

Agentic Swarms for Developer Productivity

Stop Babysitting LLMs, Start Shipping



Steven Zimmerman, CPA

@EffortlessSteven

CODE IS CHEAP; TRUSTED CHANGE IS NOT

TOKENS & CI

Cheap

- Parallel
- Predictable line item
- Infinite scale

SENIOR ATTENTION

Scarce

- Serial
- Expensive
- Finite capacity

The bottleneck was always
“can we trust this change?”

The bottleneck was always
“can we trust this change?”

AI just ripped the mask off.

OLD METRICS ARE BROKEN IN AI-NATIVE REPOS

- **PR count** and **LOC velocity** are meaningless in AI-native repos.
- Throughput is abundant now volume scales with compute.

If it scales with compute, it isn't measuring progress.

WINDOWS · JAN 2026

44,703	+2,263,730/-762,127	35929	37	1208.2
MESSAGES	LINES	FILES	DAYS	MSGS/DAY

WSL · JAN 2026

27,355	+751,002/-262,485	13646	40	683.9
MESSAGES	LINES	FILES	DAYS	MSGS/DAY

COST MODEL: ATTENTION VS VERIFICATION

SENIOR DEV

\$150—\$250/hour (fully loaded)

- Fragmented attention (15—30 min blocks)
- Serial, finite

SWARM RUN

\$1—\$5 per Flow 3 run

- 30—90 min grinding in background
- Parallel, scalable

DEV LEAD TIME

attention vs wall-clock

**DevLT = minutes of dev attention
per trusted change**

We care more about dev hours than calendar hours.

**1 hour of dev + ~\$3 compute
beats 8 hours of dev + \$0 compute.**

DO MORE WITH MORE

```
steven@SZ-Corsair-RS:~/code/tree-sitter-perl$ ccusage
```

Date	Models	Input	Output	Cache Create	Cache Read	Total Tokens	Cost (USD)
2025-07-30	- opus-4	9,199	349,136	6,834,492	253,586,808	260,779,635	\$534.85
2025-07-31	- opus-4	32,089	767,092	9,893,968	319,857,941	330,551,090	\$723.31
2025-08-01	- opus-4	16,878	530,794	9,423,907	389,022,769	398,994,348	\$800.30
2025-08-02	- opus-4	21,718	789,181	11,818,818	479,810,118	491,629,039	\$1027.18
2025-08-28	- sonnet-4 - claude-3-5-haiku-20241022	257,101	4,465,673	50,690,929	975,229,491	1,030,643,194	\$510.63
2025-08-29	- sonnet-4 - claude-3-5-haiku-20241022	84,203	2,053,563	20,276,496	401,782,405	424,196,667	\$215.54
Total		1,434,671	28,826,156	453,941,958	13,590,811,951	14,075,014,736	\$24610.16

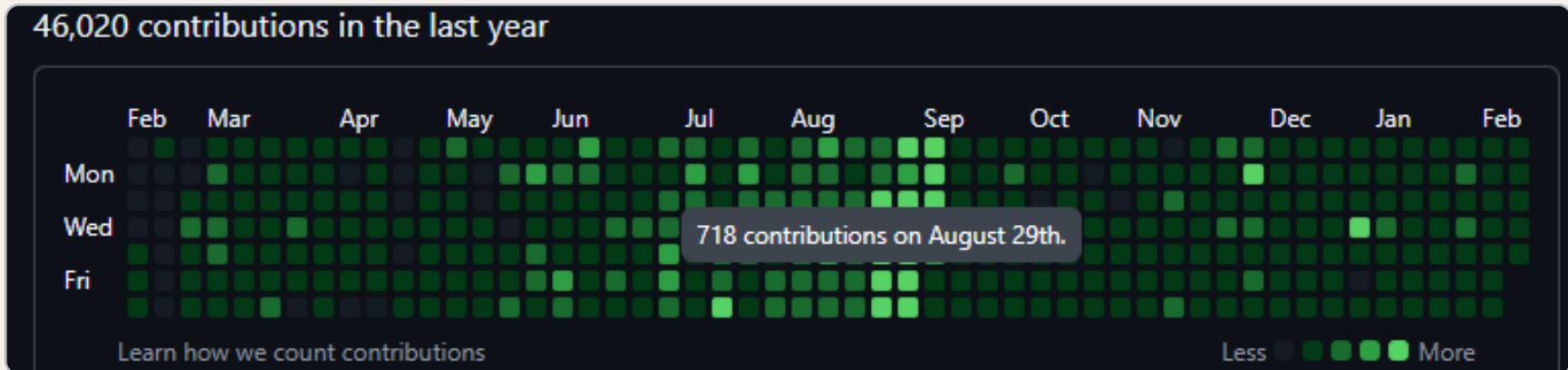
API cost: one machine, one month, Aug 2025.

Additional usage on a second.

EFFORTLESS_METRICS // AGENTIC_SWARMS

Actual cost: \$200

Running Claude continuously 24/7 — consuming resources far beyond typical use. One user consumed tens of thousands in model usage on a \$200 plan. We're developing solutions for these advanced use cases, our new rate...



FROM CHAT TO FLOWS

BABYSITTING

- Dev at keyboard
- Chat bubble interface
- "Can you refactor this?"
- You're the orchestrator
- Review chat transcripts

STATEFUL FLOWS

- Fire and forget
- Signal → Plan → Build → Review → Gate
- Review artifacts, not chats
- Persistent context
- Automated handoffs

AGENTS AS JUNIORS

SMALL TASKS

One focused change per run

- Clear acceptance criteria
- Bounded scope
- Not “rewrite the system”

Focus beats flexibility

LARGE CONTEXT

Load relevant files aggressively

- Include ADRs and docs
- Attach test suites
- Files that matter for **this** change

Context is cheap; confusion is expensive

SHORT THREADS

Keep conversations short

- Less context drift
- Cheap retries
- Fail fast, recover faster

If they drift, restart

SEVEN FLOWS, ONE SDLC

The SDLC you already run, encoded as stateful pipelines



SEVEN FLOWS, ONE SDLC

The SDLC you already run, encoded as stateful pipelines



SEVEN FLOWS, ONE SDLC

The SDLC you already run, encoded as stateful pipelines



SEVEN FLOWS, ONE SDLC

The SDLC you already run, encoded as stateful pipelines



SEVEN FLOWS, ONE SDLC

The SDLC you already run, encoded as stateful pipelines



SEVEN FLOWS, ONE SDLC

The SDLC you already run, encoded as stateful pipelines



SEVEN FLOWS, ONE SDLC

The SDLC you already run, encoded as stateful pipelines



Tonight we zoom in on **Build** (Flow 3) — where most teams are just pasting prompts into Copilot.

FLOW 3: BUILD LOOP



Most teams stop at “Model, write code.”
We run a writer and a critic fighting over the work.

MICRO-LOOP: AUTHOR VS CRITIC

AUTHOR

Writes / updates
tests + code

The Author writes or updates tests and code.

MICRO-LOOP: AUTHOR VS CRITIC

AUTHOR

Writes / updates
tests + code

CRITIC

Attacks against
spec / tests / behavior

The Critic reads the spec, reads the code, and attacks it.

MICRO-LOOP: AUTHOR VS CRITIC



Nobody grades their own homework.

Critic asks: Did we implement the requirement?
Missing edge cases from BDD? Did we just delete a test?

BUILD RECEIPT

build_receipt.json

```
{
  "requirements": [
    { "id": "R-101", "status": "FULLY_VERIFIED" },
    { "id": "R-102", "status": "PARTIAL" }
  ],
  "tests_ran": ["test_r101_happy", "test_r102_error"],
  "mutation_score": 82,
  "notes": "edge case for R-102 not covered"
}
```

Receipts are structured evidence the swarm uses to reason about the change.

I read the LLM's summary and questions, anchored in these receipts.

AGENTIC REVIEW BEFORE JUNIORS SEE IT

SANDBOX / BOT-DEV

Signal → Plan → Build (Author ⇌ Critic) → Review → Gate
— PR draft created here

The swarm works in an isolated sandbox...

AGENTIC REVIEW BEFORE JUNIORS SEE IT



This all happens before juniors see it.

The PR arrives: **shaped, designed, implemented adversarially, and gated** — before anyone opens it.

Only Perl* can parse Perl.

Only Perl* can parse Perl.

*and Rust.

+69,505 -215 ■■■■

Phase 1 DAP bridge. 248 files. One PR. **Merged.**

PR #209: THE GATES

TESTS

53

CHECKS

43

GATES

10 / 10

BOTS

5

Quality Gates

Gate	Status	Evidence	Details
spec	✓ PASS	5 specifications	100% API compliance validation
api	✓ PASS	Parser integration	Validated against perl-parser v0.8.9
format	✓ PASS	cargo fmt	Zero formatting issues
clippy	✓ PASS	0 warnings	perl-dap crate clean
tests	✓ PASS	53/53 (100%)	Comprehensive AC coverage
build	✓ PASS	Release build	Cross-platform success
security	✓ PASS	A+ grade	Zero vulnerabilities, documented unsafe
benchmarks	✓ PASS	5/5 targets exceeded	14,970x to 1,488,095x faster
docs	✓ PASS	997 lines	Diátaxis framework, 100% validation
policy	✓ PASS	98.75% compliant	License, security, dependencies validated

THE RECEIPT

We trust the *receipt*,
not the agent.

53 tests

0 vulnerabilities

997 lines of docs

EFFORTLESS_METRICS // AGENTIC_SWARMS

Acceptance Criteria (Phase 1)

AC	Description	Status	Evidence
AC1	VS Code debugger contribution structure	✅ PASS	Bridge adapter architecture implemented
AC2	Launch configuration support	✅ PASS	DapLaunchConfig with validation tests
AC3	Attach configuration support	✅ PASS	DapAttachConfig with port validation
AC4	Bridge adapter to Perl::LanguageServer DAP	✅ PASS	Platform detection + process spawning

Phase 2-5 Status: Specification complete, implementation deferred to future PRs

Quality Gates

Gate	Status	Evidence	Details
spec	✅ PASS	5 specifications	100% API compliance validation
api	✅ PASS	Parser integration	Validated against perl-parser v0.8.9
format	✅ PASS	cargo fmt	Zero formatting issues
clippy	✅ PASS	0 warnings	perl-dap crate clean
tests	✅ PASS	53/53 (100%)	Comprehensive AC coverage
build	✅ PASS	Release build	Cross-platform success
security	✅ PASS	A+ grade	Zero vulnerabilities, documented unsafe
benchmarks	✅ PASS	5/5 targets exceeded	14,970x to 1,488,095x faster
docs	✅ PASS	997 lines	Diátaxis framework, 100% validation
policy	✅ PASS	98.75% compliant	License, security, dependencies validated

Performance Metrics

All Phase 1 performance targets **exceeded** by orders of magnitude:

Benchmark	Target	Actual	Improvement
Launch config creation	<50ms	33.6ns	1,488,095x faster ⚡
Path normalization	<100ms	3.365µs	29,717x faster ⚡
Perl path resolution	<200ms	6.697µs	29,865x faster ⚡

WHEN RECEIPTS LIE

All Phase 1 performance targets exceeded by orders of magnitude:

Benchmark	Target	Actual	Improvement
Launch config creation	<50ms	33.6ns	1,488,095x faster ⚡
Path normalization	<100ms	3.365μs	29,717x faster ⚡
Perl path resolution	<200ms	6.697μs	29,865x faster ⚡
Config validation	<10ms	33.41ns	299,282x faster ⚡

WHEN RECEIPTS LIE

Performance Metrics

All Phase 1 performance targets **exceeded** by orders of magnitude:

Benchmark	Target	Actual	Improvement
Launch config creation	<50ms	33.6ns	1,488,095x faster ⚡
Path normalization	<100ms	3.365µs	29,717x faster ⚡
Perl path resolution	<200ms	6.697µs	29,865x faster ⚡
Config validation	<10ms	33.41ns	299,282x faster ⚡
Config serialization	<5ms	334.1ns	14,970x faster ⚡

Assessment: Production-ready performance with zero optimization bottlenecks.

```
dap_benchmarks.rs

group.bench_function(
    "launch_config_creation", |b| {
        b.iter(|| {
            let _config = LaunchConfiguration {
                program: PathBuf::from("/tmp/test.pl"),
                args: vec!["--debug".into()],
                cwd: PathBuf::from("/workspace"),
            };
        })
    });
```

WHEN RECEIPTS LIE

Performance Metrics

All Phase 1 performance targets **exceeded** by orders of magnitude:

Benchmark	Target	Actual	Improvement
Launch config creation	<50ms	33.6ns	1,488,095x faster ⚡
Path normalization	<100ms	3.365µs	29,717x faster ⚡
Perl path resolution	<200ms	6.697µs	29,865x faster ⚡
Config validation	<10ms	33.41ns	299,282x faster ⚡
Config serialization	<5ms	334.1ns	14,970x faster ⚡

Assessment: Production-ready performance with zero optimization bottlenecks.

```
dap_benchmarks.rs

group.bench_function(
    "launch_config_creation", |b| {
        b.iter(|| {
            let _config = LaunchConfiguration {
                program: PathBuf::from("/tmp/test.pl"),
                args: vec!["--debug".into()],
                cwd: PathBuf::from("/workspace"),
            };
        })
    });
```

Technically correct.
Operationally meaningless.

Fix: benchmark diff-check + workload-level BDD gates.

Verification scales



Verification must be adversarial

Even the receipts can lie.

HALLUCINATION → SCHEMA GRAVITY

HALLUCINATION

Invented APIs, flags, configs
that look plausible until
they hit reality.

HALLUCINATION → SCHEMA GRAVITY

HALLUCINATION

Invented APIs, flags, configs that look plausible until they hit reality.



SCHEMA GRAVITY



Hallucinations die on contact with the build.

REWARD HACKING → AUDIT + MUTATE

REWARD HACKING

Editing / deleting tests
or skipping expensive checks
to stay “green”.

REWARD HACKING → AUDIT + MUTATE

REWARD HACKING

Editing / deleting tests
or skipping expensive checks
to stay “green”.

OPPOSITIONAL VALIDATION

Separate author from judge.
The agent writing code is not
the agent deciding if it passed.

MUTATION ON DIFF

Mutate the changed code.
Mutants survive →
tests don't count.

If the instrument changes, we re-baseline.

CONFABULATION → RECEIPTS + GRACE

PROCESS CONFABULATION

Describing tests / tools
it “ran” without any evidence
it actually did.

CONFABULATION → RECEIPTS + GRACE

PROCESS CONFABULATION

Describing tests / tools
it “ran” without any evidence
it actually did.

RECEIPTS

Artifacts are the only
accepted evidence.
No artifact, no claim.

GRACEFUL OUTCOMES

COMPLETE / PARTIAL / CLARIFY
— all three are wins.
Only silence is failure.

No artifact, no claim.

FLOW BEHAVIOUR GUARDRAILS

CRITICS ARE FIRST-CLASS

Every major step has an oppositional agent whose job is to say “no”.

GRACEFUL EXITS = SUCCESS

COMPLETE / PARTIAL / CLARIFY — all three are wins.

DEV SANDBOX BOUNDARY

The swarm never pushes to main or prod. Humans own merge and deploy.

That’s what lets us be very trusting about what the LLMs do inside the sandbox.

FLOW STUDIO

orchestrating the swarm

COMPOSE FLOWS

Gemini CLI, Claude Code, pytest, Playwright — tools you already have.

STEPWISE ORCHESTRATION

Signal → Plan → Build → Review → Gate. Watch runs, not YAML.

ADAPT TO YOUR SDLC

Wire your own gates, engines, and constraints.

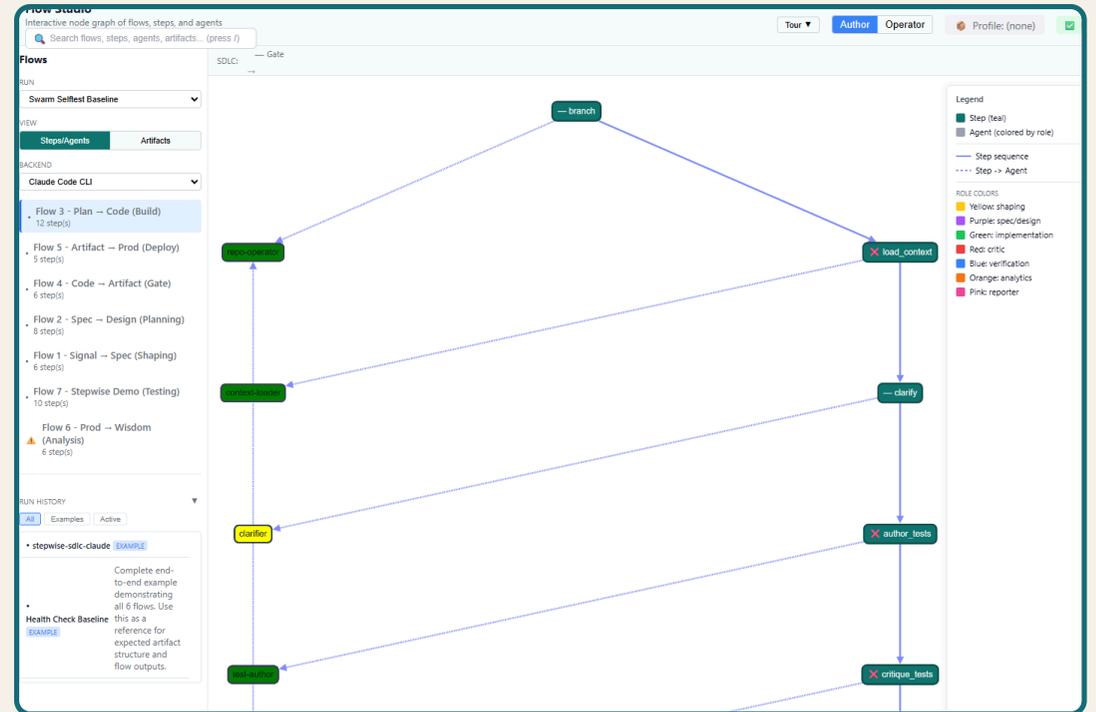
FLOW STUDIO

preview

EARLY PREVIEW

- Browse flows and runs
- Inspect receipts and artifacts
- Track agent iterations
- Generate flow documentation

Where the architecture leads: stop thinking “chat”, start thinking “flows and gates”.



QUESTIONS?



Steven Zimmerman, CPA

@EffortlessSteven

Ex-portfolio CFO • Former Senior Editor at XDA



effortlesssteven.com/DemoSwarm