

# Another RAND AI effort?

---

- RAND's Center on AI, Security, and Technology (CAST) Management desires to accelerate the pace at which it can produce high quality research products
- LLMs are able to automate and accelerate research workflows. RAND has yet to integrate AI and realize these gains
- There is a significant amount of effort in this space
  - Google has a Co-Scientist system and an “empirical software” optimizer
  - Sakana has an AI Scientist (with 2 published blind peer-reviewed papers)
  - All major models can perform web-based “deep research”
- Our goal is to make sure we can leverage ALL of these innovations and incorporate them into a system that is highly robust and reliable
- To facilitate this, we've started by focusing on a mechanical peer review system that we can trust
  - To confidently and responsibly use these tools, CAST (and RAND) need to “own” a reliable and quick system for evaluating the robustness of our process, findings, recommendations, etc.

# Rand-AI Reviewer

*Automating Document Review Process*

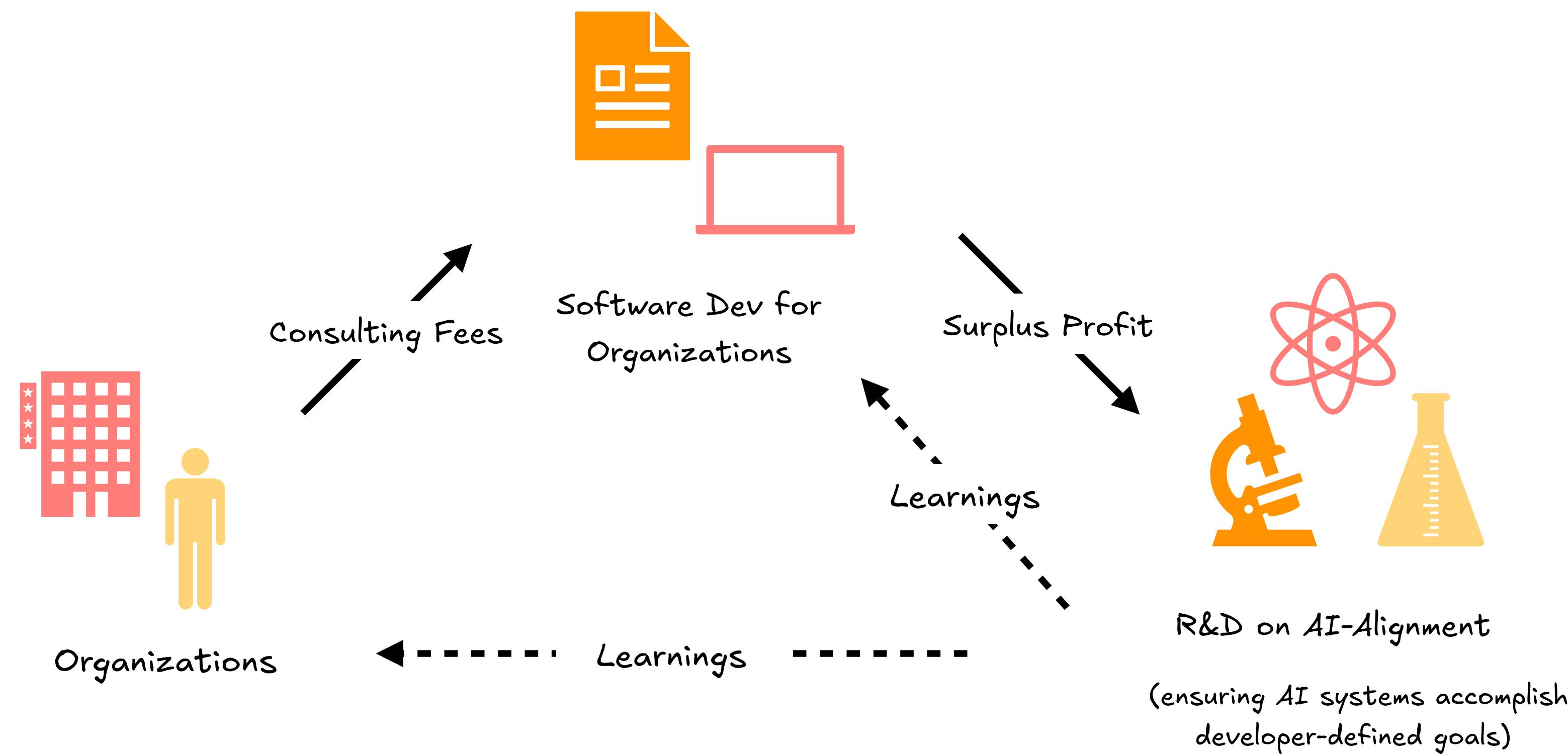
*Rand Corporation – AE Studio Collaboration – Tuesday November 4, 2025*

# Who is doing this work?



AE.STUDIO

Consulting firm focused on software development and AI alignment



# Who is doing this work?

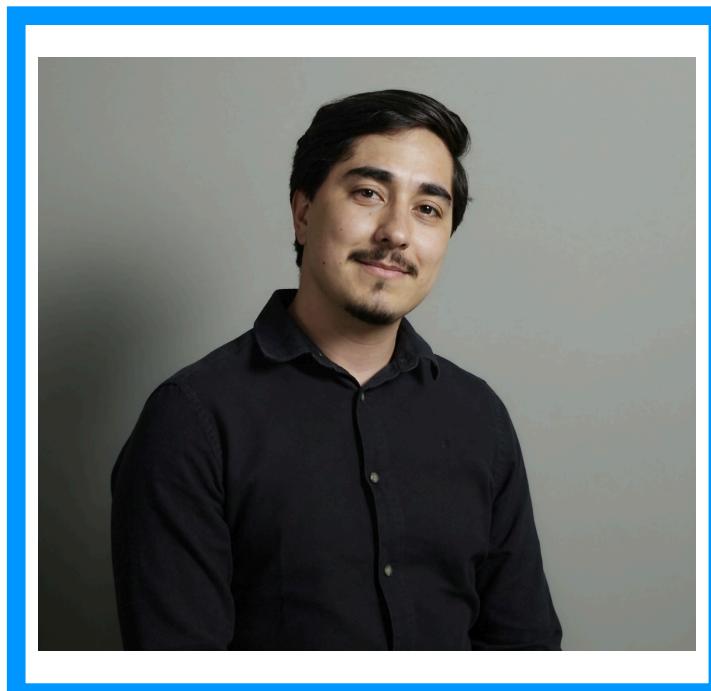
.....



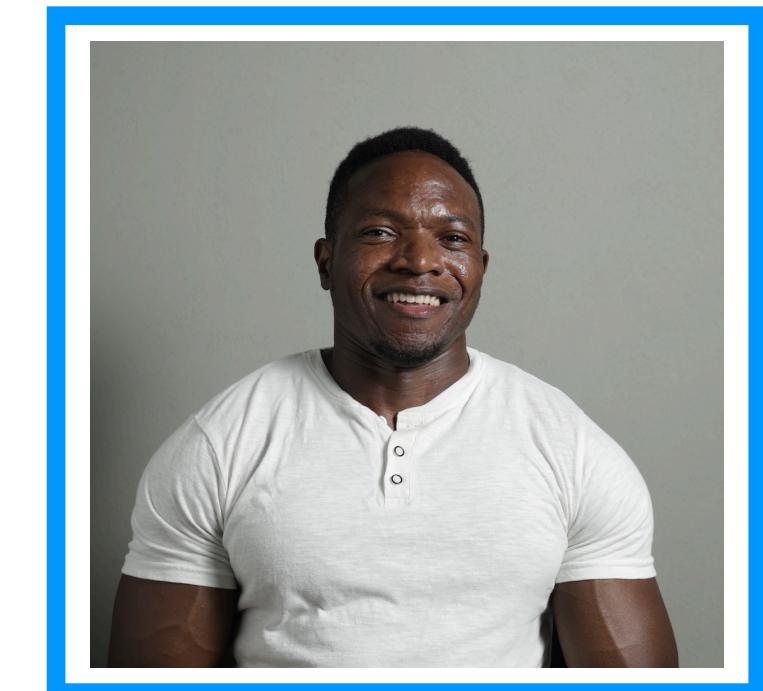
RAND-AI Reviewer Team



Carlos Bonetti  
(Technical Product Manager)  
(BS, Computer Science)



Ricardo Schieck  
(Full Cycle Engineer)  
(BS, Computer Software Engineering)



Mobolaji Williams  
(Data Scientist)  
(PhD, Physics)

# Presentation Roadmap

---

I. Motivation

II. Video Demo

III. System Details

- .i Architecture

- .ii Agents and LangGraph

- .iii Design principles

IV. Gaps and Next Steps

V. Live Demo & Q&A

# Presentation Roadmap

## I. Motivation

## II. Video Demo

## III. System Details

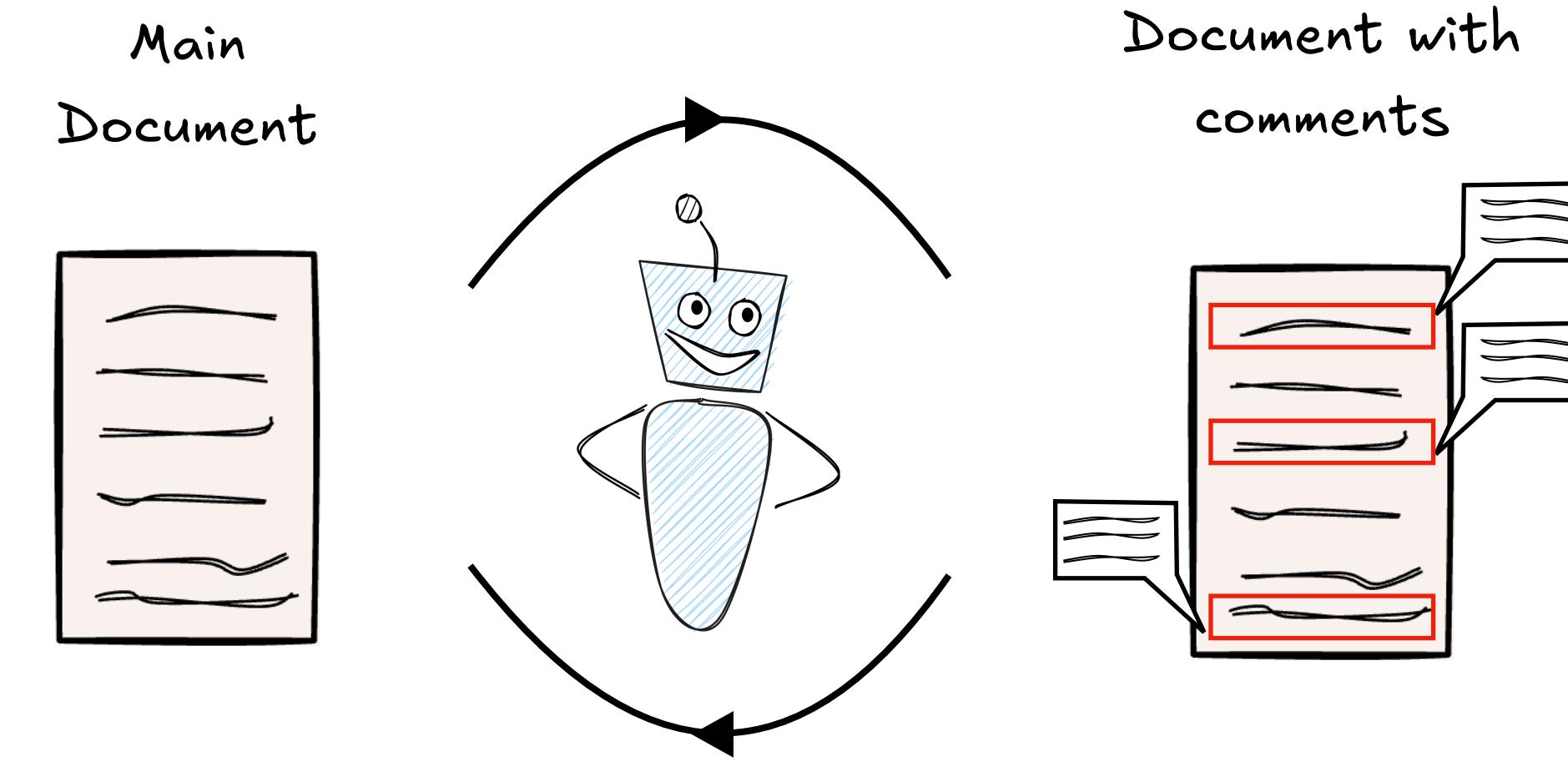
- .i Architecture

- .ii Agents and LangGraph

- .iii Design principles

## IV. Gaps and Next Steps

## V. Live Demo & Q&A



# Motivation: Embarrassments of LLM Capabilities

---

## Deloitte to Refund Government After Admitting AI-Generated Errors in \$440K Report

<https://colitco.com/deloitte-refund-ai-errors-government-report/>

### Consequences

#### Micro

Embarrassing for Deloitte  
(and financially costly;  
\$440,000 refund for services)

#### Macro

Welfare discussions were based on  
false information; Further erosion  
of belief in efficacy of LLMs

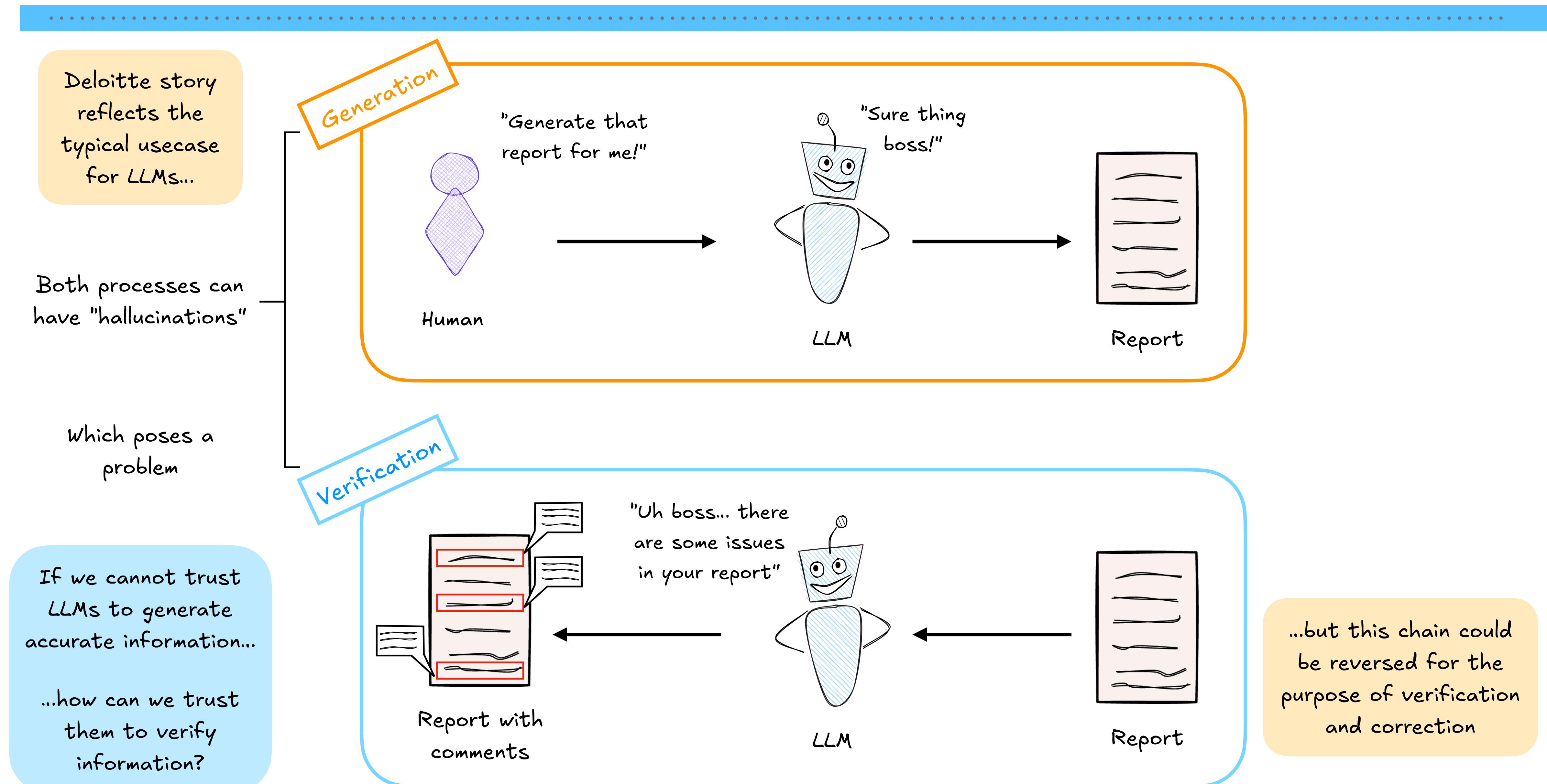
Australian govt commissioned a report  
to examine compliance issues of welfare  
recipients'



Consultants used LLMs for their reports  
which led to "false quotes", fabricated  
"experts" and "references"

**Conclusion:** LLM-generated information  
cannot be trusted without verification

# Motivation: Reversing Generation



# Motivation: LLM Systems & Peer Review

- (Motivating Question)-

How can we build an AI system  
that we can trust to conduct  
the main tasks of peer review?

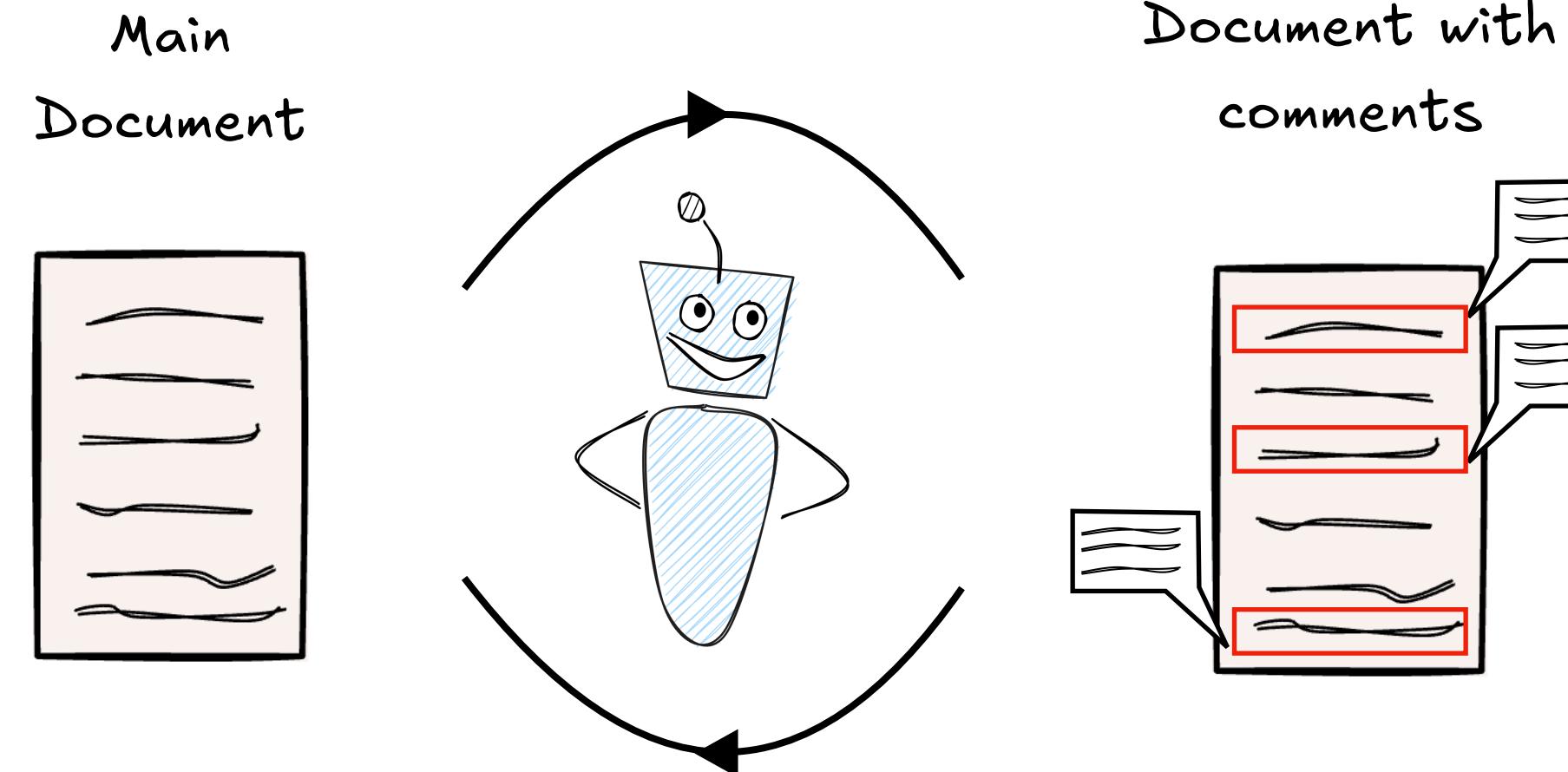
We need to narrow the problem  
because peer review can check  
for a lot of things....

Syntax, grammar,  
formatting

Existing  
efforts at  
RAND for this

Statement veracity  
and context

Overall argument  
coherence, soundness



NOTE!

There is a lot of existing work  
in this area (<https://arxiv.org/pdf/2501.10326.pdf> "LLMs for automated peer review: a survey (2025)")

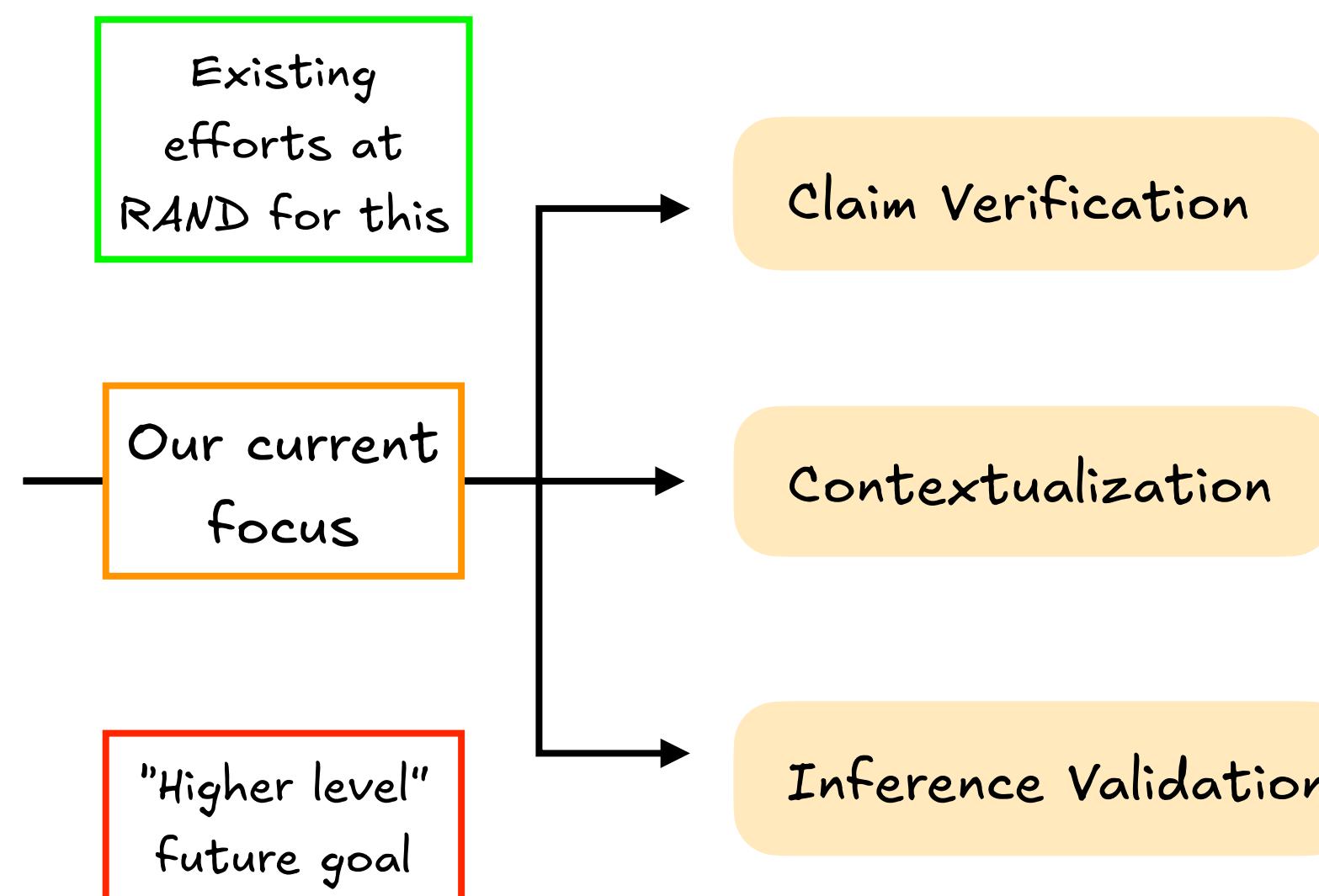
We did not want to  
add just another  
system but to...

- checking whether  
sources support claims

- checking whether there  
is additional literature  
that provides context

- checking whether  
arguments are valid

...understand the design  
principles within in any such  
trustworthy "LLM for  
peer-review" system

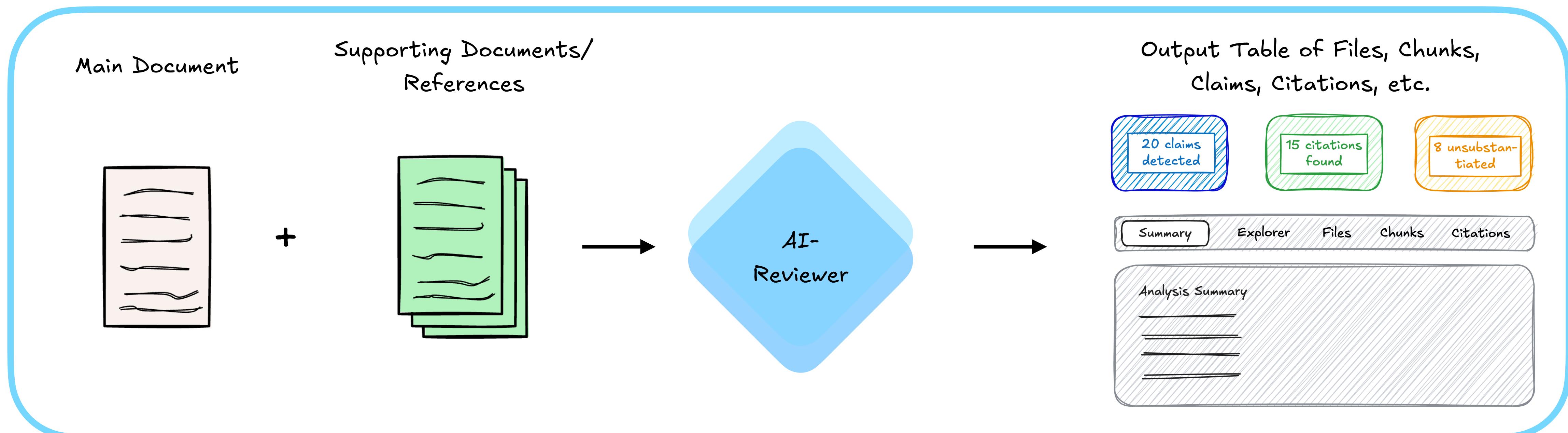


# Motivation: Application Introduction

- (Motivating Question)-

How can we build an AI system  
that we can trust to conduct  
the main tasks of peer review?

- (General Processing Pipeline) -



# Presentation Roadmap



## I. Motivation

### II. Video Demo

### III. System Details

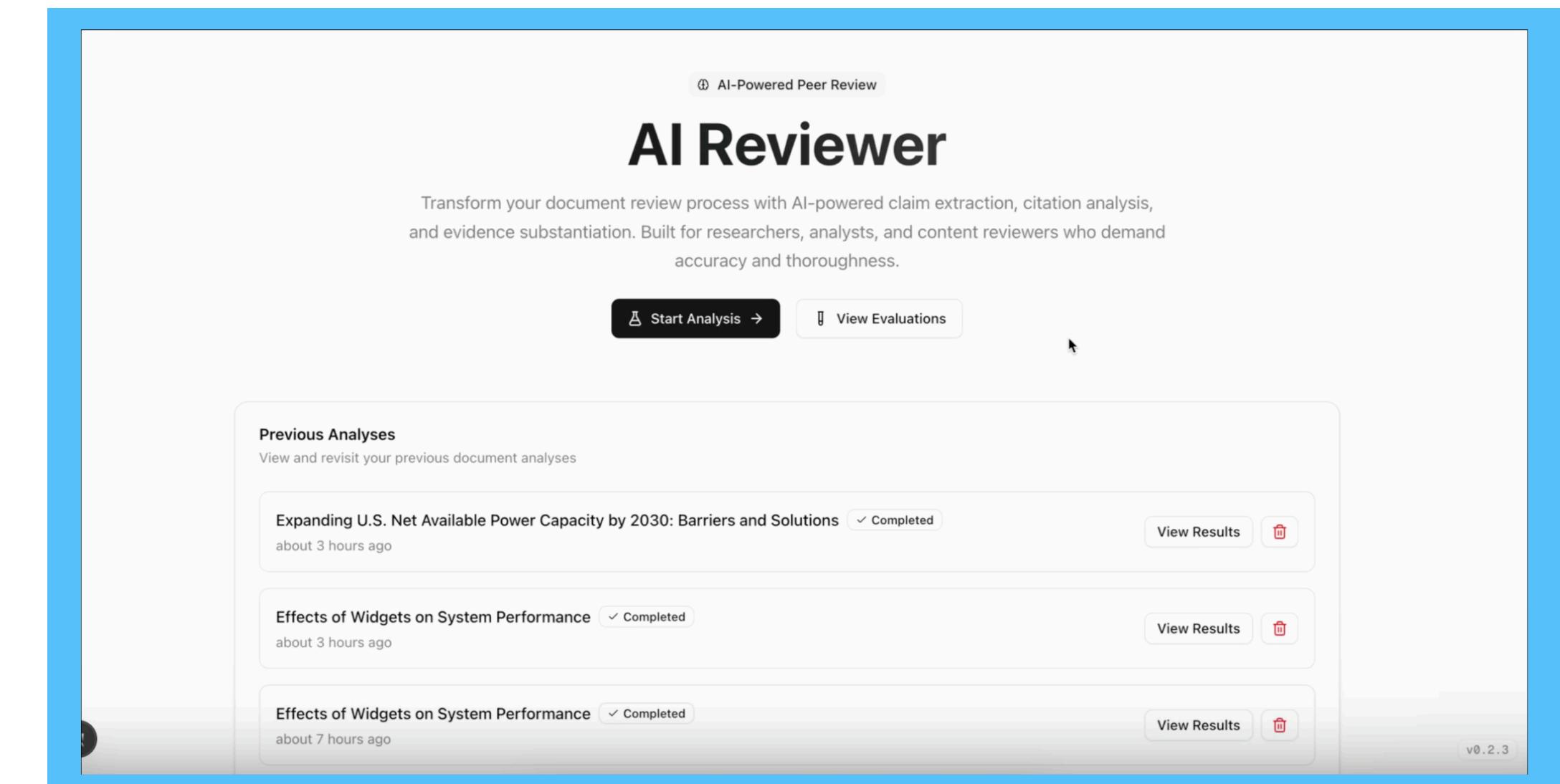
- .i Architecture

- .ii Agents and LangGraph

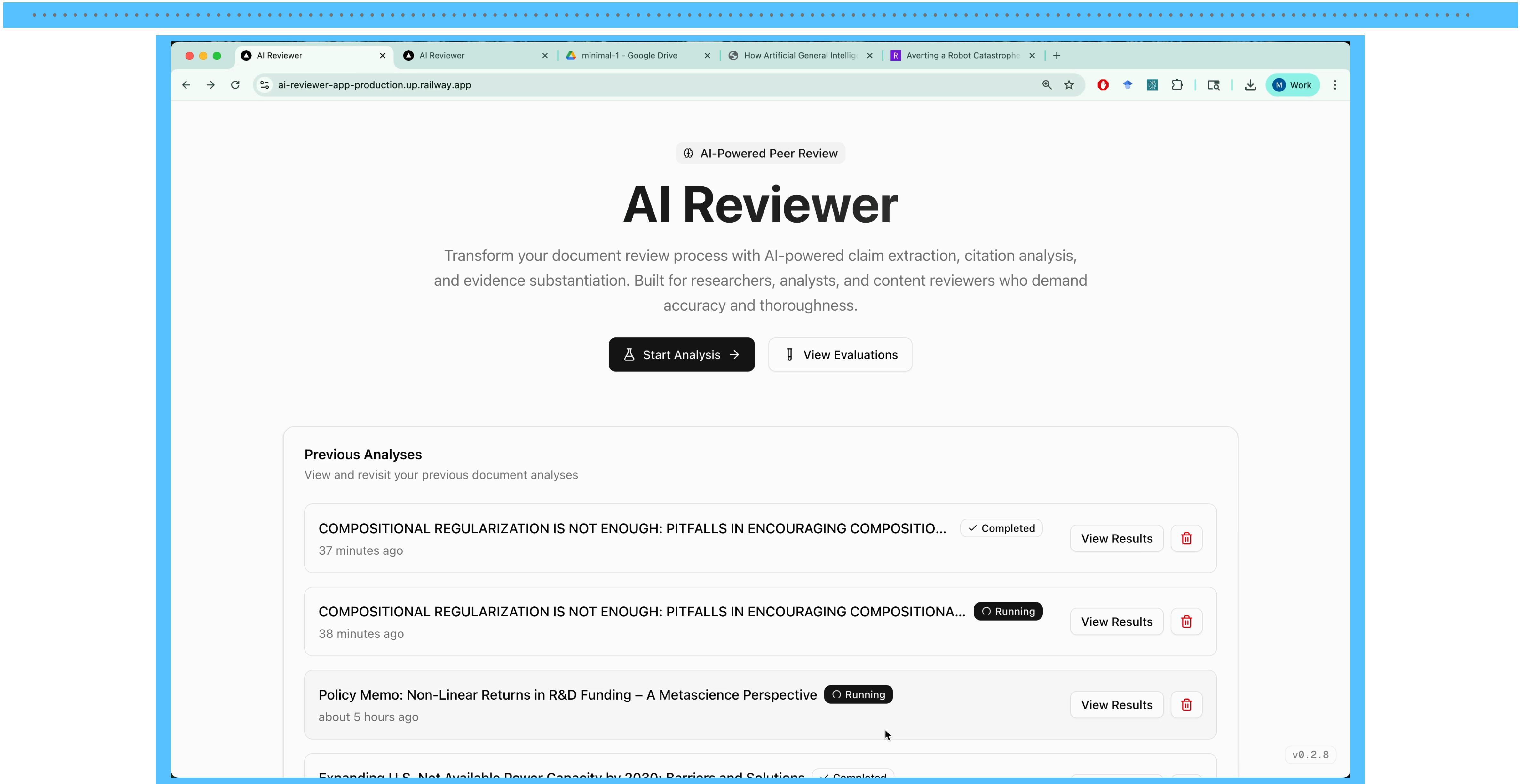
- .iii Design principles

### IV. Gaps and Next Steps

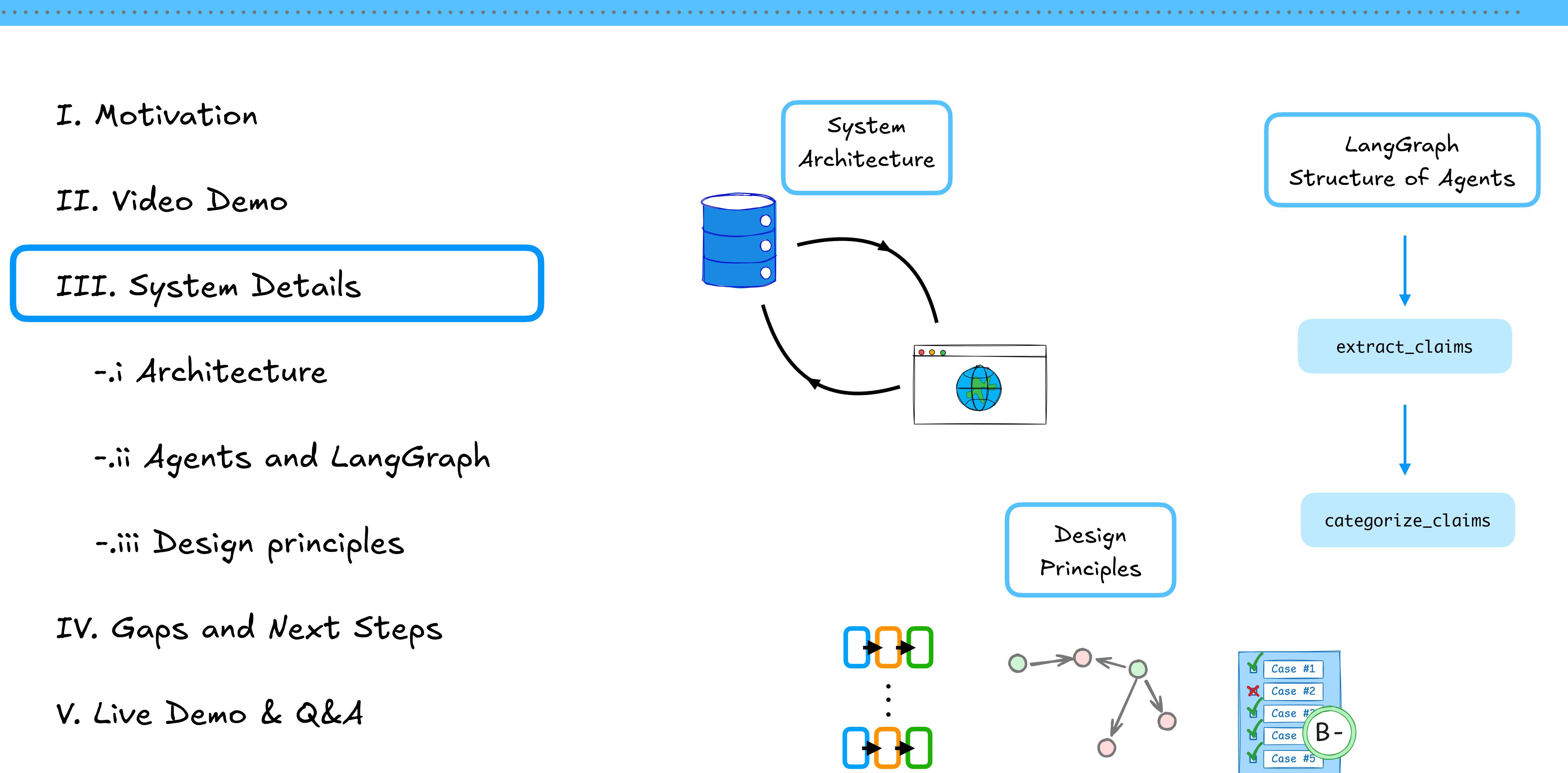
### V. Live Demo & Q&A



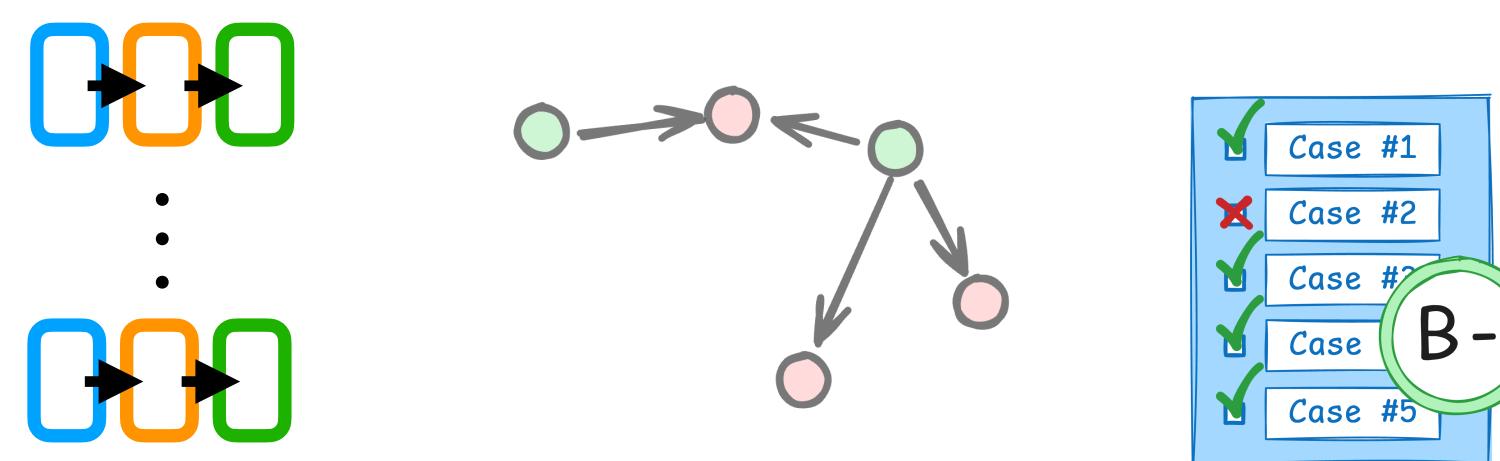
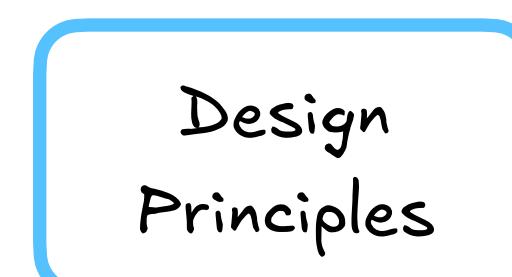
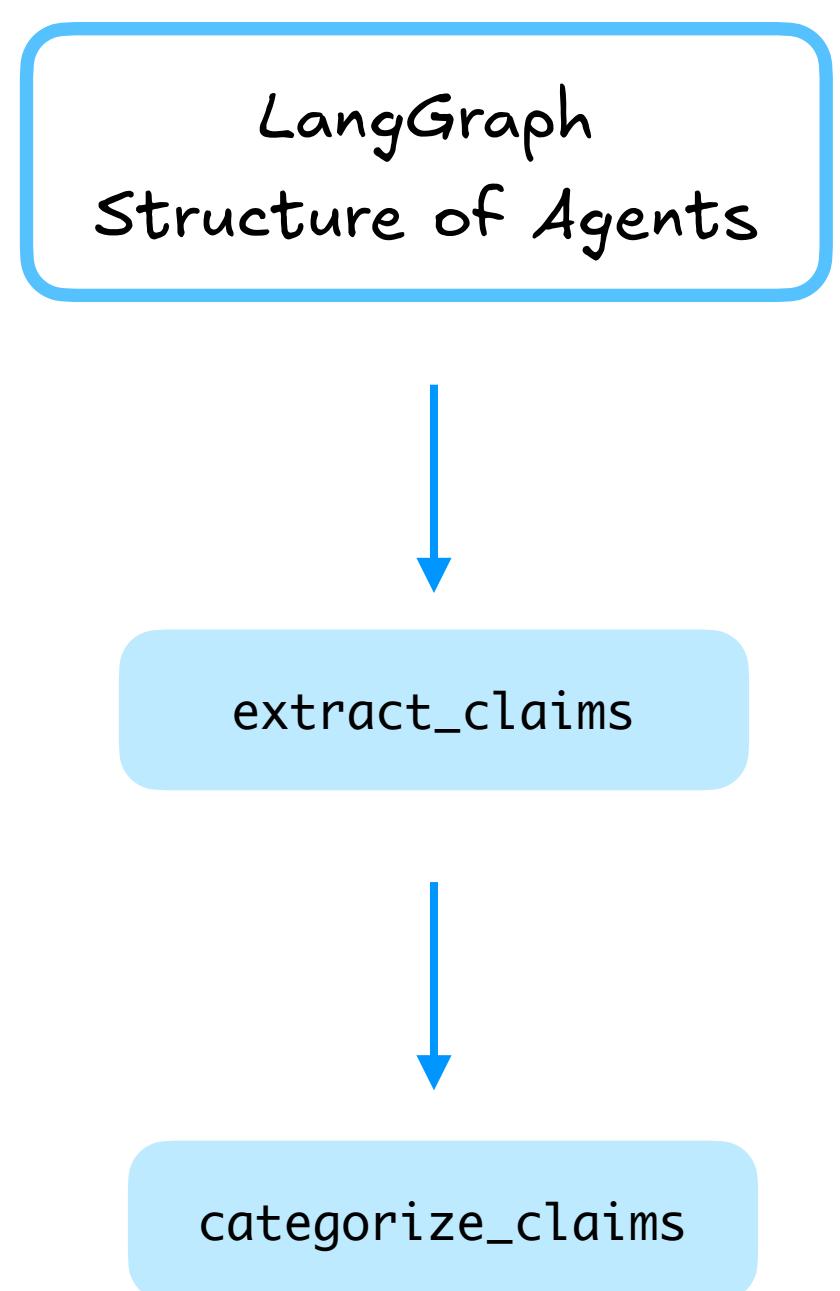
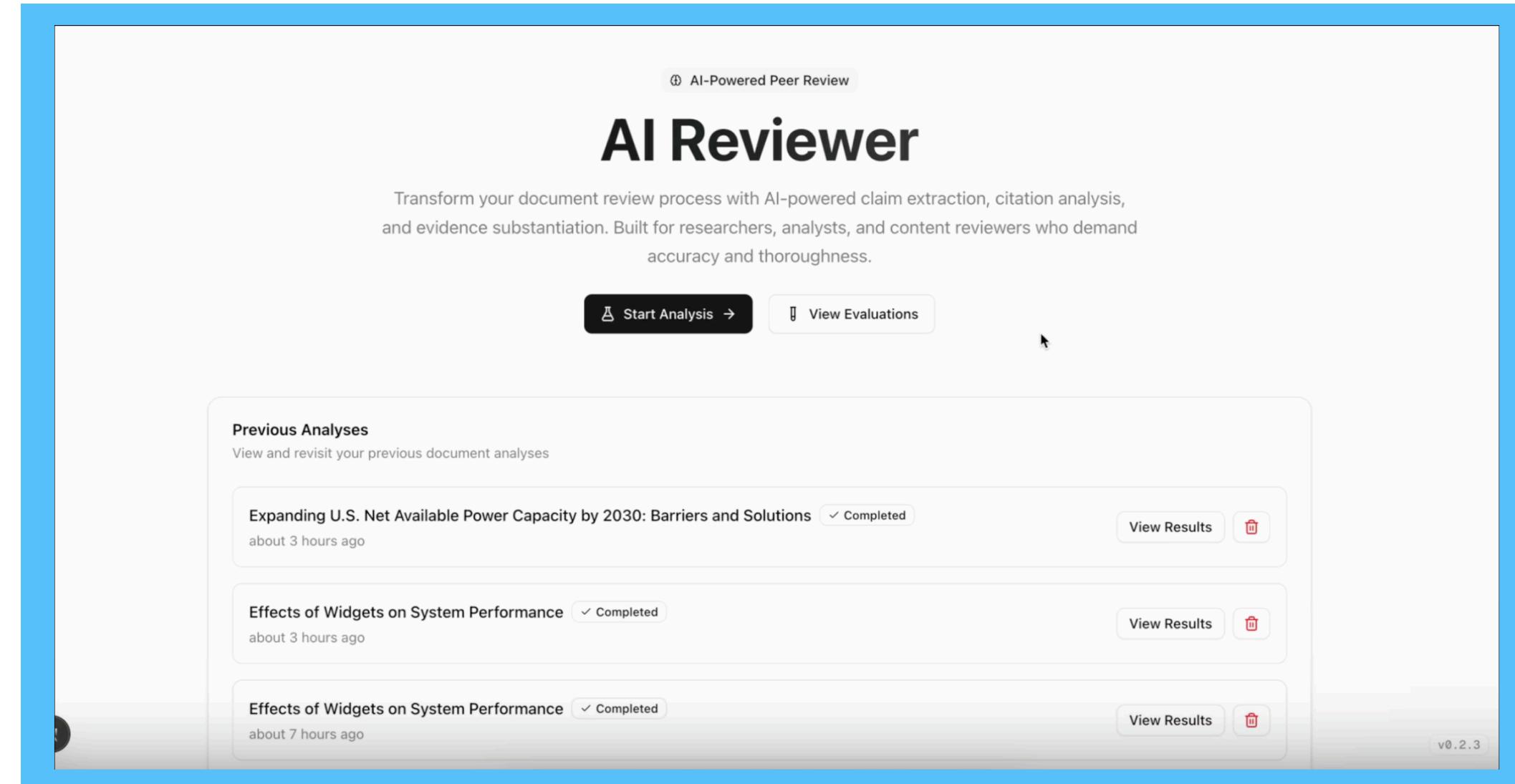
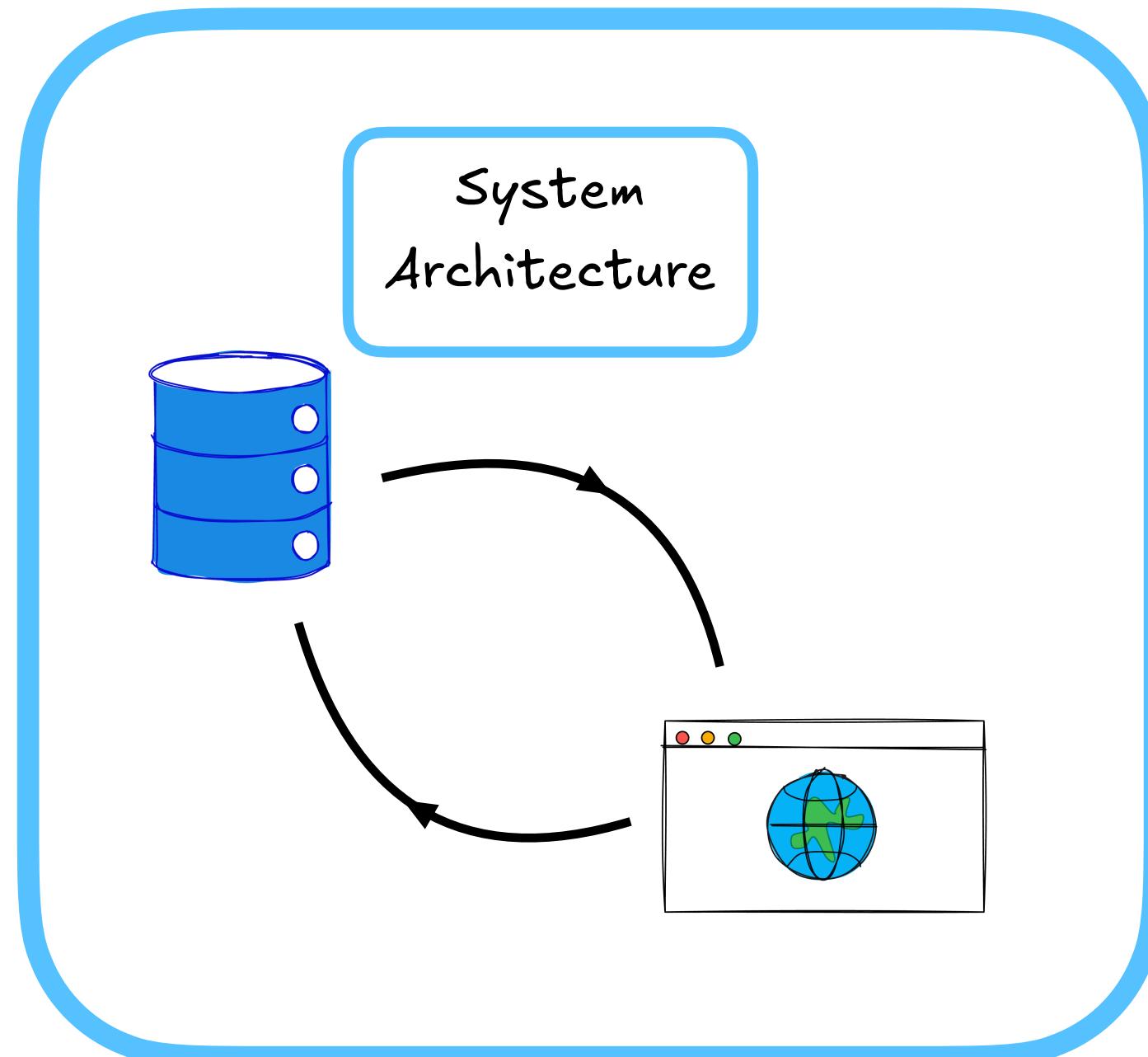
# Video Demo of Application



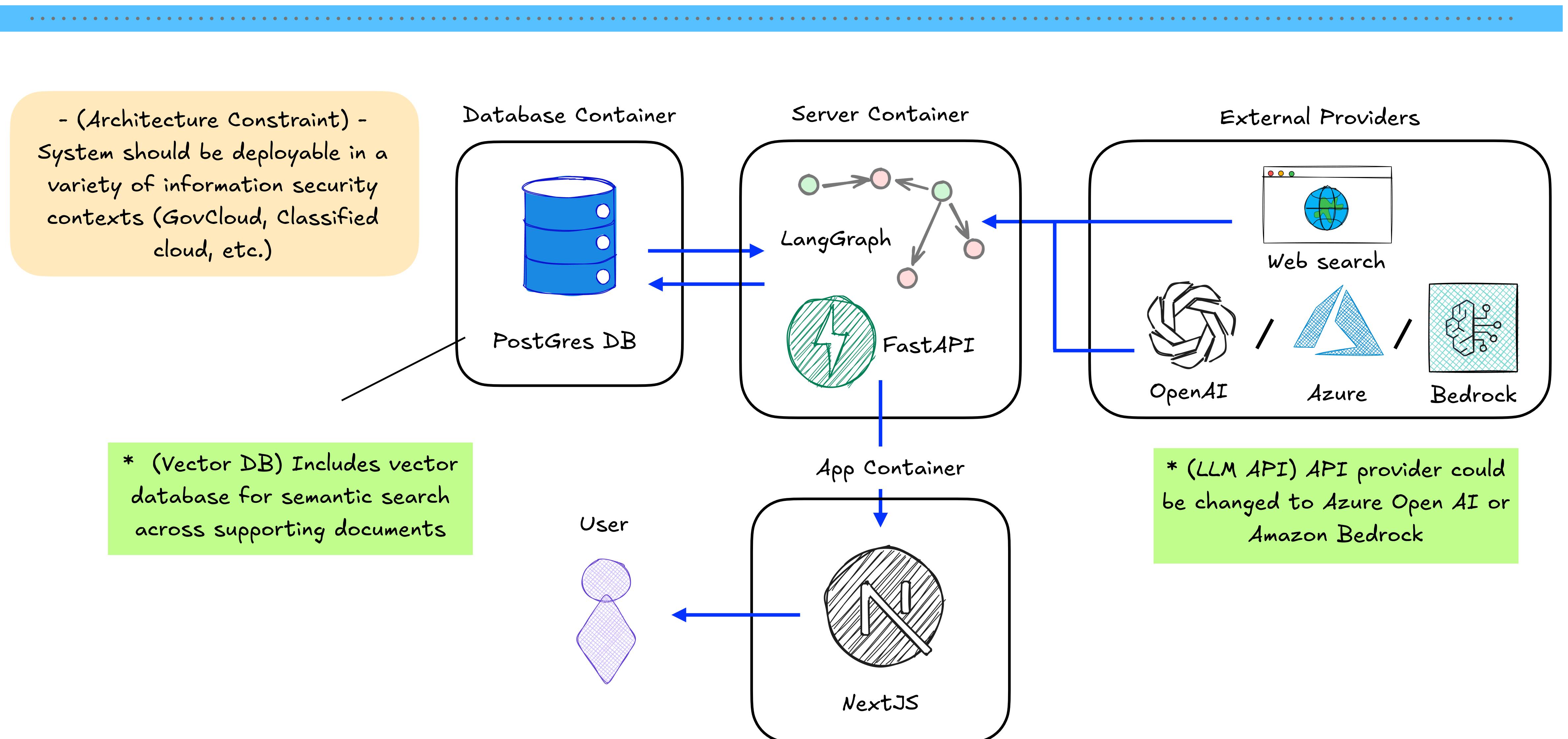
# Presentation Roadmap



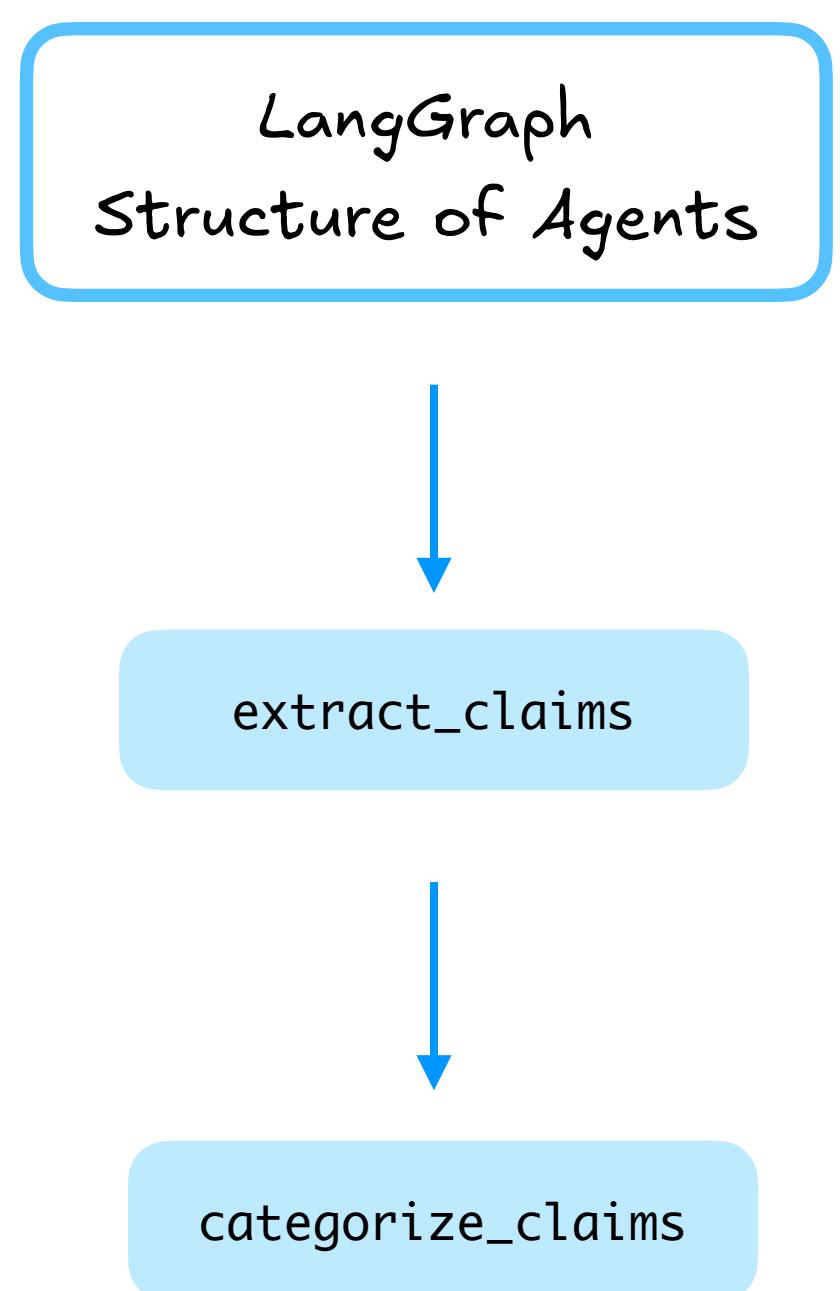
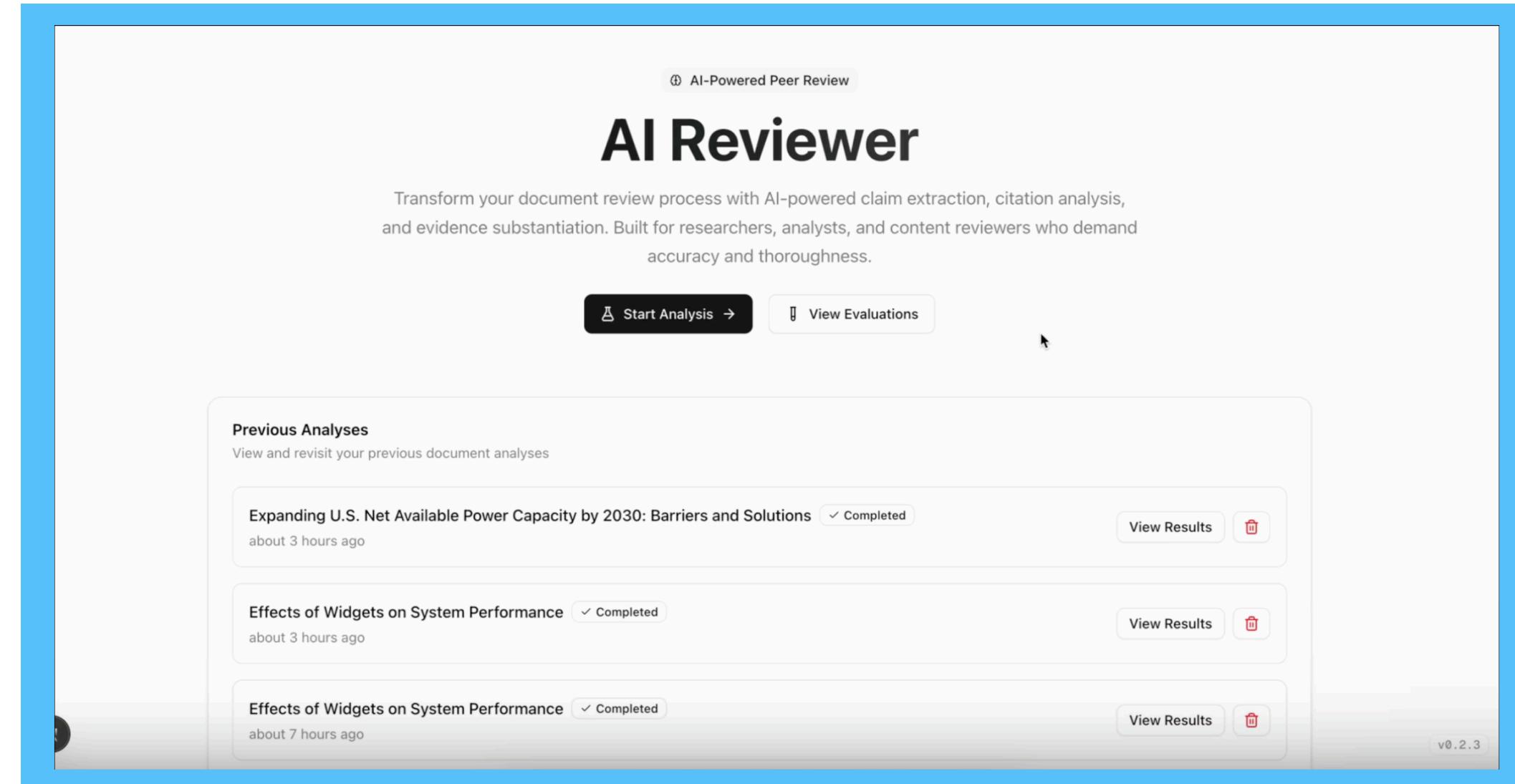
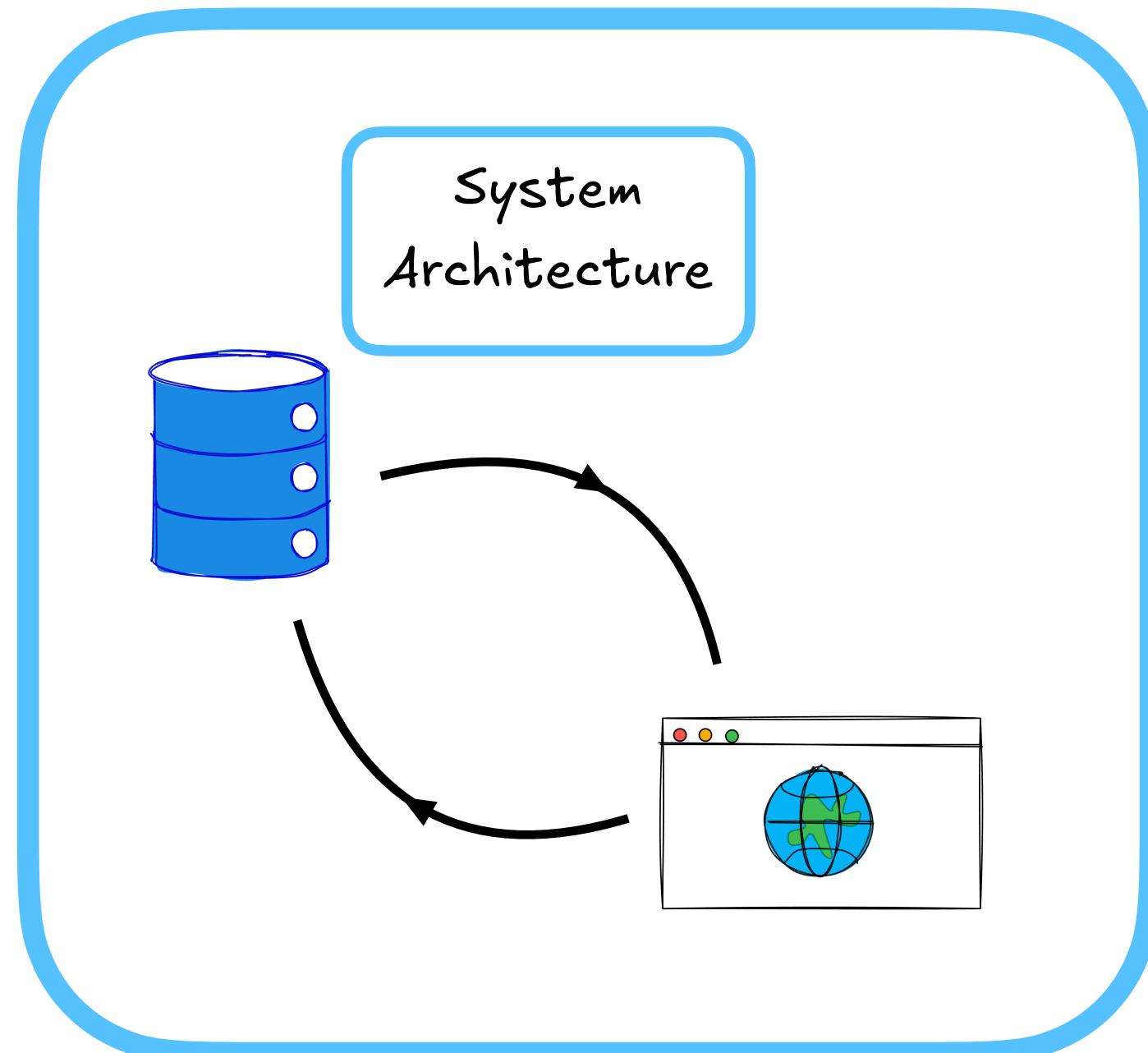
# System Details: Architecture



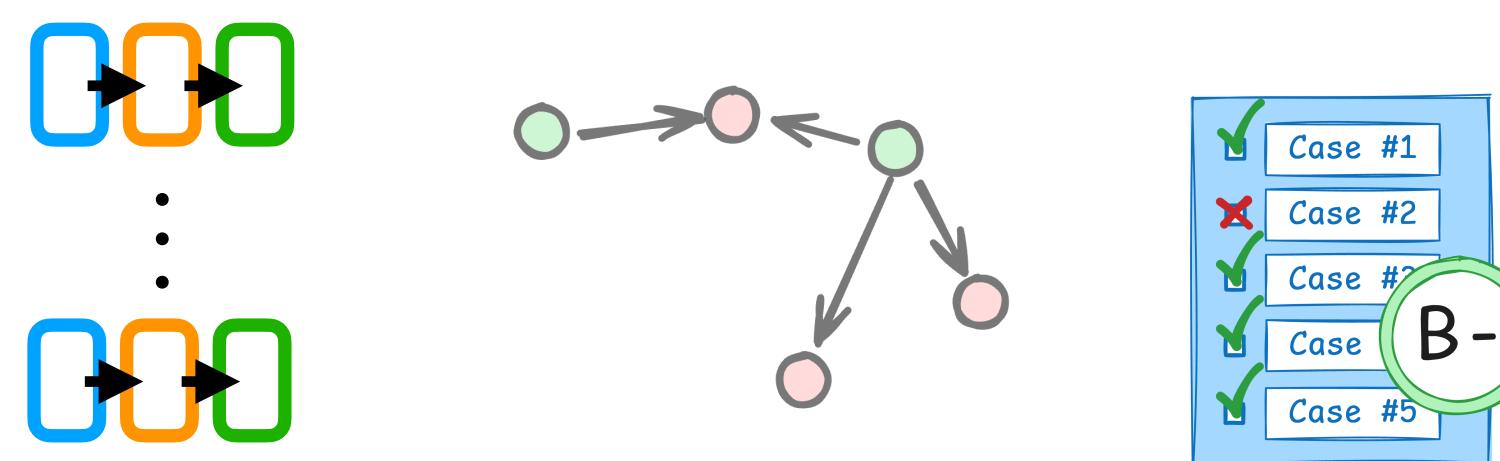
# System Details: Architecture



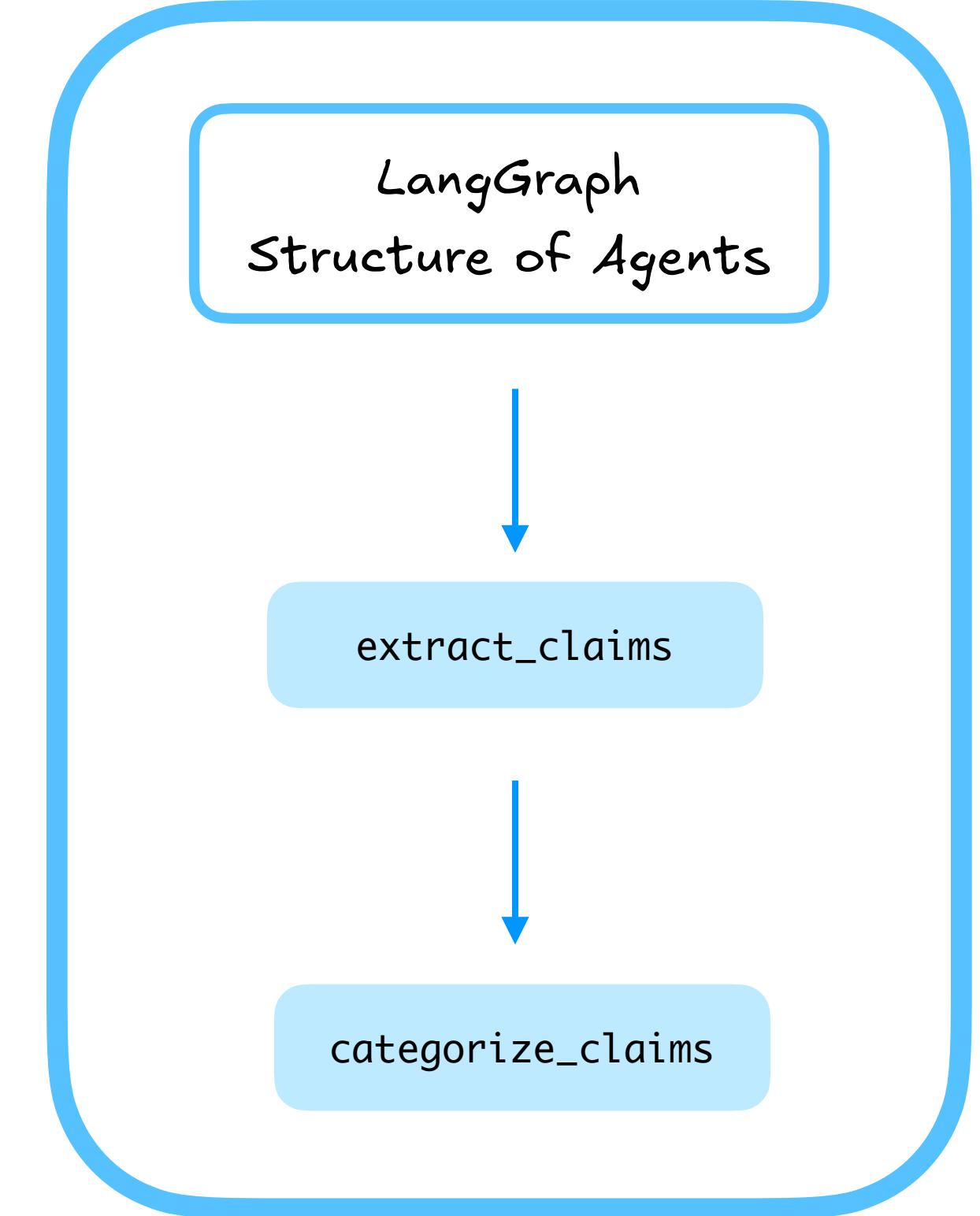
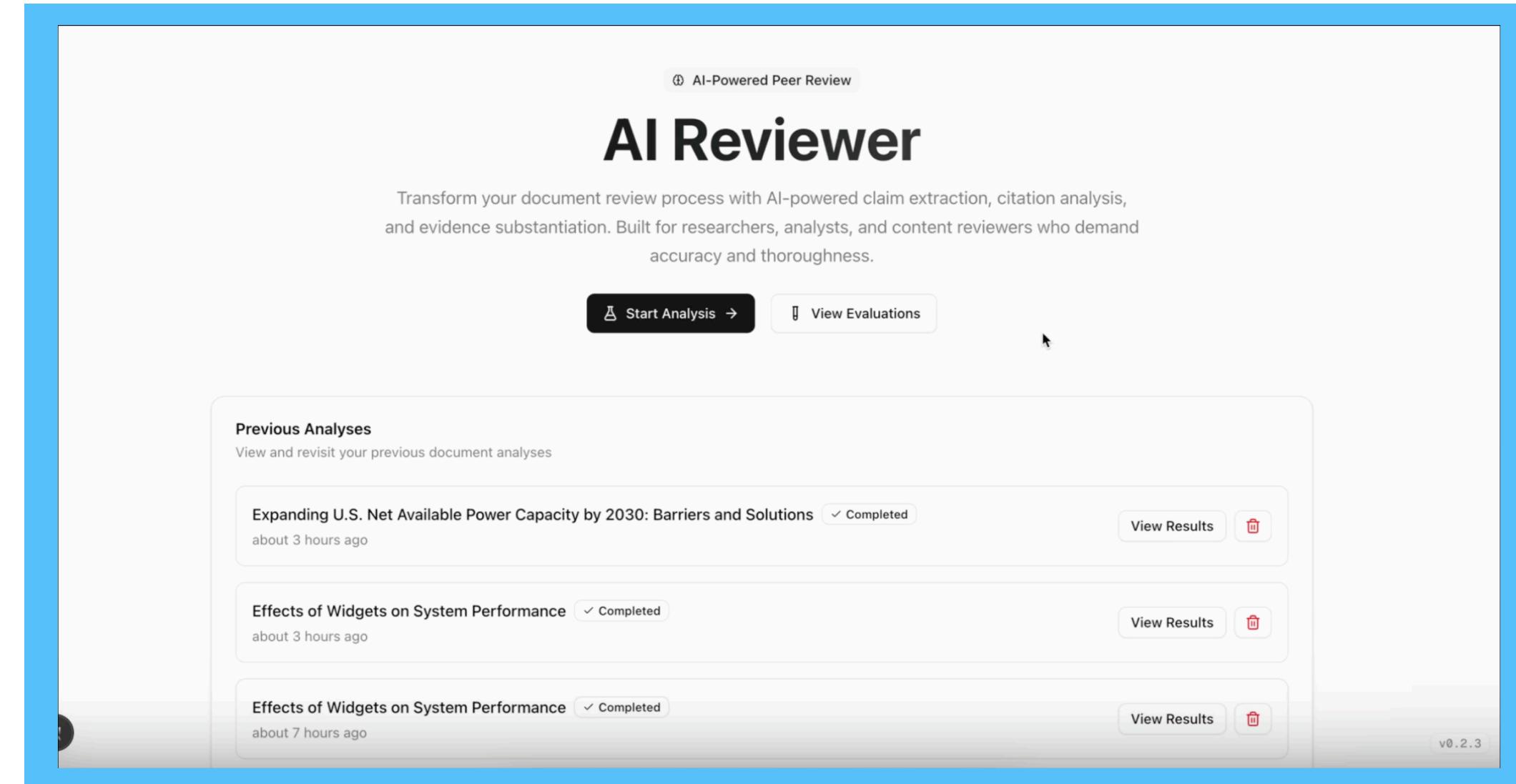
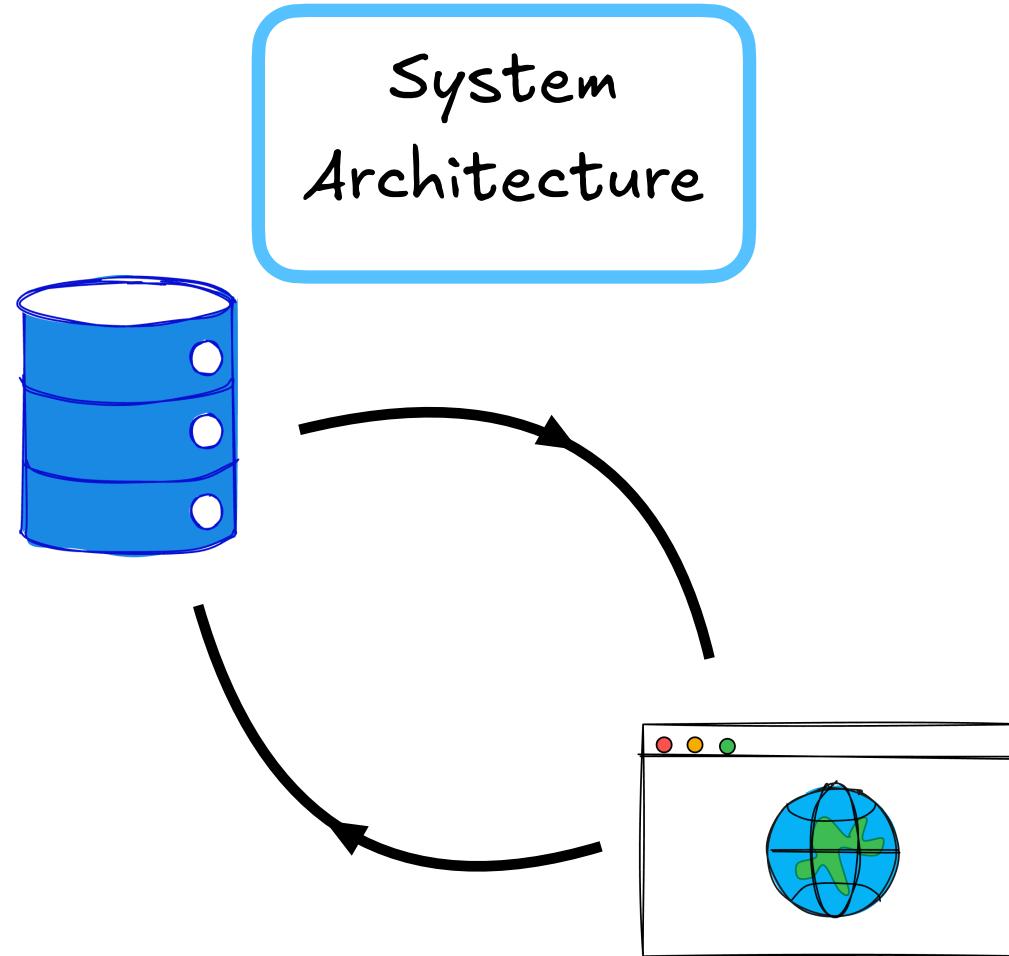
# System Details: Architecture



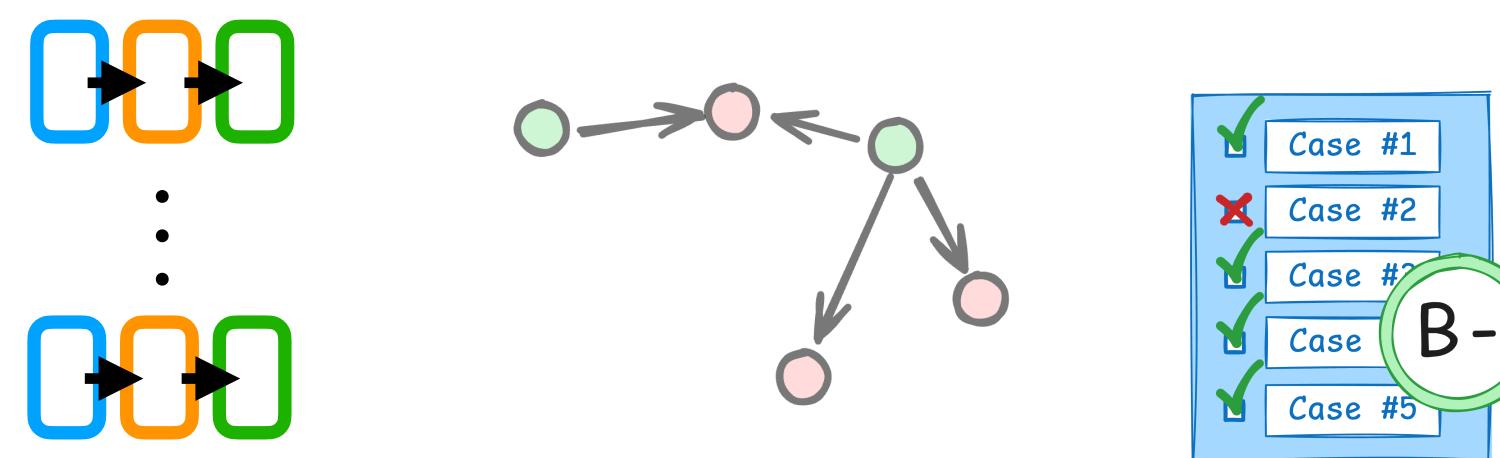
Design  
Principles



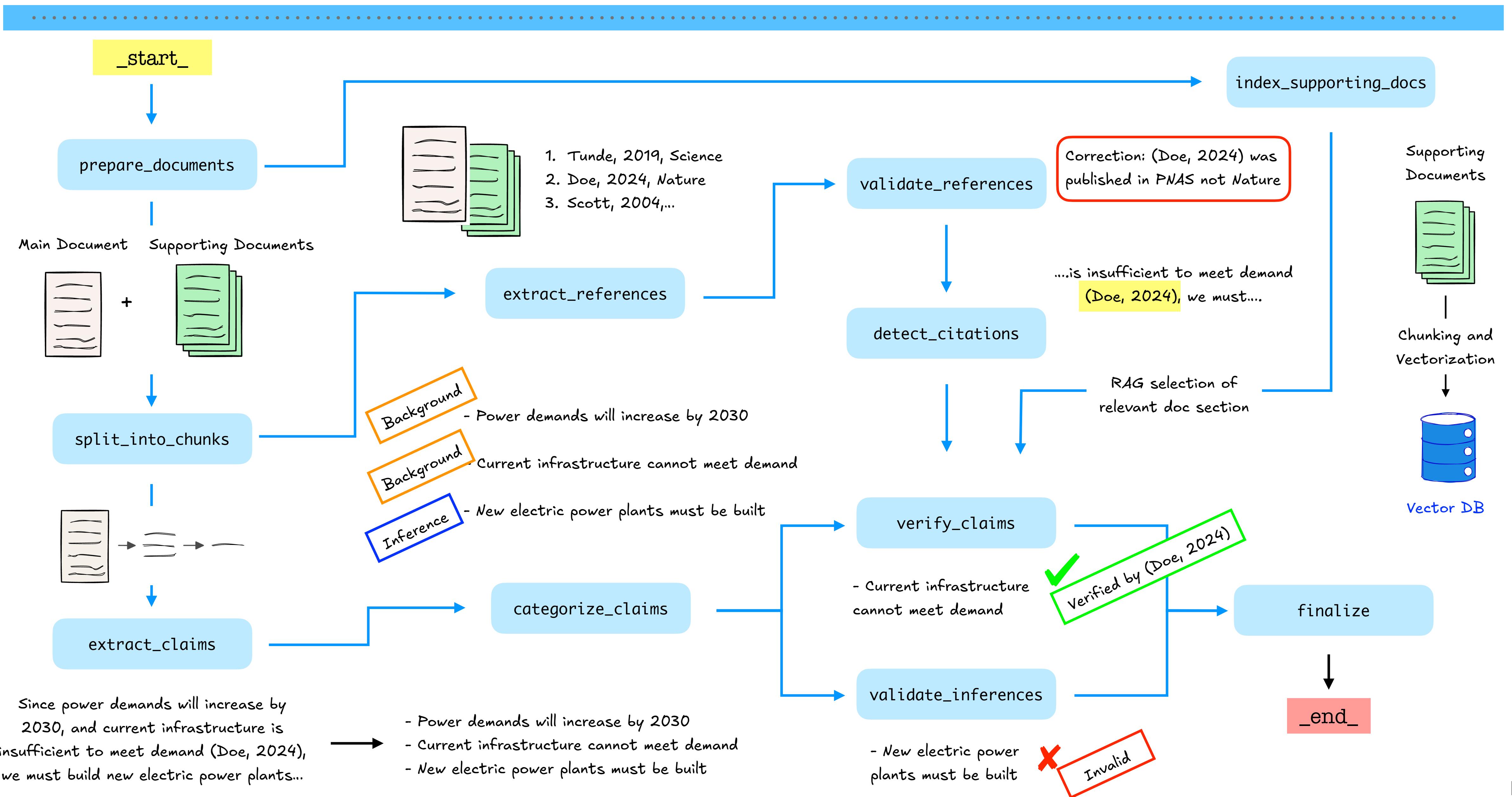
# System Details: LangGraph Structure



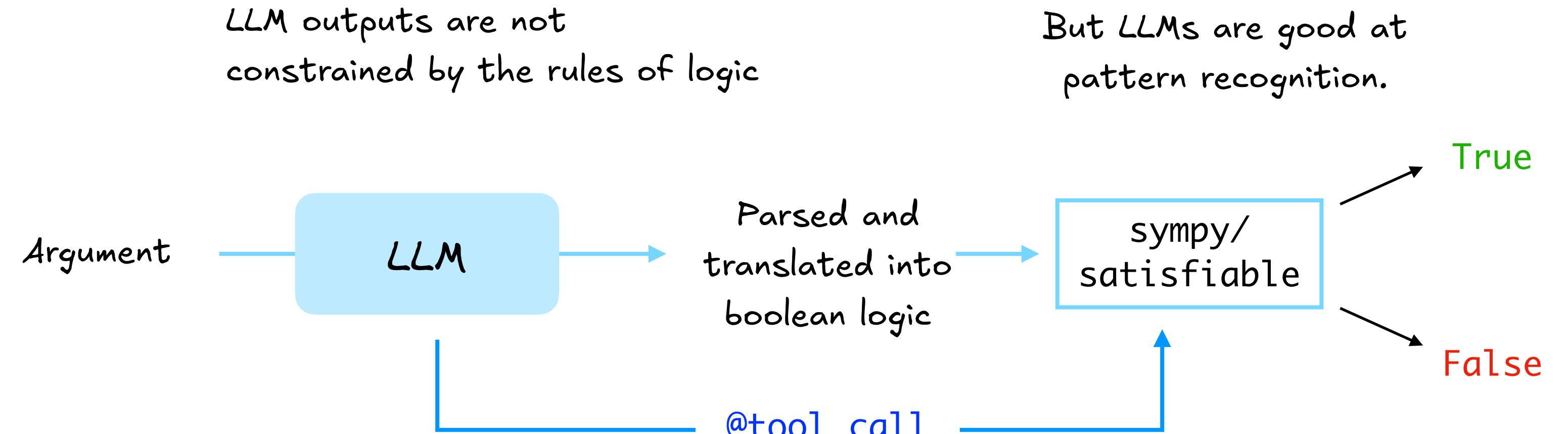
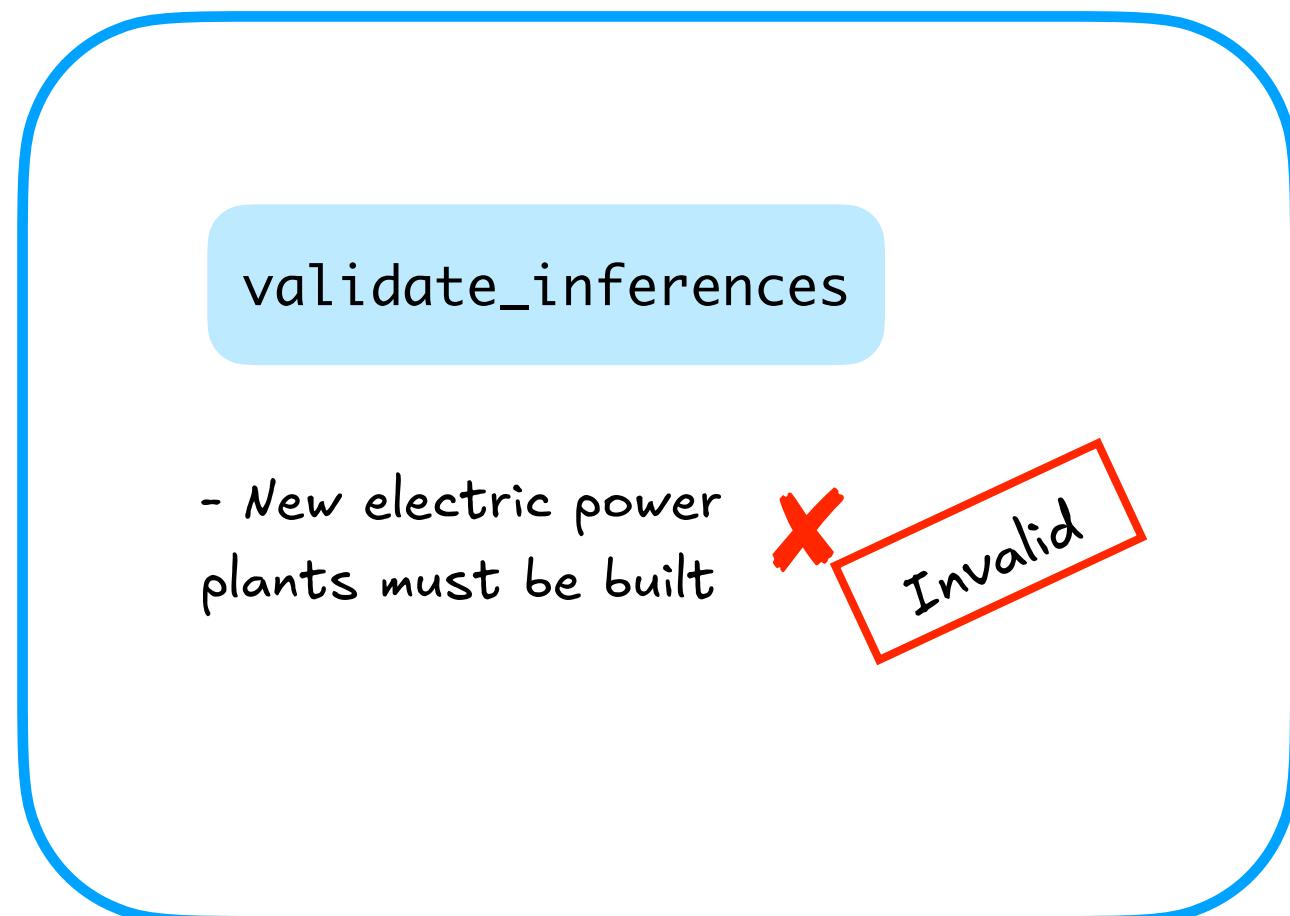
Design  
Principles



# System Details: LangGraph Structure



# System Details: Symbolic Logic and Inference Evaluation



Classic Example

Statement	Logic Symbol
i. Socrates is a man	P
ii. If socrates is a man, he is mortal	P implies Q
iii. Socrates is mortal	Q

Valid ✓  
Argument

It is impossible for "P" and "P implies Q" to be true while "Q" is false.

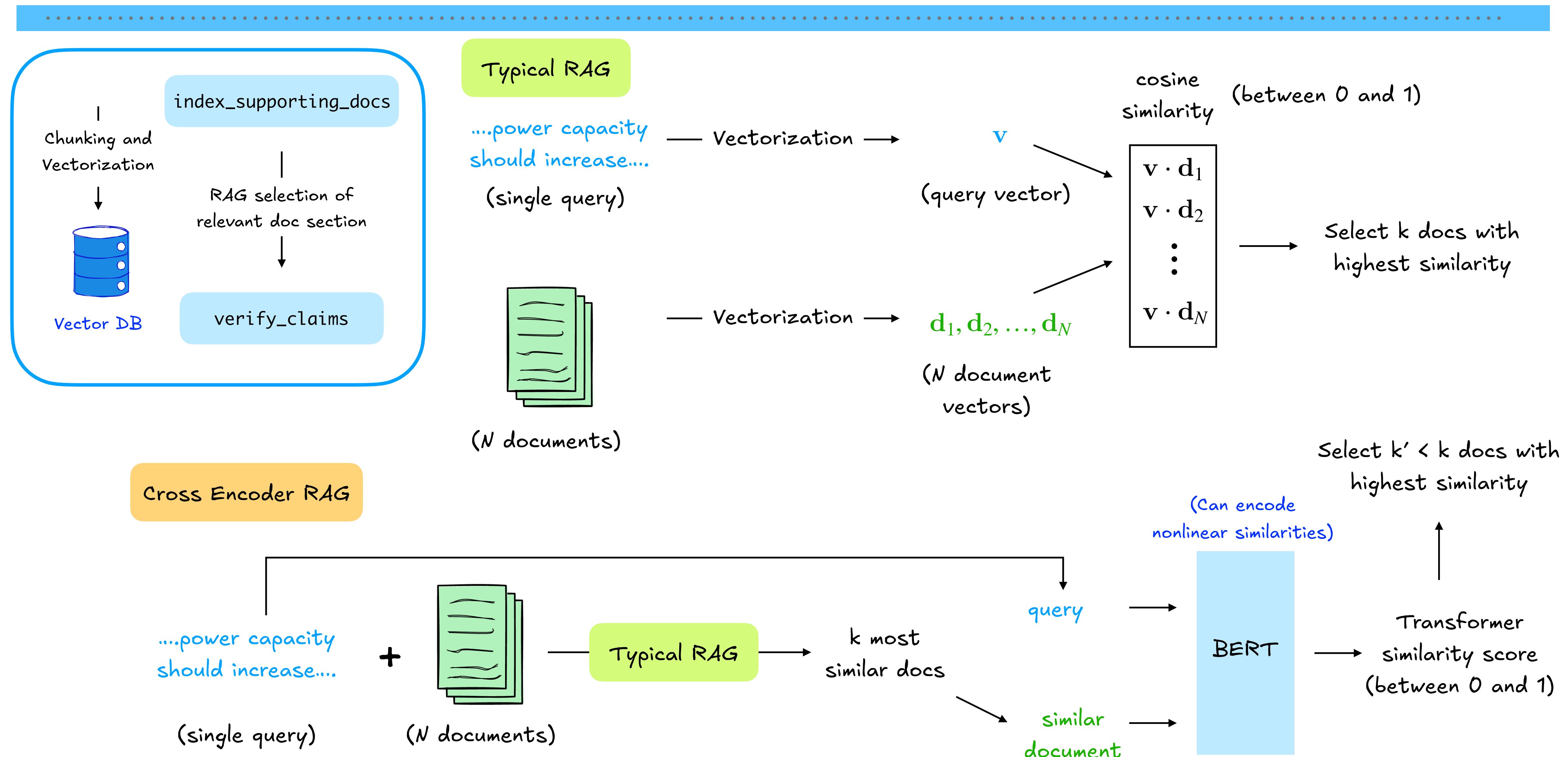
Applied Example

Statement	Logic Symbol
i. Power demands will increase by 2030	P
ii. Current infrastructure cannot meet demand	Q
iii. New electric power plants must be built	R

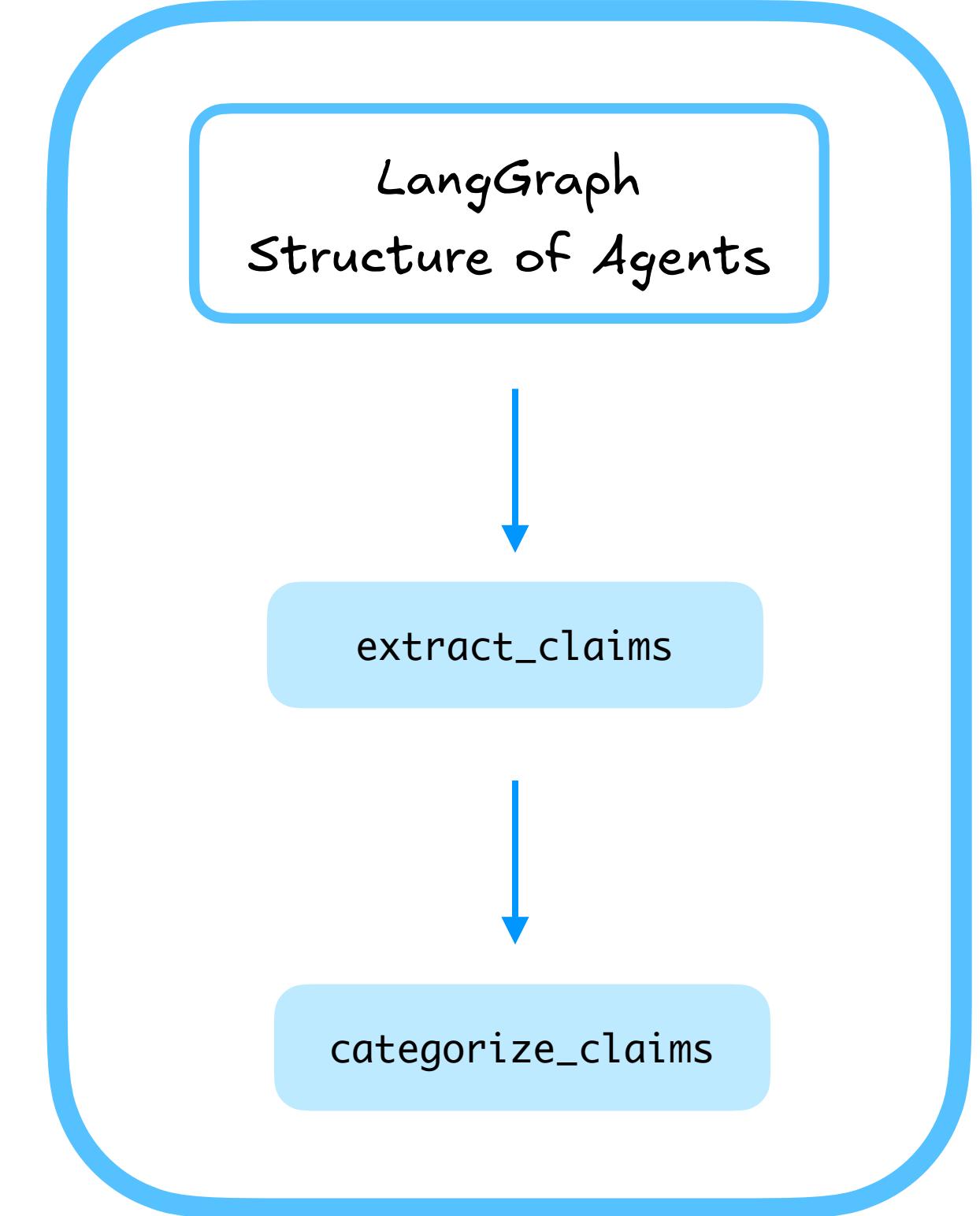
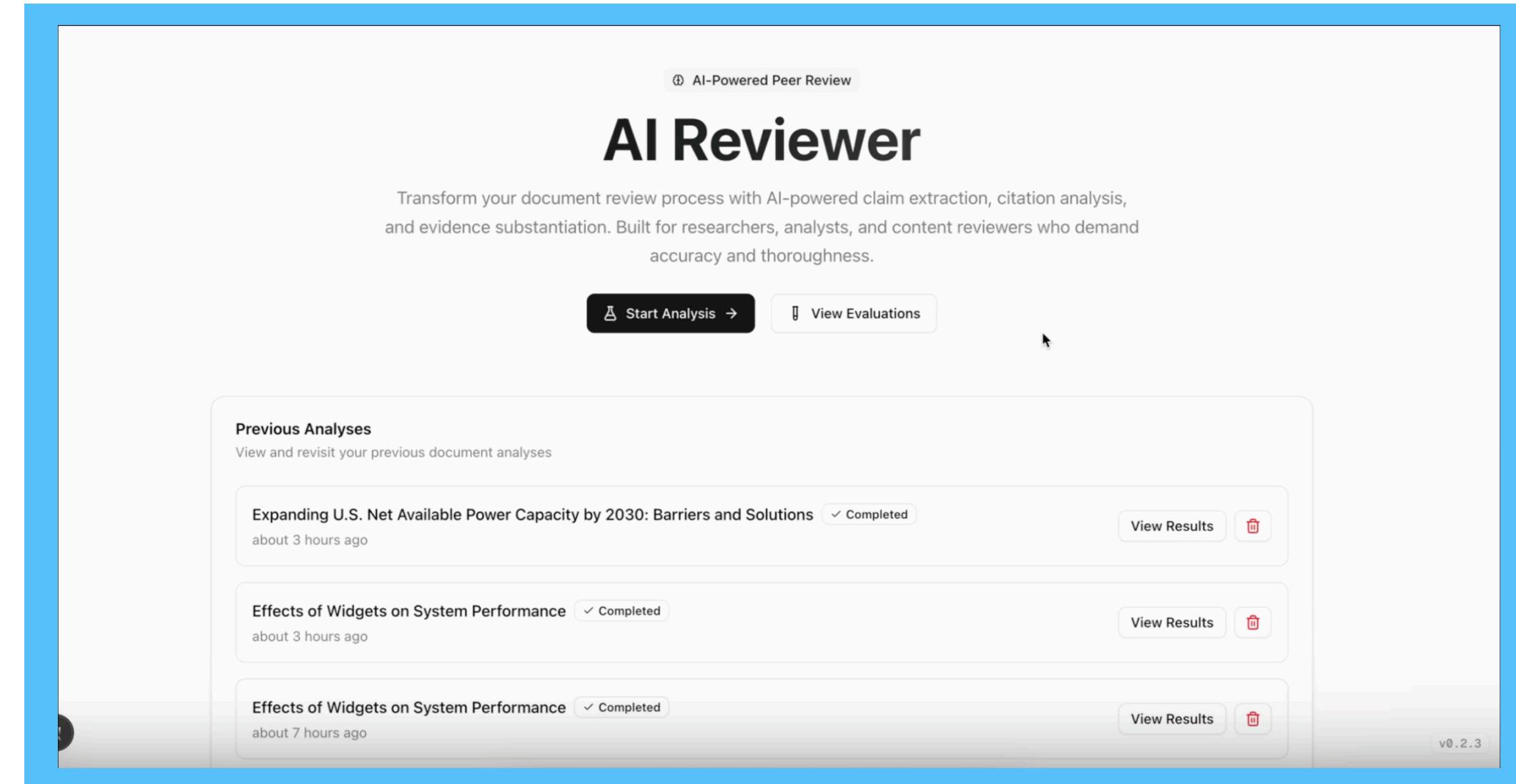
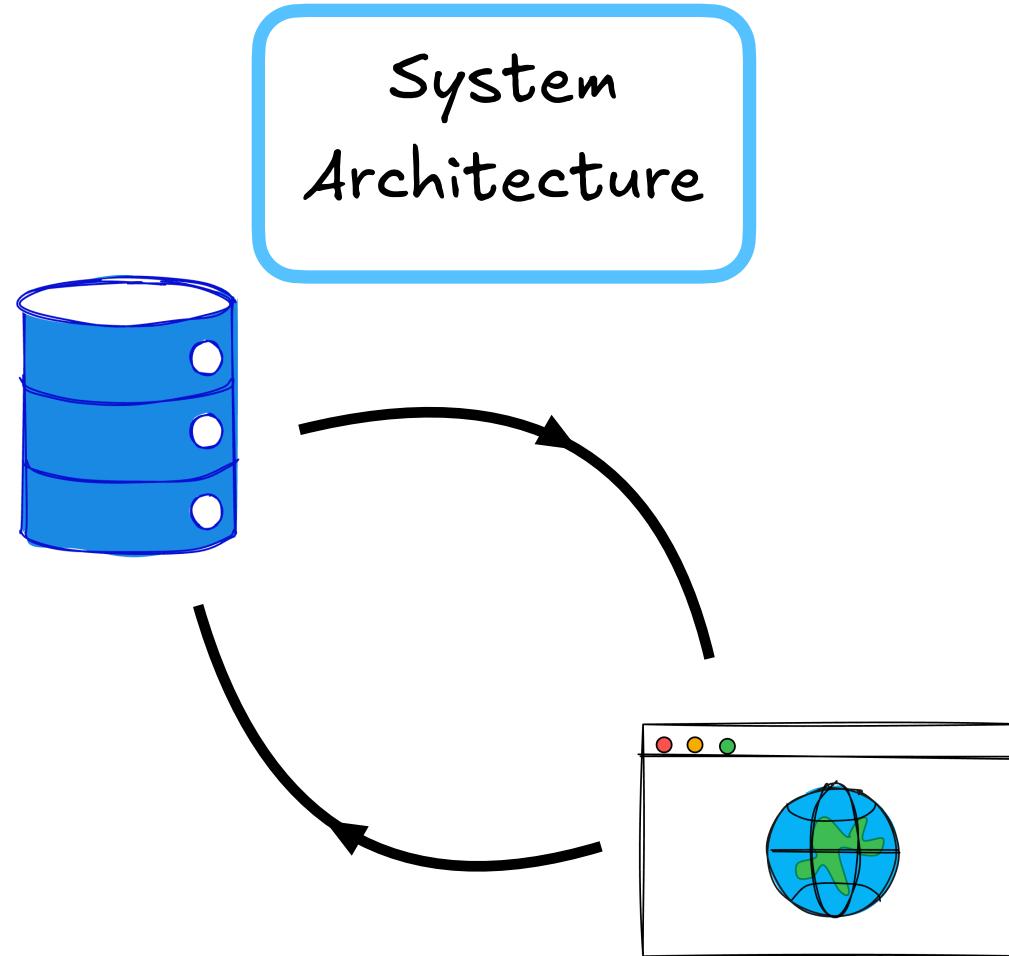
It is possible for P and Q to be true while R is false.

Invalid ✗  
argument

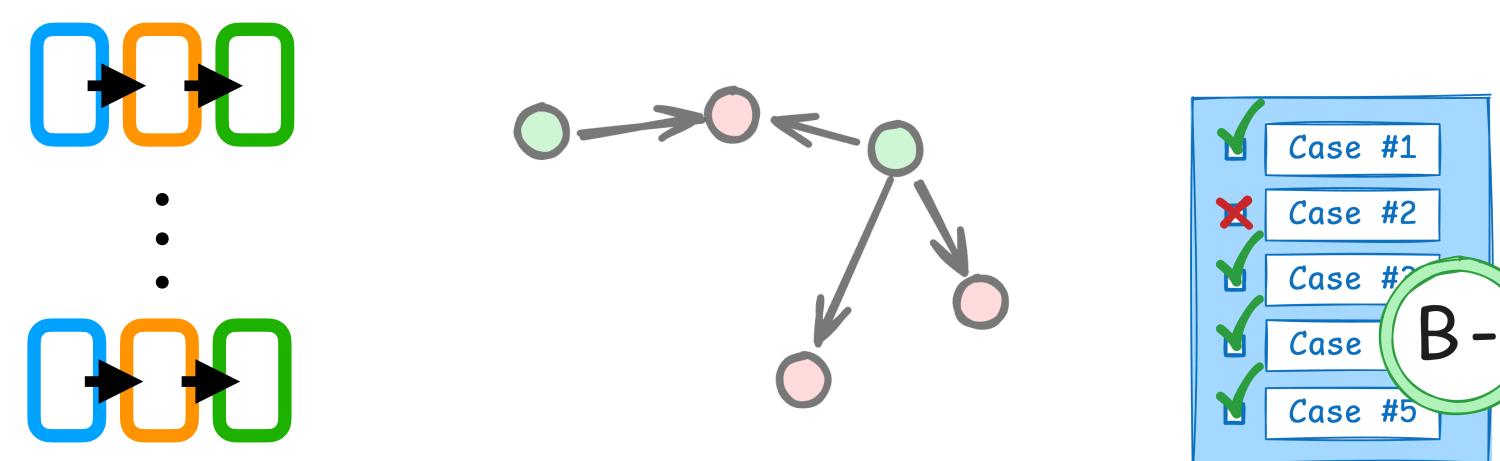
# System Details: Verify Citations and Cross Encoder RAG



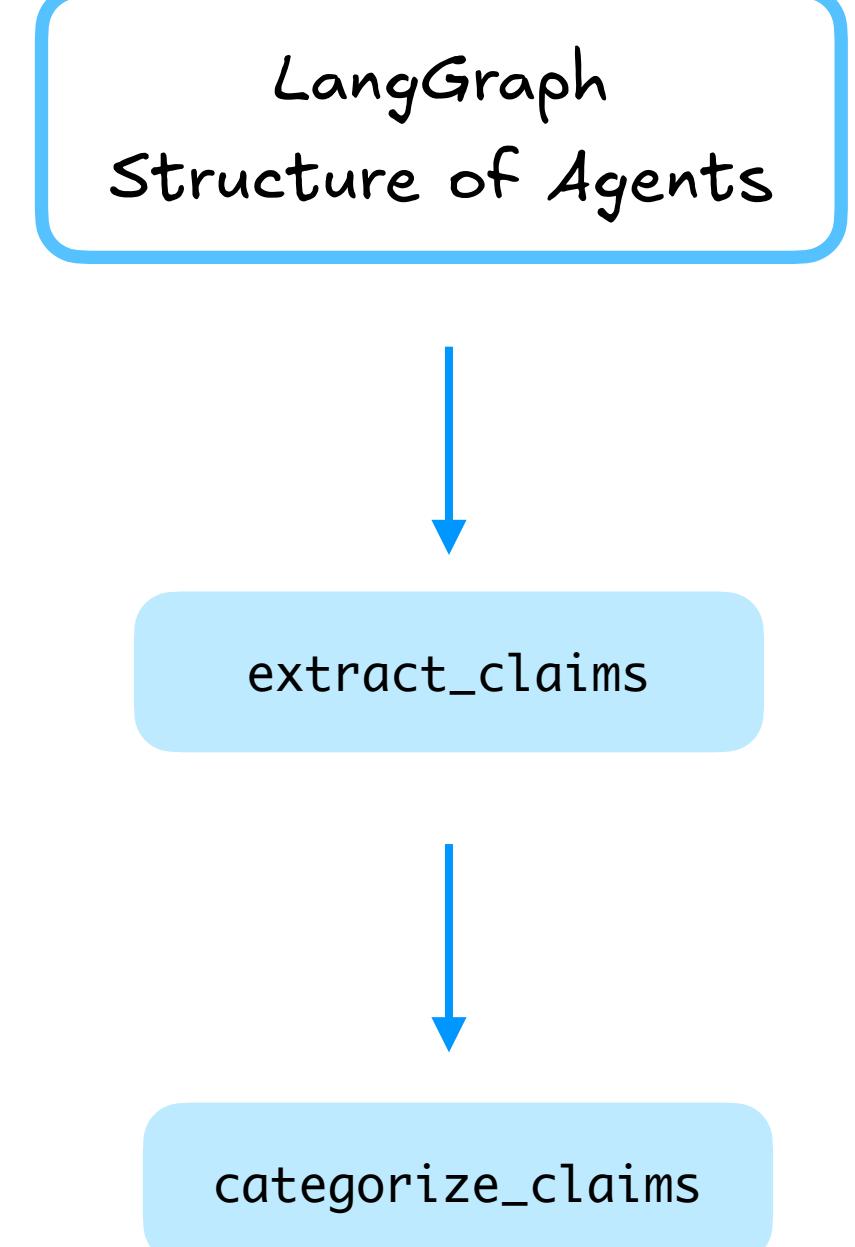
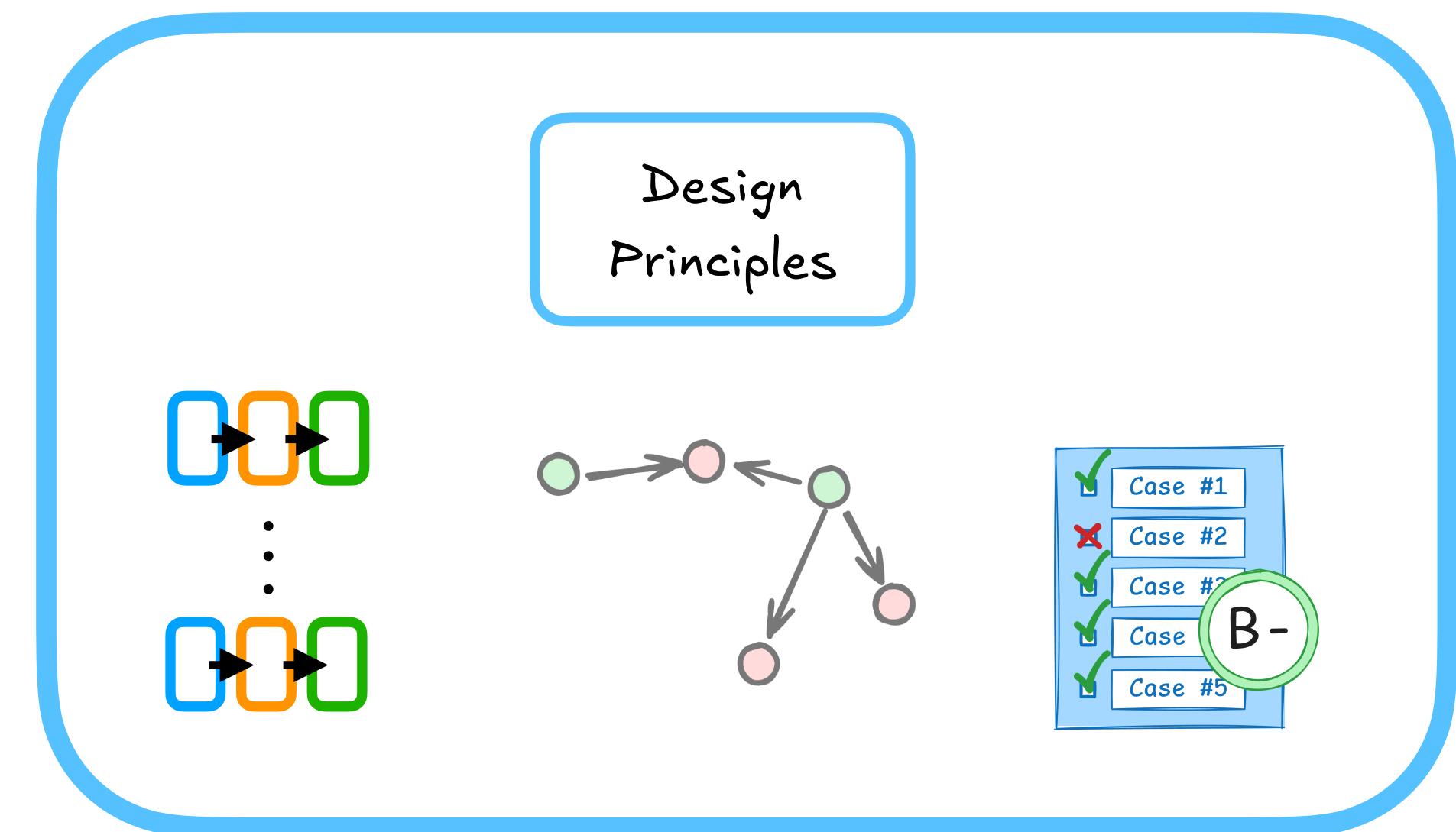
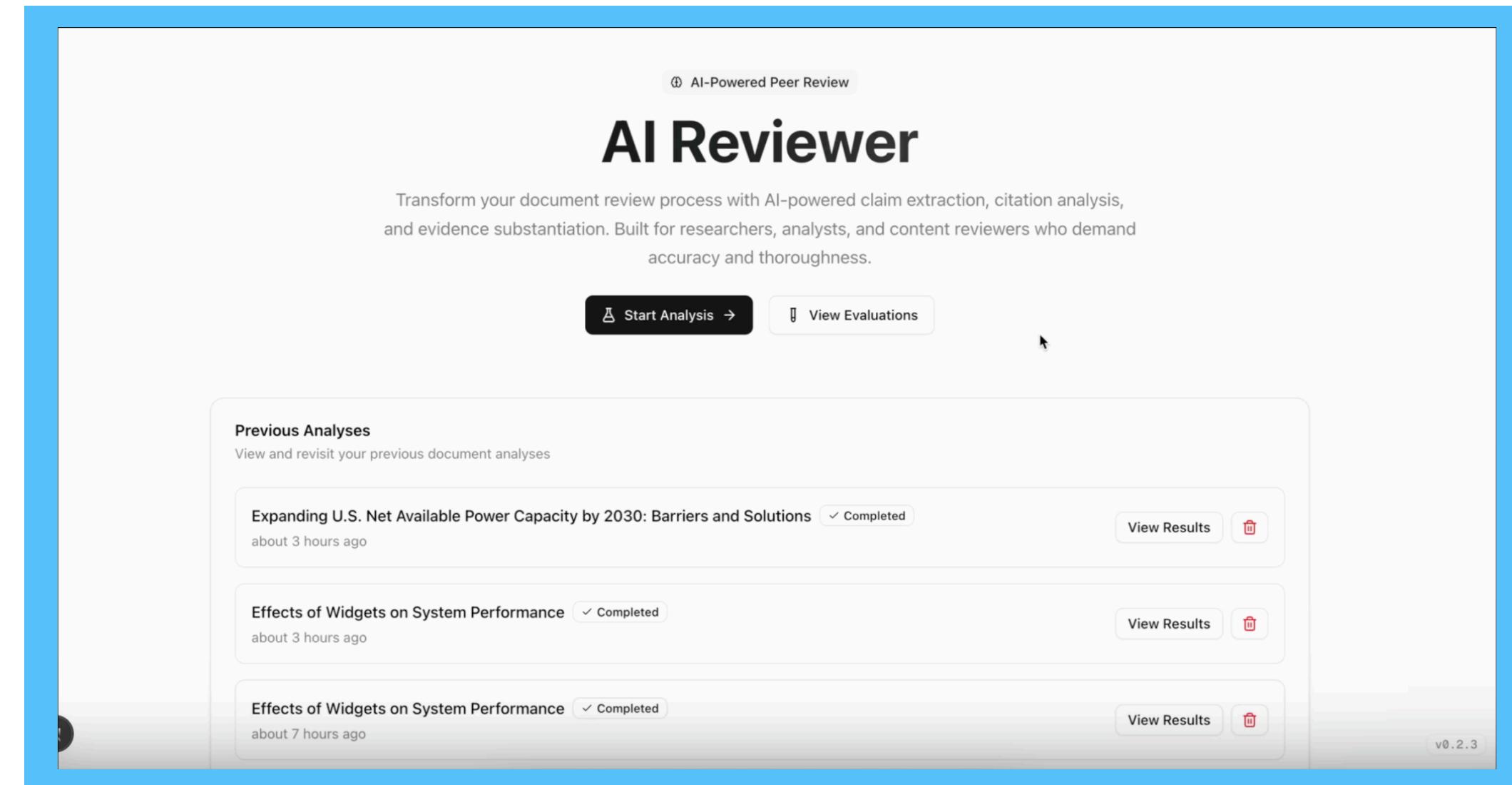
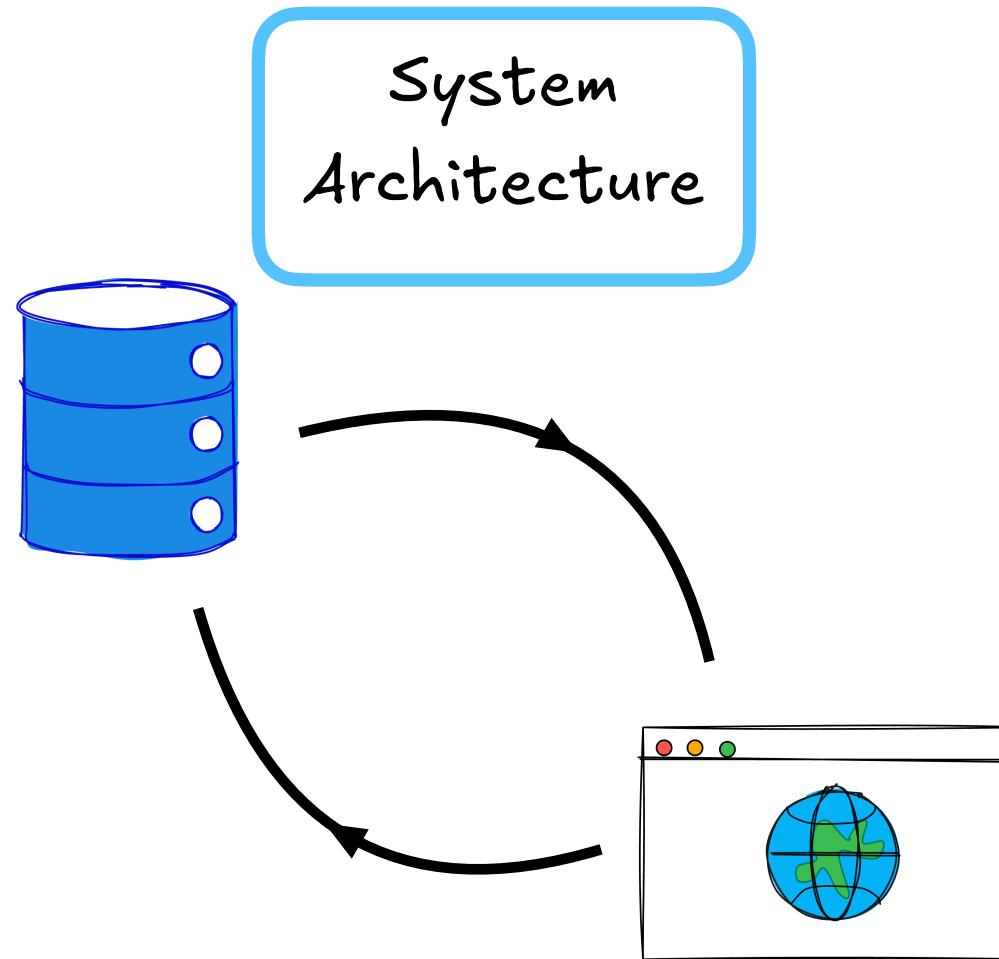
# System Details: LangGraph Structure



Design  
Principles



# System Details: Design Principles



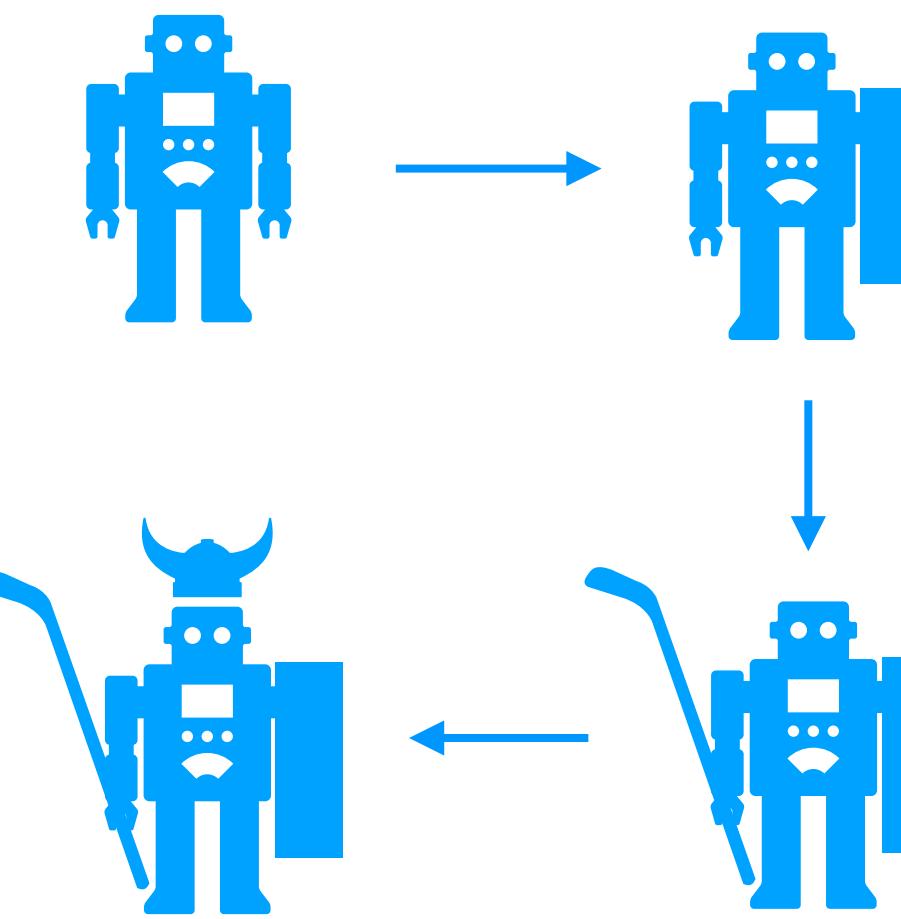
# System Details: Design Principles

Why are design principles necessary?

"Design Principles"

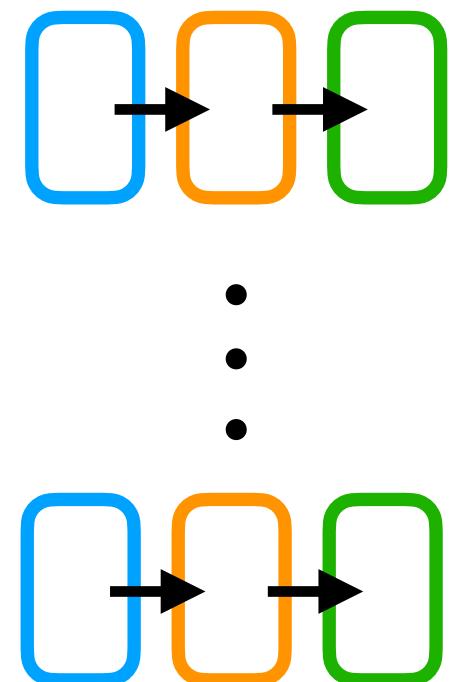
No system is final...

...so the created system must be able to (easily) handle future modifications

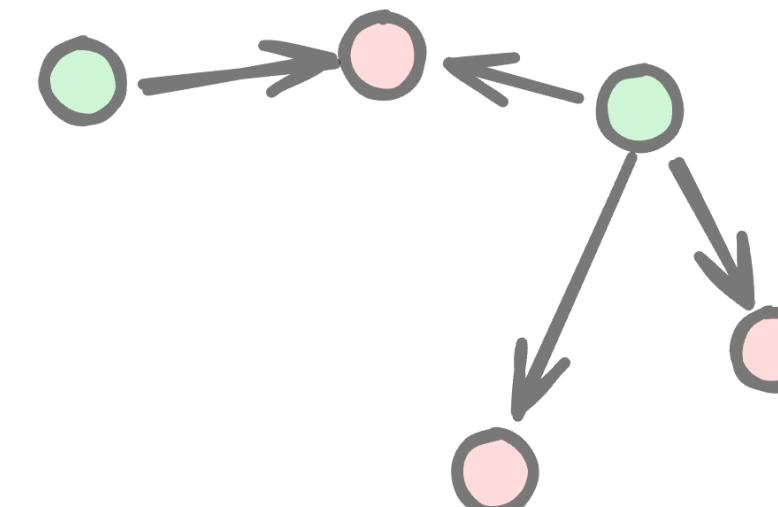


Design principles help ensure the system is built in a way that looks forward to possible change

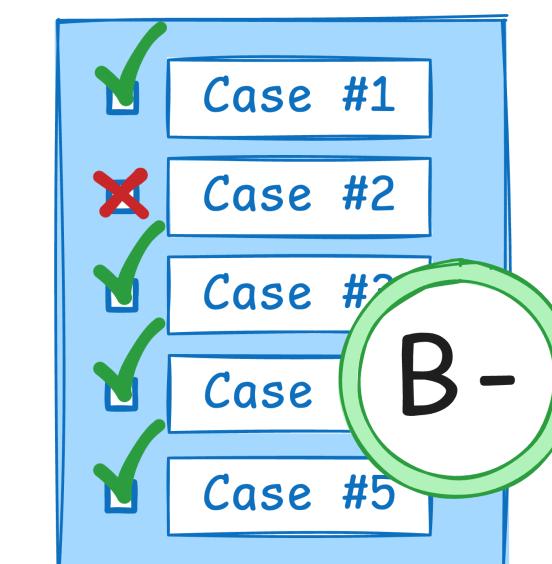
- (Uniformity) -



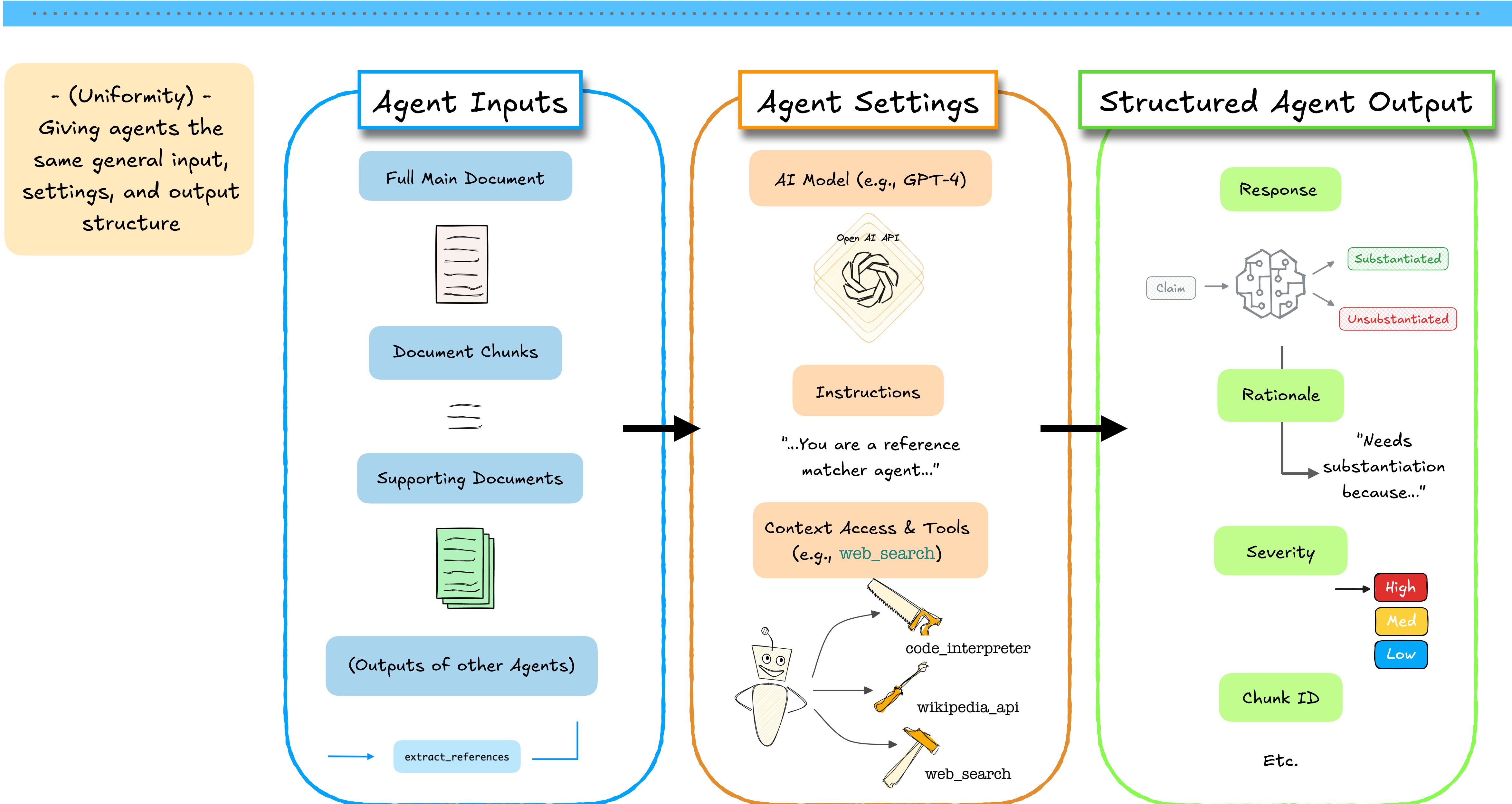
- (Granularity) -



- (Visibility) -

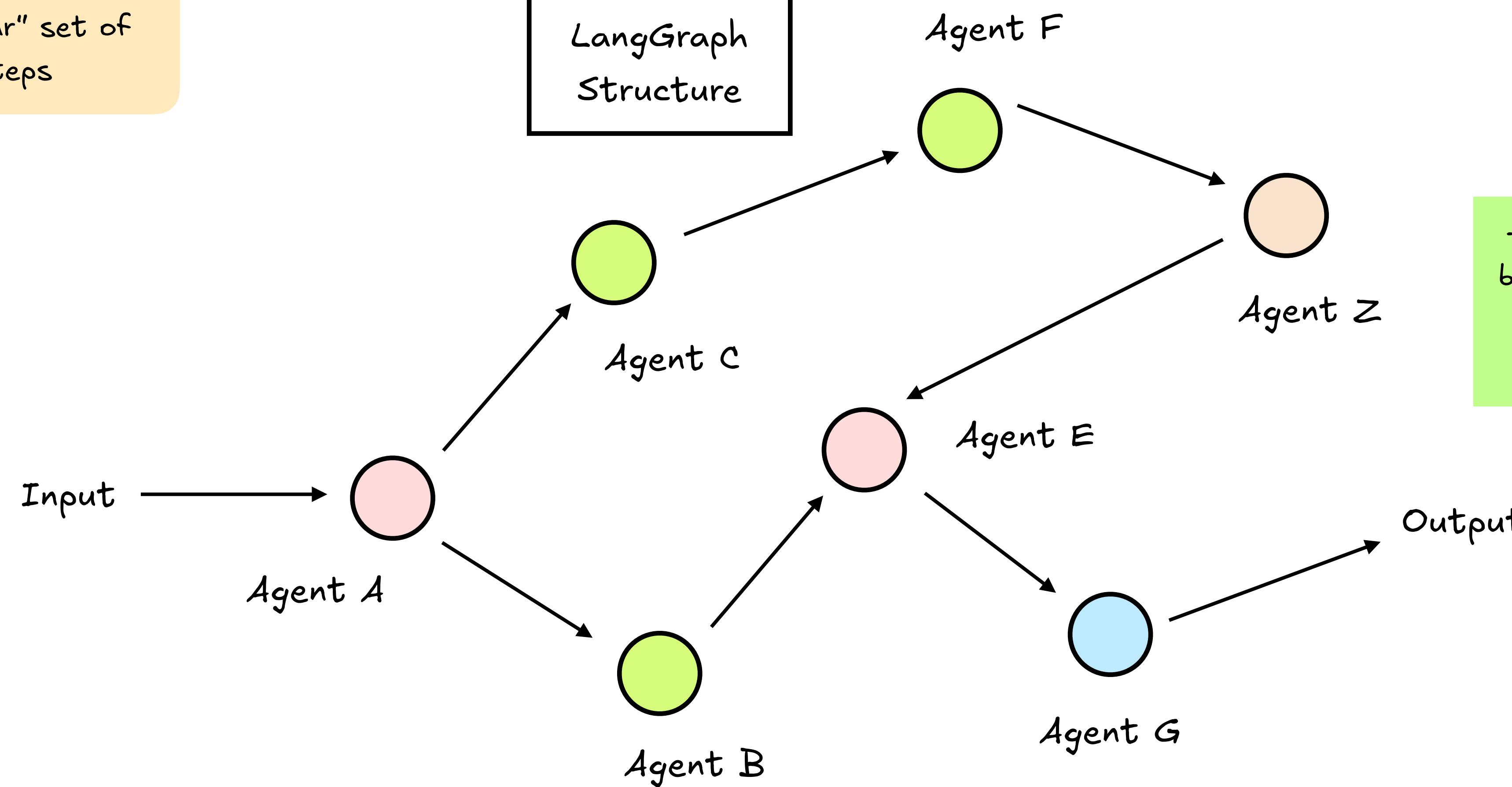


# Design Principle: Uniformity of Information Flow



# Design Principle: Granularity of Actions

- (Granularity) -  
Large tasks are  
broken into a  
"granular" set of  
steps



- All agents have the same basic structure so they can be easily added, replaced, subdivided as needed

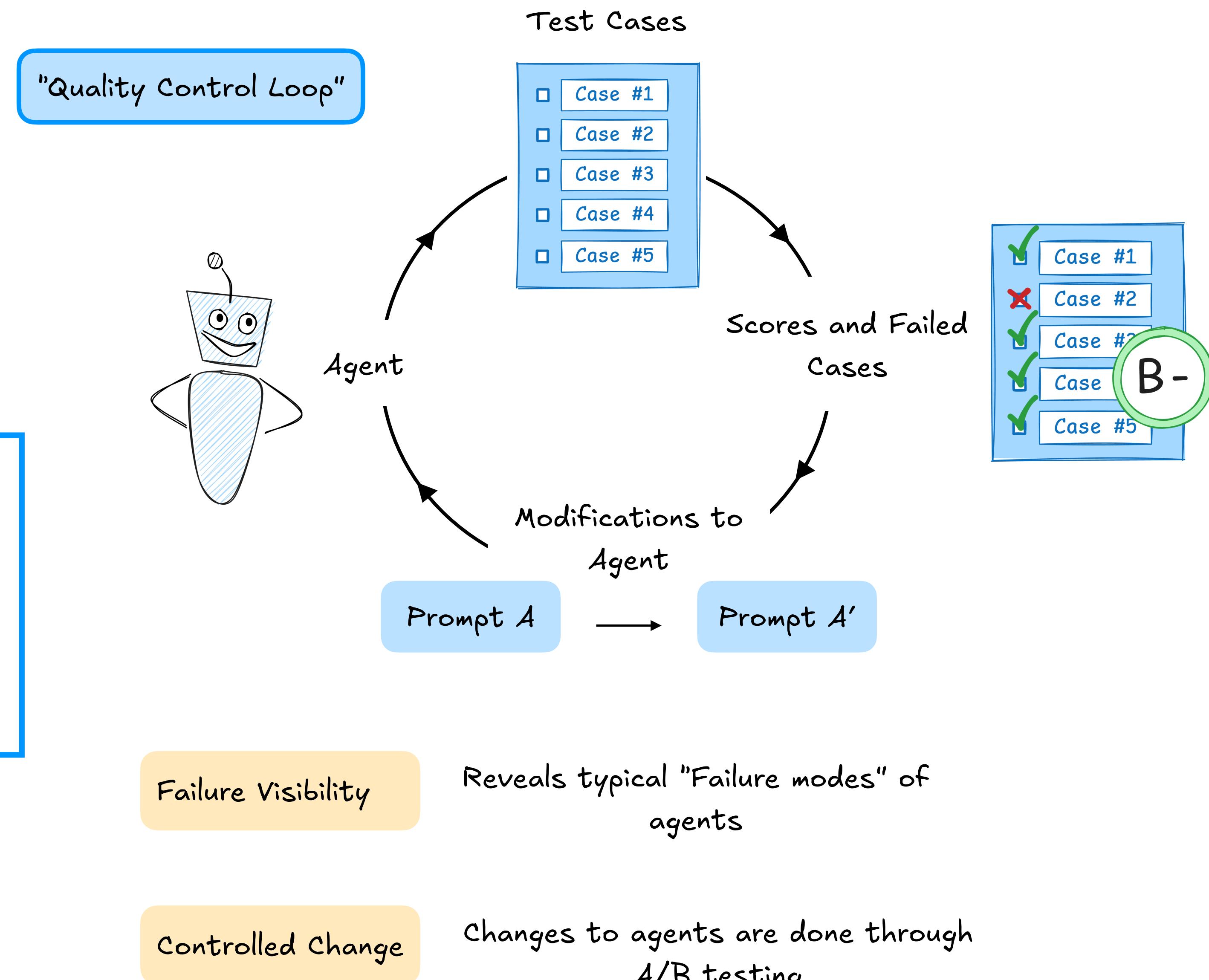
# Design Principle: Visibility of Performance

- (Visibility) -  
Test cases are created for agents to understand failure modes AND effects of changes

\*How often the same tests pass or fail

Num Test Cases	Pass Rate	Consistency Prob
Agent A	20	80%
Agent A'	20	85%

(  
gives us two main things →



# Presentation Roadmap

.....

I. Motivation

II. Video Demo

III. System Details

- .i Architecture

- .ii Agents and LangGraph

- .iii Design principles

IV. Gaps and Next Steps

V. Live Demo & Q&A

what the system  
cannot do now

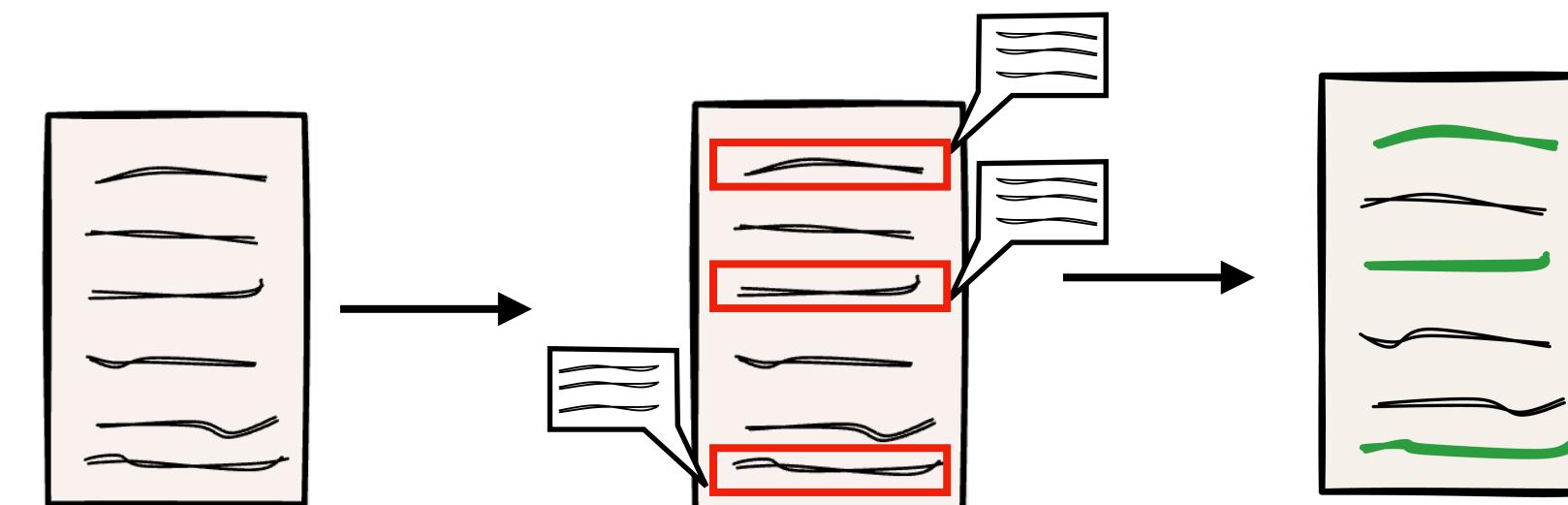
&

What we are planning  
to build in the future

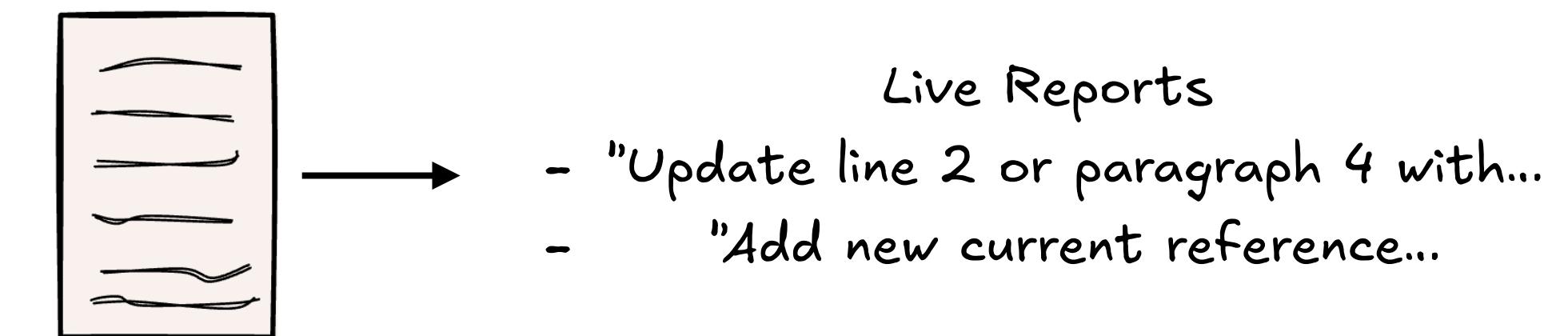
# Gaps and Next Steps

What the system  
cannot do now  
&  
What we are planning  
to build in the future

Modify text "in-line"  
without user copy-pasting



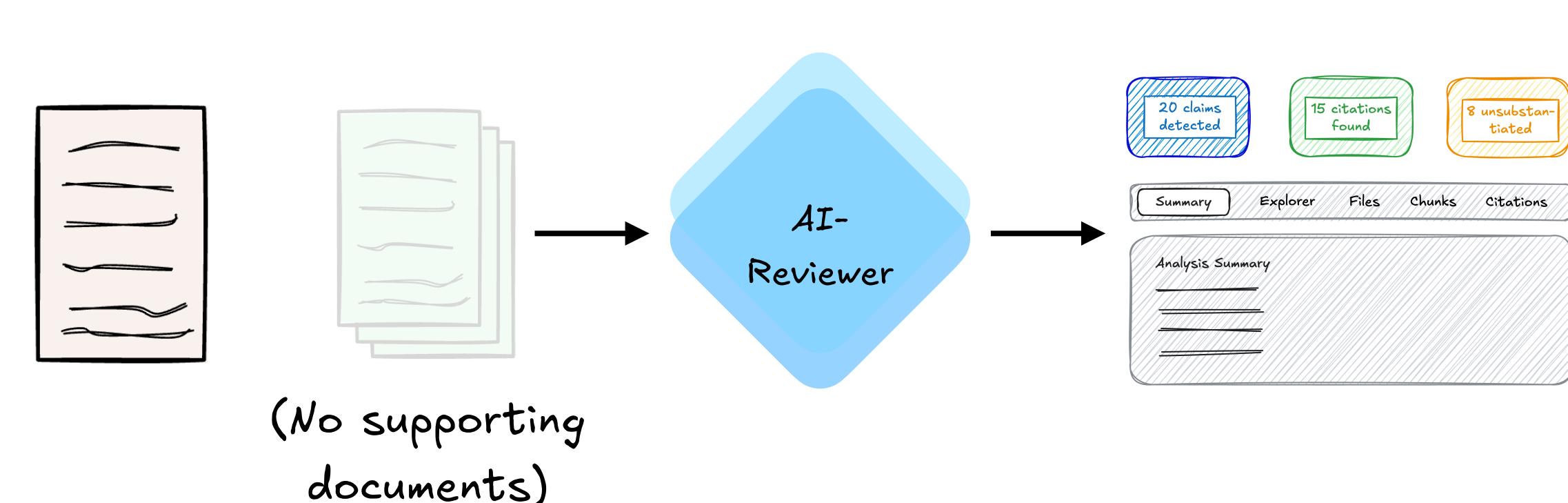
Live Reports Addendum



Higher Level Coherence

"The overall argument of the paper  
is....there are gaps in this argument  
from...."

Conduct PDF/paper  
searchers to find references  
without user input



# Links before Live Demo

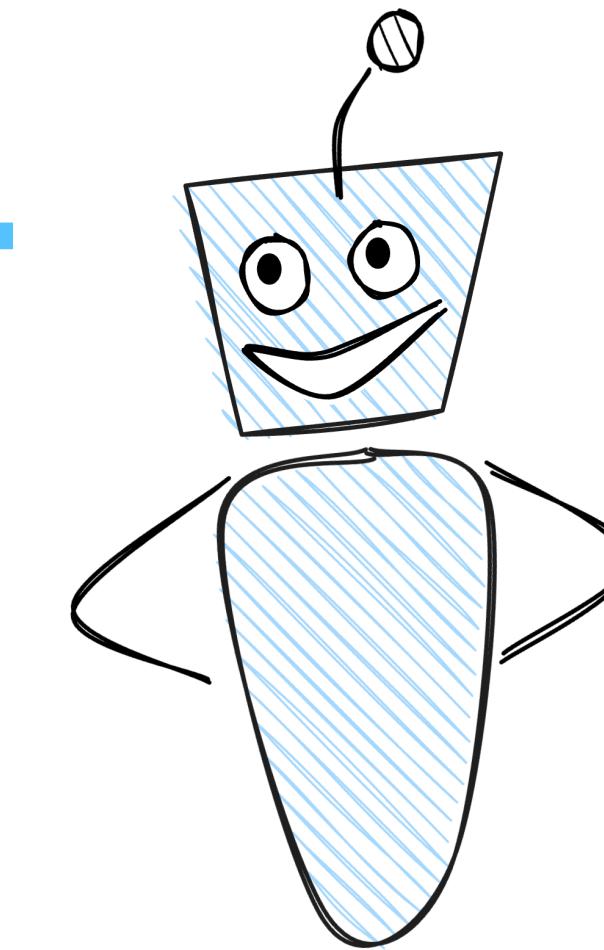


IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

(<https://xkcd.com/1425/>)

## Relevant Links

- AE Studio
  - (<https://ae.studio/>)
- AI Reviewer Github
  - (<https://github.com/agencyenterprise/ai-reviewer>)
- Slides
  - (Linked within repo)



# Presentation Roadmap

.....

I. Motivation

II. Video Demo

III. System Details

- .i Architecture

- .ii Agents and LangGraph

- .iii Design principles

IV. Gaps and Next Steps

V. Live Demo & Q&A

<redacted>

## Live Demo

---

<redacted>

## Extra Slides

