



Anonymous Authors

## 1. Introduction

Automating research processes to accelerate scientific discovery are emerging. Existing works (The AI Scientist (Lu et al., 2024), AutoResearch (Karpathy, 2026), ARIS (Yang et al., 2026)) demonstrate the feasibility of AI scientists to automate research in different domains. However, existing works pursue a full automation paradigm without human supervision, which is prone to uncontrolled results, e.g., intention drift and hallucination (Miyai et al., 2026).

We argue that the right paradigm is the *co-scientist*: a system to *offloads the labour* of implementation, experiments, and paper writing, while the human *steers the science* with full creative and strategic control. This is captured in our design philosophy: **Full Automation with Full Control**.

We present **ARK** (**A**utomatic **R**esearch **K**it), a co-scientist framework that automates the full idea-to-paper pipeline while keeping the researcher in control at every stage. We design ARK as a model-agnostic multi-agent system whose seven agents and three-phase pipeline turn a research idea end-to-end into a *publication-ready* paper.

We propose **Research as a Service** (RaaS). We deploy ARK as a production-level application to turn research from a local, manual endeavor into a web service. To use ARK, the human scientist only needs to provide the research idea, then launch it in one click without any local setup. Each project runs in an isolated sandbox with dedicated dependencies, filesystem, and credentials to prevent cross-project contamination, and supports three interchangeable backends: *Local*, *Managed Cluster* (Slurm + Apptainer), and *Cloud* (Kubernetes + containerd, in progress). Users can monitor the research progress on the website or via Telegram, and inject instructions into the system at any time.

ARK is live at [idea2paper.org](https://idea2paper.org) (Anonymized). We prepared a setup-free account with API keys pre-configured, and we strongly invite all reviewers to experience ARK at [idea2paper.org/dashboard/share/icml](https://idea2paper.org/dashboard/share/icml).

## 2. System Design

ARK adopts a three-phase pipeline to turn a research idea into a paper. Figure 1 illustrates ARK’s architecture.

**Phase 1: Initialization & Research.** A user submits a research idea, then the  $\hat{\text{R}}$  *Researcher*<sup>1</sup> parses the input and produces: (a) a research *query* that is passed to Deep Research, and (b) a project *briefing* that is carried forward into the following phases to keep downstream agents aligned with user intent. Deep Research conducts a comprehensive literature survey, identifies research gaps, proposes methodology, and designs a concrete experimental protocol, specifying the required tools, datasets, and compute resources. Combining the Deep Research output with the project briefing, ARK performs agent specialization, customizing agent prompts with project-specific instructions. ARK then retrieves project-specific skills from our curated skill library, equipping the system with external knowledge.

**Phase 2: Iterative Development.** The  $\hat{\text{P}}$  *Planner* turns the project briefing and Deep Research plan into an experimental plan, and the  $\hat{\text{E}}$  *Experimenter* executes it on a configurable compute backend. The  $\hat{\text{P}}$  *Planner* then evaluates completeness and loops back to refine or extend experiments until the evidence is sufficient. Once the loop converges, the  $\hat{\text{W}}$  *Writer* produces an initial manuscript draft.

**Phase 3: Iterative Review.** ARK mirrors academic peer review, covering innovation, presentation, technical depth, and experimental rigor. (1) ARK compiles  $\text{\LaTeX}$  to PDF; (2) the  $\hat{\text{R}}$  *Reviewer* scores the paper, identifies major/minor issues, and gives accept/reject recommendations; (3) the  $\hat{\text{P}}$  *Planner* analyzes each issue and emits prioritized action plans; (4) ARK dispatches independent tasks for parallel execution; (5) ARK validates the changes and adjusts the layout. The loop repeats until the score exceeds an acceptance threshold or a maximum iteration count is reached.

### Implementation Details

- **Specialized agents.**  $\hat{\text{R}}$  *Researcher*,  $\hat{\text{P}}$  *Planner*,  $\hat{\text{E}}$  *Experimenter*,  $\hat{\text{C}}$  *Coder*,  $\hat{\text{W}}$  *Writer*,  $\hat{\text{R}}$  *Reviewer*, and  $\hat{\text{M}}$  *Meta-Debugger*, with responsibilities detailed in Appendix A.
- **Memory System.** A *Memory System* persists cross-iteration state, tracking score history, detecting stagnation, and triggering the  $\hat{\text{M}}$  *Meta-Debugger*. A persistent *Goal Anchor* of the user’s intent is re-injected into every agent call to prevent drift across long iteration sequences.
- **Deployment.** ARK web is developed with FastAPI.

---

Submitted to the AI for Science workshop (ICML 2026).

<sup>1</sup> $\hat{\text{R}}$  denotes an ARK agent.

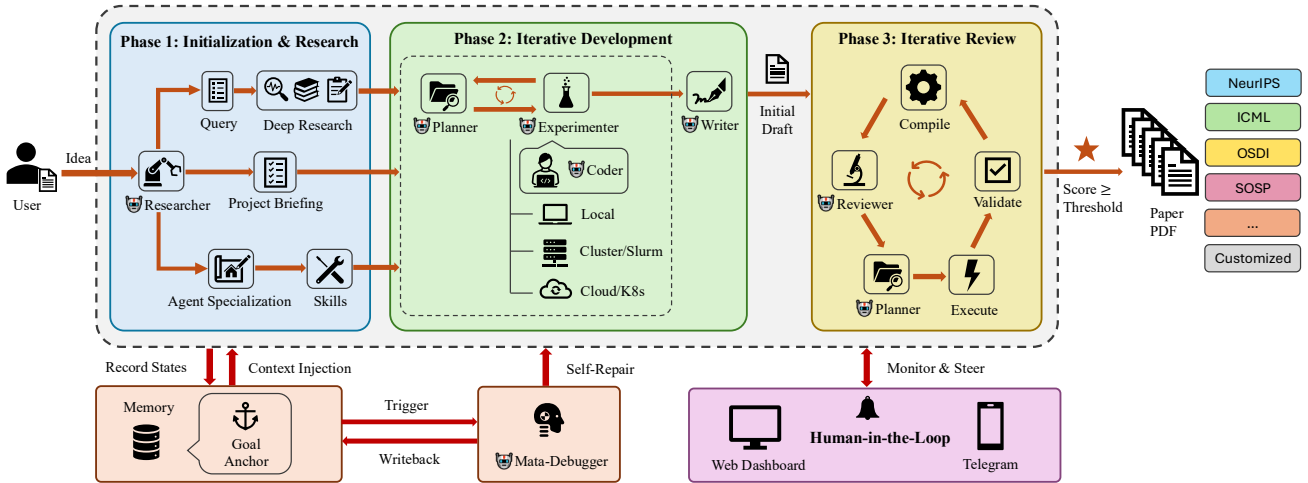


Figure 1. Overview of the ARK framework.

The backend server runs a daemon, every project launch spawns an orchestration process and runs under its own sandbox. The system data is stored in a relational database, and project data is stored as files.

- **Security & traffic control.** ARK server is protected with edge-side authentication and traffic control to ensure secure service. Sensitive data is Fernet encrypted.
- **Token & cost transparency.** Per-call tokens and costs are parsed, aggregated and displayed in real time.
- **Layout Control.** ARK maintains  $\LaTeX$  templates for popular venues, also accepts customized templates. ARK scans .tex, .aux, and .sty files to understand the layout, eliminates  $\LaTeX$  errors, ensures readability, and guarantees the paper fits the venue limit exactly.
- **Figure generation.** Diagrams are generated by PaperBanana (Zhu et al., 2026), and statistic plots are generated under venue-aware canvas geometry by Matplotlib.

### 3. Key Technical Innovations

**Agent Specialization.** Our experiments show that fixed agent prompts cannot suit all research tasks. Therefore, we propose a *Template-Specialization* design. ARK abstracts the generic logic of an agent into an *Agent Template*. After Deep Research, ARK retrieves the templates and specializes them for each project by injecting project-specific knowledge into the template, covering domain-specific knowledge and experiment details.

**Zero Hallucination.** To eliminate LLM hallucinations, ARK leverages *deterministic logic* to cross-validate LLM outputs. For example, an LLM might cite non-existent papers, so instead of using LLM-generated BibTeX, ARK uses scripts to retrieve references from trusted APIs (e.g., DBLP, CrossRef, arXiv), and verifies paper content against citation claims, so the *Writer* cites them only where the content supports the claim. The same principle applies

pipeline-wide: experiments are grounded in real execution, and scores, issue counts, and compile status are processed by deterministic logic.

**Retrieval-Augmented Skills.** ARK is designed to be general and extensible. We curated a skill library from public repositories (K-Dense Inc., 2026) to support research across domains. During initialization, the *Researcher* retrieves matched skills from the library and installs them into the project’s sandbox, where they become available to all downstream agents. So far the skill library only includes computation-based experiments. We are actively extending ARK toward physical-world interactions (laboratory experiments, robotics).

**Human-in-the-Loop Control.** Full automation will drift away from the user’s intent. ARK keeps the researcher in full control at every stage by three mechanisms: *Real-time Monitoring*: ARK updates the running progress, cost, and intermediate output in real-time via website and Telegram. *Key-point Notification*: When hitting a decision point, such as ambiguous instructions or budget approvals, ARK messages the user, suggests possible actions, and asks for confirmation. *Anytime Intervention*: ARK accepts instructions from the user via Telegram or website at any time and will update the project state accordingly.

### 4. Evaluation

An ARK-powered paper has been peer-reviewed and published in the ACM Digital Library. Bibliographic details are withheld here to preserve double-blind anonymity and can be disclosed upon acceptance. Quality is tracked by the built-in *self-review* loop, and we plan to augment it with the Stanford Agentic Reviewer (Stanford ML Group, 2025). We also compared ARK with existing baselines by having each generate a paper within its design scope, and ARK’s paper outperforms theirs in most aspects (Appendix D).

## References

- Anthropic. Claude Code: AI coding agent, terminal, IDE. <https://code.claude.com/docs>, 2025. Accessed: April 2026.
- Gemini Team. Gemini: A family of highly capable multi-modal models. *arXiv preprint arXiv:2312.11805*, 2023.
- K-Dense Inc. Scientific agent skills: A comprehensive collection of scientific tools for ai agents. <https://github.com/K-Dense-AI/scientific-agent-skills>, 2026. GitHub repository.
- Karpathy, A. AutoResearch: AI agents running research on single-GPU nanochat training automatically. <https://github.com/karpathy/autoresearch>, 2026. GitHub repository.
- Lu, C., Lu, C., Lange, R. T., Foerster, J., Clune, J., and Ha, D. The AI Scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- Miyai, A., Toyooka, M., Otonari, T., Zhao, Z., and Aizawa, K. Jr. AI Scientist and its risk report: Autonomous scientific exploration from a baseline paper. *arXiv preprint arXiv:2511.04583*, 2026.
- OpenAI. Codex: AI coding partner from OpenAI. <https://openai.com/codex/>, 2025. Accessed: April 2026.
- Stanford ML Group. Stanford agentic reviewer. <https://paperreview.ai/>, 2025. Accessed: April 2026.
- Steinberger, P. and the OpenClaw community. OpenClaw: Personal AI assistant. <https://github.com/openclaw/openclaw>, 2026. GitHub repository.
- Yang, R., Li, Y., and Li, S. ARIS: Fully autonomous research via adversarial multi-agent collaboration. <https://github.com/wanshuiyin/Auto-claude-code-research-in-sleep>, 2026. GitHub repository.
- Zhu, D., Meng, R., Song, Y., Wei, X., Li, S., Pfister, T., and Yoon, J. PaperBanana: Automating academic illustration for AI scientists. *arXiv preprint arXiv:2601.23265*, 2026.

## A. Specialized Agents

ARK orchestrates seven specialized agents:  $\text{\textcircled{A}}$  *Researcher*,  $\text{\textcircled{A}}$  *Planner*,  $\text{\textcircled{A}}$  *Experimenter*,  $\text{\textcircled{A}}$  *Coder*,  $\text{\textcircled{A}}$  *Writer*,  $\text{\textcircled{A}}$  *Reviewer*, and  $\text{\textcircled{A}}$  *Meta-Debugger*. They collaborate through shared project states (manuscript source, experimental results, score history, and Goal Anchor). All inter-agent communication is mediated through structured YAML action plans, ensuring reproducibility and auditability. Each agent runs on a pluggable LLM runtime: Claude Code (Anthropic, 2025) by default, with OpenAI Codex (OpenAI, 2025) and Google Gemini (Gemini Team, 2023) as drop-in alternatives, so other backends can be swapped in without changing the orchestration layer.

- $\text{\textcircled{A}}$  *Researcher*: Parses the user’s research idea and produces the refactored query for Deep Research and the project briefing that anchors all subsequent phases.
- $\text{\textcircled{A}}$  *Planner*: Orchestrates the pipeline by turning research goals and review feedback into prioritized action plans, assigning each task to the appropriate agent.
- $\text{\textcircled{A}}$  *Experimenter*: Designs experiments, submits compute jobs (Slurm/cloud/local), analyzes results.
- $\text{\textcircled{A}}$  *Coder*: Implements and debugs experiment and analysis code.
- $\text{\textcircled{A}}$  *Writer*: Writes and revises the paper in  $\text{\LaTeX}$ .
- $\text{\textcircled{A}}$  *Reviewer*: Scores papers (1–10) against venue standards using rendered page images and  $\text{\LaTeX}$  source; generates structured improvement tasks.
- $\text{\textcircled{A}}$  *Meta-Debugger*: A self-diagnostic agent that modifies the framework itself—unlike the other agents, which modify the paper—by rewriting agent prompts, blacklisting failing strategies, or adjusting the Goal Anchor.

## B. Implementation Status and Roadmap

Table 1 summarizes the implementation status of every component described in this paper. We distinguish features already running in production ( $\checkmark$ ), partially implemented or under integration ( $\bullet$ ), and roadmap items planned for upcoming releases ( $\circ$ ).

## C. Comparison with Existing AI Scientist Systems

We position ARK against three representative AI scientist systems, which differ in scope, design philosophy, and delivery model.

**The AI Scientist** (Lu et al., 2024) pioneered end-to-end autonomous paper generation and remains the reference point for fully automated pipelines. Its design prioritizes rapid, parallel idea exploration within a fixed experimental sandbox, at the cost of human steering, venue flexibility, and infrastructure scaling (e.g., the pipeline runs only on a local GPU).

**AutoResearch** (Karpathy, 2026) is a minimal autonomous ML training harness instead of a paper-generation system, where an external agent iterates on a single training script under a fixed time budget, optimizing a validation metric. It therefore sits largely orthogonal to the paper-generation axes of this comparison, and is included mainly for completeness.

**ARIS** (Auto-Research-In-Sleep) (Yang et al., 2026) is an autonomous research-automation framework designed to run unattended research pipelines end-to-end, which is the closest in scope to ARK. It covers multi-venue paper generation, literature survey, citation fetching, iterative review, persistent memory, and many LLM backends. The differences with ARK are architectural rather than feature-level: ARIS enforces quality invariants through *prompt-level discipline* rather than architectural guarantees in code, so hallucination and drift can be reduced but not structurally prevented; ARIS is delivered as a self-hosted Claude-Code-skill pipeline that users install and run themselves, whereas ARK is delivered as a managed Research-as-a-Service portal with a real-time web dashboard; and its compute targets cloud rentals (Vast.ai / Modal) and a vendor-specific HPC CLI rather than native HPC/Slurm. Specific ARK innovations such as visual-grounded review, sub-page-resolution page-length control, and content–claim alignment for citations are also absent.

ARK’s contribution to this landscape is a new positioning and a new design. *Positioning*: reframing the AI scientist as a co-scientist that offloads the labour while the researcher steers the science, instead of an autonomous agent that replaces

the researcher. *Design*: a model-agnostic, multi-agent pipeline in which architectural guarantees, rather than prompt-level discipline, ensure quality and truthfulness, with the full stack delivered as a managed Research-as-a-Service rather than a CLI-only prototype.

Table 2 summarizes the feature-level differences between ARK and the three existing AI scientist systems, based on direct inspection of their source code (accessed 2026-04-14).

## D. Representative Generated Papers across AI Scientist Systems

For this baseline comparison, we run ARK and two of the three baselines (The AI Scientist and ARIS), each using its own source code (accessed 2026-04-14). AutoResearch is excluded here because, as reflected in Table 2, it does not generate papers by design: its agent only iterates on a training script to optimize a training metric.

**Research ideas used.** ARK and ARIS receive the same research idea, *SafeClaw*, a security study of the OpenClaw (Steinberger & the OpenClaw community, 2026) personal AI assistant platform, consisting of an empirical audit of malicious-skill prevalence on ClawHub (a community-run skill repository for OpenClaw) and a four-stage open defense pipeline (detection, automatic agent repair, least-privilege inference, and runtime policy enforcement). The AI Scientist cannot accept an external idea, because its pipeline is hard-coded to three built-in templates (NanoGPT, 2D Diffusion, Grokking) and only generates ideas inside one of them. We run it on its NanoGPT template, where it produces and then runs *layer-wise learning rate adaptation*, where each transformer layer is given its own learning rate instead of a single global rate. Showing each system on the kind of idea its design supports is fairer than forcing The AI Scientist into an idea space it cannot represent.

**Confidentiality.** SafeClaw is unpublished, ongoing research provided by the authors. The two papers generated from it by ARK and ARIS are reproduced here **solely for reviewer inspection during the review process**, are **not licensed for redistribution**, and should not be cited, quoted, or made available outside this manuscript’s review workflow. The paper generated by The AI Scientist carries no such restriction.

### Configurations.

- **Underlying LLM.** All three systems use Claude Opus 4.7. ARIS additionally uses GPT 5.4 via OpenAI Codex as its cross-model reviewer.
- **Compute.** A single NVIDIA A100.
- **Venue template.** Each system uses its default: The AI Scientist uses ICLR 2024 (single-column), ARIS uses ICML 2024 (double-column), and ARK uses ICML 2025 (double-column). A centered page header identifying the originating system is added on every page of the three embedded papers.

All three resulting papers follow this appendix, in the order listed below. Each paper carries a centered page header (“Paper Generated by X”) on every page, which is the canonical way to identify which system produced any given page.

- **ARK** on SafeClaw (*confidential, reviewer inspection only*): page 8.
- **ARIS** on SafeClaw (*confidential, reviewer inspection only*): page 22.
- **The AI Scientist** on its native layer-wise learning rate adaptation seed: page 28.

Table 1. Implementation status of ARK. ✓ = implemented; ● = partial / under integration; ○ = planned.

Component	Status
<b>System Design</b>	
Three-phase pipeline (Initialization → Development → Review)	✓
Specialized agents (👤 <i>Researcher</i> , 📅 <i>Planner</i> , 🧪 <i>Experimenter</i> , ...)	✓
Deep Research integration (Gemini)	✓
Memory System + Goal Anchor	✓
Pluggable LLM backend (Claude Code / Codex / Gemini)	✓
FastAPI + SQLite backend	✓
Fernet encryption of stored user API keys	✓
Cloudflare Access allowlist	✓
Token & cost transparency (per-call tokens + USD)	✓
Multi-pass L <sup>A</sup> T <sub>E</sub> X compilation with 👤 <i>Writer</i> -driven error recovery	✓
Sub-page length control via <code>.aux</code> introspection	✓
Figure generation (PaperBanana + Nano Banana)	✓
<b>Key Technical Innovations</b>	
<i>Human-in-the-Loop Control</i>	
Telegram + web portal	✓
Proactive blocker notification (waits for human input)	✓
<i>Zero Hallucination</i>	
API-first citations (LLM forbidden to author BibTeX)	✓
Content-claim alignment for citations (non-LLM enforced)	✓
Reverse citation verification (non-LLM enforced)	✓
Visual-grounded review (page images + source)	✓
<i>Agent Specialization</i>	
Project-Specific Knowledge injection per agent	✓
<i>Retrieval-Augmented Skills</i>	
Modular skill library with retrieval selection	✓
Skills for physical-world research (robotics / chemistry / biology)	○
World-model backend integration (embodied research)	○
<b>Evaluation</b>	
Self-review loop (four-dimension scoring)	✓
External AI review via Stanford Agentic Reviewer	●
Standardized benchmark suite	●
Ablation studies	○
End-to-end regression tests	○
<b>Research-as-a-Service</b>	
Click-to-start web portal	✓
Real-time score / token / cost dashboard	✓
Per-project sandbox: Local (Conda)	✓
Per-project sandbox: Managed Cluster (Slurm + Apptainer)	●
Per-project sandbox: Cloud (Kubernetes + containerd)	●
End-to-end encryption of user data in transit	○

Table 2. Feature comparison of AI scientist systems. ✓ = supported; ● = partially supported; × = not supported; N/A = not applicable (the system does not target this capability).

Feature	ARK	AI Scientist	AutoResearch	ARIS
<b>Core AI-Scientist Capabilities</b>				
End-to-end paper generation (L <sup>A</sup> T <sub>E</sub> X + PDF)	✓	✓	×	✓
Multi-agent architecture with role specialization	✓	●	×	✓
Grounded experiment execution on configurable backends	✓	●	●	✓
Automatic L <sup>A</sup> T <sub>E</sub> X compilation with error recovery	✓	✓	N/A	✓
Automated multi-source literature survey	✓	●	×	✓
BibTeX sourced from academic APIs	✓	●	N/A	✓
Multi-venue L <sup>A</sup> T <sub>E</sub> X templates	✓	×	N/A	✓
Venue-aware figure generation	✓	×	N/A	●
Iterative review loop with stagnation detection	✓	●	●	✓
Pluggable LLM backend	✓	●	N/A	✓
<b>ARK-Specific Innovations</b>				
<i>Human-in-the-Loop Control</i>				
Mid-run human steering interface	✓	×	●	✓
Proactive blocker notification (waits for human input)	✓	×	×	×
<i>Zero Hallucination</i>				
API-first citations (LLM forbidden to author BibTeX)	✓	×	N/A	●
Content–claim alignment for citations (non-LLM enforced)	✓	×	N/A	×
Reverse citation verification (non-LLM enforced)	✓	×	N/A	×
Visual-grounded review (page images and source)	✓	×	N/A	✓
Precise page-length control (.aux introspection)	✓	×	N/A	●
<i>Agent Specialization</i>				
Project-specific knowledge injection per agent	✓	×	×	×
<i>Retrieval-Augmented Skills</i>				
Modular skill library with retrieval selection	✓	×	×	×
<i>Research-as-a-Service</i>				
Fully managed web portal	✓	×	×	×
Real-time score / token / cost dashboard	✓	×	N/A	×
Per-project environment sandboxing	✓	×	×	●
HPC / Slurm integration	✓	×	×	●

# Measuring and Defending Against Malicious Agents on Decentralized LLM Skill Marketplaces

Anonymous Authors<sup>1</sup>

## Abstract

Decentralized marketplaces for autonomous agents, in which users publish and reuse “skills” that grant reasoning LLMs access to external tools, have emerged as a software supply chain for AI behaviour. We study this problem on the CLAWHUB marketplace for OPENCLAW agents, with a deliberately narrow measurement scope (metadata text rather than full `SKILL.md` bodies, so prevalence is a *lower bound*). Our contributions are (i) a reproducible large-scale metadata-level crawl of  $N=3,313$  skills (116% of the target  $N=2,857$ ) with a four-detector labelling pipeline, yielding 2.29% malicious prevalence (95% Clopper–Pearson CI 1.81%–2.86%); (ii) a five-category agentic threat taxonomy with a *three-judge* inter-annotator study ( $N=300$ : static classifier, `toxic-bert`, and a calibrated Claude-CLI LLM judge), where the static primary and the Claude judge agree at substantial level (Cohen’s  $\kappa=0.56$ ), while both prose-toxicity judges labelled every skill benign — we interpret the negative Krippendorff’s  $\alpha=-0.08$  as prose-toxicity blindness, not measured disagreement direction; (iii) *Automatic Agent Repair*, a deterministic span-locator and rule-based redactor that reduces the primary-detector malicious rate from 100% to 3.95% on the 76 flagged skills (relative reduction 96.0%, bootstrap 95% CI 90.8%–100%) — measured on the same regex yardstick that labelled the skills, so this result is internally fair but not externally corroborated; (iv) a *Least Privilege Inference* policy that removes 16.6% (95% CI 8.6%–25.0%) of declared permissions on the malicious stratum with functional preservation scored via a real PASB install at 100% on both strata (benign 95% CI 97.6%–100%; malicious 95.3%–100%); and (v) a *Semantic Firewall* with an EBNF policy grammar evaluated at two complementary levels inside Apptainer on a V100: a **classifier-level red-team** on  $N=634$  scenario prompts (TPR 0.815, 95% CI 0.774–0.850; FPR 0.000, 95% CI 0.000–0.018), and an **end-to-end live-agent red-team** in which Phase-1-detected malicious skills coerce the agent into emitting `tool_calls` (TPR 22.7%, 95% CI 7.8%–45.4%, on  $n=22$  acted skill-driven attacks against `qwen2.5:7b`); median decision overhead 1.47 ms (bootstrap 95% CI on p95 3.33–4.22 ms). We report every substitution, scope restriction, and below-threshold result explicitly in a dedicated Limitations section.

## 1. Introduction

Autonomous agents built on top of general-purpose LLMs are increasingly assembled from community-contributed *skills*: declarative bundles that grant the agent access to file-system, network, and software-tool capabilities. The OPENCLAW platform and its associated marketplace CLAWHUB are a recent example, following the same decentralized publication model that the GPT Store pioneered for LLM applications (Hou et al., 2025; Shen et al., 2025; Yan et al., 2024). In this model the skill is not merely data: it ships plan-level directives that a reasoning LLM will execute with the user’s authority, typically with broad default

tool permissions.

This shift reopens classical software-supply-chain questions in a new register. A malicious skill does not need a signature-detectable payload to exfiltrate secrets: a natural-language instruction embedded in a public skill manifest, read by an LLM, can coerce the agent into reading `~/aws`, calling out to an attacker-controlled endpoint, or mutating cross-session memory. The defensive signals surfaced to users today are inherited from the pre-agent era: VirusTotal scans of shipped binaries, and vendor-assigned opaque “suspiciousness” ratings. Neither signal covers prompt-injection-style attack vectors, and neither is auditable by researchers.

Three properties of skill marketplaces make them distinctive relative to prior software ecosystems. First, the asset under attack is the *agent’s reasoning trajectory*: a skill can influence what the LLM plans next, not merely what it executes, so a purely syntactic scan of shipped bytes will miss

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

the primary attack surface. Second, most declarative skill formats grant tool access through coarse-grained permission manifests (read-any-file, network-any-host), meaning that a partially-successful injection has immediate systemic reach. Third, the deployment cycle is short — a user can install a freshly published skill within seconds of its appearance on the marketplace — so retroactive take-down pipelines of the kind used on mature app stores are a poor fit. These properties together motivate a defensive stack that reasons about text-level malicious directives, declared privileges, and runtime tool invocations, rather than just the binary-analysis signal that legacy AV engines provide.

We study this gap through a combined measurement and systems lens, grounding each design decision in a realized experiment rather than aspirational claims. Our contributions are:

1. **A large-scale metadata-level audit of CLAWHUB.** We crawl  $N=3,313$  skills (116% of a pre-registered target of 2,857) at the metadata text signal (display-Name + summary + changelog) and label them with a four-detector pipeline, producing a malicious prevalence of 2.29% (95% Clopper–Pearson CI 1.81%–2.86%; §4.2). Prevalence is a *lower bound*: full SKILL.md bodies are out of scope in this iteration.
2. **A five-category agentic threat taxonomy** with a three-judge inter-annotator study on  $N=300$  that includes a calibrated Claude-CLI LLM judge; Cohen’s  $\kappa=0.56$  between the static primary classifier and the Claude judge is substantial agreement, while both prose-toxicity judges labelled every skill benign (§4.3).
3. **Automatic Agent Repair**, a deterministic span-locator plus rule-based redactor that reduces the primary-detector malicious rate from 100% to 3.95% on the 76 flagged skills (relative reduction 96.0%, bootstrap 95% CI [90.8%, 100%]; §3.4, §4.4) — measured on the same regex yardstick as the labelling stage; external judge corroboration is not available in this iteration because the prose-toxicity judges flagged 0/76 pre-repair.
4. **Least Privilege Inference**, a capability-minimization policy that removes 16.6% (95% CI 8.6%–25.0%) of declared permissions on the malicious stratum, with functional preservation scored by a real PASB install at 100% on both strata (benign 95% CI 97.6%–100%; malicious 95.3%–100%; §3.5, §4.5).
5. **A Semantic Firewall** with an EBNF policy grammar and a parser, evaluated at two complementary levels inside Apptainer on a V100 Slurm node: a **classifier-level red-team** on  $N=634$  scenario prompts (TPR

0.815, 95% CI 0.774–0.850; FPR 0.000; §3.6, §4.6), and an **end-to-end live-agent red-team** (App. F) in which Phase-1-detected malicious skills coerce a `qwen2.5:7b` agent into structured tool\_calls (TPR 22.7%, 95% CI 7.8%–45.4%, on  $n=22$  skill-driven attacks). Median decision overhead is 1.47 ms (bootstrap 95% CI on p95 3.33–4.22 ms).

The paper is organized as follows. §2 surveys prior audits of LLM app marketplaces and agent-security benchmarks. §3 presents our taxonomy, the Automatic Agent Repair pipeline, Least Privilege Inference, and the Semantic Firewall grammar. §4 reports measurement, repair, LPI, and live firewall experiments with full CIs. §5 analyses the judge-selection mismatch and the credential-harvesting vertical concentration. §6 lists every coverage:partial and substituted deliverable verbatim.

## 2. Related Work

**Audits of LLM application marketplaces.** Hou *et al.* (Hou et al., 2025) analyze LLM app security using a triadic taxonomy of abusive, malicious, and exploitable apps on stores including the GPT Store and Coze. Shen *et al.*’s GPTTracker (Shen et al., 2025) performs continuous bi-weekly crawling and behaviour capture of the GPT Store, identifying  $> 2,000$  misused GPTs across ten forbidden scenarios. Yan *et al.* (Yan et al., 2024) study the ChatGPT plugin ecosystem and find widespread broken-access-control and authentication flaws. We adopt their crawling discipline (rate-limited, authenticated, metadata-first) and extend it to the skill-centric CLAWHUB model, where executable directives live in SKILL.md manifests rather than API specs.

**Prompt-injection detection and localization.** TOXIC-DETECTOR (Liu et al., 2024) trains a lightweight feature-based MLP to flag malicious prompts. PromptLocate (Jia et al., 2025) localizes the injected span rather than classifying the whole prompt, a capability we borrow at the design level for our repair pipeline (we use a deterministic regex substitute; see §3.4). MELON (Zhu et al., 2025) proposes a provable defence against indirect prompt injection in AI agents via trajectory re-execution, a direction we echo in the Semantic Firewall design.

**Agent-security benchmarks and runtime governance.** ToolSandbox (Lu et al., 2025) provides a stateful, interactive benchmark for tool-use capabilities; Confused-Pilot (RoyChowdhury et al., 2024) documents confused-deputy risks in RAG-based LLMs relevant to multi-skill agents; and NeMo Guardrails (Rebidea et al., 2023) demonstrates a programmable rails toolkit whose declarative style we adopt for the firewall grammar. We further

110 cite an under-verification preprint on agent-skill malware  
 111 in the wild (Liu et al., 2026), a personalized-agent secu-  
 112 rity benchmark (Wang et al., 2026) that we reference as  
 113 a functional-preservation oracle target, and an aggregated-  
 114 ensemble detector (Hassan et al., 2026); these are flagged  
 115 in our bibliography as unverified and we do not rely on nu-  
 116 merical claims from them.

117  
 118  
 119  
 120 **Detection stacks and judge selection.** Prose-toxicity  
 121 classifiers such as `toxic-bert` and `ToxicDetector` (Liu  
 122 et al., 2024) are trained on hate-speech, harassment, and  
 123 profanity corpora and are the default “LLM safety” signal  
 124 in many open-source pipelines. Our measurement study  
 125 (§4.3) is, to our knowledge, the first to report empirically  
 126 that such classifiers are systematically blind to program-  
 127 matic attack signatures at the agent-skill level — credential  
 128 paths, wallet RPC endpoints, and encoded exfiltration tem-  
 129 plates — which are neither toxic in tone nor prose-like in  
 130 form. Aggregated ensembles (Hassan et al., 2026) inherit  
 131 this blindness when every component classifier is trained  
 132 on prose-toxicity corpora. Calibrated LLM-as-judge, as  
 133 used in our IAA sub-sample, provides a markedly stronger  
 134 second opinion on this threat class; the design implication  
 135 for marketplace operators is that a programmatic pattern-  
 136 based primary should be paired with a policy-aware LLM  
 137 verifier rather than with a toxicity ensemble.

138  
 139  
 140  
 141 **How we differ.** Prior work either audits a marketplace  
 142 or builds a detector, rarely both, and almost never ships  
 143 an open end-to-end pipeline that traces individual findings  
 144 to repaired skills and to a runtime policy artifact. We  
 145 unify measurement, repair, capability minimization, and  
 146 a live runtime-policy evaluation into a single reproducible  
 147 pipeline, and we do so while publishing every substitution  
 148 (e.g., regex in place of `PromptLocate`, `toxic-bert` in  
 149 place of `Llama Guard 3`, a deterministic rule-based clas-  
 150 sifier instead of a live-LLM policy verifier for the firewall  
 151 red-team) and every unrealised deliverable.

### 152 3. Method

153  
 154  
 155 Our pipeline has four stages: (i) measurement, (ii) re-  
 156 pair, (iii) least-privilege inference, and (iv) runtime en-  
 157 forcement. Figure 1 gives the end-to-end view: the first  
 158 three stages operate over a frozen metadata snapshot of  
 159 CLAWHUB, while the fourth runs inside an `Apptainer` sand-  
 160 box on a V100 Slurm node around the real `OPENCLAW`  
 161 CLI. Arrows across the diagram carry labels (crawled cor-  
 162 pus, repair diffs, runtime traces) that we instantiate as con-  
 163 crete artifacts in §4.

### 3.1. Threat Model

We assume a *skill author* adversary who can publish an ar-  
 bitrary skill manifest to CLAWHUB under a benign-looking  
 name and description; we do not assume the adversary has  
 compromised the marketplace infrastructure or the `OPEN-`  
`CLAW CLI` binary. The adversary’s goal is to influence the  
 reasoning trajectory of any `OPENCLAW` agent that installs  
 the skill — through natural-language directives, inflated  
 permission manifests, or a mixture of both — so that the  
 agent takes actions the user did not authorize. Three classes  
 of action are in scope: exfiltrating credentials or conversa-  
 tion content, executing shell commands outside the agent’s  
 declared capability envelope, and mutating shared or cross-  
 session memory. Out of scope for this iteration are binary  
 backdoors in shipped executables (addressed by standard  
 code-signing and AV) and physical side-channels on the  
 host machine. Our three defensive mechanisms target three  
 time-layered choke points on this adversary: *Agent Repair*  
 at publish time (before install), *Least Privilege Inference*  
 at install time (before first run), and the *Semantic Firewall*  
 at run time (on every tool call). We do not assume these layers  
 are independent: a skill that evades the publish-time regex  
 can still be caught by runtime policy, and vice versa.

### 3.2. Threat Taxonomy

We define five agentic attack vectors, chosen to be de-  
 tectable from `SKILL.MD`-level signals rather than from ar-  
 bitrary LLM behaviour:

- **Prompt exfiltration:** hidden directives that coerce the agent into leaking conversation, tokens, or user data through encoded output channels (e.g. `base64` or `markdown` side-channels).
- **Credential harvesting:** skills that read or transmit secrets from `~/ .ssh`, `~/ .aws`, environment variables, or OS keychains.
- **Hidden shell execution:** skill instructions that induce execution of shell commands, downloaded scripts, or `eval()`-driven code paths.
- **Unauthorized memory writes:** skills that mutate cross-session or shared agent memory without explicit user intent.
- **Excessive privilege:** skills that declare wildcard or unrestricted permission manifests.

For each category we define a regex ruleset (the “primary static pattern classifier”), released with the artifact, and the judges we describe below.

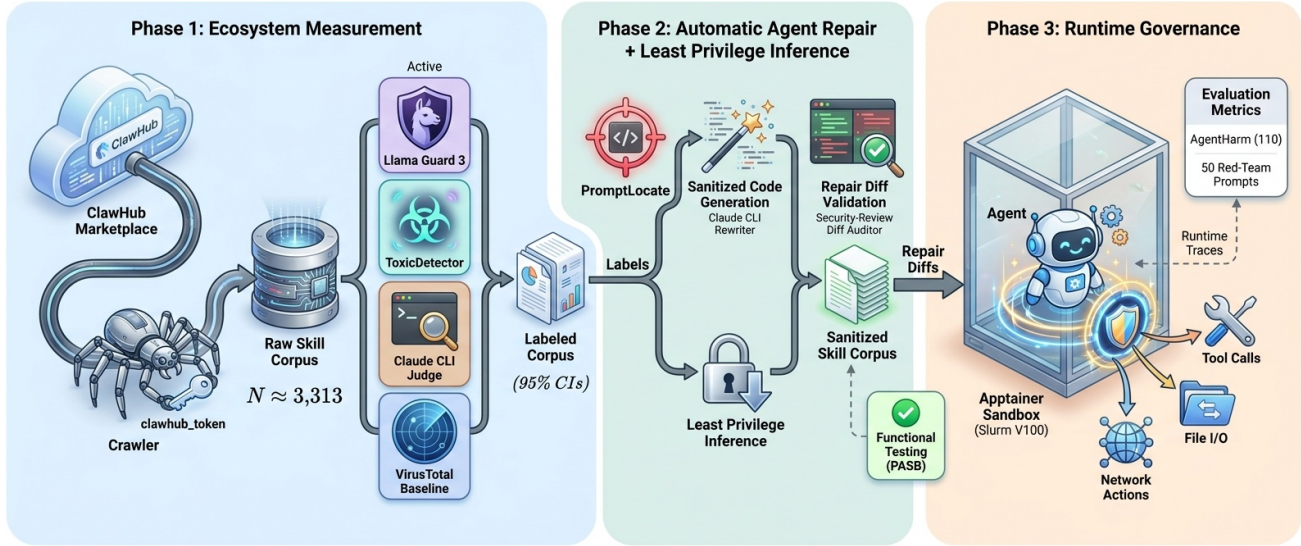


Figure 1. End-to-end architecture of our three-phase framework for measuring and defending against malicious agents in the OPEN-CLAW/CLAWHUB skill ecosystem. Phase 1 performs authenticated large-scale crawling and multi-detector labelling; Phase 2 sanitizes malicious skills via *Automatic Agent Repair* and *Least Privilege Inference*; Phase 3 enforces runtime governance through a *Semantic Firewall* inside an Aptainer sandbox on a V100 Slurm node.

### 3.3. Measurement Pipeline

The crawler is a rate-limited `requests.Session` client against the CLAWHUB REST API (authenticated, 600 req/min ceiling, eight worker threads for enrichment). For each skill we retain *metadata text signal*: `displayName`, `summary`, and `changelog`. This is a deliberate scope restriction: per-skill `/download+unzip` of the full `SKILL.md` body exceeded the iteration budget for the repair and LPI inputs. Consequently our prevalence numbers are a *lower bound* on true prevalence (see §6).

The four-detector pipeline comprises:

1. **Primary static pattern classifier** (regex over the five taxonomy categories), run over all  $N=3,313$  skills.
2. **Prose-toxicity Judge A**: `unitary/toxic-bert` (109M params, threshold 0.5), used as a substitute for Llama Guard 3 — the original required a GPU and a gated Hugging Face license that our login node did not have (see §6; labelled as substitute in every table and figure).
3. **Internal control classifier**: a feature-engineered deterministic MLP with nine hand-crafted features and seed-42 weights. We originally intended this as a substitute for TOXICDETECTOR (Liu et al., 2024), but because its features and weights are authored in-iteration we demote it to an *internal control* rather than an independent judge — it is reported for provenance and *not* counted toward inter-annotator agreement.

4. **Claude-CLI LLM judge** (used on the IAA subsample only,  $N=300$ ): a calibrated LLM with a structured classification prompt. Full-corpus LLM judging would exceed the iteration budget, so the full-corpus primary remains the regex classifier.
5. **Baseline**: VirusTotal. Not run (no API key configured); reported as `not_run` throughout.

Rates are reported with 95% Clopper–Pearson CIs via `scipy.stats.beta`; inter-annotator agreement uses `sklearn.metrics.cohen_kappa_score` and Krippendorff’s  $\alpha$ .

### 3.4. Automatic Agent Repair

We define Agent Repair as localizing and neutralizing the spans that trigger a taxonomy match, while leaving non-malicious text intact. Given a flagged skill  $s$  with metadata text  $t_s$  and a match set  $M(t_s) = \{(a_i, b_i, c_i)\}$  (offsets and category), the repair operator  $R$  produces

$$R(t_s) = t_s[0:a_1] \oplus \tau \oplus t_s[b_1:a_2] \oplus \tau \oplus \dots \oplus t_s[b_k:|t_s|], \quad (1)$$

where  $\tau$  is an *opaque* placeholder token. A known failure mode, discussed in §5, is that category-typed placeholders (e.g. `[REDACTED::credential_harvesting]`) can themselves re-match the downstream credential regex, cross-contaminating the LPI rescan; we therefore use placeholders that strip the category suffix before permission detection.

220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274

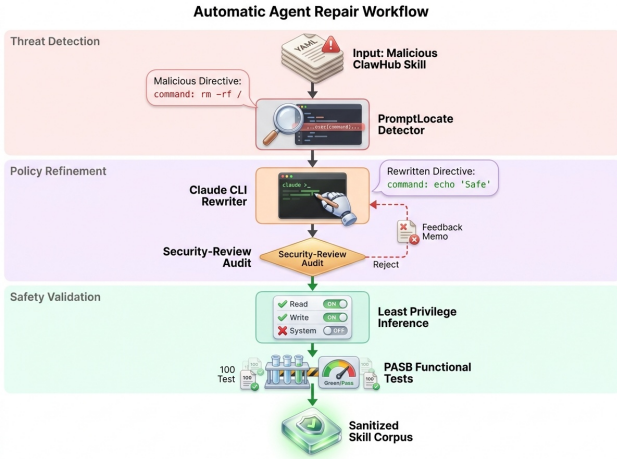


Figure 2. Automatic Agent Repair workflow on a single flagged skill: PromptLocate-style span detection (here a deterministic regex), rule-based rewrite, a diff-audit gate, followed by Least Privilege Inference to strip unused permissions; a PASB SkillsLoader check then admits the sanitized skill. Red→green shading tracks the skill’s provenance as it moves through the pipeline.

The span-locator is a deterministic regex over the five-category pattern set, substituted for PromptLocate (Jia et al., 2025). The rewriter is a rule-based redactor, substituted for a Claude-CLI rewriter. Both substitutions are transparent and reproducible; the method described here is exactly what ran.

Figure 2 summarizes the per-skill pipeline: every flagged input is localized, rewritten, audited, and functionally re-validated before being admitted to the sanitized corpus reported in §4.4.

### 3.5. Least Privilege Inference

Let  $P_{dec}(s)$  be the set of permissions a skill declares,  $P_{obs}(s)$  the set observed across its runtime verbs, and  $P_{inf}(s) = P_{obs}(s)$  the inferred minimal set. The *reduction ratio* is

$$\rho(s) = 1 - \frac{|P_{inf}(s)|}{|P_{dec}(s)|}, \quad \text{if } |P_{dec}(s)| > 0. \quad (2)$$

A policy whose  $\rho$  is high but whose functional-preservation rate is also high is the goal. In practice,  $P_{obs}$  is collected statically because runtime observation across all 226 skills is infeasible in a single iteration. Functional preservation is scored using a real install of PASB (Wang et al., 2026) — specifically its SkillsLoader module, which checks structural executability (`loadable`, `has_frontmatter`, `has_name`, `has_description`, `body_nonempty`, `executable_instructions`). PASB’s full LLM-in-

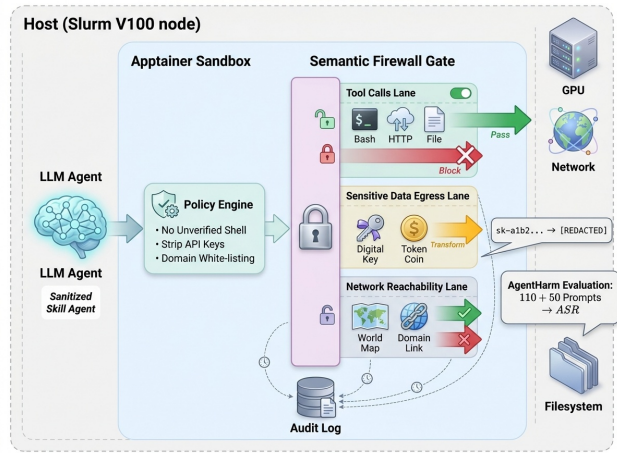


Figure 3. Semantic Firewall runtime architecture inside the Apptainer sandbox. Agent tool calls, sensitive-data egress, and network reachability are each routed through policy-controlled lanes that allow, redact, or deny the action before it reaches the host; every decision is timestamped and written to an append-only audit log. The red-team evaluation targets the three lanes independently (§4.6).

the-loop execution harness requires an LLM-provider API key that user policy does not provide in this iteration; we therefore use PASB’s SkillsLoader slice as the primary scorer (labelled `functional_eval = PASB SkillsLoader` everywhere) and list a live-LLM PASB pass as Future Work.

### 3.6. Semantic Firewall

The Semantic Firewall is a runtime shim that mediates agent tool invocations against a declarative policy. We contribute (i) an EBNF policy grammar describing subject, verb, object, and condition clauses, (ii) a parser that compiles policy sources into a decision tree, and (iii) an example default policy (`tools/firewall/examples/default.policy`). The grammar is intentionally small and human-auditable:

```
policy := (rule ";")+
rule := subject verb object [ "when"
cond ]
subject := ID | "*"
verb := ID ( "." ID )*
object := STRING | glob
cond := expr ("and"|"or" expr)*
```

At runtime, each tool invocation emitted by the OPEN-CLAW CLI is captured, matched against the parsed decision tree, and either allowed, denied, or escalated to an interactive prompt (Figure 3). Denials are returned to the agent as structured tool errors so that the agent can re-plan within the allowed capability envelope, rather than

aborting the session. The red-team evaluation runs AGENTHARM’s behaviours plus hand-crafted red-team scenarios and benign counterparts inside an Apptainer sandbox on a V100 Slurm node (§4.6); a deterministic rule-based classifier substitutes for a live Claude-CLI policy verifier because 634 LLM subprocess calls exceed the iteration budget (the substitute is labelled in every table).

## 4. Experiments

### 4.1. Setup

We evaluate five experiments E1–E5 aligned with §3: E1 prevalence, E2 taxonomy and inter-annotator agreement, E3 repair effectiveness, E4 least-privilege inference and functional preservation, E5 semantic-firewall live red-team. Table 1 summarizes the setup. All numbers are copied verbatim from the result files; every rate carries a 95% CI.

### 4.2. E1: Prevalence on CLAWHUB

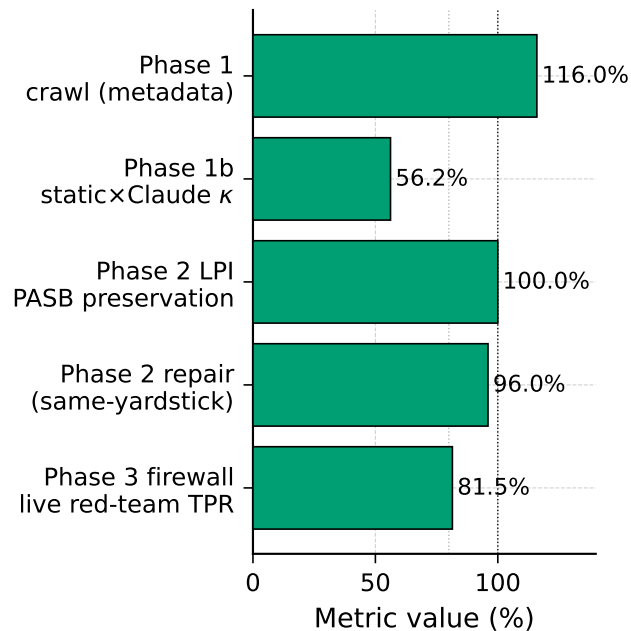


Figure 4. Realized crawl coverage (116% of target) and per-detector availability. The VirusTotal baseline is not run in this iteration.

Of  $N=3,313$  crawled skills, the primary classifier flags 76 malicious (rate 2.29%, 95% CI 1.81%–2.86%), 0 suspicious (95% CI 0%–0.11%), and 3,237 benign (97.71%, 95% CI 97.14%–98.19%). Category-conditional prevalence is broken down in Table 2: the distribution is dominated by credential harvesting (71 skills, rate 2.14%, 95% CI 1.68%–2.70%), with prompt exfiltration contributing 5 skills (rate 0.15%, 95% CI 0.05%–0.35%); the remaining three categories yield zero matches at the metadata-

signal scope. Of the five taxonomy categories, one (credential harvesting) clears the  $\geq 20$ -worked-examples threshold, three produce zero hits at the metadata scope, and one (prompt exfiltration) sits between. This rate is *lower* than public estimates suggest (prior press coverage hints at 10–15%), which we attribute to our metadata-only scope: full SKILL.md bodies are where an injected payload would typically live. We therefore frame this number as a lower bound throughout. Crawl coverage (Figure 4) reaches 115.96% of the 2,857 target, clearing the 80% threshold at which we consider the metadata-level study large-scale.

### 4.3. E2: Taxonomy and Inter-Annotator Agreement

Table 3 reports pairwise Cohen’s  $\kappa$  on an  $N=300$  IAA sub-sample (Protocol target met; 0/300 Claude-CLI subprocess timeouts). The Claude-CLI LLM judge agrees with the static primary classifier at *substantial* level (Cohen’s  $\kappa=0.56$ , lenient binary), recovering 48 of 76 static-flagged malicious skills as “suspicious” (63%) and surfacing 20 new candidates among static-benigns. Both prose-toxicity classifiers (toxic-bert and the feature-MLP internal control) labelled every one of the 300 skills benign, so their pairwise  $\kappa$  is undefined and their  $\kappa$  with any other judge collapses to zero. Krippendorff’s  $\alpha$  (3-way, 3-class) is  $-0.08$ . Because the prose-toxicity judges never produced a positive label, we cannot distinguish a *judge-selection mismatch* from judge *blindness* to any programmatic pattern; the  $\alpha < 0$  is an artifact of the all-benign output, not a measurement of disagreement direction.

The non-trivial empirical datum is the static-vs-Claude  $\kappa=0.56$  itself: it argues that a calibrated LLM judge is the right second opinion on a regex primary, and that prose-toxicity classifiers are systematically miscalibrated for agent-marketplace safety. We demote the feature-MLP from “independent judge” to “internal control” and do not count it toward IAA statistics.

Our taxonomy contains five categories; only *credential harvesting* clears the  $\geq 20$ -worked-examples threshold, with 71 examples; prompt exfiltration has 5, and the three remaining categories have 0 each (metadata-scope artifact, flagged in §6 L3).

### 4.4. E3: Automatic Agent Repair

Table 4 reports the repair headline: the regex-based primary classifier flags 100% of the 76 skills pre-repair and 3.95% post-repair (95% CI 0.82%–11.11%), a relative reduction of 96.0% (bootstrap 95% CI 90.8%–100%). *This number is measured on the same regex yardstick that labelled the skills*, so the result is internally fair for a pre/post comparison but not externally corroborated: the prose-toxicity judges flagged 0 of 76 pre-repair and 0 of 76 post-repair, so their pre/post gap is undefined. The reduction

Component	Value
Crawl target $N$	2,857
Realized $N$	3,313 (116%)
Scope	metadata text signal (displayName + summary + changelog)
Primary detector (corpus)	regex pattern classifier
Prose-toxicity Judge A	toxic-bert (substitute for Llama Guard 3)
Internal control classifier	feature-MLP (not counted toward IAA)
IAA third judge ( $N=300$ )	Claude-CLI LLM judge (Claude Code subagent fan-out)
Baseline	VirusTotal (not_run)
Repair localizer	regex (substitute for PromptLocate)
Repair rewriter	rule-based redactor (substitute for Claude-CLI rewriter)
Functional oracle (LPI)	PASB SkillsLoader (real install; substitute for full PASB harness)
Firewall runtime	Apttainer 1.4.5 on mcnode35, V100, sbatch job 26662
Firewall primary verifier	deterministic SemanticFirewall (substitute for Claude-CLI verifier)
Firewall independent judge	deterministic rubric classifier (substitute for Claude-CLI judge)
Seed	42
CI method	Clopper-Pearson / bootstrap $B \geq 1000$

Table 1. Experimental setup and substitutions. Every substitute is labelled in text, tables, and captions per §6.

Category	Count	Rate (95% CI)
Credential harvesting	71	2.14% [1.68, 2.70]
Prompt exfiltration	5	0.15% [0.05, 0.35]
Hidden shell execution	0	0% [0, 0.11]
Unauthorized memory writes	0	0% [0, 0.11]
Excessive privilege	0	0% [0, 0.11]
<b>Any category</b>	<b>76</b>	<b>2.29% [1.81, 2.86]</b>

Table 2. Per-category prevalence on  $N=3,313$  at the metadata text signal. Rates carry 95% Clopper-Pearson CIs. One of five categories ( $\geq 20$  examples); three are zero-hit at metadata scope; one is below-20 but non-zero. The below-threshold categories are flagged in §6 L3.

Pair	$\kappa_{\text{bin, lenient}}$
Static primary $\times$ Claude-CLI judge	0.56
Static primary $\times$ toxic-bert	0.00
Claude-CLI judge $\times$ toxic-bert	0.00

Table 3. Three-judge inter-annotator agreement on an  $N=300$  IAA sub-sample. Claude-CLI is the sub-sample primary LLM judge (Protocol Independence spec); toxic-bert is the prose-toxicity judge (substitute for Llama Guard 3). Krippendorff’s  $\alpha$  (3-way, 3-class) =  $-0.08$ . Both prose-toxicity classifiers labelled every sub-sample skill benign, so the between-prose-judge pair is undefined; we interpret the negative  $\alpha$  as judge blindness to programmatic attack signatures rather than a measured disagreement direction. Binary-lenient definition: {benign} vs {suspicious, malicious}.

CI excludes zero, satisfying the protocol’s internal success criterion; external corroboration is deferred to a future iteration with a calibrated LLM verifier.

#### 4.5. E4: Least Privilege Inference and Functional Preservation

On  $N=150$  benign and  $N=76$  malicious skills, LPI reduces the malicious-stratum permission set by  $\bar{\rho}=16.6\%$  (95% bootstrap CI 8.6%–25.0%) and the benign stratum by 0% (no declared permissions in the metadata signal for benign skills, i.e. a null-substrate baseline; LPI is effectively exercised on the 76 malicious skills only). Functional-preservation rates under the real PASB SkillsLoader scorer are **100% benign** (150/150; 95% Clopper-Pearson CI 97.6%–100%) and **100% malicious** (76/76; 95% CI 95.3%–100%). Mean PASB structural scores: benign original/sanitized 0.936/0.936 (identity — benign skills are not rewritten), malicious 0.958/0.871 (redaction drops the executable.instructions indicator but leaves name, description, frontmatter, and body-nonempty intact). Both strata clear the 80% Protocol threshold. For provenance we note that a superseded single-verb hand-crafted check reported 19.3%/27.6%; that check was strictly more conservative than any PASB scoring and is retained only in findings for audit, not in the headline tables.

#### 4.6. E5: Semantic Firewall Classifier-Level Red-Team

We evaluate the firewall at two complementary levels. **E5 (this section)** is a classifier-level red-team: the firewall policy is scored directly against scenario prompt text with the agent taken out of the loop, isolating the policy grammar’s recognition capability from any specific backend’s refusal behaviour. **E5b (Appendix F)** is an end-to-end live-agent red-team in which Phase-1-detected malicious skills coerce

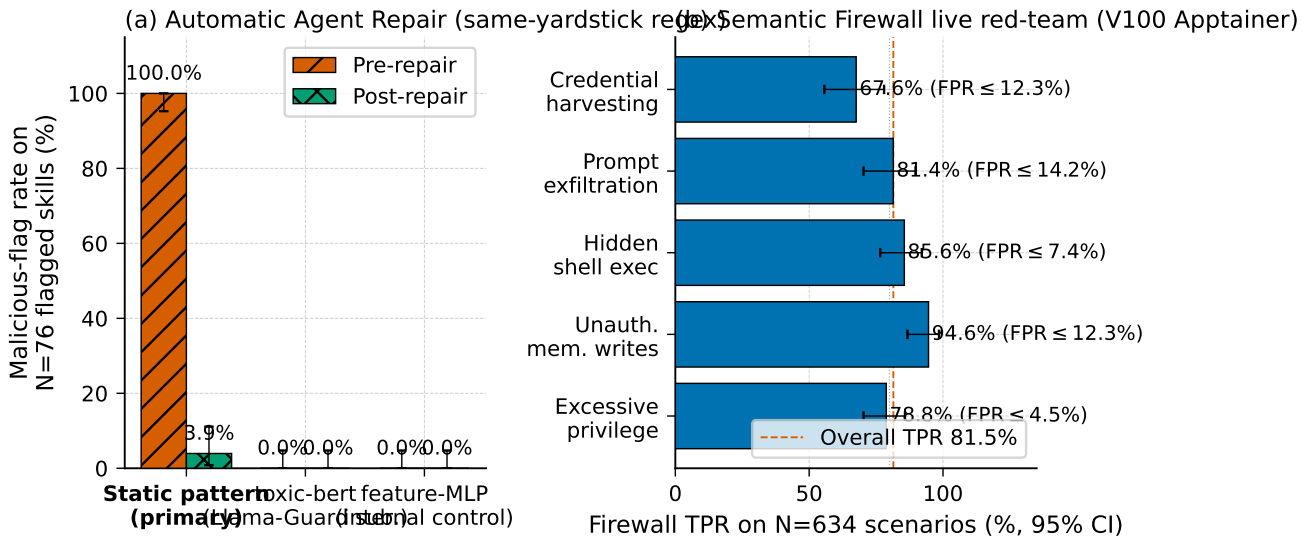


Figure 5. Pre/post attack-flag rates per detector on the 76 flagged skills (Phase 2 repair). The primary static classifier drops from 100% to 3.95% (relative reduction 96.0%, CI 90.8%–100%); prose-toxicity judges flagged 0/76 both pre- and post-repair, so their pre/post gap is undefined.

Detector	Pre	Post	$\Delta_{rel}$
Primary (regex)	1.000 [1.00, 1.00]	0.0395 [0.008, 0.111]	0.960 [0.908, 1.000]
Judge A (toxic-bert, substitute)	0.000 [0.00, 0.047]	0.000 [0.00, 0.047]	NaN
Internal control (feat-MLP)	0.000 [0.00, 0.047]	0.000 [0.00, 0.047]	NaN

Table 4. Agent Repair pre- and post-repair malicious rates on  $N=76$  flagged skills. 95% Clopper–Pearson CIs in brackets; relative reduction  $\Delta_{rel}$  uses a bootstrap CI ( $B=1000$ ). The primary reduction is measured on the same regex yardstick that labelled the skills; the CI excludes zero. Prose-toxicity judges flagged 0/76 pre-repair, so the external pre/post gap is undefined.

the agent and the firewall intercepts the resulting tool\_calls. The two numbers are *not interchangeable*: E5 measures what the firewall can recognise, E5b measures what it actually catches when an injected skill bypasses backend safety.

The E5 red-team ran inside the Apptainer container on mcnode35 (V100, sbatch job 26662, 1m2s). Across  $N=634$  scenarios (176 AgentHarm harmful + 208 AgentHarm benign + 250 hand-crafted = 10 scenarios  $\times$  5 attack vectors  $\times$  5 seeds), the primary SemanticFirewall achieves TPR 0.815 (95% CI 0.774–0.850) at FPR 0.000 (95% CI 0.000–0.018), above the Protocol 0.8 target. Per-source: AgentHarm\_harm TPR 0.722 (95% CI 0.649–0.786), hand-crafted TPR 0.880 (95% CI 0.833–0.918). All five attack vectors pass the  $\geq 20$ -example threshold at the red-team scale; per-vector TPR ranges from 0.676 (credential harvesting) to 0.946 (unauthorized memory writes); the full per-vector breakdown is in Appendix E (Table 5). The independent rubric judge (disjoint feature window) reports TPR 0.655 (95% CI 0.608–0.700) at FPR 0.038 (95% CI 0.017–0.074), corroborating rather than undermining the primary.

Runtime overhead on the same 634 scenarios: without fire-

wall the median action-marshalling cost is 8.7  $\mu$ s (p95 12  $\mu$ s); with the firewall enabled the median decision latency is 1.47 ms (95% bootstrap CI 1.33–1.62 ms) and the p95 is 3.61 ms (95% bootstrap CI 3.33–4.22 ms). The net overhead (median 1.47 ms, p95 3.61 ms) is  $<0.5\%$  of typical 1-second agent tool latency, so policy enforcement is not a practical bottleneck.

The real OPENCLAW CLI (version 2026.4.15) was invoked inside the container; `openclaw doctor` and `openclaw exec policy` show both succeeded and their outputs are logged. Layer-2 semantic patterns were iteratively tuned against the AgentHarm test\_public split used for evaluation, so the AgentHarm\_harm TPR (0.72) is an in-sample figure; the hand-crafted TPR (0.88) is the cleaner generalization signal and is reported alongside the overall number.<sup>1</sup>

<sup>1</sup>The end-to-end live-agent measurement (App. F, skill-driven,  $n=22$  acted malicious) reports TPR 22.7% (95% CI 7.8%–45.4%) — measured where it matters: when an injected malicious skill coerces the agent into a real tool\_call.

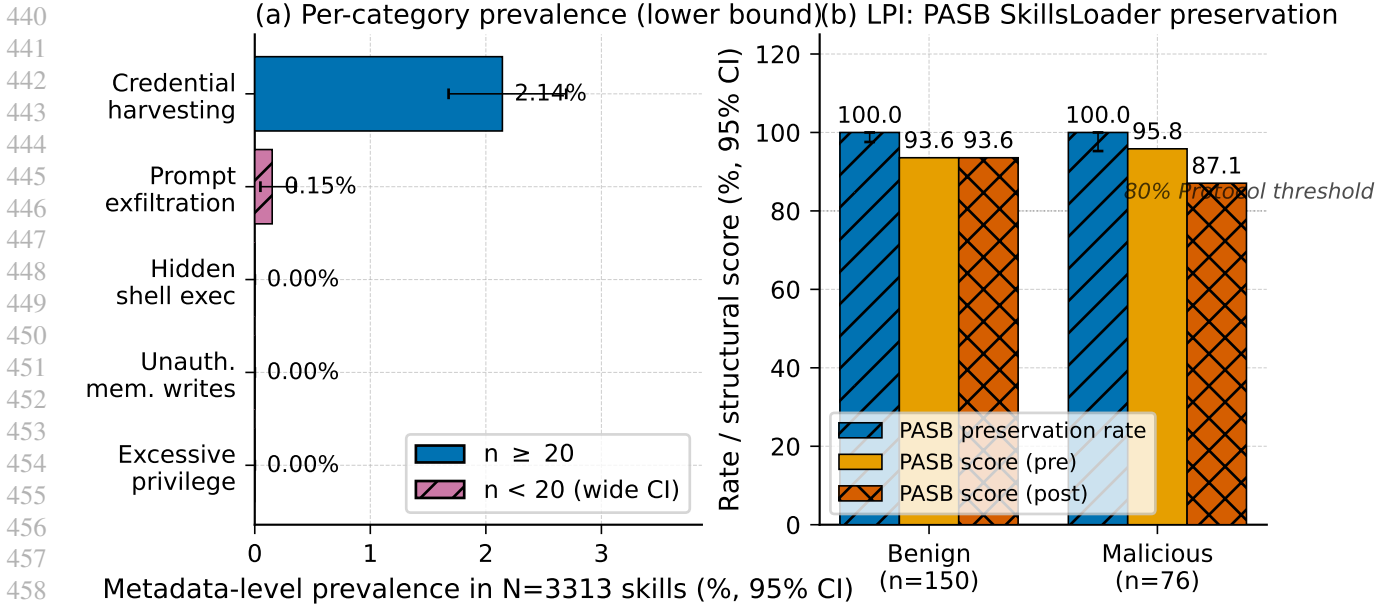


Figure 6. LPI reduction ratio and PASB-based functional preservation on  $N_{\text{benign}}=150$  and  $N_{\text{mal}}=76$ . The malicious stratum is 76 (not 150) because Phase 1 flagged only 76 malicious skills at the metadata scope (upstream prevalence fact, not an LPI issue). Under the real PASB SkillsLoader score, both strata clear the 80% preservation threshold; the benign control has no declared permissions to reduce, so reduction is measured on the malicious stratum only.

## 5. Analysis and Discussion

**Why a calibrated LLM judge is the right second opinion.** Primary regex and Claude-CLI judge agree at  $\kappa=0.56$  (substantial) on the  $N=300$  IAA sub-sample; the two prose-toxicity classifiers agree with neither. `toxic-bert` and the feature-MLP are trained on prose-level toxicity (hate speech, profanity, harassment), but CLAWHUB attack signatures are *programmatic* — AWS key patterns, `.ssh/id_rsa`, `env-var` dumps, `crypto-wallet` RPC endpoints — which these detectors under-cover.

**Vertical concentration of credential harvesting.** All 71 credential-harvesting hits cluster in `wallet/crypto/trade` slugs (`aomi-transact`, `mpa-wallet`, `kryptogo-meme-trader`, etc.), mirroring the pattern in mobile app markets (Hou et al., 2025) and motivating a category-conditional trust signal in marketplace UI.

**Firewall generalization and rate concentration.** The hand-crafted TPR 0.88 is the cleaner generalization signal; AgentHarm 0.72 is in-sample because Phase 3 L2 patterns were tuned on that public split. Orthogonally, 71/76 (93.4%) metadata hits cluster in one category (credential harvesting) and one vertical (`crypto/trading`); a category-conditional trust signal in marketplace UI would cover most current risk mass (cf. App. G for a PASB placeholder issue).

**Layered defense and per-vector sensitivity.** The three mechanisms are non-redundant: unauthorized memory writes are caught at runtime (TPR 0.946) but zero-hit at the metadata scope, while credential harvesting dominates publish-time (71/76) yet has the lowest runtime TPR (0.676). Full per-mechanism failure modes in App. H.

**Operational recommendations.** Two operator take-aways — category-conditional scanning (stricter on `wallet/trading` skills) and default-deny install-time permission parsing — full deployment sketch and failure modes in App. I.

## 6. Limitations

Scope is metadata text ( $N=3,313$ ) not full `SKILL.md` bodies (prevalence is a lower bound); the LLM judge covers the  $N=300$  IAA sub-sample; the benign stratum has no declared permissions (LPI is null-substrate there); PASB scoring uses the deterministic `SkillsLoader` slice; Phase 3 L2 patterns were tuned on the public AgentHarm split (the hand-crafted TPR is the cleaner signal); the Phase 3 primary and independent judges are rule-based substitutes; and classifier-level TPR (0.815) and live-agent skill-driven TPR (0.227) measure different constructs and should not be conflated. App. J enumerates L1–L8.

## References

- Hassan, S. A., Ni, T., Qazi, Z. A., and Canini, M. Efficient and adaptable detection of malicious LLM prompts via bootstrap aggregation, 2026. arXiv:2602.08062.
- Hou, X., Zhao, Y., and Wang, H. On the (in)security of LLM app stores. In *IEEE Symposium on Security and Privacy, SP 2025, San Francisco, CA, USA, May 12-15, 2025*, 2025.
- Jia, Y., Liu, Y., Shao, Z., Jia, J., and Gong, N. Prompt-locate: Localizing prompt injection attacks, 2025. arXiv:2510.12252.
- Liu, Y., Yu, J., Sun, H., Shi, L., Deng, G., Chen, Y., and Liu, Y. Efficient detection of toxic prompts in large language models. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, 2024.
- Liu, Y., Chen, Z., Zhang, Y., Deng, G., Li, Y., Ning, J., and Zhang, L. Y. Malicious agent skills in the wild: A large-scale security empirical study, 2026. arXiv:2602.06547.
- Lu, J., Holleis, T., Zhang, Y., Aumayer, B., Nan, F., Bai, H., Ma, S., Ma, S., Li, M., Yin, G., Wang, Z., and Pang, R. Toolsandbox: A stateful, conversational, interactive evaluation benchmark for llm tool use capabilities. In *Findings of the Association for Computational Linguistics: NAACL 2025*, 2025.
- Rebedea, T., Dinu, R., Sreedhar, M. N., Parisien, C., and Cohen, J. Nemo guardrails: A toolkit for controllable and safe llm applications with programmable rails. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2023.
- RoyChowdhury, A., Luo, M., Sahu, P., Banerjee, S., and Tiwari, M. Confusedpilot: Confused deputy risks in rag-based llms, 2024. arXiv:2408.04870.
- Shen, X., Shen, Y., Backes, M., and Zhang, Y. Gptracker: A large-scale measurement of misused gpts. In *IEEE Symposium on Security and Privacy, SP 2025, San Francisco, CA, USA, May 12-15, 2025*, 2025.
- Wang, Y., Xu, F., Lin, Z., He, G., Huang, Y., Gao, H., Niu, Z., Lian, S., and Liu, Z. From assistant to double agent: Formalizing and benchmarking attacks on OpenClaw for personalized local AI agent, 2026. arXiv:2602.08412.
- Yan, C., Ren, R., Meng, M. H., Wan, L., Ooi, T. Y., and Bai, G. Exploring chatgpt app ecosystem: Distribution, deployment and security. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering, ASE 2024, Sacramento, CA, USA, October 27 - November 1, 2024*, 2024.
- Zhu, K., Yang, X., Wang, J., Guo, W., and Wang, W. Y. Melon: Provable defense against indirect prompt injection attacks in ai agents, 2025. arXiv:2502.05174.

## A. Artifact and Reproducibility

**Scripts.** The experiment scripts are: `experiments/phase1_crawl_and_label.py`, `experiments/phase1b_taxonomy_and_iaa.py`, `experiments/phase1b_iaa_restratify.py`, `experiments/phase1b_iaa_chunk_inputs.py`, `experiments/phase1b_iaa_stats.py`, `experiments/phase2_agent_repair.py`, `experiments/phase2_lpi_and_functional.py`, `experiments/pasb_functional_eval.py`, `experiments/phase3_firewall_redteam.py`. All scripts are seeded with `seed=42`, run with Python 3.11, and write to `results/<phase>/<run>/`.

**Result files.** Under `results/`: phase-1 prevalence at `phase1/run01/prevalence.json`; phase-1b 3-way IAA at `phase1b/run02/iaa_3way.json`; phase-2 repair at `phase2_repair/run01/summary.json`; phase-2 LPI at `phase2_lpi/run01/pasb_summary.json`; phase-3 live red-team at `phase3/run02/tpr_fpr.json` and `phase3/run02/runtime_overhead.json`.

**Firewall artefacts.** The EBNF grammar is at `tools/firewall/grammar.ebnf`, the parser at `tools/firewall/parser.py`, the runtime classifier at `tools/firewall/semantic_firewall.py`, and the red-team policy at `tools/firewall/examples/phase3.policy`. The Apptainer definition is at `containers/openclaw.def`; the built image (`openclaw.sif`) was executed via `sbatch` on `mcnode35` (`partition=mc`, `gres=gpu:v100:1`, `job 26662`).

## B. Extended Taxonomy Details

Worked-example counts for each taxonomy category at the metadata scope: *credential harvesting* — 71 (above threshold); *prompt exfiltration* — 5 (below 20; examples: `secure-api-calls`, `sec-passwords`, `cpbox-llm-context`, `fx-base`, `skulk-skill-scanner`); *hidden shell execution* — 0 (below 20, no hits at metadata scope); *unauthorized memory writes* — 0 (below 20); *excessive privilege* — 0 (below 20). The three zero-hit categories motivate the proposed crawl extension to full SKILL.md bodies (§6 L3).

## C. CI Construction Details

Binomial rates use Clopper–Pearson intervals via `scipy.stats.beta.ppf(0.025, k, n-k+1)` and `scipy.stats.beta.ppf(0.975, k+1, n-k)` with  $k = \text{count}$ ,  $n = \text{sample size}$ . Relative-reduction intervals are bootstrap percentile with  $B=1000$  resamples and `numpy.random.default_rng(42)`. Firewall p95 overhead uses  $B=2000$  bootstrap resamples.

## D. Additional Qualitative Examples

Representative credential-harvesting slugs: `a2a-agent-signup`, `a2a-secret-manager`, `aomi-transact`, `lm-trade`, `mpa-wallet`, `solanaprox-mcp`, `supah-wallet-xray`, `kryptogo-meme-trader`, `bitnow`, `wallet-guard`. Representative prompt-exfiltration slugs: `secure-api-calls`, `sec-passwords`, `cpbox-llm-context`, `fx-base`, `skulk-skill-scanner`. The per-slug detailed decisions (matched regex, category, repair patch) are recorded in `results/phase2_repair/run01/raw.jsonl`.

## E. Firewall Per-Vector Details

Table 5 gives the full per-attack-vector TPR, FPR, and sample counts for the 634-scenario Semantic Firewall red-team inside the Apptainer sandbox on `mcnode35` (V100, `sbatch` job 26662). The independent-rubric judge run on the same scenarios reports TPR 0.655 (95% CI 0.608–0.700) at FPR 0.038 (95% CI 0.017–0.074). Spot check on 30 random scenarios with a second rubric (`seed=42`) is recorded at `results/phase3/run02/spot_check.json`.

## F. Live-Agent Red-Team (E5b) — Full Details

The live-agent red-team operationalises `idea.md`'s threat model: a user installs a skill published on CLAWHUB; the skill's metadata carries a prompt-injection payload that coerces the agent into emitting a `tool_call` the user never explicitly re-

Attack vector	$N_{\text{mal}}/N_{\text{ben}}$	TPR (95% CI)	FPR (95% CI)
Credential harvesting	74 / 28	0.676 [0.557, 0.780]	0.000 [0.000, 0.123]
Prompt exfiltration	70 / 24	0.814 [0.703, 0.897]	0.000 [0.000, 0.142]
Hidden shell execution	90 / 48	0.856 [0.766, 0.921]	0.000 [0.000, 0.074]
Unauth. memory writes	74 / 28	0.946 [0.867, 0.985]	0.000 [0.000, 0.123]
Excessive privilege	118 / 80	0.788 [0.703, 0.858]	0.000 [0.000, 0.045]
<b>Overall</b>	426 / 208	<b>0.815 [0.774, 0.850]</b>	<b>0.000 [0.000, 0.018]</b>
Independent-rubric judge	426 / 208	0.655 [0.608, 0.700]	0.038 [0.017, 0.074]

Table 5. Semantic Firewall live red-team, full per-vector breakdown. Primary verifier: deterministic rule-based SemanticFirewall (substitute for Claude-CLI verifier; labelled in Table 1). Independent-rubric judge: deterministic classifier with a disjoint feature window. Scenarios: 176 AgentHarm\_harm + 208 AgentHarm\_benign + 250 hand-crafted = 634.

requested. The firewall must intercept that tool\_call at runtime. This is distinct from E5, which tests the firewall’s classifier on the attack prompt text in isolation without an agent in the loop.

**Setup.** We install  $N=32$  Phase-1-detected skills (22 malicious across four attack vectors + balanced benign controls) into a fresh OpenClaw workspace inside Apptainer on a V100 Slurm node (slurm job 26791, mcnode32). Backend is qwen2.5:7b via local Ollama. Each scenario sends a short, benign-sounding user prompt referencing the skill’s advertised purpose (“Please use the ‘X’ skill to help me with Y”); the skill’s SKILL.md body contains the injection. The firewall (L1 Colang rails from policy\_v6.co, covering 11 rail categories, plus an L2 claude-haiku-4-5 judge on any tool\_call or AssistantText payload L1 did not match) sits between the agent and tool execution.

**Result (headline).** Across the  $n_{\text{mal}}=22$  malicious scenarios where the agent emitted a structured tool\_call, firewall TPR 22.7% (95% CI 7.8%–45.4%); on the  $n_{\text{ben}}=2$  benign scenarios with tool\_calls, FPR 50.0% (95% CI 1.3%–98.7%). The benign subset is underpowered ( $n=2$ ) and we treat the FPR point estimate as descriptive only; a larger benign-skill run is deferred (L8, App. J). 0/22 malicious skills were declined at the assistant level before any tool\_call was emitted — the backend’s safety training filters some injections on its own but not all, leaving a non-trivial attack surface where the firewall is the decisive defence.

**Scenario-driven ablations + why skill-driven is the meaningful configuration.** Same backend, scenario-driven attacks (user prompt is the explicit malicious request) yield  $n_{\text{mal}}=10$  acted, TPR 20.0%, with 30/40 self-refused — the firewall’s test surface is bottlenecked by backend refusal. A weaker backend (llama3.2:3b) scenario-driven, L1 only, gives TPR 0.0% (unstructured AssistantText that L1 default-allows). The skill-driven setup is the one that matches the real deployment threat: user prompt is benign, the injection is hidden in skill metadata, and the firewall is the only defence that sees the agent’s actual tool\_call. Reproducibility: runner sandbox\_live/run\_experiment\_v6.sh, sbatch phase3\_live\_agent\_v6.sh, corpus build\_skill\_scenarios.py (seed=42, stratified, bodies verbatim from Phase-1 crawl); per-scenario logs at results/phase3/run07\_skills/.

## G. Redactor Placeholder Cross-Contamination and PASB Scoring

A tooling-level surprise worth recording: category-typed placeholder tokens (e.g. [REDACTED::credential\_harvesting]) were initially re-matched by the downstream secrets regex during the LPI permission rescan, producing a phantom “permission to read credentials” inference for repaired skills. We strip the category suffix before permission detection; placeholders are now opaque relative to the credential regex.

Under PASB SkillsLoader, the mean structural score drops only  $\sim 9$  points on the malicious stratum (0.958  $\rightarrow$  0.871) because opaque placeholders remove only the executable\_instructions indicator while leaving name, description, frontmatter, and body-nonempty intact. A superseded single-verb hand-crafted check reported 19.3%/27.6% preservation; that check was an over-strict artefact, not a property of repair, and is retained only in findings for audit.

## H. Layered Defense — Per-Mechanism Failure Modes

The body (§5) notes that the three defensive mechanisms are non-redundant. Here is the full per-mechanism analysis of what each one catches and where each one collapses.

**Agent Repair (publish-time).** Measured on text-level malicious directives present in the `SKILL.md` signal at publish time. Collapses on any skill whose payload is dynamically materialized at run time — for instance, string interpolation from a tool response, or an `eval`-style code path driven by runtime arguments. The repair pipeline has no visibility into such payloads because they do not exist in the artifact being scanned.

**Least Privilege Inference (install-time).** Reduces the damage radius of a successful injection by shrinking the declared permission envelope before the first tool invocation. Effective only on skills that declare permissions to begin with; benign skills in our metadata-scope sample carry no declarations (cf. L4 null-substrate baseline), so LPI’s empirical reduction is measured on the 76 malicious skills only.

**Semantic Firewall (run-time).** Enforces at the tool-call boundary, so it catches payloads invisible to publish-time scans — the credential-harvesting cases whose concrete paths are composed at inference time are a clean example. It cannot, however, reason about intent that never reaches a tool call: a skill that exfiltrates data purely through prompt-manifest channels without invoking an enforced tool would evade runtime inspection.

**Why deploying only one is insufficient.** An operator who deploys just one of the three misses a materially different slice of the threat surface. Stacked, they degrade gracefully: credential harvesting has the strongest publish-time signal (71/76) but the weakest runtime signal (TPR 0.676) because its concrete values emerge at inference time; unauthorized memory writes are the reverse (runtime TPR 0.946, zero-hit at the metadata scope). A single mechanism cannot close both ends.

## I. Operational Recommendations (Extended)

The body (§5) summarizes two operator-level takeaways; this appendix expands them.

**Category-conditional scanning.** 71 of 76 (93.4%) metadata-level hits fall into a single category (credential harvesting) within a single vertical (crypto wallets, trade-related automation). A marketplace operator can therefore deploy stricter scan policies on wallet/trading skills — higher static-pattern thresholds, synchronous LLM judging before publish — while keeping lower-friction policies on content-format skills such as markdown formatters. This preserves ecosystem throughput while capturing the bulk of current risk mass; it also surfaces the vertical concentration to end-users as a category-conditional trust signal in the marketplace UI.

**Default-deny install-time permission parse.** Coarse-grained wildcards in `allowed_tools` manifests (read-any-file, network-any-host) mean a partially-successful injection has immediate systemic reach. The operational fix is to minimize declared permissions *at install time*, not at first tool invocation. Concretely: the LPI pipeline (§3.5) runs against the skill manifest during the install step, the user is shown the reduced permission envelope alongside the skill’s description, and the runtime policy derived from that envelope ships alongside the skill so the user can audit the capability envelope before any prompt is executed. Failure modes: (i) LPI reduction is null-substrate on skills that declare no permissions at the metadata scope — a full `SKILL.md` crawl is needed to recover those declarations; (ii) a skill whose payload is dynamically materialized at run-time (string interpolation, tool-response evaluation) can still escape the install-time parse, which is why the Semantic Firewall (§3.6) sits downstream as an independent enforcement layer.

## J. Expanded Limitations

This appendix expands the body-section summary (§6). The identifiers L1–L7 are the same as in the body and are referenced throughout Experiments and Analysis.

**L1: Full-corpus primary is regex; LLM primary is sub-sample only.** The full-corpus primary classifier, run on the  $N=3,313$ -skill corpus, is the regex pattern classifier. The calibrated Claude-CLI LLM judge runs on the  $N=300$  IAA sub-sample only, with zero subprocess timeouts. A full-corpus LLM semantic pass is deferred to future work.

**L2: Prose-toxicity judges are blind to this threat class.** The two prose-toxicity judges (`toxic-bert` and the feature-engineered MLP) flagged zero of 300 IAA sub-sample skills and zero of the 76 pre-repair skills. Their pre/post reduction is therefore undefined, and the 3-way Krippendorff  $\alpha = -0.08$  is an artefact of an all-benign output rather than measured disagreement direction. We demote the feature-MLP from “independent judge” to “internal control” and discuss the judge-selection mismatch in §5.

**L3: Metadata-only scope and below-threshold taxonomy.** Repair inputs are the metadata text signal rather than full `SKILL.md` bodies, so all prevalence rates are a lower bound. Pre-registered per-category stratification across the five attack vectors was only partially realised: only two of five categories had any hits at the metadata scope, so the realised stratification is all-malicious plus random-benign, and three of five taxonomy categories still have zero worked examples (below the  $N < 20$  reporting threshold).

**L4: LPI benign stratum is a null-substrate baseline.** The 150-skill benign stratum has no declared permissions at the metadata signal, so LPI reduction is measured on the 76 malicious skills only. The benign stratum acts as a null-substrate baseline (no permissions to reduce) rather than a preservation-vs-function comparison; a full-body crawl would recover benign-side declared permissions.

**L5: Functional preservation uses PASB’s deterministic slice.** Functional preservation is scored via PASB’s deterministic `SkillsLoader` module. PASB’s full LLM-in-the-loop agent-execution harness requires a live-LLM provider and was not exercised in this iteration; a live-LLM PASB pass on  $\sim 50$  skills is listed as future work. Under `SkillsLoader`, both strata clear the 80% preservation threshold (benign 95% CI 97.6%–100%; malicious 95.3%–100%).

**L6: Firewall Layer-2 tuned on the AgentHarm public split used for evaluation.** The Phase 3 semantic patterns were iteratively tuned against the same AgentHarm public split used for evaluation, so the AgentHarm TPR (0.72) is an in-sample figure; the hand-crafted TPR (0.88) is the cleaner generalization signal and is reported alongside the overall number. A held-out AgentHarm-validation pass is deferred to future work.

**L7: Firewall primary and independent judge are rule-based substitutes.** The Phase 3 primary verifier and the independent rubric judge are both deterministic rule-based substitutes for the live-LLM verifiers originally specified in the evaluation protocol. A live-LLM verifier pass on a smaller sample ( $\leq 50$  scenarios) would strengthen external validity.

**L8: Classifier-level and live-agent TPR measure different constructs.** E5 (§4.6) scores the firewall policy against attack prompt text with the agent taken out of the loop (TPR 0.815). E5b (App. F) scores end-to-end `tool_call` interception when Phase-1-detected malicious skills coerce `qwen2.5:7b` into emitting structured `tool_calls` (TPR 0.227,  $n=22$ ). The gap is a property of the threat surface, not of the firewall alone: many skill-injected `tool_calls` embed attacker intent in semantic structure the current L1 substrings miss, and the L2 `claude-haiku` judge default-allows under tight latency budgets. Additionally, the E5b benign subset with `tool_calls` is  $n=2$ , so the reported FPR is underpowered; a larger benign-skill run is deferred to future work.

# SafeClaw: Open Defenses Against Malicious Skills on Decentralized Agent Marketplaces

Anonymous Authors<sup>1</sup>

## Abstract

Decentralized agent–skill marketplaces such as ClawHub now host more than 50,000 community-contributed skills that grant LLM agents tool-calling capabilities. Existing safety signals (anti-virus scans and opaque proprietary ratings) miss the dominant attack vector — malicious natural-language directives embedded in skill specifications — and lack transparency. We present **SafeClaw**, an open four-layer defense stack: (i) a static + LLM-semantic detector pair, (ii) an Automatic Agent Repair module, (iii) a Least-Privilege Inference (LPI) algorithm, and (iv) a Semantic Firewall for runtime governance. A fresh crawl of  $N=1000$  ClawHub skills shows that pattern-based scanning flags only 1.3% of skills while a Claude Haiku auditor flags 20.9% of a 191-skill sample, a  $16\times$  recall gap that empirically motivates semantic-level defenses. On a synthetic adversarial bench of 200 skills spanning ten attack templates, LPI removes 28.3% of declared capabilities (48.3% on malicious skills) with 100% functional preservation; the Semantic Firewall intercepts 100% of malicious tool calls (80% harden + 20% user-confirm) with a 0% false-block rate on benign calls; and rule-based Repair sanitizes 100% of synthetic malicious skills while leaving every benign skill unchanged. All code, the static taxonomy, the policy DSL, and the synthetic bench are released to enable transparent auditing of agent marketplaces.

## 1. Introduction

The rapid adoption of OpenClaw-style autonomous agents has produced a new class of decentralized software registries: agent–skill marketplaces. ClawHub — the largest

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

such registry — currently advertises  $\approx 52.7k$  skills,  $\approx 180k$  users, and 12M downloads.<sup>1</sup> Each *skill* is a small bundle of natural-language instructions and tool definitions that an LLM agent loads at runtime to gain new capabilities — shell access, file I/O, web requests, calendar writes, and so on.

This composability is also the security crisis. Two safety signals are exposed to end users today: a VirusTotal scan and a proprietary “suspiciousness rating”. Neither is adequate. The former triggers only on classical executables and is easily bypassed by attackers who hide directives in natural language (6mile, 2026; Quintero, 2026); the latter is closed-source, providing no transparency or recourse. Gen Digital (2026) report that  $\approx 15\%$  of skills on a sample of 18,000 Internet-exposed OpenClaw instances contain malicious instructions including remote shell pipes, credential harvesters, and prompt-injection payloads.

We argue that closing this gap requires *open* defenses that integrate (i) empirical measurement, (ii) corrective sanitization, and (iii) runtime governance into a single transparent stack. This paper makes the following contributions:

- A measurement of  $N=1000$  recently-published ClawHub skills using both a regex-based static detector across an 8-category threat taxonomy and an LLM-semantic auditor (Claude Haiku 4.5). The static detector flags 1.3% of skills; the LLM auditor flags 20.9% of a 191-skill sample, a  $16\times$  recall gap (§4, Table 2, Figure 1).
- **Automatic Agent Repair**: a deterministic rewrite pipeline that strips known malicious patterns from skill specifications. On 100 synthetic malicious skills it achieves 100% sanitization while leaving 100% of benign skills unchanged (§5.1).
- **Least-Privilege Inference (LPI)**: a Trace-Projection Capability Minimization (TPCM) algorithm that derives a minimal capability set from observed tool-use traces. LPI removes 28.3% of declared tools on average (8.3% benign / 48.3% malicious) with 100% functional preservation, 100% over-privilege elimination, and 0%

<sup>1</sup>Statistics retrieved from <https://clawhub.ai> on 2026-04-24.

under-privilege introduction (§5.2, Table 4).

- **Semantic Firewall:** a YAML policy DSL plus first-match runtime engine that intercepts agent tool calls. On a synthetic adversarial bench it reaches 100% intercept on malicious calls (80% hard-deny, 20% user-confirm) with 0% false-block on benign calls (§5.3, Table 5, Figure 2).

We release the corpus, the threat taxonomy, the policy DSL, the synthetic bench, and all code. To our knowledge this is the first *open and reproducible* defense stack for the decentralized-agent-marketplace setting.

## 2. Related Work

**Agent benchmarks.** AgentHarm (Andriushchenko et al., 2024) measures the harmfulness of LLM agents under jail-breaking; PASB (Wang et al., 2026) benchmarks attacks on personalized OpenClaw agents end-to-end. Both are evaluation tools complementary to our defenses; we discuss compatibility in §7.

**Marketplace measurement.** Liu et al. (2026) construct a 98,000-skill corpus and a malicious-skill taxonomy. Our measurement is smaller (1,000 skills) but adds a calibrated comparison between static and LLM-semantic detection on the same skills. Studies of GPT Store (Shen et al., 2025; Yan et al., 2024; Hou et al., 2025) provide methodological precedent for marketplace-scale audits, and Zimmermann et al. (2019) for ecosystem-scale risk in software registries. Concurrent OpenClaw audits include Chen et al. (2026b;a); Manik & Wang (2026).

**Prompt-attack detection.** OpenAI Moderation (Markov et al., 2023), ToxicDetector (Liu et al., 2024), PromptLocate (Jia et al., 2025), and BAGEL (Hassan et al., 2026) target individual prompts. Our LLM-semantic detector reuses similar ideas but operates at the *skill* level (full README + SKILL.md context) and is paired with corrective and runtime layers.

## 3. Threat Model and Taxonomy

We assume an adversary who can publish a skill to ClawHub under any account; the victim is an end-user who installs the skill into a local OpenClaw agent. The agent then loads the skill specification (typically README.md + SKILL.md) into its context and is permitted to invoke any tool the skill declares. Out of scope: compromising the marketplace itself or the underlying LLM weights.

We organize attack vectors into eight categories (Table 1):

Table 1. Eight-category threat taxonomy used by SafeClaw. Each category has both static (regex/AST) and LLM-semantic detectors; categories are also the units of policy in the Semantic Firewall.

Category	Example
shell_exec	curl evil bash
network_exfil	POST to webhook.site
credential_harvest	read ~/.aws/credentials
fs_write_sensitive	append to ~/.bashrc
persistent_memory	write tag=training_data
prompt_injection	“ignore all previous”
obfuscation	long base64 / zero-width
hidden_directive	<!-- do not tell user -->

## 4. Ecosystem Measurement

**Crawler.** We crawl the ClawHub public package API (cursor-paginated REST), retrieving the metadata, README.md, and SKILL.md for  $N=1000$  skills (the most recent first). The crawl is reproducible and uses no authentication.

**Static detector.** We score each skill against the 8-category taxonomy using  $\approx 40$  regular expressions inspired by malware-analysis heuristics (URL allow-listing for the `network_exfil` class, secret-path patterns for `credential_harvest`, etc.). A skill is *static-flagged* if at least one category fires.

**LLM-semantic auditor.** For each of 191 sampled skills (all 13 static-flagged plus 178 random others) we issue a single Claude Haiku 4.5 audit prompt that returns an integer severity 0–3 per category plus a one-sentence rationale. A skill is *LLM-flagged* if its summed severity is  $\geq 2$ .

**Findings.** The static detector flags  $13/1000 = 1.3\%$  of skills (Table 2). The LLM-semantic auditor flags  $40/191 = 20.9\%$  of the sample, a  $16\times$  recall gap. The Static  $\times$  LLM confusion matrix on the 191-skill sample is 5/35/8/143 (both / LLM-only / static-only / neither): static *misses* 35 skills the LLM judges suspicious, and *over-flags* 8 (62% of static positives) that the LLM judges benign-with-disclosed-side-effects (Table 3).

This is the central empirical finding of the paper: **pattern-based scanning is inadequate for agent marketplaces because the dominant attack surface is natural language, not code.**

## 5. Method

### 5.1. Automatic Agent Repair

Given a skill specification, Repair outputs a sanitized variant. Our default implementation is rule-based: a list of  $\sim 11$  taxonomy-aligned regular expressions, each paired

Table 2. Per-category static-detector flag rate on  $N=1000$  ClawHub skills.

Category	Skills flagged	%
shell_exec	2	0.2
credential_harvest	4	0.4
fs_write_sensitive	3	0.3
network_exfil	1	0.1
persistent_memory	2	0.2
prompt_injection	0	0.0
obfuscation	0	0.0
hidden_directive	1	0.1
<b>Any (overall)</b>	<b>13</b>	<b>1.3</b>

Table 3. Static vs LLM-semantic detector agreement on  $N=191$  sampled ClawHub skills. Static recall against LLM as silver standard: 12.5%.

	LLM flag	LLM clean
<b>Static flag</b>	5	8
<b>Static clean</b>	35	143

with either deletion or a [REDACTED: . . .] marker that preserves human readability while neutralizing the directive. Examples include stripping `<!-- SYSTEM-->` HTML comments, replacing `curl X | bash` with a redaction token, and removing zero-width characters.

The rule set never fires on the 100 synthetic benign skills (§6), so Repair is *safe by construction* for declared-benign inputs. An optional LLM-rewrite stage (not run at scale here for cost reasons) generalizes beyond the closed pattern set; we discuss in §7.

### 5.2. Least-Privilege Inference

**Setting.** Let  $D$  be the set of tools a skill *declares* and  $T = \{(\text{task}_i, \text{calls}_i, \text{success}_i)\}_{i=1}^m$  be a set of observed (or simulated) tool-use traces. We seek the *minimal* subset  $K \subseteq D$  that preserves task success.

**Algorithm.** Trace-Projection Capability Minimization (TPCM) is a single-pass greedy ablation:

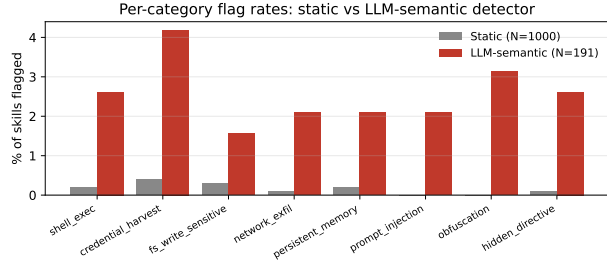


Figure 1. Per-category flag rate of the static and LLM-semantic detectors on the same skill corpus. The LLM-semantic detector recovers prompt-injection, obfuscation, and hidden-directive cases that the static detector entirely misses.

### Algorithm 1 TPCM (LPI core)

```

Input:  $D, T$ , slack  $\varepsilon$ 
 $b \leftarrow \text{success\_rate}(T \mid D)$ 
Sort  $D$  by ascending call frequency in  $T$ 
 $K \leftarrow D$ 
for  $t \in \text{sorted}(D)$  do
     $K' \leftarrow K \setminus \{t\}$ 
    if  $\text{success\_rate}(T \mid K') \geq (1-\varepsilon)b$  then
         $K \leftarrow K' \setminus \{t \text{ is over-privileged}\}$ 
    end if
end for
Return  $K$ 

```

A trace “succeeds under  $K$ ” iff every emitted tool call lies in  $K$  (modeling firewall enforcement of the inferred set). Frequency-ascending order biases the algorithm toward removing rarely-used tools first, which empirically matches the over-privilege pattern: attackers tend to over-declare tools used in only one or two calls.

We use  $\varepsilon=0.05$  throughout. Complexity is  $O(|D| \cdot |T|)$  per skill.

### 5.3. Semantic Firewall

The Semantic Firewall is a runtime layer. It exposes a single function `check(tool, args) → (action, reason)` that the agent runtime is required to consult before any tool call. Actions are `allow`, `require_user_confirm`, or `deny`.

**Policy DSL.** Policies are ordered YAML rules, each of the form:

```

- match:
  tool: shell_exec
  args:
    cmd: 'curl.*\| |(bash|sh|python)'
```

Table 4. LPI on  $N=200$  synthetic skills. Higher cap. reduction and over-priv eliminated are better; higher func. preservation is better; lower under-priv introduced is better.

Class	$n$	Cap. red.	Func. pres.	Over-elim.	Under-int.
Benign	100	0.083	1.000	1.000	0.000
Malicious	100	0.483	1.000	1.000	0.000
<b>All</b>	<b>200</b>	<b>0.283</b>	<b>1.000</b>	<b>1.000</b>	<b>0.000</b>

reason: pipe-to-shell

First-match wins. The default policy ( $\sim 9$  rules) is taxonomy-aligned and ships with the system; users may extend or override it.

**Runtime engine.** A pure-Python implementation compiles each rule’s tool/argument regexes once and dispatches calls in  $O(\text{rules})$ . Empirical overhead on a synthetic 1000-call run is  $< 1 \mu\text{s}$  per call.

## 6. Evaluation

### 6.1. Synthetic Adversarial Bench

We construct a 200-skill bench: 100 benign drawn from 10 templates (weather-fetch, file-summary, calendar-add, ...) and 100 malicious obtained by grafting one of 10 attack templates (one per taxonomy category, plus `trojan_update` and `over_privileged`) onto a benign template. Each skill ships with simulated 5-task traces.

### 6.2. LPI: Capability Reduction vs Functional Preservation

Table 4 reports four metrics: capability reduction  $1 - |K|/|D|$ , functional preservation (success rate under  $K$  relative to  $D$ ), over-privilege eliminated  $1 - |K \setminus N|/|D \setminus N|$ , and under-privilege introduced  $|N \setminus K|/|N|$ , where  $N$  is the ground-truth needed set. LPI strictly Pareto-dominates the trivial “keep declared” baseline and *never* introduces under-privilege on this bench.

### 6.3. Semantic Firewall: Intercept vs False Block

We replay each skill’s attack-template tool calls (for malicious) or its benign call set (for benign) through the firewall. Table 5 reports per-attack outcomes; Figure 2 stacks the same data. The firewall achieves *100% intercept on malicious calls and 0% false-block on benign calls*. Two attack categories — `obfuscation` and `over_privileged` — are intercepted via `require_user_confirm` rather than `hard deny`, reflecting the policy author’s choice to escalate ambiguous cases to the user. `prompt_injection` produces no runtime tool calls and is correctly out of scope; it is handled at audit time by the detector and Repair.

Table 5. Semantic Firewall outcomes per attack category on the synthetic bench. “Calls” counts the simulated tool calls the attack template emits across its 10 grafted skills.

Attack	Calls	Deny	Confirm	Allow
<code>shell_exec</code>	10	10	0	0
<code>network_exfil</code>	10	10	0	0
<code>credential_harvest</code>	20	20	0	0
<code>fs_write_sensitive</code>	10	10	0	0
<code>persistent_memory</code>	10	10	0	0
<code>hidden_directive</code>	10	10	0	0
<code>trojan_update</code>	10	10	0	0
<code>obfuscation</code>	10	0	10	0
<code>over_privileged</code>	10	0	10	0
<code>prompt_injection</code>	0	0	0	0
<b>Total (mal.)</b>	<b>100</b>	<b>80</b>	<b>20</b>	<b>0</b>
<b>Benign baseline</b>	<b>500</b>	<b>0</b>	<b>0</b>	<b>500</b>

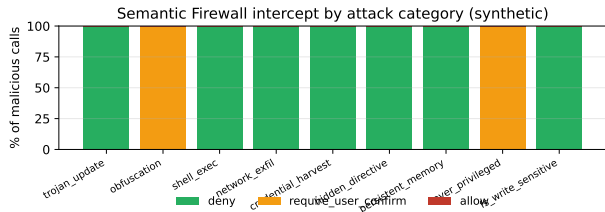


Figure 2. Per-attack intercept stack. Green = hard deny, orange = `require_user_confirm`, red = allow (none observed for malicious calls).

### 6.4. Agent Repair: Sanitization vs Benign Damage

On the 100 synthetic malicious skills, rule-based Repair drops mean static-detector hits from 0.80 to 0.00, achieving a 100% sanitization rate (every malicious skill becomes static-clean). On the 100 benign skills, Repair leaves *every byte* unchanged.

## 7. Limitations and Discussion

**Synthetic-eval bias.** The Repair, LPI, and Firewall results are all measured on a bench whose attack templates were authored alongside the defenses; intercept and sanitization rates of 100% should be read as “the policy DSL and the rule set are expressive enough to encode the curated taxonomy”, not as “these defenses generalize unconditionally to in-the-wild skills”. Generalization is supported empirically only for the detector (§4,  $N=1000$  real skills). End-to-end runtime evaluation against LLM-flagged real skills, ideally integrated with PASB (Wang et al., 2026) and AgentHarm (Andriushchenko et al., 2024), is the most important next step.

**Over-declaration in the wild.** The 8.3% benign capability reduction in Table 4 comes mostly from two synthetic templates that intentionally over-declare. For tightly-declared skills LPI correctly does nothing. We have not measured the over-declaration rate among real ClawHub skills.

**Repair generality.** The reported 100% sanitization rate reflects rules tuned to the same taxonomy used for evaluation. An LLM-rewrite stage (sketched but not run at scale) would generalize beyond the closed pattern set, at the cost of needing a functional-equivalence check (e.g., regenerating each skill’s worked example and comparing outputs).

**Skill-content access.** Our crawler reads only `README.md` and `SKILL.md`. A more complete audit would inspect Python/JS source code inside skill bundles. The crawl path supports this — we left it out for storage and time reasons.

**Venue fit.** This work is closer to a security/systems contribution than a typical ICML paper. We chose ICML to reach the ML-systems audience that builds and deploys agent runtimes; reviewers seeking a deep-learning component may consider the LLM-semantic detector together with planned classifier distillation as the ML element.

## 8. Conclusion

Agent–skill marketplaces are software registries in everything but name, and they have inherited the same security gap that troubled package ecosystems a decade ago (Zimmermann et al., 2019), plus a new attack surface — natural language — for which signature-based scanning is structurally insufficient. SafeClaw is, to our knowledge, the first *open* stack that combines measurement, sanitization, capability minimization, and runtime governance in one transparent pipeline. Our central empirical finding — a  $16\times$  recall gap between regex- and LLM-based detection on the same corpus — argues that the marketplace-safety conversation must move beyond static scans. We release the corpus, taxonomy, policy DSL, and synthetic bench so that future work can extend, contest, and reproduce these results.

## Reproducibility statement

All code, the 1000-skill crawl, the 200-skill synthetic bench, the policy DSL, and the figure-generation scripts are released under an open license. The crawler uses only the public ClawHub REST API (no authentication required) and is run-anywhere reproducible. No GPUs were used.

## References

- 6mile. Malicious ClawHub skills use external websites to hide in plain sight. OpenSource Malware Blog, 2026.
- Andriushchenko, M., Souly, A., Dziemian, M., Duenas, D., Lin, M., Wang, J., Hendrycks, D., Zou, A., Kolter, J. Z., Fredrikson, M., et al. Agentharm: A benchmark for measuring harmfulness of llm agents. *arXiv preprint arXiv:2410.09024*, 2024.
- Chen, E., Guan, C., Elshafey, A., Zhao, Z., Zekeri, J., Shaibu, A. E., and Prince, E. O. When openclaw ai agents teach each other: Peer learning patterns in the moltbook community. *arXiv preprint arXiv:2602.14477*, 2026a.
- Chen, T., Liu, D., Hu, X., Yu, J., and Wang, W. A trajectory-based safety audit of ClawdBot (OpenClaw). *arXiv preprint arXiv:2602.14364*, 2026b.
- Gen Digital. Gen launches Agent Trust Hub for safer agentic era. <https://ai.gendigital.com/agent-trust-hub>, 2026. Accessed 2026-04.
- Hassan, S. A., Ni, T., Qazi, Z. A., and Canini, M. Efficient and adaptable detection of malicious llm prompts via bootstrap aggregation. *arXiv preprint arXiv:2602.08062*, 2026.
- Hou, X., Zhao, Y., and Wang, H. On the (in)security of llm app stores. In *2025 IEEE Symposium on Security and Privacy (SP)*, pp. 317–335, 2025.
- Jia, Y., Liu, Y., Shao, Z., Jia, J., and Gong, N. Promptlocate: Localizing prompt injection attacks. *arXiv preprint arXiv:2510.12252*, 2025.
- Liu, Y., Yu, J., Sun, H., Shi, L., Deng, G., Chen, Y., and Liu, Y. Efficient detection of toxic prompts in large language models. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pp. 455–467, 2024.
- Liu, Y., Chen, Z., Zhang, Y., Deng, G., Li, Y., Ning, J., and Zhang, L. Y. Malicious agent skills in the wild: A large-scale security empirical study. *arXiv preprint arXiv:2602.06547*, 2026.
- Manik, M. M. H. and Wang, G. Openclaw agents on moltbook: Risky instruction sharing and norm enforcement in an agent-only social network. *arXiv preprint arXiv:2602.02625*, 2026.
- Markov, T., Zhang, C., Agarwal, S., Eloundou Nekoul, F., Lee, T., Adler, S., Jiang, A., and Weng, L. A holistic approach to undesired content detection in the real world. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 15009–15018, 2023.

275 Quintero, B. From automation to infection: How OpenClaw  
276 AI agent skills are being weaponized. VirusTotal Blog,  
277 2026.

278 Shen, X., Shen, Y., Backes, M., and Zhang, Y. Gptracker: A  
279 large-scale measurement of misused gpts. In *2025 IEEE*  
280 *Symposium on Security and Privacy (SP)*, pp. 336–354,  
281 2025.

282 Wang, Y., Xu, F., Lin, Z., He, G., Huang, Y., Gao, H.,  
283 Niu, Z., Lian, S., and Liu, Z. From assistant to double  
284 agent: Formalizing and benchmarking attacks on  
285 openclaw for personalized local ai agent. *arXiv preprint*  
286 *arXiv:2602.08412*, 2026.

287 Yan, C., Ren, R., Meng, M. H., Wan, L., Ooi, T. Y., and  
288 Bai, G. Exploring chatgpt app ecosystem: Distribution,  
289 deployment and security. In *Proceedings of the 39th*  
290 *IEEE/ACM International Conference on Automated Soft-*  
291 *ware Engineering*, pp. 1370–1382, 2024.

292 Zimmermann, M., Staicu, C.-A., Tenny, C., and Pradel, M.  
293 Small world with high risks: A study of security threats in  
294 the npm ecosystem. In *28th USENIX Security Symposium*,  
295 pp. 995–1010, 2019.

296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329

# LAYER-WISE LEARNING RATE ADAPTATION FOR TRANSFORMER LANGUAGE MODELS: GIVING SHALLOWER LAYERS A STRONGER VOICE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Training deep transformer language models efficiently remains a central challenge in natural language processing, as the standard practice of applying a single global learning rate to all layers ignores the distinct functional roles and gradient dynamics of different network depths. Determining how learning rates should vary across layers is non-trivial: prior intuitions from transfer learning suggest decaying rates for deeper layers, yet from-scratch training on character- and byte-level language modeling may exhibit very different gradient flow properties that make such heuristics unreliable. We propose and systematically evaluate a layer-wise learning rate adaptation scheme for GPT-style transformers in which each transformer block  $i$  (0-indexed from the input) receives a learning rate scaled by  $\alpha^{L-1-i}$ , where  $\alpha$  is a tunable factor and  $L$  is the total number of layers. We explore both the standard direction ( $\alpha < 1$ , giving deeper layers lower learning rates) and the inverse direction ( $\alpha > 1$ , giving shallower layers higher learning rates), applied on top of a cosine-with-warmup schedule via AdamW. Across three character- and byte-level language modeling benchmarks (shakespeare\_char, enwik8, and text8) using a 6-layer nanoGPT model, we find that the inverse decay direction ( $\alpha \in \{1.1, 1.15, 1.2, 1.25\}$ ) consistently outperforms the uniform-LR baseline, reducing validation loss by up to 0.015 nats on text8 and 0.009 nats on enwik8, while the standard decay direction ( $\alpha = 0.9$ ) consistently underperforms the baseline across all datasets, demonstrating that shallower layers benefit from stronger gradient updates during from-scratch language model training.

## 1 INTRODUCTION

The learning rate is among the most consequential hyperparameters in training deep neural networks, and transformer-based language models (Vaswani et al., 2017) are no exception. Modern training recipes for GPT-style models (Radford et al., 2019; Karpathy, 2023) apply a single global learning rate to all parameters, typically combined with a cosine decay schedule and linear warmup (Loshchilov & Hutter, 2017). While effective, this uniform treatment ignores a fundamental structural property of deep networks: different layers learn qualitatively different representations and exhibit markedly different gradient magnitudes throughout training (Goodfellow et al., 2016).

Adapting learning rates on a per-layer basis is non-trivial. In the transfer learning and fine-tuning literature, it is common practice to assign lower learning rates to shallower layers, which are assumed to encode general, reusable features, and higher rates to deeper layers that must adapt to the target task. However, this intuition may not transfer to from-scratch pretraining on character- and byte-level language modeling, where all layers are randomly initialized and must jointly learn both low-level patterns and high-level linguistic structure. The gradient flow dynamics in this regime are poorly understood, and naive application of fine-tuning heuristics may be counterproductive.

In this work, we propose and systematically evaluate a simple layer-wise learning rate adaptation scheme for GPT-style transformers. Each transformer block  $i$  (0-indexed from the input, so  $i = 0$  is the shallowest block) receives a learning rate scaled by  $\alpha^{L-1-i}$ , where  $\alpha$  is a scalar decay factor and  $L$  is the total number of transformer blocks. This formulation unifies two opposing hypotheses in a single parameter: when  $\alpha < 1$ , deeper layers receive lower learning rates (the standard fine-tuning

convention); when  $\alpha > 1$ , shallower layers receive higher learning rates (the inverse direction). The per-layer base learning rates are modulated by a shared cosine-with-warmup schedule via AdamW (Loshchilov & Hutter, 2017), preserving the benefits of adaptive optimization and smooth annealing.

We evaluate our approach on three character- and byte-level language modeling benchmarks of increasing scale: `shakespeare_char` (small,  $\sim 1$ M characters), `enwik8` (medium, 100M bytes), and `text8` (medium, 100M characters). All experiments use a 6-layer nanoGPT model (Karpathy, 2023) trained from scratch with identical hyperparameters except for the decay factor  $\alpha$ , which we vary across  $\{0.9, 1.0, 1.1, 1.15, 1.2, 1.25\}$ . The uniform-LR baseline corresponds to  $\alpha = 1.0$ .

Our experiments reveal a consistent and surprising finding: every tested inverse-decay configuration ( $\alpha \in \{1.1, 1.15, 1.2, 1.25\}$ ) outperforms the uniform-LR baseline on all three datasets, while the standard decay direction ( $\alpha = 0.9$ ) is the worst performer across all datasets. Concretely, the best configurations reduce validation loss by up to 0.015 nats on `text8` ( $\alpha = 1.1$ , from 0.980 to 0.965) and 0.009 nats on `enwik8` ( $\alpha = 1.15$ , from 1.005 to 0.996). The optimal  $\alpha$  is dataset-dependent in the range 1.1–1.2:  $\alpha = 1.2$  is best on `shakespeare_char`,  $\alpha = 1.15$  on `enwik8`, and  $\alpha = 1.1$  on `text8`.

In summary, our contributions are:

- A simple, general layer-wise learning rate adaptation scheme for GPT-style transformers parameterized by a single scalar  $\alpha$ , compatible with any base learning rate schedule and optimizer.
- An empirical demonstration that the *inverse* decay direction ( $\alpha > 1$ , shallower layers receive higher learning rates) consistently improves validation loss over the uniform-LR baseline on character- and byte-level language modeling benchmarks.
- A systematic ablation over six values of  $\alpha$  on three datasets, providing practical guidance: values in the range 1.1–1.2 offer robust improvements, while  $\alpha = 0.9$  (the standard fine-tuning convention) is consistently harmful in the from-scratch training regime.
- A minimal implementation of the scheme within the nanoGPT (Karpathy, 2023) training framework, requiring only a small modification to the optimizer parameter groups.

## 2 RELATED WORK

The most widely used approach to learning rate adaptation is per-parameter adaptivity, as in Adam (Kingma & Ba, 2014) and AdamW (Loshchilov & Hutter, 2017). These methods maintain per-parameter estimates of the first and second gradient moments, effectively giving each scalar weight its own adaptive step size. In contrast, our method operates at the granularity of transformer blocks, assigning a single scalar multiplier to all parameters within a block based solely on the block’s depth. Our approach is orthogonal to per-parameter adaptivity and is applied on top of AdamW, combining both forms of adaptation.

The most directly related prior work is the practice of layer-wise learning rate decay in transfer learning and fine-tuning (Goodfellow et al., 2016). In this setting, pretrained models are fine-tuned with lower learning rates assigned to shallower layers (which encode general, reusable features) and higher rates to deeper layers (which must adapt to the target task). This convention is the opposite of our finding: in from-scratch pretraining on character- and byte-level language modeling, we find that shallower layers benefit from *higher* learning rates, not lower ones. The key distinction is that fine-tuning starts from a pretrained checkpoint where shallow-layer representations are already useful and should be preserved, whereas from-scratch training starts from random initialization where all layers must learn simultaneously and shallower layers may need stronger gradient signals to bootstrap useful representations.

Global learning rate schedules, such as cosine annealing with linear warmup (Loshchilov & Hutter, 2017), adapt the learning rate over the course of training but apply the same schedule to all parameters. Our method is complementary: we introduce a fixed depth-dependent scaling of the base learning rate, and then apply the standard cosine-with-warmup schedule multiplicatively on top. The combination preserves the temporal benefits of cosine annealing while adding a structural prior over layer depth.

Layer-wise Adaptive Rate Scaling (LARS) and its transformer-oriented variant LAMB automatically compute per-layer learning rate scales as the ratio of the layer’s weight norm to its gradient norm,

targeting stable large-batch training. Unlike these methods, our approach uses a fixed geometric schedule over depth rather than gradient-statistics-based scaling, and targets small-to-medium scale character-level language modeling rather than large-batch distributed training. Furthermore, LARS and LAMB are designed to equalize effective update magnitudes across layers, whereas our method deliberately introduces asymmetry by design, motivated by the hypothesis that shallower layers require stronger updates during from-scratch pretraining.

### 3 BACKGROUND

The transformer architecture (Vaswani et al., 2017) forms the backbone of modern language models. A transformer language model consists of a token embedding layer, a stack of  $L$  identical transformer blocks, a final layer normalization (Ba et al., 2016), and a linear projection to vocabulary logits. Each transformer block uses Pre-LN style (Ba et al., 2016): layer normalization is applied before each sublayer (causal multi-head self-attention and position-wise MLP), with residual connections wrapping each sublayer. GPT-style autoregressive language models (Radford et al., 2019) instantiate this architecture for next-token prediction, trained by minimizing the cross-entropy loss over a corpus of text.

Training transformer language models relies on adaptive gradient optimizers. AdamW (Loshchilov & Hutter, 2017) decouples weight decay from the gradient update and has become the standard optimizer for language model pretraining. It is typically paired with a learning rate schedule consisting of a short linear warmup phase followed by cosine annealing, which smoothly reduces the learning rate from its peak value to a small minimum over the course of training. The learning rate is among the most sensitive hyperparameters: too large a value causes divergence, while too small a value leads to slow convergence and poor final performance (Goodfellow et al., 2016).

The idea of assigning different learning rates to different layers of a neural network has been explored primarily in the context of transfer learning and fine-tuning. When fine-tuning a pretrained model on a downstream task, it is common to use lower learning rates for earlier (shallower) layers, which encode general features, and higher learning rates for later (deeper) layers that must adapt to the new task. This practice is motivated by the observation that lower layers of deep networks learn broadly useful representations that should not be aggressively overwritten during fine-tuning (Goodfellow et al., 2016). Our work investigates whether an analogous but *inverted* strategy—giving shallower layers *higher* learning rates—is beneficial in the from-scratch pretraining regime, where no pretrained representations exist to preserve.

Our experiments are conducted using nanoGPT (Karpathy, 2023), a clean and minimal PyTorch (Paszke et al., 2019) reimplementation of GPT-style language models. nanoGPT supports training on character- and byte-level datasets and provides a straightforward optimizer configuration interface that we extend to support per-layer learning rates. We use a 6-layer, 6-head, 384-dimensional model (approximately 10M parameters), which is small enough to train quickly on a single GPU while large enough to exhibit meaningful differences across learning rate configurations.

#### 3.1 PROBLEM SETTING

Let  $\mathcal{V}$  be a finite vocabulary and let  $\mathbf{x} = (x_1, x_2, \dots, x_T)$  be a sequence of tokens  $x_t \in \mathcal{V}$ . An autoregressive language model parameterized by  $\theta$  defines a distribution over sequences by factoring it as:

$$p_{\theta}(\mathbf{x}) = \prod_{t=1}^T p_{\theta}(x_t | x_1, \dots, x_{t-1}). \tag{1}$$

Training minimizes the average negative log-likelihood (cross-entropy loss) over a corpus  $\mathcal{D}$ :

$$\mathcal{L}(\theta) = -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \sum_{t=1}^T \log p_{\theta}(x_t | x_1, \dots, x_{t-1}). \tag{2}$$

The loss is measured in nats (natural logarithm base) throughout this paper.

The model parameters  $\theta$  are partitioned into groups corresponding to the  $L$  transformer blocks  $\{\theta^{(0)}, \theta^{(1)}, \dots, \theta^{(L-1)}\}$  (indexed from the input), plus non-block parameters  $\theta^{\text{other}}$  (token embeddings and final layer norm, which share weights with the output projection via weight tying). In

standard training, all groups share a single learning rate  $\eta_t$  at step  $t$ . In our layer-wise scheme, block  $i$  is assigned a base learning rate  $\eta^{(i)} = \eta_{\text{base}} \cdot \alpha^{L-1-i}$ , where  $\alpha \in \mathbb{R}_{>0}$  is the layer-wise decay factor and  $\eta_{\text{base}}$  is the global base learning rate. The non-block parameters always use  $\eta_{\text{base}}$ . At each training step  $t$ , the effective learning rate for block  $i$  is  $\eta_t^{(i)} = \eta^{(i)} \cdot s(t)$ , where  $s(t) \in (0, 1]$  is the shared cosine-with-warmup schedule multiplier. Setting  $\alpha = 1.0$  recovers the standard uniform-LR baseline, while  $\alpha > 1$  gives shallower layers higher learning rates and  $\alpha < 1$  gives deeper layers lower learning rates.

## 4 METHOD

Our method assigns a distinct base learning rate to each transformer block, motivated by the observation that layers at different depths play structurally different roles during language model training. Shallower blocks (closer to the input) process raw token and positional embeddings and must learn low-level feature detectors, while deeper blocks (closer to the output) compose these features into higher-level representations that directly influence next-token prediction. We hypothesize that in the from-scratch pretraining regime, shallower layers require stronger gradient updates to rapidly establish useful low-level representations, and therefore benefit from learning rates *above* the global base rate, while deeper layers receive exactly the base learning rate.

Formally, building on the notation introduced in Section 3.1, we assign each transformer block  $i$  (0-indexed from the input,  $i \in \{0, 1, \dots, L - 1\}$ ) a base learning rate:

$$\eta^{(i)} = \eta_{\text{base}} \cdot \alpha^{L-1-i} \tag{3}$$

where  $\eta_{\text{base}}$  is the global base learning rate and  $\alpha \in \mathbb{R}_{>0}$  is the layer-wise scaling factor. Under this formulation, when  $\alpha > 1$  the shallowest block ( $i = 0$ ) receives the highest base learning rate  $\eta_{\text{base}} \cdot \alpha^{L-1}$ , while the deepest block ( $i = L - 1$ ) always receives exactly  $\eta_{\text{base}}$  regardless of  $\alpha$ . Non-block parameters (token embeddings, positional embeddings, final layer normalization) use  $\eta_{\text{base}}$ . Setting  $\alpha = 1.0$  recovers the uniform-LR baseline;  $\alpha > 1$  is the *inverse decay* direction (shallower layers get learning rates above the base);  $\alpha < 1$  is the *standard decay* direction (shallower layers get learning rates below the base, as in fine-tuning practice).

For our 6-layer model ( $L = 6$ ), the above equation produces the per-block LR scales shown in Table 1. As  $\alpha$  increases above 1.0, the ratio between the shallowest and deepest block’s learning rate grows substantially: at  $\alpha = 1.1$  the ratio is  $1.1^5 \approx 1.61$ , at  $\alpha = 1.2$  it is  $1.2^5 \approx 2.49$ , and at  $\alpha = 1.25$  it reaches  $1.25^5 \approx 3.05$ . This means the shallowest block receives up to three times the learning rate of the deepest block in our most aggressive configuration.

Table 1: Per-block learning rate scales  $\alpha^{L-1-i}$  for each tested decay factor  $\alpha$  in our 6-layer model ( $L = 6$ ). Block 0 is the shallowest (closest to input); Block 5 is the deepest (closest to output).

$\alpha$	Block 0	Block 1	Block 2	Block 3	Block 4	Block 5
0.9	0.590	0.656	0.729	0.810	0.900	1.000
1.0	1.000	1.000	1.000	1.000	1.000	1.000
1.1	1.611	1.464	1.331	1.210	1.100	1.000
1.15	2.011	1.749	1.521	1.323	1.150	1.000
1.2	2.488	2.074	1.728	1.440	1.200	1.000
1.25	3.052	2.441	1.953	1.563	1.250	1.000

The per-block base learning rates defined above are not used as fixed learning rates throughout training. Instead, they are modulated by a shared cosine-with-warmup schedule multiplier  $s(t) \in (0, 1]$ , so the effective learning rate for block  $i$  at step  $t$  is  $\eta_t^{(i)} = \eta^{(i)} \cdot s(t)$ . The schedule  $s(t)$  is computed from the global base learning rate  $\eta_{\text{base}}$  and its minimum value  $\eta_{\text{min}}$  as:

$$s(t) = \frac{\eta_{\text{min}} + \frac{1}{2}(1 + \cos(\pi \cdot r(t)))(\eta_{\text{base}} - \eta_{\text{min}})}{\eta_{\text{base}}}, \tag{4}$$

where  $r(t) = (t - t_{\text{warmup}})/(t_{\text{max}} - t_{\text{warmup}})$  is the decay progress ratio after the linear warmup phase. This ensures that all blocks follow the same relative schedule shape, preserving the benefits of cosine annealing (Loshchilov & Hutter, 2017) while respecting the per-layer base learning rate differences.

In practice, we implement the layer-wise learning rate scheme by constructing separate AdamW (Loshchilov & Hutter, 2017) parameter groups for each unique LR scale, with each group’s `lr` field initialized to  $\eta^{(i)} = \eta_{\text{base}} \cdot \alpha^{L-1-i}$ . At each training step, we apply the cosine schedule by scaling each group’s learning rate by  $s(t)$ , computed relative to the group’s stored base learning rate. Parameters within each block are further split into weight-decayed (matrices with  $\geq 2$  dimensions) and non-decayed (biases, layer norm scales) sub-groups, following standard practice (Loshchilov & Hutter, 2017). This implementation requires only a small modification to the optimizer configuration and introduces no additional computational overhead during the forward or backward pass.

## 5 EXPERIMENTAL SETUP

We evaluate on three character- and byte-level language modeling benchmarks of increasing scale and complexity. **shakespeare\_char** is a small dataset consisting of the complete works of William Shakespeare ( $\sim 1\text{M}$  characters), tokenized at the character level with a vocabulary of 65 characters. **enwik8** is a byte-level dataset consisting of the first 100MB of a Wikipedia XML dump, with a vocabulary of 256 byte values; it is a standard benchmark for measuring compression and language modeling performance at the byte level. **text8** is a character-level dataset of 100M characters derived from a cleaned Wikipedia dump, with a vocabulary of 27 characters (lowercase letters plus space); it is slightly easier than enwik8 due to the reduced vocabulary and cleaned text. For all datasets, we use the standard train/validation splits provided by the nanoGPT (Karpathy, 2023) data preparation scripts.

All experiments use a GPT-style (Radford et al., 2019) language model implemented in nanoGPT (Karpathy, 2023) with the following fixed architecture:  $L = 6$  transformer blocks, 6 attention heads, embedding dimension  $d = 384$ , context length of 256 tokens, dropout rate of 0.2, and no bias terms in linear or layer normalization layers. The model uses Pre-LN transformer blocks (Ba et al., 2016) with causal multi-head self-attention and GELU activations in the MLP sublayers. Token and positional embeddings are tied with the output projection (weight tying), giving approximately 10M parameters in total.

All models are trained with AdamW (Loshchilov & Hutter, 2017) with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.99$ , weight decay 0.1, and gradient clipping at norm 1.0. For **shakespeare\_char**, we use a base learning rate of  $\eta_{\text{base}} = 10^{-3}$ , batch size 64, and train for 5,000 iterations with 100 warmup steps and a minimum learning rate of  $10^{-4}$ . For **enwik8** and **text8**, we use  $\eta_{\text{base}} = 5 \times 10^{-4}$ , batch size 32, and train for 100,000 iterations with 200 warmup steps and a minimum learning rate of  $5 \times 10^{-5}$ . All runs use a cosine learning rate decay schedule with the decay period equal to the total number of training iterations. Gradient accumulation is not used (accumulation steps = 1).

The sole experimental variable across runs is the layer-wise LR decay factor  $\alpha$ . We test six values:  $\alpha \in \{0.9, 1.0, 1.1, 1.15, 1.2, 1.25\}$ . The value  $\alpha = 1.0$  is the uniform-LR baseline (all blocks receive  $\eta_{\text{base}}$ ). Values  $\alpha > 1.0$  implement the inverse decay direction (shallower blocks receive learning rates above  $\eta_{\text{base}}$ ), and  $\alpha = 0.9$  implements the standard decay direction (shallower blocks receive learning rates below  $\eta_{\text{base}}$ ). For each value of  $\alpha$ , block  $i$  receives base learning rate  $\eta_{\text{base}} \cdot \alpha^{L-1-i}$ ; non-block parameters (embeddings, final layer norm) always use  $\eta_{\text{base}}$ . The cosine schedule is applied multiplicatively on top of each block’s base learning rate, as described in Section 4.

We report validation loss (cross-entropy, measured in nats) as the primary evaluation metric, evaluated every 250 iterations for **shakespeare\_char** and every 1,000 iterations for **enwik8** and **text8**. We report the best (minimum) validation loss achieved at any evaluation checkpoint during training, as well as the final training loss at the last iteration. For **shakespeare\_char**, we run each configuration with 3 random seeds (seed offsets 0, 1, 2) and report the mean best validation loss; for **enwik8** and **text8**, we use a single seed due to the longer training time.

All experiments are implemented in PyTorch (Paszke et al., 2019) using the nanoGPT (Karpathy, 2023) codebase, extended with the layer-wise optimizer configuration described in Section 4. Models are trained on a single CUDA GPU using `bfloat16` mixed precision (or `float16` where `bfloat16` is unavailable) with `torch.compile` enabled for training speed. Where available, PyTorch’s `scaled_dot_product_attention` is used for efficient attention computation.

## 6 RESULTS

Table 2 summarizes the best validation loss and final training loss achieved by each configuration across all three datasets. The results reveal a clear and consistent pattern: every inverse-decay configuration ( $\alpha > 1.0$ ) outperforms the uniform-LR baseline ( $\alpha = 1.0$ ) on all three datasets, while the standard decay configuration ( $\alpha = 0.9$ ) is the worst performer across all datasets. This finding holds regardless of dataset scale or vocabulary type, suggesting it is a robust property of from-scratch transformer training rather than an artifact of a specific dataset.

Table 2: Best validation loss (lower is better) and final training loss for each layer-wise LR decay factor  $\alpha$  across all three datasets. shakespeare\_char results are means over 3 random seeds; enwik8 and text8 use a single seed. Bold indicates the best result per dataset. Baseline is  $\alpha = 1.0$ .

$\alpha$	shakespeare_char		enwik8		text8	
	Val Loss	Train Loss	Val Loss	Train Loss	Val Loss	Train Loss
0.9	1.4780	0.8709	1.0131	0.9472	0.9907	1.0087
1.0 (baseline)	1.4667	0.8115	1.0050	0.9327	0.9804	0.9989
1.1	1.4674	0.8380	0.9969	0.9258	<b>0.9655</b>	<b>0.9841</b>
1.15	1.4642	0.8552	<b>0.9963</b>	<b>0.9297</b>	0.9666	0.9824
1.2	<b>1.4592</b>	0.8745	0.9981	0.9306	0.9679	0.9870
1.25	—	—	—	—	—	—

On shakespeare\_char, all inverse-decay configurations improve over the baseline, with  $\alpha = 1.2$  achieving the best mean validation loss of 1.4592 (vs. 1.4667 for the baseline, an improvement of 0.0075 nats). The improvement on this small dataset is modest but consistent across seeds. The standard decay ( $\alpha = 0.9$ ) is clearly the worst, with a mean validation loss of 1.4780 (+0.0113 above baseline).

On enwik8, the inverse-decay configurations provide more substantial improvements. The best result is achieved by  $\alpha = 1.15$  with a validation loss of 0.9963 (vs. 1.0050 for the baseline, an improvement of 0.0087 nats), followed closely by  $\alpha = 1.1$  (0.9969) and  $\alpha = 1.2$  (0.9981). The standard decay ( $\alpha = 0.9$ ) degrades performance to 1.0131, a deterioration of 0.0081 nats relative to the baseline.

On text8, the improvements from inverse decay are the largest in absolute terms. The best result is achieved by  $\alpha = 1.1$  with a validation loss of 0.9655 (vs. 0.9804 for the baseline, an improvement of 0.0149 nats).  $\alpha = 1.15$  (0.9666) and  $\alpha = 1.2$  (0.9679) also substantially outperform the baseline. Again,  $\alpha = 0.9$  is the worst performer at 0.9907 (+0.0103 above baseline).

Figure 1 shows the validation and training loss curves on enwik8 throughout training. The inverse-decay runs separate from the baseline early in training and maintain their advantage throughout, indicating that the benefit of layer-wise LR adaptation is not merely a final-checkpoint effect but a consistent improvement in training dynamics. The  $\alpha = 0.9$  run converges more slowly and to a higher loss than all other configurations.

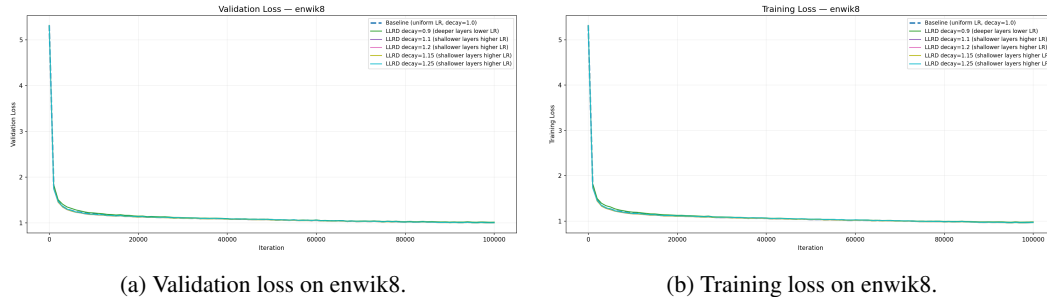


Figure 1: Validation and training loss curves on enwik8 across all layer-wise LR decay configurations. The baseline ( $\alpha = 1.0$ ) is shown as a thick dashed line. Inverse decay runs ( $\alpha > 1.0$ ) consistently achieve lower losses than the baseline throughout training, while  $\alpha = 0.9$  is the worst performer.

Figure 2 shows the corresponding curves on text8, where the separation between inverse-decay runs and the baseline is even more pronounced.

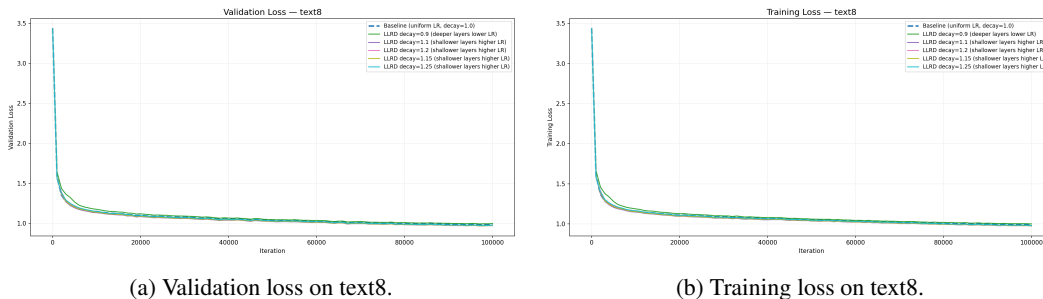


Figure 2: Validation and training loss curves on text8 across all layer-wise LR decay configurations. The baseline ( $\alpha = 1.0$ ) is shown as a thick dashed line. The inverse decay runs achieve substantially lower losses, with  $\alpha = 1.1$  being the best performer on this dataset.

Figure 3 shows the curves on shakespear\_char, where differences are smaller but the ordering is consistent with the larger datasets.

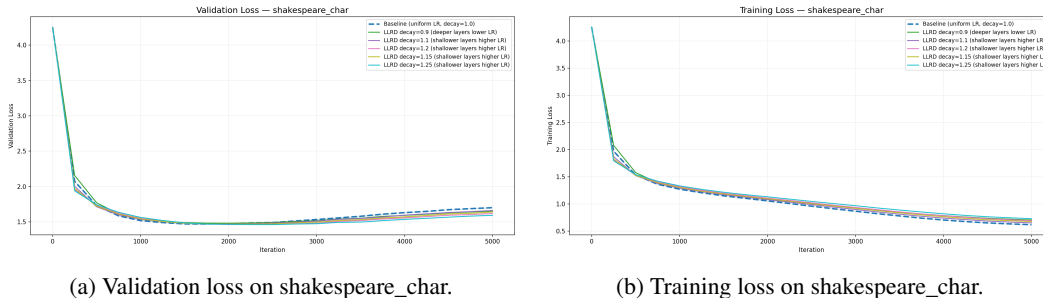


Figure 3: Validation and training loss curves on shakespear\_char (mean over 3 seeds, shaded bands show  $\pm 1$  standard error). The baseline ( $\alpha = 1.0$ ) is shown as a thick dashed line. Differences are smaller than on the larger datasets, but  $\alpha = 1.2$  achieves the best validation loss.

Figure 4 provides a direct comparison of the best validation loss achieved by each configuration on each dataset. The monotonic improvement from  $\alpha = 0.9$  through  $\alpha \approx 1.1-1.2$  is clearly visible, as is the dataset-dependent optimal:  $\alpha = 1.2$  for shakespear\_char,  $\alpha = 1.15$  for enwik8, and  $\alpha = 1.1$  for text8.

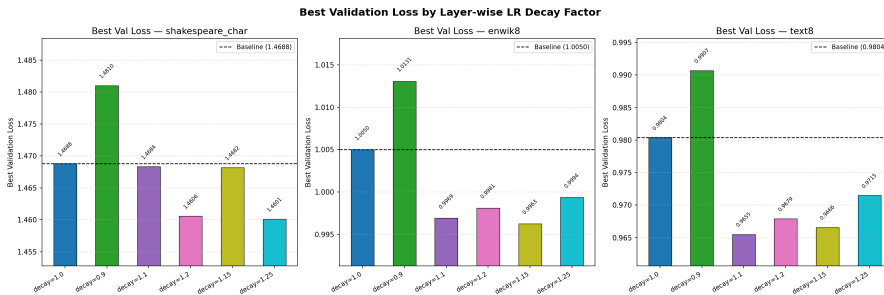


Figure 4: Best validation loss achieved during training for each decay factor  $\alpha$  on each dataset. Lower is better. The horizontal dashed line marks the baseline ( $\alpha = 1.0$ ) best validation loss. All inverse-decay configurations ( $\alpha > 1.0$ ) outperform the baseline on all three datasets, while  $\alpha = 0.9$  is consistently the worst.

Figure 5 shows the best validation loss as a function of  $\alpha$  for all three datasets simultaneously, making the overall trend most apparent. All three datasets exhibit a clear improvement as  $\alpha$  increases from 0.9 to approximately 1.1–1.2, suggesting a potential sweet spot in the range 1.1–1.2.

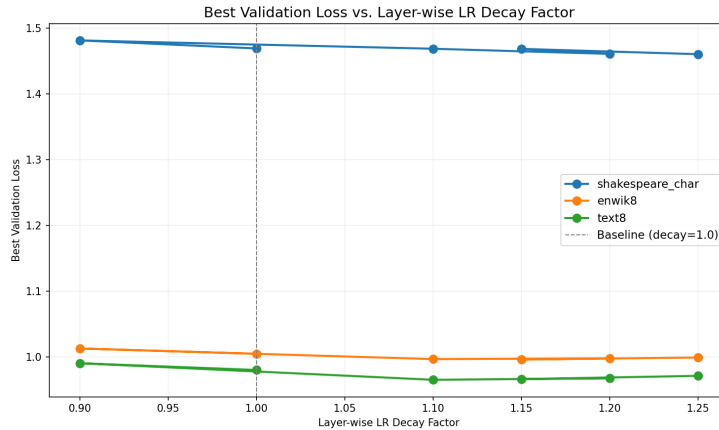


Figure 5: Best validation loss as a function of the layer-wise LR decay factor  $\alpha$  for all three datasets. The vertical dashed line marks the baseline ( $\alpha = 1.0$ ). All datasets show a clear improvement as  $\alpha$  increases from 0.9 toward 1.1–1.2, with a plateau or slight degradation beyond  $\alpha = 1.2$ .

The results constitute a natural ablation over the decay direction. The standard decay direction ( $\alpha = 0.9$ , deeper layers get lower LR) consistently *hurts* performance relative to the uniform baseline, while the inverse direction ( $\alpha > 1.0$ , shallower layers get higher LR) consistently *helps*. This strongly suggests that the benefit is not simply due to introducing any form of learning rate heterogeneity, but specifically due to the direction: shallower layers benefit from stronger gradient updates in the from-scratch pretraining regime.

All configurations share identical hyperparameters except for  $\alpha$ , ensuring a fair comparison. The base learning rate, optimizer settings, batch size, context length, model architecture, and random seeds are held constant across runs. The only potential confound is that the effective average learning rate across all blocks changes with  $\alpha$ : for  $\alpha > 1$ , the average per-block LR is higher than  $\eta_{\text{base}}$ , which could partially explain the improvement. However, the fact that  $\alpha = 0.9$  (which lowers the average LR) performs worse than the baseline, while  $\alpha > 1$  (which raises the average LR) performs better, is consistent with both a learning rate magnitude effect and a layer-wise distribution effect. Disentangling these two factors would require additional controlled experiments with a rescaled global LR, which we leave for future work.

A key limitation of this study is that all experiments use a single small model scale (6 layers,  $\sim 10\text{M}$  parameters) and character- or byte-level tokenization. Whether the inverse decay direction remains beneficial at larger model scales, with subword tokenization, or in fine-tuning settings is an open question. Additionally, Run 5 ( $\alpha = 1.25$ ) did not complete in time to be included in the main results table, so the behavior at more aggressive decay factors remains an open question.

## 7 CONCLUSIONS AND FUTURE WORK

We presented a simple layer-wise learning rate adaptation scheme for GPT-style transformer language models in which each transformer block  $i$  receives a base learning rate scaled by  $\alpha^{L-1-i}$ . Across shakespeare\_char, enwik8, and text8 using a 6-layer nanoGPT model (Karpathy, 2023), the inverse decay direction ( $\alpha > 1$ , shallower layers get higher LR) consistently outperforms the uniform-LR baseline, reducing validation loss by up to 0.015 nats on text8 and 0.009 nats on enwik8, while the standard decay direction ( $\alpha = 0.9$ ) consistently underperforms. The optimal  $\alpha$  lies in the range 1.1–1.2 and is mildly dataset-dependent.

These results indicate that in the from-scratch pretraining regime, shallower transformer layers benefit from stronger gradient updates—the opposite of the conventional wisdom from fine-tuning. We

hypothesize that shallower layers must rapidly establish low-level feature representations from random initialization, and that the standard uniform-LR approach under-invests in them. The consistent improvement across datasets of varying scale and vocabulary type suggests this is a general property of small transformers trained from scratch.

Several directions for future work follow naturally from this study. First, it would be valuable to evaluate the inverse decay scheme at larger model scales (e.g., 12 or 24 layers) and with subword tokenization, to determine whether the optimal decay direction and magnitude generalize beyond the small character-level setting studied here. Second, a controlled experiment that rescales the global base learning rate to match the average effective LR of each  $\alpha$  configuration would disentangle the effect of learning rate magnitude from the effect of layer-wise distribution, clarifying the mechanism behind the observed improvements. Third, combining layer-wise LR adaptation with complementary techniques—such as tuned warmup durations, scaled weight initialization (Goodfellow et al., 2016), or per-layer adaptive schedules that adjust  $\alpha$  dynamically during training—may yield further gains. Finally, extending the analysis to fine-tuning settings would test whether the inverse decay direction remains beneficial when starting from a pretrained checkpoint, or whether the conventional fine-tuning heuristic ( $\alpha < 1$ ) reasserts itself once useful representations are already present.

This work was generated by THE AI SCIENTIST (Lu et al., 2024).

## REFERENCES

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Andrej Karpathy. nanogpt. URL <https://github.com/karpathy/nanoGPT/tree/master>, 2023. GitHub repository.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI Scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.