

Conf42 Large Language Models (LLMs) 2026

# Building Reliable LLM Systems

Managing Decision Latency and Cognitive Load  
in the AI Era

---

**Lakshmi Priya Gopalsamy**

Independent Researcher & Technology Lead, Software Engineering

IEEE Senior Member · Sigma Xi Full Member · ACM Professional Member

# The bottleneck has shifted.

Your LLM generates code in **4 seconds**.  
Your approval loop takes  
**4 days**.

## AI Execution Speed

Code generation

~4s

Test generation

~6s

Doc generation

~3s

Analysis

~8s

## Human Decision Loop

Model approval

2–4 days

Prompt review

3–5 days

Data governance check

5–10 days

Safety validation

1–2 weeks

# A reframe.

*Reliability and governance are increasingly limited not by model capability — but by decision latency and cognitive load.*

## DECISION LATENCY

The elapsed time between an AI output being ready and a human making a consequential decision about it.

### Found in:

- Model approvals
- Prompt reviews
- Data governance checks
- Safety validations
- Incident response loops

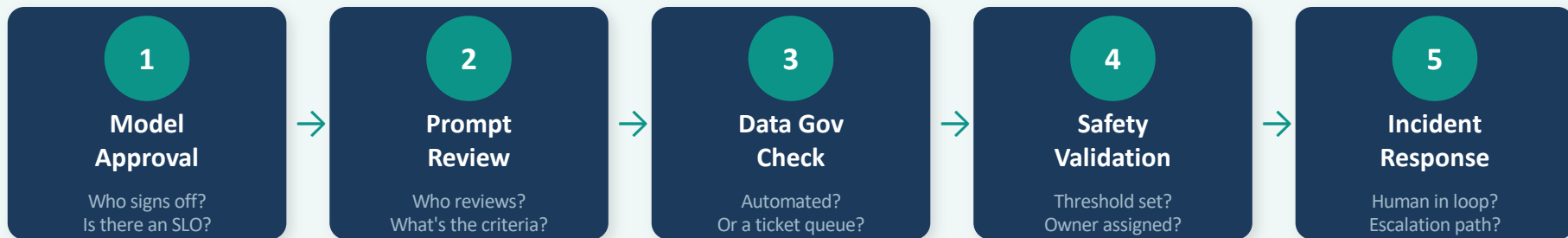
## COGNITIVE LOAD

The accumulated mental overhead engineers carry when managing LLM systems at scale and velocity.

### Found in:

- Model selection decisions
- Evaluation metric design
- Hallucination mitigation
- Prompt engineering
- Security + compliance management

# Where decision latency lives.



## The invisible queue

Each step above is a queue. Every queue has latency. Most organizations have never measured it. AI initiatives stall not because the model failed — but because nobody defined who owns the approval decision or what the SLO is for the human in the loop.

## Pattern from enterprise practice:

When we shifted from centralized enforcement to self-service governance with automated guardrails, teams that previously waited 5–10 days for governance approvals were unblocked. The latency didn't disappear — it moved into the system, where it could be measured, tuned, and owned.

# Where cognitive load lives.

The stack engineers are managing simultaneously in LLM-enabled systems:

1

## Model Selection

Which model for which task?  
How to evaluate?

2

## Evaluation Metrics

What does 'good' mean? Who decides?

3

## Hallucination Mitigation

Detection, grounding, human-in-loop design

4

## Prompt Engineering

Versioning, testing, regression detection

5

## Security Boundaries

Prompt injection, data leakage, access control

6

## Compliance Requirements

Data residency, audit trails, model documentation

Each layer is manageable in isolation. All six simultaneously, at the pace LLMs enable, without deliberate design — leads to burnout and fragile systems.

# Apply the distributed systems lens.

## Traditional Reliability Thinking

Machines fail

- Design for redundancy
- Measure: error rate, latency, availability
- SLOs on system components



## LLM Reliability Thinking

Humans bottleneck

- Design for decision throughput
- Measure: decision latency, cognitive load
- SLOs on human loops too

## Reapply the vocabulary you already know:

### Error Budgets

How many human review failures per sprint before the AI initiative slows down?

### Toil

Manual governance steps that scale linearly with AI adoption — automate them.

### SLOs on Human Loops

P95 for decision turnaround. If you don't measure it, you can't improve it.

# Explicit Decision Boundaries

Not every AI output needs a human decision. Define which ones do — and at what level.

## PRE-AUTHORIZED

### No review needed

Low-stakes outputs · Reversible decisions ·  
Well-tested prompt templates · Internal  
tooling suggestions

Target latency: < 1 second

## ASYNC REVIEW

### Review within SLO

Medium-stakes outputs · Customer-facing  
but correctable · New prompt patterns ·  
Models within approved families

Target latency: < 4 hours

## EXPLICIT APPROVAL

### Synchronous gate

Irreversible decisions · High-risk outputs ·  
New model families · Data classification  
changes · Public-facing AI

Target latency: Defined SLO

*Analogy: circuit breakers in distributed systems — known failure mode triggers known response, no human required in the hot path.*

# Clear AI Platform Ownership Models

*Every LLM in production needs an owner. Not a committee — a team with a pager and an SLO.*

## The common anti-pattern

- AI initiative starts with a cross-functional 'AI committee'
- No single team owns reliability for any model in production
- When something breaks: who does the incident go to?
- Answer: everyone, which means no one
- Result: accumulated hidden risk, stalled response

## The ownership model

- Each LLM has an owning team with explicit accountability
- SLO defined: availability, response quality, decision latency
- On-call rotation includes AI system events
- Architecture decisions require owner sign-off
- Ownership documented, reviewed quarterly

**Ambiguous ownership is how AI initiatives accumulate hidden risk.**

### Decision Latency

When the owner is clear, approval loops have a destination and SLO.

### Cognitive Load

Ownership boundaries define who carries which concerns.

### Incident Confusion

No ambiguity about who gets paged and what they're responsible for.

# Policy-as-Code Guardrails

Move governance decisions out of human review loops and into automated enforcement.

## Before: Manual Review Loop

AI Output

→ Slack thread

→ **Human review (days)**

→ Back to team

→ Maybe approved



## After: Policy in the Pipeline

AI Output

→ Prompt boundary check (ms)

→ Data classification rule (ms)

→ Output safety threshold (ms)

→ **Automated decision + audit log**

## What to encode as policy:

### Prompt Boundary Checks

Input length, injection pattern detection, topic restriction

### Output Safety Thresholds

Confidence floors, topic filters, format validators

### Data Classification Rules

PII detection, data residency enforcement, access scope

### Audit Trail Generation

Every decision logged with model version, timestamp, policy version

# Opinionated Golden Paths for LLM Adoption

*Reduce cognitive load by making the right choice the easy choice.*

## What a Golden Path Encodes

- Approved model catalog with version pinning
- Standard evaluation framework and metrics
- Pre-built prompt templates with safety checks
- Security boundary configuration by data class
- Observability hooks + alerting defaults
- Compliance documentation templates
- Escalation paths and on-call runbooks

## The adoption incentive:

If the golden path is faster than building from scratch, adoption isn't a culture problem — it's math.

## Teams can still diverge:

Golden paths aren't lock-in. They're sensible defaults. Teams opt out with explicit justification — which creates a useful forcing function for architectural review.

*Analogy: a Unified Compute Gateway that abstracts infrastructure complexity — teams submit jobs without knowing the underlying machinery. Same principle: abstract the hard choices into validated defaults.*

# The compounding effect.

Apply all four patterns together. The outcome is not additive — it's multiplicative.



## Sustainable AI Velocity

Fast enough to compete. Reliable enough to trust. Governed well enough to sustain.

### ↓ Decision latency

From days to hours or minutes

### ↓ Cognitive load

Fewer decisions per engineer per sprint

### ↑ Reliability

Fewer incidents from governance gaps

### ↑ Team retention

Engineers who aren't burned out ship better systems

# What to do Monday.

Four actions. No new budget required.

01

## Map your decision loops

Draw every point where a human waits on an AI output — or an AI waits on a human. Measure current latency. You cannot improve what you haven't made visible.

02

## Assign an owner to every LLM in production

Name. Team. SLO. Incident escalation path. If you can't answer those four things for any model currently running, that's your highest-priority action.

03

## Identify one governance decision to automate this sprint

Pick the highest-volume manual review step. Write it as policy. Deploy it. Measure the latency delta. Build from there.

04

## Establish a golden path for your most common LLM use case

Document what 'good' looks like: approved models, evaluation criteria, security boundaries, observability hooks. Make it easier to adopt than to ignore.

The teams that will build reliable LLM systems  
at scale are the ones who treat

**human decision capacity**

as a design constraint — not an afterthought.

---

[linkedin.com/in/lakshmipriyagopalsamy](https://www.linkedin.com/in/lakshmipriyagopalsamy)

[lakshmipriyagopalsamy.com](https://lakshmipriyagopalsamy.com)

**Questions & discussion welcome.**

Let's continue the conversation.