

Does Wikipedia Make Code Better?

A Load Balancer Case Study

Timothy Brinded

February 2026

Abstract

We test whether augmenting an LLM coding agent with access to Wikipedia improves the quality of generated code. A single design-rich programming problem—a graceful-degradation load balancer—is implemented by Claude Code under eight conditions: one control and seven Wikipedia-augmented variants with different research strategies. Solutions are evaluated by a deterministic behavioral benchmark (7 fault scenarios, weighted 0–100) and a blinded LLM judge (design quality). Two conditions separate from the pack: the *flaneur* (95.24) and *biomimetic* (94.77) conditions score ~ 6 points above control (89.39), while the remaining five Wikipedia conditions cluster within 0.6 points of control. However, the experiment’s apparatus reveals as much as its results: the *flaneur* did not follow its random-walk instructions—it searched on-topic articles instead—yet scored highest. The blinded judge and benchmark diverge *dramatically* on this condition: the judge ranks it below control (−6 margin) while the benchmark ranks it first. Biomimetic is the only condition that followed its cross-domain research instructions, produced a structurally distinct architecture, and scored well. We discuss the apparatus limitations these findings expose, including research-strategy enforcement, parametric knowledge contamination, and the compute-time confound.

1 Introduction

Large language models generate competent code from training data alone. But can cross-domain knowledge—the kind found in encyclopedic references—improve the *design* of generated software? We test a simple hypothesis: giving a coding

agent access to Wikipedia articles alongside a programming problem produces better-designed code than coding without external references.

We study a single problem: implementing a load balancer with graceful degradation. Load balancers are design-rich systems where multiple valid architectures exist, failure modes are subtle, and cross-domain analogies (TCP congestion control, biological homeostasis, power grid protection) can yield structural insights.

Eight conditions are evaluated. The **control** receives only the problem specification. Seven Wikipedia-augmented conditions each use a different research strategy: direct research (*explicit*), ambient exposure (*subtle*), historical precedent (*reflective*), undirected exploration (*flaneur*), cross-domain convergence (*consilience*), biological analogy (*biomimetic*), and adversarial critique (*contrarian*). Table 1 summarizes all eight conditions.

The results are mixed—and more interesting for their failures than their successes. Two conditions clearly outperform control, but the experiment’s apparatus reveals gaps that complicate causal claims. We present both the findings and an honest accounting of what the experimental design can and cannot tell us.

2 Method

2.1 Three-Stage Pipeline

Each Wikipedia condition runs through a three-stage pipeline:

Stage 1: Retrieval. A research subagent (**wiki-explorer**) explores a local corpus of ~ 7 million Wikipedia articles using the same tools available for code—`ripgrep` for cross-article full-text search, `file` reads for deep inspection. The sub-

Table 1: Experimental conditions and research strategies.

Condition	Research Strategy
control	No Wikipedia access
explicit	Instructed to research Wikipedia
subtle	Wikipedia tools available, not mentioned
reflective	Seek historical precedent and past failures
flaneur	Undirected random Wikipedia exploration
consilience	Find convergent patterns across 5+ domains
biomimetic	Research only biological analogies
contrarian	Stress-test the obvious approach first

agent’s system prompt encodes the research persona (e.g., “take a random walk” for flaneur, “find biological analogies” for biomimetic). Retrieval runs until the subagent terminates or a timeout is reached.

Stage 2: Synthesis. The research findings are passed to a synthesis step that distills the retrieved articles into a structured research brief. This brief accompanies the problem specification in the next stage.

Stage 3: Coding. A fresh Claude Code session receives the problem specification, the research brief (if any), and implements the solution. The coding agent has no direct access to Wikipedia—only the synthesized brief. It implements an `AbstractLoadBalancer` base class with four required methods: `on_backend_added`, `route_request`, `on_request_complete`, and `on_tick`. Backend health is inferred through an opaque `BackendHandle` protocol—the agent cannot query health directly, only observe request outcomes and probe results.

The control condition skips stages 1 and 2, receiving only the problem specification.

2.2 Conditions

The **control** receives the problem specification and nothing else. All seven Wikipedia conditions receive the same specification plus a research brief produced by the retrieval and synthesis stages. The conditions

differ only in the system prompt framing given to the research subagent:

- **Explicit:** “Research Wikipedia for load balancer design patterns.”
- **Subtle:** Tools available silently; no mention in prompt.
- **Reflective:** “Find historical failures that inform resilient design.”
- **Flaneur:** “Take a random walk through Wikipedia; let connections emerge.”
- **Consilience:** “Find the same pattern in 5+ unrelated domains.”
- **Biomimetic:** “Research only biological systems as design analogies.”
- **Contrarian:** “Before coding, find reasons the obvious approach fails.”

2.3 Evaluation

Solutions are evaluated by three independent methods.

Deterministic Benchmark. Seven fault-injection scenarios run against each implementation, producing scores on $[0, 1]$ per scenario. An aggregate score is computed as the weighted average scaled to $[0, 100]$.

Blinded LLM Judge. Each Wikipedia condition is paired against control in a blinded A/B comparison (randomized assignment to “Solution A” and “Solution B”). The judge scores six dimensions: correctness ($\times 3$), design ($\times 2$), robustness ($\times 2$), algorithmic novelty ($\times 2$), cross-domain insight ($\times 1$), and proportionality ($\times 1$). Maximum weighted total: 110 points.

Quantitative Metrics. Duration, API cost, lines of code, and subagent token usage are recorded for each run.

2.4 Benchmark Design

Table 2 lists the seven benchmark scenarios and their weights. Priority Protection carries the highest weight (3.0) because correctly shedding low-priority traffic during overload is the core requirement. Degradation Detection (2.0) tests whether the agent can sense a backend that is slowly failing—a subtle scenario that separates sophisticated health tracking from naive approaches.

Table 2: Benchmark scenarios and weights.

Scenario	Tests	Weight
Steady State	Even distribution, full success	1.0
Degradation	Detect slowly failing backend	2.0
Priority	Shed low-priority traffic first	3.0
Recovery	Smooth ramp-up after failure	2.0
Cascading	Isolate failure, protect healthy	2.0
Flapping	Handle on/off oscillation	1.5
Asymmetric	Route by backend speed	1.0

2.5 Apparatus Controls and Limitations

Before presenting results, we document what this experiment controls for and what it does not.

What we controlled.

- *Retrieval verification*: each research subagent produces a manifest of articles read from disk, confirming that Wikipedia content was actually consumed.
- *Comment stripping*: all code is evaluated by the benchmark as-is; comments and docstrings do not affect scores.
- *Position randomization*: the judge sees solutions in randomized A/B order, with no condition labels.
- *Identical contract*: all conditions implement the same `AbstractLoadBalancer` interface against the same `BackendHandle` protocol.

What we did not control.

- *Parametric knowledge contamination*: the coding agent (Claude) has been trained on load-balancing literature. Wikipedia may confirm or refine existing knowledge rather than introduce genuinely new information. We cannot distinguish “the agent learned this from Wikipedia” from “the agent already knew this and Wikipedia reinforced it.”
- *Research strategy enforcement*: the retrieval manifests verify *which* articles were read, not *how* they were found. An agent instructed to random-walk can instead search directly for relevant articles.

As we will see, this is exactly what happened with the *flaneur* condition.

- *Compute-time confound*: research conditions consume additional wall-clock time and API tokens during retrieval and synthesis. A condition that spends 1,200 seconds on research and 500 seconds coding has had more total “thinking time” than control, which spent 468 seconds coding. We cannot fully separate the effect of Wikipedia content from the effect of additional compute.
- *The subtle failure*: the subtle condition was designed to test whether an agent would organically discover and use Wikipedia tools. It did not—making it effectively a second control with lower performance, not a meaningful test of ambient exposure.

3 Results

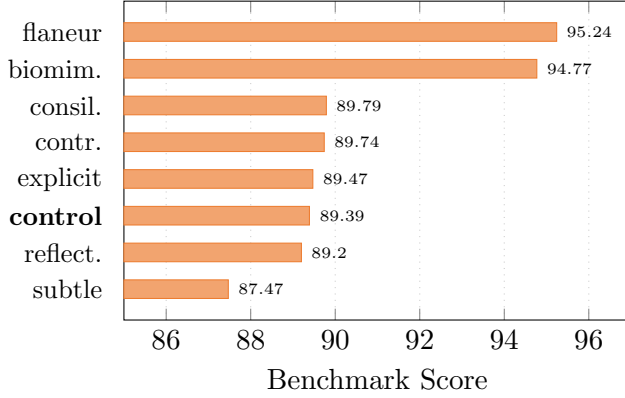
3.1 Benchmark Rankings

Figure 1a shows aggregate benchmark scores ranked by performance. Two conditions clearly separate from the pack: *flaneur* (95.24) and *biomimetic* (94.77) score ~ 5 –6 points above control (89.39). The remaining five Wikipedia conditions cluster within 0.6 points of control, forming an undifferentiated middle band from 89.20 (reflective) to 89.79 (consilience). Subtle (87.47) is the only condition that clearly underperforms control.

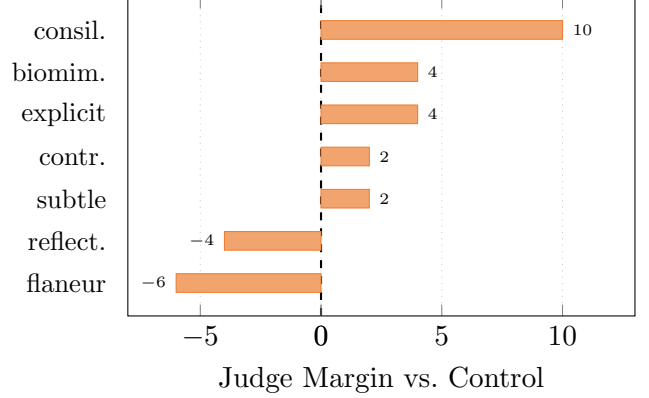
3.2 Scenario Heatmap

Table 3 shows per-scenario scores for all conditions, sorted by aggregate rank. The top two conditions dominate the same scenarios: cascading failure (*flaneur* 1.00, *biomimetic* 0.98) and asymmetric routing (*flaneur* 1.00, *biomimetic* 0.99). These are the scenarios where the middle pack struggles most—control scores 0.87 and 0.85 respectively. Degradation detection also separates the leaders (*flaneur* 1.00, *biomimetic* 0.99) from control (0.85). Recovery smoothness is the most uniform scenario, with all conditions scoring 0.80–0.85.

Two conditions show notably weak asymmetric routing: consilience (0.52) and explicit (0.53). Both use weighted-random selection that distributes load roughly evenly regardless of backend speed—a failure mode that costs them on this scenario but is masked by strong performance elsewhere.



(a) Aggregate benchmark scores (0–100). Control baseline in bold.



(b) Blinded LLM judge margin vs. control (out of 110). Positive = Wikipedia condition wins.

Figure 1: Benchmark scores (left) and judge margins (right). The two measures diverge sharply: flaneur ranks 1st on benchmarks but loses to control on the judge (−6). Consilience ranks 3rd on benchmarks but wins the largest judge margin (+10).

Table 3: Per-scenario benchmark scores (0.0–1.0), sorted by aggregate rank. Color key: ≥ 0.90 (green), 0.80–0.89 (light green), 0.70–0.79 (yellow), 0.50–0.69 (orange), < 0.50 (red).

	Steady (w = 1.0)	Degrade (w = 2.0)	Priority (w = 3.0)	Recovery (w = 2.0)	Cascade (w = 2.0)	Flapping (w = 1.5)	Asymm. (w = 1.0)	Agg.
flaneur	0.95	1.00	1.00	0.85	1.00	0.84	1.00	95.24
biomimetic	0.97	0.99	1.00	0.84	0.98	0.84	0.99	94.77
consilience	0.97	0.96	0.96	0.81	0.96	0.91	0.52	89.79
contrarian	0.95	0.94	0.99	0.83	0.88	0.82	0.79	89.74
explicit	0.97	0.92	0.99	0.83	0.93	0.90	0.53	89.47
control	1.00	0.85	0.96	0.83	0.87	0.90	0.85	89.39
reflective	0.96	0.90	0.99	0.83	0.86	0.83	0.80	89.20
subtle	1.00	0.83	0.98	0.80	0.82	0.87	0.77	87.47

3.3 Judge Rankings

Figure 1b shows the blinded LLM judge margin for each condition versus control. Consilience achieves the largest margin (+10 points), with the judge praising its AIMD recovery model, weighted-random selection for proportional degradation, and data-driven shed tiers. Biomimetic and explicit tie at +4. Contrarian and subtle each win by +2.

Two conditions *lose* to control: reflective (−4) and flaneur (−6). The judge found real code-quality issues in flaneur—a mislabeled routing function and a no-op recovery transition—that do not affect benchmark scores but reflect lower design clarity. Control scored between 60 and 68 across its seven pairings, introducing meaningful noise in the mar-

gins.

3.4 Judge vs. Benchmark Divergence

The judge and benchmark diverge more dramatically than in previous work. The most striking case:

- **Flaneur:** Benchmark rank 1st (95.24); Judge margin −6 (control wins). The benchmark rewards flaneur’s near-perfect cascading failure isolation (1.00) and asymmetric routing (1.00). The judge penalizes code-quality issues and awards no cross-domain insight credit—because flaneur’s research was on-topic, not cross-domain.
- **Consilience:** Benchmark rank 3rd (89.79, within 0.4 of control); Judge rank 1st (+10 margin). The judge praised AIMD recovery and weighted-

random selection. The benchmark shows consilience is dragged down by poor asymmetric routing (0.52)—its weighted-random selection ignores backend speed.

- **Subtle:** Benchmark rank 8th (87.47); Judge margin +2. The judge slightly preferred subtle’s smooth weighted round-robin, but the benchmark penalizes its slower degradation detection (0.83) and weaker cascading failure handling (0.82).

The judge evaluates architecture and design reasoning. The benchmark evaluates whether the code works under fault injection. Each captures a different dimension of quality, and they can disagree sharply on the same implementation.

3.5 Cost and Efficiency

Figure 2a plots benchmark score against API cost. Table 4 shows the full breakdown. Subtle is the cheapest overall (\$0.68, 3:03) but scores lowest. Control is the most cost-efficient at \$1.43 for 89.39 points (62.5 score/\$). Among Wikipedia conditions, flaneur achieves the best return—\$1.94 for 95.24 points (49.1 score/\$), the cheapest research condition and the highest score. Biomimetic costs \$3.96 for 94.77 points. Consilience is the most expensive at \$4.01 for only 0.4 points above control.

Research time dominates cost for most conditions. Flaneur’s research completed in 426 seconds with 37 retrieval turns—the fastest research condition. Consilience’s research took 1,275 seconds with 79 turns and 769 seconds of synthesis, reflecting its instruction to find convergent patterns across five or more domains.

Table 4: Cost, efficiency, and code metrics. Score/\$ = benchmark score per dollar of API cost. Research time includes retrieval and synthesis.

Condition	Time	Cost	LOC	Res.	Sc./\$
subtle	3:03	\$0.68	407	—	128.6
control	7:48	\$1.43	414	—	62.5
flaneur	9:45	\$1.94	446	7:06	49.1
reflective	12:06	\$2.55	428	8:42	35.0
explicit	13:52	\$2.65	543	8:53	33.8
contrarian	14:02	\$3.05	530	10:11	29.4
biomimetic	20:14	\$3.96	728	11:58	23.9
consilience	31:09	\$4.01	437	21:15	22.4

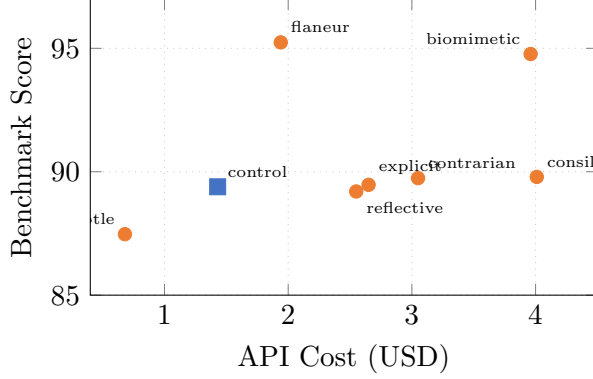
3.6 Wikipedia Usage

Table 5 shows research effort versus knowledge extraction. Six conditions actively retrieved Wikipedia articles; control and subtle did not. All six active conditions read 9–12 articles each—a narrow range that suggests the retrieval pipeline converges on a similar volume regardless of persona.

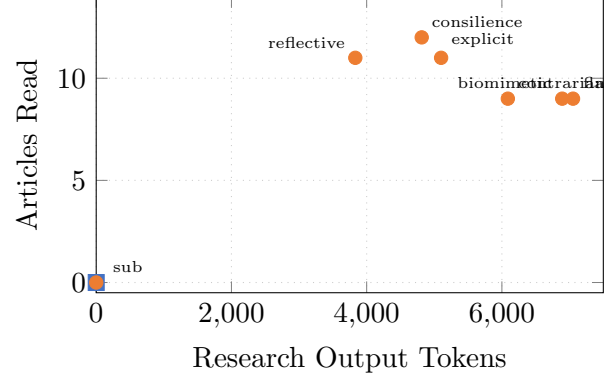
The *content* of the articles varies dramatically. Biomimetic read exclusively biological articles (allostasis, homeostasis, apoptosis, torpor, optimal foraging theory). Contrarian focused on failure modes (systemic risk, demand response, C10k problem). But flaneur—instructed to random-walk—read fault tolerance, thundering herd problem, consistent hashing, weighted round-robin, and deficit round-robin. Every article is directly relevant to load balancing. The flaneur *defected from its research persona*, searching on-topic articles instead of exploring randomly (see Section 5).

Table 5: Wikipedia usage per condition. “Tokens” = research subagent output tokens; “Articles” = articles read from retrieval manifest.

Cond.	Tokens	Art.	Key Sources
control	0	0	—
subtle	0	0	Tools not discovered
explicit	5,101	11	Load shedding, DiffServ
reflective	3,832	11	Circuit breaker, C10k
flaneur	7,052	9	Fault tolerance, WRR
consilience	4,814	12	TCP cong., HRA
biomimetic	6,892	9	Allostasis, torpor
contrarian	6,088	9	Systemic risk, C10k



(a) Benchmark score vs. API cost. Control (blue square); flaneur delivers the highest score at $1.4\times$ control’s cost.



(b) Research effort vs. knowledge extraction. All active conditions read 9–12 articles; token usage varies $2\times$.

Figure 2: Cost–performance trade-off (left) and Wikipedia research patterns (right).

4 What Each Agent Built

Table 6 summarizes the architectural choices made by each condition. Despite identical problem specifications, the eight agents produced meaningfully different designs.

Control (414 LOC) implements weighted least-connections routing with EMA health tracking ($\alpha=0.4$). Priority-based shedding activates at fixed capacity thresholds (background $<70\%$, normal $<40\%$, critical $<15\%$). An asymmetric weight ramp ($+0.15/\text{tick up}$, $-0.4/\text{tick down}$) prevents thundering herd on recovery. A solid, conventional baseline.

Flaneur (446 LOC) uses weighted round-robin with a scoring function: $\text{weight}/(1+\text{avg.latency}/100)$, normalizing latency relative to a 100ms baseline. Failed requests inject a synthetic 500ms latency penalty into the EMA, making health tracking highly responsive. Its “no-volt release” recovery—fast ramp-down ($\times 0.5$) and slow ramp-up ($+0.10/\text{tick}$)—mirrors power-grid demand-response patterns. The implementation excels at cascading failure (1.00) and asymmetric routing (1.00) because latency-weighted scoring naturally routes away from slow backends.

Biomimetic (728 LOC) is the most architecturally distinctive. Routing uses optimal foraging theory: each backend’s “profitability” is health/latency (energy per handling

time), and the router selects from the top-50% profitability tier. A four-state machine (HEALTHY \rightarrow DEGRADED \rightarrow TORPID \rightarrow RECOVERING) governs backend lifecycle—torpid backends receive probes only, and recovery caps requests at $\max(1, \text{round}(\text{progress}\times 5))$. This is the only condition where biological analogies produced structural—not merely terminological—differences in the architecture.

Consilience (437 LOC) implements AIMD-based weighted random selection. Multiplicative decrease ($\times 0.5$) on failure, additive increase ($+0.05/\text{tick}$) on recovery, but recovery is gated through probes: only probe outcomes trigger additive increase, rate-limiting recovery to one tick at a time. Despite referencing five convergent patterns (TCP congestion, power-grid UFLS, ALOHA, high-redundancy actuation, DiffServ), its asymmetric routing fails (0.52)—weighted-random distributes load evenly regardless of backend speed.

Contrarian (530 LOC) combines weighted-random routing with active-request tracking (a least-connections variant). Multi-tier latency sensitivity distinguishes “degraded” (300ms) from “severe” (800ms) with different multiplicative decreases. Its multi-stage shedding escalates through four thresholds. The adversarial research strategy produced careful fault detection but no structural innovation.

Explicit (543 LOC) implements weighted-random selection with graduated continuous health

Table 6: Architectural comparison across conditions. WLC = weighted least-connections; WRR = weighted round-robin; SWRR = smooth weighted round-robin; AIMD = additive increase / multiplicative decrease; EMA/EWMA = (exponential) weighted moving average.

Condition	Routing	Health Tracking	Shedding	Recovery	LOC
control	WLC	EMA ($\alpha=0.4$)	Priority tiers	Asym. ramp (+0.15/−0.4)	414
flaneur	Weighted RR	EWMA + synthetic latency	Demand-response	No-volt release ($\times 0.5/+0.10$)	446
biomimetic	Foraging (E/h)	Trail pheromone EMA	Allostatic redistrib.	4-state torpor machine	728
consilience	Weighted random	AIMD weight	Tiered escalation	Probe-gated (+0.05)	437
contrarian	Weighted LC	EMA + latency tiers	Multi-stage UFLS	AIMD ($\times 0.5/+0.05$)	530
explicit	Weighted random	Graduated continuous	Capacity-coupled	Linear (+0.05)	543
reflective	SWRR (Nginx)	EMA + sliding window	Priority tiers	Asym. ramp (+0.15/−0.3)	428
subtle	WLC	Graduated continuous	Staged tiers	Asym. ramp (+0.15/−0.4)	407

scoring—latency is linearly interpolated between 50ms (healthy) and 500ms (failed), avoiding binary up/down decisions. Linear recovery (+0.05/tick) avoids oscillation in small- N deployments.

Reflective (428 LOC) uses smooth weighted round-robin (Nginx’s SWRR algorithm), the only condition with deterministic routing. A three-state enum (HEALTHY/DEGRADED/UNHEALTHY) with dual-threshold transitions provides clean state management.

Subtle (407 LOC) is architecturally near-identical to control: weighted least-connections with graduated health scoring and staged shedding. It never discovered the Wikipedia tools during its run, confirming the retrieval manifest’s zero-article record.

5 Discussion

Six observations emerge from this experiment.

1. Cross-domain research can improve functional benchmarks by ~ 6 points—but the effect is concentrated. Flaneur and biomimetic score 5–6 points above control, but the gap comes almost entirely from three scenarios: cascading failure, asymmetric routing, and degradation detection. On recovery, priority protection, and flapping, they perform comparably to the middle pack. The improvement is real but narrow—specific failure scenarios where latency-aware or profitability-based routing excels, not a general quality lift.

2. The flaneur defected from its instructions—and this is the most interesting finding. The flaneur was instructed

to “take a random walk through Wikipedia; let connections emerge.” Its retrieval manifest shows it read: fault tolerance, thundering herd problem, exponential backoff, consistent hashing, weighted round-robin. Every article is directly relevant to load balancing. The agent ignored its creative persona and searched for what it already knew was relevant.

This is simultaneously an *apparatus flaw* and a *finding about agent behavior*. As an apparatus flaw: we cannot claim the flaneur condition tests undirected exploration, because undirected exploration didn’t happen. As a behavioral finding: when an LLM agent is given a creative research instruction but knows the problem domain, it gravitates toward on-topic material. The retrieval verification catches *whether* articles were read, not *how* they were found. Enforcing research strategy—as opposed to research access—is an open problem for agent experiment design.

3. Biomimetic is the only genuine cross-domain success. Of the seven Wikipedia conditions, biomimetic is the only one that (a) followed its cross-domain research instructions (reading articles on allostasis, homeostasis, apoptosis, torpor, and optimal foraging theory), (b) produced a structurally distinct architecture (four-state torpor machine, foraging-profitability routing), and (c) scored well on benchmarks (94.77, rank 2). However, it is also the most expensive (\$3.96, $2.8\times$ control) and produced the most code (728 LOC vs. 414 control). The structural translation from biological concepts to software design appears to require significant

additional compute.

4. The judge and benchmark measure orthogonal things. The flaneur divergence is the starkest example: benchmark rank 1st, judge rank last (−6). But the pattern extends further. Consilience wins the largest judge margin (+10) while sitting at 89.79 on the benchmark—0.4 points from control. The judge rewards architectural ambition, cross-domain insight, and design reasoning. The benchmark rewards whether the code actually isolates failures and routes traffic correctly. Neither is wrong, but for production evaluation, *behavioral benchmarks should take precedence*—a load balancer that isolates cascading failures matters more than one the judge finds architecturally interesting.

5. The middle pack is noise. Five conditions (consilience, contrarian, explicit, control, reflective) cluster in a 0.6-point band from 89.20 to 89.79. With $N=1$ per condition, we cannot distinguish this from variance. The experiment has statistical power to detect flaneur and biomimetic’s ~ 6 -point gap but cannot make meaningful claims about conditions separated by less than a point.

6. The compute-time confound is real but not decisive. Research conditions consume 3–31 minutes of total time versus control’s 8 minutes. However, flaneur’s *coding* stage took only 159 seconds—less than subtle (183s) and far less than control (468s)—yet produced the highest score. Biomimetic’s coding stage (496s) was comparable to control’s. The confound is real—more total compute is available—but the coding-time data suggests the research *content*, not merely the extra thinking time, drives the effect for the top two conditions.

Limitations. This is a single-problem ($N=1$), single-model (Claude), single-judge experiment. Each condition was run once, so we cannot estimate within-condition variance. The LLM judge’s control scores varied between 60 and 68 across seven pairings, introducing noise in the margin estimates. The flaneur defection means we tested six research personas, not seven—the flaneur result is interesting but does not test undirected exploration as designed. Parametric knowledge contamination means Wikipedia may be reinforcing what the model already knows rather than introducing genuinely new

information.

6 Conclusion

Wikipedia access can improve LLM-generated code, but the effect is narrower and more conditional than a headline score suggests. Two of seven Wikipedia conditions clearly outperformed control—flaneur by 5.85 points and biomimetic by 5.38—while five clustered within 0.6 points of the baseline. The improvement concentrates in specific failure scenarios (cascading failure, asymmetric routing, degradation detection), not a general quality lift.

The experiment’s most interesting finding may be its flaws. The flaneur did not random-walk—it searched on-topic. The subtle condition did not discover its tools. The judge and benchmark disagree sharply on the same code. These failures expose real challenges in agent experiment design: enforcing research strategy compliance, separating Wikipedia’s contribution from parametric knowledge, and choosing evaluation methods that align with the quality dimensions that matter.

Biomimetic stands out as the only condition that followed its cross-domain instructions, produced a structurally distinct architecture, and scored well—suggesting that when structural translation from another domain actually happens, the results can be meaningful. Whether this transfers to other problem types, at larger N , remains an open question.