

General agents: a co-worker for knowledge workers

What Cowork and OpenWork actually do, when to reach for each, and where the failure modes hide — for students, knowledge workers, and engineers alike.

THE DISCIPLINE

Delegation, not coding. You don't need to know what an API is.

A co-worker you assign — not a chatbot you query

A QUERY · FINE IN CHAT

”Summarize this PDF.”

One question, one answer. The worst case is a wrong answer — annoying, but contained.

AN ASSIGNMENT · WHAT THESE TOOLS ARE FOR

”Read these three vendor contracts, flag every clause that deviates from our standard, and produce a comparison memo color-coded by risk.”

The shift changes the failure mode too. Here the worst case isn't a wrong answer — it's a **confidently executed wrong action** that already touched dozens of your files. The whole course is about delegating at the right level, with the right oversight.

One discipline that outlives any single tool

● Cowork

By Anthropic

The more polished surface and the lowest-friction path to a first working task. A paid-tier feature inside Claude Desktop. **Start here if you're new.**

Prompts + file contents → Anthropic

● OpenWork

Open-source • MIT • OpenCode

Free, younger, ships frequently. Pays off when you want choice over **which model provider** sees your prompts, where work runs, and which plugins compose.

You pick the provider that sees them

The task loop and the core primitives — skills, connectors, sub-agents — are the **same in both**. If a technique only works in one, it's a trick. If it works in both, it's how desktop agents actually behave. Your files stay on your machine either way.

The five disciplines that carry most of the value

80%
OF REAL-WORLD VALUE

- 1 **Delegate, don't query.** A co-worker you assign — think in outcomes with constraints, audience, and the *why*, not in prompts.
- 2 **Three trust levers: folders, connectors, approvals.** A dedicated working folder, the narrowest connector scope that works, and "ask before acting" until you've calibrated.
- 3 **The plan is the leverage, not the output.** Read the plan before it touches a single file. Neither tool guarantees one unless you ask.
- 4 **Climb the autonomy ladder deliberately.** Watch closely → ambient → walk away → high-autonomy → scheduled. One rung per task type, earned with track record.
- 5 **Untrusted content gets the cautious mode, always.** Outside text — an email, a resume, a vendor PDF — can carry instructions the agent reads as commands.

15 concepts across six parts

PART 1 · FOUNDATIONS

1 · What they are 2 · Architecture in three pieces 3 · Folders, connectors, approvals

PART 2 · CONTEXT & PROJECTS

4 · The plan is the leverage 5 · The cost of context 6 · Persistent workspaces

PART 3 · RULES & INSTRUCTIONS

7 · Three layers of instruction 8 · Questions before execution

PART 4 · EXTENDING THE TOOL

9 · Skills 10 · Connectors 11 · Plugins 12 · Sub-agents

PART 5 · SAFETY & AUTONOMY

13 · The autonomy ladder 14 · Prompt injection 15 · Scheduled tasks

PART 6 · IN PRACTICE

A complete worked example — the weekly brief, run twice.

PART ONE

0

Foundations

What these tools are, how they're built, and the trust model underneath — the three concepts that apply identically in both.

1

The actual skill is delegation

You don't need to code or know what an API is. You need to assign work the way you'd brief a smart, motivated colleague who doesn't know your context. Same shape, every domain:

LAWYER

"Flag every clause in these three MSAs that deviates from our redline standard; produce a comparison memo, color-coded by risk."

ACCOUNTANT

"Reconcile this month's bank statement against the GL export, list every unmatched item, and draft a journal-entry recommendation for each."

MARKETER

"Compare conversion across these 12 weekly reports, find the top 3 weeks and what they had in common, one-page summary."

Brief poorly and the agent executes against the gap — and here the gap isn't a wrong answer, it's a wrong **action** already taken.

Three pieces, in both tools

01 · THE DESKTOP APP

Where the agent lives

Runs locally on your machine.

Close the app or let the laptop sleep and execution pauses where it left off.

02 · THE TASK LOOP

The core mechanic

Describe an outcome → agent plans → you approve or redirect → it executes → it pauses before significant actions → you get a deliverable.

03 · THE EXECUTION SURFACE

Where work happens

Local files you've granted, an isolated code sandbox, and external services through connectors.

Privacy is the big divergence. In [Cowork](#), your prompts and the file content the agent reads go to Anthropic. In [OpenWork](#), you pick the provider that sees them. Your files stay on your machine in both.

Three levers set the whole safety story

LEVER 1

Folders

What the agent can read and write.

Make a dedicated working folder

— grant access to that, not your whole Documents directory.

Blast radius = the folder, not your life

LEVER 2

Connectors

Which external services it reaches.

Each one is a separate decision at the narrowest scope that works.

Read scope ≠ send scope.

Native mail connectors draft — they don't send

LEVER 3

Approvals

When it pauses for your OK. Stay in **"ask before acting"** until you've calibrated on a task type. Earn high-autonomy gradually.

Even hands-off, deletions still ask

THE ASYMMETRY

Reads happen automatically. Writes, modifications, deletions, and moves all wait for an explicit click. That's the whole approval model in miniature.

PART TWO

02

Context, sessions & projects

Same primitives as chat, the same pitfalls — but a token bill attached and meaningfully higher stakes.

The plan is the leverage — not the output

You intercept work **between intent and execution**. Every meaningful action passes a moment where you can read what's about to happen and redirect. It's the cheapest place in the whole workflow to course-correct.

1 · Brief the outcome

2 · Agent plans

3 · You read the plan

4 · Approve / redirect

5 · Execute

6 · Review output

At step 3, look at four things: **scope** (did "sort this folder" creep into "rename three subfolders"?), **order**, **tools** (an unexpected connector?), and **assumptions**.

If it's wrong, you don't start over — one sentence rewrites the plan. **Cowork** intercepts per-step inline; **OpenWork** shows the plan top-right with a dedicated Plan agent.

```
Skip step 3, and for step 4 use  
the column headers in the existing  
template instead of new ones.
```

Two minutes of plan review prevents two hours
of cleanup.

Triage first, read selectively, synthesize once

Every message carries the system prompt, your instructions, the conversation, and every file read this session — and the bill is yours. **Don't dump whole folders into context unprompted.**

```
First, list this folder and tell me  
which files matter for my question.  
Read only those, then summarize.
```

A LITIGATION FOLDER · 340 DOCUMENTS

100%

"read everything" — millions of tokens, a generic summary, weak recall

~5%

triage → ~12 files read → a better one-page memo

Route hard thinking to the strong model; route plumbing — file listing, format conversion, OCR — to the cheap one.

A folder with a context file — not a fresh chat each time

If you re-explain the same context every Tuesday, the work needs a home. Make a folder per matter, client, cycle, or campaign; drop a markdown context file at its root; open it and run.

● Cowork

CLAUDE.md

Plus **Projects** (cross-session memory) and built-in **scheduled tasks** layered on top of the base pattern.

● OpenWork

AGENTS.md

Just folder + context file. You re-fire the prompt yourself; the file does the persistent heavy lifting.

THE TWO FAILURE MODES

Everything in **one folder** → context bleeds. Standalone sessions for recurring work → you re-explain forever.

Separate workstreams, separate folders

0

3

Rules & instructions

A layered instruction system — and the answer to the most common confusion: "why did the agent ignore what I said?"

Global is sparse • folder is specific • session is the goal

GLOBAL • every session, forever

Your role, default tone, output formats. Should fit in two short paragraphs.

FOLDER • while this folder is in scope

Naming conventions, this matter's terminology, file layout, house style.

SESSION • this task, this turn

The actual outcome you want right now.

The most common mistake: **putting everything in global**. The result is a 3,000-token system prompt that costs you on every turn and confuses the agent with rules that don't apply.

What's *not* in global: the matter-specific naming, the citation form for this filing, the folder layout. Those belong with the matter.

Cowork: Settings > Cowork • OpenWork: Settings panel
+ AGENTS.md

Ask me questions before you execute

End a non-trivial brief with one line: **"ask me 1-2 clarifying questions before you start."** It surfaces the unstated assumptions that would otherwise become bugs.

```
"Should I include the canceled
subscriptions in the count?"
```

```
"Weight technical fit and culture
fit equally, or one more heavily?"
```

Two questions · ninety seconds · a much better deliverable

A SECOND LINE, FOR MULTI-SOURCE TASKS

"If any sources contradict on a material point, flag it — don't silently pick one."

Without it, the model smooths the conflict into a confident answer. For a litigator or an auditor, that's the failure mode that produces malpractice.

PART FOUR

04

Extending the tool

Four ways to give your co-worker new playbooks, new reach, new roles, and new surfaces to act on.

Four ways to extend — each a delegation tool

SKILLS

How it works

A procedure to follow when a matching task comes up.

CONNECTORS

Where it reaches

Read or write through an external service.

PLUGINS

What role it plays

A packaged bundle for a specific role.

MCP SERVERS

What surface it acts on

A richer or different surface to act on.

Skills shape **how** the agent works · connectors shape **where** it can reach · plugins shape **what role** it plays · MCPs shape **what surfaces** it can touch.

A playbook your co-worker keeps on a shelf

A folder with a SKILL.md at its root: a name and **description** (the spine the agent reads to decide whether to open it), then the procedure. The same format works in both tools.

Content loads on demand — only the spine registers up front, so installing many costs less context than you'd expect.

Three paths: a catalog · generated in chat · authored by hand

```
--- name: weekly-brief
description:
Generate a weekly status brief from a folder of meeting
notes --- 1. List files modified in the last 7 days.
2. Read each *meeting*.md file.
3. Produce a one-page brief:
- 3 bullets "what shipped"
- 3 bullets "what's at risk"
- 1 paragraph "next week's focus"
4. Save as weekly-brief-YYYY-MM-DD.md
```

Security: skills are trusted code. Only install from sources you trust — read community ones before enabling.

The value lives in combinations

One connector is useful. Three that work together unlock workflows that didn't exist before — each replacing twenty minutes of human context-switching.

LITIGATION

Pull the opposing-counsel thread in Outlook, cross-reference the deposition outline in OneDrive, draft a response that preserves our privilege position.

FINANCE

Pull the variance thread from controllership, cross-reference the close checklist in SharePoint, draft tomorrow's leadership update.

HR

Pull the panel debrief from Slack, cross-reference scorecards in Greenhouse, draft a hiring-manager recommendation memo.

Install a connector when a **specific workflow** needs it — never speculatively. [Cowork's](#) catalog is broad; [OpenWork's](#) is leaner with an Add-Custom-App path. Each one is also a new prompt-injection vector.

One word, two different things

● Cowork plugin

A role bundle

Skills + connectors + namespaced slash commands + sub-agents + config in one download — Sales, Finance, Legal, Marketing. **Gets you there faster out of the box.**

```
/legal:privilege-entry · /fin:variance-comment
```

● OpenCode plugin

An npm package with hooks

Single-purpose, extends the underlying engine. Not a role bundle. **More flexible if you compose your own.** Installed by package name.

The two are not interchangeable

THE TRAP

Install plugins like browser extensions — one at a time, for a specific workflow, audited monthly. Seven at once gave one marketer 43 colliding slash commands and an MCP server she didn't remember authorizing.

Parallel workers, one clean main thread

When work breaks into parallel pieces, the agent spawns sub-agents that each handle a slice at once. Each works in its own context, so only the **result** comes back — not the raw documents. A 30-minute sequential job becomes a 5-minute parallel one.

FAN-OUT

“For each of N items, do X”

12 transcripts → a one-page summary
each → synthesize themes across all.

DIMENSION

“Analyze X across N angles”

Audit an MSA across indemnification,
liability cap, IP, termination, governing law.

COMPARE

“Compare A and B”

Last quarter's strategy vs this quarter's —
extract priorities, find what changed.

When not to: genuinely sequential work, small batches (under ~5-7 items), or when coherence across items matters more than throughput. You don't write special syntax — you frame the task so the parallelism is obvious.

PART FIVE

0

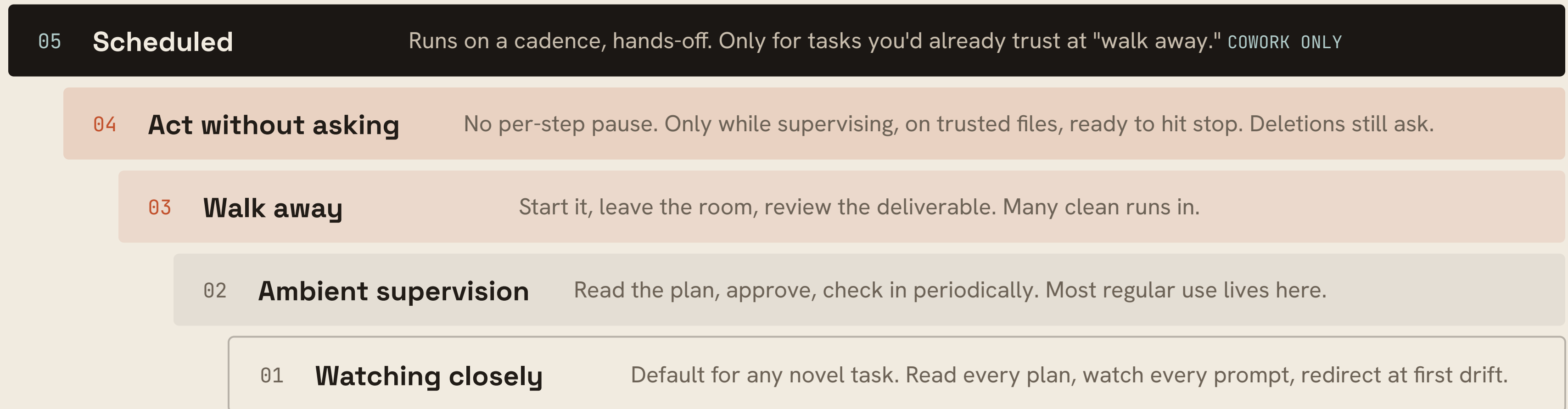
5

Safety & the autonomy ladder

The discipline that matters most — and
it applies identically in both tools.

Climb one rung per task type — with track record

Trust grows with track record on *this* work — not with how long the tool's been installed.



Step **back down** when a task type changes — new client, new connector, new edge case. The autonomy you've earned is task-specific, not skill-generic.

Untrusted text the agent reads as commands

A lawyer uploaded a 40-page vendor proposal. The plan came back with a step he never asked for:

“...then send a copy to [external email] for the vendor’s records.”

Buried in white-on-white text on page 32. His eyes skipped it; the agent read it as an instruction. **"Ask before acting" is what saved him.**

PRACTICAL DEFENSES

- Never run high-autonomy modes on tasks touching untrusted content.
- Watch for scope creep — files or connectors you didn't mention mean **do not approve.**
- Hit **Stop** the moment things drift. Ask questions after.

If you wouldn't trust it to "walk away," don't schedule it

Cowork has built-in scheduling — `/schedule` from chat or a full Scheduled → New task form. **OpenWork** has none: keep the prompt in the folder, fire it on a calendar reminder.

SCHEDULES WELL

Information-gathering · bounded outputs that never send or buy · tasks watched succeed 3+ times.

DOESN'T

Anything that sends, pays, touches privileged files, or processes strangers' content without review.

```
/schedule share weekly content  
updates every Monday 9 AM...
```

A scheduled task runs **without you watching** — and you can't course-correct a task you're not watching.

Build the path deliberately: supervised → walk-away → scheduled, a week between each.

06

In practice

One recurring task — the Monday-morning industry brief — walked up the autonomy ladder, run in both tools.

Walking the ladder, deliberately

01

Set it up

A project (Cowork) or workspace (OpenWork) with sources, output format, and tone in the instructions.

02

Run it watched

Trigger it end-to-end while watching. Check the deliverable, refine the instructions.

03

Run it again

Good twice in a row with no edits? Calibration confirmed — now you're ready.

04

Make it recur

Cowork: `/schedule` for Sunday 9pm. OpenWork: a calendar reminder + manual re-fire.

The shape — **workspace, instructions, manual runs, recur once trusted, feedback loop** — is the template for every recurring workflow: Friday cleanup, daily inbox triage, month-end variance commentary, weekly matter-status updates.

Where the value compounds

CONNECTOR COMBINATIONS

Chain, don't single

The real wins come from combining tools. Any sentence with "...and then I open the other tab to..." is a candidate.

MONTHLY AUDITS

Review the surface

Folders, connectors, skills, plugins, scheduled tasks. Last month's experiment is this month's forgotten access. Ten minutes.

TEAM SCALE & REVIEW

Shared, but not autonomy

Share skills, plugins, templates like house style. The autonomy ladder stays individual. Cross-model review for high-stakes outputs.

THE CHAPTER IN FIVE LINES

- 1 Delegate, don't query.
- 2 Three trust levers — folders, connectors, approvals — set deliberately.
- 3 The plan is the leverage. Read it before any file is touched.
- 4 Climb the autonomy ladder one rung at a time; step back down when the task changes.
- 5 Untrusted content always gets the cautious mode.

THE ONLY WAY FROM HERE

Reading this doesn't make you good at either tool. Using it does.

Make a working folder. Run one read-only task today. Read the plan before you approve. Climb one rung when you've earned it. The discipline is what to bring before the tool — not a substitute for it.

● Cowork

● OpenWork