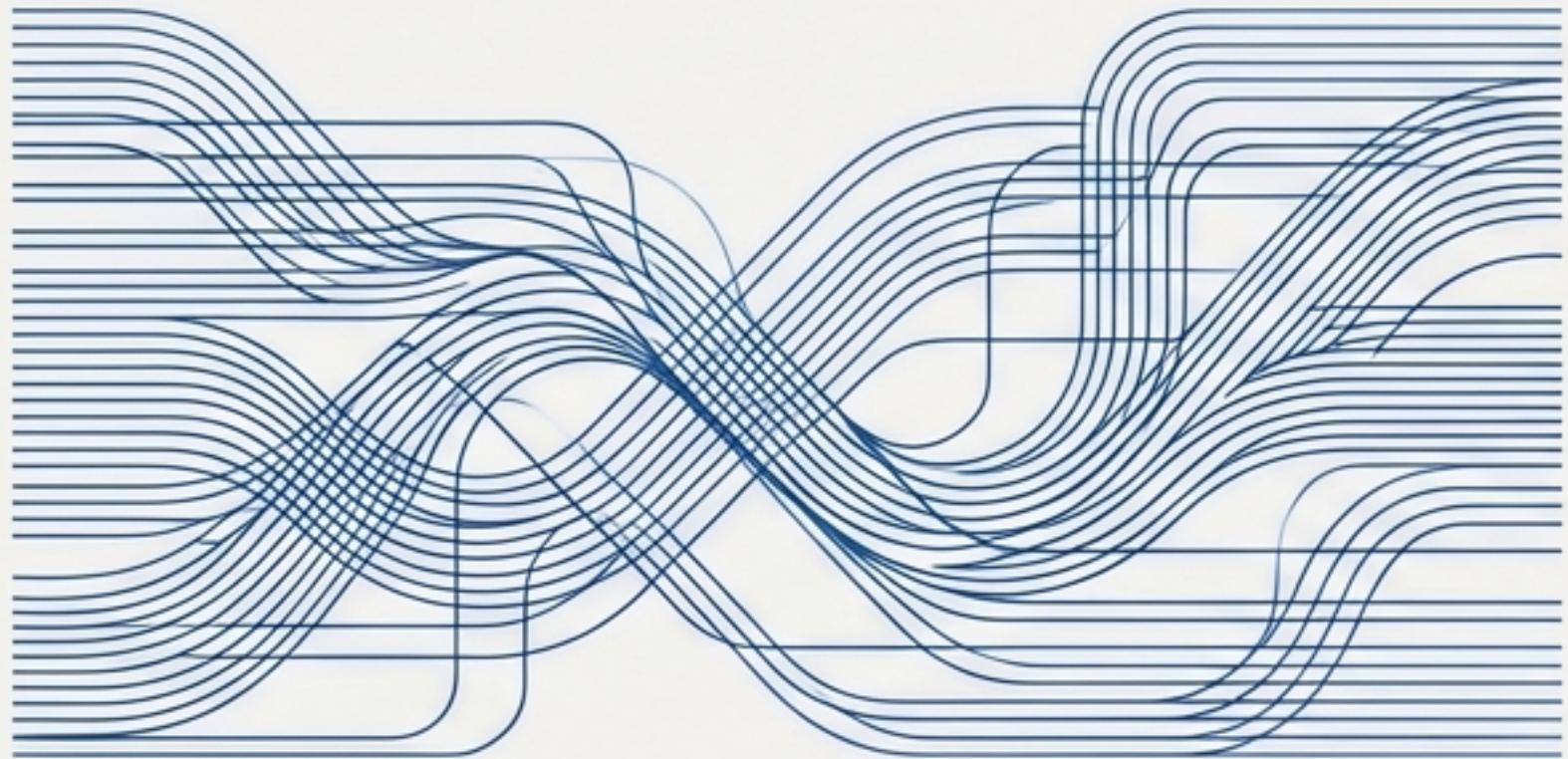


Agents on Rails

Structured Paths for Safe and
Explainable GenAI Autonomy

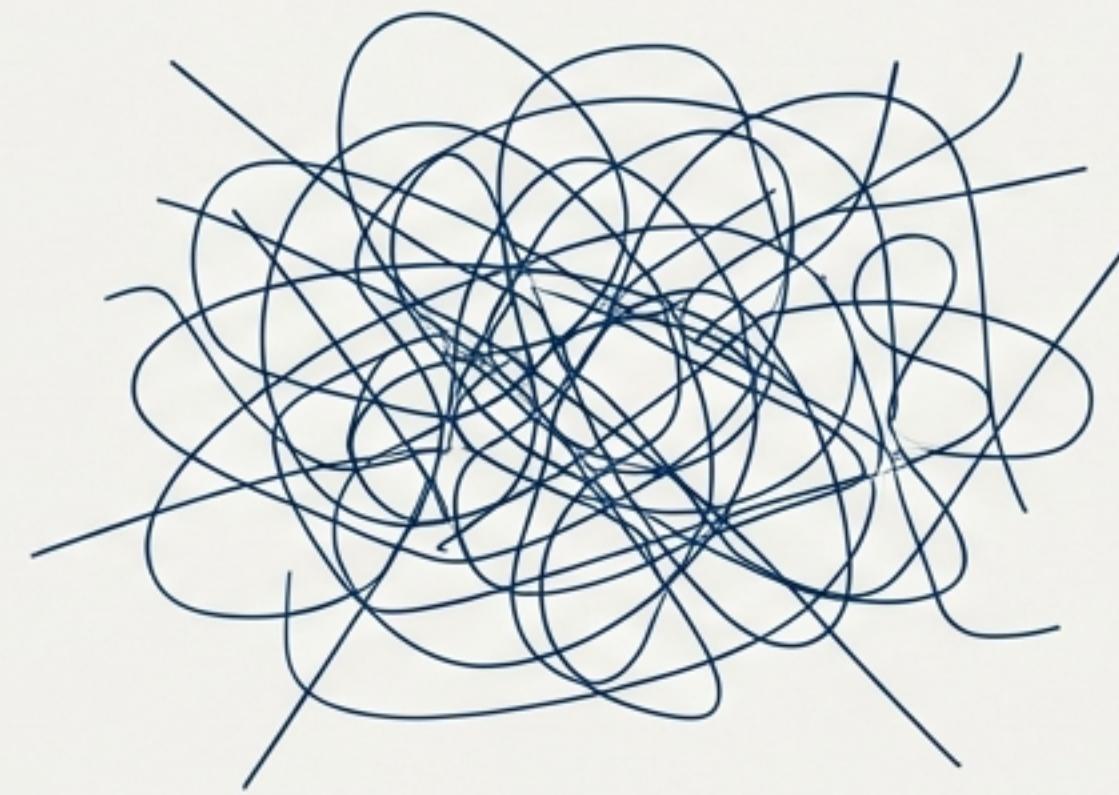


The Promise vs. The Peril of GenAI Agents



The Promise

Generative AI agents hold great promise for automating complex tasks by making independent decisions.

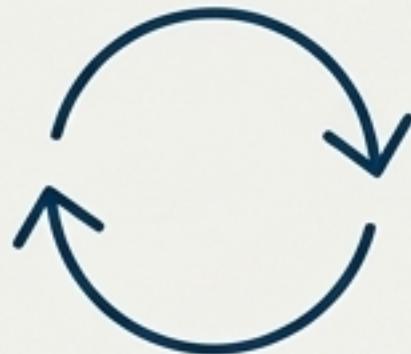


The Peril

But when given free rein, even powerful LLMs can go ‘off the rails’—producing irrelevant or harmful outputs, or veering away from intended goals.

We are at a critical juncture: we must choose structured capability over unbounded chaos to realise the promise of AI.

Why Free-Form Agents Are “Radically Less Reliable”



Lack of Focus & Looping

Agents frequently meander or loop infinitely, revisiting the same step instead of making forward progress. LLM errors compound over iterations, leading them astray from the goal.



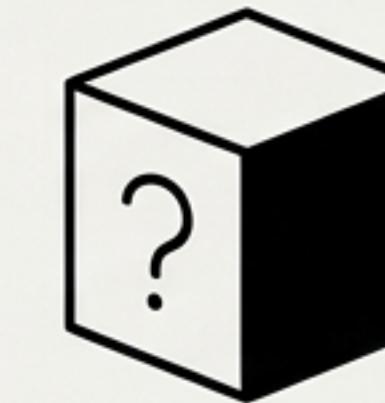
Security Vulnerabilities

Open-ended agents are highly vulnerable to prompt injection, where malicious instructions hidden in inputs can trick the agent into executing unintended actions.



Unpredictable Outputs

The probabilistic nature of LLMs means decisions and outputs differ from run to run, undermining reliability and making it nearly impossible to guarantee outcomes or reproduce results.



Lack of Explainability

Decisions are opaque. We are left asking, ‘Why did the AI do that?’ with no clear audit trail to trace the decision back to inputs or intermediate reasoning.

“After an initial burst of hype, most users abandoned these agents... due to their inability to ‘reliably [complete] multi-step reasoning’ tasks.”

The Way Forward: From Free-Roaming Car to High-Speed Train



Unbounded Agent:

High flexibility, but unpredictable, unpredictable, unsafe, and easily lost.

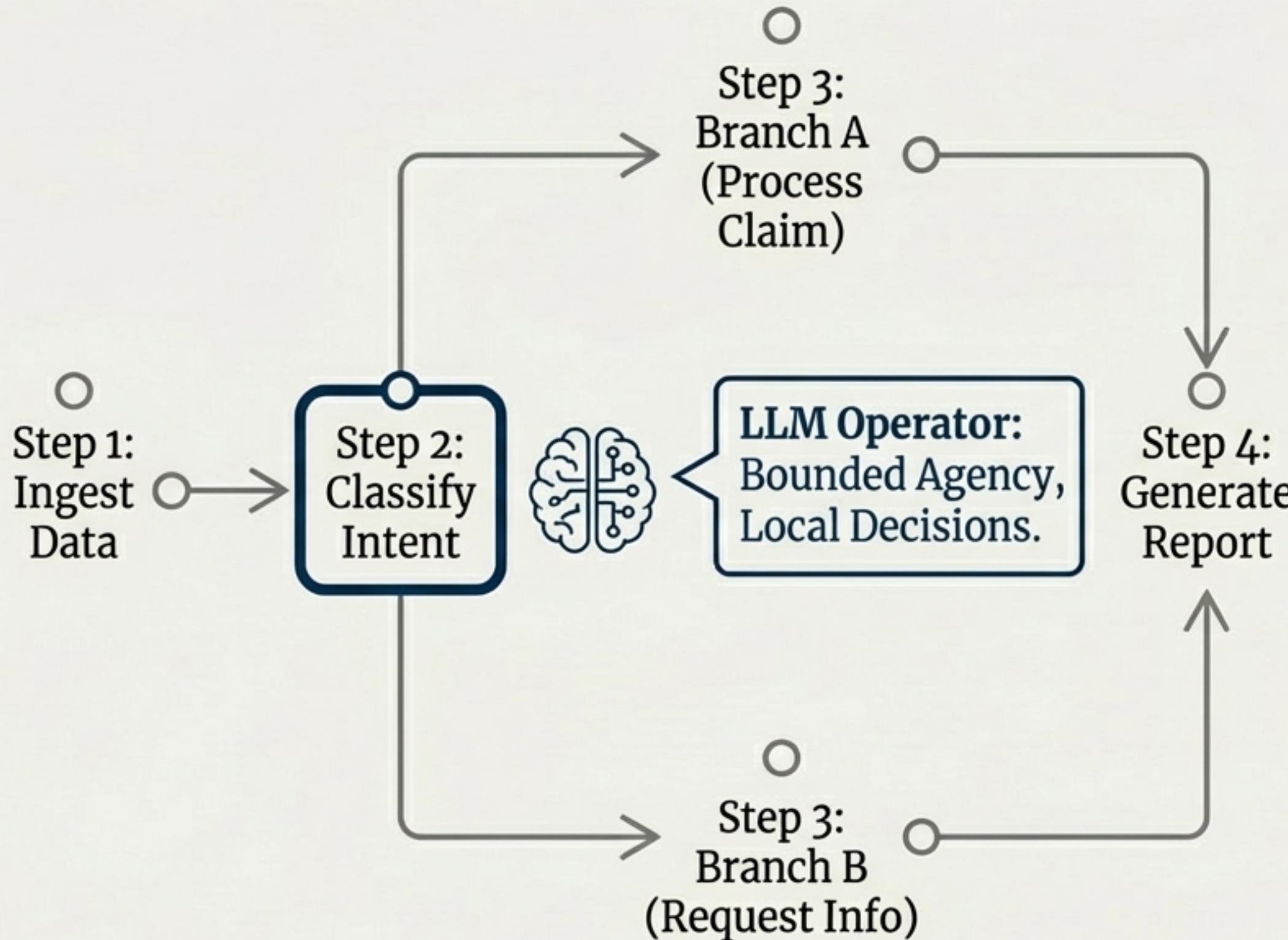


Agent on Rails:

High reliability, efficiency, and safety on a predetermined, version-controlled route.

“AI Automation is like a train on a track... highly efficient and reliable for its predetermined route.” – Rajat Jain

Combining LLM Intelligence with Deterministic Control



The Track (Deterministic Workflow)

The overall sequence of actions is fixed in a schema (e.g., a DAG). All possible branches are enumerated. The business logic lives in this schema, not in the LLM.

The Conductor (LLM Operator)

The LLM is invoked at specific steps to perform bounded, well-defined tasks (e.g., extract, classify, summarise).

The LLM is a 'scripted actor' for the macro-level flow, not the director.

Pillar 1: Enforce Structure with Strongly Typed Schemas

Instead of free-form text, every LLM interaction must adhere to a predefined, structured format (e.g., JSON). The LLM's job is to “populate specific fields,” not just generate prose.

Schema Contract

```
{  
  "type": "object",  
  "properties": {  
    "action": { "type": "string" },  
    "data": { "type": "number" }  
  },  
  "required": ["action", "data"]  
}
```

Validated Output

```
{  
  "action": "classify_intent",  
  "data": 12345  
}
```



Consistency & Predictability: Greatly reduces output variability from the LLM.

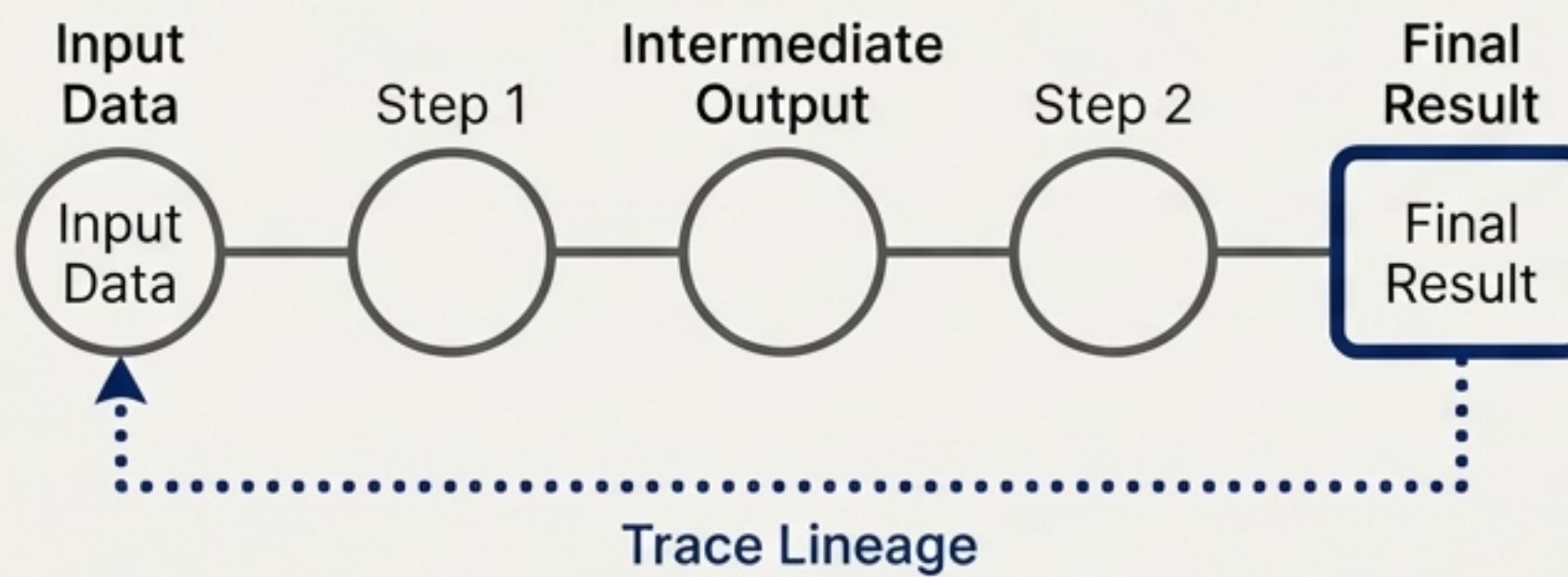
Automatic Validation: The system can immediately detect, retry, or reject malformed outputs.

Seamless Integration: Downstream systems can reliably consume well-formed data.

We treat the LLM like a microservice that must adhere to a strict API contract.

Pillar 2: Achieve Explainability with Knowledge Graph Provenance

Log every LLM call, input, output, and decision into a structured repository. This creates a complete, queryable trace of the agent's journey, turning an opaque process into a transparent one.

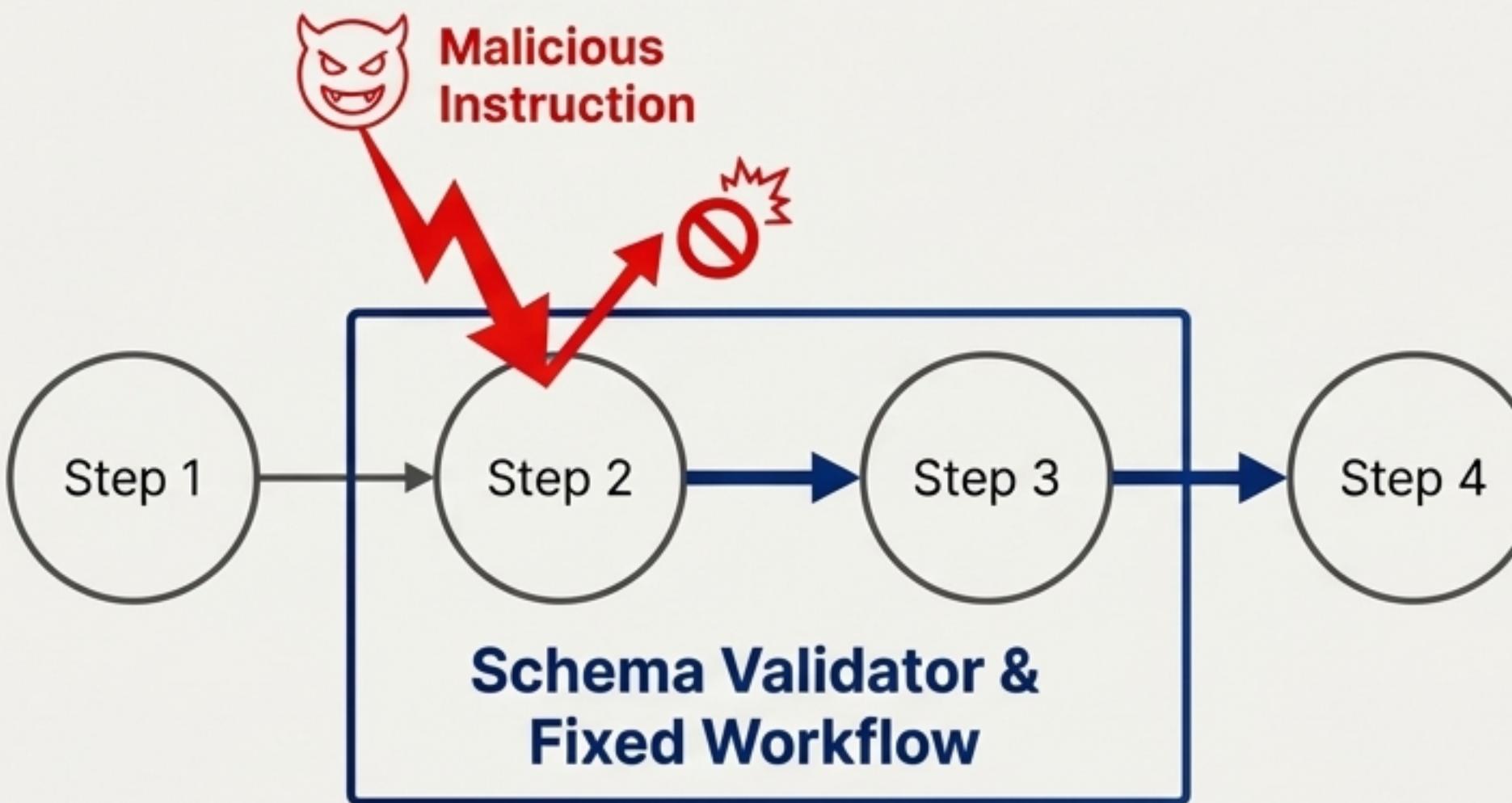


- **Full Transparency:** Answer 'Why did the AI do that?' with logged, auditable evidence. The chain of outputs forms a 'provenance trail'.
- **Rapid Debugging:** Pinpoint the exact step, inputs, and outputs where an error occurred. Enables unit testing of each step.
- **Audit & Compliance:** Inherently produces the detailed audit trails required for regulated industries.

We turn an opaque process into an open, auditable one, building trust by making every decision accountable.

Pillar 3: Ensure Safety with a Reduced Blast Radius

The agent cannot perform any action not explicitly defined in the workflow schema. The impact of any error or malicious input is confined to the current, sandboxed step.



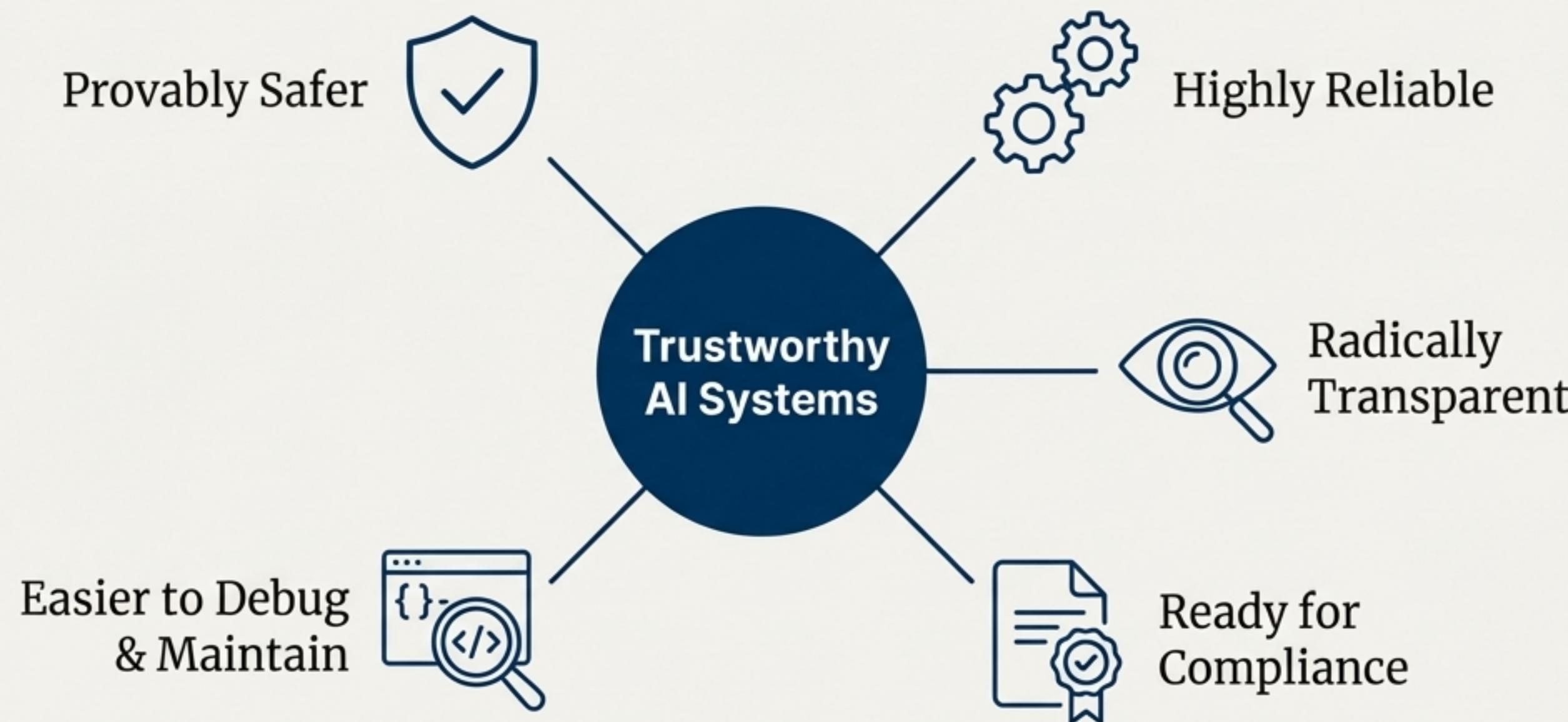
- **Mitigates Prompt Injection:** Malicious commands either fail schema validation or have ‘no track’ to run on. The request goes nowhere.
- **Controlled Tool Use:** Tool calls are predefined, permissioned, and auditable within the workflow.
- **Macro-Level Determinism:** The overall process is predictable and can be version-controlled like code.

“Bottom line: keep the core deterministic, let the LLM do scoped extraction, and enforce schemas plus tests.”

From Chaos to Control: A Systematic Solution

The Problem	The “Agents on Rails” Solution	The Outcome
Lack of Focus & Looping	Deterministic Workflow Schema	Guided, goal-oriented progress
Unpredictable Outputs	Pillar 1: Strongly Typed Schemas	Consistent & verifiable results
Lack of Explainability	Pillar 2: Knowledge Graph Provenance	Transparent & auditable decisions
Security Vulnerabilities	Pillar 3: Reduced Blast Radius	Contained risk & attack mitigation

The Business Impact: Delivering Innovation with Confidence



This architecture allows businesses to leverage advanced AI capabilities while meeting the rigorous standards of enterprise-grade software.

The Future of Reliable AI is Hybrid

We unlock powerful capabilities without sacrificing oversight by combining deterministic structure with probabilistic intelligence. This is a form of *guided independence*.

Acknowledging Trade-offs:

This approach requires upfront design and knowledge of the task structure.

It is best suited for known workflows, not pure open-ended exploration.

“We hope this paper serves as a blueprint for building AI systems that are not only advanced in capability but also safe, transparent, and grounded in the deterministic logic of their creators.”

