SnT 2025
8 SEPTEMBER ONLINE DAY
9 TO 12 SEPTEMBER AT HOFBURG PALACE, VIENNA & ONLINE
CTBT: SCIENCE AND TECHNOLOGY CONFERENCE

P4.3-576

# Evaluating AI Code Assistants for Enhanced Software Development in CTBTO: A Modular Approach

Evangelos Dellis

Preparatory Commission for the Comprehensive Nuclear-Test-Ban Treaty Organization (CTBTO),
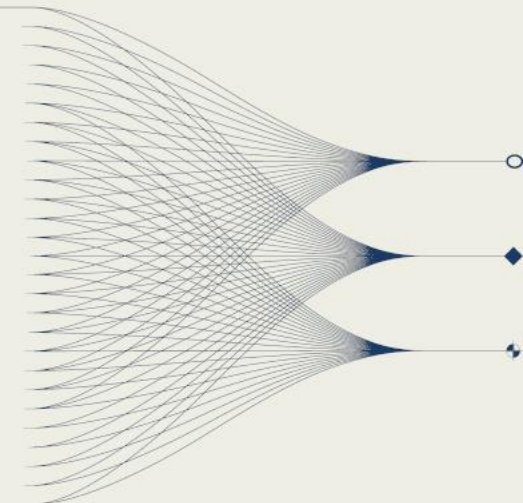P.O. Box 1200, 1400 Vienna, Austria

CTBTO | PUTTING AN END TO NUCLEAR EXPLOSIONS
PREPARATORY COMMISSION

## INTRODUCTION AND MAIN RESULTS

AI code assistants have emerged as a promising solution, offering capabilities such as autocomplete, code refactoring and context-aware suggestions. However, evaluating the effectiveness of these AI-powered tools in the CTBTO context is crucial to ensure their reliability, security, and compliance with organizational standards.

In this work we present the results of the evaluation of open-source AI code assistants in CTBTO, focusing on a modular approach that incorporates multiple components, including autocomplete models, chat models, local and remote context engines, filtering mechanisms, and training engines.

SnT 2025
CTBT: SCIENCE AND TECHNOLOGY CONFERENCE

8 SEPTEMBER
ONLINE DAY
9 TO 12 SEPTEMBER
AT HOFBURG PALACE, VIENNA & ONLINE

# Evaluating AI Code Assistants for Enhanced Software Development in CTBTO: A Modular Approach

Evangelos Dellis

**P4.3-576**

## Introduction and defitinions

An AI coding Assistant/Agent is a **tool** that uses a Large Language Model (**LLM**) that is finetuned for coding tasks and for "tool-use", more particularly:

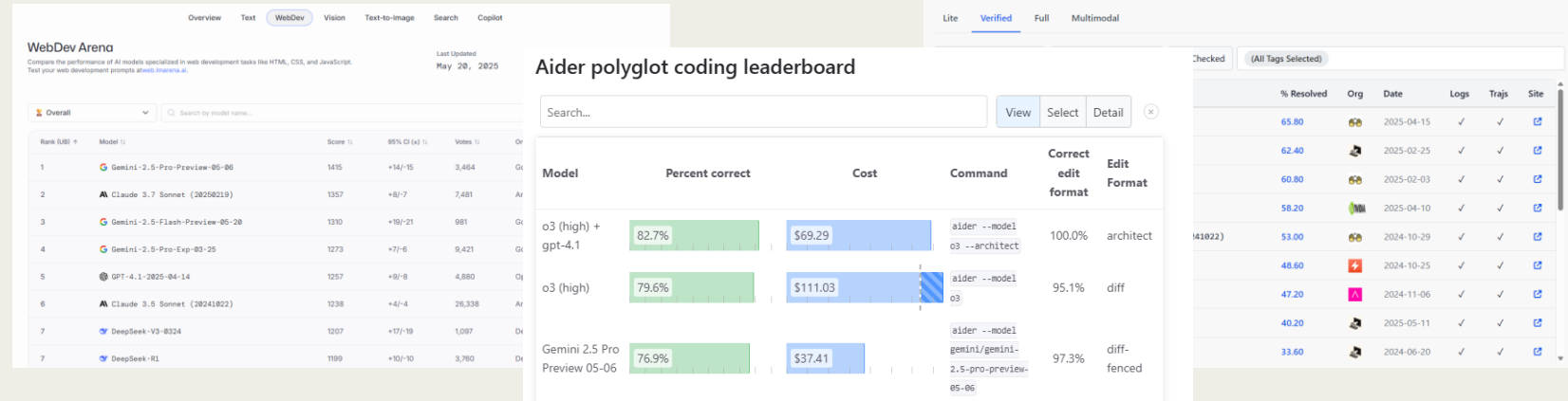**AI Coding Assistants** (higher predictability – lower agency)**:**
- Allow developers to work faster by automating repetitive or time-consuming tasks
- Solve simple, pre-defined tasks but struggle with novel or open-ended problems
- **Uses an IDE** where the developer is driving the coding

**AI Coding Agents** (higher agency – lower predictability):
- Agents are typically implemented as an LLM performing actions (via tool-calling) based on environmental feedback **in a loop**.
- Handle complex, multi-step problem-solving **tasks**, often across dynamic environments.
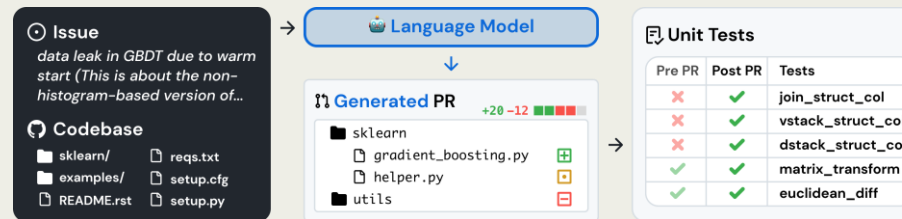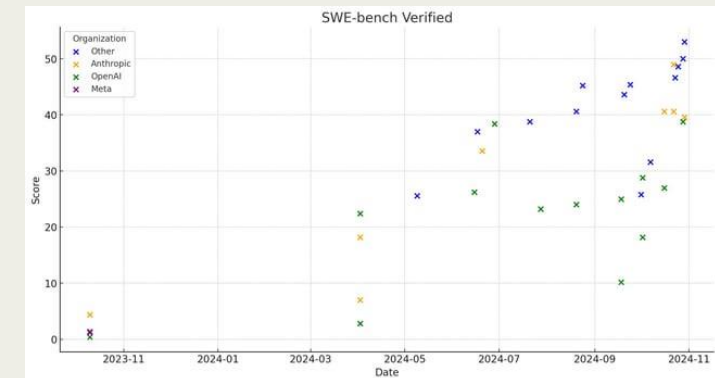


## Benchmarks and Leaderboards



**SWE-bench:** A benchmark used to evaluate the ability of large language models (LLMs) to solve real-world software engineering problems

- Given a codebase and an issue, a language model is tasked with generating a patch that **resolves** the issue.
- Results are **verifiable**!



**Benchmarks:** SWE-bench, Codeforces, HumanEval, MBPP, LiveBench Coding, LiveCodeBench & SciCode


SWE-bench Verified

CTBTO
PREPARATORY COMMISSION
PUTTING AN END TO NUCLEAR EXPLOSIONS

# Evaluating AI Code Assistants for Enhanced Software Development in CTBTO: A Modular Approach
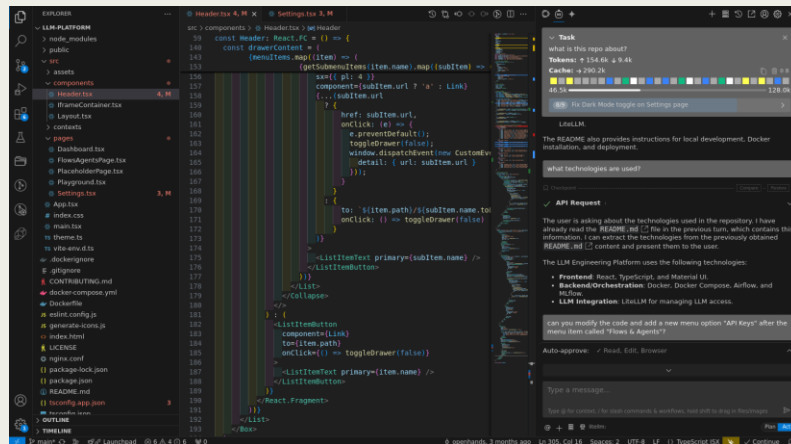
Evangelos Dellis

## AI IDE Coding Assistants

An **AI IDE** integrates code-generating LLMs directly into your existing development environment (like VS Code or JetBrains), helping you write, refactor, and understand code.
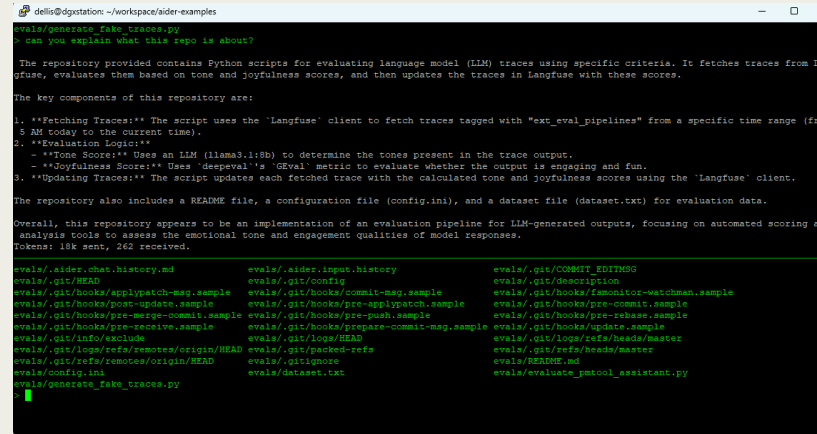
o Designed to **assist** rather than replace developers: code suggestions, completion, inline explanations, and debugging hints.

o Runs locally or via plugin, with real-time code context and seamless integration into your workflow.

We use vLLM for inference which is suitable for multiuser environments. We installed the SOTA open-source coding LLM for copilots (codestral:22b) locally in our GPU infrastructure.



## CLI-based Agents



**CLI-Based Agents** are agents you interact with through your terminal.

o They wrap LLM reasoning into a command-line interface, helping you navigate codebases, answer questions, or modify code—all from within your local dev environment.
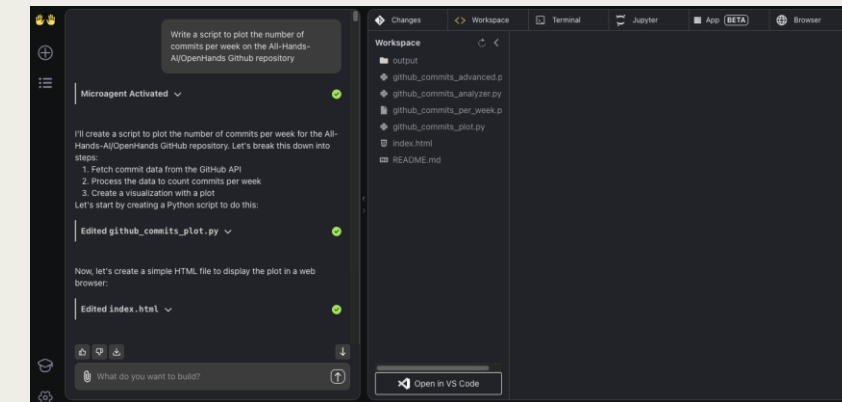
**Examples:** Claude Code, Codex CLI (open-source), Qwen Code and Aider (open-source)

Both Aider and Codex CLI can work with local docs, codebase from GitHub/Gitlab, can consult documentation, connect to Jira, etc.

## SWE-Agent or Teammate Agents

A SWE-Agent is an **asynchronous** software engineering agent that you can **delegate tasks** to and receive a fully tested PR ready to be reviewed and merged in main branch. Examples include Devin, Codex, Jules, **OpenHands** (open-source)

o It is a cloud-based software engineering agent that is running on **public** or **private** cloud (sandboxed) → spins-up new docker containers in the cloud!

o The task can be anything from fixing a bug, review code, do refactors, and implement new functionalities in response to user feedback.

o Integration with Github/Gitlab: ability to clone the repo and the ability to push a pull request

CTBTO PREPARATORY COMMISSION | PUTTING AN END TO NUCLEAR EXPLOSIONS

# Evaluating AI Code Assistants for Enhanced Software Development in CTBTO: A Modular Approach

Evangelos Dellis

P4.3-576

## A guide on how to use an autonomous SWE-Agent

You specify a prompt, and the agent goes to work in its own environment and after a few minutes, the agent gives you back a diff. The environment may contain:

- o Set package versions (like python 3.12)
- o Environment variables and Secrets
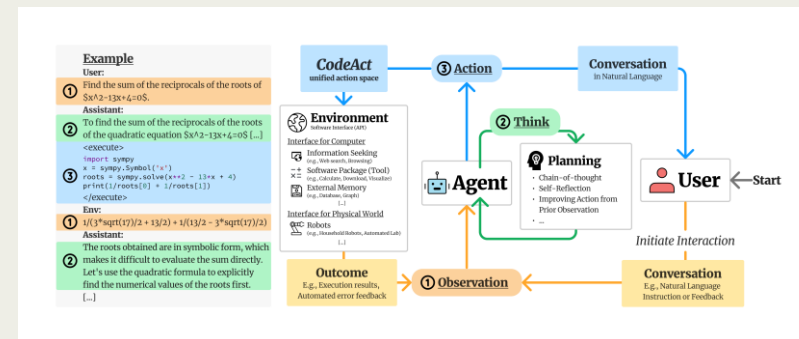- o Setup scripts (pip install -r requirements.txt)

You can execute tasks/prompts in either ask mode or code mode:

a) **Ask mode:** the agent clones a read-only version of your repo, booting faster and giving you follow-up tasks → Ask mode for brainstorming, audits, or architecture questions.

b) **Code mode:** the agent creates a full-fledged environment that the agent can run and test against. → Code mode for when you want automated refactors, tests, or fixes applied

**Development Workflow:**

1. Connect to GitHub/Gitlab **repo** and select a branch
2. Submit a **new task** (implement a new feature, fix a bug, do refactor, review code, etc)
3. A new **Docker container** based upon a base image is created. The repo is cloned at the desired branch and run any setup scripts you have from the specified working directory (install dependencies, compile, required software to run unit tests, etc).
4. The agent then **runs terminal commands** in a loop. It writes code, runs tests, and attempts to check its work. The agent attempts to honor any specified lint or test commands you've defined in an AGENTS.md file. The agent does not have access to any special tools outside of the terminal or CLI tools you provide.
5. When the agent completes your task, it presents a diff or a set of follow-up tasks. You can choose to **open a PR** or respond with follow-up comments to ask for more changes → a senior developer will review/accept the PR.

## Conclusions



**AI IDEs:** We evaluated VS Code Continue vs Cline: Both support autocomplete, code refactor, access to documents, connect to DBs, connect to GitHub/Gitlab, Jira, etc. → augments developer's work (copilot)

**CLI-based Agents:** Both tools can work with local docs, Gitlab codebases, can consult documentation, connect to Jira, etc → good Linux CLI skills are crucial

**Teammate Agents (SWE-Agent):** SWE-Agents are experimental coding agents that helps you fix bugs and build new features. They integrate with GitHub/Gitlab, understands your codebase, and works asynchronously. We evaluated OpenHands with different LLMs and for different scenarios → excels in greenfield development!

CTBTO PREPARATORY COMMISSION | PUTTING AN END TO NUCLEAR EXPLOSIONS