

Claw-Eval: A Framework That Tests Whether AI Agents Are Safe, Not Just Successful

Kabui, Charles

2026-04-19

[Read at ToKnow.ai](#)

Claw-Eval: Testing AI Agents For Safety, Not Just Success

Open framework that scores side effects, not only task completion

- 326★**
GitHub stars in day one
#2 paper on Hugging Face
- 4 axes**
Complexity, horizon,
side effects, reproducibility
- Local**
No cloud or GPU cluster
required to run

An agent that deletes files to finish a task scores poorly, even if the goal was met

April 19, 2026

ToKnow.ai

The Claw-Eval team released an open framework for evaluating autonomous AI agents that goes beyond pass/fail on fixed tasks. Posted to [arXiv on April 7, 2026](#), it reached #2 Paper of the Day on [Hugging Face](#) and pulled 326 GitHub stars in its first day. The framework has four design choices that distinguish it from existing agent benchmarks like [SWE-bench](#) or [WebArena](#): configurable complexity (so you can map exactly where an agent breaks down

instead of getting one binary score), scalable horizons (single-step actions through long-horizon multi-step workflows), lightweight environments that run locally without cloud infrastructure, and scoring that explicitly tracks unintended side effects. An agent that completes a task by deleting important files, making irreversible changes, or otherwise damaging the environment scores poorly even if the goal was technically achieved.

This matters because companies are deploying coding agents, web agents, and enterprise automation agents faster than evaluation methodology has caught up. Most current benchmarks ask “did the task get done?” and stop there. A coding agent that ships a feature by silently rewriting the test suite passes. A browser agent that books a flight by also clicking through three popups passes. Claw-Eval lets a team measure those failure modes before deployment, and the lightweight design means a small team can run rigorous evaluations without a GPU cluster.

Trust, not raw capability, is now the bottleneck for putting agents into production. Frameworks that score side effects and reproducibility shift the conversation from “can it do the task” to “can we ship it.”

Read More: [GrandCode beat every human in live Codeforces competitions](#) using a multi-agent system trained with a different RL approach.

Sources:

- [Claw-Eval paper \(arXiv\)](#)
- [Claw-Eval GitHub](#)
- [Hugging Face Daily Papers](#)
- [SWE-bench](#)
- [WebArena](#)

Disclaimer: For information only. Accuracy or completeness not guaranteed. Illegal use prohibited. Not professional advice or solicitation. Read more: [/terms-of-service](#)