```
ror_mod = modifier_ob.
mirror object to mirror
mirror_mod.mirror_object
"peration = "MIRROR_X":
mirror_mod.use_x = True
mirror_mod.use_y = False
mirror_mod.use_z = False
 de Agents — Code
 irror mod.use x = False
 generation &
  ob select= 1
  debugging
  er ob.select=1
  mta.objects[one.name].se
  int("please select exaction
    OPERATOR CLASSES ----
  (vpes.Operator):
   X mirror to the selecte
  ject.mirror_mirror_x"
  FFOF X"
 ext.active_object is not
```

Learning Objectives

1

Define code agents and contrast with completion-only tools 2

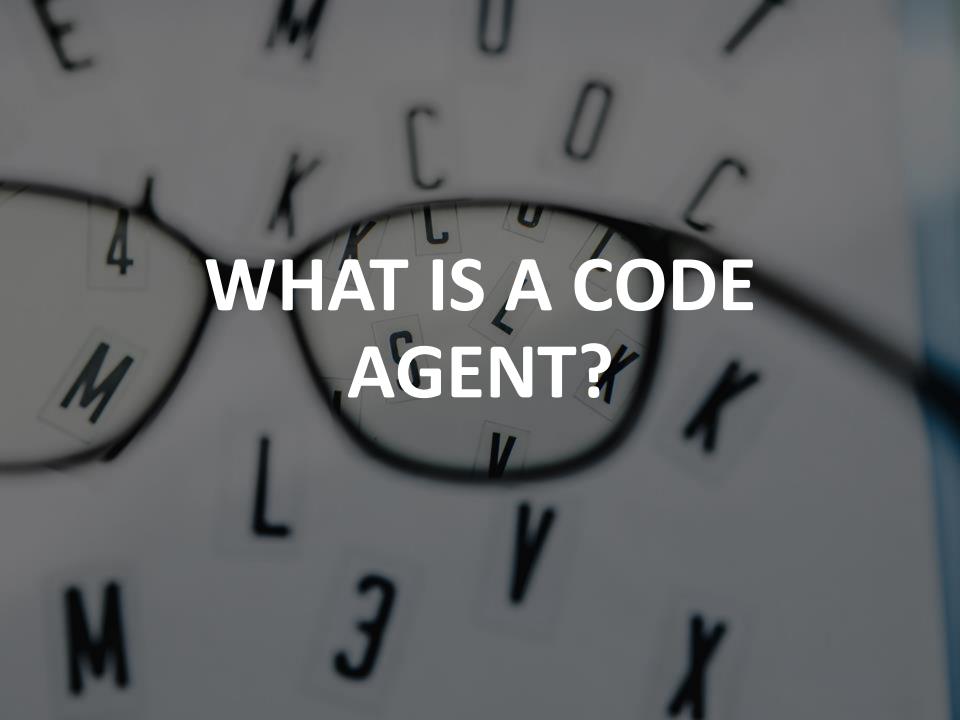
Diagram architecture & trust boundaries (ACI, tools, memory, harness)

3

Design retrieval/context plans; budget tokens effectively 4

Execute a closedloop debugging cycle; capture artifacts 5

Interpret benchmark results; pick meaningful metrics

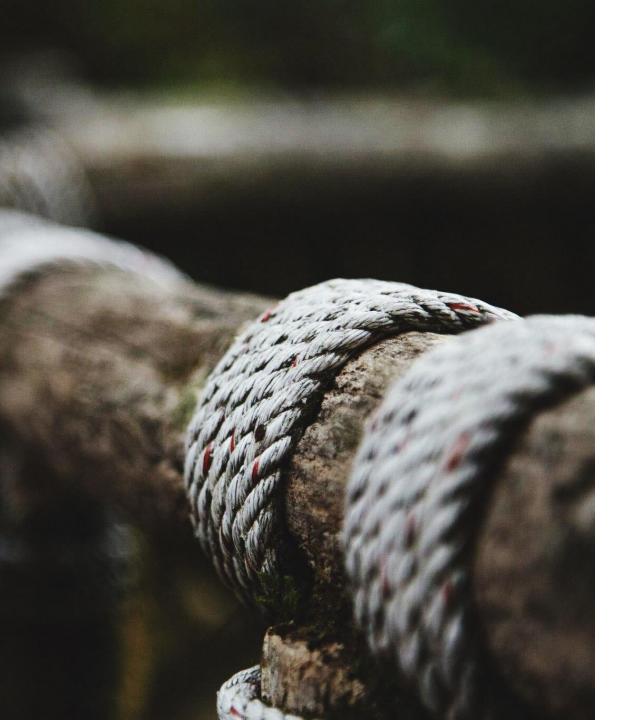


Definition & Scope

Autonomous system that plans, edits, runs, and evaluates code changes

Inputs: issue/ticket, repo files, docs; Outputs: diffs/PRs & test results Environment: editor + shell + test harness + Version Control System/Continuous Integration via Agent Computer Interface

Not just code completion; **closed loop** with execution feedback



Why "closed loop" is essential

Closed loop: the agent doesn't just generate code and stop, it

- runs the code/tests,
- observes results,
- and edits again until acceptance criteria are met

Benefits of "closed loop"

- Grounding in reality: Execution results (pass/fail, stack traces)
 cut through hallucinations and guesswork.
- 2. Objective success signal: Continuous Integration/test outcomes define "done," not the agent's confidence.
- 3. Smaller, safer diffs: Iteration encourages minimal patches with clear effects.
- **4. Better debugging:** Failures become *data* the agent can use to localize bugs.
- **5.** Reproducibility & auditability: Each cycle leaves artifacts (diffs, logs, Cl runs).
- 6. Safety hooks: The loop runs inside a sandbox with gates (coverage, SAST, secret scans), so risky edits get blocked early.

Benefits of 'closed loop

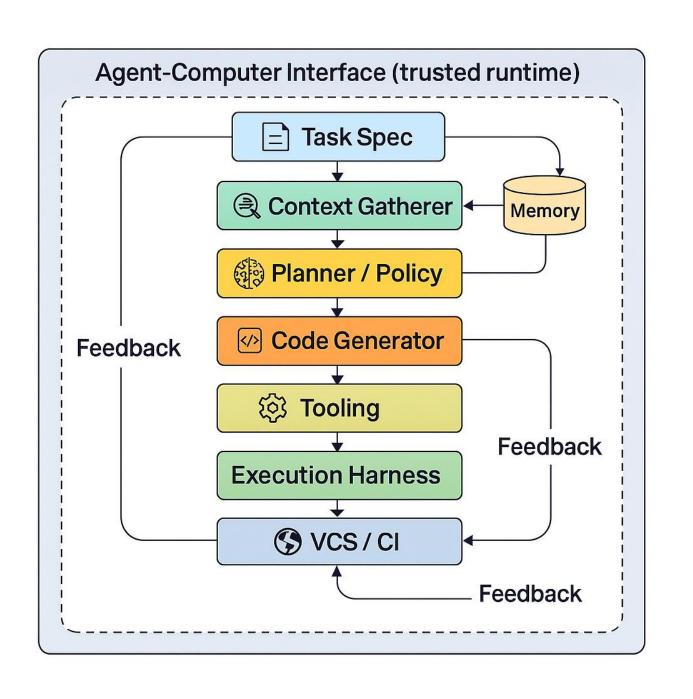


Reproducibility & auditability

Reproducibility & auditability

Mental model (pseudocode)

```
while not acceptance criteria met:
  plan = analyze(issue, context)
  patch = propose minimal diff(plan)
  apply(patch)
  results = run harness() # tests, lint, build, coverage
  if results.green:
    open or update PR(patch, results)
    break
  else:
    diagnose(results) # read failing tests/tracebacks
    refine plan()
```



Context Strategies for Large Repos

- Embed & retrieve: files, symbols, docs, issues (rank by task relevance)
- Static signals: call graph, imports, types, ownership metadata
- Budget: prioritize high-signal snippets; iterative retrieval
- Determinism: cache retrieval results for reproducible runs

Generation Patterns

- Spec-first: restate requirements; define acceptance tests
- Scaffold-first: create stubs/interfaces; fill post-tests
- Diff-first: **unified patches** to minimize churn
- Decompose: small, verifiable steps vs. monolithic changes

Debugging Loop (execution-guided)

- Run tests/build; capture errors, logs, coverage
- Localize fault (binary search, delta debugging)
- Hypothesize cause; **minimal fix**; regenerate
- Selective re-runs; stop criteria: tests pass + coverage ≥ baseline

Evaluation & Benchmarks

- SWE-bench & SWE-bench Verified: real issues/PRs in multi-file repos
- SWE-agent ACI improves agent performance on SWE-bench (reported 12.5% pass@1)
- DebugBench/DebugEval: bug localization/repair tasks
- Metrics: build-green rate, time-to-green, diff churn, coverage Δ , flake rate

In-Class Activity (10–15m)

Artifact	What to produce
Plan	Subgoals + acceptance tests
Retrieval set	Top files/symbols to inspect (with rationale)
Patch outline	1–3 edits with expected test impact
Risks	2 risks + quick mitigations

Reading — Lecture 1

- SWE-bench (<u>arXiv 2310.06770</u>); SWE-bench Verified (OpenAl, 2024/25)
- SWE-agent (<u>arXiv 2405.15793</u>; NeurIPS'24)
- OpenHands (<u>arXiv 2407.16741</u>; OpenReview)
- AutoCodeRover (<u>arXiv 2404.05427</u>)
- <u>Self-Debug</u> (ICLR'24); <u>Reflexion</u> (NeurIPS'23)
- DebugBench/DebugEval (2024–25)