



Equipping Your Digital Armor

World 4: Safe and Responsible AI Use

Defense Enables Offense



The Goal

Build safe default habits for everyday AI use.



The Reality

AI assistants are powerful tools, but they lack human context for privacy and risk.



The Mindset

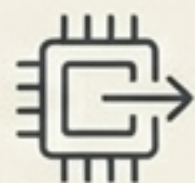
Safety rules are not barriers. They are the digital armor that allows you to use Claude confidently without exposing sensitive information.

The Quarantine List: Do Not Paste

Never enter these into AI tools unless you have explicit permission and understand the underlying privacy rules.



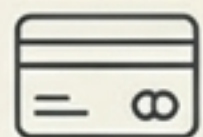
Passwords



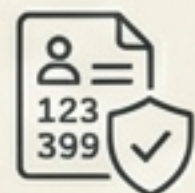
API keys



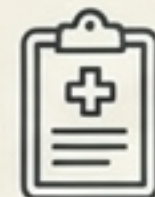
Private customer records



Payment card information



Social insurance / SSNs



Medical records



Legal documents with private facts



Confidential business plans

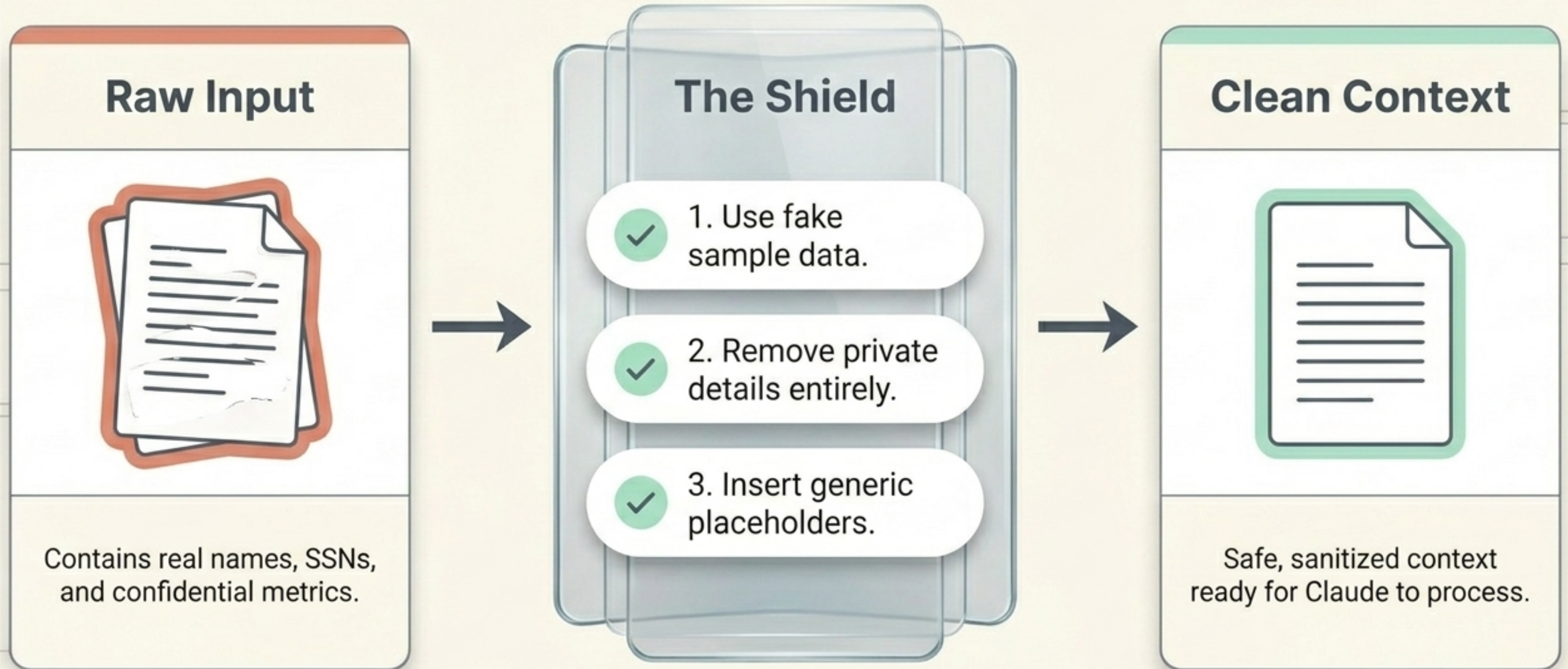


Internal profit, margin, markup, or cost data

The Data Triage Board

Safe to Share	Redact First	Never Paste
<p data-bbox="236 784 1096 1014">A public website link</p> <p data-bbox="236 1061 1096 1292">A general question about writing better prompts</p> <p data-bbox="236 1339 1096 1570">A company policy that is already public</p>	<p data-bbox="1236 784 2095 1014">A fake sample email (based on a real scenario)</p> <p data-bbox="1236 1061 2095 1405">Internal drafts (with sensitive names/data removed)</p>	<p data-bbox="2235 784 3095 1014">A private contract</p> <p data-bbox="2235 1061 3095 1292">A customer phone number</p> <p data-bbox="2235 1339 3095 1570">Payment card details</p>

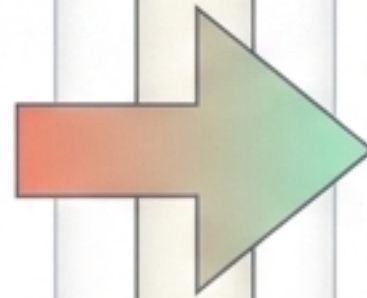
The Scrubbing Engine



Tactical Redaction in Action

Risky Prompt

Help me rewrite this email to **John Doe** at **Acme Corp**. His account is past due, and his API key key is **sk-1234abcd**. His credit card ending in **4321** failed.



Scrubbed Prompt

Help me rewrite this email to **[CUSTOMER_NAME]** at **[COMPANY]**. Their account is past due, and their API key **[API_KEY_HE]** needs rotation. Their payment method ending in **[XXXX]** failed.

The Abstraction Workflow

Core Insight: You don't feed the AI the data; you feed the AI the structure.

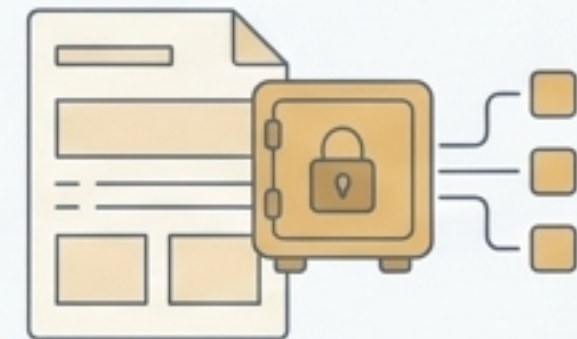
Step 1: AI Task

Ask Claude to generate a blank template, checklist, or structural framework.



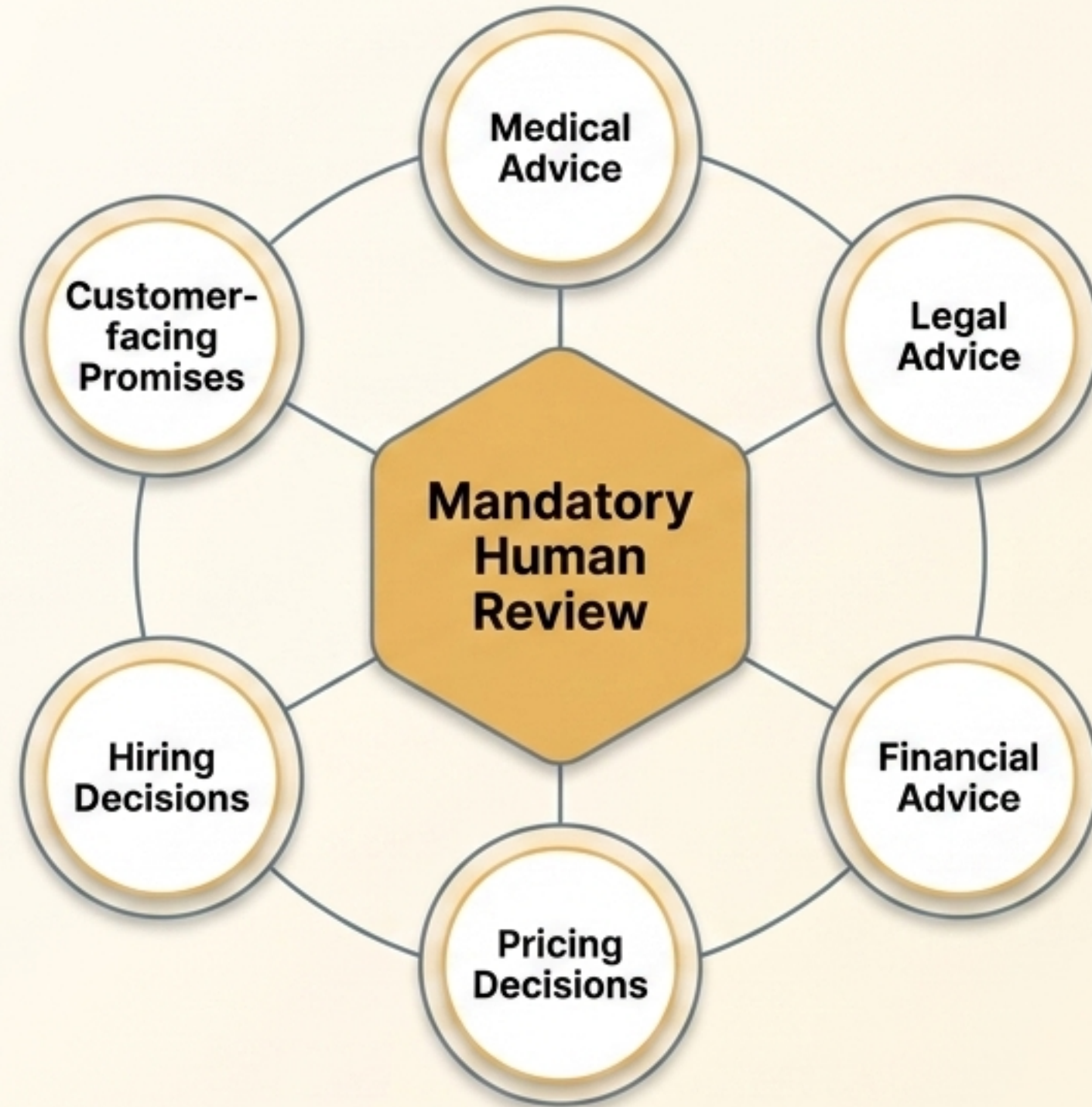
Step 2: Human Task

Take the AI-generated template offline. Securely merge your private, confidential data into the document on your own trusted system.



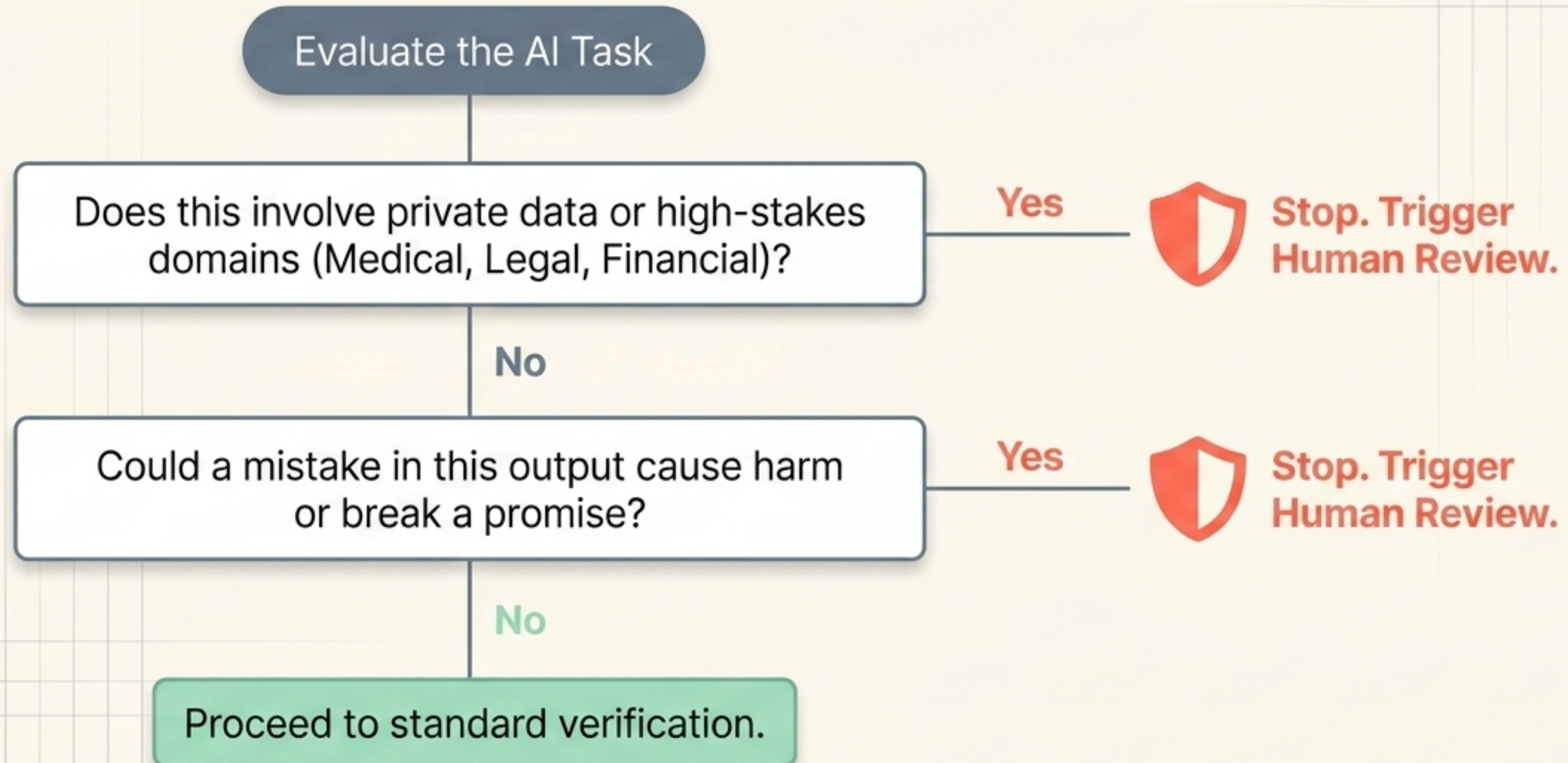
Takeaway: Keep the private data on your side of the screen.

High-Stakes Domains



Guiding Principle: Anything that could harm someone if wrong requires strict human oversight. AI assists; humans decide.

The Oversight Decision Tree



The Verification Rule



Confidence is not correctness.

AI assistants can suffer from “Hallucinations”—giving answers that sound confident but are wrong, unsupported, outdated, or completely made up.

The 5-Step Verification Checklist

- 1 Identify the core claim.
- 2 Check against a trusted source.
- 3 Verify all names, dates, and numbers.
- 4 Ask: Could a mistake here cause harm?
- 5 Get human review for high-stakes topics.

World 4 Checkpoint Cleared

The Red List (Unlocked)

You know not to paste passwords, API keys, private customer data, or sensitive business records.

The Tactical Shield (Unlocked)

You know how to use fake sample data, redaction, and [PLACEHOLDERS] to scrub prompts.

The Verification Habit (Unlocked)

You know that confidence isn't correctness, and how to trigger human review for high-stakes decisions.

[Next Objective: World 5 - Real-World Workflows] →