# Building an Open-Source AI Coding Agent

# When do you think AI will write 90% of all new code?

A.  Less than 5 Years
B.  More than 5 Years
C.  Never

# AI may automate software development tomorrow.



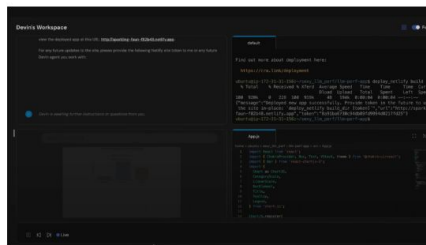**The 'Magic' Breakthrough That Got Friedman and Gross to Bet $100 Million on a Coding Startup**

**Devin: AI Software Engineer that Codes Entire Projects from Single Prompt**

Devin, an autonomous AI agent, can plan and execute complex software engineering tasks

Ben Wodecki, Jr. Editor
March 13, 2024

2 Min Read

Cognition describes Devin as a "teammate" – not to replace human engineers   COGNITION

**At a Glance**

AI startup Cognition develops an AI software engineering platform that can automate entire projects.

**Latest News**

Verticals

**FDA Clears AI Tool for Detecting Cancer Signs in Bone Marrow**
by Ben Wodecki

Apr 22, 2024                    2 Min Read

**Amazon Music Launches Maestro: AI Tool Generates Playlists From Text**
by Ben Wodecki

Subscribe to the AI Agenda newsletter

But we can automate the parts of development that suck, today.

Software Development Lifecycle today is laborious, disruptive and slow.

# LLMs have a lot of potential to automate the SDLC...



**Automated Unit Test Improvement using Large Language Models at Meta**

Nadia Alshahwan
Jubin Chheda
Anastasia Finegenova
Beliz Gokkaya
Mark Harman
Inna Harper
Alexandru Marginean
Shubho Sengupta
Eddy Wang
Meta Platforms Inc.,
Menlo Park, California, USA

**Unit Test Generation**



**LLaMA-Reviewer: Advancing Code Review Automation with Large Language Models through Parameter-Efficient Fine-Tuning**

Junyi Lu, Lei Yu, Xiaojia Li, Li Yang, Chun Zuo

**Code Review**



**Automated Program Repair in the Era of Large Pre-trained Language Models**

Chunqiu Steven Xia
University of Illinois
Urbana-Champaign
chunqiu2@illinois.edu

Yuxiang Wei
University of Illinois
Urbana-Champaign
ywei40@illinois.edu

Lingming Zhang
University of Illinois
Urbana-Champaign
lingming@illinois.edu

**Bug Fixing**

But they need to be integrated seamlessly,

with deep context,

and developer insights.

While preserving privacy,

and offering flexibility.

# AI Agents

Software that interacts with the environment, collects data, and use the data to perform self-determined tasks to meet predetermined goals.

# Introducing Patchwork

An open-source framework that effortlessly integrates into and automates your SDLC tasks, while giving you complete flexibility and control.

# Patchwork Overview

# Demo





https://github.com/patched-codes/patchwork

https://patched.codes

# SWE-Bench



**SWE-bench: Can Language Models Resolve Real-World GitHub Issues?**

https://arxiv.org/abs/2310.06770

## Leaderboard

| | Lite | Verified | Full | | | | | |

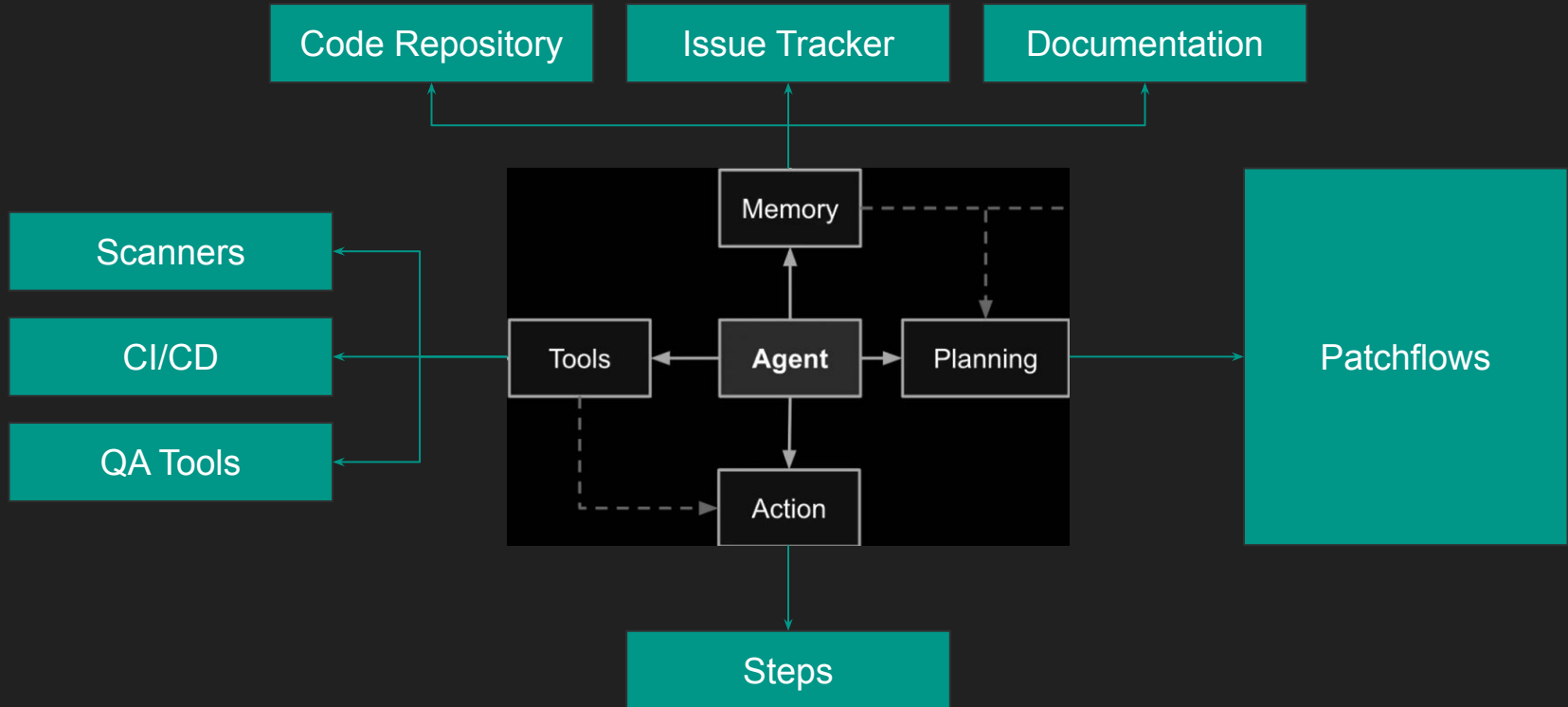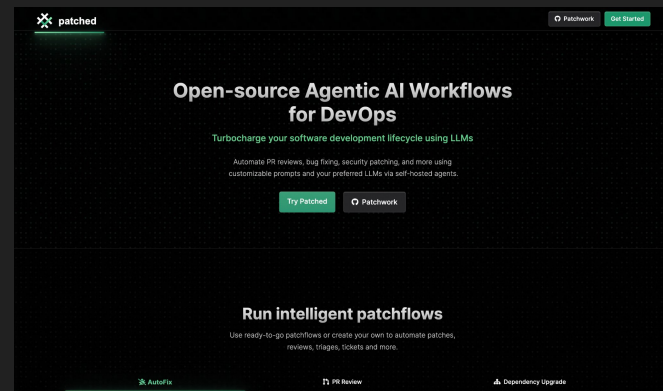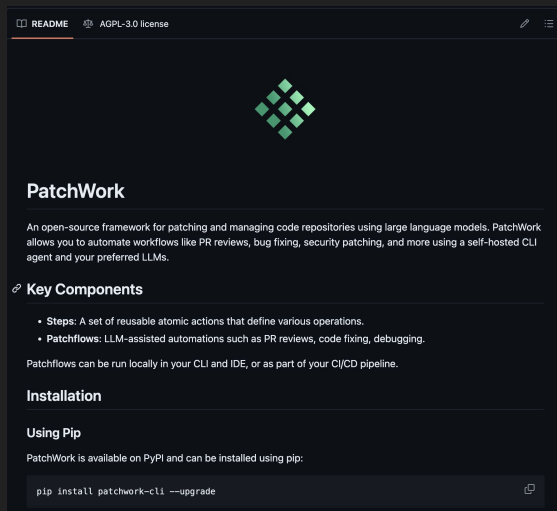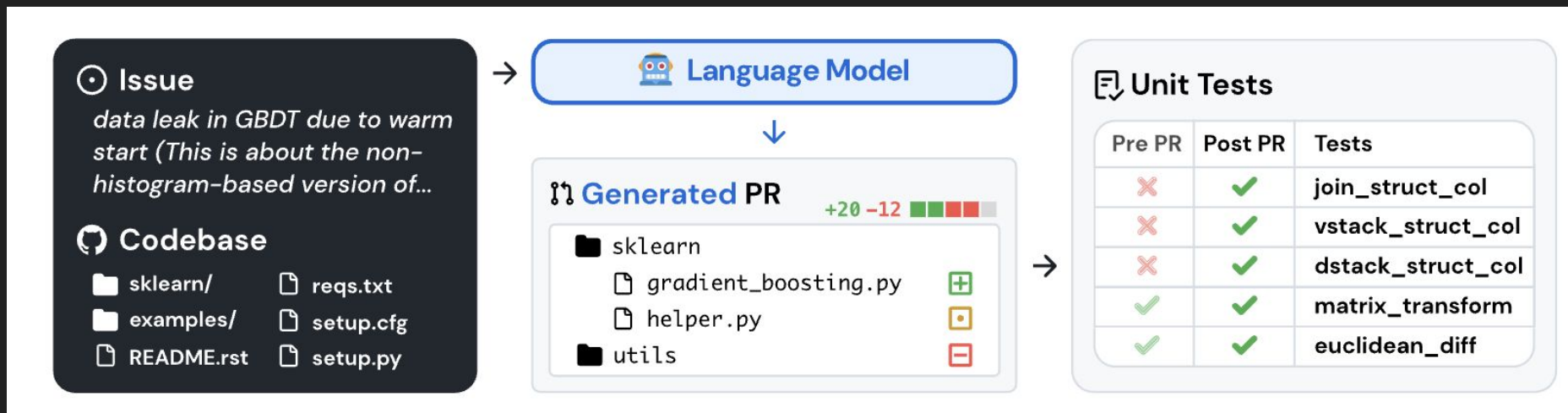| Model | % Resolved | Org | Date | Logs | Trajs | Site |
|---|---|---|---|---|---|---|
| NEW 🥇 Blackbox AI Agent | 49.00 | – | 2024-12-20 | ✓ | ✓ | 🔗 |
| NEW 🥈 Gru(2024-12-08) | 48.67 | 🔲 | 2024-12-08 | – | – | 🔗 |
| NEW 🥉 Globant Code Fixer Agent | 48.33 | G▶ | 2024-11-27 | ✓ | ✓ | 🔗 |
| NEW devlo | 47.33 | ◉ | 2024-11-22 | ✓ | ✓ | 🔗 |
| NEW 🤠 Kodu-v1 + Claude-3.5 Sonnet (20241022) | 44.67 | 🔴 | 2024-12-07 | ✓ | ✓ | 🔗 |
| 🤠 ✅ OpenHands + CodeAct v2.1 (claude-3-5-sonnet-20241022) | 41.67 | 👀 | 2024-10-25 | ✓ | ✓ | 🔗 |
| NEW 🤠 PatchKitty-0.9 + Claude-3.5 Sonnet (20241022) | 41.33 | 🌐 | 2024-12-20 | ✓ | ✓ | – |
| 🤠 Composio SWE-Kit (2024-10-30) | 41.00 | ⚡ | 2024-10-30 | ✓ | ✓ | 🔗 |
| NEW 🤠 Agentless-1.5 + Claude-3.5 Sonnet (20241022) | 40.67 | ⌛ | 2024-12-02 | ✓ | ✓ | 🔗 |
| Bytedance MarsCode Agent | 39.33 | 〽 | 2024-09-12 | ✓ | ✓ | 🔗 |
| NEW 🤠 ✅ Moatless Tools + Claude 3.5 Sonnet (20241022) | 38.33 | – | 2024-11-17 | ✓ | ✓ | 🔗 |
| Honeycomb | 38.33 | 🔶 | 2024-08-20 | ✓ | ✓ | 🔗 |
| NEW 🤠 Patched.Codes Patchwork | 37.00 | ✖ | 2025-01-04 | ✓ | ✓ | 🔗 |
| 🤠 ✅ AppMap Navie v2 | 36.00 | ◢ | 2024-11-13 | ✓ | ✓ | 🔗 |
| NEW CodeFuse-AAIS | 35.67 | CODEFUSE | 2025-01-04 | ✓ | ✓ | 🔗 |
| Gru(2024-08-11) | 35.67 | 🔲 | 2024-08-11 | ✓ | ✓ | 🔗 |

https://www.swebench.com/

Thank You!