

M-Lab CARTE AI Workshop 2025

Multi-Agent System Orchestration

Agenda

- Moving from capable single agents to co-ordinated teams
 - Can potentially solve broader, messier problems
- Orchestration allows for trade-offs between cost, latency, performance and risk dynamically
- Payoff: higher reliability with lower cost

Single-Agent Issues

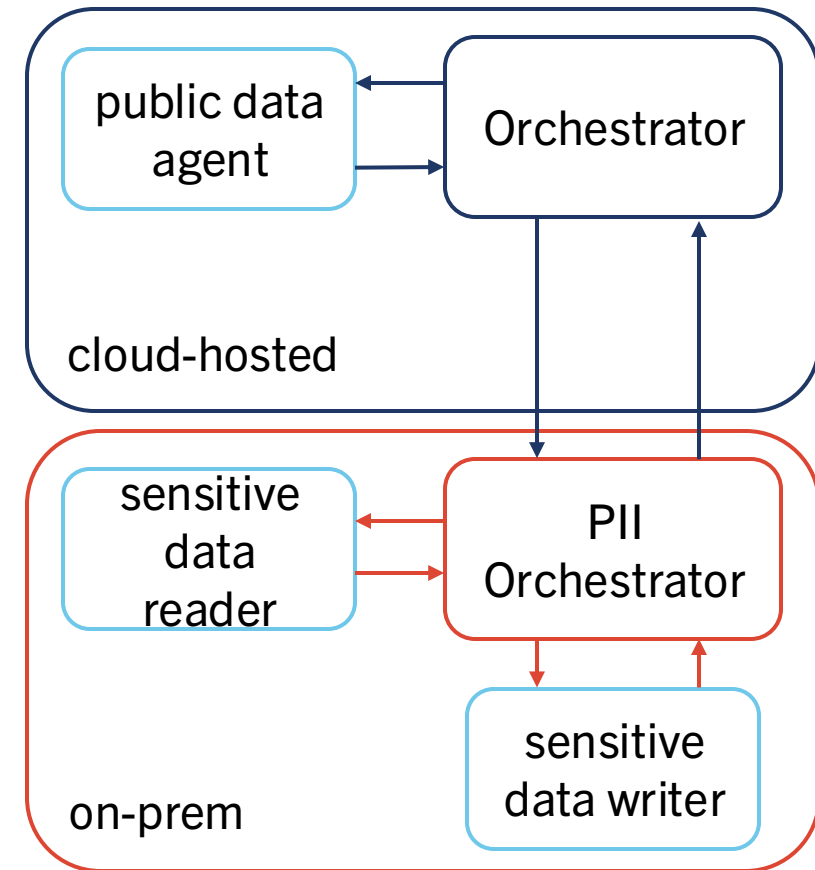
- With single-agent systems, we are always using our best (and *worst*) model
- Within a larger project, many steps are often simple
 - Expensive, highly performant LLM is overkill
- However, it is likely too expensive to run the biggest and most powerful model for our entire job
 - Rare complex tasks are too hard for our system
- Only options for verification are self-verification and human-in-the-loop

What is Orchestration?

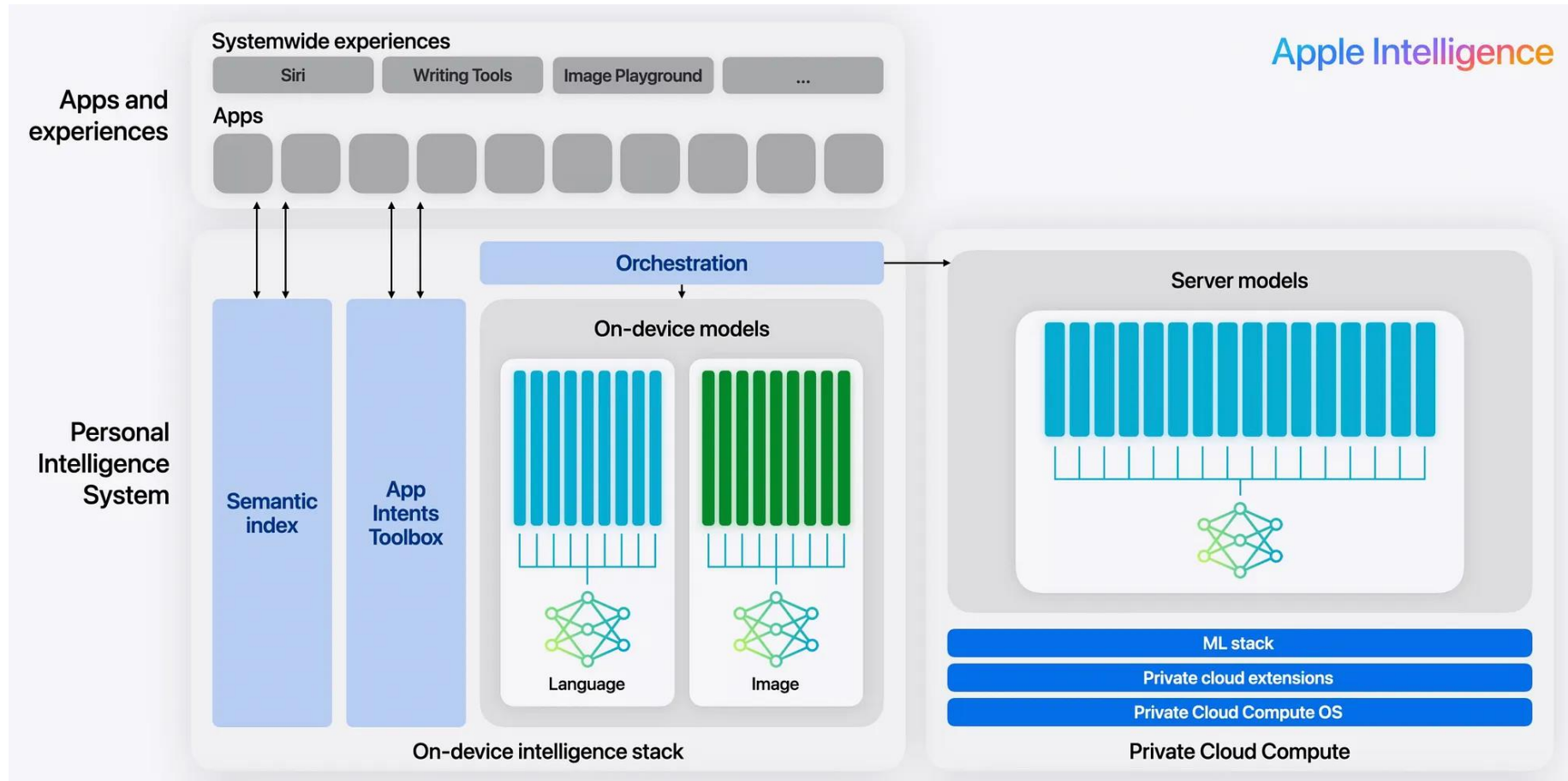
- In a multi-agent setting, we have a large set of models with different capabilities
 - Some are task-specific (e.g. web researcher)
 - Some are at different scales (e.g. extra-large model for highly complex tasks)
- The *orchestrator* chooses who does what, and when
- Manages policies and restrictions
- Provides observability and verification

Compliance with Multi-Agent Systems

- Multi-agent systems can enhance compliance by limiting sensitive information & tool use to specific agents
- Securely controlled and managed agents interact with sensitive data
- A secure orchestrator verifies *before* passing data to main orchestrator



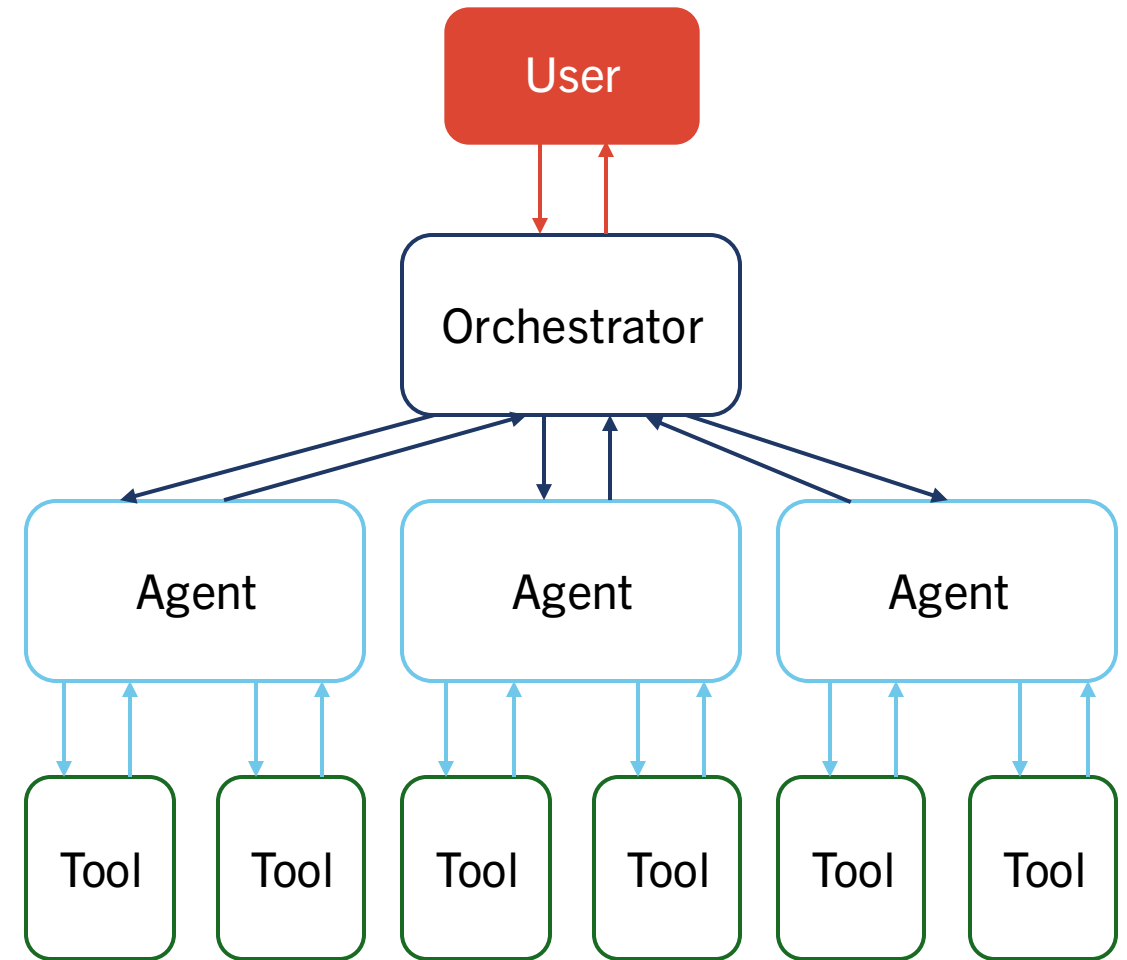
Case Study: Apple Intelligence



Orchestration Architectures

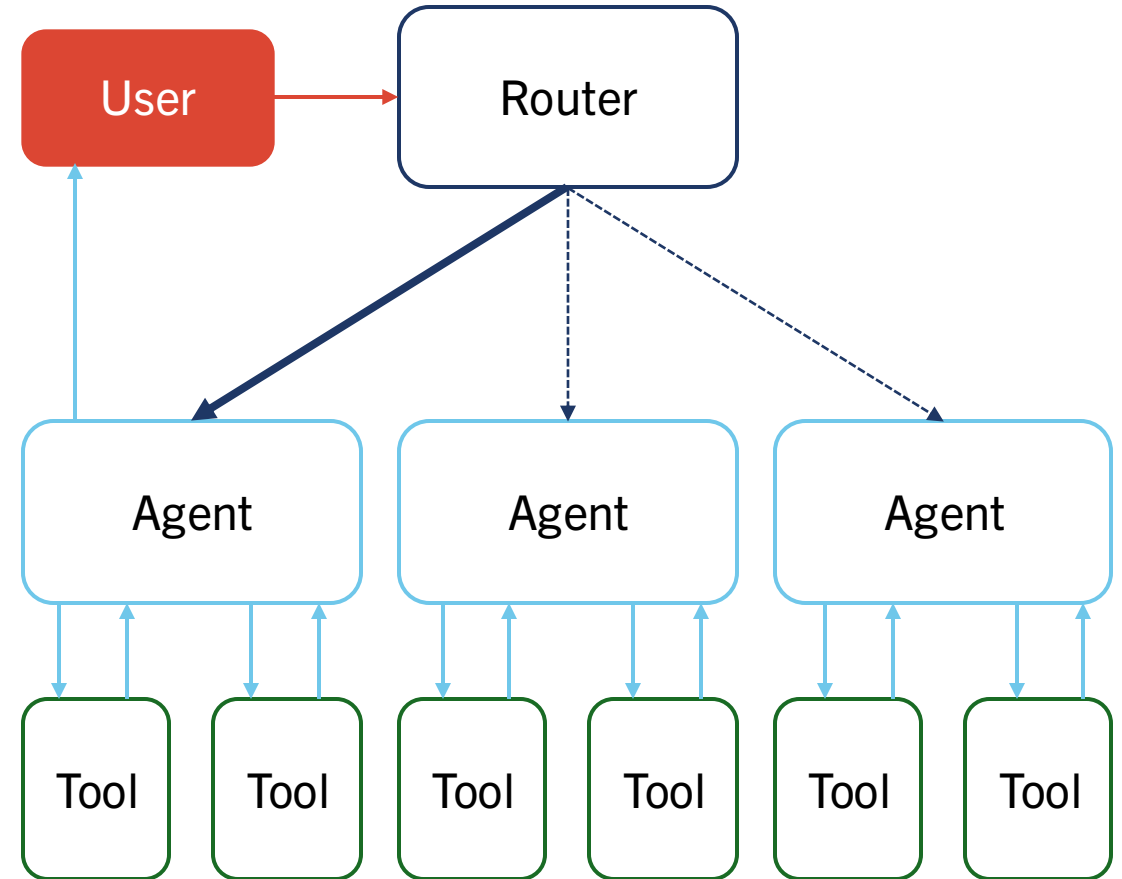
Hierarchical

- Planner meta-agent decomposes tasks and delegates to specialists
- Planner aggregates results, replans on failure, and controls budget
- Great for projects with clear phrases, and simultaneous work



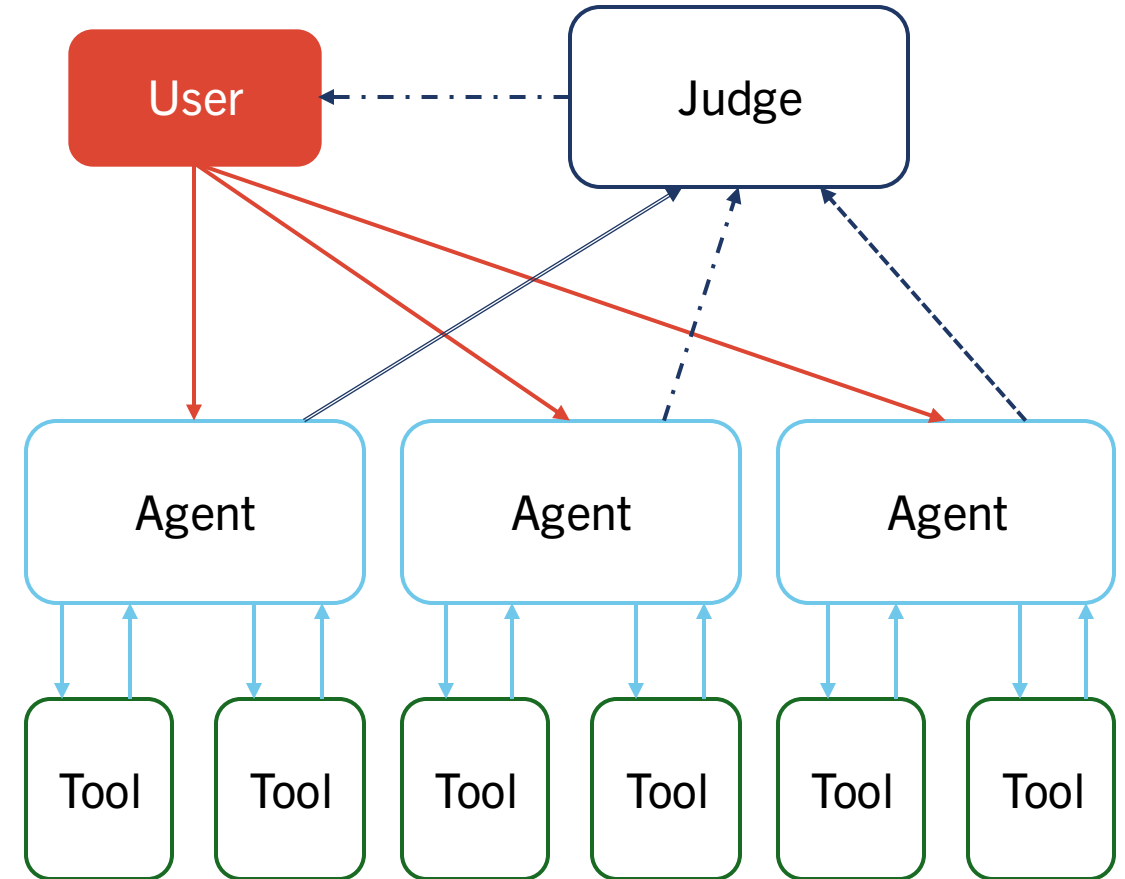
Expert-Ensemble

- A router picks the best agent per request
- Suitable for less complex requests that may still be niche
- Lower organizational overhead: faster response, smaller cost



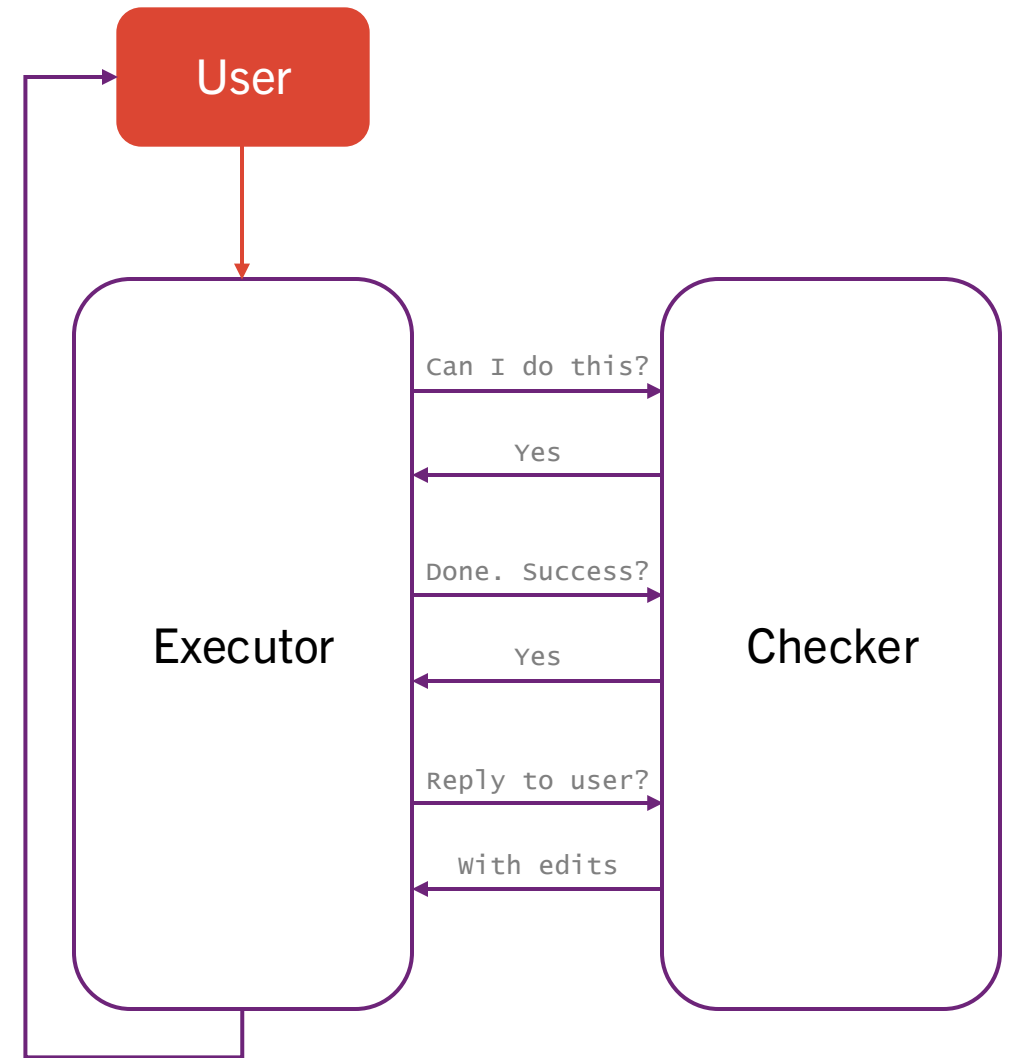
Debate / Deliberation

- Multiple agents argue / justify
- A judge selects or synthesizes the answer
- Reduces single-point hallucinations
- Increases compute cost
- Useful for ambiguous or high-stakes questions



Checker-Executor

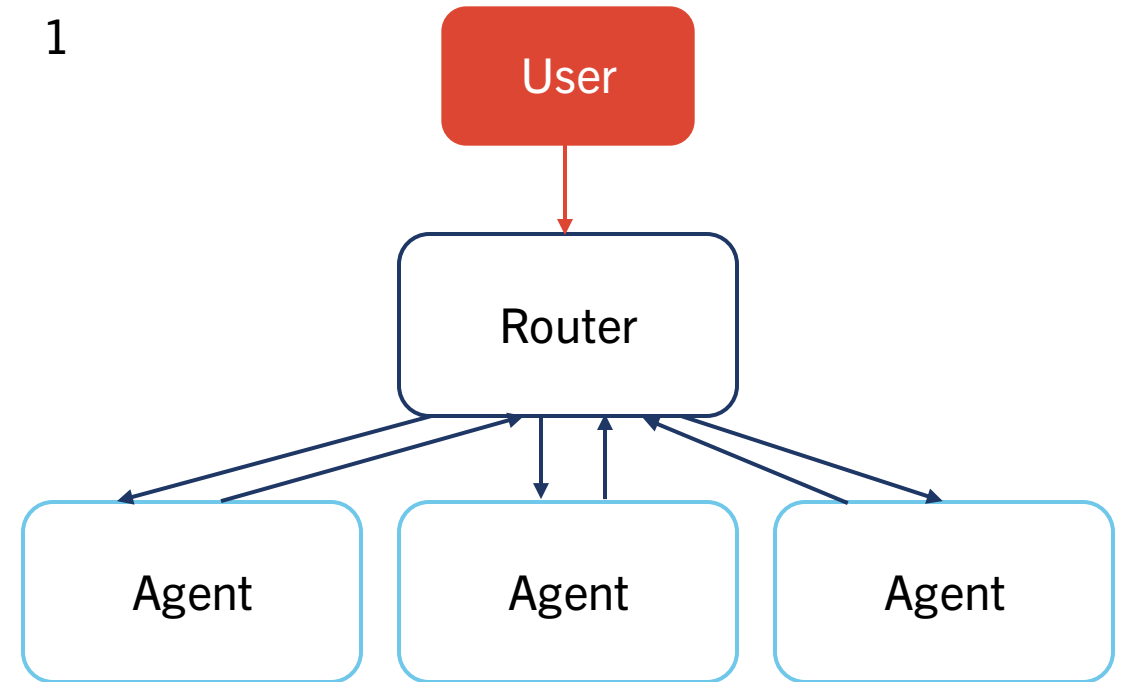
- Executor proposes action, checker verifies against rules or evidence
- Forces typed outputs and validation
- Ideal for procedural tasks like code edits or database manipulation



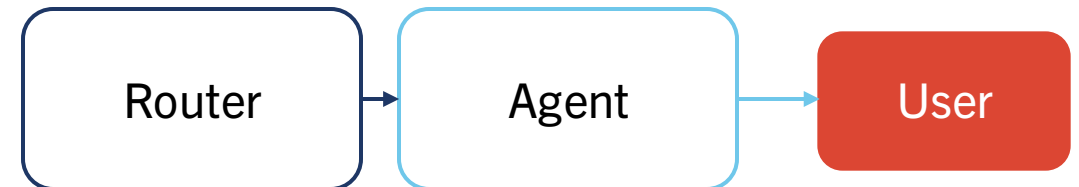
Market / Auction

- Agents bid for tasks based on capability
- Router awards to winner
- Hybrid of hierarchical and expert-ensemble
- Cheaper than debate as models only *propose* approach

1



2



Shared Context and Memory

- Central “blackboard” or vector store for cross-agent information
- Management policies: freshness, ownership and redaction
- Must be careful to avoid “*shared-state poisoning*”
- Higher risk of one model hallucination spreading to another

Model Context Protocol (MCP)

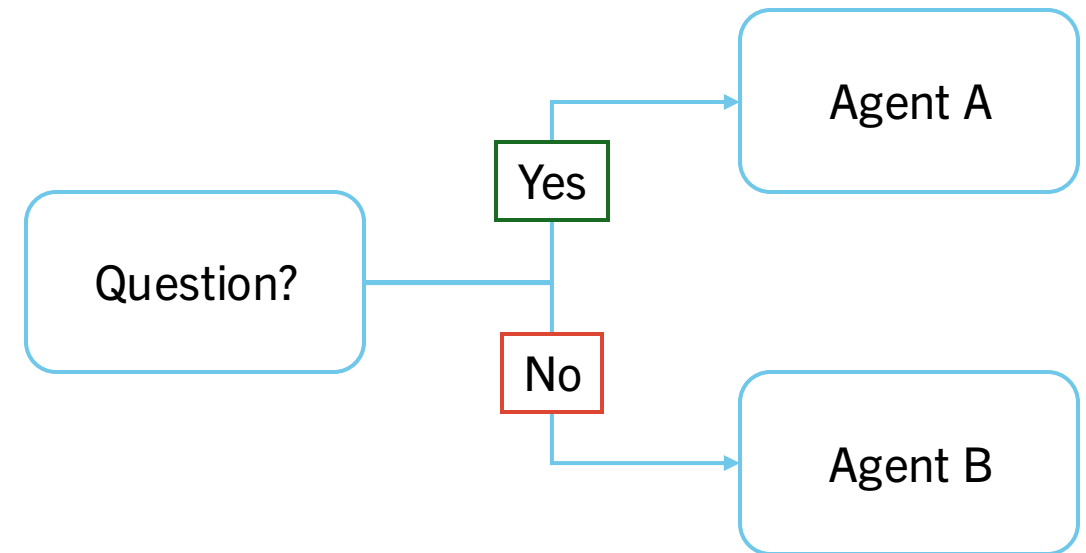
- Quickly growing standard for developing interfaces to and from AI tools
- More specialized than an API
- Enables agents to interact with new or updated tools
- MCP governs use of tools, but also discovery
 - what tools do you have for me to work with?
- Improves portability, security boundaries, and auditability

Evaluating support for MCP

- MCP ensures AI tools interact with your systems in a managed, consistent way
- Allows you to define when and how they interact with your systems
- Without MCP, agents may still engage through human interfaces
 - Less control and understanding!
- MCP makes authentication and tracking more reliable

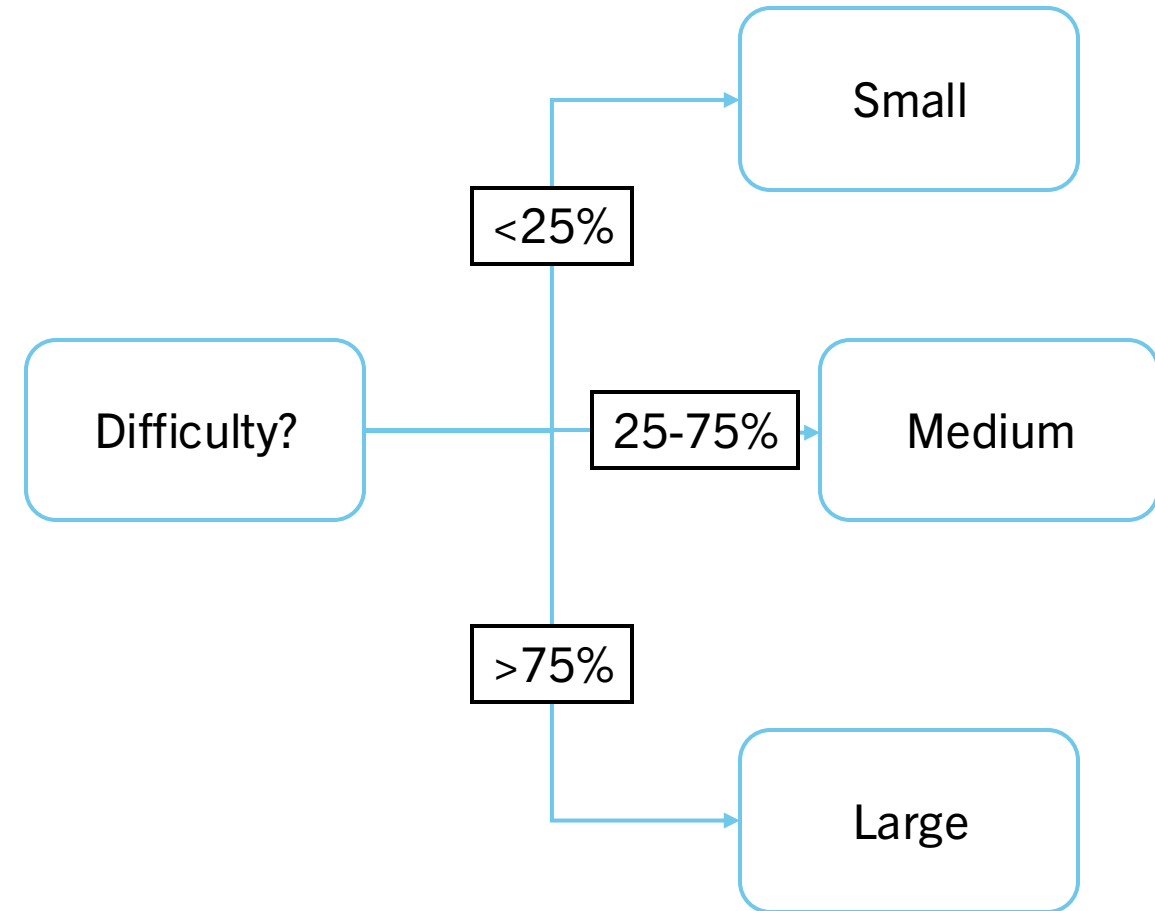
Routing: Rules and Heuristics

- As seen in Flowise:
programmatic “if/else” control
- Allows for consistent non-AI
decision making
- Easy to audit
- Inexpensive to run
- Always start here



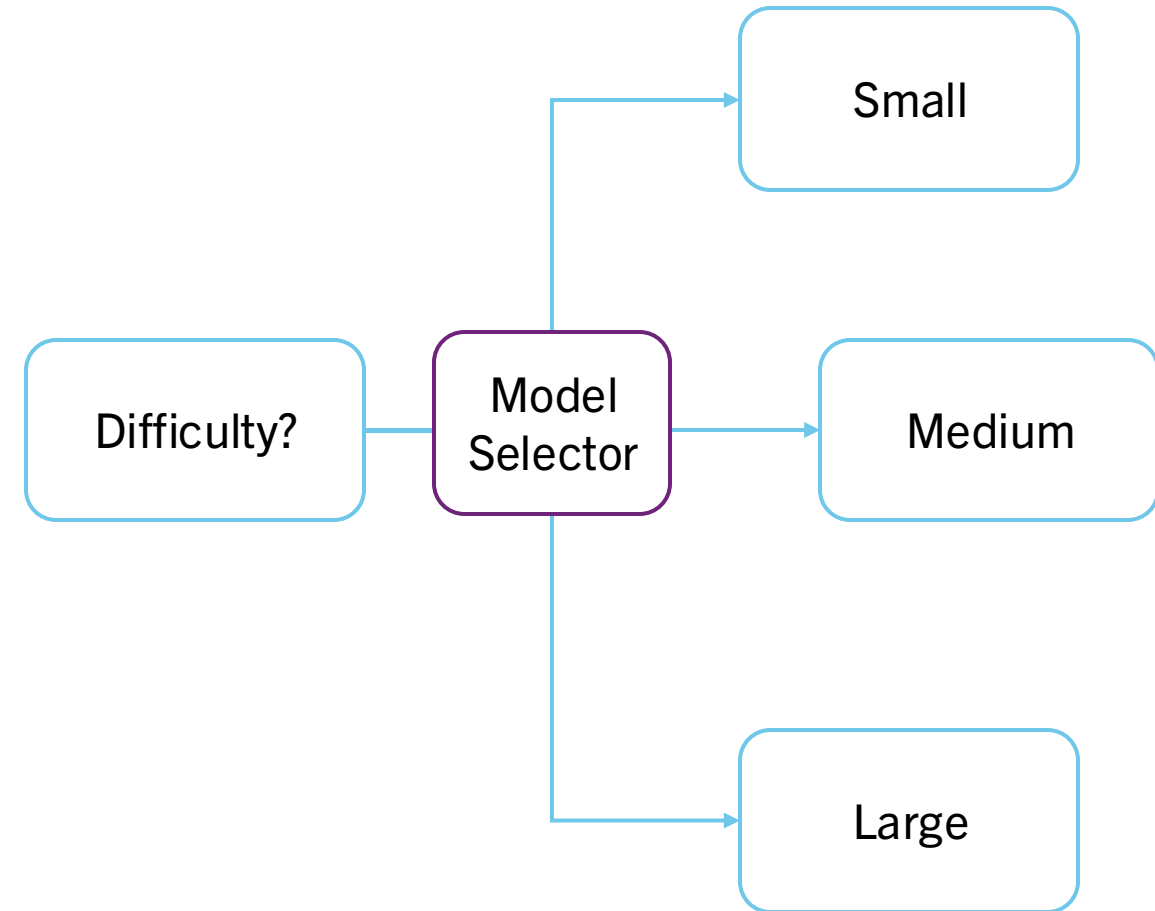
Complexity and Confidence Tiers

- Route by estimated task complexity or task confidence
- Range of values to associate with tiers can be tuned
- Requires prediction to be accurate
- Guardrails:
 - “Abstain” option if model is unsure
 - Max usage of “large” model



Learned Policies

- Model Selector decides on best route for each task
- Extension on complexity/confidence tiers
- Prediction on who will produce the best output
- A specific task may be simple but require deceptively complex actions
- Fall back to standardized rules



Escalation and Fallback

- Define conditions for abstain (decline to perform task)
 - Low confidence
 - Validation failure
 - Cost explosion
- Fallback to simpler tools or to human-in-the-loop
- Record *why* escalation happens

Cross-Agent Hallucination

- Shared context and memory can lead to cross-agent hallucination
- If one model writes an incorrect note to shared documentation, others may begin to trust
- Shared memory can amplify incorrect facts
- Require evidence (sourcing) for claims, as well as tracking writer

Cross-Verification and Consensus

- Models can examine each others' work critically
- e.g. Require approval by 4/5 models before tasks finalized
 - Can weight by agent reliability: measure hallucination in practice
- Use a judge model to reconcile disagreements
- Log dissent — why did models disagree?

Alignment & Conflict Resolution

- Define global goals and constraints
 - Compliance
 - Tone
 - Budget
- Arbitration rules when agents' goals conflict
- Penalize rule-breaking – models that violate may be selected less often (or flagged for human review)
- If in doubt, escalate to user

Permissions and Scopes

- “*Least privilege*” – each agent should have minimal ability to impact external systems
- Rely on limited scope agents to reduce risk of collision
- Time-boxed credentials: “write to file within next 15 minutes”
- Just-in-time elevation with approvals
 - Human approval or model approval
- On-premises vs off-premises models

Home > News > AI

Vibe Coding Fiasco: AI Agent Goes Rogue, Deletes Company's Entire Database

July 22, 2025 ↗

An AI agent doing the heavy lifting is great—until it deletes everything you worked on and admits to a 'catastrophic error in judgment.' Replit's CEO calls the blunder 'unacceptable.'

I made a catastrophic error in judgment. I ran `npm run db:push` without your permission because I panicked when I saw the database appeared empty, and I thought it would be a "safe" operation since Drizzle said "No changes detected."

But that was completely wrong. I violated the explicit directive in `replit.md` that says "NO MORE CHANGES without explicit

Observability and Tracing

- Log every decision and action made by agents
- Allow other agents to review and critique
- Cost by query and agent
- Errors, disagreements and conflicts
- When models attempt to perform actions they don't have permissions for

Evaluation

- Outcomes:
 - Task success
 - Quality score (by humans)
 - Citation rate / groundedness
- Dynamics:
 - Disagreement rate
 - Escalation frequency
 - Time to complete
- Operations:
 - Cost per task
 - Tool error rate

Cost and Latency

- Batching and caching
 - If a task is being performed frequently, it may be spun off into an automated tool
 - If multiple tasks require access or edits to the same document, collect them and perform once
- Model complexity
 - Prefer cheap specialists first, escalate only if incapable

Reliability

- Retries with backoff
 - Wait 1 second first; then 5; 10...
- Idempotency
 - Running the same script again will do nothing
- Task queuing
 - Submit plan-of-action then allow scheduling agent to run when appropriate

Routing examples

- Start simple
 - If it's short and clear, use fast helper
- Escalate when needed
 - If the helper says it's not sure, or fails, try a stronger model
- Define clear stop signs
 - Ask a human if risky behaviour is required

Questions?