











The 2026 GenAI Stack



The new AI engineer map.

Not just prompts.
Not just RAG.
Not just agents.

 A production GenAI system now needs:



	MODELS	Choose the right model for the right task
	CONTEXT	Give the model the right information
	TOOLS	Let the model act in the real world
	EVALS	Measure quality, safety, and performance
	TRACES	See what happened, step by step
	GUARDRAILS	Control risk and prevent bad outcomes
	HUMANS	Add judgment and approvals where needed
	OPERATIONS	Run, monitor, and improve in production

 Save this. You'll need it. 

The shift



The skill has moved from:



"Can I get a good answer?"

to:

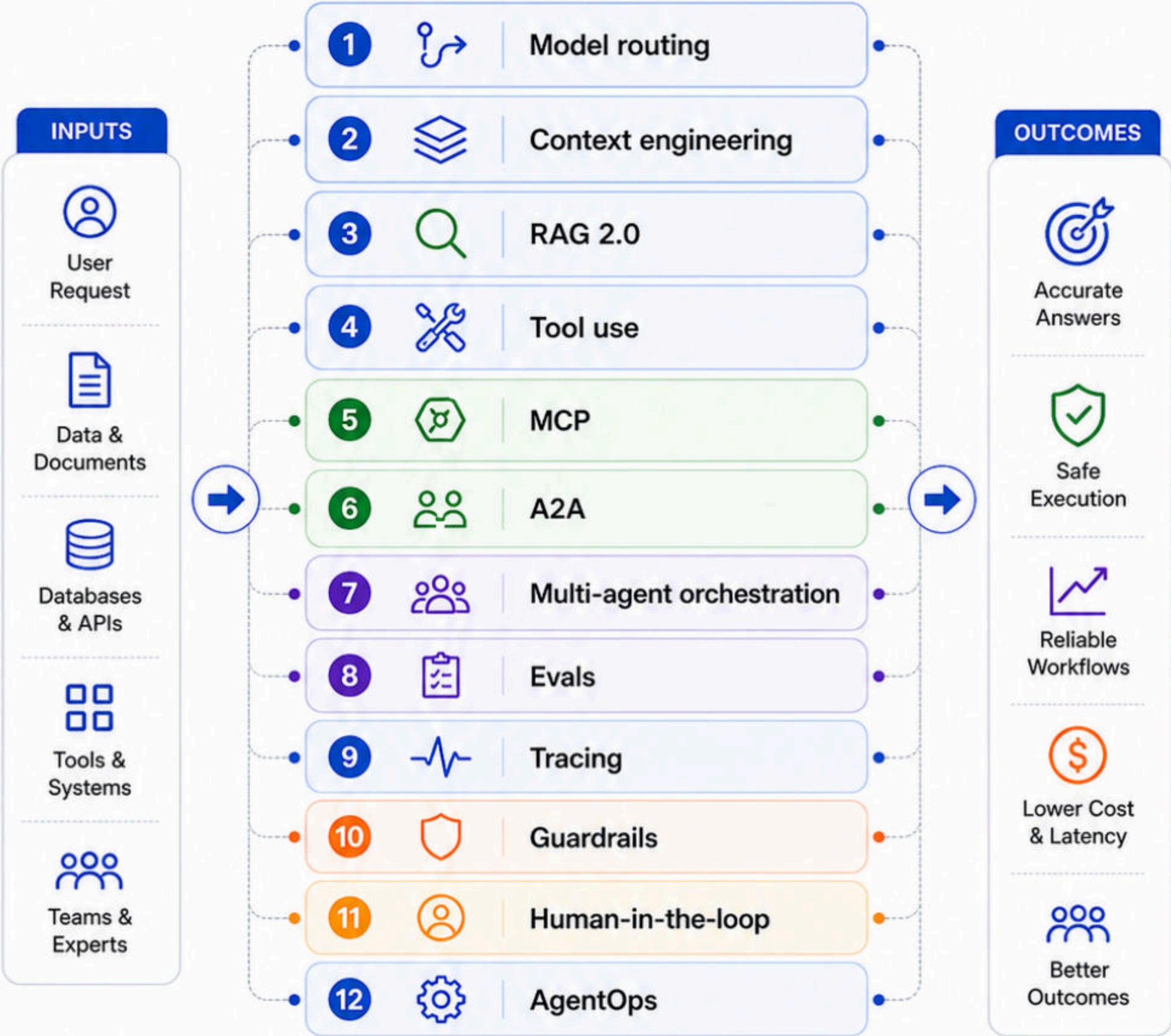


"Can this system complete real tasks reliably, safely, and repeatedly?"




Save this. You'll need it.

The 2026 GenAI Stack

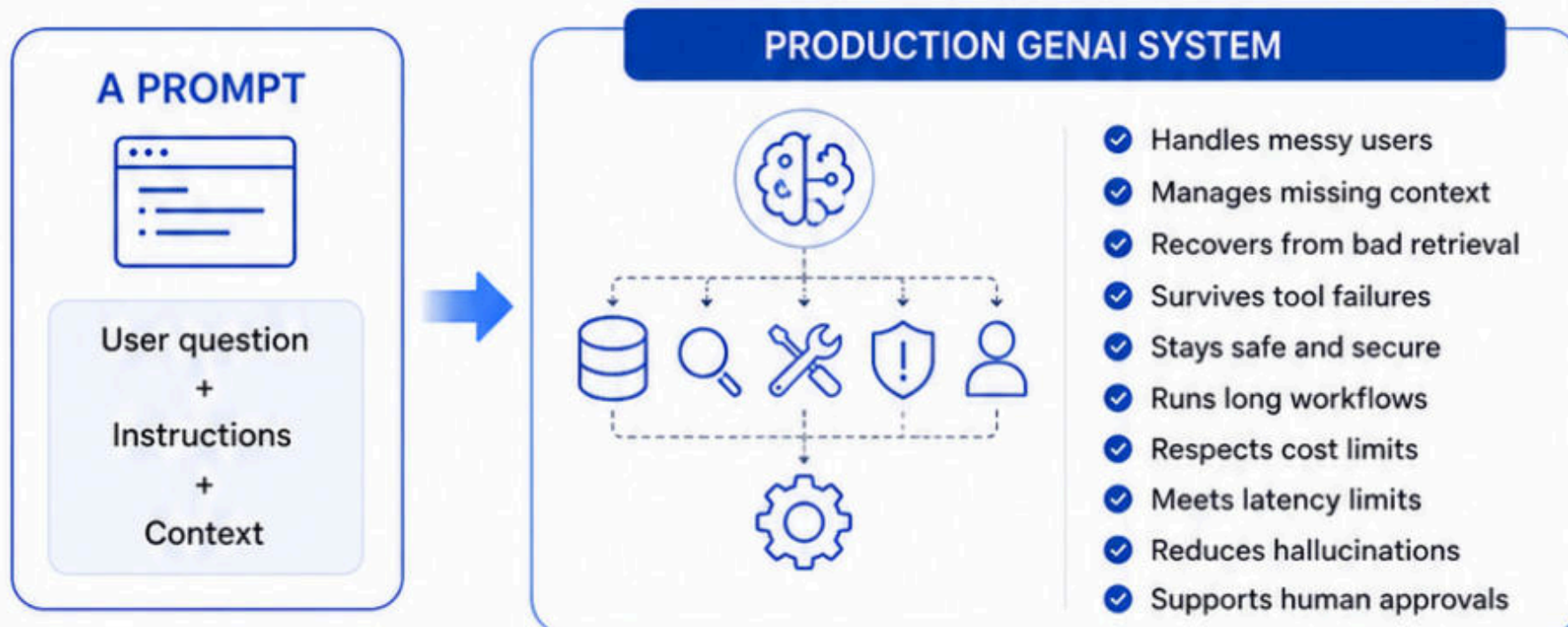
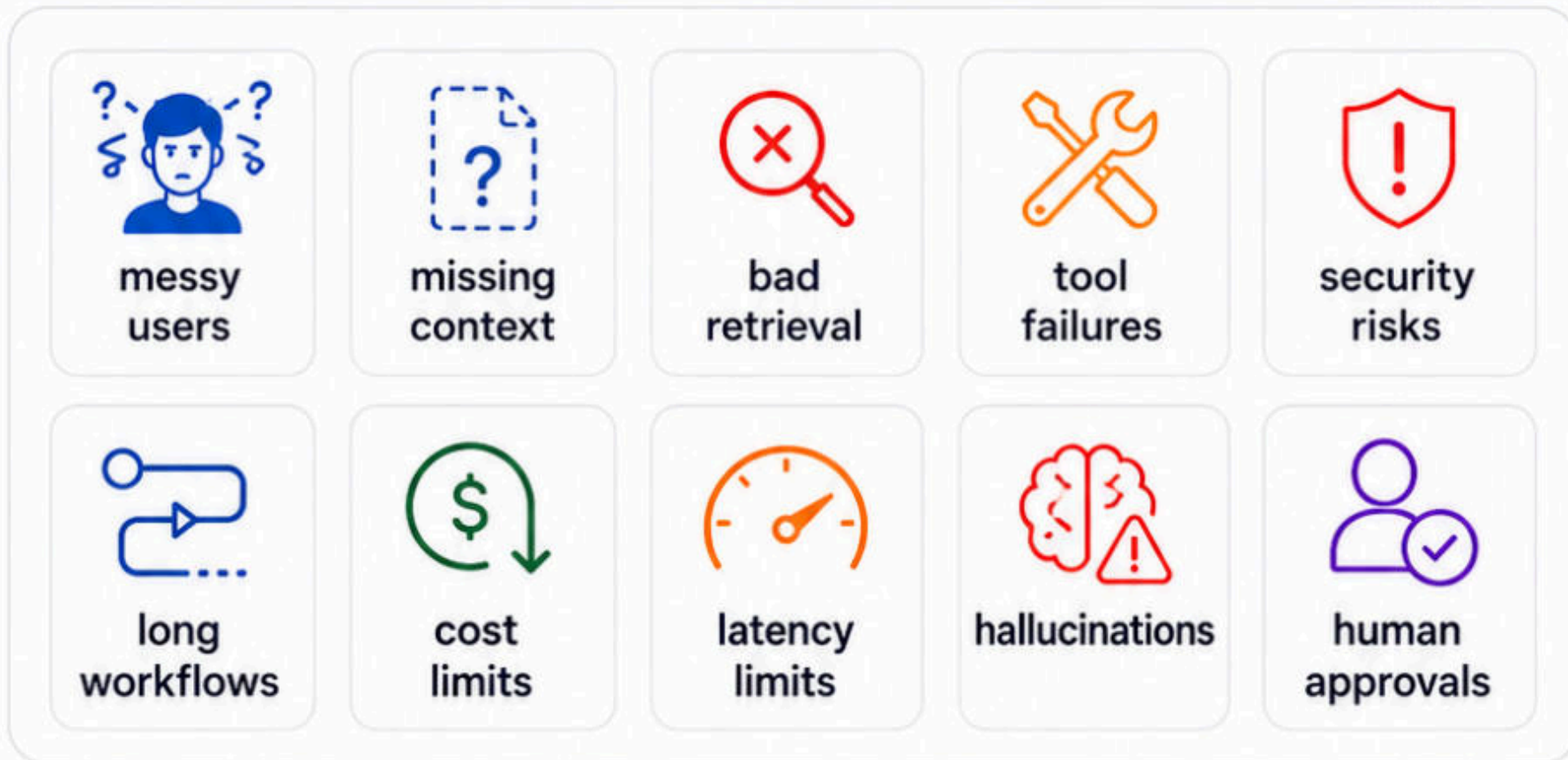


 **This is the new AI engineer map.**

 Save this slide — you'll need it.

Why prompting is not enough anymore

A prompt can improve one answer.
But production systems need to handle:



Prompting is still useful.
But it is no longer the whole stack.



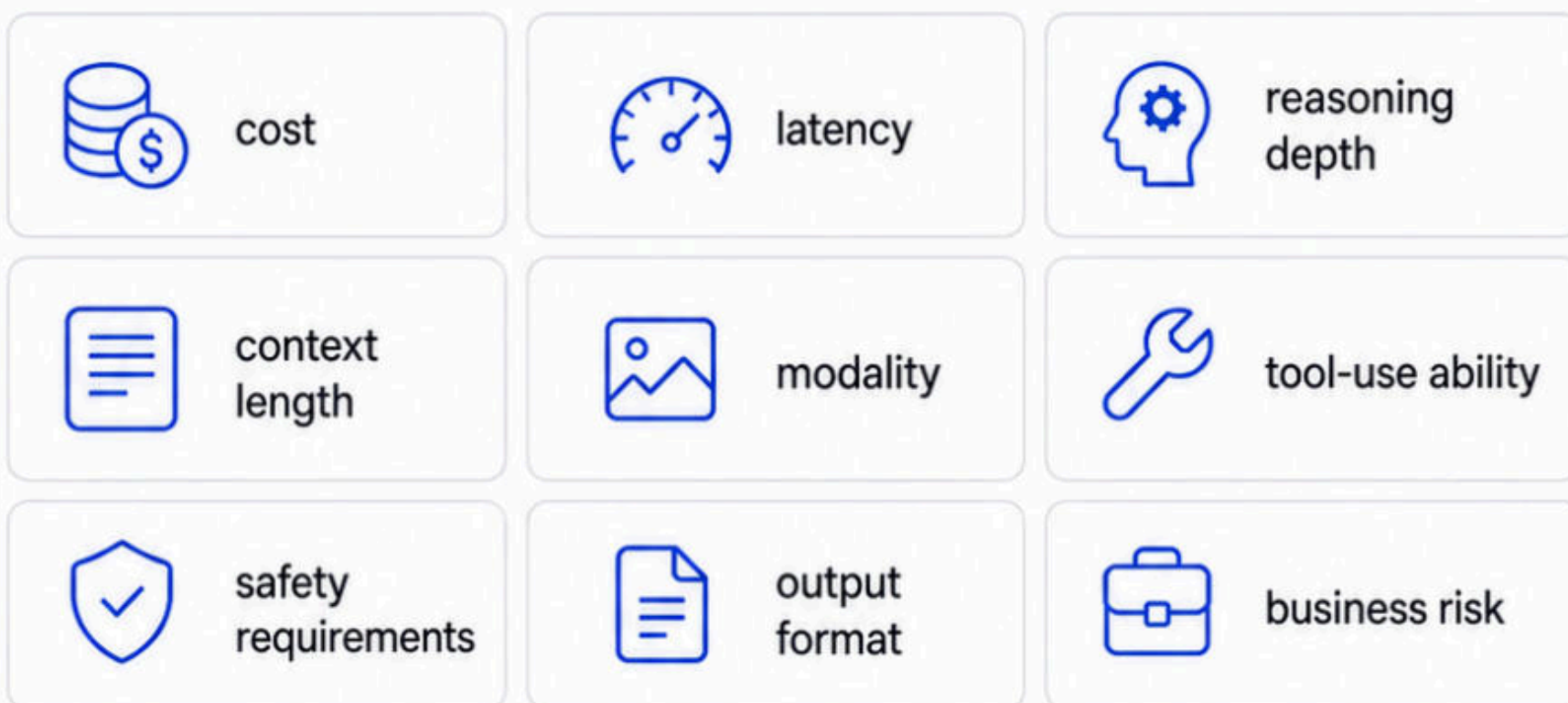
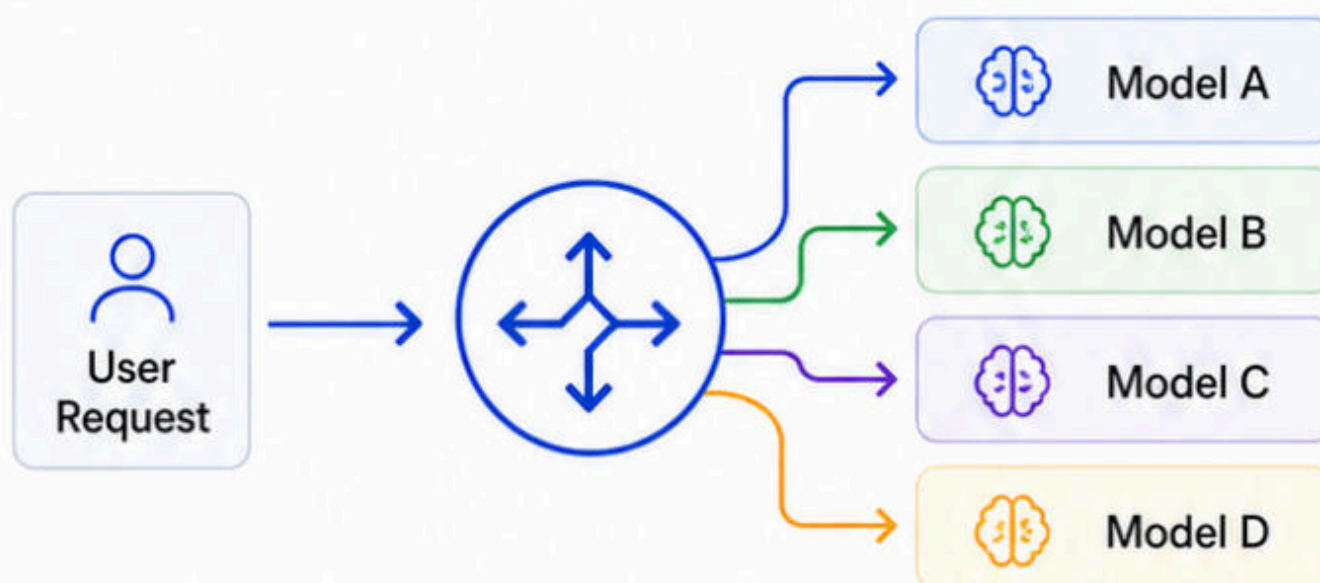
Save this slide for quick revision!

Layer 1: Model Routing



One model for everything is **lazy architecture**.

In 2026, smart GenAI systems **route tasks to different models** based on:



✗ The question is not:

Which model is best?

VS

✓ The question is:















Which model is best for this task?



Save this. You'll need it.

Model Routing: Practical Map

Use different models for different jobs:

JOB	MODEL TYPE
 Simple classification →  small, fast model	
 Complex reasoning →  stronger reasoning model	
 Code generation →  coding-specialized model	
 Summarization →  low-latency model	
 Vision tasks →  multimodal model	
 Sensitive decisions →  safer, controlled model	
 Final answer →  high-quality response model	

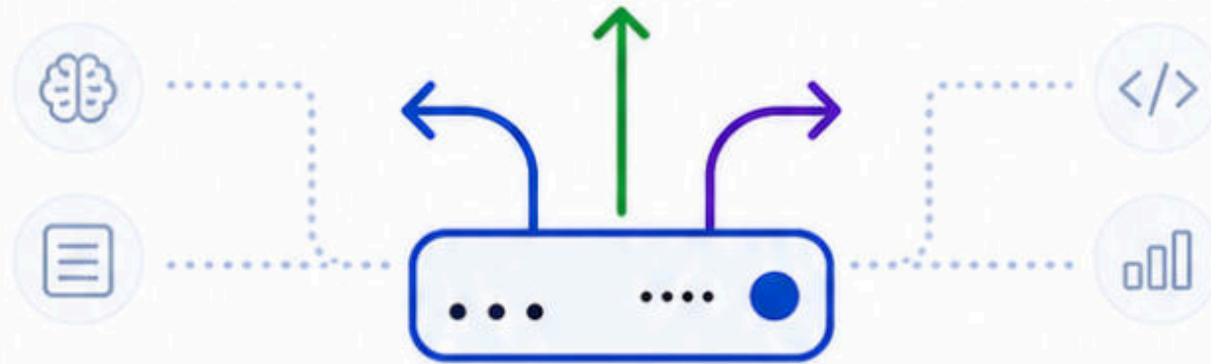


The best GenAI apps are not model-loyal.
They are **outcome-loyal.**



Save this slide for quick revision!

Model Routing: What to track



Before routing models, define:

- ✓ task type
- ✓ input complexity
- ✓ required accuracy
- ✓ allowed latency
- ✓ max cost
- ✓ safety risk
- ✓ tool requirement
- ✓ fallback model
- ✓ escalation rule

Bad routing gives you:

- high cost
- slow output
- unpredictable quality

Good routing gives you:

- speed
- quality
- control

Save this slide for quick revision!

Layer 2: Context Engineering



Prompt engineering asks:

“What should I say to the model?”

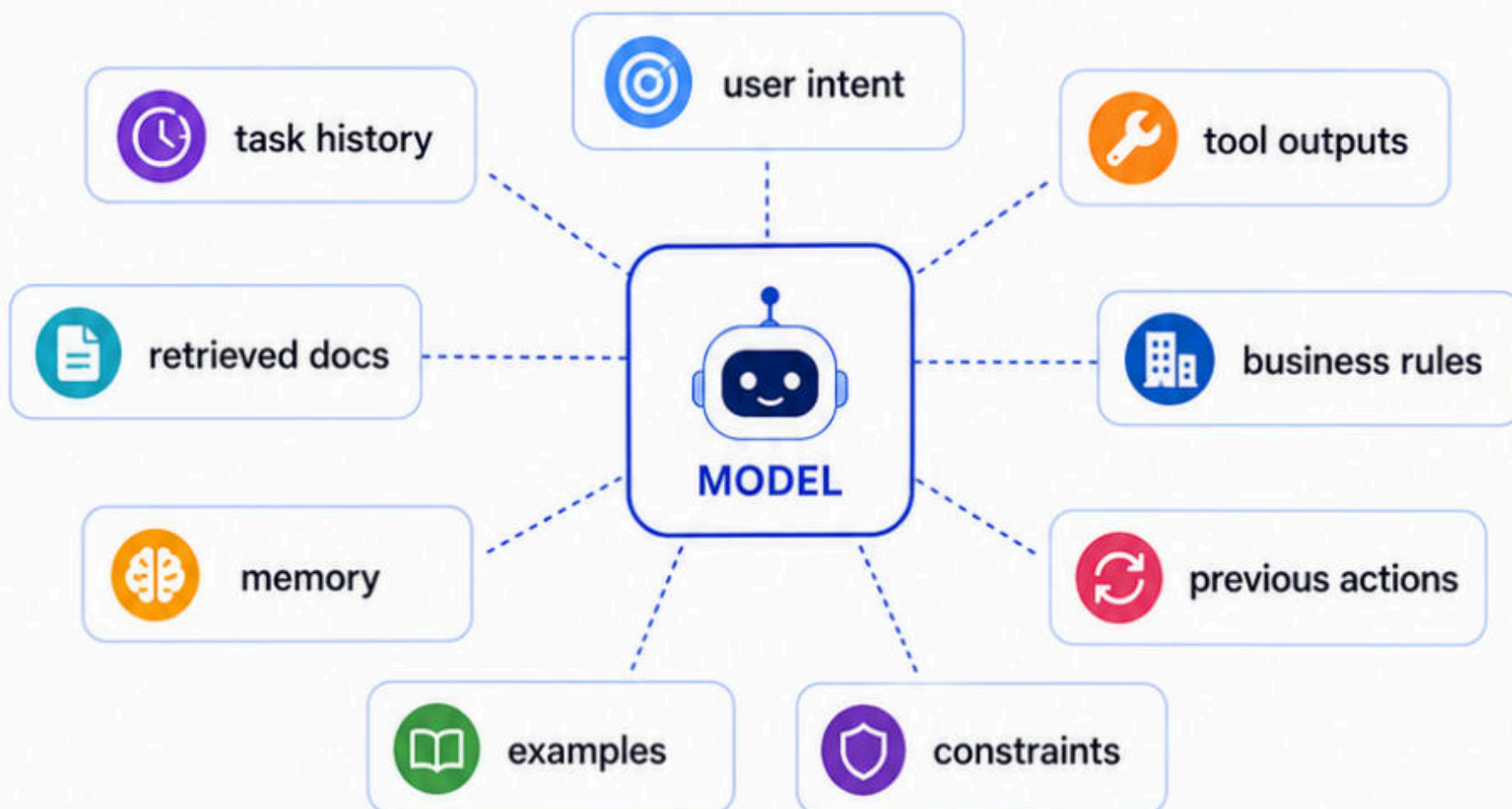


Context engineering asks:

“What should the model know before it answers?”



This includes:

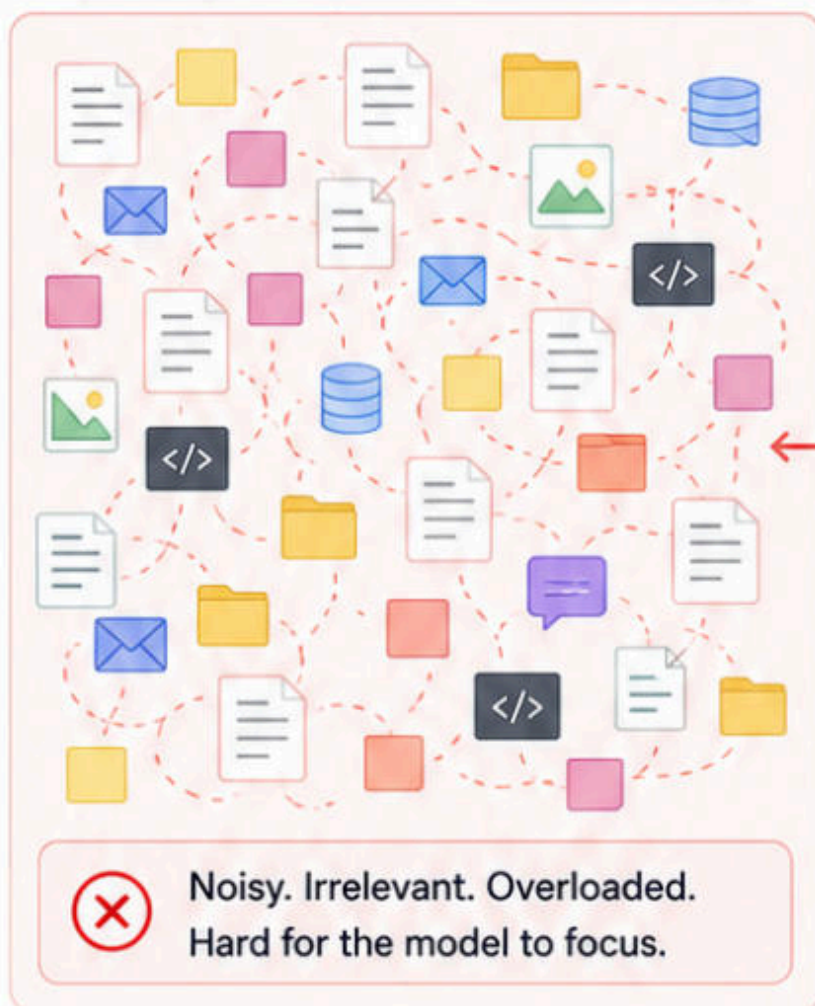


In 2026, context is the real prompt.

Context Engineering: The mistake

Most teams **dump more context**.
Better teams **select better context**.

MORE CONTEXT (THE MISTAKE)



BETTER CONTEXT (THE PRACTICE)



GREAT SYSTEMS KNOW:



More tokens \neq better intelligence.

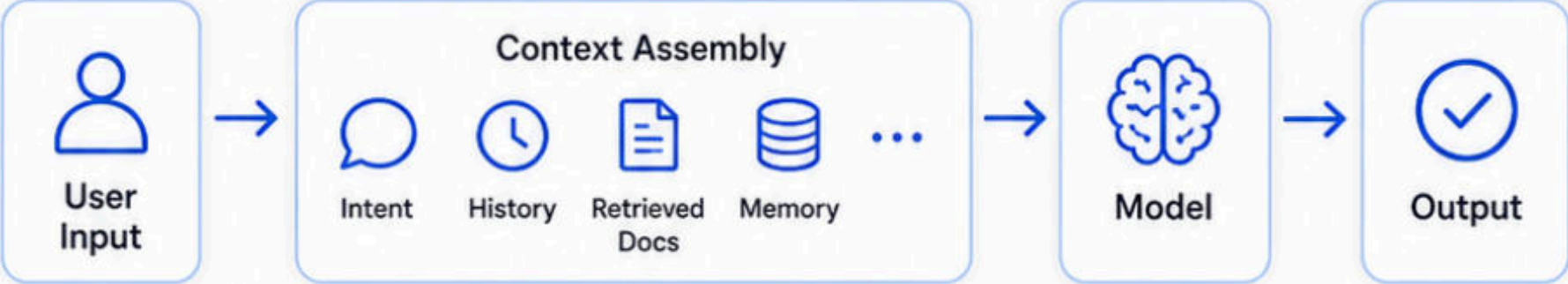



Better context = better output.













Save this slide for quick revision!

Context Engineering Checklist

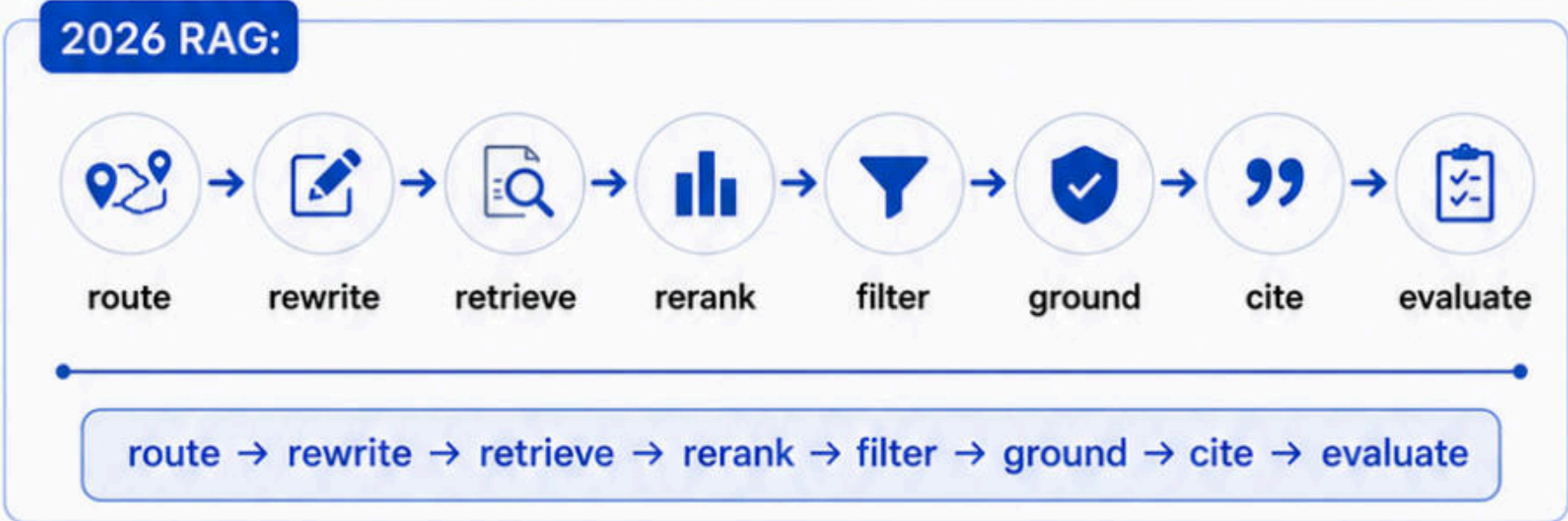


 **Before sending context to the model, ask:**

-  Is this relevant?
-  Is this too much?
-  Is this fresh?
-  Is this structured?
-  Is this trusted?
-  Is this user-specific?
-  Is this allowed?
-  Is this task-specific?
-  Is this complete?
-  Is this safe to expose?

 **Context quality is system quality.** 

Layer 3: RAG 2.0



- Advanced RAG proves:**
- ✓ where the answer came from
 - ✓ why the source is relevant
 - ✓ whether the source is fresh
 - ✓ whether the model stayed grounded
 - ✓ whether the final answer is faithful

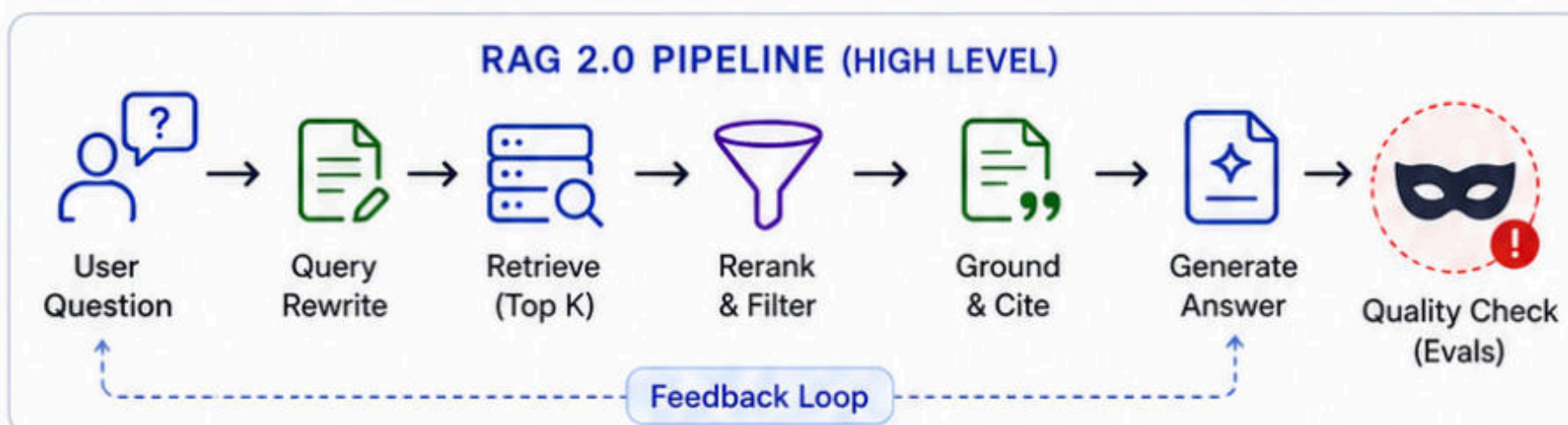
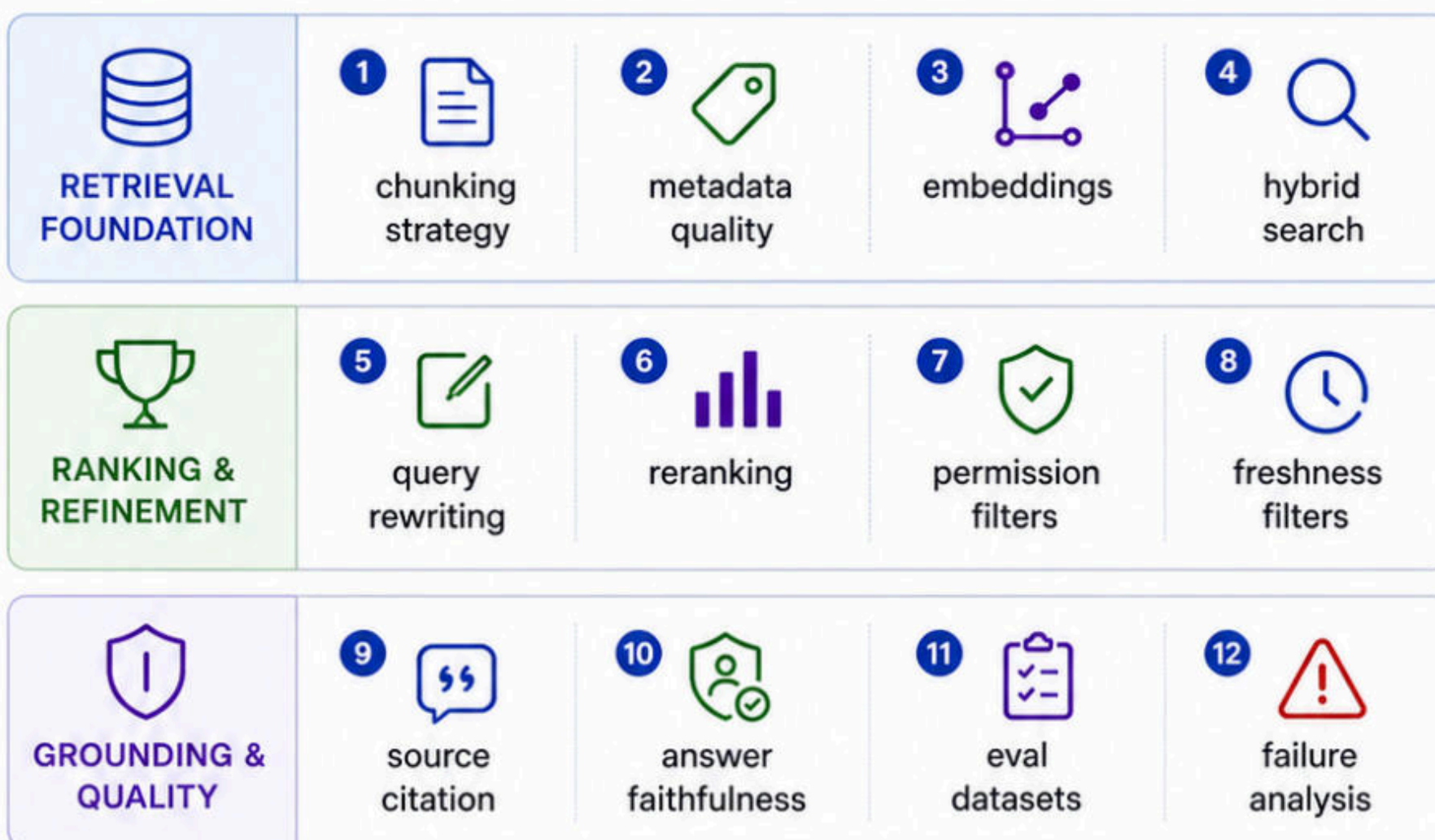
 **RAG is not dead.** |  **Lazy RAG is dead.**

RAG 2.0:

What actually matters



Your RAG quality depends on:



Most RAG failures are retrieval failures wearing a generation mask.












Save this slide for quick revision!

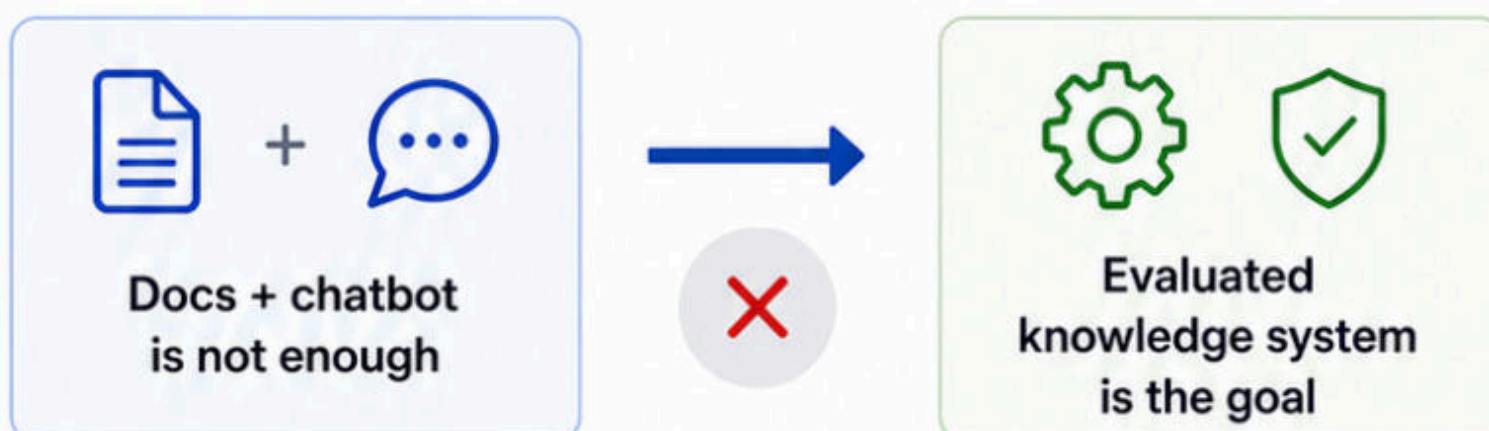
RAG 2.0:

Production questions

Before shipping RAG, ask:

<input type="checkbox"/> Can it handle vague queries? 	<input type="checkbox"/> Can it handle conflicting sources? 	<input type="checkbox"/> Can it avoid stale documents? 
<input type="checkbox"/> Can it respect user permissions? 	<input type="checkbox"/> Can it say "I don't know"? 	<input type="checkbox"/> Can it cite correctly? 
<input type="checkbox"/> Can it detect missing context? 	<input type="checkbox"/> Can it recover after bad retrieval? 	<input type="checkbox"/> Can it be evaluated automatically? 

From a document chatbot...



A chatbot with documents is **not** automatically a knowledge system.

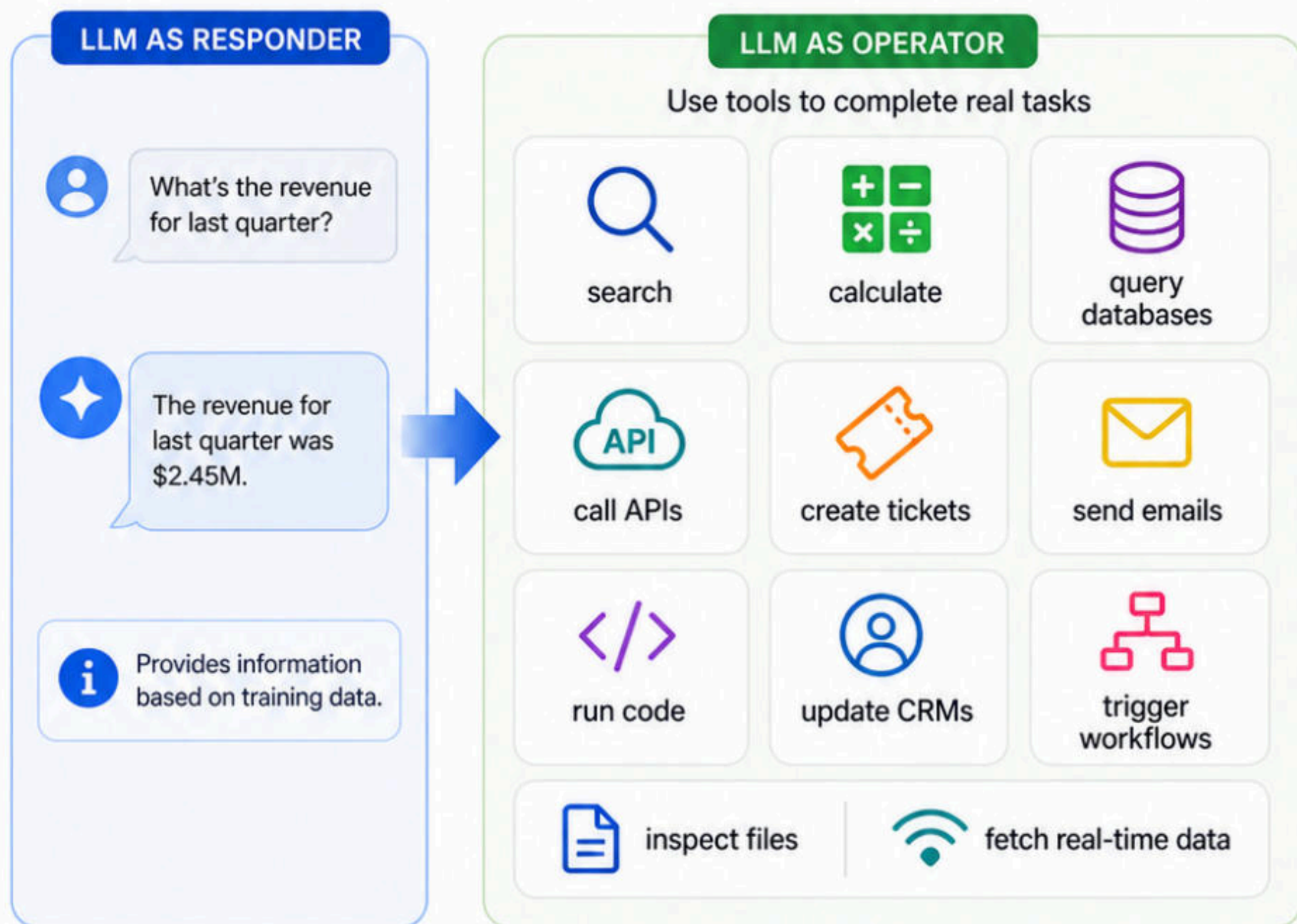


Save this slide for quick revision!

Layer 4: Tool Use

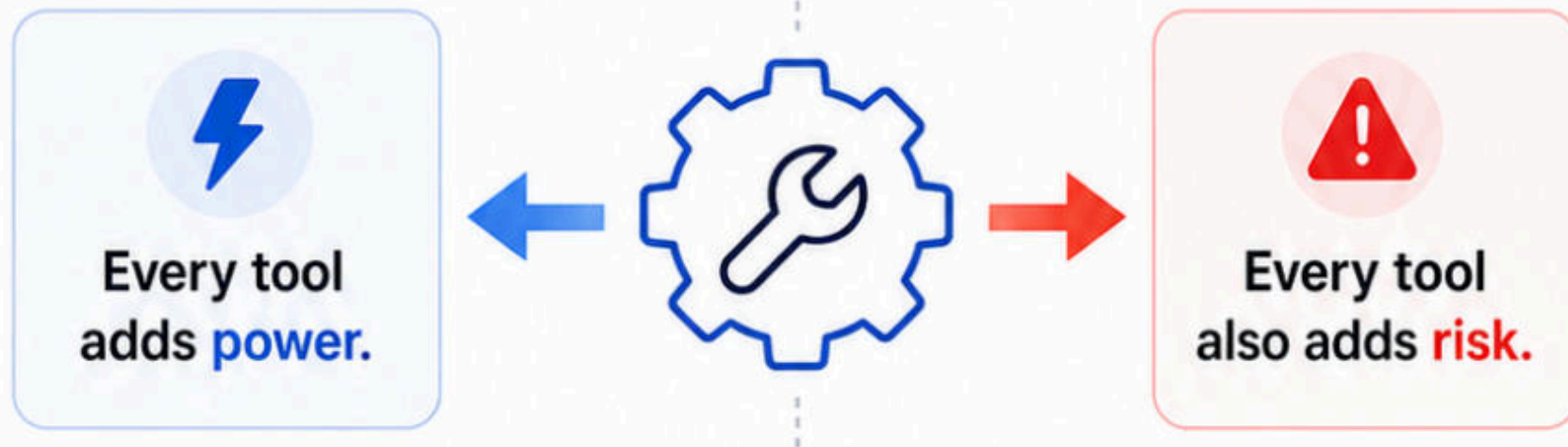
LLMs talk.
Tools let them act.

With **tools**, a GenAI system can:

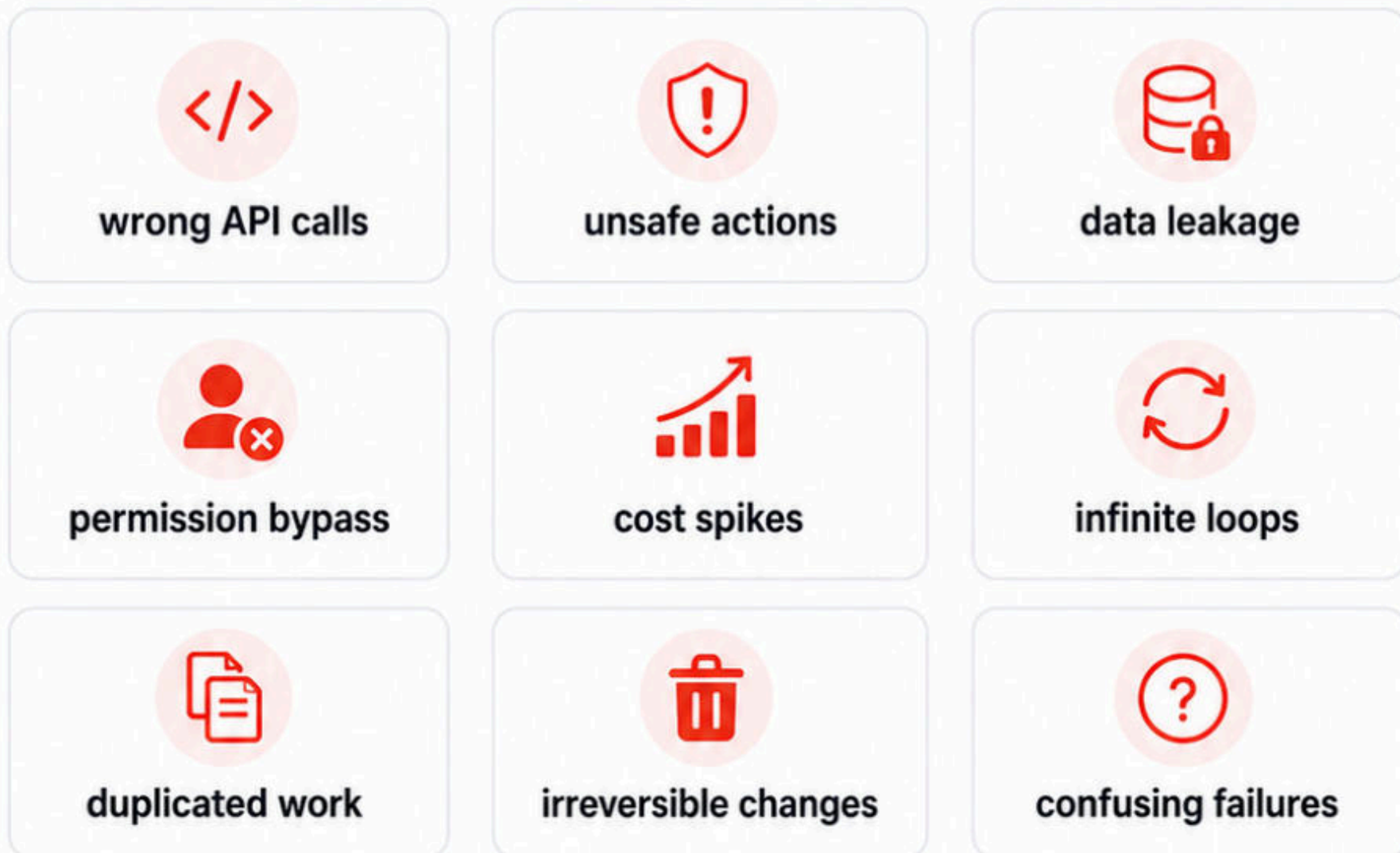


Tool use turns a model from
a **responder** into an **operator**.

Tool Use: The hidden risk



! Bad tool design causes:



The model is not the only thing that needs design.


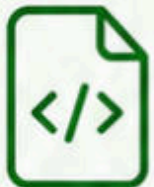


The tools need design too.





Tool Use Checklist

For every tool, define:

 Identity	1 clear name <input type="checkbox"/>
	2 clear description <input type="checkbox"/>
 Contract	3 input schema <input type="checkbox"/>
	4 output schema <input type="checkbox"/>
	5 permission level <input type="checkbox"/>
	6 failure behavior <input type="checkbox"/>
 Safety	7 timeout limit <input type="checkbox"/>
	8 retry policy <input type="checkbox"/>
 Operations	9 audit log <input type="checkbox"/>
	10 approval requirement <input type="checkbox"/>
	11 cost impact <input type="checkbox"/>
	12 rollback path <input type="checkbox"/>



A vague tool is a dangerous tool.



Save this slide for quick revision!

Layer 5: MCP

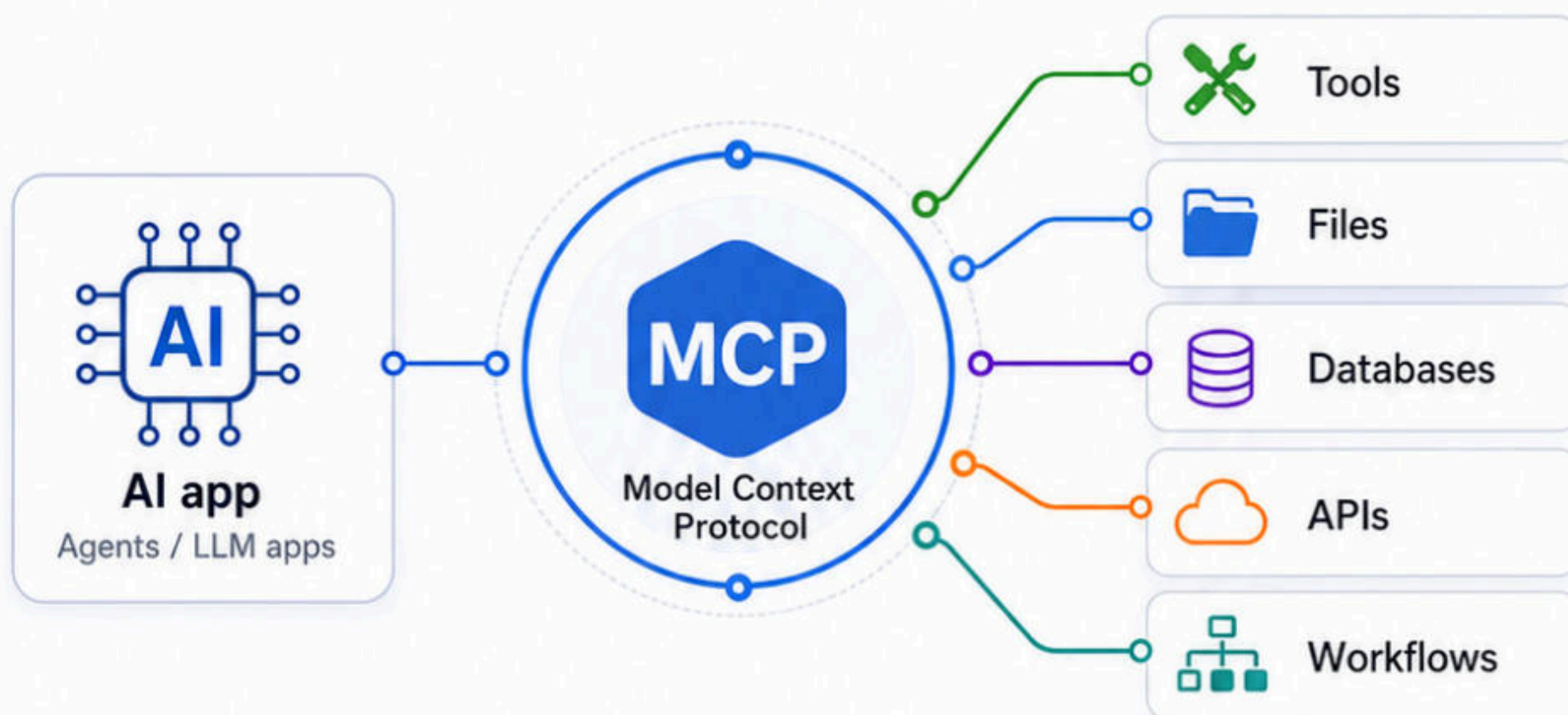
Model Context Protocol.

MCP is a standard way to connect AI applications with external systems.



Think:

AI app → MCP → tools, files, databases, APIs, workflows



Why it matters:



less custom integration



reusable tool & connectors



cleaner context access



better agent workflows



easier enterprise adoption



more standardized tool use

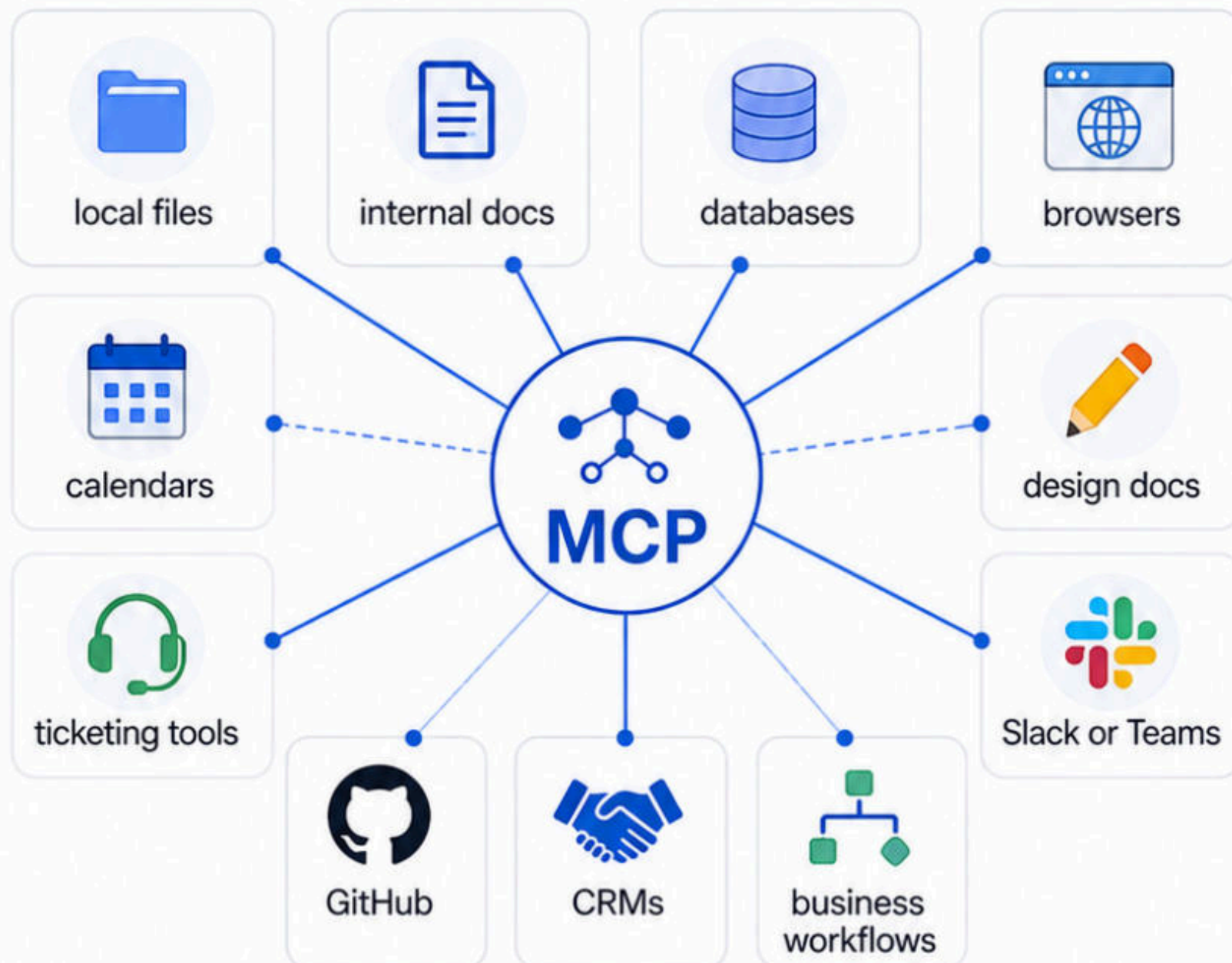


MCP is becoming one of the **key** integration layers for agentic systems.



MCP: What GenAI engineers should know

MCP is useful when agents **need access to:**



THE OLD WAY

One **chatbot**.



✘ Limited context. Limited impact.

THE 2026 WAY

AI systems **connected** to your **actual work tools**.



✔ More context. More capability. Real impact.

The future is **not** one chatbot.

The future is **AI systems** connected to your **actual work tools**.



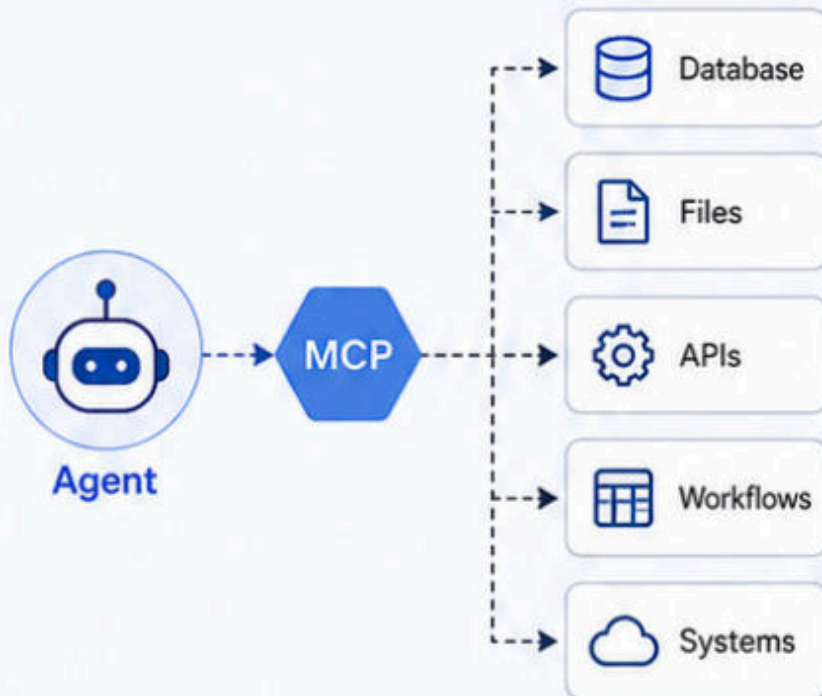
Save this slide for quick revision!

Layer 6: A2A

Agent2Agent.

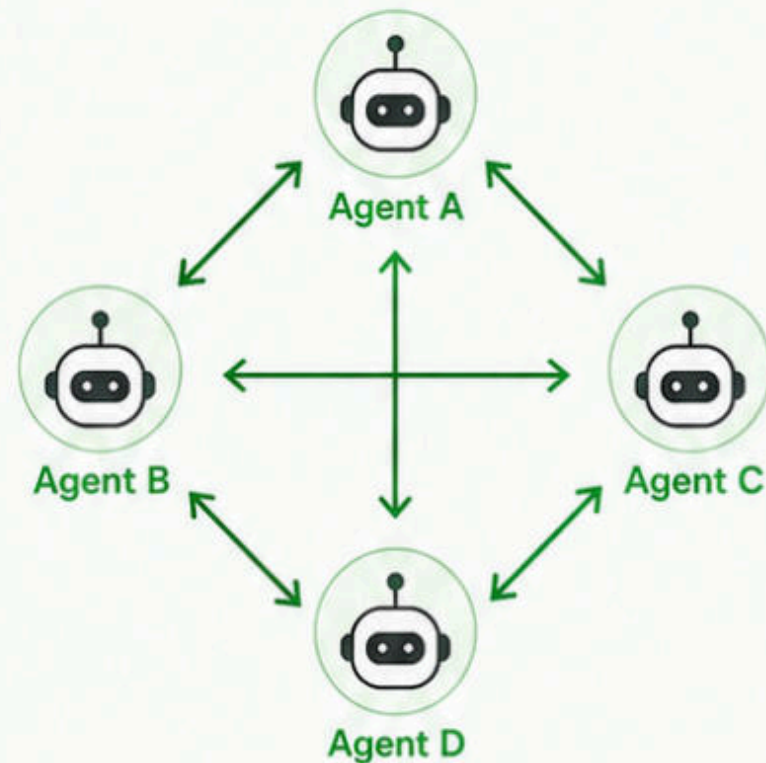
MCP connects:

agents → tools/data



A2A connects:

agents → other agents



This matters when different agents need to:

- communicate
- share task state
- exchange information
- delegate work
- coordinate actions
- collaborate across platforms
- complete enterprise workflows



A2A is about agent interoperability.



Save this slide for quick revision!

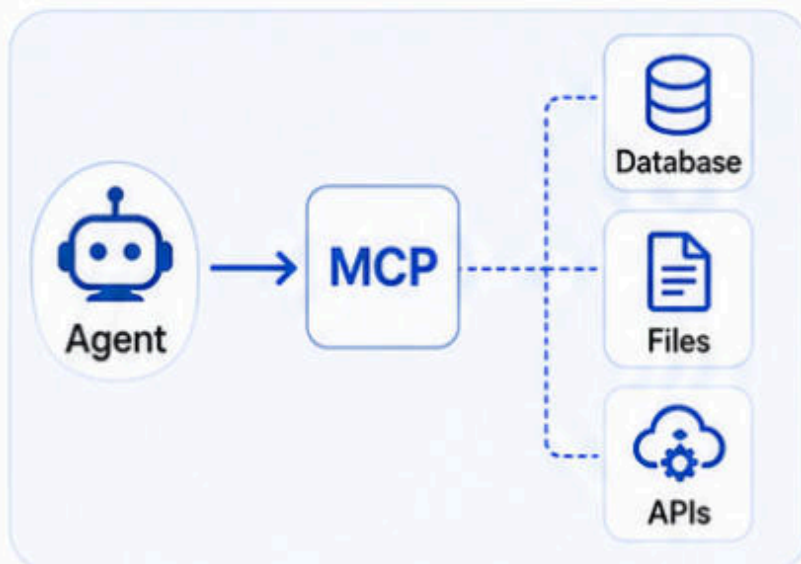
MCP vs A2A

Don't confuse them.

MCP

MCP connects:
agent to tools/data

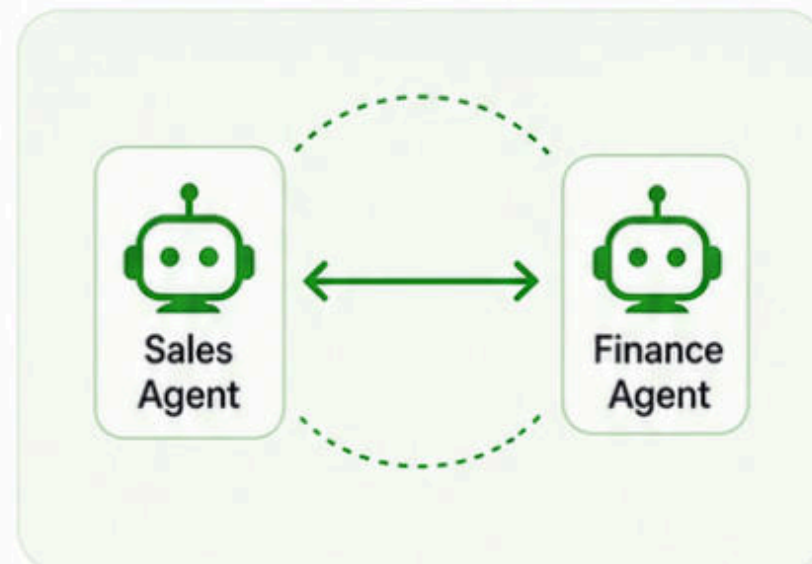
Example:
agent reads database



A2A

A2A connects:
agent to agent

Example:
sales agent talks to finance agent



Simple way to remember:



MCP = tool access

A2A = agent collaboration

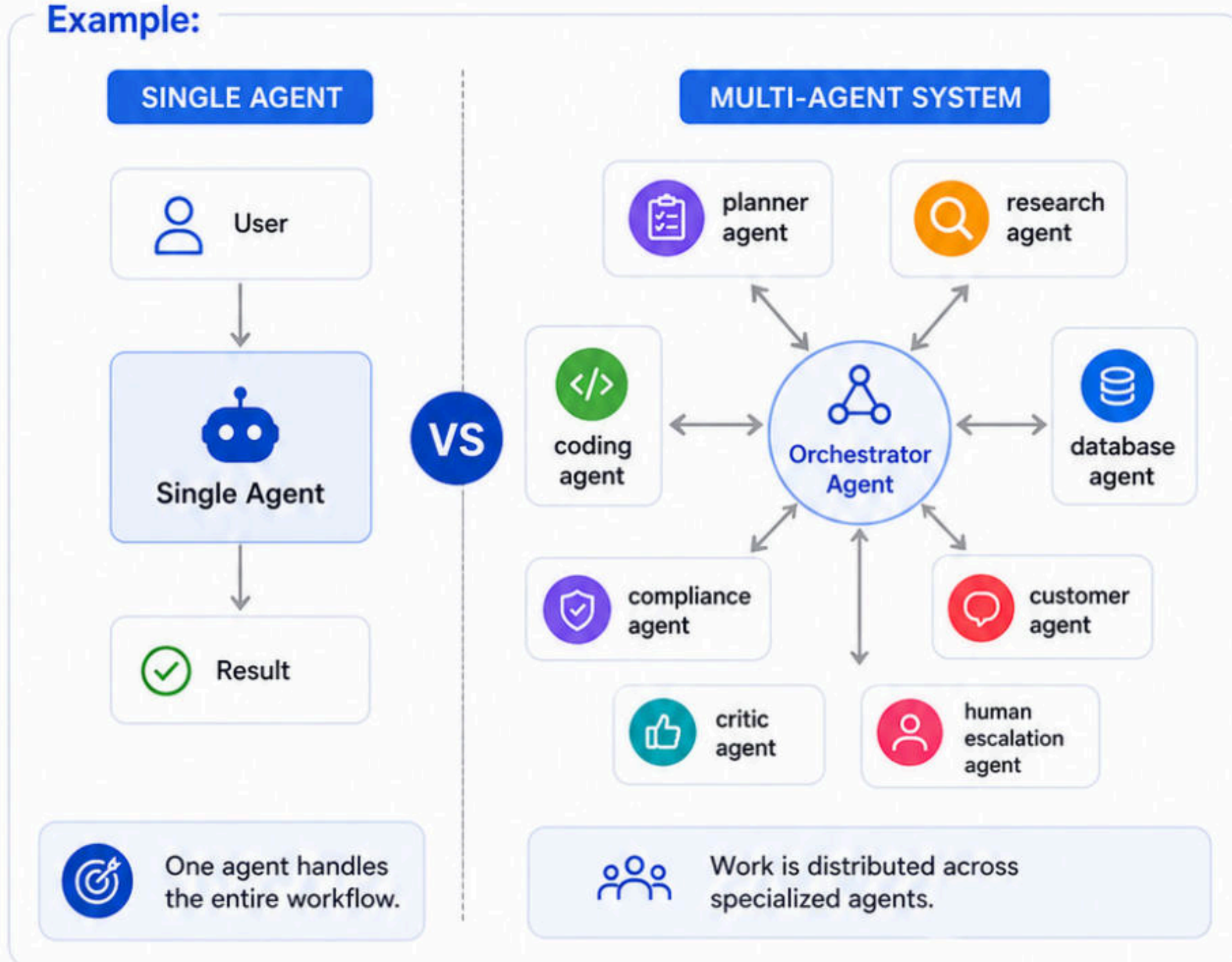


Both matter because
agentic systems are becoming more distributed.

Layer 7: Multi-Agent Orchestration

-  A single agent can do **one workflow**.
-  A multi-agent system can **divide work across specialists**.

Example:



But more agents do not automatically mean better intelligence.

Sometimes they **only create more chaos**.

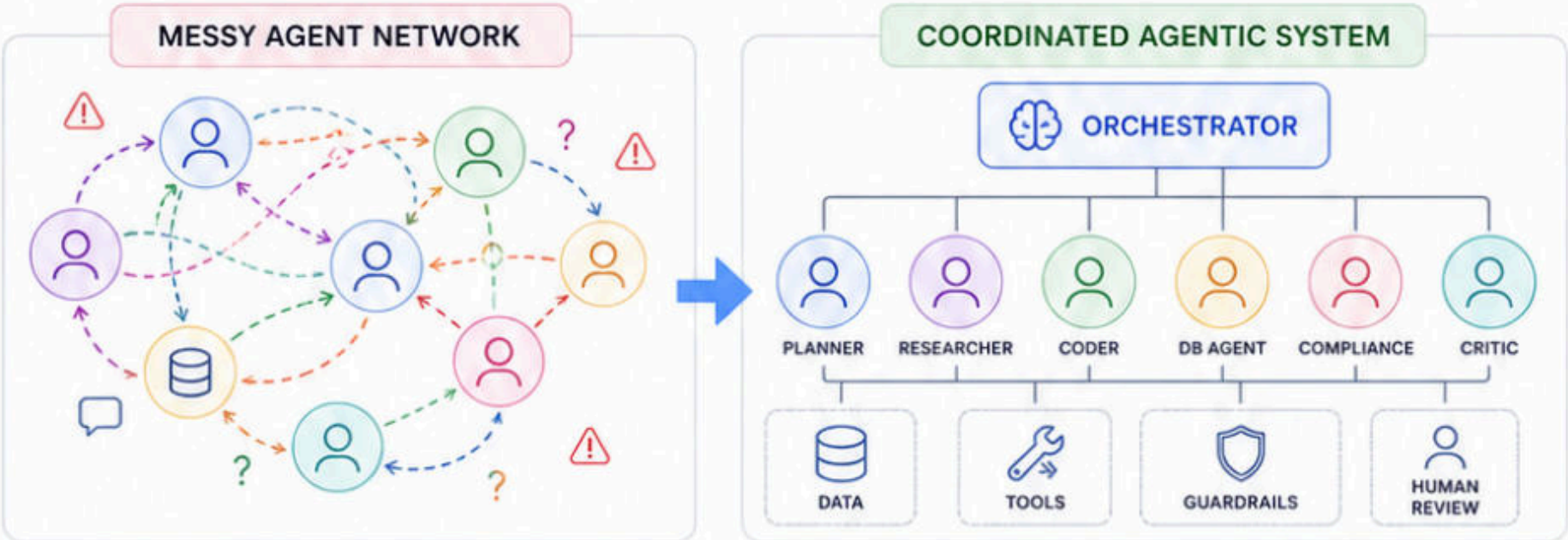


Save this slide for quick revision!

Multi-Agent Failure Modes

Multi-agent systems fail when:











- 01 agents duplicate work
- 02 agents disagree silently
- 03 one agent trusts bad context
- 04 handoffs lose information
- 05 no one owns the final answer
- 06 tool calls conflict
- 07 cost grows unpredictably
- 08 traces become unreadable
- 09 loops never stop
- 10 humans cannot intervene



 **The harder part is not creating agents. The harder part is coordinating them.**

Multi-Agent Design Checklist

Before adding another agent, ask:

- 1  Does this need a separate role?
- 2  What is this agent responsible for?
- 3  What can it access?
- 4  What can it not do?
- 5  Who reviews its output?
- 6  When does it hand off?
- 7  When does it stop?
- 8  What happens if it fails?
- 9  What does it cost?
- 10  How will we trace it?



Do not add agents because it sounds advanced.



Add agents because the workflow needs **specialization**.



Save this slide for quick revision!

Layer 8: Evals

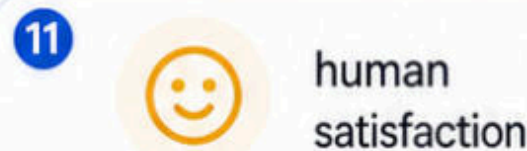
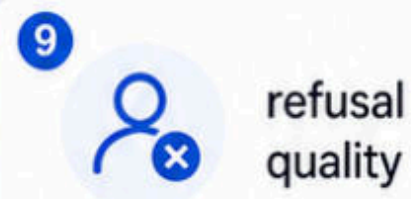
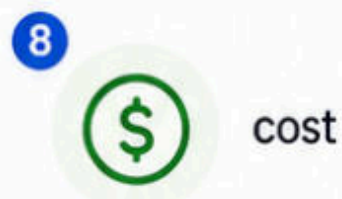
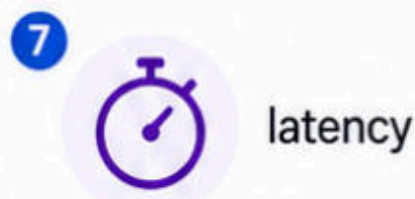
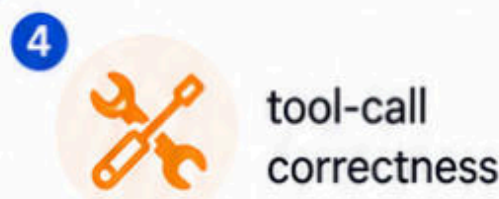
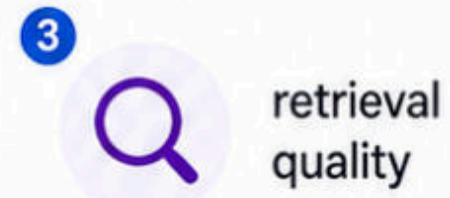
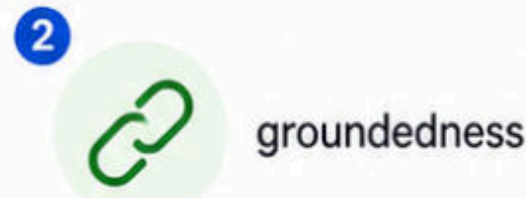
Prompting gives you
an output.



Evals give you
confidence.



In 2026, GenAI systems need **evals** for:



No evals = no production confidence.













Save this slide for quick revision!



Evals: The minimum useful suite

For a serious LLM app, build:

- 
Golden test cases
- 
Edge cases
- 
Adversarial prompts
- 
Retrieval evals
- 
Tool-call evals
- 
Hallucination checks
- 
Safety checks
- 
Regression tests
- 
Cost benchmarks
- 
Human review samples



**Your eval suite is your
GenAI test suite.**



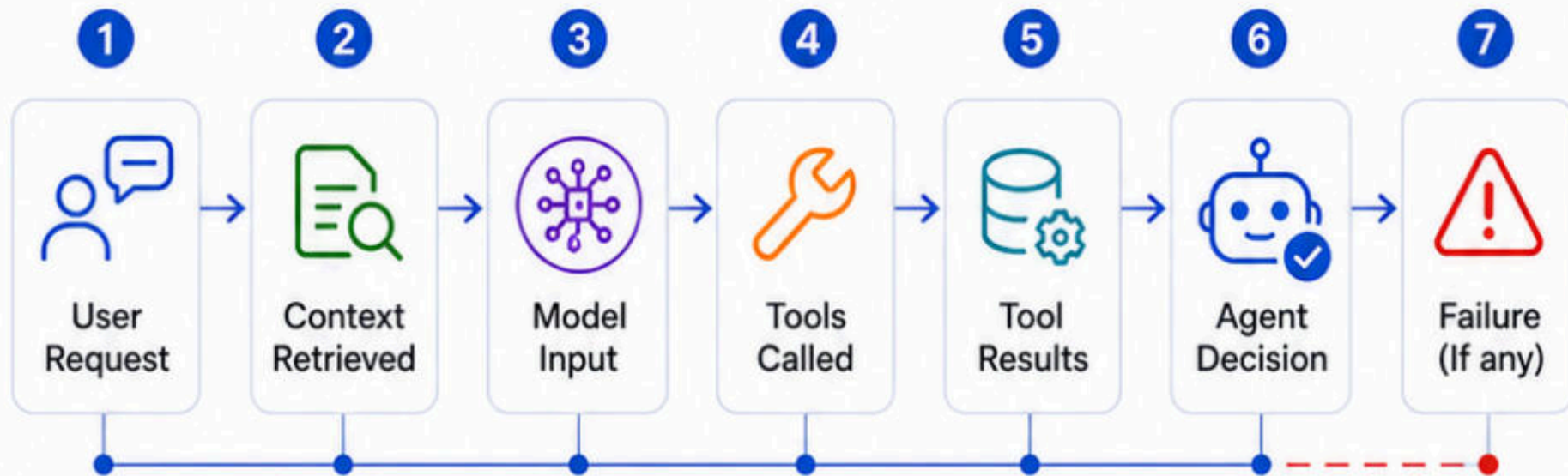
**Would you ship code without tests?
Then don't ship prompts without evals.**



Save this slide for quick revision!

Layer 9: Tracing

When an agent fails, you need to see the path.



Tracing answers:




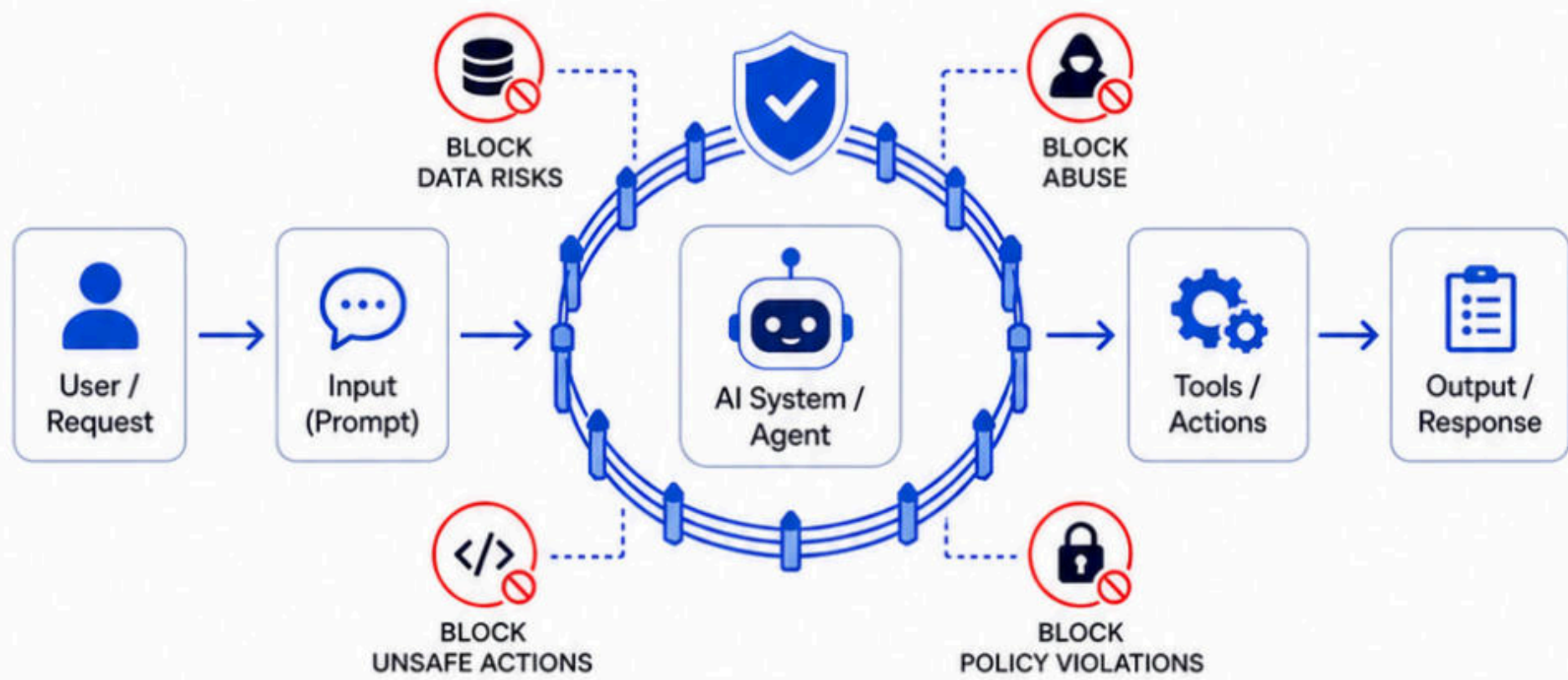
No traces = no debugging.









Save this slide for quick revision!

Layer 10: Guardrails

 Guardrails define what the system should **never** do.



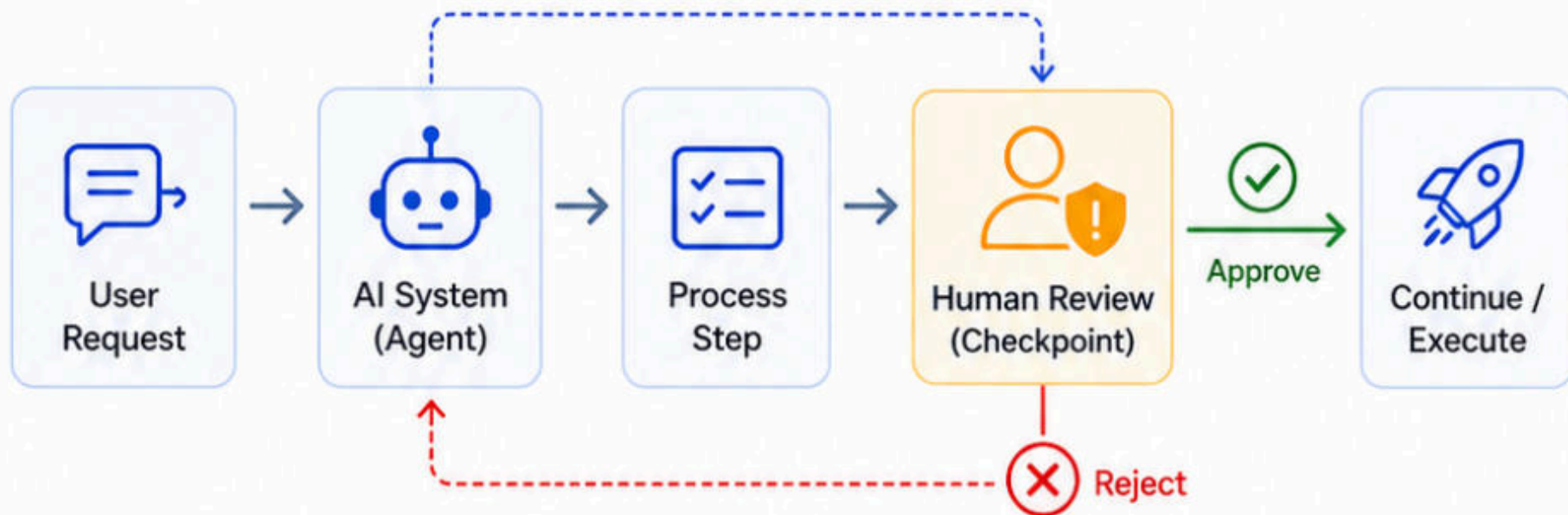
You need guardrails around:

- 1  data access
- 2  tool permissions
- 3  prompt injection
- 4  unsafe actions
- 5  financial actions
- 6  medical/legal claims
- 7  Customer communication
- 8  private data
- 9  compliance workflows
- 10  autonomous decisions

 **The more autonomy you give an agent, the more guardrails you need.**

Layer 11: Human-in-the-loop

Humans are not a weakness in agentic systems.
They are a control layer.



Use human review for:

1



high-risk
decisions

2



uncertain
outputs

3



irreversible
actions

4



financial
approvals

5



legal/compliance
steps

6



customer-
sensitive
messages

7



low-confidence
answers

8



policy
violations

9



escalation
cases



The best systems do not remove humans.
They involve humans at the right moment.



Save this slide for quick revision!

Layer 12: AgentOps



DevOps

DevOps operates **software**.



MLOps

MLOps operates **models**.



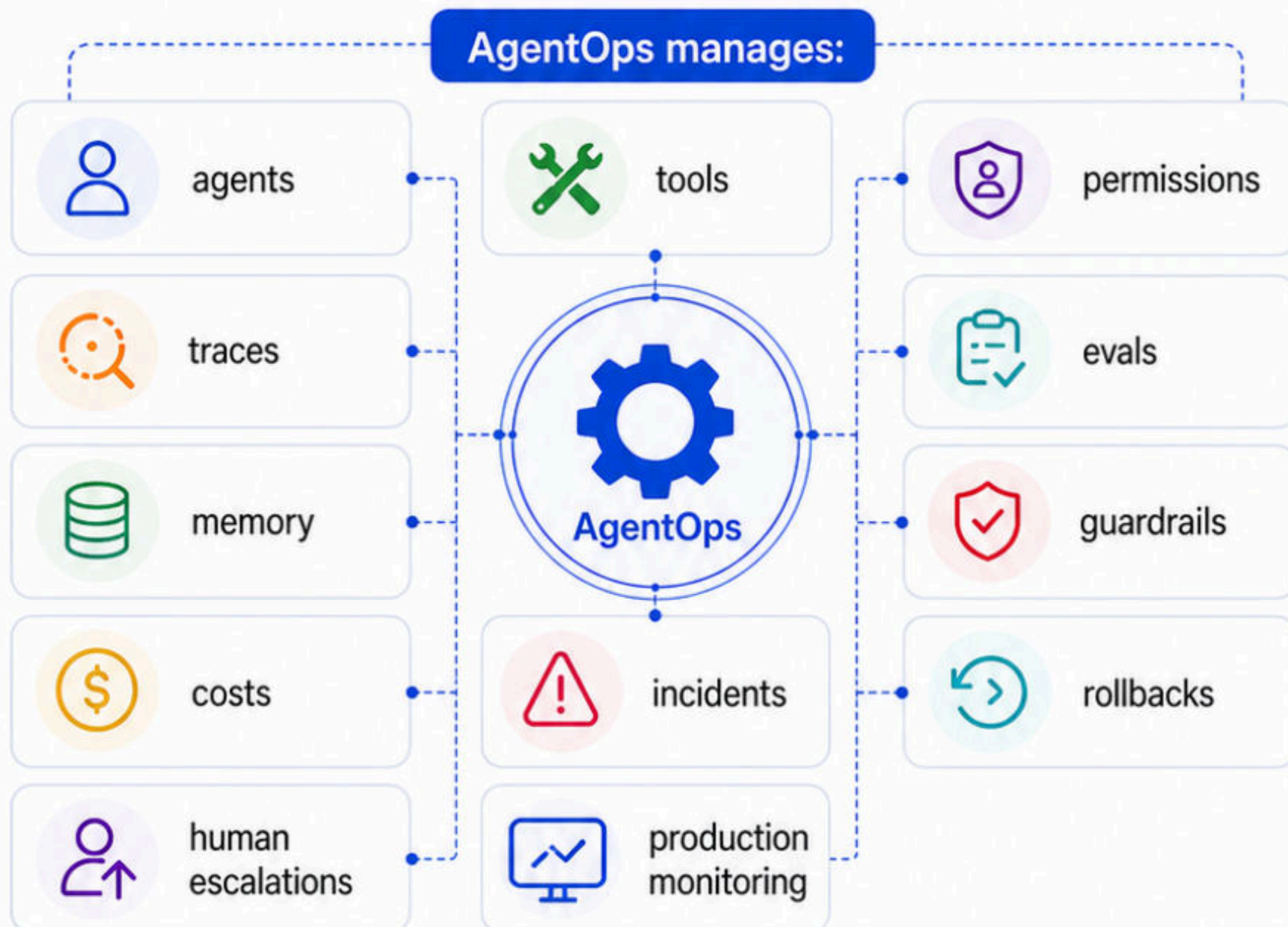
LLMOps

LLMOps operates **LLM apps**.



AgentOps

AgentOps operates **autonomous workflows**.









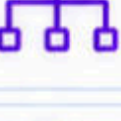
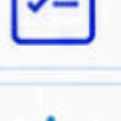




AgentOps is where GenAI becomes real production infrastructure.



Save this slide for quick revision!

The final map

The 2026 GenAI Stack:

-  **Model routing** | chooses the right model.
-  **Context engineering** | gives it the right information.
-  **RAG 2.0** | grounds it in trusted knowledge.
-  **Tool use** | lets it act.
-  **MCP** | connects it to systems.
-  **A2A** | connects it to other agents.
-  **Orchestration** | coordinates workflows.
-  **Evals** | test quality.
-  **Tracing** | explains failures.
-  **Guardrails** | control risk.
-  **HITL** | adds judgment.
-  **AgentOps** | runs it in production.



Save this slide for your next build!



The 2026 AI Engineer Mindset

OLD MINDSET



Do not just ask:

Can the model
answer?

- ✗ Focus on one answer
- ✗ Optimizes for model output
- ✗ Works in demos

NEW MINDSET



Ask:

Can the system complete
the task **reliably, safely,**
repeatedly, and
affordably?

- ✓ Focus on end-to-end outcomes
- ✓ Optimizes for system reliability
- ✓ Works in production

That is the difference between:

demo AI



- ✗ Looks impressive
- ✗ Brittle and inconsistent
- ✗ Fails in real-world edge cases
- ✗ Hard to trust

and

production AI



- ✓ Solves real problems
- ✓ Reliable and consistent
- ✓ Handles edge cases
- ✓ Trusted and scalable



You don't just build AI.
You build **reliable AI systems.**



Save this slide for quick revision!



Save this



The future belongs to people who can turn LLMs into **reliable systems**.

- ⊗ Not just impressive prompts.
- ⊗ Not just beautiful demos.
- ⊗ Not just another chatbot.

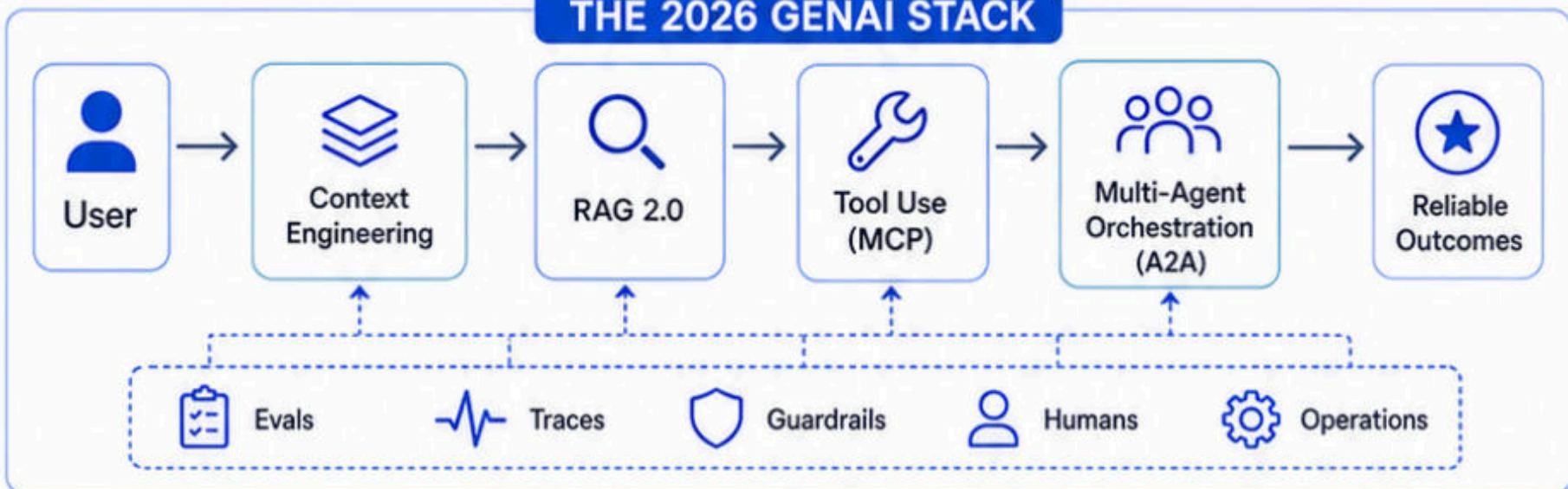


Reliable systems.

With:



THE 2026 GENAI STACK



That is the **2026 GenAI stack**.



Save this for your future self!