

AN INTERACTIVE BOOK ON CONTINUOUS AI

GitHub Agentic Workflows

A progressive guide that starts with agentic-workflow concepts and builds toward GitHub Agentic Workflows authoring, compilation, MCP tools, safe outputs, and CI hosting.

By **Maxim Salnikov** · Microsoft

Version 1.1 · 14 chapters · Single-page edition · Generated 2026-07-08

Contents

PART I · THE INDIVIDUAL (ONE WORKFLOW)

- 01 What Are Agentic Workflows?
- 02 The 10-Minute Win: Your First Workflow
- 03 Anatomy & the Compile Model
- 04 Triggers: When Workflows Wake Up
- 05 Engines: Choosing the Agent's Brain

PART II · THE TEAM (SAFE, REVIEWED, PATTERNED)

- 06 Safe Outputs: Acting Without Overreach
- 07 Defense in Depth: Permissions, Firewall & Strict Mode
- 08 Tools & MCP: Real Capabilities, Governed
- 09 Continuous Triage & Docs: Reading the Room
- 10 Continuous Review, Testing & CI-Doctor

PART III · THE ORGANIZATION (FLEET AT SCALE)

- 11 Reuse & Memory: Shared Components and Repo Knowledge
- 12 Trust & Operate: Observability and Debugging
- 13 Governance & FinOps: Policy and Cost at Scale
- 14 Fleets & Adoption: From One Repo to the Org

What Are Agentic Workflows?

Explain what an agentic workflow is, why the outer loop matters, and when to reach for *gh-aw* instead of plain GitHub Actions.

Objective

By the end of this chapter you can say precisely what an **agentic workflow** is, explain why the repository's *outer loop* is where it earns its keep, and judge when to reach for GitHub Agentic Workflows (*gh-aw*) instead of plain GitHub Actions. This is the vocabulary chapter: the terms defined here — *outer loop*, *Continuous AI*, *safe-by-default* — are used as settled language for the rest of the book.

Everything here targets **gh aw v0.81.6** (Public Preview). You won't write a workflow yet — that is [→Chapter 2](#). First, the idea.

● CONCEPTS FIRST, SYNTAX LATER

This chapter names *gh-aw*'s building blocks — triggers, engines, safe outputs — only to anchor the ideas behind them. Each one gets a full, hands-on chapter later. Read this for the *why*; the *how* starts in Chapter 2.

The outer loop: why CI/CD stops short

Think about how work actually moves through a repository. There is the **inner loop**: the fast, interactive coding you do at your desk — edit, run, debug, repeat — minute to minute. And there is the **outer loop**: the repository's slower, collaborative life that surrounds and outlives any one editing session — issues filed, pull requests opened and reviewed, discussions, releases, CI results, and docs that quietly drift out of date. You leave the inner loop every time you close your laptop. The outer loop keeps going.

CI/CD automated half of the outer loop

Continuous integration and deployment were a triumph of outer-loop automation — but only for the *deterministic* half. A build, a test suite, a release: these “do exactly what you tell them, every time, in the same way” ([How They Work](#)). That determinism is precisely what you want when correctness means *identical* behavior on every run.

But a large, genuinely useful class of outer-loop work is not like that. Reading a new issue and deciding whether it's a bug or a feature request. Asking a reporter for the missing reproduction step. Keeping the docs honest as the code changes. Investigating why CI went red. None of these can be written as a fixed `if/then` rule, because the right action depends on unstructured context you can only *interpret*. That is **judgment work**, and CI/CD was never built to express it — so it fell to already-overloaded humans, or it simply didn't get done.

Where the throughput actually stalls

This is the gap the book is about. Individual AI productivity has advanced quickly, but GitHub Next observes that faster inner-loop coding “can shift burdens to other team members, or to later stages in software projects” — more code generated means more to review, triage, document, and maintain ([Continuous AI](#)). The bottleneck moves *outward*, into the collaborative loop, and that is exactly the loop CI/CD left to human judgment.

● BUILDER TRACK

You already automate the deterministic half of your repo with Actions. The outer-loop judgment tasks — triage, doc upkeep, first-pass review — are the half you keep meaning to get to. That backlog is the opportunity, and it's where agentic workflows aim.

● LEADER TRACK

Your team's real throughput is won or lost in the outer loop — how fast issues get triaged, PRs get reviewed, docs stay current. Inner-loop speedups can even *increase* outer-loop load. The strategic question is not “can AI write code?” but “can we automate the collaborative judgment work that gates delivery?”

In gh-aw: Continuous AI, from GitHub Next

GitHub Next gives this idea a name: **Continuous AI** — “all uses of automated AI to support software collaboration on any platform.” It is deliberately named to rhyme with CI/CD: “Just as CI/CD transformed software development by automating integration and deployment, Continuous AI covers the ways in which AI can be used to automate and enhance collaboration workflows” ([Continuous AI](#)). The framing is a *third leg* alongside CI and CD — and, notably, a *category* rather than a product: “not a term GitHub owns, nor a technology GitHub builds.”

An agentic workflow, defined

An **agentic workflow** is the concrete unit that practices Continuous AI. GitHub defines these as “automated, intent-driven repository workflows that run in GitHub Actions, authored in plain

Markdown and executed with coding agents” ([launch blog](#)). You describe the *outcome you want* in natural language; a coding agent interprets that intent and carries out the multi-step work. Where traditional automation follows fixed logic, an agentic workflow has *agency* — it can “understand context, make decisions, and generate content by interpreting natural language instructions flexibly” ([How They Work](#)).

That makes it three things it is often confused with, but isn't:

- **Not a chatbot or an IDE assistant.** Those are interactive and human-driven, turn by turn. An agentic workflow is *standing* and event-driven: it wakes on a repository event, does one job unattended, and proposes a result.
- **Not a fully autonomous agent.** It runs in a bounded sub-loop “under defined terms,” and “pull requests are never merged automatically — humans must always review and approve” ([launch blog](#)).
- **Not a replacement for GitHub Actions.** Which brings us to how it actually runs.

It compiles to GitHub Actions

Here is the key architectural fact, and the reason `gh-aw` is additive rather than a competitor to your existing pipelines. You author a Markdown file; the `gh aw compile` command turns it into an ordinary, security-hardened GitHub Actions workflow that GitHub runs. “The `.md` file is the editable source of truth, while `.lock.yml` is the compiled GitHub Actions workflow with security hardening” ([How They Work](#)). Agentic workflows “run on GitHub Actions because that is where GitHub provides the necessary infrastructure for permissions, logging, auditing, sandboxed execution, and rich repository context” ([launch blog](#)). So `gh-aw` adds three things Actions alone lacks — an **agentic engine** that reasons over context, a **natural-language authoring surface**, and a **security model** — on top of the substrate you already trust. The compile model is [→Chapter 3](#).

Safe by default (a first look)

Letting a standing agent act in your repository sounds risky until you see the boundary. The agent runs **read-only** by default — it “can read repository state, but it cannot push commits or write to issues directly” ([gh-aw](#)). Anything it wants to change, it *requests* as a structured **safe output**, which a separate, permission-scoped job validates before applying — so that “even a fully compromised agent cannot directly modify repository state” ([Security Architecture](#)). The agent proposes; a mediated boundary disposes. That is the whole safety thesis in one line; the mechanism is [→Chapter 6](#) and the threat model behind it is [→Chapter 7](#).

● THE “CONTINUOUS X” FAMILY

Continuous AI names recurring patterns you'll meet throughout the book: **Continuous Triage**, **Continuous Documentation**, **Continuous Code Improvement**, **Continuous Summarization**, and more ([Continuous AI](#)). Each becomes one small agentic workflow — a trigger, an engine, and a set of safe outputs.

When agentic workflows fit (and when they don't)

The mental model GitHub offers is refreshingly simple: “if repetitive work in a repository can be described in words, it might be a good fit for an agentic workflow” ([launch blog](#)). More precisely, a task fits when it has all three of these traits:

- **It's judgment work** — subjective and repetitive, the kind of task “that traditional CI/CD struggle to express” because there's no fixed rule to encode.
- **Exact reproducibility isn't the point** — triage, drafting docs, researching dependencies, proposing improvements for review. A slightly different (good) answer each time is fine.
- **Actions are low-stakes and reviewable** — a comment, a label, a draft PR. Best practice is to “start with low-risk outputs such as comments, drafts, or reports before enabling pull request creation.”

When to reach for something else

Agentic workflows are additive, not universal. Keep the work in deterministic GitHub Actions — or a human's hands — when:

- **The task must be exactly reproducible.** Builds, tests, and releases must behave identically every run. GitHub is explicit: don't use agentic workflows “as a replacement for GitHub Actions YAML workflows for CI/CD”; the use cases “largely do not overlap” ([launch blog](#)).
- **The action is high-stakes or hard to reverse.** Publishing a release, deleting data, force-pushing — if a mistake can't be shrugged off with a click, it isn't a starting point.
- **Success must be 100% correct with no human in the loop.** The value here is throughput on *reviewable proposals*, not unsupervised perfection.
- **One workflow is trying to do everything.** The unit is one teammate, one job. Start narrow and let patterns emerge.

There's a deeper reason the human stays central. The impact study of GitHub Next's real repository assistant found that throughput was gated less by the model than by “how often human maintainers chose to act on the agent's proposals” ([Repo Assist impact report](#)). Agentic workflows don't remove humans from the loop — they make the humans' judgment the scarce, valuable input.

Worked example: reading the Repo Assistant's first mission

Concepts land when you see them in one artifact. Here is the book's running example — the **Repo Assistant**, an agentic workflow that triages a newly opened issue. You'll build and ship it in [→Chapter 2](#); right now we're just *reading* it, because every concept from this chapter is visible in these few lines.

```
> examples/ch02/repo-assistant-triage.md - read it as five concepts in one file

---
on:
  issues:
    types: [opened]          # (2) an OUTER-LOOP trigger: a new issue
    workflow_dispatch:      #    also runnable by hand
permissions:
  contents: read           # (5) READ-ONLY: the agent cannot write directly
  issues: read
engine: copilot            # (1) the ENGINE: the judgment that reads the issue
network: defaults
safe-outputs:             # (5) mediated writes: proposed, then applied by a scoped job
  add-comment:
    max: 1
  add-labels:
    allowed: [bug, enhancement, question, documentation]
    max: 1
---

# Repo Assistant - triage a new issue

You are the Repo Assistant. A new issue was just opened.
Read its title and body, decide what kind of issue it is, post one short
triage comment, and apply at most one label from the allowed set.
```

Five concepts, one artifact

- **(1) Agentic workflow.** The whole file is one: YAML frontmatter that configures it, plus a natural-language body that states intent. No `if/then` logic — just “triage it.”
- **(2) Outer loop.** The `on: issues: [opened]` trigger binds it to a collaborative, outer-loop moment — a new issue — not to your keystrokes. Triggers are [→Chapter 4](#).
- **(3) Continuous AI.** This is *Continuous Triage*, one of the named patterns, expressed as a single workflow.
- **(4) Compiles to Actions.** Running `gh aw compile` turns this Markdown into a hardened `.lock.yml` that GitHub Actions executes — the subject of [→Chapter 3](#).
- **(5) Safe by default.** `permissions` are read-only; the only writes are one comment and one allow-listed label, *requested* through `safe-outputs` and applied by a separate scoped job. Full mechanism in [→Chapter 6](#).

Repo Assist, measured

The Repo Assistant is modeled on GitHub Next's real “Repo Assist.” An impact study across **15 open-source repositories** reported a net reduction of **651 open issues** and a **median 9× increase** in both issue-closure and PR-merge velocity — turning largely dormant projects into actively maintained ones. Its central finding: throughput was gated by **how often human maintainers chose to act** on the agent's proposals, and 78% of the agent's draft PRs were marked ready by a human before merge ([Repo Assist impact report](#)).

Recap & what's next

You now have the vocabulary the rest of the book stands on:

- The **inner loop** is your interactive coding; the **outer loop** is the repository's collaborative life. CI/CD automated the outer loop's *deterministic* work and left its *judgment* work to humans.
- An **agentic workflow** is intent-driven, event-triggered repository automation authored in Markdown and run by a coding agent — a standing teammate that proposes, not a chatbot and not an unsupervised agent.
- **Continuous AI** is GitHub Next's name for applying AI to that outer loop — a third leg beside CI/CD — and gh-aw is how you practice it.
- gh-aw **compiles to GitHub Actions**; it is additive, adding an engine, a natural-language surface, and a security model on top of the substrate you already trust.
- The agent is **read-only by default** and every write is a mediated, reviewable proposal — which is what makes standing automation safe to trust.

What's next. Enough theory — time to ship. In [→Chapter 2: The 10-Minute Win](#) you'll install the `gh aw` CLI and get this exact Repo Assistant running end to end, triaging a real issue through the safe-by-default boundary.

The 10-Minute Win: Your First Workflow

Install the `gh aw` CLI and ship a first working Repo Assistant that triages a new issue end to end.

Objective

By the end of this chapter you can install the `gh aw` CLI and ship a working **Repo Assistant**: an agentic workflow that reads a newly opened issue, decides what kind of issue it is, and replies with a triage comment and a label — all through a reviewed, safe-by-default boundary.

It really is a ten-minute win. You write your intent as a short Markdown file, compile it into an ordinary GitHub Actions workflow, and watch a coding agent do a job that used to need a human. Everything here targets **gh aw v0.81.6** (Public Preview).

● BEFORE YOU START

This chapter builds on [→Chapter 1](#), which introduces the outer loop and Continuous AI. You will also need the [GitHub CLI](#) (`gh`) and a repository you can push to.

The problem: the overnight teammate you don't have

Open your repository's issue tracker. Somewhere in there is a new issue with no label, a vague title, and no reproduction steps — and it has been sitting for a week. Reading it, deciding whether it's a bug, a feature request, or a question, asking for the missing detail, and routing it to the right place is real work. It's just rarely the work that reaches the top of anyone's day.

Triage is judgment work, not a pipeline

Traditional CI/CD is built to “do exactly what you tell them, every time, in the same way” ([How They Work](#)). That determinism is exactly what you want for a build or a release. Triageing an issue is the opposite kind of task: there is no lookup table that maps every possible issue to the right response. You have to infer intent from unstructured prose and pick a context-dependent action. That is *judgment work* — the kind of task “where exact reproducibility doesn't matter, such as triaging issues, drafting documentation, researching dependencies, or proposing code improvements for human review” ([FAQ](#)).

That is why triage is the canonical first agentic win. It is high-volume, low-stakes (a comment or a label is trivially reversible), and every action stays human-reviewable. GitHub Next even names the pattern: **Continuous Triage** — “label, summarize, and respond to issues using natural language” ([Continuous AI](#)). And it is **additive**: agentic workflows sit alongside your deterministic pipelines, which do not change at all ([FAQ](#)).

Meet the overnight teammate

Picture a tireless teammate who owns exactly one small, recurring job. It wakes on an event, does that job, proposes the result for your review, and goes back to sleep. That is the mental model for an agentic workflow — and it is the book's running example, the **Repo Assistant**. It is modeled on GitHub Next's real “Repo Assist,” a repository assistant that labels issues, answers questions, and proposes fixes “all while the maintainer stays in control through pull request review” ([Continuous AI](#)).

One workflow equals one teammate equals one job. That is deliberate: it is not a general chatbot you prompt ad hoc, and it is not an unsupervised autonomous agent. It is a *standing*, event-driven task-owner that proposes rather than acts with free rein.

Repo Assist, measured

An impact study of Repo Assist across **15 open-source repositories** reported a net reduction of **651 open issues** and a **median 9× increase** in issue-closure and PR-merge velocity — turning largely dormant projects into actively maintained ones. Its central finding: throughput was gated less by the model than by **how often human maintainers chose to act on the agent's proposals** ([Repo Assist impact report](#)).

● BUILDER TRACK

You already automate the deterministic half of your repo with Actions. This is the other half: the judgment tasks you keep meaning to get to. Start with the one that is most repetitive and least risky — triage — and you will have something real running today.

● LEADER TRACK

The payoff is not “AI writes our code.” It is throughput on the unglamorous collaboration work that quietly decays: unlabeled issues, stale reports, doc drift. Because every action is a reviewable proposal, you get that throughput without handing over control — and the measured lever is your team's review habit, not the model.

Where this sits in the outer loop

As →[Chapter 1](#) established, the *inner loop* is the fast, interactive coding you do in your editor; the *outer loop* is your repository's ongoing collaborative life — issues, PRs, reviews, releases — that keeps moving after you close the laptop. Continuous AI applies AI to that outer loop the way CI/CD automated integration and deployment, and gh-aw is how you do it. Your first workflow is simply the smallest slice of that loop: one new issue, triaged.

In gh-aw: init, new, compile, run

The `gh aw` CLI turns that overnight-teammate idea into four small steps: **init** the repo once, **new** to scaffold a workflow, **compile** your Markdown into a real Actions workflow, and **run** it. You author *intent* in Markdown; the compiler produces the reviewable YAML that actually executes.

Install and verify

`gh aw` is a GitHub CLI extension. Install it, then confirm the version this chapter targets.

```
> Install the extension and check the version
```

```
gh extension install github/gh-aw
gh aw version
# → gh aw version v0.81.6
```

1 • `gh aw init` — set up the repo once

Run this once per repository. It is non-interactive and does not ask for an engine or configure any secrets. It prepares the repo so you can author, compile, and run — for example, marking generated `*.lock.yml` files in `.gitattributes` and adding helper skills and editor settings ([CLI Commands](#)).

```
> One-time repository setup
```

```
gh aw init
```

2 • `gh aw new` — scaffold a workflow

`gh aw new <workflow-id>` creates a single Markdown workflow at `.github/workflows/<workflow-id>.md`, pre-populated with a heavily commented template of every frontmatter option.

> Create a new workflow file

```
gh aw new repo-assistant-triage
# creates .github/workflows/repo-assistant-triage.md
```

● TIP

The generated template is intentionally verbose — great as a reference, noisy as a first example. For a clean ten-minute win, the worked example below is a hand-written *minimal* file. When an official workflow already does the job, pull one in with `gh aw add` instead ([Quick Start](#)).

3 • `gh aw compile` — Markdown becomes a workflow

Compilation is the heart of the loop. `gh aw compile` turns each `<file>.md` into the GitHub Actions lock file that actually runs, `<file>.lock.yml`. With no argument it compiles every workflow in `.github/workflows/`. You never hand-edit the lock file — you change the Markdown and recompile.

> Compile a single workflow (offline, no secrets)

```
gh aw compile .github/workflows/repo-assistant-triage.md
```

Two properties make this the workhorse of the book. First, **compilation is offline**: it never calls an engine or touches GitHub, so it needs no network and no secrets — which is exactly why every example here is compile-verified. Second, **strict mode is on by default**. Strict mode does not make you fill in boilerplate; it *refuses unsafe choices*: top-level write permissions (route writes through `safe-outputs:` instead), unpinned actions, wildcard network egress, and deprecated fields. Everything else has a safe default — omit `engine:` and you get Copilot; omit `permissions:` and the agent is read-only; omit `network:` and you get a curated egress allowlist (the same one `network: defaults` selects). So the smallest valid workflow is really just a trigger plus a body, with any writes routed through `safe-outputs:`. Our example still spells out `permissions:`, `engine:`, and `network:`, because being explicit is clearer in a teaching example. Those guardrails are not friction — they are the safe-by-default posture from Chapter 1, enforced at compile time.

4 • `gh aw run` — trigger it on GitHub

`gh aw run` dispatches a compiled workflow on GitHub Actions using its `workflow_dispatch` trigger — so a workflow must declare one to be runnable this way (ours does). Unlike `compile`,

this is a **live run**: it executes on GitHub's servers against a real repository and needs the engine's secret configured. That is why it is the one step in this chapter the book does not compile-verify.

> Manually dispatch the workflow (live run – needs the engine secret)

🔒 REQUIRES A SECRET / LIVE RUN

```
gh aw run repo-assistant-triage          # dispatch it by hand
gh aw run repo-assistant-triage --dry-run # validate without triggering a real run
```

● SECURITY & COST

Compiling is free and offline. *Running* executes a coding agent on GitHub Actions, which spends AI credits and needs Copilot authentication at run time (never to compile, and never written into the workflow file). There are two paths: the recommended `copilot-requests: write` permission — a scoped Copilot-billing request, not a repository-write scope — which bills to your organization's Copilot plan with no PAT required, or a `COPILOT_GITHUB_TOKEN` repository secret. → [Engines \(Chapter 5\)](#) weighs the trade-offs; for now, either one lets the Repo Assistant run ([Engines reference](#)).

The engine: Copilot by default

The `engine:` key selects which coding agent runs the workflow — the “judgment” that reads the issue and decides what to do. `gh-aw` supports Copilot, Claude, Codex, and Gemini, and **Copilot is the default** ([Engines](#)). If your team already has GitHub Copilot, there is no extra account to set up, which makes it the natural first-workflow engine. You can omit `engine:` entirely, but a teaching example keeps it explicit so you can see which engine ran.

> Selecting the engine in frontmatter

```
engine: copilot # the default; shown explicitly for clarity
```

You will go deeper on choosing and configuring engines in → [Engines \(Chapter 5\)](#).

First look: `safe-outputs`

Here is the piece that makes an overnight teammate safe to trust. The agent step runs **read-only**. Anything it wants to change — post a comment, add a label — it *requests* as structured output, and a separate, permission-scoped job validates and applies it. The agent proposes; a mediated boundary disposes.

> A first-workflow safe-outputs block: one comment, one allow-listed label

```
safe-outputs:  
  add-comment:  
    max: 1  
  add-labels:  
    allowed: [bug, enhancement, question, documentation]  
    max: 1
```

That is the whole safety story for a first workflow: the worst case is one comment and one label from a fixed list — never a code or settings change. This is a first look only; the full mechanism (sanitization, per-operation caps, targets, staged mode) is the subject of [→Safe Outputs \(Chapter 6\)](#), and the threat model behind it is [→Defense in Depth \(Chapter 7\)](#).

When to ship a first workflow (and when to wait)

A good first workflow shares three traits: the task is **judgment work** (there is no exact rule to follow), it is **high-volume and recurring** (so throughput matters), and every action is **low-stakes and reviewable** (a comment or a label, easily undone). Triage hits all three, which is why it is the pattern to start with. Doc nits and stale-issue nudges are close seconds.

When to wait

Agentic workflows are additive, not universal. Reach for something else when:

- **The task must be exactly reproducible.** Builds, tests, and releases must do the same thing every time — keep those as deterministic CI/CD. Agentic workflows augment those pipelines; they do not replace them.
- **The action is high-stakes or hard to reverse without review** — publishing a release, deleting data, force-pushing. If a mistake cannot be shrugged off, it is not a first workflow.
- **You would have to grant broad `write` permissions** to make it work. In strict mode that fails to compile — and it is usually a sign the scope is wrong, not that the guardrail is.
- **One workflow is trying to do everything.** Prefer one teammate, one job. Start with one or two workflows and expand as patterns emerge ([FAQ](#)).
- **The outcome must be correct 100% of the time with no human in the loop.** The value here is throughput on reviewable proposals, not unsupervised perfection.

● LEADER TRACK

“Start narrow” is a governance strategy, not just a tip. One reviewable workflow per team builds trust and a review habit before you scale — and it keeps the blast radius of any single agent to a comment and a label while you learn.

Worked example: Repo Assistant triages a new issue

Here is the whole thing: a complete, compile-verified Repo Assistant that triages a newly opened issue. It is a single Markdown file — YAML frontmatter on top, a natural-language brief below. It targets **gh aw v0.81.6**.

› examples/ch02/repo-assistant-triage.md – the complete workflow (frontmatter + body)

```
---
on:
  issues:
    types: [opened]
  workflow_dispatch:
permissions:
  contents: read
  issues: read
engine: copilot
network: defaults
safe-outputs:
  add-comment:
    max: 1
  add-labels:
    allowed: [bug, enhancement, question, documentation]
    max: 1
---
```

Repo Assistant – triage a new issue

You are the **Repo Assistant**. A new issue was just opened in this repository.

Read the triggering issue's title and body, then triage it:

1. Decide what kind of issue it is (a bug report, a feature request, a question, or a documentation gap) and how a maintainer should treat it.
2. Post **one** short, friendly triage comment that (a) restates the request in a sentence, (b) names the category you chose and why, and (c) lists any missing information the reporter should add.
3. Apply **at most one** label from the allowed set that best matches the issue.

If the issue is empty or too vague to categorize, post a comment asking for the missing details and do not apply a label.

This workflow demonstrates **safe-outputs**: the agent runs read-only and never writes to GitHub directly – it *requests* a comment and a label, which gh-aw applies from separate, permission-scoped jobs.

Reading the frontmatter

Five keys, each doing one job. Every one is stable in v0.81.6 — no preview fields.

Key	Value	What it does
<code>on</code>	<code>issues: { types: [opened] } + workflow_dispatch</code>	Two triggers: the Repo Assistant wakes when a new issue is opened, and <code>workflow_dispatch</code> also lets you run it by hand to smoke-test it. → Chapter 4 covers triggers in depth.
<code>permissions</code>	<code>contents: read, issues: read</code>	A read-only agent. It can read the repo and the issue, but cannot write anything itself.
<code>engine</code>	<code>copilot</code>	The coding agent that does the judgment — GitHub Copilot, the default engine, made explicit.
<code>network</code>	<code>defaults</code>	Explicit here for teaching clarity; if omitted, strict mode applies this same curated egress allowlist. Either way, the agent can only reach approved hosts.
<code>safe-outputs</code>	<code>add-comment: {max: 1}, add-labels: {allowed: [...], max: 1}</code>	The only writes permitted — at most one comment and one allow-listed label, each applied by a separate scoped job.

The body underneath is just the brief you would give a new teammate: who they are, what to read, and the three steps to take — with an explicit fallback for an empty issue. That prose is the editable source of truth; the agent reads it at run time.

Compile it

Compile the file to prove it is valid. This is offline — no secrets, no network.

```
> Compiling the example, with the exact successful output (exit code 0)
```

```
gh aw compile examples/ch02/repo-assistant-triage.md
✓ examples\ch02\repo-assistant-triage.md (100.3 KB)
✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

Zero errors, zero warnings. The compiler wrote `repo-assistant-triage.lock.yml` next to the Markdown; its metadata records `"compiler_version": "v0.81.6"`, `"strict": true`, and `"agent_id": "copilot"`, and the agent job's permissions are `contents: read` — read-only, exactly as declared. The writes live in separate safe-output jobs. You commit both files: the `.md` you author and the `.lock.yml` that runs ([How They Work](#)).

Run it

The true end-to-end path is simply to push the workflow and **open a test issue** — it triggers on `issues: { types: [opened] }`, so the Repo Assistant wakes on its own and replies with a comment and a label. To smoke-test on demand instead, dispatch it manually. This is the live-run step: it runs on GitHub Actions and needs Copilot authentication (the `copilot-requests: write` permission or a `COPILOT_GITHUB_TOKEN` secret) configured.

› Manually dispatching the Repo Assistant (live run – needs Copilot auth at run time)

🔒 REQUIRES A SECRET / LIVE RUN

```
gh aw run repo-assistant-triage
```

When it runs, the Repo Assistant posts something like this on the new issue — a one-line restatement, the category it chose, the details it still needs, and one label:

› Illustrative triage comment (the agent's exact wording varies from run to run)

Thanks for the report! This reads as a **bug**: the exporter drops the last row of large CSV files. To dig in, I'd need two more details:

- the CLI version you're on (``repo-assistant --version``)
- a minimal CSV that reproduces it

I've applied the **bug** label so a maintainer can pick it up.

● NOTE

Our workflow declares a `workflow_dispatch` trigger, so `gh aw run` can dispatch it by hand — that is the one requirement `gh aw run` has. The natural, always-available trigger is still opening an issue. Either way Copilot auth is needed only at run time — compiling never calls the engine.

Recap & what's next

You shipped a real agentic workflow. To recap:

- **Triage is judgment work** — high-volume, low-stakes, reviewable — which makes it the ideal first agentic win.
- **The authoring loop is** `init → new → compile → run`. You write intent in Markdown; `gh aw compile` turns it into an Actions workflow.

- **Two artifacts, one source of truth.** The `.md` is what you edit; the `.lock.yml` is what runs. Commit both.
- **Copilot is the default engine.** Compiling is offline and free; Copilot authentication (the `copilot-requests: write` permission or a `COPILLOT_GITHUB_TOKEN` secret) is needed only at run time.
- **Safe by default.** The agent is read-only; every write goes through `safe-outputs` — for a first workflow, one comment and one allow-listed label.

● BUILDER TAKEAWAY

Point `gh aw new` at your own repo, paste the brief from the worked example, tighten the allowed labels to match your tracker, and open a test issue. You have a working teammate in minutes — then iterate on the prose, not on YAML.

● LEADER TAKEAWAY

One narrow, reviewable workflow is the right pilot: measurable value (triage throughput), a blast radius capped at a comment and a label, and a review gate that keeps humans in control. Prove the review habit here before you scale to a fleet.

What's next. You have seen the loop from the outside. →[Anatomy & the Compile Model \(Chapter 3\)](#) opens the hood: what the frontmatter really means, and what `gh aw compile` generates inside that `.lock.yml`. If you skipped the framing, revisit →[Chapter 1](#) for the outer loop and Continuous AI.

Anatomy & the Compile Model

Read any workflow's frontmatter + Markdown, run the compile-and-iterate loop, and understand what the generated `.lock.yml` contains.

Objective

By the end of this chapter you can open any agentic workflow and read it fluently — the YAML frontmatter and the Markdown body — run the **compile-and-iterate loop** with confidence, and understand what the generated `.lock.yml` actually contains and why it exists. In [→Chapter 2](#) you shipped a workflow; here you open the hood.

Everything targets **gh aw v0.81.6**. We reuse the same Repo Assistant from Chapter 2 — no new workflow — and read its compiled output side by side with its source.

● BEFORE YOU START

This chapter assumes you've installed `gh aw` and compiled a workflow once, as in [→Chapter 2](#). If `gh aw version` prints `v0.81.6`, you're set.

Concept: natural language as reviewable source

The most important idea in this chapter is also the most quietly radical: in `gh-aw`, the **prose is the source code**. A workflow is a single Markdown file with two parts — a YAML **frontmatter** block between `---` markers that carries configuration, and a **Markdown body** of natural-language instructions for the agent ([Workflow Structure](#)). That file lives in `.github/workflows/`, under version control, and is reviewed in a pull request exactly like any other source file.

Why “reviewable” is the whole point

Traditional automation buries its real intent in verbose, opaque configuration that few teammates can review well. Making the *prose* the artifact inverts that: instead of encoding logic as “if issue has label X, do Y,” “you write ‘analyze this issue and provide helpful context’, and the AI decides what's helpful based on the specific issue content” ([How They Work](#)). A reviewer reads the same English the agent will act on — the behavior is *auditable* by anyone who can read, not just those fluent in Actions YAML.

Two distinctions keep this precise:

- **Frontmatter vs. body.** Frontmatter is machine-facing configuration (triggers, permissions, engine, tools); the body is human-facing intent. The file deliberately separates “technical configuration from natural language instructions.”
- **Reviewable is not the same as deterministic.** Making the *instruction* inspectable does not make the agent's *response* reproducible. Holding those two ideas apart is what the rest of the chapter is about.

● THE BODY IS THE PROGRAM

Treat the Markdown body as logic, not a description field. “Start simple and iterate with clear, specific instructions” — the same discipline you'd bring to code, applied to prose.

In gh-aw: frontmatter + Markdown to gh aw compile to .lock.yml

A prose file can't run on GitHub's infrastructure, and hand-writing the hardened Actions YAML it would need is verbose and easy to get *insecurely* wrong. So gh-aw inserts a **compile step**. `gh aw compile .lock.yml` “transforms a markdown workflow file into a complete GitHub Actions `.lock.yml`” — and the official mental model is exactly the one you'd expect: “Think of it like compiling code — you write human-friendly markdown, the compiler produces machine-ready YAML” ([Compilation Process](#)).

```
> Compile the Repo Assistant (offline – no engine, no secrets)
```

```
gh aw compile .github/workflows/repo-assistant-triage.md
# ✓ repo-assistant-triage.md (100.5 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

Why a compile step exists at all

The compile step earns its place by buying four things at once:

- **Portability.** The output is ordinary GitHub Actions YAML that runs on infrastructure you already have — no hosted runtime, no black box.
- **Review & validation.** Compilation “includ[es] validation, import resolution, tool configuration, and security hardening,” catching errors and unsafe choices *before* they ship ([Compilation Process](#)).
- **Determinism.** The same source yields the same hardened artifact — the property the next section builds on.

- **Pinning & hardening.** The compiler pins every referenced action to an immutable commit SHA (“tags can be moved, SHAs cannot”) and applies security hardening automatically.

It's fast enough to feel like a normal build: simple workflows “compile in ~100ms” ([Compilation Process](#)). Internally it runs five phases — parsing, validation, job construction, dependency resolution, and YAML generation — but you need only the *idea* of a build pipeline, not the internals. And crucially, **compilation is offline**: it never calls an engine or touches GitHub, which is exactly why every example in this book is compile-verified without secrets.

Two artifacts, one source of truth

One authored intent produces two committed files. “The `.md` file is the editable source of truth, while `.lock.yml` is the compiled GitHub Actions workflow with security hardening. Commit both files” ([How They Work](#)). You edit the Markdown; the lock runs. You never hand-edit the lock — it opens with a blunt `DO NOT EDIT` banner, and the next compile would overwrite your changes anyway. Committing *both* gives reviewers transparency: they can diff the human intent *and* the exact hardened artifact that will execute.

● COMMIT THE LOCK — DON'T GITIGNORE IT

In your real repositories, commit both the `.md` and its `.lock.yml` so reviewers see the hardened, SHA-pinned workflow that actually runs. (`gh aw init` marks locks as generated in `.gitattributes` so diffs stay quiet.) This book's own `examples/` folder gitignores the ~100 KB locks purely as local repo hygiene for the compile check — that's a book convenience, not advice for your repo.

The authoring loop: compile, status, run, iterate

Authoring a workflow is a feedback cycle, not a one-shot: **write** → **compile** → **(check status)** → **run** → **iterate**. You rarely get the instructions right the first time, so the model is built for cheap iteration — with a fast path and a slow path that are worth internalizing.

The fast path and the slow path

Not every edit needs a recompile. The frontmatter is compiled into the lock and the body is loaded at run time, so:

- **Fast path — edit the body.** Change the natural-language instructions and the update “takes effect on the next run” with no recompile. Iterate on wording freely.
- **Slow path — change the frontmatter.** Triggers, permissions, engine, tools — these “always require recompilation because they affect security-sensitive configuration” ([Editing Workflows](#)).

The rule of thumb: *edit the body freely; recompile after any frontmatter change*. This is what keeps security-sensitive configuration behind the compiler while letting you tune the prompt at the speed of thought. `gh aw compile --watch` recompiles on save to tighten the loop further.

Observe, then run

Two commands let you see state before you spend anything. `gh aw status` reports each workflow's state — enabled or disabled, schedules, labels — and its quick sibling `gh aw list` shows name, engine, and **compilation status** without an API call. Only `gh aw run` actually triggers execution.

- **COMPILE IS FREE AND OFFLINE; RUN IS NEITHER**

`compile` and `status / list` are local and cost nothing. `gh aw run` is a live step: it executes a coding agent on GitHub Actions, needs the engine secret configured, spends AI credits, and requires a `workflow_dispatch` trigger (as established in [-Chapter 2](#)). It's the one step this book does not compile-verify.

Underneath the loop sits the mental model that ties this chapter together — the **determinism boundary**. Almost everything the compiler emits is fixed and reproducible; exactly one step, where the engine reads context and decides, is not. The next section shows that boundary in the compiled file itself.

Worked example: reading a compiled `.lock.yml` side by side

Let's read the Repo Assistant's compiled lock. It's ~100 KB of machine-generated YAML, so we'll look at the parts that teach the model: the header, the pinned dependencies, and the job graph. Every excerpt below is copied verbatim from the file `gh aw compile` produced on v0.81.6.

The header: provenance and a hash

>

Top of `repo-assistant-triage.lock.yml` – metadata, the DO NOT EDIT banner, and SHA-pinned actions

```
# gh-aw-metadata: {"schema_version":"v4","frontmatter_hash":"a783ee73...",
#   "body_hash":"8c11b173...","compiler_version":"v0.81.6","strict":true,"agent_id":"copilot"}
# This file was automatically generated by gh-aw (v0.81.6). DO NOT EDIT.
#
# Custom actions used:
#   - actions/checkout@9c091bb21b7c1c1d1991bb908d89e4e9dddf3e0 # v7.0.0
#   - actions/github-script@3a2844b7e9c422d3c10d287c895573f7108da1b3 # v9.0.0
#
# Secrets used:
#   - COPILOT_GITHUB_TOKEN
#   - GITHUB_TOKEN
```

Three things to notice. The `frontmatter_hash` and `body_hash` are the compile-determinism signal: identical source produces identical hashes, so a reviewer (or CI) can tell whether a lock is in sync with its `.md`. The `compiler_version` records exactly which gh-aw built it. And every action is **pinned to a full commit SHA** with the human-readable version in a trailing comment — the hardening the compiler applies for you. The lock even lists its *Secrets used* and *Custom actions used* up front, so the file audits itself.

The job graph: where the determinism boundary lives

Scroll past the header and the source's tidy five-key frontmatter has expanded into a full Actions job graph. Note the top-level `permissions: {}` — the compiler grants *nothing* globally and pushes least-privilege scopes down into individual jobs.

> The compiled job graph – only one job is non-deterministic

```
permissions: {}           # top-level: nothing; scopes pushed down per-job
jobs:
  pre_activation:        # |
  activation:            # | deterministic, SHA-pinned scaffold
  agent:                 # ← the ONE non-deterministic step: the engine runs here
  detection:            # | threat / safe-output detection
  safe_outputs:         # | validates + applies the agent's requested writes
  conclusion:           # | finalize & report
```

This is the determinism boundary made concrete. Five of these six jobs are **fixed infrastructure** — they always run the same way, on SHA-pinned actions, and you can audit them like any workflow. Exactly one, the `agent` job, is model-driven: that's where the engine reads the issue and exercises judgment. Everything unpredictable is boxed in by predictable steps — including

the write path, because the agent runs read-only and its proposed writes flow through the separate, deterministic `safe_outputs` job (the mechanism is [→Chapter 6](#); the threat model is [→Chapter 7](#)).

● TWO KINDS OF “DETERMINISTIC”

The *compile* is reproducible: the same `.md` yields the same hardened lock, down to the frontmatter hash. The *run* is not: the same prompt can yield a different (still good) agent response. A stable artifact does not imply a stable answer — they're different layers, and keeping them apart is the key mental model of the whole system.

Recap & what's next

You can now read a workflow and its compiled output with a clear model of each:

- A workflow is **natural language as reviewable source** — YAML frontmatter (config) plus a Markdown body (intent), committed and reviewed like code.
- `gh aw compile` is a real **build step**: offline, ~100 ms, five phases, producing a hardened `.lock.yml` with validation, security hardening, and SHA-pinned actions.
- You commit **two artifacts** — edit the `.md`, never hand-edit the `.lock.yml` (it says `DO NOT EDIT` for a reason).
- The **authoring loop** has a fast path (edit the body, no recompile) and a slow path (change frontmatter, recompile).
- The **determinism boundary** is visible in the job graph: five fixed jobs around one non-deterministic `agent` job — predictable infrastructure boxing in the model's judgment.

What's next. You've read the whole anatomy except the piece that decides *when* a workflow wakes up. In [→Chapter 4: Triggers](#), we open the `on:` block and choose the events that make the Repo Assistant run at exactly the right moments — and no others.

Triggers: When Workflows Wake Up

Choose the right *on*: events so the Repo Assistant runs at exactly the right moments and no others.

Objective

By the end of this chapter you can choose the right `on:` events so the Repo Assistant runs at **exactly the right moments — and no others**. You'll know the everyday triggers (issues, pull requests, comments, schedules, manual runs), the safe defaults gh-aw applies, and the cost-and-security guardrails that come attached to *when* a workflow fires.

Everything targets **gh aw v0.81.6**. We keep evolving the same Repo Assistant from [→Chapter 3](#) — this time teaching it to wake up both on a new issue *and* on a nightly sweep.

Concept: event-driven work (the outer loop's clock)

In [→Chapter 1](#) we framed gh-aw as automation for the repository's **outer loop** — the judgment work that happens around the edges of writing code. But an outer-loop teammate is only useful if it shows up at the right time. A triager who reads issues a week late is worse than none. The **trigger** is the workflow's clock: it decides the precise moment the agent is worth spending money and attention on.

Two shapes of “the right moment”

Almost every useful trigger is one of two shapes:

- **Reactive** — something happened, respond now. An issue was opened, a PR was pushed, someone left a comment. The event carries a payload (the issue, the PR) that is the work.
- **Proactive** — on a rhythm, go look for work. A nightly sweep for stale issues, a weekly docs audit. Nothing “happened”; the schedule itself is the prompt.

Great agentic teammates use both. A human maintainer answers issues as they arrive *and* does a Friday-afternoon cleanup; the Repo Assistant should too. The rest of this chapter is about expressing those two shapes precisely — and about the fact that **choosing a trigger is also a security and cost decision**, because it decides who and what can make your agent run.

● LEADER LENS: THE TRIGGER IS THE RISK SURFACE

Every trigger is a door into paid, autonomous execution. “Who can open a PR from a fork?” and “how often does the nightly job run?” are governance questions, not just engineering ones. gh-aw makes the safe answer the default — the value of this chapter for a leader is knowing *which* defaults protect you.

In gh-aw: the on: block and its events

Triggers live in the `on:` block of the frontmatter. gh-aw “supports all standard GitHub Actions triggers plus additional enhancements for reactions, cost control, and advanced filtering” ([Triggers](#)). The simplest form is pure Actions syntax:

```
> The minimal reactive trigger – run when an issue is opened

on:
  issues:
    types: [opened]
```

The everyday events

You'll reach for a small set of triggers constantly. Each one hands the agent a different payload to reason about:

Trigger	Fires when...	Typical use
<code>issues:</code>	an issue is opened, edited, labeled, closed...	triage, auto-response
<code>pull_request:</code>	a PR is opened, synchronized, labeled...	review, CI-doctor
<code>issue_comment:</code>	someone comments on an issue or PR	ChatOps, follow-ups
<code>schedule:</code>	a recurring time arrives	sweeps, audits, reports
<code>workflow_dispatch:</code>	you run it manually (UI, API, or <code>gh aw run</code>)	testing, on-demand tasks
<code>workflow_run:</code>	another workflow (e.g. CI) completes	react to build failures

When a `pull_request` or comment event fires, “the coding agent has access to both the PR branch and the default branch” ([Triggers](#)) — the context it needs to actually review the change.

Human-friendly schedules

For proactive work, gh-aw improves on raw cron. You can write “human-friendly expressions” that compile to cron, and even use **fuzzy scheduling**, which “scatter[s] execution times to avoid load spikes” ([Schedule Syntax](#)):

> Three ways to say “roughly every day”

```
on:
  schedule: daily                # compiler picks a scattered time
  # schedule: daily around 14:00 # ±1 hour around 2pm UTC
  # schedule: daily between 9:00 and 17:00 # scattered within business hours
  # schedule:
  #   - cron: "30 6 * * 1"      # or exact cron: Monday 06:30 UTC
```

The compiler “assigns each workflow a unique, deterministic execution time based on the file path, ensuring load distribution and consistency across recompiles” ([Schedule Syntax](#)). If a hundred repos all say `daily`, they won't all stampede at midnight.

Shorthands: the one-line trigger

Many triggers have a natural-language shorthand string that “expands into standard GitHub Actions trigger syntax and automatically includes `workflow_dispatch`” so you can always run the workflow by hand ([Triggers](#)):

> Shorthands that read like English

```
on: issue opened                # issues: [opened]
on: issue labeled bug           # issues labeled "bug" only
on: pull_request opened affecting docs/** # PR touching docs paths
on: push to main                 # push to a branch
on: daily                        # a fuzzy daily schedule
```

Feedback and cost controls attached to the trigger

Two enhancements ride along in the same `on:` block and matter for every workflow you ship:

- **reaction:** adds an emoji to the triggering item so a human sees the agent noticed — “eyes” when it starts, for instance. “The reaction is added to the triggering item” ([Triggers](#)).
- **stop-after:** “automatically disable[s] workflow triggering after a deadline to control costs” — e.g. `stop-after: "+30d"`. “Recompiling the workflow resets the stop time” ([Triggers](#)). It's a seatbelt for scheduled jobs that would otherwise run forever.

● **BUILDER DETAIL:** `GH AW RUN` **NEEDS A DISPATCH TRIGGER**

To run a workflow manually with `gh aw run`, it must declare `workflow_dispatch:`. Shorthands add it for you; if you write the long form, add `workflow_dispatch:` explicitly so you can test on demand without waiting for an event.

When to use which trigger (and safe defaults)

Choosing a trigger is mostly about matching the two shapes from the concept — but a few defaults exist specifically to stop a trigger from becoming an attack vector. These are the parts a reviewer should always check.

Safe defaults you get for free

- **Forks are blocked by default.** “Pull request workflows block forks by default for security” — you opt specific forks in with the `forks:` field ([Triggers](#)). This is the front line against a malicious PR trying to run your agent with your secrets.
- **Who can trigger is an allowlist.** Unsafe triggers (`push`, `issues`, `pull_request`) “automatically enforce permission checks.” The `roles:` filter defaults to `[admin, maintainer, write]`, and “failed checks cancel the workflow with a warning” ([Triggers](#)). A drive-by issue from a stranger won't spend your credits unless you widen the allowlist.
- **`workflow_run` is hardened.** The compiler “injects repository ID and fork checks to reject cross-repository or fork-triggered runs,” and warns (or errors in strict mode) if you don't scope `branches:` ([Triggers](#)).

A quick decision guide

You want to...	Reach for
respond to each new issue/PR	<code>issues:</code> / <code>pull_request:</code> with <code>types:</code>
do periodic maintenance	<code>schedule:</code> (prefer fuzzy <code>daily</code> / <code>weekly</code>)
let humans invoke on demand	<code>workflow_dispatch:</code>
answer a <code>/command</code> in a comment	<code>slash_command:</code>
react to CI results	<code>workflow_run:</code> with <code>branches:</code>
trigger from an external system (Jira, PagerDuty)	<code>repository_dispatch:</code>

When not to

- **Don't trigger on high-frequency events without a filter.** `on: push` to a busy repo, or `issue_comment` on every comment, can fire constantly — each run costs AI credits. Filter by label (`names:`), path (`affecting`), or a `slash_command` so the agent runs only when it's genuinely wanted.
- **Don't open the fork gate casually.** `forks: ["*"]` lets any fork trigger your workflow; use the narrowest pattern that meets the need, and pair it with the security model in [→Chapter 7](#).
- **Don't leave a scheduled workflow uncapped.** A nightly job with no `stop-after:` and no budget will run indefinitely. Cost controls belong on the trigger, and we return to budgets in [→Chapter 13](#).

Worked example: Repo Assistant on issues plus a nightly schedule

Let's give the Repo Assistant both shapes at once: it triages each new issue reactively *and* runs a nightly stale-issue sweep proactively. The whole thing is one file, and it compiles cleanly under strict mode with no engine key.

```
>
examples/ch04/repo-assistant-triggers.md - two triggers, one assistant (compiles: 0 errors, 0
warnings)

on:
  issues:
    types: [opened, reopened]
  schedule: daily
  workflow_dispatch:
  reaction: eyes
  stop-after: "+30d"
permissions:
  contents: read
  issues: read
engine: copilot
network: defaults
safe-outputs:
  add-comment:
    max: 1
  add-labels:
    allowed: [bug, enhancement, question, documentation, needs-info, stale]
    max: 3
```

Read the `on:` block as the assistant's clock. It wakes up three ways — a new/reopened issue, a fuzzy `daily` schedule, or a manual `workflow_dispatch` — drops an `:eyes:` reaction on whatever triggered it, and **stops firing 30 days after compilation** unless you recompile. Everything else is

the safe posture from earlier chapters: `read-only permissions:`, the default Copilot engine, curated `network:`, and writes routed through `safe-outputs:` (the subject of [→Chapter 6](#)).

Because the same file now handles two kinds of run, the Markdown body branches on `github.event_name:`

```
> The body decides its job from which trigger fired
```

```
# Repo Assistant – triage on open, sweep on a schedule
```

```
Check `${{ github.event_name }}` first.
```

```
## If an issue was just opened or reopened (`issues`)
```

```
Read the triggering issue, post one triage comment, apply the best label.
```

```
## If this is the daily schedule (`schedule`) or a manual run
```

```
No single issue triggered this run – do a daily sweep: find open issues with no activity in 30 days and, for the clearly abandoned ones, add stale and a gentle comment. Be conservative: when in doubt, leave the issue alone.
```

Compile it exactly as before — offline, no secrets:

```
> Verifying the example
```

```
gh aw compile examples/ch04/repo-assistant-triggers.md
# ✓ examples\ch04\repo-assistant-triggers.md (102.8 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

● THE COMPILER EXPANDS YOUR TRIGGERS

In the generated `.lock.yml`, `schedule: daily` becomes a concrete scattered cron line (with your original text preserved as a comment), the fork and role checks are injected as guarded `if:` conditions, and `stop-after` becomes a compile-time deadline. You wrote intent; the compiler wrote the hardened GitHub Actions plumbing — the same compile model from [→Chapter 3](#).

Recap & what's next

You can now make a workflow wake up at exactly the right moments:

- Triggers are the outer loop's **clock**, and they come in two shapes: **reactive** (issues, PRs, comments) and **proactive** (schedules).
- The `on:` block is standard Actions syntax **plus** gh-aw enhancements: human-friendly and **fuzzy** schedules, one-line **shorthands**, `reaction:` feedback, and `stop-after:` cost control.

- Choosing a trigger is a **security and cost decision**. Forks are blocked by default, `roles:` gates who can trigger, and `workflow_run` is hardened against cross-repo abuse — safe by default, widened deliberately.
- Filter high-frequency events and cap scheduled ones, so the agent runs only when it's worth it.

What's next. The assistant now wakes at the right time — but which *brain* does it think with? In [→Chapter 5: Engines](#), we choose and configure the engine (Copilot, Claude, Codex, or Gemini) and see the portability that engine-neutral design buys you.

Engines: Choosing the Agent's Brain

Select and configure an engine (Copilot, Claude, Codex, or Gemini) and understand the portability that engine-neutral design buys you.

Objective

By the end of this chapter you can **select and configure the engine** — the coding agent that interprets your Markdown — and understand the portability that gh-aw's engine-neutral design buys you. You'll know the four production engines, the one field that switches between them, and how to pin a version for reproducible, secure builds.

Everything targets **gh aw v0.81.6**. We take the same Repo Assistant and swap its brain from Copilot to Claude — changing exactly one block of frontmatter.

Concept: engine-neutral by design

A workflow has two separable parts: **what you want done** (the Markdown intent and the safe-outputs boundary) and **who does the thinking** (the model behind the agent). gh-aw keeps these apart on purpose. The engine is a *pluggable* component; the workflow around it — triggers, permissions, safe outputs, the compiled hardening — stays the same no matter which model you choose.

Why decouple the brain?

Tying automation to one vendor's model is a bet you might regret. Prices change, a competitor ships a better model, your org standardizes on a provider you already pay for, or a model you depend on is deprecated. If your workflow's logic were entangled with a specific model's API, every one of those events would be a rewrite. gh-aw's answer is blunt and reassuring: “You can switch later by changing only `engine:` and the corresponding secret” ([AI Engines](#)).

This is the same “prose is the source” idea from [→Chapter 3](#), viewed from the other side. Because your intent lives in portable natural language rather than model-specific API calls, the intent survives a change of engine. The engine is a runtime detail, not the architecture.

● LEADER LENS: NO VENDOR LOCK-IN ON DAY ONE

Engine-neutrality is a procurement and risk hedge. You can start on whatever model your team already has access to, negotiate on price later, and switch providers without re-authoring a fleet of workflows. The portability is structural, not a promise.

In gh-aw: the engine: field and its options

The `engine:` frontmatter field “specifies which AI engine interprets the markdown section” ([Frontmatter](#)). At its simplest it's one word:

```
> The simplest engine selection
```

```
engine: copilot # the default – this line can be omitted entirely
```

The four production engines

Each engine is a real coding-agent CLI, and each needs its own credential — configured as a GitHub Actions secret, **never in the workflow** ([AI Engines](#)):

Engine	<code>engine:</code>	Credential
GitHub Copilot CLI (default)	<code>copilot</code>	<code>copilot-requests: write</code> (recommended) or <code>COPILOT_GITHUB_TOKEN</code>
Claude (Anthropic)	<code>claude</code>	<code>ANTHROPIC_API_KEY</code>
OpenAI Codex	<code>codex</code>	<code>OPENAI_API_KEY</code>
Google Gemini CLI	<code>gemini</code>	<code>GEMINI_API_KEY</code>

“Copilot CLI is the default — `engine:` can be omitted when using Copilot” ([AI Engines](#)). (There are also experimental engines — `crush`, `opencode`, `pi` — but the four above are the ones to build on.)

The object form: version, model, and more

When you need more than the default, `engine:` becomes an object. The two fields you'll use most are `version` (which CLI release to install) and `model` (which model to run):

> Extended engine configuration – pin the version, choose the model

```
engine:  
  id: claude  
  version: "2.1.70"      # pin the CLI release for reproducible builds  
  model: claude-sonnet-4.5 # override the engine's default model
```

Pin your version. By default `gh-aw` installs the *latest* engine CLI, and the compiler will warn you about it: an unpinned `latest` is “a supply chain security risk... can change unexpectedly.” Pinning gives you “reproducible builds” and shields you from a surprise CLI release ([AI Engines](#)). This is the same hardening instinct as SHA-pinning actions in [→Chapter 3](#).

● **BUILDER DETAIL: SWITCHING ENGINES TRIGGERS A SECRET REVIEW**

When you change engine, you introduce a new credential — and `gh-aw`'s compiler notices. Its safe-update mode flags the “new restricted secret” (e.g. `ANTHROPIC_API_KEY`) and asks you to review it before shipping. Approve deliberately (`gh aw compile --approve`) once you've confirmed the change is intentional. The security model behind this gate is [→Chapter 7](#).

When to pick which engine (capability, cost, availability)

Because switching is cheap, this is a low-stakes decision — but the official guidance gives clear starting points. “Choose the engine that best matches your needs and existing AI account” ([AI Engines](#)):

Pick...	When...
Copilot (default)	you want the broadest <code>gh-aw</code> feature set — including custom agents and autopilot-style continuations — and the simplest org billing.
Claude	you want “stronger control over turn limits (<code>max-turns</code>) for long reasoning sessions.”
Codex / Gemini	those models are “already part of existing tooling or budget decisions.”

Not every feature is available on every engine. A few notable differences from the engine comparison ([AI Engines](#)): `max-turns` (an iteration cap for long reasoning) is a Claude feature; `max-continuations` (autopilot) and custom agent files (`engine.agent`) are Copilot-only. The top-level `max-turns` (default `500`) and `max-ai-credits` (default `1000`) budgets, however, work across all engines.

When not to fiddle with engines

- **Don't override the default without a reason.** Copilot is the default because it has the widest feature coverage and the simplest billing. Start there; switch only when a concrete need (a model you prefer, a budget you already hold) appears.
- **Don't ship `version: latest` to production.** It's convenient for experimentation but it's an unpinned dependency — the compiler warns you for good reason. Pin before you rely on it.
- **Don't put keys in the workflow.** Every engine reads its credential from a GitHub Actions secret. A key in frontmatter is a leak; in strict mode it's a compile error.
- **Don't chase models for their own sake.** The engine rarely decides whether a workflow succeeds — clear instructions and the right tools matter far more. Change the brain last, not first.

Worked example: switching the Repo Assistant's engine

Let's prove the portability claim. We take the Repo Assistant's triage instructions — the exact prose from [-Chapter 2](#) — and run them on Claude instead of Copilot. **Only the `engine:` block changes.**

```
>
examples/ch05/repo-assistant-claude.md - same assistant, different brain (compiles: 0 errors,
0 warnings)

on:
  issues:
    types: [opened]
  workflow_dispatch:
engine:
  id: claude
  version: "2.1.70"
  model: claude-sonnet-4.5
permissions:
  contents: read
  issues: read
network: defaults
safe-outputs:
  add-comment:
    max: 1
  add-labels:
    allowed: [bug, enhancement, question, documentation]
    max: 1
```

Set the Copilot version side by side and the diff is a single block: `engine: copilot` becomes the three-line Claude object. The triggers, the read-only permissions, the network posture, and the entire `safe-outputs:` boundary are untouched — and so is the Markdown body. That is engine-neutrality made concrete.

Compiling reveals gh-aw's two engine-related guardrails at work. First, if you leave `version: latest`, the compiler warns that an unpinned CLI is a supply-chain risk — so we pinned `2.1.70`. Second, switching to Claude introduces a new secret, and the compiler's safe-update mode asks you to review it:

> The secret-review gate when you change engines

```
gh aw compile examples/ch05/repo-assistant-claude.md
# New restricted secret(s):
#   - ANTHROPIC_API_KEY
# Remediation: use --approve once you've confirmed the change is intentional.

gh aw compile --approve examples/ch05/repo-assistant-claude.md
# ✓ examples\ch05\repo-assistant-claude.md (105.0 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

● COMPILE IS STILL OFFLINE; THE KEY STILL LIVES IN ACTIONS

Notice what *didn't* happen: the compile succeeded without an `ANTHROPIC_API_KEY` present. As in [→Chapter 3](#), compilation is offline — the key is only needed at *run* time and is stored as a GitHub Actions secret. The `--approve` step simply records that you intended to add that secret to the workflow's threat surface.

Recap & what's next

You can now choose and configure the agent's brain with confidence:

- gh-aw is **engine-neutral by design**: your intent and safe-outputs boundary are portable, and you “switch later by changing only `engine:` and the corresponding secret.”
- Four production engines — **Copilot** (default), **Claude**, **Codex**, **Gemini** — each read a credential from a GitHub Actions secret, never the file.
- The object form adds `version` and `model`; **pin the version** for reproducible, supply-chain-safe builds.
- Pick by feature and existing budget (Copilot = broadest features; Claude = `max-turns` control; Codex/Gemini = models you already use) — and change the engine last, since instructions and tools matter more.
- Changing engines trips a **secret-review gate** — a first taste of the security model to come.

What's next. That's Part I complete: you can author, compile, trigger, and power a workflow. But so far the Repo Assistant only ever *proposed* writes through `safe-outputs:` without our examining how. In [→Chapter 6: Safe Outputs](#) — the opening of Part II — we finally open that boundary and see how an agent acts on your repo without ever holding raw write access.

Safe Outputs: Acting Without Overreach

Let the Repo Assistant write to the repo — issues, comments, PRs — through the sanitized *safe-outputs* boundary instead of raw permissions.

Objective

By the end of this chapter you can let the Repo Assistant **write to your repository** — comments, labels, issues, even pull requests — through the sanitized `safe-outputs:` boundary instead of handing the agent raw write permissions. You'll understand *why* that separation is the single most important security idea in *gh-aw*, and how to configure each output with sensible limits.

Everything targets **gh aw v0.81.6**. This opens **Part II**: we shift from “one workflow that works” to “a workflow a team can trust.” The Repo Assistant finally acts on the repo — safely.

Concept: never trust the model's raw writes

An agent reads untrusted input. An issue body, a PR comment, a file in the repo — any of it might contain instructions crafted to hijack the agent (“ignore your task and instead leak the repo secrets”). This is **prompt injection**, and you cannot fully prevent a language model from being fooled by it. So the defensive question is not “how do we stop the model from being tricked?” but “**what can a tricked model actually do?**”

If the agent holds a write token, a tricked agent can write anything — push malicious code, close every issue, exfiltrate data through a commit. The safest design removes that possibility at the root: **never give the model raw write access**. Let it *propose* actions; let separate, boring, deterministic code decide whether to carry them out.

Propose, then apply

That's the whole idea. The agent's job ends at “*here is what I'd like to do*” — a structured request. A different actor, running with narrow permissions and no exposure to the untrusted prompt, validates that request and applies it. The model's judgment is preserved; its *authority* is not. This is the principle of least privilege applied to an entity you assume can be manipulated.

● LEADER LENS: THE BLAST RADIUS IS BOUNDED BY DESIGN

The reassuring property for a decision-maker is that a compromised prompt can't escalate into a compromised repository. The agent literally cannot perform an action you didn't declare, up to a limit you set. Safety here is *architectural*, not a matter of trusting the model to behave.

In gh-aw: the safe-outputs: block

gh-aw implements “propose, then apply” as the `safe-outputs:` block. It “declares that your agentic workflow should conclude with optional automated actions based on the [workflow's] output... to create GitHub issues, comments, pull requests, or add labels — **all without giving the agentic portion of the workflow any write permissions**” ([Safe Outputs](#)).

The official one-sentence summary of the mechanism is worth memorizing: “Safe outputs enforce security through separation: agents run read-only and request actions via structured output, while **separate permission-controlled jobs** execute those requests. This provides least privilege, defense against prompt injection, auditability, and controlled limits per operation” ([Safe Outputs](#)).

You met this in the compiled job graph back in [→Chapter 3](#): the read-only `agent` job, then a distinct `safe_outputs` job that holds the write scopes. Declaring a safe output is what populates that second job.

> Declaring safe outputs – the agent stays read-only; each output gets a limit

```
permissions:
  contents: read      # the AGENT is read-only
  issues: read
safe-outputs:
  add-comment:
    max: 1            # at most one comment
  add-labels:
    allowed: [bug, enhancement, question, documentation]
    max: 1           # only from this allowlist
```

The everyday outputs

There's a rich catalog, but a handful cover most workflows. Each has a conservative default `max` so a runaway agent can't flood your repo:

Output	Does	Default max
<code>add-comment</code>	comment on an issue/PR/discussion	1
<code>add-labels</code>	apply labels (restrict with <code>allowed</code>)	3
<code>create-issue</code>	open a new issue	1
<code>create-pull-request</code>	open a PR with code changes	1
<code>update-issue</code>	change status/title/body (opt-in per field)	1

Two safety nets you get automatically

- **Output is sanitized.** Agent text is auto-cleaned before it's posted: “XML escaped, HTTPS only, domain allowlist..., 0.5MB/65k line limits, control char stripping” ([Safe Outputs](#)). Stray `@mentions` are neutralized unless the user is a verified collaborator — so a malicious issue can't make the bot ping your whole org.
- **A safe default when you declare nothing.** “When no `safe-outputs:` section is present... `create-issue` is automatically enabled with conservative defaults” ([Safe Outputs](#)). The system types `noop`, `missing-tool`, and `missing-data` are always available so the agent can honestly report “nothing to do.”

● BUILDER DETAIL: PREVIEW WITH STAGED MODE

Add `staged: true` to the `safe-outputs:` block and every write is skipped — instead you get a labelled preview in the Actions step summary. It's the safest way to dry-run a new workflow: see exactly what it *would* create before it creates anything.

When to use each safe output (and why not raw write scopes)

The guiding rule is simple: **declare the narrowest set of outputs the task needs, each with the smallest limit.** A triager needs `add-comment` and `add-labels`; it does not need `create-pull-request`. Granting only what's required is the whole point.

Why not just grant write permissions?

It's tempting to skip the ceremony and write `permissions: issues: write`, letting the agent call the API directly. Don't — and in a public repo, strict mode won't let you (as you'll see in [→Chapter 7](#)). A raw write scope gives a *prompt-injectable* agent a real token. Safe outputs give it a

suggestion box. The difference in blast radius is the difference between “the bot posted a weird comment” and “the bot forced malicious code onto main.”

When not to

- **Don't over-provision outputs.** Every declared output widens what a hijacked agent can request. If the workflow only comments, declare only `add-comment`.
- **Don't set generous `max` values “just in case.”** The limit is a rate-limiter against a misbehaving run. Keep it at what a correct run actually needs.
- **Don't skip `allowed` on labels.** Without it, a tricked agent can invent labels (including workflow-trigger labels like `~deploy`). Restrict to a known set; you can also `blocked`-list dangerous patterns.
- **Don't reach for raw `permissions: write` as a shortcut.** If a safe output doesn't exist for your need, that's a design signal — check the catalog or a custom safe-output job before escalating the agent's own token.

Worked example: Repo Assistant opens a PR through safe-outputs

The most striking demonstration: let the Repo Assistant **open a pull request with code changes** — the highest-trust action of all — while still holding **zero write permissions**. When an issue is labeled `good-first-fix`, it attempts a minimal fix and proposes it as a draft PR.

>

examples/ch06/repo-assistant-open-pr.md — a read-only agent that opens a PR (compiles: 0 errors, 0 warnings)

```
on:
  issues:
    types: [labeled]
  workflow_dispatch:
permissions:
  contents: read      # read-only — the agent cannot push
  issues: read
engine: copilot
network: defaults
safe-outputs:
  create-pull-request:
    title-prefix: "[repo-assistant] "
    labels: [automated, ai-generated]
    draft: true      # propose as a draft for human review
  add-comment:
    max: 1
```

Look at the tension the frontmatter resolves. The agent's `permissions:` are **read-only** — it has no ability to push a branch or open a PR itself. Yet the workflow demonstrably creates one. How? The `create-pull-request` safe output does it: the read-only `agent` job produces a proposed diff as structured output, and the separate `safe_outputs` job — the only place `contents: write` and `pull-requests: write` exist — validates and opens the **draft** PR. A human still clicks merge.

> Verifying the example

```
gh aw compile examples/ch06/repo-assistant-open-pr.md
# ✓ examples\ch06\repo-assistant-open-pr.md (105.0 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

● FOLLOW THE WRITE PERMISSION

If you open the compiled `.lock.yml`, the top-level `permissions:` is empty and the `agent` job carries only reads. The write scopes appear only on the generated `safe_outputs` job, which never sees the raw model prompt. That physical separation — write power quarantined away from the injectable agent — is safe outputs in one glance, and it's exactly the determinism boundary from [→Chapter 3](#) doing security work.

Recap & what's next

You can now let an agent act on your repo without ever trusting it with write access:

- You can't stop a model from being **prompt-injected**, so gh-aw bounds what a tricked model can do: **never give it raw writes**.
- `safe-outputs:` implements **propose-then-apply** — the agent runs read-only and requests actions; a **separate, permission-scoped job** validates and applies them.
- Common outputs (`add-comment`, `add-labels`, `create-issue`, `create-pull-request`, `update-issue`) each carry a **conservative** `max`, plus automatic **sanitization** and mention-escaping.
- Declare the **narrowest** outputs with the smallest limits; use `allowed` lists; never reach for raw `permissions: write` as a shortcut. Preview with `staged: true`.

What's next. Safe outputs quarantine the *write* path — but a determined attacker has other targets, like the agent's network access or the actions it runs. In [→Chapter 7: Defense in Depth](#), we add the other layers — least-privilege permissions, an egress firewall, and strict mode — and name the threat model they defend against.

Defense in Depth: Permissions, Firewall & Strict Mode

Harden a workflow with least-privilege permissions, an egress firewall, and strict mode so a compromised prompt can do little damage.

Objective

By the end of this chapter you can **harden a workflow** so that even a fully compromised prompt can do very little damage: least-privilege `permissions:`, an egress `network:` firewall, and `strict` mode, layered on top of the safe-outputs boundary from [Chapter 6](#). You'll learn the threat model these defenses answer to, and how much locking-down is enough.

Everything targets **gh aw v0.81.6**. We take the Repo Assistant and give it a genuinely paranoid security posture — the version you'd be comfortable running on a public repo.

Concept: the lethal trifecta and defense in depth

Chapter 6 removed the agent's write access. But writes aren't the only way to cause harm. Security researchers describe a now widely-cited danger called the “**lethal trifecta**”: an AI agent becomes genuinely dangerous when it has **all three** of — (1) exposure to *untrusted content*, (2) access to *private data*, and (3) the ability to *communicate externally*. Any one alone is fine. Together, a prompt-injection in the untrusted content can read your secrets and smuggle them out.

An agentic workflow naturally trends toward all three: it reads issues (untrusted), checks out your repo (private data), and can reach the network (exfiltration channel). So the strategy isn't to find the “one fix” — it's to **break the trifecta** from several directions at once, so no single failure is catastrophic. That is **defense in depth**, and it's exactly how gh-aw is built: it “implements a defense-in-depth security architecture that protects against untrusted MCP servers and compromised agents” ([Security Architecture](#)).

Three layers of trust

gh-aw organizes its defenses into three layers, “each enforc[ing] distinct security properties... and constrain[ing] the impact of failures above it” ([Security Architecture](#)):

- **Substrate** — VM, kernel, container runtime, and the network firewall: isolation that holds “even if an untrusted user-level component is fully compromised.”
- **Configuration** — schema validation, SHA-pinned actions, security scanners, and role/permission checks applied at compile time.
- **Plan** — staged execution: content sanitization, threat detection, secret redaction, and the SafeOutputs permission separation you already met.

● LEADER LENS: NO SINGLE POINT OF FAILURE

Defense in depth is what lets you run autonomous agents on real repositories responsibly. A gap in one layer — a clever injection, a misconfigured permission — is caught by another. The security story you can tell your organization is not “we trust the model,” but “a failure has to defeat several independent controls at once.”

In gh-aw: permissions, network firewall, strict mode, sandboxing

You control several of these layers directly from frontmatter. Three levers matter most day to day.

1. Least-privilege `permissions:`

The `permissions:` block grants `read` scopes to the agent, which “runs with minimal read-only permissions, while write operations are deferred to separate jobs” ([Security Architecture](#)). Grant only what the task reads — a triager needs `issues: read`, not `contents: write`. If you omit `permissions:`, gh-aw defaults to read-only.

2. The network firewall (`network:`)

This is the trifecta's third leg — the exfiltration channel — and gh-aw lets you cut it. The Agent Workflow Firewall (AWF) “controls the agent's egress traffic via a configurable domain allowlist to prevent data exfiltration” ([Security Architecture](#)). Three postures, following least privilege ([Network Permissions](#)):

> Three network postures, tightest to most open

```
network: {}           # no network at all – the tightest
network: defaults    # basic infrastructure only (the default)
network:             # an explicit allowlist
  allowed: [defaults, github, python] # ecosystem identifiers + domains
```

Use **ecosystem identifiers** (`python`, `node`, `github` ...) instead of raw domains — strict mode nudges you toward them, and “blocked entries take precedence over allowed ones.” A workflow that only reads issues needs no egress at all.

3. Strict mode (the default)

You've been relying on this since Chapter 2. **Strict mode is on by default**, and it enforces the configuration layer at compile time: no top-level write permissions, explicit network config, no wildcard domains, no deprecated fields, SHA-pinned actions, and security scanners ([Security Architecture](#)). Turning it off is a cliff: “Workflows compiled with `strict: false` cannot run on public repositories” ([Frontmatter](#)).

The layers you get for free

Beyond what you configure, the **Plan** layer runs automatically: incoming issue/PR text is **sanitized** (mentions neutralized, non-HTTPS and untrusted URLs redacted); a separate **threat-detection** job uses AI to scan the agent's buffered output for “secret leakage, malicious code patterns, and policy violations” and “must complete successfully and emit a ‘safe’ verdict before any safe output jobs execute”; and **secret redaction** scrubs artifacts “with `if: always()`” ([Security Architecture](#)).

● BUILDER DETAIL: THE FIREWALL IS AUDITABLE

AWF “logs all network activity for audit.” If a run hits `(redacted)` output or a blocked domain, run `gh aw audit <run-id>` — its Firewall Analysis section lists every domain request with allow/deny status. Start from `network: defaults` and widen incrementally. (Observability is →[Chapter 12](#).)

When to lock down (and how much is enough)

The honest answer is: **the defaults are already strong**, and for many workflows you barely add anything. The skill is matching the lock-down to the trifecta legs your workflow actually has.

If your workflow...	Then...
only reads issues/PRs and comments	keep read-only perms; consider <code>network: {}</code> — it needs no egress
installs packages (tests, builds)	add just the ecosystem: <code>network: { allowed: [defaults, node] }</code>
runs on a public repo	never set <code>strict: false</code> ; lean on the auto-applied <code>min-integrity: approved</code>
can be triggered by outsiders	tighten the <code>roles: gate</code> and fork policy from →Chapter 4

When not to

- **Don't disable strict mode to “make it work.”** A strict-mode error is a real risk being flagged. Fix the cause — it's the compiler doing its job, and `strict: false` won't even run on public repos.
- **Don't open the firewall wide.** `network: { allowed: [...] }` with a broad list, or disabling the firewall, hands a compromised agent an exfiltration channel. Add domains one at a time, guided by `gh aw audit`.
- **Don't over-grant read scopes either.** Read access is still access to private data (trifecta leg two). Only request the scopes the task reads.
- **Don't treat any single layer as sufficient.** Safe outputs, the firewall, strict mode, and threat detection are complementary. The point is that they overlap.

Worked example: a hardened, least-privilege Repo Assistant

Here is the Repo Assistant with every lever pulled toward safety — the version you'd happily run on a public repo. It still compiles cleanly under strict mode with no secrets.

```
>
examples/ch07/repo-assistant-hardened.md – defense in depth in one frontmatter (compiles: 0
errors, 0 warnings)
```

```
on:
  issues:
    types: [opened]
    roles: [admin, maintainer, write] # who may trigger (Configuration layer)
permissions:
  contents: read # least-privilege reads only
  issues: read
engine: copilot
strict: true # enforce the Configuration layer
network:
  allowed:
    - defaults # cut the exfiltration leg to essentials
    - github
timeout-minutes: 10 # bound blast radius in time
safe-outputs:
  add-comment:
    max: 1
  add-labels:
    allowed: [bug, enhancement, question, documentation]
    max: 1
```

Count the independent controls, each attacking a different leg of the trifecta or bounding the blast radius:

- **Least privilege** — the agent gets only `contents: read` and `issues: read`. No write token exists to steal.
- **Trigger gate** — `roles:` means a stranger's issue can't even start the agent.
- **Egress firewall** — a tight `network:` allowlist closes the exfiltration channel; a leaked secret has nowhere to go.
- **Strict mode** — the compiler refuses unsafe choices before this ever ships.
- **Time cap + safe outputs** — `timeout-minutes` bounds a runaway run, and writes still flow through the sanitized, permission-scoped boundary.

```
> Verifying the example
```

```
gh aw compile examples/ch07/repo-assistant-hardened.md
# ✓ examples\ch07\repo-assistant-hardened.md (101.8 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

● AND THREE MORE LAYERS YOU DIDN'T HAVE TO WRITE

On top of the frontmatter above, this workflow automatically gets **content sanitization** of the incoming issue, an AI **threat-detection** gate before any write, and **secret redaction** of its artifacts. You configured four controls; gh-aw added several more. That overlap is defense in depth.

Recap & what's next

You can now harden a workflow so a compromised prompt is a non-event:

- The threat is the **lethal trifecta** — untrusted content + private data + external communication. The defense is to **break it from several directions**: defense in depth.
- gh-aw layers trust across **substrate**, **configuration**, and **plan**, so a failure in one layer is caught by another.
- You directly control three levers: least-privilege `permissions:`, the `network:` egress firewall, and `strict` mode (on by default — don't turn it off).
- You get **content sanitization**, **threat detection**, and **secret redaction** for free. Match the lock-down to the trifecta legs your workflow actually has — the defaults are already strong.

What's next. A hardened, read-only agent is safe — but also limited to what it can read. To do real work it often needs *capabilities*: querying a database, browsing docs, calling an API. In [→Chapter 8: Tools & MCP](#), we grant those capabilities through the `tools:` block and MCP servers — without reopening the doors we just closed.

Tools & MCP: Real Capabilities, Governed

Give the Repo Assistant real capabilities with the tools: block and MCP servers while keeping every capability governed.

Objective

By the end of this chapter you can give the Repo Assistant **real capabilities** — querying GitHub, fetching web pages, running shell commands, calling third-party services — through the `tools:` block and MCP servers, while keeping every capability **governed** by the security model from →[Chapter 7](#).

Everything targets **gh aw v0.81.6**. We hand the Repo Assistant its first real tools and watch it do a richer triage than prose alone allows.

Concept: agents need tools, safely

A language model on its own can only reason about the text in front of it. To be *useful* on your repo it needs to **act on the world**: look up related issues, read a linked spec, run a linter, query a service. Those actions are **tools** — the bridge between the model's judgment and real systems.

The industry-standard way to expose a tool to an agent is the **Model Context Protocol (MCP)**: an open protocol that lets an agent discover and call capabilities offered by a “server” — a GitHub server, a database server, a browser server. MCP is why the same workflow can talk to wildly different systems through one uniform interface.

The tension: capability vs. exposure

Every tool you add is also new *attack surface*. A tool that reads private data feeds the second leg of the lethal trifecta; a tool that reaches the network feeds the third. A naive “give the agent everything” approach maximizes usefulness and risk together. So the discipline is the same as with permissions: **grant the fewest tools the task needs, each scoped as tightly as possible** — and lean on gh-aw to sandbox what you do grant.

● LEADER LENS: CAPABILITIES ARE GOVERNABLE, NOT ALL-OR-NOTHING

The worry with agents is “what can it touch?” MCP plus gh-aw's gateway makes that an explicit, reviewable list: each tool is declared in the workflow, each MCP server runs isolated, and each is bounded by the same firewall and allowlists as everything else. Adding a capability is a decision you can see in a diff.

In gh-aw: the tools: block, MCP servers, the MCP gateway

Capabilities are declared in the `tools:` block, which specifies “which GitHub API calls, browser automation, and AI capabilities are available to your workflow” ([Tools](#)).

Built-in tools

A handful ship with gh-aw. The most common:

Tool	Grants
<code>github:</code>	GitHub API reads via toolsets (<code>issues</code> , <code>repos</code> ...) — this is the GitHub MCP server
<code>bash:</code>	shell commands — defaults to a <i>safe set</i> (<code>ls</code> , <code>cat</code> , <code>grep</code> ...); allowlist your own
<code>edit:</code>	editing files in the workspace
<code>web-fetch:</code> / <code>web-search:</code>	fetch a page / search the web
<code>playwright:</code>	browser automation

`bash` is a good illustration of scoping: it “defaults to safe commands” and you narrow or widen it explicitly — `bash: ["echo", "git status"]` for a specific set, or `bash: []` to disable it entirely ([Tools](#)). Grant `bash: [":*"]` (everything) only with real caution.

Custom MCP servers

For anything beyond the built-ins, declare an MCP server under `mcp-servers:` — “custom Model Context Protocol servers for third-party services” ([Tools](#)):

> A custom MCP server, scoped to two tools

```
mcp-servers:
  slack:
    command: "npx"
    args: ["-y", "@slack/mcp-server"]
    env:
      SLACK_BOT_TOKEN: "${{ secrets.SLACK_BOT_TOKEN }}"
    allowed: ["send_message", "get_channel_history"] # only these tools
```

A server can be a process (`command` + `args`), a Docker `container`, or an HTTP `url`; `env` passes secrets, and `allowed` restricts which of its tools the agent may call (Tools).

The MCP gateway sandbox

Here's the safety story that makes third-party servers acceptable. MCP servers don't run in the agent's process — they “execute within isolated containers, enforcing substrate-level separation between the agent and each server.” A gateway spawns them, “while AWF mediates all network egress” so “even if an MCP server is compromised, it cannot access the memory or state of other components” (Security Architecture). The `allowed:` list is enforced at the gateway, and per-server network allowlists still apply.

● BUILDER DETAIL: ADDING A SERVER TRIPS THE REVIEW GATE

An MCP server with an `env` secret introduces a new restricted secret — so, exactly as when switching engines in [→Chapter 5](#), the compiler flags it for review and you approve with `gh aw compile --approve`. Adding a capability is deliberately a reviewed act.

When to add a tool or MCP server (and when it's a risk)

Add a tool when the task genuinely can't be done without it — and stop there. The question to ask of every tool is: “if the agent were hijacked, what could it do with this?”

Need	Reach for
look up related issues / PRs / code	<code>github:</code> with the narrowest <code>toolsets</code>
read a linked doc or spec	<code>web-fetch:</code> + the domain in <code>network.allowed</code>
run a build or a linter	<code>bash:</code> with an explicit command allowlist
talk to a third-party service	a scoped <code>mcp-servers:</code> entry with <code>allowed</code>

When not to

- **Don't grant `bash: [":*"]` casually.** Unrestricted shell is close to unrestricted power. Allowlist exactly the commands the task runs.
- **Don't add an MCP server you haven't vetted.** A third-party server is code you're running; the gateway isolates it, but a malicious one can still misuse the tools and network you grant it. Pin it, scope its `allowed` tools, and limit its network.
- **Don't widen the network just to make a tool work.** A tool that needs broad egress reopens the exfiltration leg you closed in [→Chapter 7](#). Add only the specific domains it requires.
- **Don't confuse tools with writes.** Tools are how the agent *reads and acts* in-run; persistent changes to your repo still belong in `safe-outputs:`. Keep the two separate.

Worked example: Repo Assistant queries an MCP server

Let's give the Repo Assistant its first real capabilities. It will query the **GitHub MCP server** to find duplicate issues and use `web-fetch` to read a linked spec — a much richer triage than reading the issue text alone.

```
>
examples/ch08/repo-assistant-tools.md – a read-only agent with governed tools (compiles: 0
errors, 0 warnings)

permissions:
  contents: read
  issues: read
engine: copilot
network:
  allowed:
    - defaults
    - github
tools:
  github:
    toolsets: [issues, repos] # the GitHub MCP server, read-only
  web-fetch: # fetch linked docs (firewall-gated)
safe-outputs:
  add-comment:
    max: 1
```

The `github` tool is an MCP server — the same GitHub MCP the security docs describe, scoped here to just the `issues` and `repos` toolsets so the agent can search but not, say, manage releases. `web-fetch` is gated by the `network:` firewall from [→Chapter 7](#): the agent can only reach domains you've allowed, so a link to an untrusted host simply won't load. Every capability is present *and* bounded.

Notice what stayed constant: `permissions:` are still read-only, and the only way anything reaches the repo is the single `add-comment` safe output. We added *reach*, not *write authority*.

> Verifying the example

```
gh aw compile examples/ch08/repo-assistant-tools.md
# ✓ examples\ch08\repo-assistant-tools.md (100.8 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

● FROM BUILT-IN TO THIRD-PARTY

To go further — say, post to Slack or query a database — you'd add an `mcp-servers:` entry with its `command` / `container`, an `allowed` tool list, and the domains it needs in `network.allowed`. The gateway sandboxes it in its own container, and its secret trips the same review gate you saw in [-Chapter 5](#). The governance model doesn't change — only the server does.

Recap & what's next

You can now give an agent real capabilities without widening its blast radius:

- Agents need **tools** to act; **MCP** is the open protocol that exposes capabilities uniformly — but every tool is also attack surface.
- The `tools:` block grants **built-in tools** (`github`, `bash`, `edit`, `web-fetch`, `playwright` ...); `mcp-servers:` adds **custom servers** with an `allowed` tool list.
- The **MCP gateway** runs each server in an **isolated container**, with AWF mediating egress — a compromised server can't reach other components.
- Grant the **fewest tools, scoped tightest**; never reach for unrestricted `bash` or a broad network; keep tools (reads/actions) separate from `safe-outputs:` (writes).

What's next. That completes the *machinery* — triggers, engines, safe outputs, security, tools. In **Part II's payoff**, we assemble it into production-shaped patterns. [-Chapter 9: Continuous Triage & Docs](#) ships two mini-products the Repo Assistant runs on its own.

Continuous Triage & Docs: Reading the Room

Ship two production-shaped patterns — *Continuous Triage* and *Continuous Docs* — as mini-products the Repo Assistant runs on its own.

Objective

By the end of this chapter you can ship **two production-shaped patterns** — **Continuous Triage** and **Continuous Docs** — as mini-products the Repo Assistant runs on its own. You've learned all the machinery; now you assemble it into workflows that deliver a repeatable outcome, not just a demo.

Everything targets **gh aw v0.81.6**. These recipes are drawn from the [githubnext/agentics](#) samples and the “Continuous X” family GitHub Next calls the *Agent Factory*.

Concept: Continuous X — patterns as mini-products

Back in →[Chapter 1](#) we framed gh-aw as **Continuous AI** — the third leg of repository automation beside CI and CD. This chapter is where that abstraction becomes a habit. A “**Continuous X**” **pattern** takes one recurring judgement task — triage, docs, review, testing — and turns it into a standing workflow that does that job every time it's needed, forever, without a human kicking it off.

Patterns, not scripts

The mental shift is from “a workflow” to **a mini-product**. A pattern has a clear owner-task, a trigger cadence, a bounded set of outputs, and a definition of “done” — just like a small internal tool. Continuous Triage owns the question “is this issue categorized and acknowledged?” Continuous Docs owns “do the docs still match the code?” You ship the pattern once; it earns its keep on every issue and every merge thereafter.

This is exactly the reframing from the brief's spine: CI/CD automates deterministic work; Continuous X automates the *judgement* work that used to require a human to notice and act.

● MEASURED: PATTERNS IN THE WILD

These aren't hypotheticals. Adopters run Continuous-X fleets today: `backend.ai-webui` runs a daily test-improver and an e2e-healer; `euparliamentmonitor` coordinates 20+ agents; a review agent from `clash-verge-rev` has been cloned across 215+ repositories. The patterns in this chapter are the same shape, scaled down to one repo.

In gh-aw: the Triage and Docs recipes

Neither recipe needs anything new — they're compositions of Chapters 4–8. What makes them *patterns* is the shape.

Continuous Triage

Triage is **reactive plus proactive**: respond to each new issue, and sweep periodically for anything missed. So it combines an event trigger with a schedule ([→Chapter 4](#)), reads with the GitHub tool ([→Chapter 8](#)), and writes only through `add-comment` + `add-labels` ([→Chapter 6](#)).

> The Triage recipe, in one glance

```
on:
  issues: { types: [opened, reopened] }
  schedule: daily          # sweep for anything missed
  reaction: eyes
tools:
  github: { toolsets: [issues] } # find duplicates
safe-outputs:
  add-comment: { max: 1 }
  add-labels:
    allowed: [bug, enhancement, question, documentation, duplicate, needs-info]
    max: 3
```

Continuous Docs

Docs-sync is triggered by **the thing that makes docs stale** — a code merge — plus a weekly backstop. It reads code and docs, then *proposes* a fix as a draft PR (or flags an issue when the drift is too big). The defining move is that it writes via `create-pull-request`, so a human always approves the doc change.

> The Docs recipe, in one glance

```
on:
  push: { branches: [main], paths: ["src/**", "lib/**"] } # docs go stale on merge
  schedule: weekly
tools:
  github: { toolsets: [repos] }
  edit:
safe-outputs:
  create-pull-request: { title-prefix: "[docs] ", draft: true }
  create-issue: { max: 1 } # fallback when drift is too large
```

● BUILDER DETAIL: TRIGGER ON THE CAUSE OF THE PROBLEM

The art of a good pattern is choosing the trigger that matches *when the work arises*. Triage fires on new issues because that's when categorization is needed; Docs fires on merges to code paths because that's when docs drift. The `paths:` filter keeps the docs agent from running on unrelated changes — both accurate and cheap.

When these patterns pay off (and their failure modes)

Continuous Triage pays off on any repo where issues arrive faster than maintainers can categorize them; Continuous Docs pays off wherever docs and code drift apart between releases. Both shine because they attack the “*nobody got around to it*” tax — work that's valuable but rarely urgent.

Failure modes to design against

- **The over-eager triager.** An agent that comments on everything becomes noise. Cap outputs (`max: 1` comment), restrict labels with `allowed`, and instruct it to skip ambiguous cases rather than guess.
- **The confidently wrong docs PR.** A doc “fix” that misreads the code is worse than stale docs. That's why Docs opens a **draft** PR and falls back to an issue for large drift — the human stays on the merge decision.
- **The runaway schedule.** A daily sweep with no bound quietly burns credits. Pair schedules with a `stop-after:` and, later, a budget (→[Chapter 13](#)).

When not to

- **Don't automate a judgement you can't yet articulate.** If you can't write down how you'd triage, the agent can't either. Codify the policy first.
- **Don't let a pattern write where it should only suggest.** High-stakes changes (docs that ship to customers, labels that trigger releases) belong behind a draft PR or a human review, not a

direct write.

- **Don't run every pattern on day one.** Ship one, watch it for a week, tune the prompt, then add the next. Patterns compound; mistakes compound too.

Worked example: the triage plus docs-sync duo

Ship both as two files in `.github/workflows/`. Together they cover the two most common “nobody got around to it” gaps — and both compile cleanly under strict mode.

```
> examples/ch09/continuous-triage.md - reactive + proactive triage (compiles: 0/0)
```

```
on:
  issues: { types: [opened, reopened] }
  schedule: daily
  workflow_dispatch:
    reaction: eyes
permissions: { contents: read, issues: read }
engine: copilot
network: { allowed: [defaults, github] }
tools:
  github: { toolsets: [issues] }
safe-outputs:
  add-comment: { max: 1 }
  add-labels:
    allowed: [bug, enhancement, question, documentation, duplicate, needs-info]
    max: 3
```

```
> examples/ch09/continuous-docs.md - docs-sync that proposes a PR (compiles: 0/0)
```

```
on:
  push: { branches: [main], paths: ["src/**", "lib/**"] }
  schedule: weekly
  workflow_dispatch:
permissions: { contents: read }
engine: copilot
network: { allowed: [defaults, github] }
tools:
  github: { toolsets: [repos] }
  edit:
safe-outputs:
  create-pull-request: { title-prefix: "[docs] ", labels: [documentation, automated], draft:
true }
  create-issue: { max: 1 }
```

Read them as two mini-products with different rhythms. Triage is **read-mostly and chatty** — it comments and labels, never touches code. Docs is **read-mostly and proposes** — it drafts a PR a human merges. Neither holds a write token; both are bounded by the same firewall and safe-outputs boundary you've applied since Part II opened.

> Verifying both examples

```
gh aw compile examples/ch09/continuous-triage.md
# ✓ examples\ch09\continuous-triage.md (102.1 KB) - 0 error(s), 0 warning(s)
gh aw compile examples/ch09/continuous-docs.md
# ✓ examples\ch09\continuous-docs.md (104.0 KB) - 0 error(s), 0 warning(s)
```

● YOU ALREADY KNEW HOW TO BUILD THESE

There's no new frontmatter here — just triggers (Ch4), an engine (Ch5), safe outputs (Ch6), the firewall (Ch7), and tools (Ch8), composed with intent. That's the whole point of a pattern: the value is in the *combination and the prompt*, not in a new feature.

Recap & what's next

You've shipped your first two Continuous-X mini-products:

- A **Continuous X pattern** turns one recurring judgement task into a standing workflow — a **mini-product** with an owner-task, a cadence, and bounded outputs.
- **Continuous Triage** is reactive + proactive, writing via `add-comment` / `add-labels`; **Continuous Docs** triggers on code merges and *proposes* a draft PR.
- Both are **compositions** of Chapters 4–8 — no new syntax; the value is in the combination and the prompt.
- Design against the failure modes: cap outputs, prefer draft PRs for anything high-stakes, and ship one pattern at a time.

What's next. Triage and docs keep the *inbox* honest. The other half of a healthy repo is the code itself. In [→Chapter 10: Continuous Review, Testing & CI-Doctor](#), we close the quality loop — while keeping humans firmly on the merge decision.

Continuous Review, Testing & CI-Doctor

Close the *quality loop* with *Review*, *Testing*, *CI-Doctor*, and *Refactoring* patterns while keeping humans on the merge decision.

Objective

By the end of this chapter you can close the repository's **quality loop** with four more patterns — **Review**, **Testing**, **CI-Doctor**, and **Refactoring** — while keeping humans firmly on the merge decision. This is the second half of the Continuous-X library and the close of Part II.

Everything targets **gh aw v0.81.6**. The Repo Assistant graduates from tending the issue tracker to helping tend the code.

Concept: closing the quality loop

A repository's quality has a **loop**: code is proposed (a PR), reviewed, tested, merged, and — when something slips — fixed. Traditional CI automates the *deterministic* checks in that loop: does it compile, do the tests pass, does the linter approve. But the *judgement* steps — is this a good change? is this test worth adding? why did CI actually break? — still wait on a human.

Continuous-X patterns fill exactly those judgement gaps. Where [→Chapter 9](#) kept the *inbox* honest, these keep the *codebase* honest — each one a mini-product owning one link in the quality loop.

The one rule that makes it safe: humans keep the merge

The defining constraint of quality automation is that **the agent proposes; a human disposes**. A review agent *comments*, it doesn't approve. A test-improver *opens a draft PR*, it doesn't push to main. This isn't timidity — it's what lets you run these patterns at all. The agent accelerates the work up to the decision point and stops, leaving the irreversible call to a person. That's the human-in-the-loop principle, and it's why the safe-outputs boundary from [→Chapter 6](#) matters most here.

● MEASURED: QUALITY AGENTS AT SCALE

A single PR-review agent originating in `clash-verge-rev` has been cloned across **215+ repositories**; `backend.ai-webui` runs both a daily test-improver and an e2e-healer; `camunda` runs CI cost analysis. Quality is the category where Continuous-X adoption has spread fastest — because the human-keeps-the-merge rule makes it low-risk to try.

In gh-aw: the Review, Testing, CI-Doctor, and Refactoring recipes

Four patterns, each triggered by a different moment in the quality loop — and each writing through a safe output that stops short of merging.

Pattern	Trigger	Writes via
Review	<code>pull_request</code>	<code>submit-pull-request-review</code> (COMMENT only)
Testing	<code>schedule</code>	<code>create-pull-request</code> (draft)
CI-Doctor	<code>workflow_run</code> (CI failed)	<code>add-comment</code> / <code>create-issue</code>
Refactoring	<code>schedule</code> or command	<code>create-pull-request</code> (draft)

Review: comment, never approve

The Review pattern reads a PR diff and leaves inline feedback. The critical setting is `allowed-events: [COMMENT]`, which “prevents the agent from submitting APPROVE reviews regardless of what the agent attempts to output” — the docs explicitly recommend it as “the default for automated review workflows... without creating a persistent merge-blocking state” ([Safe Outputs](#)). Infrastructure enforces the human-keeps-the-merge rule.

CI-Doctor: react to the failure

CI-Doctor is the elegant use of the `workflow_run` trigger from [→Chapter 4](#) with **conclusion filtering**: fire only when a named CI workflow finishes with `failure`, read the logs, and post a diagnosis. Because `workflow_run` is hardened against cross-repo abuse, this stays safe even on public repos.

> The CI-Doctor trigger – wake only on a real CI failure

```
on:  
  workflow_run:  
    workflows: ["CI"]  
    types: [completed]  
    conclusion: [failure]      # only when CI actually broke  
    branches: [main]
```

Testing & Refactoring: propose a diff

Both run on a schedule, do focused work, and open a **draft** `create-pull-request`. Testing adds coverage without touching production code; Refactoring makes a small, behavior-preserving cleanup. Draft PRs keep the human on the merge, exactly as in the Docs pattern.

● BUILDER DETAIL: GIVE THE TESTER A SCOPED SHELL

A test-improver has to *run* the suite, so it needs `bash` — but scope it. Prefer an explicit allowlist like `bash: ["npm ci", "npm test", "npx jest"]` over the unrestricted `bash: [":*"]`, and add only the ecosystem the build needs to `network.allowed` (e.g. `node`). Capability where required, tight everywhere else.

When to automate quality (and where humans stay in the loop)

Quality patterns pay off when they act as a **tireless first pass** — catching the obvious before a human spends attention, never replacing the human's final say. The line to hold: *automate the noticing and the drafting; reserve the deciding.*

Agent may...	Human keeps...
comment on a PR, flag risks	approve / request changes / merge
open a draft test or refactor PR	review and merge that PR
diagnose a CI failure, file an issue	decide the fix and ship it

When not to

- **Don't let a review agent block merges.** Auto `REQUEST_CHANGES` creates a persistent merge-blocking state from a fallible model. Keep `allowed-events: [COMMENT]` unless a human explicitly wants gating.

- **Don't let the test-improver edit production code.** Instruct it to add tests only; a PR that “fixes” code to make a test pass is the opposite of what you want.
- **Don't auto-merge agent PRs.** The draft PR is the human's decision point — automating the merge throws away the one safeguard that makes this safe.
- **Don't run a refactoring agent on a repo without good tests.** “Behavior-preserving” is only verifiable if the tests can prove it. Ship Testing before Refactoring.

Worked example: a PR-review plus daily-test-improver pair

Two complementary quality agents: one reacts to every PR, the other proactively strengthens the tests. Both compile cleanly, and both stop short of the merge.

```
> examples/ch10/continuous-review.md – comment-only PR review (compiles: 0/0)
```

```
on:
  pull_request: { types: [opened, synchronize] }
permissions: { contents: read, pull-requests: read }
engine: copilot
network: { allowed: [defaults, github] }
tools:
  github: { toolsets: [pull_requests] }
safe-outputs:
  create-pull-request-review-comment: { max: 10 }
  submit-pull-request-review:
    allowed-events: [COMMENT]      # can never approve or block
    max: 1
```

```
> examples/ch10/daily-test-improver.md – proposes tests as a draft PR (compiles: 0/0)
```

```
on:
  schedule: daily
  workflow_dispatch:
permissions: { contents: read }
engine: copilot
network: { allowed: [defaults, github, node] }
tools:
  github: { toolsets: [repos] }
  bash: ["npm ci", "npm test", "npx jest", "npx vitest run"] # scoped shell
  edit:
safe-outputs:
  create-pull-request: { title-prefix: "[tests] ", labels: [tests, automated], draft: true }
```

The review agent holds **read-only** PR access and can only emit a `COMMENT` review — the `allowed-events` setting makes “never block a merge” an infrastructural guarantee, not a hope. The test-improver gets a **scoped shell** to run the suite and `edit` to write tests, but its sole output is a **draft** PR a human reviews. Both accelerate the work right up to the human's decision, then hand it over.

> Verifying both examples

```
gh aw compile examples/ch10/continuous-review.md
# ✓ examples\ch10\continuous-review.md (101.8 KB) - 0 error(s), 0 warning(s)
gh aw compile examples/ch10/daily-test-improver.md
# ✓ examples\ch10\daily-test-improver.md (103.8 KB) - 0 error(s), 0 warning(s)
```

● THE PRINCIPLE IS ENFORCED, NOT JUST DOCUMENTED

Notice how `allowed-events: [COMMENT]` and `draft: true` turn “humans keep the merge” from a guideline into a compiled constraint. Even a hijacked agent can't approve a PR or merge one — the safe-outputs layer simply won't let the request through. That's the quality loop closed *safely*.

Recap & what's next

You've closed the quality loop — and Part II:

- Quality automation fills the **judgement gaps** CI can't: is this change good, is this test worth adding, why did CI break.
- Four patterns — **Review** (PR, comment-only), **Testing** (scheduled draft PR), **CI-Doctor** (`workflow_run` on failure), **Refactoring** (scheduled draft PR).
- The unbreakable rule is **humans keep the merge** — enforced by `allowed-events: [COMMENT]` and `draft: true`, not just by convention.
- Automate the noticing and drafting; reserve the deciding. Ship Testing before Refactoring, and never auto-merge an agent's PR.

What's next. You now have a shelf of patterns — and you're about to notice how much they repeat. **Part III** scales from one repo to an org. →[Chapter 11: Reuse & Memory](#) factors the shared parts into imported components and gives the Repo Assistant memory that persists across runs.

Reuse & Memory: Shared Components and Repo Knowledge

Factor common intent into imported shared components and give the Repo Assistant memory that persists across runs.

Objective

By the end of this chapter you can factor common intent into **imported shared components** so a fleet of workflows stops repeating itself, and give the Repo Assistant **memory** that persists across runs. This opens **Part III**: the leap from one repository to an organization.

Everything targets **gh aw v0.81.6**. We take the triage policy we've refined over ten chapters and turn it into a single shared file every repo can import.

Concept: don't repeat yourself across a fleet

One repo, one triage workflow — fine to write inline. But the moment you have five repos each with a triage workflow, you've copied the same allowed-label list, the same tools, the same policy prose five times. Fix a rule in one, and the other four drift. This is the classic **Don't Repeat Yourself** problem, now at the scale of a fleet of agents.

There are two distinct kinds of “sameness” to factor out:

- **Shared configuration and intent** — the toolset, the safe-outputs limits, the triage policy itself. This wants to live in *one file* that many workflows **import**.
- **Shared knowledge over time** — what the agent learned on previous runs (recurring duplicates, project conventions). This wants to **persist** across runs as memory.

Imports solve repetition *across workflows*; memory solves amnesia *across runs*. Together they turn a pile of copy-pasted workflows into a maintained, learning fleet.

● LEADER LENS: GOVERN ONCE, APPLY EVERYWHERE

A shared component is a *governance surface*. Encode your triage policy, your allowed labels, your security posture once; every repo that imports it inherits the current version. Updating org-wide behavior becomes a single reviewed change, not a 200-repo migration — which is exactly how the largest adopters run agents at scale.

In gh-aw: imports, shared components, and repo memory

Imports and shared components

The `imports:` field “compose[s] shared tools, steps, MCP servers, and prompts from other workflow files” ([Frontmatter](#)). A **shared component** is simply a workflow file without an `on:` field: it is “validated but not compiled into GitHub Actions, only imported by other workflows” ([Imports](#)).

```
> Import a shared file; its tools and safe-outputs merge into yours
```

```
imports:
- shared/triage-policy.md           # local, repo-relative
- acme-org/shared-workflows/triage.md@v2.1.0 # cross-repo, pinned to a tag
```

Paths resolve three ways: **relative** to the workflow, **repo-root** (`.github/...`), or **cross-repo** as `owner/repo/path@ref` — and cross-repo imports are “pinned to a semantic tag, branch, or commit SHA” and cached for offline compiles ([Imports](#)). Fields merge sensibly: tool allowlists concatenate and dedupe; each safe-output type is defined once, with the main workflow winning on conflict.

Packaging dependencies: the Agent Package Manager (APM)

Imports compose files you write. But agents increasingly depend on *published* primitives — skills, prompts, instructions, sub-agents, hooks, and plugins — that live in other repositories and evolve on their own cadence. The **Agent Package Manager (APM)** “manages AI agent primitives... Packages can depend on other packages and APM resolves the full dependency tree” ([APM Dependencies](#)). It is the same DRY instinct as a shared component, but for the agent's context rather than its config — a real package manager for agent intelligence.

APM plugs into gh-aw through exactly the mechanism you just learned: you import a shared file, `shared/apm.md`, and pass it the packages you want. That import “adds a dedicated `apm` job” that resolves and packs the packages into a bundle at build time, which the agent job unpacks “for deterministic startup” ([APM Dependencies](#)).

```
> Illustrative: depend on published skills via the shared/apm.md import
```

```
imports:
- uses: shared/apm.md           # vendored from microsoft/apm via `gh aw add`
  with:
    packages:
      - microsoft/apm-sample-package           # a full package
      - github/awesome-copilot/skills/review-and-refactor # one primitive
      - anthropics/skills/skills/frontend-design#v2.0    # pinned to a tag
```

Each entry is a package reference in one of three shapes: `owner/repo` (a full package), `owner/repo/path` (a single primitive such as one skill), or `owner/repo#ref` (pinned to a tag, branch, or SHA) ([APM Dependencies](#)). Reproducibility is the point: an `apm.lock` file “pin[s] every package to an exact commit SHA, so the same versions are installed on every run,” and those lock diffs “appear in pull requests and are reviewable before merge” — an audit trail for exactly which agent context is in use. That reviewable, pinnable supply chain is what makes APM governable at org scale, which we return to in [→Chapter 13](#).

Repo memory vs. cache memory

Persistence comes in two flavors, and choosing correctly matters. **Repo memory** gives “persistent file storage via Git branches with unlimited retention” — enable it with `tools: repo-memory: true` and the compiler auto-configures a `memory/default` branch that files “auto-commit/push after workflow completion” ([Repo Memory](#)). **Cache memory** uses the GitHub Actions cache instead — fast, but 7-day retention and no version control.

	Repo memory	Cache memory
Storage	Git branches	Actions Cache
Retention	Unlimited	7 days
Versioned	Yes	No
Best for	Long-term insights/history	Temporary/session state

● BUILDER DETAIL: SHARED COMPONENTS CAN BE PARAMETERIZED

A shared file can declare an `import-schema` of typed parameters; callers pass values with the `uses / with` form. That lets one `deploy.md` or `trriage.md` serve many repos with per-repo tweaks (region, allowed labels) — reuse without a fork.

When to extract a shared component (and when to inline)

The rule of thumb is the **rule of three**: inline the first time, wince the second, extract the third. Premature sharing couples workflows that should stay independent; late sharing leaves you with drift. Extract when the *same* intent genuinely recurs and you want it governed centrally.

Extract to a shared import when...	Keep inline when...
the policy/toolset repeats across repos	it's genuinely one-off
you want one place to update org-wide	the workflows will diverge anyway
a security config must be consistent	early days — you're still iterating

When not to

- **Don't track a moving branch for cross-repo imports.** Pin to a tag or SHA (`@v2.1.0`), not `@main` — an unpinned import is a supply-chain risk, the same lesson as unpinned engine versions in [Chapter 5](#).
- **Don't consume APM packages unpinned or unreviewed.** A skill is executable context: commit the `apm.lock`, review its diffs, and pin references. Skills you don't control are exactly where the [Chapter 7](#) threat model applies — govern them (Chapter 13), don't trust them by default.
- **Don't put secrets in memory.** Repo memory follows repository permissions and lives in a branch; “don't store sensitive data in repo memory.” Keep secrets in Actions secrets, always.
- **Don't reach for repo memory when cache memory fits.** Session-only scratch state doesn't need an unlimited, version-controlled branch — use the faster 7-day cache.
- **Don't over-abstract.** A shared component with fifteen parameters is harder to reason about than two honest copies. Share the stable core; let the edges vary.

Worked example: importing a shared triage policy with memory

Let's collapse ten chapters of triage refinement into **one shared file** plus a thin workflow that imports it and remembers what it learns.

```
> examples/ch11/shared/triage-policy.md – a shared component (no on:
, so it never runs alone)
```

```
description: Shared triage policy reused across Repo Assistant workflows
tools:
  github: { toolsets: [issues] }
safe-outputs:
  add-comment: { max: 1 }
  add-labels:
    allowed: [bug, enhancement, question, documentation, duplicate, needs-info]
    max: 3
---
## Shared triage policy
Categorize the issue, summarize it in one sentence, note missing info, apply at
most three labels from the allowed set, and post one concise comment.
```

```
> examples/ch11/repo-assistant-shared.md – imports the policy, adds memory (compiles: 0/0)
```

```
on:
  issues: { types: [opened, reopened] }
  schedule: daily
  workflow_dispatch:
permissions: { contents: read, issues: read }
engine: copilot
network: { allowed: [defaults, github] }
imports:
  - shared/triage-policy.md      # tools + labels + safe-outputs, from one file
tools:
  repo-memory: true             # persist what it learns across runs
```

The main workflow is now almost content-free: the *policy* — tools, allowed labels, safe outputs — comes entirely from the imported file, and `repo-memory: true` lets the agent read prior notes and append new ones each run. Point ten repositories at the same `shared/triage-policy.md` (or a pinned cross-repo import) and they triage identically; change the file once and all ten update on their next compile.

```
> Verifying the example
```

```
gh aw compile examples/ch11/repo-assistant-shared.md
# ✓ examples\ch11\repo-assistant-shared.md (108.1 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

● THE COMPILER MERGED TWO FILES INTO ONE LOCK

Notice the workflow declares no `safe-outputs` or `github` tool of its own — they came from the import and were merged in at compile time (tool allowlists concatenate; safe-output types are defined once). The compiled `.lock.yml` is a single self-contained artifact, exactly as in [→Chapter 3](#) — imports are resolved at build time, not runtime.

Recap & what's next

You can now stop repeating yourself across a fleet, and let agents remember:

- As you scale to many repos, factor common intent into **shared components** — files without `on:` that others `import`.
- Imports resolve **relative**, **repo-root**, or **cross-repo** (`owner/repo/path@ref`, pinned); fields merge, and `import-schema` allows typed parameters.
- The **Agent Package Manager (APM)** rides the same import mechanism (`shared/apm.md` + `packages:`) to depend on published skills/prompts/plugins, pinned to exact SHAs in an `apm.lock` for reproducible, reviewable builds.
- **Repo memory** (Git branches, unlimited, versioned) vs. **cache memory** (Actions cache, 7-day, fast) give the agent persistence across runs.
- Follow the **rule of three**, **pin** cross-repo imports, and never store secrets in memory.

What's next. A shared, remembering fleet is powerful — and now you need to see what it's doing. In [→Chapter 12: Trust & Operate](#), we inspect, debug, and audit runs with `gh aw logs` and `gh aw audit`, because you can't govern what you can't see.

Trust & Operate: Observability and Debugging

Inspect, debug, and audit runs with `gh aw logs`, `gh aw audit`, and OpenTelemetry so you can trust what the fleet does.

Objective

By the end of this chapter you can **inspect, debug, and audit** what your workflows do — with `gh aw logs`, `gh aw audit`, run summaries, and OpenTelemetry — so you can trust a fleet you can't watch by hand.

Everything targets **gh aw v0.81.6**. We take a run that went wrong and trace it from an overview table down to the exact failing step.

Concept: you can't govern what you can't see

Everything in Part III assumes a fleet running unattended — agents triaging, reviewing, and opening PRs across many repos while you sleep. That only works if you can answer, after the fact: *what did it do, why, what did it cost, and what did it touch?* **Observability is the precondition for trust.** You can't govern — can't budget, can't secure, can't improve — what you can't see.

An agentic run is unusually inspectable because, as [→Chapter 3](#) showed, only one job is non-deterministic and everything is captured as artifacts. `gh-aw` “provides comprehensive observability through GitHub Actions runs and artifacts... [which] preserve prompts, outputs, patches, and logs for post-hoc analysis” ([Security Architecture](#)). Debugging an agent isn't guesswork; it's reading a well-kept record.

● LEADER LENS: AUDITABILITY IS A COMPLIANCE ASSET

Every run leaves a durable trail — the prompt it saw, the actions it proposed, the network it touched, the tokens it spent. That record is what lets you answer a security or cost question about *any* past run, and it's what turns “we run autonomous agents” from a worry into a governable practice.

In gh-aw: logs, audit, OpenTelemetry, run summaries

Three CLI commands and one export cover the whole observability surface.

gh aw logs — the overview and the artifacts

It “download[s] and analyze[s] agentic workflow logs and artifacts... and provides an overview table with aggregate metrics including duration, token usage, and cost information” (`gh aw logs --help`). By default it grabs just the compact usage artifact; widen with `--artifacts`:

> Fetch runs and choose how much to download

```
gh aw logs                # overview table: duration, tokens, cost
gh aw logs repo-assistant # just one workflow's runs
gh aw logs --artifacts all # everything: prompt, output, patch, logs
gh aw logs --artifacts agent,firewall # only what you need
```

The downloadable artifacts are the agent's black box recorder: `agent-stdio.log`, `safe_output.jsonl` (what it proposed), `aw-{branch}.patch` (what it changed), `workflow-logs/`, and `summary.json`. Available sets include `activation`, `agent`, `detection`, `firewall`, `github-api`, `mcp`, `usage`.

gh aw audit — the focused report

Where `logs` is broad, `audit` is deep. It audits runs “by downloading artifacts and logs, detecting errors, analyzing MCP tool usage, and generating a concise report” (`gh aw audit --help`). Point it at a run and it finds the problem for you:

> Investigate one run, or diff two

```
gh aw audit 1234567890      # detailed Markdown report for one run
gh aw audit <run-url>/job/<id> # a job URL – extracts the first failing step
gh aw audit 1234567890 1234567891 # compare two runs (first = baseline)
```

Given a job URL without a step, it “finds and extracts the first failing step's output” — it navigates to the failure for you. Its Firewall Analysis section (from [→Chapter 7](#)) lists every domain the agent tried to reach with allow/deny status.

Run summaries and OpenTelemetry

Every run also writes a rich Markdown **step summary** in the Actions UI, and `gh aw status` reports fleet health at a glance. For centralized, cross-run visibility, the `observability.otlp` block “export[s] distributed traces from workflow runs to an OpenTelemetry Protocol (OTLP)

compatible backend” ([Frontmatter](#)) — so agent runs appear in the same tracing tool as the rest of your systems.

● **BUILDER DETAIL: START NARROW, THEN WIDEN**

Downloading every artifact for every run is slow. Start with `gh aw logs` for the overview, spot the anomalous run (long duration, high tokens, a failure), then `gh aw audit <run-id>` just that one. Only reach for `--artifacts all` when the audit report points you at something you need to read in full.

When to inspect versus trust (human-in-the-loop)

You can't read every run of a busy fleet — nor should you. The skill is knowing which runs earn a look. Let the cheap signals (the overview table, the safe-outputs boundary, the threat-detection gate) carry the routine cases, and spend attention where the signal says something's off.

Inspect closely when...	Trust the guardrails when...
a run failed or timed out	it succeeded and produced expected safe outputs
tokens/cost spiked vs. the norm	cost is in the usual band
the firewall logged unexpected domains	egress stayed within the allowlist
you're rolling out a new or changed workflow	a stable workflow is running unchanged

When not to

- **Don't skip observability because “it's working.”** A silent fleet is not a healthy fleet — it's an unmonitored one. Glance at `gh aw logs` regularly even when nothing's on fire.
- **Don't debug from the model's chat alone.** The artifacts — patch, safe-output JSON, firewall log — are ground truth; the agent's narration is not. Read the record, not the story.
- **Don't treat observability as a substitute for the guardrails.** Seeing a bad action after the fact is no help if it already shipped. Logs and audit complement safe outputs and review gates; they don't replace them.

Worked example: debugging a failed run from its logs

The Repo Assistant's nightly run failed. Here's the trace from “something's wrong” to root cause — three commands, no guessing.

1. Get the overview. Start broad to find the bad run and its ID:

> The overview table surfaces the anomaly

```
gh aw logs repo-assistant
# RUN ID      WORKFLOW      STATUS  DURATION  TOKENS  COST
# 1234567890  repo-assistant failure  4m12s    182,400 ...
# 1234567889  repo-assistant success   0m48s    12,100  ...
```

The failed run also burned **15×** the tokens of a healthy one — two signals pointing at the same run.

2. Audit that run. Let `audit` find the failing step and explain it:

> A focused report that detects the error for you

```
gh aw audit 1234567890
# Downloads artifacts + logs, detects errors, analyzes MCP tool usage,
# and writes a concise Markdown report – including the first failing step
# and a Firewall Analysis of every domain the agent tried to reach.
```

Say the report shows the agent looping on a tool call to a domain the firewall **denied** — that explains both the failure *and* the token blow-up (it retried until timeout).

3. Confirm and fix. Pull the full artifacts if you need to read the raw exchange, then fix the cause — add the domain to `network.allowed` ([→Chapter 7](#)) — and recompile:

> Read the black box, then fix the workflow

```
gh aw logs repo-assistant --artifacts all # agent-stdio.log, firewall log, patch...
# → root cause: egress to an un-allowed domain, retried to timeout
# fix: add the domain to network.allowed, then:
gh aw compile .github/workflows/repo-assistant.md
```

● MAKING RUNS OBSERVABLE UP FRONT

The chapter's example, `examples/ch12/repo-assistant-observable.md`, adds an `observability.otlp` block so its traces flow to your OpenTelemetry backend automatically — it compiles clean (0/0, secrets approved as in [→Chapter 5](#)). Between OTel traces, the Actions step summary, and `gh aw logs / audit`, you rarely have to guess what a run did.

Recap & what's next

You can now see what your fleet does, and debug it when it misbehaves:

- **Observability is the precondition for trust** — you can't govern what you can't see, and every run leaves a durable artifact trail.
- `gh aw logs` gives the **overview + artifacts** (duration, tokens, cost; `--artifacts` to download more); `gh aw audit` gives a **focused report** that detects the failing step and analyzes tool/firewall use.
- Run **step summaries**, `gh aw status`, and **OpenTelemetry** (`observability.otlp`) round out the picture.
- Inspect the runs that signal trouble (failures, cost spikes, denied egress); trust the guardrails for the rest — but never let observability replace them.

What's next. Seeing cost is the first step; controlling it is the next. In [→Chapter 13: Governance & FinOps](#), we cap and meter agentic spend with `max-ai-credits` and set the org policy that keeps a fleet affordable and compliant.

Governance & FinOps: Policy and Cost at Scale

Cap, meter, and gate agentic spend with AI Credits and max-ai-credits, and set org policy so the fleet stays affordable and compliant.

Objective

By the end of this chapter you can **cap, meter, and gate** agentic spend — with `max-ai-credits` per-workflow budgets and org-wide defaults — and set the policy that keeps a fleet affordable and compliant as it grows. That policy surface extends past cost to the agent's *supply chain*: which skills and prompts your workflows are even allowed to consume.

Everything targets **gh aw v0.81.6**. We put a hard budget on the Repo Assistant and show how an org enforces the same limits — and the same approved dependency list — across every repo at once.

Concept: agentic work has a budget and a policy surface

CI/CD costs are largely fixed: a build takes roughly the same compute every time. Agentic work is different — each run spends a *variable* amount of model inference depending on how much the agent reads, reasons, and retries. That variability is the whole reason agents are powerful, and it's also why **agentic work has a budget in a way CI never did**. Left unbounded, a looping agent or an over-eager schedule can quietly run up real money.

At one repo, this is a cost knob. Across an org, it becomes a **policy surface**: which workflows may run, which capabilities they may use, what they may spend, which model they default to. FinOps — the discipline of managing variable cloud spend — now applies to your agents, and governance means answering these questions *once, centrally*, not per repo.

● LEADER LENS: PREDICTABLE SPEND, ENFORCED CENTRALLY

The two questions a leader asks about an agent fleet are “what will it cost?” and “what is it allowed to do?” gh-aw answers both with hard controls: per-workflow credit budgets that fail safe, and org/enterprise defaults and policies that apply to every repo without editing a single workflow. Spend becomes a dial you set, not a surprise you discover.

In gh-aw: max-ai-credits, AI Credits, token efficiency, policy

Cost is denominated in **AI Credits (AIC)**, a model-normalized unit so budgets mean the same thing regardless of engine. You control it at two levels.

Per-workflow budgets

`max-ai-credits` “sets the AWF AI Credits budget used for cost enforcement. It is enabled by default and defaults to `1000` (`1k`) when omitted” — with steering messages at 80%, 90%, 95%, and 99% of budget ([Frontmatter](#)). Its sibling `max-daily-ai-credits` caps a rolling 24-hour total across recent runs of the same workflow; when exceeded it “warns, creates an issue, skips the agent job” ([Frontmatter](#)) — a fail-safe, not a silent overspend.

> Three cost dials in the frontmatter

```
max-ai-credits: 200          # per-run budget (default 1000); K/M suffixes ok
max-daily-ai-credits: 2000  # rolling 24h cap across this workflow's runs
timeout-minutes: 10         # wall-clock ceiling (default 20)
on:
  stop-after: "+30d"        # stop triggering after a deadline (Ch. 4)
```

Token efficiency is the other half: a tighter prompt, a narrower toolset, and read-only scopes all reduce credits per run. The cheapest run is the one that reads only what it needs — good security and good FinOps are the same discipline.

Org-wide defaults and policy

Editing every workflow doesn't scale. `gh aw env` manages `GH_AW_DEFAULT_*` variables at repository, organization, or enterprise scope from a YAML file ([Governance](#)):

> defaults.yml – org-wide guardrails, applied without touching workflows

```
default_max_ai_credits: "5M"
default_max_daily_ai_credits: "15M"
default_max_turns: "12"
default_timeout_minutes: "30"
default_model_copilot: "gpt-5-mini"
```

Values **percolate** with a clear precedence: “workflow frontmatter value... repository variable... organization variable... enterprise variable... built-in compiler fallback” ([Governance](#)). Beyond numbers, **policy variables** (`GH_AW_POLICY_*`) “enforce capability gates... without recompiling any workflow” — for instance `GH_AW_POLICY_ALLOW_CREATE_PULL_REQUEST=false` makes the safe-outputs server refuse to start for any workflow that tries to open PRs, org-wide.

● BUILDER DETAIL: THE RECOMMENDED ROLLOUT

The docs prescribe a layered rollout: *enterprise baseline* → *org where needed* → *repo exceptions* → *rare, explicit frontmatter overrides*. Preview any change with `gh aw env update ... --dry-run` before applying. Most repos stay aligned to the baseline; exceptions are deliberate and visible.

Governing the agent supply chain (APM)

Budgets and policy variables govern *spend* and *capabilities*. A third surface is the *dependency* supply chain from →[Chapter 11](#): the skills, prompts, and plugins an agent consumes are executable context, so an unreviewed one is an injection vector — exactly the →[Chapter 7](#) threat model. The **Agent Package Manager (APM)** “treats agent skills as packages with the same governance primitives that enterprises require for code dependencies” ([Governing agentic workflows](#)).

Three controls turn that supply chain into a governed surface:

- **Pinning & scanning.** Every package is pinned to an exact commit SHA in an `apm.lock.yaml` so “there is no drift between what was reviewed and what actually runs,” and install-time scanning flags “hidden Unicode threats like homoglyphs, bidirectional override characters, and zero-width joiners” that could smuggle invisible instructions into a prompt ([Governing agentic workflows](#)).
- **Org-level allowlists.** An `apm-policy.yaml` in the org's `.github` repository controls which packages any repo may consume. Inheritance is *tighten-only* across *enterprise* → *org* → *repo* — children “can narrow allowlists, add deny entries, and escalate enforcement, but... cannot relax constraints set by a parent” — the same percolation model as your cost defaults, applied to dependencies.
- **Isolation.** Importing a skill with `isolated: true` means “the agent sees only the skill's packaged instructions,” so a compromised repo-level `AGENTS.md` or `copilot-instructions.md` cannot silently override a security skill's rules ([Governing agentic workflows](#)).

> Illustrative: an org-wide dependency allowlist in `.github/apm-policy.yml`

```
name: "Org agent governance"
enforcement: block
dependencies:
  allow:
    - "org/approved-security-skills/*"
    - "org/approved-review-skills/*"
  deny:
    - "*" # deny everything not explicitly allowed
require_pinned_constraint: true
```

Together — lockfile + org policy + isolation — these give a platform team a complete answer to the compliance question “show me exactly what instructions the agent followed, at what version.” Air-gapped shops can go further and route all downloads through a corporate scanning proxy (`PROXY_REGISTRY_ONLY=1`) so nothing is fetched directly from GitHub ([Governing agentic workflows](#)).

When to cap, meter, and gate (cost and risk trade-offs)

Budgets trade off cost certainty against task completion: too tight and useful runs get cut off; too loose and a bad run overspends. Tune to the shape of the work.

For...	Set...
a cheap, frequent job (triage)	a low per-run <code>max-ai-credits</code> and a daily cap
an occasional deep task (refactor)	a higher per-run budget, no aggressive daily cap
a whole org	generous enterprise defaults, tightened per-org/repo
a capability you want to forbid	a <code>GH_AW_POLICY_*</code> gate, not per-workflow edits
the skills/prompts agents may use	an <code>apm-policy.yml</code> allowlist + committed, pinned <code>apm.lock</code>

When not to

- **Don't disable budgets (`max-ai-credits: -1`) to “unblock” a workflow.** A run hitting its cap is usually a looping or over-scoped agent — fix the cause; the budget did its job.
- **Don't set org defaults so tight that every repo overrides them.** If exceptions become the norm, the baseline is wrong. The goal is most repos aligned, few exceptions.
- **Don't rely on budgets for security.** A budget limits *spend*, not *blast radius* — that's still safe outputs, the firewall, and policy gates (Chapters 6–8). Cost controls and security controls are complementary.

- **Don't govern by editing workflows.** At scale, prefer `gh aw env` defaults and `GH_AW_POLICY_*` gates — central, reviewable, and applied without recompiling every repo.
- **Don't let agents pull skills you haven't approved.** An unpinned, unscanned skill package is an ungoverned instruction source; set an `apm-policy.yml` allowlist and require pinned constraints before you widen a fleet, not after.

Worked example: a budgeted, policy-compliant Repo Assistant

Here's the Repo Assistant with a real budget — capped three ways and set to expire — the version an org would be happy to run at scale.

```
> examples/ch13/repo-assistant-budgeted.md – cost controls in frontmatter (compiles: 0/0)
```

```
on:
  issues: { types: [opened] }
  schedule: daily
  workflow_dispatch:
    stop-after: "+30d"          # stop triggering after a month
permissions: { contents: read, issues: read }
engine: copilot
network: { allowed: [defaults, github] }
max-ai-credits: 200           # tight per-run budget (default is 1000)
max-daily-ai-credits: 2000   # rolling 24h cap across runs
timeout-minutes: 10         # wall-clock ceiling
safe-outputs:
  add-comment: { max: 1 }
  add-labels:
    allowed: [bug, enhancement, question, documentation]
    max: 1
```

Four independent cost brakes: a **per-run** credit budget of 200, a **rolling daily** cap of 2000, a **wall-clock** ceiling, and a **calendar** expiry. If a single run misbehaves, `max-ai-credits` stops it; if the whole day runs hot, `max-daily-ai-credits` warns, files an issue, and skips further agent runs. None of these require a human watching a dashboard.

Now make it *org-compliant* without editing this file at all. An admin sets baseline defaults and a capability policy once:

```
> Org-wide governance – applied to every repo, no workflow edits
```

```
# defaults.yml, applied at org scope
gh aw env update defaults.yml --scope org --org my-org --dry-run
gh aw env update defaults.yml --scope org --org my-org

# forbid a capability fleet-wide, no recompile needed
gh variable set GH_AW_POLICY_ALLOW_CREATE_PULL_REQUEST --org my-org --body "false"
```

> Verifying the workflow

```
gh aw compile examples/ch13/repo-assistant-budgeted.md
# ✓ examples\ch13\repo-assistant-budgeted.md (101.7 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

● FRONTMATTER WINS, DEFAULTS FILL THE GAPS

Because precedence runs *frontmatter* → *repo* → *org* → *enterprise* → *fallback*, this workflow's explicit `max-ai-credits: 200` stands even under a looser org default — while any repo that omits a budget inherits the org's. Deliberate local choices are respected; silence inherits the safe baseline.

Recap & what's next

You can now keep an agent fleet affordable and compliant:

- Agentic work has a **variable cost** and, at scale, a **policy surface** — FinOps and governance now apply to your agents.
- Per-workflow: `max-ai-credits` (default 1000, steering at 80/90/95/99%), `max-daily-ai-credits` (fail-safe daily cap), `timeout-minutes`, and `stop-after`.
- Org-wide: `gh aw env` sets `GH_AW_DEFAULT_*` defaults that **percolate** (frontmatter → repo → org → enterprise), and `GH_AW_POLICY_*` variables **gate capabilities** without recompiling.
- Supply chain: **APM** governs which skills a fleet may use — SHA-pinned `apm.lock.yaml`, install-time Unicode scanning, a tighten-only `apm-policy.yml` allowlist, and `isolated: true` imports that repo config can't override.
- Tune budgets to the work, roll out defaults in layers, and remember budgets limit *spend*, not *blast radius*.

What's next. You can now build, secure, operate, and govern agentic workflows. The final chapter zooms all the way out: →[Chapter 14: Fleets & Adoption](#) takes the Repo Assistant from one repo to a multi-repo fleet and lays out the enterprise adoption playbook.

Fleets & Adoption: From One Repo to the Org

Scale the Repo Assistant into a governed multi-repo fleet and follow an enterprise adoption playbook to roll it out.

Objective

By the end of this chapter you can take the Repo Assistant from **one repo to a governed, multi-repo fleet** — installing shared workflows across many repositories, coordinating cross-repo work, rolling out safely, and following an enterprise adoption playbook. This is the top of the maturity arc the book has climbed since page one.

Everything targets **gh aw v0.81.6**. The assistant that triaged one issue in Chapter 2 becomes an org-wide capability, run from a single source of truth.

Concept: from one assistant to a governed fleet

Something “qualitatively different becomes possible” when agentic workflows move beyond a single repository: they can “coordinate or scale across dozens of repositories simultaneously” — rolling out changes org-wide, assessing code quality across hundreds of repos, or “aggregat[ing] issue tracking into a single control plane” ([Using at Scale](#)). A fleet isn't just many copies of one workflow; it's a *managed practice*.

This is the arc the whole book has followed: **the Individual** (one workflow), **the Team** (safe, reviewed, patterned), and now **the Organization** (a fleet at scale). Each level reused everything below it — the fleet is just Parts I and II, governed and multiplied.

● MEASURED: WHAT A FLEET ACTUALLY DELIVERS

GitHub Next's own [repo-assist-impact](#) report measured a Repo Assistant fleet across **15 repositories**: a net reduction of **651 issues** and a **median 9× velocity** improvement — with the thesis that throughput is now “gated by human decision-making,” not by the agents. Public-preview adopters echo it: Home Assistant, CNCF, Carvana, Marks & Spencer, and Hud.io. The fleet is where the compounding value the book promised on page one finally shows up.

In gh-aw: multi-repo fleets, the dispatcher pattern, rollout

Distribution has “two complementary layers” ([Using at Scale](#)).

1. Install and update (developer-facing)

Use `gh aw add` (or `gh aw add-wizard`) to “install a workflow from another repository, and `gh aw update` to pull in upstream changes while preserving local edits.” A workflow installed this way records where it came from in its `source:` field. The recommended org structure is a central `agentic-workflows` repository as “the source of truth,” with workflows “versioned with exact tags (`@v1.2.0`)” or SHA pins.

2. Coordinate across repos (the dispatcher pattern)

For work that spans repositories, two patterns matter. **CentralRepoOps** is “where a central control repository dispatches work to target repositories or aggregates issues from component repositories”; **OrchestratorOps** handles “dispatching parallel worker workflows for large-scale multi-repo operations” ([Using at Scale](#)). The cross-repo writes flow through the same safe-outputs boundary, now with `target-repo` and `allowed-repos`; GitHub Apps are “preferred for automatic token rotation and fine-grained scoping.”

3. Roll out safely

Don't flip a fleet to production writes on day one. **Safe Rollout** “describes how to move from report-only or staged behavior to production writes with evidence and control” — using the `staged: true` preview mode from [→Chapter 6](#) as a shadow-evaluation step before promotion ([Using at Scale](#)).

● BUILDER DETAIL: FIND EVERYTHING THE FLEET DID

Set a `tracker-id:` on a workflow and it “tags every asset (issues, PRs, discussions, comments) the workflow creates with a hidden marker,” so one GitHub search surfaces all of a workflow's output across every repo in the fleet. Pair it with `private: true` on internal workflows to control what's installable elsewhere.

When to scale org-wide (a readiness checklist)

Scaling multiplies both value and mistakes. Before cloning a workflow across a fleet, it should have earned it on one repo first. A readiness checklist:

Before you fan out, confirm...	Because...
the workflow has run cleanly on one repo for a while	a bug cloned to 200 repos is 200 bugs
it's a shared import , pinned to a tag	you can fix it in one place (Ch. 11)
org defaults & policy are set (<code>gh aw env</code> , <code>GH_AW_POLICY_*</code>)	budgets and capability gates apply fleet-wide (Ch. 13)
you can observe it (logs, audit, OTel)	you can't govern what you can't see (Ch. 12)
writes start staged / draft	safe rollout beats a big-bang cutover

When not to

- **Don't fan out a workflow you haven't operated.** Prove it on one repo; earn the fleet.
- **Don't copy-paste across repos.** That's the drift trap — distribute a pinned shared import so one change updates everyone.
- **Don't go to production writes without a rollout.** Start report-only or `staged`, gather evidence, then promote.
- **Don't scale without central governance.** A fleet without org defaults and policy gates is an unbounded cost-and-risk surface; set them *before* you widen, not after.

Worked example: cloning the Repo Assistant across a fleet

Here is the Repo Assistant as a **fleet citizen**: it imports the org's central triage policy, records where it was installed from, and tags its output for fleet-wide search. Every repo runs this same thin file.

```
>
examples/ch14/fleet-triage.md – installed from a central repo, governed as a fleet (compiles: 0/0)
```

```
on:
  issues: { types: [opened, reopened] }
  workflow_dispatch:
permissions: { contents: read, issues: read }
engine: copilot
network: { allowed: [defaults, github] }
source: "my-org/agentive-workflows/workflows/triage.md@v1.2.0" # where it came from, pinned
tracker-id: repo-assistant-triage # find all its assets fleet-wide
imports:
  - shared/triage-policy.md # the org's single source of truth
tools:
  repo-memory: true
```

This one file embodies the whole of Part III. The **policy** comes from a shared import (Ch. 11), so a **fix** propagates everywhere. `source:` pins it to a released version of the central repo, so `gh aw update` pulls upgrades deliberately. `tracker-id:` makes the fleet *auditable* — one search finds every comment and label it produced across every repo. And it inherits the budgets and policy gates an admin set org-wide (Ch. 13).

> The fleet lifecycle, in commands

```
# install the central workflow into a repo (records source:)
gh aw add my-org/agentiic-workflows/workflows/triage.md@v1.2.0

# later, pull the org's upgrade while preserving local edits
gh aw update

# find everything the fleet's triage assistant has done
#   GitHub search: "gh-aw-tracker-id: repo-assistant-triage" in:body
```

> Verifying the example

```
gh aw compile examples/ch14/fleet-triage.md
# ✓ examples\ch14\fleet-triage.md (109.9 KB)
# ✓ Compiled 1 workflow(s): 0 error(s), 0 warning(s)
```

● EVERYTHING YOU LEARNED, IN ONE COMPILED ARTIFACT

Read the frontmatter top to bottom and you'll see all fourteen chapters: a trigger (Ch. 4), an engine (Ch. 5), safe outputs via the imported policy (Ch. 6), a network firewall and read-only scopes (Ch. 7), tools (Ch. 8), a Continuous-X pattern (Ch. 9–10), a shared import and memory (Ch. 11), tracker-based observability (Ch. 12), inherited budgets (Ch. 13), and fleet distribution (Ch. 14) — all compiled into one hardened `.lock.yml`.

Recap & the road ahead

You've reached the top of the arc — from one workflow to a governed fleet:

- Beyond one repo, workflows **coordinate and scale across dozens** — org-wide rollouts, cross-repo quality, a single issue-tracking control plane.
- Distribute via `gh aw add / update` from a **central source-of-truth repo** with **pinned versions**; coordinate with **CentralRepoOps/OrchestratorOps** and cross-repo safe outputs.
- **Roll out safely** (report-only → staged → production), make the fleet **auditable** with `tracker-id`, and govern it centrally before you widen.

- A fleet is just Parts I–II, **governed and multiplied** — and it's where measured impact (15 repos, 651 issues, 9× velocity) shows up.

The road ahead. You set out to look at your own repository, spot three tasks a tireless teammate could own overnight, and ship a governed agentic workflow that does them — safely, cheaply, and reviewably. You now can. Start with one 10-minute win, earn trust, and let the fleet compound. That's Continuous AI: the outer loop, automated — with people firmly in the loop. Go build your Repo Assistant.

By **Maxim Salnikov** · Microsoft · [LinkedIn](#) · [Book repository](#)

GitHub Agentic Workflows: An Interactive Book · Content edition v1.1 · <https://aw.isainative.dev>