evaluation metrics.

**LLM Non-determinism.** LLMs exhibit inherent stochasticity. While we employ temperature=0 for the Delegator and low temperature (0.2) for analysis to reduce variance, complete elimination of stochasticity is not possible. Future work should include results from multiple evaluation runs with confidence intervals to establish statistical significance rigorously.

**Implementation Fidelity.** Bugs in the orchestration logic or tool wrappers could skew results. We mitigate this through extensive unit testing of the `batch_execute_tool` and manual verification of agent traces to ensure tools are being invoked and parsed correctly, but rare edge cases may persist.

**Ground Truth Quality.** Human code reviews are not a perfect oracle; reviewers miss bugs. We address this via the “Valid Extras” protocol, but there remains a risk that our workspace-backed verifier (LLM with tool access) itself makes mistakes. We audit a sample of verifier decisions to estimate this error rate, but scaling this audit is labor-intensive.

**Metric Validity.** Standard precision/recall metrics penalize agents for finding valid issues that

humans missed. Our modified metrics attempt to correct this, but rely on the accuracy of the verification step.

## 9 Ethics Statement

Our work introduces **Magistrate**, an automated code review framework. We adhere to the ACL Ethics Policy and consider the broader impact of this technology.

**Human-in-the-Loop.** Magistrate is designed to augment, not replace, human reviewers. Automated systems can hallucinate or miss subtle context (as noted in our Limitations). We explicitly position this tool as a “filter” for objective correctness, relying on human judgment for architectural and stylistic decisions. Over-reliance on automation without human oversight could lead to a degradation in code quality or the acceptance of subtle security vulnerabilities.

**Bias and Fairness.** The LLMs used in this study (Gemini, Devstral, etc.) are trained on public code repositories, which may contain biased language or non-inclusive terminology. Furthermore, the models may perform unevenly across different programming languages, potentially disadvantaging developers working in less-represented languages (though our benchmark covers 12 languages to mitigate this).

**Energy Consumption.** Multi-agent systems with iterative reasoning (like Magistrate) consume significantly more energy than single-shot inference. We analyze this trade-off in our results, highlighting that while computational cost is higher, the density of valid issues found justifies the resource usage for critical codebases. Future work will focus on optimizing token efficiency.

**Potential for Misuse.** While designed for defense (code review), the vulnerability detection capabilities could theoretically be adapted for malicious code auditing. However, we believe the defensive benefits of securing open-source software outweigh this risk, and the system requires a full repository context to function effectively.

## References

coderrabbitai. 2024. ast-grep-essentials: Community-led collection of essential ast-grep rules. <https://github.com/coderrabbitai/ast-grep-essentials>. Accessed: 2025.

514	Hanyang Guo, Xunjin Zheng, Zihan Liao, Hang Yu,	John Yang, Carlos E. Jimenez, Alexander Wettig, Kil-	566
515	Peng Di, Ziyin Zhang, and Hong-Ning Dai. 2025.	ian Lunt, Shunyu Yao, and Karthik Narasimhan.	567
516	CodeFuse-CR-Bench: A comprehensiveness-aware	2024. SWE-agent: Agent-computer interfaces enable	568
517	benchmark for end-to-end code review evaluation in	automated software engineering. <i>arXiv preprint</i>	569
518	Python projects. <i>arXiv preprint arXiv:2509.14856</i> .	<i>arXiv:2405.15793</i> .	570
519	Robert Heumüller and Frank Ortmeier. 2025. Previ-	Zhengran Zeng, Ruikai Shi, Keke Han, Yixin Li,	571
520	ously on... automating code review. <i>arXiv preprint</i>	Kaicheng Sun, Yidong Wang, Zhuohao Yu, Rui Xie,	572
521	<i>arXiv:2508.18003</i> .	Wei Ye, and Shikun Zhang. 2025. Benchmarking and	573
522	Ruida Hu, Xincheng Wang, Xin-Cheng Wen, Zhao	studying the LLM-based code review. <i>arXiv preprint</i>	574
523	Zhang, Bo Jiang, Pengfei Gao, Chao Peng, and	<i>arXiv:2509.01494</i> .	575
524	Cuiyun Gao. 2025. Benchmarking LLMs for fine-		
525	grained code review with enriched context in practice.		
526	<i>arXiv preprint arXiv:2511.07017</i> .		
527	Imen Jaoua, Oussama Ben Sghaier, and Houari Sahraoui.		
528	2025. Combining large language models with static		
529	analyzers for code review generation. <i>arXiv preprint</i>		
530	<i>arXiv:2502.06633</i> .		
531	Carlos E. Jimenez, John Yang, Alexander Wettig,		
532	Shunyu Yao, Kexin Pei, Ofir Press, and Karthik		
533	Narasimhan. 2024. SWE-bench: Can language		
534	models resolve real-world GitHub issues? <i>arXiv</i>		
535	<i>preprint arXiv:2310.06770</i> .		
536	Shuochuan Li, Dong Wang, Patanamon Thongta-		
537	nunam, Zan Wang, Jiuqiao Yu, and Junjie Chen.		
538	2025. Issue-oriented agent-based framework for au-		
539	tomated review comment generation. <i>arXiv preprint</i>		
540	<i>arXiv:2511.00517</i> .		
541	Junyi Lu, Xiaojia Li, Zihan Hua, Lei Yu, Shiqi		
542	Cheng, Li Yang, Fengjun Zhang, and Chun Zuo.		
543	2025. DeepCRCEval: Revisiting the evaluation of		
544	code review comment generation. <i>arXiv preprint</i>		
545	<i>arXiv:2412.18291</i> .		
546	Ke Ma, Zhaochen Meng, Weicheng Wang, Yu Wang,		
547	Ying Zhang, Min Zhang, and Zhiyuan Liu. 2024.		
548	Codeagent: Enhancing code generation with		
549	tool-integrated agent systems. <i>arXiv preprint</i>		
550	<i>arXiv:2402.02172</i> .		
551	Atharva Naik, Marcus Alenius, Daniel Fried, and Car-		
552	olyn Rose. 2025. CRScore: Grounding automated		
553	evaluation of code review comments in code claims		
554	and smells. <i>arXiv preprint arXiv:2409.19801</i> .		
555	Tianyang Sun and 1 others. 2025. BitsAI-CR: Au-		
556	tomated code review via LLM in practice. <i>arXiv</i>		
557	<i>preprint arXiv:2501.15134</i> .		
558	Rosalia Tufano, Ozren Dabić, Antonio Mastropaolo,		
559	Matteo Ciniselli, and Gabriele Bavota. 2024. Code		
560	review automation: Strengths and weaknesses of the		
561	state of the art. <i>arXiv preprint arXiv:2401.05136</i> .		
562	Chunqiu Steven Xia, Yinlin Deng, and Lingming		
563	Zhang. 2024. Agentless: Demystifying LLM-		
564	based software engineering agents. <i>arXiv preprint</i>		
565	<i>arXiv:2407.01489</i> .		

## A Detailed System Architecture

This appendix provides a comprehensive breakdown of the Magistrate architecture.

### A.1 Orchestration Layer

The orchestration layer manages the lifecycle of review tasks and coordinates computational resources via a FastAPI Server and Docker Orchestrator.

### A.2 Execution Layer

Worker containers provide the execution environment, built using a multi-stage Docker build process (Python + Bun runtime + system tools).

### A.3 Agent Layer

The agent layer implements the core review logic through a hierarchical multi-agent architecture. **MagistrateOrchestrator** coordinates the workflow. **Delegator** analyzes metadata to create an execution plan. **Issue Detector** analyzes batches in parallel using `batch_execute_tool`.

## B Agent I/O Schemas and Prompt Structure

We implement typed inputs/outputs for each agent via Pydantic models.

### Delegator (Planner) Prompt

**Role:** Analyze all changed files, generate a PR summary, and create optimal batches for parallel Issue Detector execution.

**Input:** List of DelegationCandidate objects with file paths, status, change metrics, and computed risk scores.

**Output:** DelegationPlan with ordered batches (each batch = 1-5 files), priority rankings, and rationale per batch. Also returns PRSummary.

```
class PRSummary(BaseModel):
    summary: str # concise PR overview
    key_changes: List[str] # bullet points
    overall_impact: str # risk assessment

class DelegationBatch(BaseModel):
    batch_name: str
    priority: int # lower = higher priority
    files: List[str] # file paths in batch
    rationale: str

class DelegationPlan(BaseModel):
    summary: PRSummary
    batches: List[DelegationBatch]
    notes: Optional[str]
```

### Issue Detector Batch Prompt

**Role:** Analyze a delegation batch of related files; detect issues with evidence; verify downstream impacts.

**Input:** Batch name, priority, rationale, list of FileChange objects, PR context.

**Output:** IssueDetectorBatchResult containing batch summary, shared analysis notes, shared tools used, and per-file FileAnalysisResult objects.

```
class FileAnalysisResult(BaseModel):
    file_path: str
    file_summary: str # 2-3 sentences
    issues: List[ReviewIssue]
    analysis_notes: str
    tools_used: List[str]

class IssueDetectorBatchResult(BaseModel):
    batch_name: str
    batch_summary: str
    shared_analysis_notes: str
    shared_tools_used: List[str]
    files: List[FileAnalysisResult]
```

## C Tool Boxes and Capabilities

### batch\_execute\_tool: Parallel Shell Execution

**Purpose:** Execute multiple shell commands concurrently for fast evidence gathering.

**Available Commands:** ripgrep (rg), ast-grep (ag/sg), sed, git, ls, find, cat, head, tail, wc, grep, awk.

**Usage:** Issue Detector constructs command lists (e.g., search for function calls, check git history, find TODOs).

**Example:**

```
batch_execute_tool([
    "rg -n 'validateToken' --glob '!
      node_modules/**' src/",
    "git log --oneline -5 src/auth.py",
    "ast-grep --lang py -p 'def $FUNC($ARGS)
      :'"
])
```

## D Dataset Statistics

Table 5 provides a detailed breakdown of the Magistrate-Bench dataset.

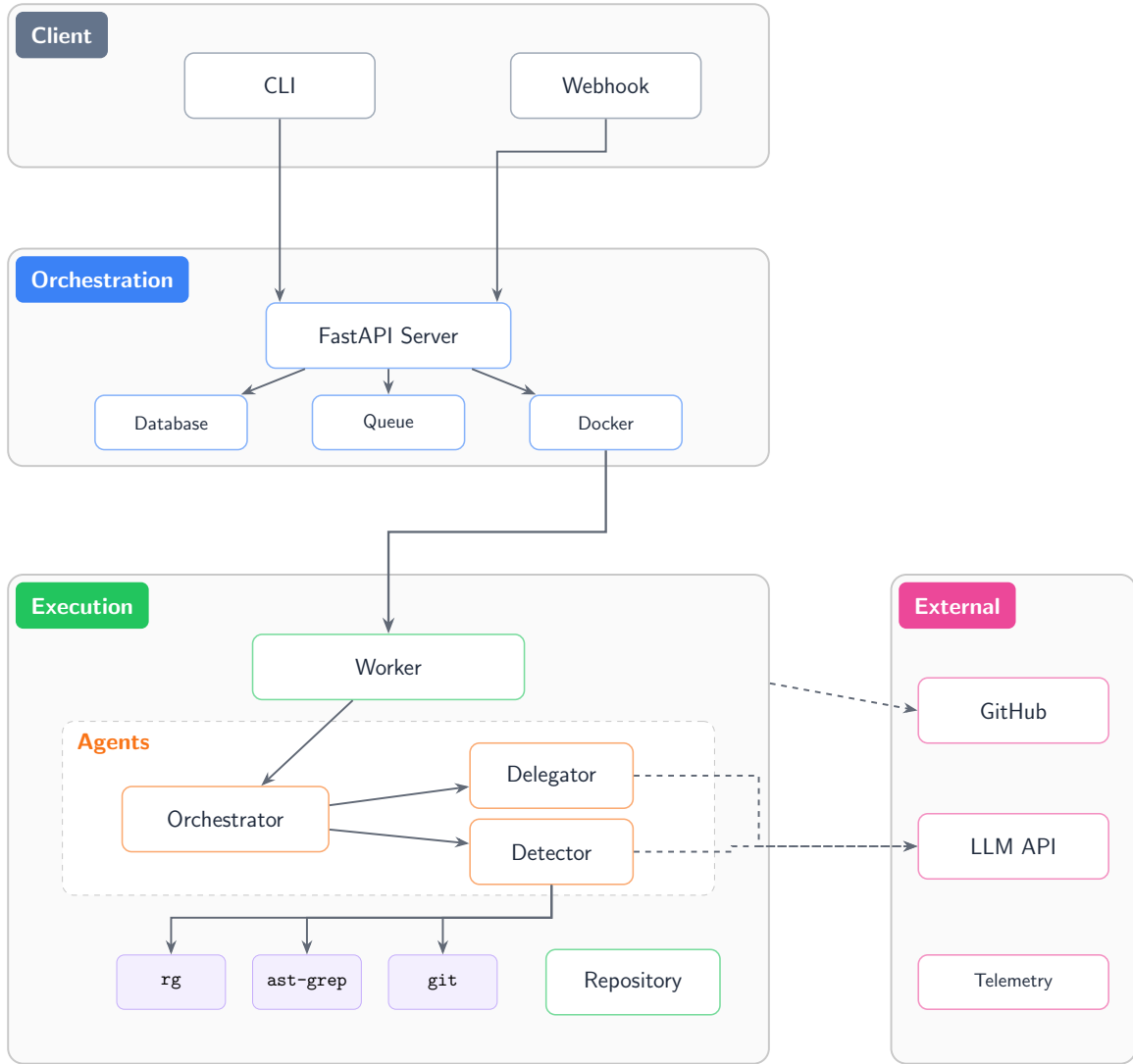


Figure 3: System architecture showing three layers: Orchestration, Execution, and Agent Layer.

General Stats		Issue Types			
Total PRs	2,042	Style	4,969		
Total Repos	245	Design	2,523		
Total Issues	12,121	Docs	1,931		
Avg. Issues/PR	5.94	Logic	1,120		
Avg. Tests/PR	0.29	Bug	973		
Language Distribution					
Python	389	Java	186	C#	98
JavaScript	272	Go	182	Swift	96
TypeScript	232	PHP	105	C++	90
Rust	207	Dart	99	Ruby	86

Table 5: Magistrate-Bench statistics showing dataset scale and issue type distribution.

## E Server-Client Architecture and System Integration

### E.1 Overview

Magistrate follows a three-layer architecture: (1) an orchestration layer (FastAPI server, SQLite

task database, Docker scheduler), (2) an execution layer (one worker container per pull request), and (3) an agent layer (Orchestrator, Delegator, and IssueDetector agents). This separation isolates infrastructure concerns from agent logic and ensures reproducible, containerized evaluation.

### E.2 Client Interaction

The review process begins when the client submits a request to POST `/api/review-pr`. The payload includes:

- repository identifier,
- commit SHA and diff metadata,
- PR title and description.

The server inserts a new task into the SQLite queue and returns a task identifier.

657	<b>E.3 Orchestration Layer</b>			
658	The FastAPI server acts as the controller for task			
659	creation, status tracking, and result retrieval. A			
660	lightweight Docker orchestrator continuously polls			
661	the task database. When a new task is available,			
662	the orchestrator provisions a fresh worker con-			
663	tainer, mounts the cached repository workspace,			
664	and provides the task metadata. No computation			
665	is performed on the host machine; all review logic			
666	executes inside the container.			
667	<b>E.4 Execution Layer</b>			
668	Each worker container executes the full pipeline:			
669	1. Extract cached repository tarball correspond-			
670	ing to the target commit.			
671	2. Run Phase 0 static prescan (via <code>ast-grep</code> ).			
672	3. Invoke the Delegator to generate the			
673	<code>PRSummary</code> and <code>DelegationPlan</code> .			
674	4. Spawn parallel <code>IssueDetector</code> runs over the			
675	assigned batches.			
676	5. Perform post-LLM verification and filtering.			
677	6. Aggregate all validated issues and write the			
678	final review. <code>json</code> .			
679	The container then uploads the results back to			
680	the server and terminates.			
681	<b>E.5 Agent Layer</b>			
682	All agents communicate using typed Pydantic mod-			
683	els to avoid drift in LLM output formats. The key			
684	structures include:			
685	• <code>DelegationCandidate</code> (file path, size, risk			
686	score),			
687	• <code>PRSummary</code> (overview, key changes, impact),			
688	• <code>DelegationPlan</code> (ordered batches with ratio-			
689	onale),			
690	• <code>FileAnalysisResult</code> (summary, issues,			
691	notes),			
692	• <code>ReviewIssue</code> (title, description, file, line,			
693	severity).			
694	All objects serialize to JSON, ensuring repro-			
695	ducibility across multiple evaluation runs.			
	<b>E.6 Tool Integration</b>			696
	The agents operate under a guarded tool palette:			697
	• <code>ast-grep</code> for syntactic pattern matching,			698
	• <code>ripgrep</code> for project-wide search,			699
	• <code>git log</code> and <code>git diff</code> for history mining,			700
	• <code>batch_execute_tool</code> for parallel shell exe-			701
	cution.			702
	These tools enable the agent to gather concrete			703
	evidence without hallucinated paths or nonexistent			704
	lines.			705
	<b>F Dataset, Models, and I/O Schemas</b>			706
	<b>F.1 Dataset Structure</b>			707
	Each example consists of:			708
	• PR metadata: commit SHA, diff hunks, file-			709
	names, change statistics,			710
	• repository workspace: extracted from cached			711
	tarballs,			712
	• ground truth issues: title, description, file, line,			713
	and human reviewer notes.			714
	The dataset is stored in JSONL format for effi-			715
	cient streaming.			716
	<b>F.2 ReviewIssue Schema</b>			717
	Agents emit issues in a strictly typed structure:			718
	<code>class ReviewIssue:</code>			719
	<code>title: str</code>			720
	<code>description: str</code>			721
	<code>file_path: str</code>			722
	<code>line: int</code>			723
	<code>severity: str</code>			724
	<b>F.3 Batch Issue Matching</b>			725
	To match agent issues to ground truth issues, we			726
	use a single-shot LLM call ( <i>BatchIssueMatcher</i> ).			727
	The model receives:			728
	• all agent issues,			729
	• all GT issues,			730
	• relevant diff and code snippets.			731
	It returns one-to-one correspondences and iden-			732
	tifies unmatched agent issues (false positives) and			733
	unmatched GT issues (false negatives).			734

735  
736  
737  
738  
739  
740

#### **F.4 Outputs and Persistence**

All outputs (summaries, batch results, detected issues, verification metadata) are written to a single review.json object per PR. This file mirrors the internal Pydantic structures and is used directly by the evaluation pipeline.

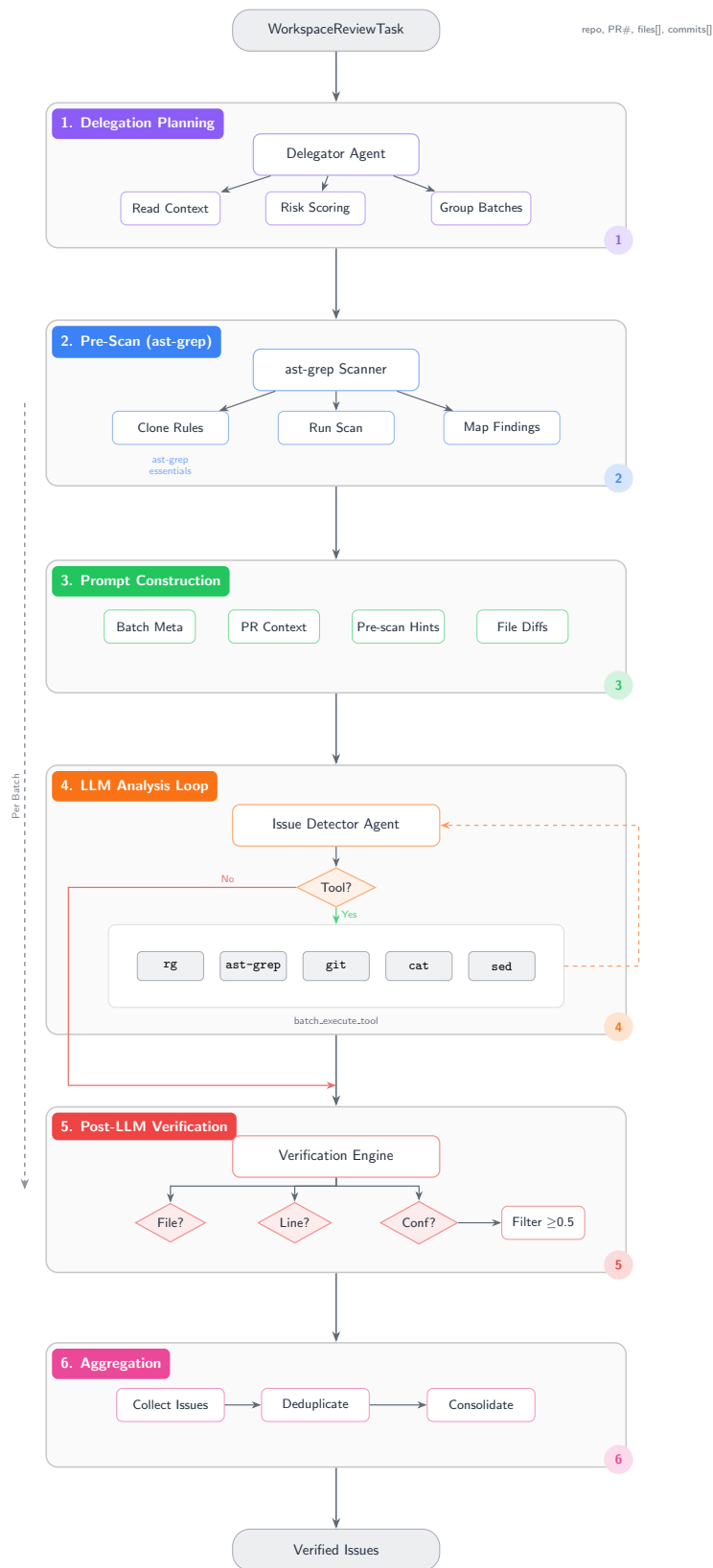


Figure 4: Issue Detector pipeline showing six phases: delegation planning, pre-scan using ast-grep essentials (coder-abbitai, 2024), prompt construction, LLM analysis with tool access, post-LLM verification, and aggregation.