

Defining the runtime that turns a model into an agent

# WHAT IS A HARNESS

WHAT IS A HARNESS • HARNESS 1.0 ARCHITECTURE • A VISUAL PRIMER

## § 01 · DEFINITION

# A harness is the *fixed architecture* that turns a model into an agent.

The modern harness was born bottom-up, out of coding agents — Cursor, Claude Code, Windsurf, Codex. These products started with a concrete problem: make an LLM write and edit real code across real repositories. In solving that problem, they independently converged on remarkably similar architectures. A while-loop that calls tools. A context manager that compresses history. A permission layer that keeps things safe. The same patterns, discovered separately, over and over again. The abstractions that define these common architectures are what we call a *harness*.

Two things make a harness fundamentally different from a framework. **First, a harness works out of the box.** You

don't configure a harness into existence. It ships as a working agent with a fixed architecture — iteration loop, context management, tool registry, permission layer — already wired together and already running. There is no assembly step.

**Second — and this is the deeper shift — a harness is not designed for humans to build agents. It is designed for the agent to accomplish almost any task.** The model reads instruction files and learns your project. It discovers available tools and composes them. It writes its own skills to extend what it can do. It spawns sub-agents when the task gets too big. The human provides the goal; the harness figures out the rest.

*"Just as computer architectures matured in the 80s, we expect harness architectures to mature over the next couple years."*

— ARIZE RESEARCH

A HARNESS IS	A HARNESS IS NOT
+ A fixed runtime architecture	- LangChain / LangGraph
+ Born from coding agents	- A prompt template
+ A while-loop around the model	- A framework for humans to assemble agents
+ Context, tools, and control-flow	- The model itself
+ Observable end-to-end	- A single API call

FIG. 1 — The harness is the layer between a pre-trained model and a useful behavior. Source: Arize Research, *What is a Harness*, 04.2026.

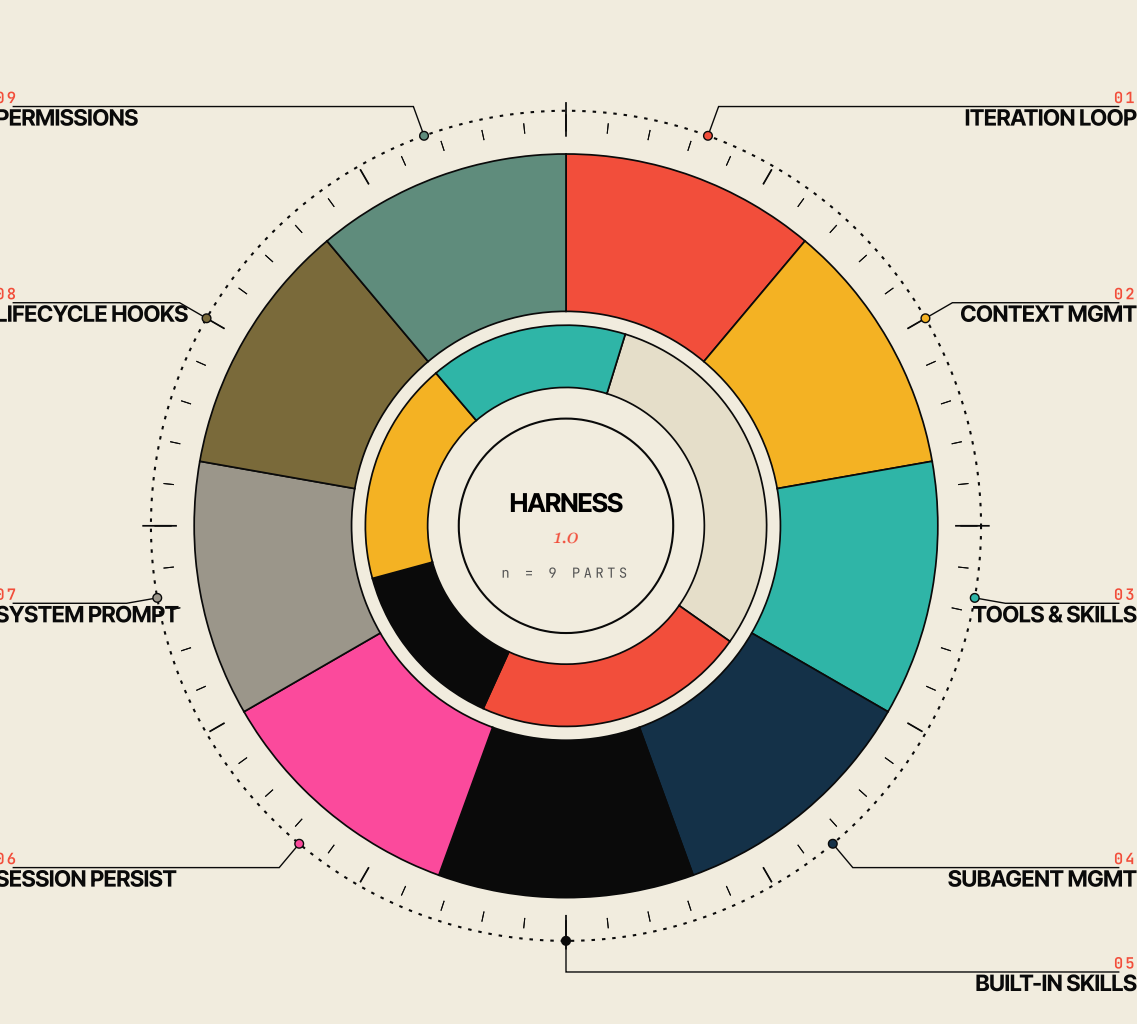
## § 02 · ANATOMY

# Nine components of the Harness 1.0 architecture.

Components fall into nine distinct areas. Every production harness we've studied implements all of them — and the quality of each is where vendors differentiate.

### THE NINE PARTS

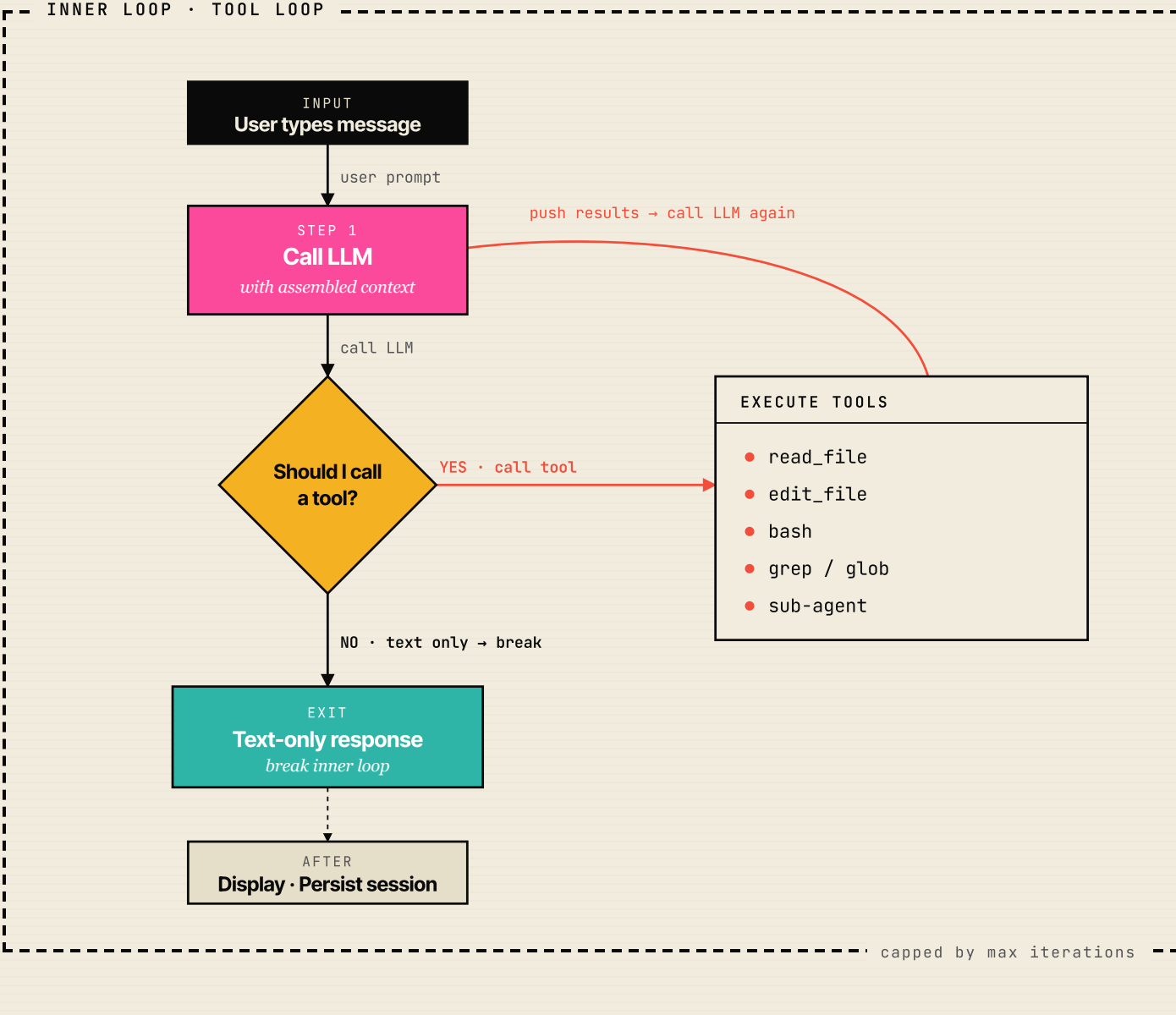
<b>01 - Outer Iteration Loop</b>	<i>while()</i>
The while-loop at the heart of the architecture: decide → call tool → push result → repeat until done.	
<b>02 - Context Management</b>	<i>tokens</i>
What enters context, how large data is compressed, how tool results are summarized before returning.	
<b>03 - Skills &amp; Tools</b>	<i>registry</i>
A registry teams can extend: build skills, add tools, manage availability, dispatch to the right executor.	
<b>04 - SubAgent Management</b>	<i>spawn</i>
Isolated child agents with their own session, restricted tools, and focused prompt. Spawn · Restrict · Collect.	
<b>05 - Built-in Skills</b>	<i>baseline</i>
The non-negotiable baseline: read/write/edit files, shell, grep, glob, plus higher-level commit/PR/test skills.	
<b>06 - Session Persistence</b>	<i>JSONL</i>
Append-only state to disk; resume exactly where a crash left off; validate, detect drift, probe health.	
<b>07 - System Prompt Assembly</b>	<i>dynamic</i>
CLAUDE.md · git · environment · tools · permissions — stitched per project with character budgets.	
<b>08 - Lifecycle Hooks</b>	<i>pre/post</i>
The extensibility seam: pre/post-tool hooks allow, deny, or modify — audit, policy, and workflow.	
<b>09 - Permission &amp; Safety</b>	<i>enforce</i>
Read-only · workspace-write · full access. Static modes + interactive approval + declarative rules.	



## § 03 · THE CORE LOOP

# The inner loop is a while-loop.

FIG. 2 — machine-paced · capped by max iterations



The while-loop is the **core architectural foundation** of the Harness 1.0 architecture.

- Decide**  
The model reads the system prompt and current context, then picks a tool — or produces a text-only answer to end the loop.
- Dispatch**  
The harness validates the call against permissions, runs the tool, captures stdout/stderr and structured results.
- Feed back**  
Tool results are pushed back into context as a `tool_result` block. The model is called again.
- Terminate**  
Loop exits on a text-only response, on max iterations, or on unrecoverable error.

## § LOOKING AHEAD · HARNESS 2.0

# Anthropic defines Harness 2.0 as a decoupling of brains and hands.

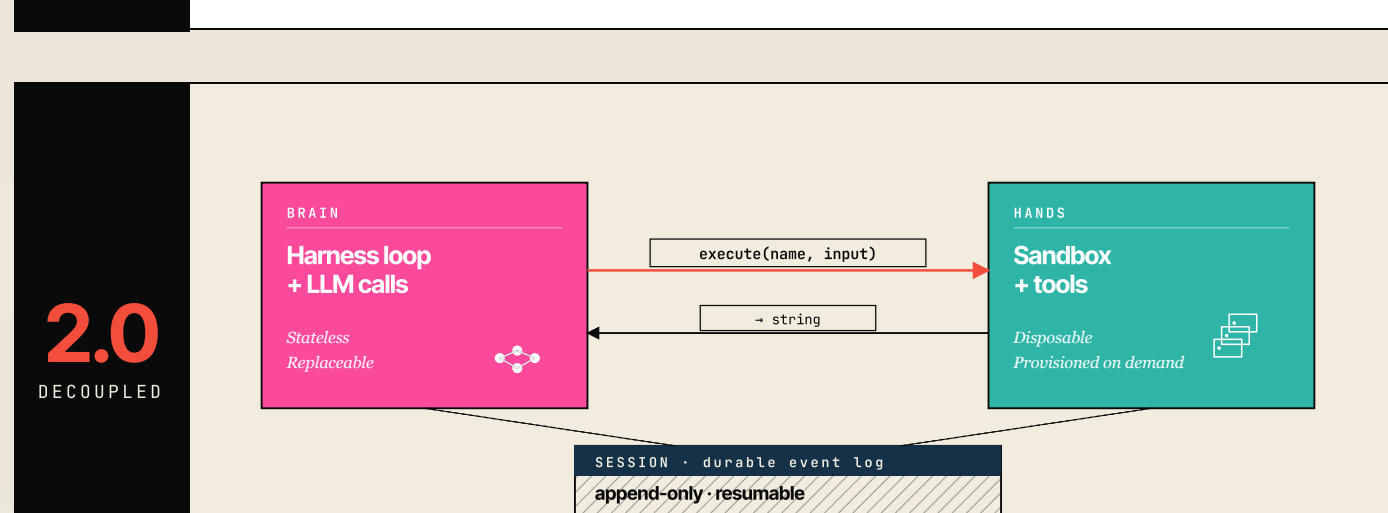
Harness 1.0 is a single process. The loop calls the LLM, executes tools, reads files, runs bash — all on one thread. It's how every coding agent works on your laptop today.

In **Harness 2.0**, the harness is just the outer loop. The **brain** is the harness loop and LLM calls — stateless, replaceable. The **hands** are sandboxes where tools run — disposable containers provisioned on demand. Between them sits the **session**: a durable event log that survives crashes on either side. The brain just calls `execute(name, input)` and gets a string back. It doesn't know if the hands are local or across the world.

*Many brains, many hands, one stable set of interfaces. This is where harness architecture is heading.*

[ANTHROPIC.COM/ENGINEERING/MANAGED-AGENTS](https://anthropic.com/engineering/managed-agents)

FIG. 5 — Decoupled architecture per Anthropic's *Managed Agents*.



N BRAINS · N HANDS · 1 STABLE SET OF INTERFACES

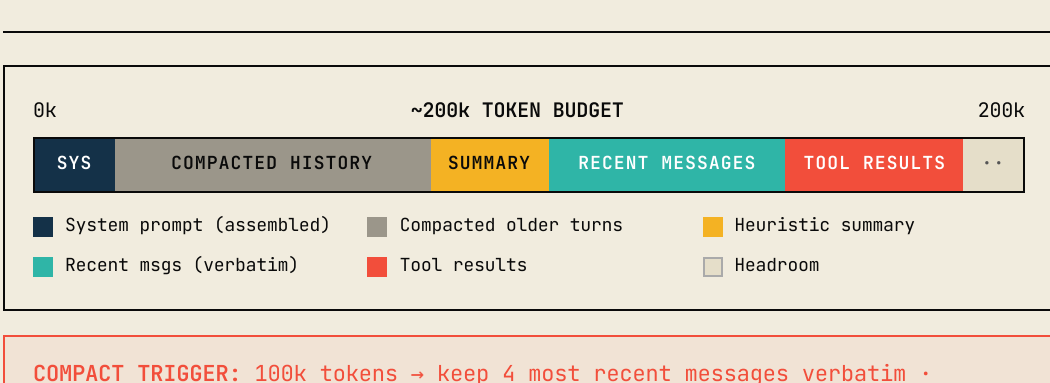
# Context, Prompts, & Sub-Agents

DIAGRAMS FOLLOW § 03.1 - 03.3

## § 03.1 · CONTEXT MANAGEMENT

# What actually fits in the context window.

Session trajectory grows every turn. The harness decides what to keep verbatim, what to summarize, and what to throw away — before the next LLM call.



**COMPACT TRIGGER:** 100k tokens — keep 4 most recent messages verbatim · summarize older turns (heuristic, not an LLM call).

### Tool-result budgets

- `read_file` — 10 MiB max, line-windowed
- `grep` — 250 lines returned
- `glob` — 100 files returned
- `bash` — 16 KB of stdout/stderr

### What streams into context

- User messages (verbatim)
- Assistant messages + reasoning blocks
- Tool results (variable, often large — hence the caps)
- Search + command output — truncated & selected

### What the harness never passes through

- Raw file contents above the line-window
- Full command output above the buffer
- Full message history above the compact trigger

FIG. 3 — Illustrative token shares · derived from *context-management* diagram (Arize Research).

## § 03.2 · SYSTEM PROMPT ASSEMBLY

# Stitched per project.

The system prompt is not a string — it is a pipeline. A static scaffold below a dynamic, per-project boundary.

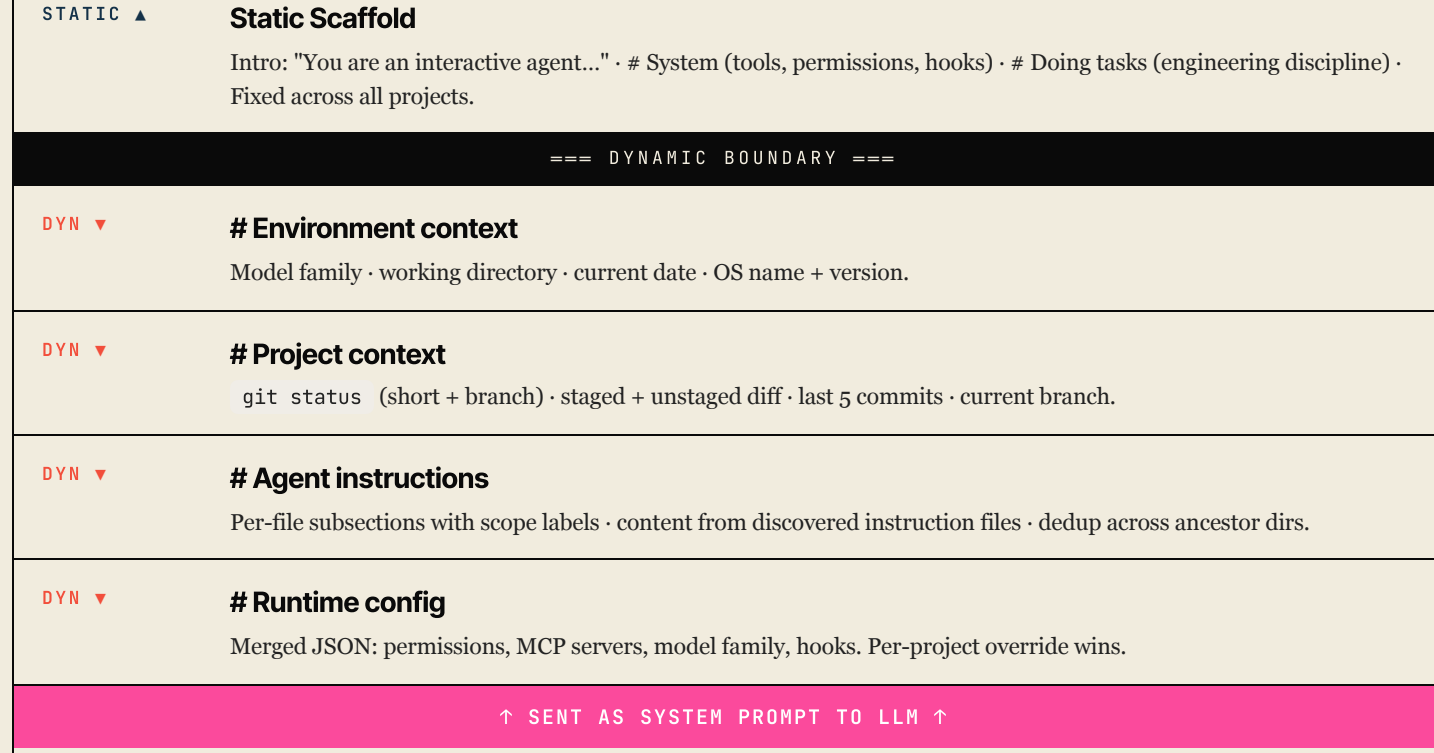
Discovery walks ancestor dirs:

- `{dir}/CLAUDE.md`
- `{dir}/AGENTS.md`
- `{dir}/.config/instructions.md`
- `{dir}/.config/rules.md`

Budgets: 4k chars per file · 12k total

Config merge order (last wins):

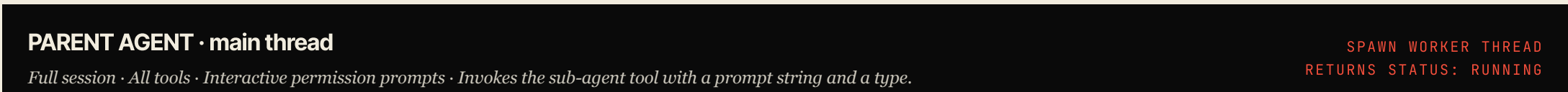
- `./config/settings.json (user)`
- `{cwd}/.config.json (project)`
- `{cwd}/.config/settings.json`



## § 03.3 · SUBAGENT MANAGEMENT

# Spawn · Restrict · Collect.

At some point a task gets too big or too parallel for a single conversation. Sub-agents intro in isolation and report results back — the parent polls a manifest on disk.



↓ SPAWN · ISOLATE · EXECUTE ↓

<b>Explore</b> (READ-ONLY)	<b>General Purpose</b> (FULL)	<b>Verification</b> (EXECUTE + READ)
Navigates the code and the web. Cannot mutate state. Best for planning and recon.	Full capabilities — but cannot spawn further sub-agents. The everyday worker.	Runs commands and reads output. Built for testing and post-change confirmation.
TOOLS: <code>read_file</code> , <code>grep</code> , <code>glob</code> , <code>web_fetch</code> , <code>web_search</code> , <code>skill</code>	TOOLS: <code>bash</code> , <code>read</code> , <code>write</code> , <code>edit</code> , <code>grep</code> , <code>glob</code> , <code>web</code> , <code>skill</code> , <code>todo</code>	TOOLS: <code>bash</code> , <code>read_file</code> , <code>grep</code> , <code>glob</code> , <code>web_fetch</code> , <code>todo</code>

**KEY CONSTRAINT** — all sub-agents run with an isolated session, an appended sub-agent directive in the system prompt, and no interactive permission prompts. The sub-agent tool is **not** in the sub-agent's own toolset — no infinite spawn. Results land in `agent-{id}.json (manifest)` and `agent-{id}.md (output)`; the parent polls.

# The Nine Components in detail

REFERENCE SHEET DIAGRAMS & COPY

## A closer look at each component.

Full architecture reference · direct from the *Harness 1.0 spec*.

<b>01 - OUTER ITERATION LOOP</b> <b>The while-loop at the core</b> The model uses the system prompt and decides what tools to call. It iterates on tools until it is finished. This is the foundation.	<b>02 - CONTEXT MANAGEMENT</b> <b>What enters the window, what gets compressed</b> How data is pulled into context, how it is simplified and compressed, how tool results are returned. Cheap to get wrong, expensive to notice.	<b>03 - SKILLS &amp; TOOLS</b> <b>Built-in tools</b> Built-in tools ( <code>read</code> , <code>edit</code> , <code>bash</code> , <code>grep</code> ) are primitives. <b>Skills</b> are the <code>SKILL.md</code> layer on top — organizational knowledge encoded as mark-down the agent invokes by name.
<b>04 - SUBAGENT MANAGEMENT</b> <b>Isolated, restricted, collected</b> A sub-agent gets its own session, a restricted tool set, and a focused system prompt. Isolation is the architectural rule.	<b>05 - BUILT-IN SKILLS</b> <b>File and code nav</b> The ops, shell, code nav — if the agent cannot read and edit files, it is not a coding agent. Higher-level built-ins ( <code>commit</code> , <code>PR</code> , <code>tests</code> ) are where vendors differentiate.	<b>06 - SESSION PERSISTENCE</b> <b>Append-only JSONL</b> Each message, each tool result, each compaction — one line at a time. Resume exactly where a crash left off; detect drift; run a health probe after compaction.
<b>07 - SYSTEM PROMPT ASSEMBLY</b> <b>Static scaffold · dynamic project context</b> Walks ancestors for <code>CLAUDE.md</code> / <code>AGENTS.md</code> . Injects git status, env, tool list. Stitched with character budgets so the window isn't blown before turn one.	<b>08 - LIFECYCLE HOOKS</b> <b>Pre- and post-tool extensibility</b> Structured protocol: JSON on stdin, exit codes for allow/deny. Hooks enforce policy: <code>never rm -rf</code> , <code>log all writes</code> , <code>approve bash that touches prod</code> .	<b>09 - PERMISSION &amp; SAFETY</b> <b>Static modes + interactive approval</b> Read-only · workspace-write · full access. Each tool declares its minimum. Bash is classified dynamically ( <code>ls</code> is read-only, <code>rm</code> isn't). Human pauses on the dangerous ones.

## § 05 · CLOSING

# The harness is the architecture. The model is the CPU.

— ARIZE RESEARCH · AI ENGINEERING DESK · 04.2026

### FURTHER READING

- Context management diagram — internal
- Harness components (9-part) — internal
- Control-loop diagram — internal
- System-prompt assembly — internal
- SubAgent management — internal

ON ARIZE.COM

- AX — AI engineering platform
- Phoenix OSS tracing
- Evaluating AI Agents (Learn)