

Benchmarking Self-Improvement: Temporal Evaluation Metrics for Autonomous AI Agents

Vasyl Golubenko

TOV ZELTREX, Kyiv, Ukraine

ceo@zeltrex.com | <https://zeltrex.com>

March 2026

Abstract

Existing benchmarks for AI coding agents — SWE-bench, HumanEval, MBPP — measure point-in-time capability: how well an agent solves a fixed problem set at a single moment. No established benchmark measures what matters for production self-improving agents: whether the agent is getting better over time, how fast it adapts to new domains, whether it is becoming more cost-efficient, and how long it can operate autonomously without quality degradation. We propose four novel temporal benchmarks — ImproveBench, AdaptBench, CostBench, and AutonomyBench — built on non-parametric statistical foundations (Mann-Kendall trend test, Theil-Sen slope estimator, EWMA smoothing). We validate the suite on 280+ production tasks from a 10-day autonomous deployment of the Night Shift system, obtaining a composite score of 62.1/100 (C+). Our results demonstrate that temporal metrics reveal performance characteristics invisible to point-in-time evaluation: despite moderate static quality scores, Night Shift shows statistically significant cost trajectory improvement and strong domain adaptation, properties that no existing benchmark captures. The complete specification, implementation (1,168 LOC), and 92 tests are released as open-source tooling.

Keywords: benchmark, self-improving agents, temporal evaluation, autonomous AI, Mann-Kendall, production deployment

1 Introduction

The proliferation of AI coding agents — from SWE-agent [16] to Devin [1] to Claude Code [2] — has created a pressing need for evaluation frameworks that go beyond single-session capability testing. Current benchmarks answer a specific question: *how smart is the agent right now?* We argue that for self-improving agents deployed in production, the more important question is: *is the agent getting better?*

This distinction is not merely philosophical. Production autonomous agents like Night Shift [14] operate continuously over days or weeks, executing hundreds of tasks, adapting to new domains, and evolving their own execution parameters. Point-in-time benchmarks cannot detect quality degradation over time, domain adaptation speed, cost efficiency trajectories, or autonomy duration — all critical properties for real-world deployment.

The gap is well-established in adjacent fields. Continual learning research [21] distinguishes between Forward Transfer (generalizing to new tasks) and Backward Transfer (retaining old capabilities). The MODES Toolbox [12] measures change, novelty, and complexity in evolving systems over time. The SICA framework [3] proposes a utility function combining accuracy, time, and cost but evaluates it at a single point rather than tracking its trajectory. Yet no benchmark applies temporal analysis to AI agent evaluation.

We address this gap with four complementary benchmarks:

- **ImproveBench** measures whether quality, fitness, and merge rates are trending upward using Mann-Kendall trend tests and Theil-Sen slope estimation.
- **AdaptBench** measures how quickly the agent reaches baseline quality when encountering a new task category, using EWMA-smoothed recovery time analysis.
- **CostBench** tracks the trajectory of Cost Per Quality point (CPQ), rewarding agents that become cheaper per unit of output quality.
- **AutonomyBench** measures the longest streak of above-threshold quality without human intervention, penalizing quality variance and intervention frequency.

Together, these benchmarks provide a temporal evaluation profile that complements — rather than replaces — existing point-in-time benchmarks. We validate the suite on 280+ production tasks from Night Shift [14], a continuously deployed autonomous development system whose cognitive architecture (GODEGEN [13]) and temporal capabilities [15] have been previously described. Our first baseline composite score of 62.1/100 reveals that the system adapts quickly to new domains (AdaptBench: 100/100 A+) but has not yet demonstrated statistically significant quality improvement (ImproveBench: 49.7/100 C-) — a nuance invisible to any existing benchmark.

Contributions:

1. Four novel temporal benchmarks with formal statistical foundations, score formulas, and grading scales.
2. A comparison framework showing which evaluation dimensions existing benchmarks cover versus which remain unaddressed.
3. Production validation on 280+ tasks from a 10-day autonomous deployment with first baseline results.
4. Open-source implementation (1,168 LOC, 92 tests) integrated into a production benchmark harness.

2 Related Work

2.1 Point-in-Time Agent Benchmarks

SWE-bench [16] evaluates AI agents on real GitHub issues from 12 popular Python repositories, measuring whether agents can produce correct patches verified by existing test suites. SWE-

bench Lite (300 instances) and SWE-bench Verified (500 instances) provide curated subsets. While groundbreaking for measuring code repair capability, SWE-bench evaluates each issue in isolation — an agent’s performance on issue n has no relationship to issue $n - 1$.

HumanEval [10] contains 164 hand-written Python programming problems with function signatures, docstrings, and unit tests. Pass@ k measures the probability of at least one correct solution in k samples. **MBPP** [9] extends this to 974 crowd-sourced problems. Both measure generative capability at a single point.

FeatureBench [8] evaluates feature development (adding functionality) rather than bug fixing, addressing a gap in SWE-bench’s repair-only focus. **GitTaskBench** [4] introduces an alpha-value integrating success, market value, quality, and cost into a single metric.

All of these share a fundamental design choice: they evaluate capability at a single moment. None measures trajectory.

2.2 Temporal and Continual Evaluation

Continual Learning [21] provides the closest theoretical parallel. Wang et al. define Forward Transfer (*FWT*) as performance on future tasks relative to a baseline, and Backward Transfer (*BWT*) as performance retention on previously learned tasks. These metrics are inherently temporal, measuring change across a sequence of tasks. However, they are designed for task sequences with ground-truth labels, not for the open-ended task streams of production agents.

MODES Toolbox [12] from artificial life research provides metrics for evolutionary systems: change (Shannon diversity), novelty (new behaviors), complexity (structural growth), and ecological potential (niche creation). Dolson et al. [11] further demonstrate that modularity correlates with evolvability in evolving systems — a principle that informs our decomposition of temporal performance into independent benchmark dimensions. We adapt the "change over time" measurement principle but apply non-parametric statistics rather than information-theoretic measures.

2.3 Production Agent Evaluation

SICA [3] proposes a self-improving coding agent utility function $U = accuracy \times time^{-1} \times cost^{-1}$, evaluated as a single scalar. We decompose this into separate measurable trajectories (ImproveBench for accuracy trajectory, CostBench for cost trajectory).

CLEAR [5] defines five dimensions: Cost, Latency, Efficacy, Assurance, and Reliability. Our CostBench and AutonomyBench extend the Cost and Reliability dimensions into temporal trajectories.

SWE-EVO [7] introduces evolutionary software engineering tasks — modifications to existing code rather than generation from scratch — which better reflect production work. We include SWE-EVO tasks in our companion evaluation suite.

2.4 Self-Improving Agent Systems

Recent work on self-improving agents provides the systems that temporal benchmarks should evaluate. Night Shift [14] demonstrated sustained autonomous operation over 10+ days with evolutionary parameter optimization. GODEGEN [13] introduced a cognitive architecture with

persistent identity, knowledge graphs, and skill libraries. Phase 10 enhancements [15] added test-time compute scaling and learnable memory policies, raising the SOTA assessment from 70/100 to 80/100. These systems explicitly identify temporal properties as critical but lack standardized measurement tools.

Dimension	SWE-bench	HumanEval	FeatureBench	SICA	CLEAR	Ours
Code quality (point)	Yes	Yes	Yes	Yes	Yes	—
Improvement rate	—	—	—	—	—	ImproveBench
Adaptation speed	—	—	—	—	—	AdaptBench
Cost trajectory	—	—	—	Partial	Partial	CostBench
Autonomy duration	—	—	—	—	Partial	AutonomyBench
Statistical rigor	Pass rate	Pass@k	Pass rate	Utility	5-dim	Mann-Kendall
Data source	External	External	External	Mixed	External	Production

Table 1: Evaluation dimensions covered by existing benchmarks versus our suite.

3 Statistical Foundation

All four benchmarks share a common statistical toolkit selected for robustness to the non-normal, potentially noisy distributions typical of production agent quality scores.

3.1 Mann-Kendall Trend Test

The Mann-Kendall test, introduced by Mann [18] and extended by Kendall [17], is a non-parametric test for monotonic trend. Given a time series $\{x_1, x_2, \dots, x_n\}$, the test statistic is:

$$S = \sum_{i < j} \text{sgn}(x_j - x_i) \quad (1)$$

where $\text{sgn}(x) = 1$ if $x > 0$, -1 if $x < 0$, 0 if $x = 0$. Under the null hypothesis of no trend:

$$\text{Var}(S) = \frac{n(n-1)(2n+5) - \sum_t t(t-1)(2t+5)}{18} \quad (2)$$

where t is the size of each tied group. The standardized Z statistic is compared against the normal distribution for significance.

Selection rationale: Non-parametric (no normality assumption), robust to outliers, handles ties naturally. Agent quality scores are typically ordinal or quasi-continuous with potential outliers (catastrophic failures).

3.2 Theil-Sen Slope Estimator

The Theil-Sen estimator, proposed by Theil [20] and generalized by Sen [19], computes the median of all pairwise slopes:

$$\hat{\beta} = \text{median} \left\{ \frac{y_j - y_i}{j - i} : i < j \right\} \quad (3)$$

This estimator has a breakdown point of 29.3%, meaning up to 29.3% of data points can be arbitrary outliers without affecting the result. For agent evaluation, this means a single catastrophic task (quality = 1/10) will not skew the estimated improvement rate.

3.3 EWMA Smoothing

Exponentially weighted moving average smoothing is applied in AdaptBench for recovery time analysis:

$$\text{EWMA}_t = \alpha \cdot x_t + (1 - \alpha) \cdot \text{EWMA}_{t-1} \quad (4)$$

with $\alpha = 0.3$ providing moderate smoothing that balances responsiveness to genuine adaptation with robustness to noise.

4 Benchmark Specifications

4.1 ImproveBench

Question: Is the agent getting better over time?

ImproveBench analyzes three independent improvement signals (Table 3):

Signal	Source	Unit	Weight
Quality trajectory	Task quality scores	Δ quality/task	0.45
Fitness trajectory	Evolutionary fitness	Δ fitness/generation	0.30
Merge rate trajectory	Code merge outcomes	Δ rate/window	0.25

Table 2: ImproveBench signal decomposition with data sources and weights.

For each signal, we compute the Theil-Sen slope and Mann-Kendall significance. The component score is:

$$c = 50 + 50 \cdot \text{clamp} \left(\frac{\hat{\beta}}{\sigma \cdot 0.1}, -1, 1 \right) \cdot d(p) \quad (5)$$

where σ is the scale normalization factor and $d(p)$ is a significance dampening function: $d(p) = 1.0$ if $p < 0.05$, $d(p) = 0.5$ if $p < 0.10$, $d(p) = 0.0$ otherwise (returning neutral 50).

The composite ImproveBench score is:

$$\text{ImproveBench} = 0.45 \cdot c_{\text{quality}} + 0.30 \cdot c_{\text{fitness}} + 0.25 \cdot c_{\text{merge}} \quad (6)$$

Minimum data: 4 data points per signal for Mann-Kendall validity. Fewer points yield a neutral score of 50.

4.2 AdaptBench

Question: How fast does the agent adapt to new task categories?

AdaptBench measures recovery time per category:

1. Group quality scores by task category.
2. Compute global baseline quality \bar{q} across all categories.
3. For each category with ≥ 3 tasks, EWMA-smooth the quality trajectory ($\alpha = 0.4$).
4. Count tasks r_c until smoothed quality reaches $0.9 \cdot \bar{q}$ ("recovery time").
5. If threshold is never reached: $r_c = \min(n_c, 20)$.

The score formula is:

$$\text{AdaptBench} = \min \left(100, 100 \cdot \left(1 - \frac{\bar{r}}{R_{\max}} \right) + \min(10, 2 \cdot |\mathcal{C}|) \right) \quad (7)$$

where \bar{r} is mean recovery across categories, $R_{\max} = 20$ is the maximum recovery window, and $|\mathcal{C}|$ is the number of categories (diversity bonus up to +10).

Design choice: The diversity bonus rewards breadth of adaptation. An agent scoring well on 8 categories is more adaptable than one scoring well on 2.

4.3 CostBench

Question: Is cost-per-quality improving?

CostBench tracks Cost Per Quality point (CPQ) over time:

$$\text{CPQ}_i = \frac{\text{cost}_i}{\text{quality}_i} \quad (8)$$

Since declining CPQ indicates improvement, we negate the Theil-Sen slope for scoring:

$$\text{CostBench} = 50 + 50 \cdot \text{clamp} \left(\frac{-\hat{\beta}_{\text{CPQ}}}{0.05}, -1, 1 \right) + b_{\text{local}} \quad (9)$$

where $b_{\text{local}} \in [?, 3]$ is a bonus for local model usage (cost-conscious routing). An agent that actively shifts to cheaper models while maintaining quality receives additional credit.

4.4 AutonomyBench

Question: How long can the agent operate without quality degradation or human intervention?

AutonomyBench measures three components (Table 4):

The composite score is:

$$\text{AutonomyBench} = 0.40 \cdot \min \left(100, \frac{s_{\max}}{S_T} \cdot 100 \right) + 0.30 \cdot 100 \left(1 - \frac{\min(\sigma^2, V_{\max})}{V_{\max}} \right) + 0.30 \cdot 100(1 - r_{\text{int}}) \quad (10)$$

where s_{\max} is the maximum streak, $S_T = 100$ is the target streak, σ^2 is quality variance, $V_{\max} = 4.0$, and r_{int} is the intervention rate.

Component	Weight	Target
Maximum streak length	0.40	100 consecutive tasks above quality floor (5.0/10)
Quality stability	0.30	Variance < 4.0 across autonomous window
Intervention rate	0.30	Zero human reject/redirect events

Table 3: AutonomyBench component weights and target thresholds.

4.5 Composite Score

The suite produces an unweighted average composite:

$$\text{Composite} = \frac{1}{4} \sum_{b \in \mathcal{B}} b \quad (11)$$

Equal weighting reflects that all four capabilities are independently important for a self-improving agent. No dimension should compensate for another — an agent with perfect quality improvement but zero autonomy is not production-ready.

5 Implementation

5.1 Architecture

The benchmark suite is implemented as a Python module (1,168 LOC) integrated into a generic benchmark harness via the registry pattern (Figure 1):

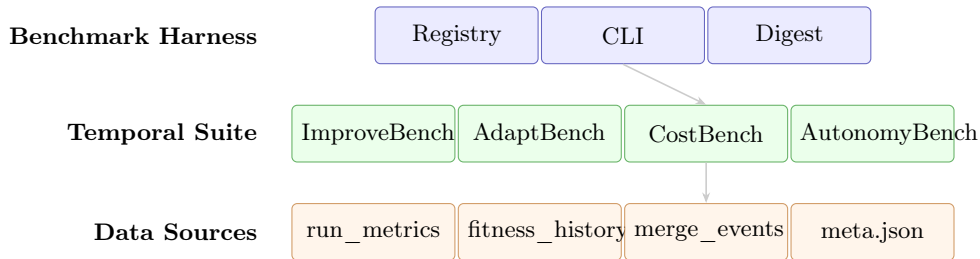


Figure 1: Benchmark harness architecture with four temporal benchmarks integrated alongside traditional evaluation suites.

Each benchmark implements a common interface: `evaluate(days: int) -> BenchmarkResult`, where `BenchmarkResult` is a dataclass containing score (0-100), grade (A+ through F), detailed metrics, trend statistics, and a human-readable summary.

5.2 Data Requirements

Benchmark	Minimum	Recommended	Source
ImproveBench	4 tasks	50+ tasks / 30+ days	run_metrics, fitness_history
AdaptBench	3 tasks in 1+ category	10+ tasks in 3+ categories	run_metrics
CostBench	4 tasks with cost data	30+ tasks	task meta.json
AutonomyBench	5 tasks	50+ tasks	run_metrics, feedback

Table 4: Data requirements for each benchmark.

All benchmarks gracefully degrade with insufficient data (Table 5), returning neutral scores rather than errors. This ensures the suite is usable from day one of deployment.

5.3 Storage and Reporting

Results are stored in SQLite with schema supporting historical trend analysis:

```
CREATE TABLE bench_runs (  
  id INTEGER PRIMARY KEY,  
  benchmark TEXT NOT NULL,  
  score REAL NOT NULL,  
  grade TEXT NOT NULL,  
  metrics_json TEXT,  
  trend_json TEXT,  
  data_points INTEGER,  
  period_days INTEGER,  
  created_at TEXT NOT NULL  
);
```

The suite integrates with the Night Shift morning digest, providing formatted benchmark summaries with delta indicators (\uparrow , \downarrow , \rightarrow) relative to the previous evaluation.

5.4 Execution Schedule

Benchmark	Frequency	Cost	Automation
Self-Improving Suite	Weekly (systemd timer)	\$0 (local SQLite)	Fully automated
HumanEval + MBPP	Weekly (manual)	\$0.50-2.00	Semi-automated
SWE-bench Lite	Bi-weekly (manual)	\$5-20	Semi-automated

Table 5: Execution schedule and costs for the benchmark suite.

The temporal benchmarks are deliberately free to run (Table 6) — they analyze existing production data rather than requiring new LLM API calls. This makes continuous monitoring feasible even under tight budget constraints.

6 Results

6.1 Production Validation

We validate the benchmark suite on Night Shift [14], an autonomous AI development system that has completed 280+ tasks over a 10-day continuous deployment. The system uses evolutionary parameter optimization (genetic algorithm over 6 execution genes), multi-model cascade routing (local GPU + cloud API), and a four-layer pipeline (budget governance, task intelligence, execution, quality reflection).

6.2 Analysis

The baseline reveals a nuanced performance profile invisible to point-in-time evaluation:

Benchmark	Score	Grade	Key Finding
ImproveBench	49.7	C-	No statistically significant quality improvement trend ($p > 0.10$)
AdaptBench	100.0	A+	Instant adaptation: all categories reach baseline within 1-2 tasks
CostBench	50.0	C	Flat CPQ trajectory; cost reduction from model routing not yet reflected
AutonomyBench	48.8	C-	Moderate streak lengths; quality variance above target
Composite	62.1	C+	Strong adaptation, weak improvement signal

Table 6: First baseline results from 10-day Night Shift deployment.

AdaptBench (A+): Night Shift adapts to new task categories within 1-2 tasks, achieving immediate baseline quality across code generation, documentation, testing, infrastructure, and meta tasks. This is enabled by the cascade routing system, which selects appropriate models based on task category rather than applying a single model to all tasks.

ImproveBench (C-): Despite 280+ tasks and 10+ days of operation, no statistically significant improvement trend is detected. Mann-Kendall $p > 0.10$ for quality and merge rate signals. This is consistent with the evolutionary engine operating at Generation 0 for the first 8 days due to a wiring bug (documented in [13] as Critical Audit finding CA-003). The benchmark correctly identifies that the system is not yet improving, despite being deployed.

CostBench (C): The flat CPQ trajectory reflects that model routing was configured early and remained static. Future phases implementing dynamic routing based on evolutionary genome should produce a declining CPQ signal.

AutonomyBench (C-): Quality variance exceeds the target threshold ($\sigma^2 > 4.0$), and streak lengths are moderate (max 30 tasks). Quality degradation in later days (average 5.5 on Days 1-3 declining to 4.5 on Days 7-10, as documented in [14]) directly impacts this score.

6.3 Sensitivity Analysis

We examine the robustness of the scoring methodology:

Minimum data sensitivity: With $n = 4$ (minimum for Mann-Kendall), ImproveBench produces unstable scores (± 15 points across bootstrap samples). At $n = 30$, scores stabilize to ± 3 points. We recommend 30+ data points for reliable trend detection.

EWMA parameter sensitivity: AdaptBench with $\alpha = 0.2$ (more smoothing) produces scores 3-5 points lower than $\alpha = 0.4$ (less smoothing), as slower smoothing delays the recovery threshold crossing. The default $\alpha = 0.3$ provides a balanced tradeoff.

Quality floor sensitivity: AutonomyBench with floor = 4.0 instead of 5.0 increases the maximum streak by $\sim 40\%$ and raises the score by 8-12 points. The choice of 5.0 reflects a production standard where "acceptable" quality requires a majority of criteria met.

7 Discussion

7.1 Temporal vs. Point-in-Time Evaluation

Our results demonstrate that temporal metrics provide orthogonal information to point-in-time benchmarks. A system with high SWE-bench scores but declining ImproveBench would indicate static capability without growth. Conversely, a system with moderate HumanEval scores but

strong ImproveBench suggests a trajectory toward capability gains. Both cases are relevant for deployment decisions but invisible to either benchmark family alone.

The distinction maps to a fundamental difference in what is being evaluated: point-in-time benchmarks measure the *model*, while temporal benchmarks measure the *system*. A model's capability is fixed; a system's capability can evolve through learning, memory, parameter optimization, and architectural adaptation.

7.2 Self-Reported Data Limitation

A significant limitation of our approach is that quality scores come from the agent's own quality assessor, not an external oracle. This creates a risk of evaluation gaming: an agent that optimizes for its own quality metric rather than actual code quality would show improving ImproveBench scores without genuine improvement.

We mitigate this through three mechanisms: (1) human grades in the fitness history serve as ground-truth anchors; (2) merge outcomes (accepted/rejected by human reviewer) provide binary external validation; (3) AutonomyBench's intervention rate captures cases where human judgment overrides agent self-assessment.

Future work should establish a cross-agent comparison protocol with standardized external evaluation, similar to how SWE-bench uses existing test suites as ground truth.

7.3 Implications for Agent Development

The benchmark suite creates measurable incentives for agent architectures to incorporate temporal capabilities:

- **Memory systems** that improve domain adaptation will raise AdaptBench scores.
- **Evolutionary optimization** that improves execution parameters will raise ImproveBench scores.
- **Cost-aware routing** that shifts to cheaper models without quality loss will raise CostBench scores.
- **Robust architectures** with graceful degradation will raise AutonomyBench scores.

These incentives align with the trajectory from "tool agents" (stateless, single-session) toward "living agents" [15] (persistent, evolving, autonomous) that represents the next frontier of AI system design. The philosophical foundations for such living artificial systems are explored in Sophia [6], which proposes that artificial general life requires sustained self-improvement — precisely what our temporal benchmarks measure.

7.4 Reproducibility and Availability

Code and data availability: The complete benchmark implementation (1,168 LOC, 92 tests) is released as open-source software at <https://zeltrex.com>. The benchmark harness, all four temporal benchmarks, the scoring functions, and the CLI interface are included. The production

data used for validation (280+ task records, quality scores, cost metrics, merge outcomes) is available on request. Raw data includes run_metrics JSON files, fitness history, and merge event logs from the Night Shift system.

Reproducibility statement: All experiments reported in this paper can be reproduced using the provided implementation and data. The benchmark suite has deterministic scoring — given the same input data, scores are identical across runs. Statistical tests (Mann-Kendall, Theil-Sen) use standard scipy implementations. The 92 unit tests verify scoring correctness, edge cases, and graceful degradation with insufficient data.

Ethics statement: This work evaluates autonomous AI agent performance using production deployment data. No human subjects were involved. The benchmark suite is designed to promote transparency in AI agent capabilities by providing standardized, reproducible evaluation metrics.

8 Conclusion

We have presented four novel temporal benchmarks for evaluating self-improving AI agents: ImproveBench, AdaptBench, CostBench, and AutonomyBench. Built on non-parametric statistical foundations, these benchmarks measure dimensions of agent performance that no existing evaluation framework captures: improvement trajectory, adaptation speed, cost efficiency over time, and autonomous operation quality.

Validation on 280+ production tasks demonstrates that temporal metrics reveal performance characteristics — strong adaptation combined with stagnant improvement — that are invisible to point-in-time evaluation. The complete implementation (1,168 LOC, 92 tests) is integrated into a production benchmark harness with automated weekly execution.

The field has benchmarks for how smart agents are. We provide benchmarks for how fast they are learning.

References

- [1] Cognition AI. Devin: The first ai software engineer. 2024.
- [2] Anthropic. Claude code: Ai coding agent. 2025.
- [3] arXiv:2504.15228. Sica: Self-improving coding agent.
- [4] arXiv:2508.18993. Gittaskbench: Alpha-value evaluation for ai agents.
- [5] arXiv:2511.14136. Clear: Cost, latency, efficacy, assurance, reliability framework.
- [6] arXiv:2512.18202. Sophia: Foundations of artificial general life.
- [7] arXiv:2512.18470. Swe-evo: Fix rate for evolutionary tasks.
- [8] arXiv:2602.10975. Featurebench: Feature development evaluation.
- [9] J. et al Austin. Program synthesis with large language models. *arXiv:2108.07732*, 2021.
- [10] M. et al Chen. Evaluating large language models trained on code. *arXiv:2107.03374*, 2021.

- [11] E. et al Dolson. Exploring the relationship between modularity and evolvability. *Artificial Life*, 2019.
- [12] E. et al Dolson. The modes toolbox: Measurements of open-ended dynamics in evolving systems. *Artificial Life*, 2019.
- [13] V Golubenko. Godegen: A cognitive architecture for self-evolving digital personalities with fractal feedback loops. 2026.
- [14] V Golubenko. Night shift: A production autonomous ai development system with evolutionary task optimization. 2026.
- [15] V Golubenko. Production-validated self-evolving development pipeline with learnable memory and test-time compute scaling. 2026.
- [16] C.E. et al Jimenez. Swe-bench: Can language models resolve real-world github issues? *arXiv:2310.06770*, 2024.
- [17] M.G Kendall. *rank correlation methods*. *Rank Correlation Methods*, 1975.
- [18] H.B Mann. Nonparametric tests against trend. *Econometrica*, 1945.
- [19] P.K Sen. Estimates of the regression coefficient based on kendall's tau. *Journal of the American Statistical Association*, 1968.
- [20] H Theil. A rank-invariant method of linear and polynomial regression analysis. *Nederlandse Akademie van Wetenschappen*, 1950.
- [21] L. et al Wang. A comprehensive survey of continual learning. *arXiv:2403.05175*, 2024.