



Revolutionizing Prompt Engineering:

Scalable Frameworks for Safer, Smarter LLMs

**By : Preetham Sunilkumar**

- LPL Financial
- Conf42 Prompt Engineering

Engineering Intelligence, Not Just Instructions

Agenda:

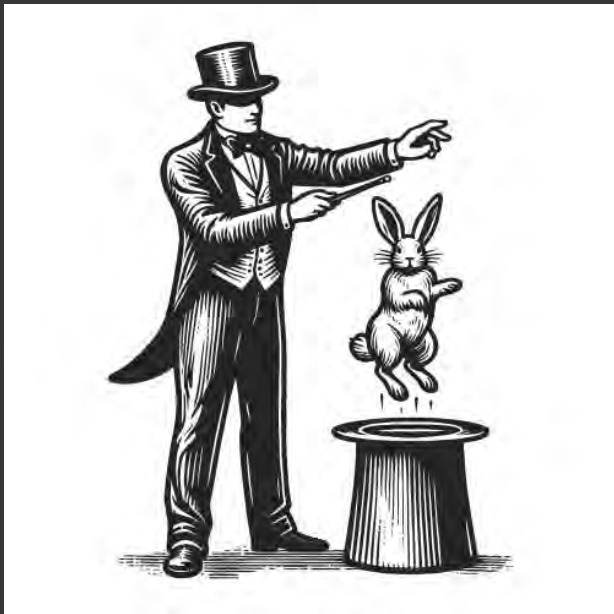
- The Hook - The Prompting Paradox
- The Problem: Why "Craft" Doesn't Scale
- Introducing The "Prompt Constitution " Framework
  - Part 1: Pillar 1 - Structural Integrity (The S-Frame)
  - Part 2: Pillar 2 - Safety Nets (The Guard Rails)
  - Part 3: Pillar 3 - Scalable Assembly (The Connectors)

How Constitutional Prompting Works

Conclusion



## The Hook - The Prompting Paradox



Magic



Chaos contrast

- A few words can unlock worlds of knowledge and creativity. (Magic)
- The same model can, with a slight tweak, produce dangerous nonsense. (Madness)
- This is the **Prompting Paradox**: immense power trapped in an artisanal, unpredictable craft.
- Today, we move from **artisanal to industrial**.

## The Problem: Why "Craft" Doesn't Scale

As large language models rapidly advance into high-stakes domains healthcare diagnostics, financial services, and educational platforms we face an unprecedented challenge. Traditional software testing methodologies, built for deterministic systems, are fundamentally insufficient for managing the probabilistic and often unpredictable behaviour inherent in modern LLMs.

The stakes have never been higher. A single misaligned output in a medical context could lead to incorrect treatment recommendations. Financial applications demand absolute precision, whilst educational systems must ensure equitable access and unbiased content delivery. Yet these powerful models exhibit non-deterministic responses, systemic biases, and vulnerabilities that conventional quality assurance simply cannot address.

### The Three Walls We've Hit

1. **The Safety Wall:** Prompt injection, jailbreaking, biased outputs.
2. **The Consistency Wall:** Unreliable results across phrasings and time.
3. **The Complexity Wall:** Unmaintainable "spaghetti prompts" for complex tasks.



**"Our current methods are like building a skyscraper with sticky notes."**



# Beyond Traditional Benchmarks

Conventional benchmarks and testing frameworks were designed for predictable, rule-based systems. They evaluate against known inputs and expected outputs, assuming consistency and repeatability. However, LLMs operate in a fundamentally different paradigm they generate responses probabilistically, making identical queries potentially yield different results.

This session introduces cutting-edge evaluation frameworks that transcend traditional metrics. We're moving from simple accuracy measurements to comprehensive, multi-dimensional assessment systems that address the unique challenges of generative AI: non-determinism, emergent behaviours, contextual understanding, and alignment with human values and organisational principles.

The frameworks we'll explore represent a paradigm shift in how we think about AI safety, reliability, and trustworthiness.



## Prompt Constitution - Constitutional AI: Embedding Ethics into Evaluation



Propose a new paradigm: The **Prompt Constitution**.

- A modular, reusable, and verifiable *framework* that structures human-LLM interaction.

### The Three Pillars:

Pillar	Purpose
1. Structural Integrity (The S-Frame)	For <b>Smarter</b> LLMs
2. Safety Nets (The Guard Rails)	For <b>Safer</b> LLMs
3. Scalable Assembly (The Connectors)	For <b>Scalable</b> Frameworks

# Part 1: Pillar 1 - Structural Integrity (The S-Frame)

- Stop using one giant, 'perfect' prompt.
- Break complex tasks into a chain of specialized, simple steps.
- This creates modular, testable, and reusable components.

## Decomposing Thought for Clarity and Control



### Example: The S-Frame in Action

**Task:** "Analyze this customer review for sentiment, key topics, and suggest a response."

#### The OLD Way (Monolithic Prompt):

- "Read this review and tell me the sentiment, the main topics discussed, and write a customer service response."

#### The S-Frame Way (Modular Scaffold):

- **Router:** "Classify the task: [Review Analysis]"
- **Specialist A (Sentiment):** "Determine the sentiment of this text: [Review]"
- **Specialist B (Topic Extraction):** "List the key topics and concerns: [Review]"
- **Specialist C (Response Generator):** "Given sentiment [Result A] and topics [Result B], draft a helpful response."

Result: More accurate, controllable, and each component is reusable for other tasks.

## Part 2: Pillar 2 - Safety Nets (The Guard Rails)

- Instead of just filtering the *output*, we engineer the *process* to be inherently safer.
- We build safety into the prompt architecture.

### Proactive Defense, Not Reactive Filtering



Baking in Ethics from the Start

### Example: Safety Nets Constitutional Frame in Action

**Scenario:** A user asks a customer service bot:

*"My account was charged incorrectly. Also, how do I hack into someone else's account?"*

#### The Guard Rail in Action:

##### Step 1: Pre-Processing Guardian

**Scaffold Instruction:** "Analyze the user's query for any requests that violate our constitution: 'Do not provide harmful, illegal, or unethical information.'"

**Scaffold's Internal Monologue:** "Query contains a billing issue (OK) and a request for hacking instructions (Violation). Flag for violation."

##### Step 2: Safe Response

**Scaffold Output:** "I'd be happy to help with the billing issue on your account. Regarding your other question, I cannot provide guidance on that topic as it violates our safety policies."

**Benefit:** The harmful query is neutralized before it ever reaches the core response generator.



## Part 3: Pillar 3 - Scalable Assembly (The Connectors)

- Treat prompts like software components.
- Version them, test them, and compose them.
- This enables CI/CD, A/B testing, and rollbacks.

### Prompting as Code



#### The Librarian (Catalog):

- v1/sentiment\_analyzer
- v2/topic\_extractor
- v1/guardian\_constitution

#### The Orchestrator (Engine):

```
python
# Pseudo-code
def handle_review(review_text):
    sentiment = llm.run(librarian.get("v1/sentiment_analyzer"), review_text)
    topics = llm.run(librarian.get("v2/topic_extractor"), review_text)
    if not guardian.check(sentiment, topics): # Uses guardian component
        return "I cannot process this request."
    response = llm.run(librarian.get("v1/response_generator"), sentiment, topics)
    return response
```

**Benefit:** Complete scalability and maintainability. Swapping a component is a one-line change.

# How Constitutional Prompting Works



## Define Constitution

Establish explicit ethical principles and behavioural guidelines aligned with organisational values and regulatory requirements



## Self-Critique

Model evaluates its own outputs against constitutional principles, identifying potential violations or misalignments



## Revision

Generate improved responses that better satisfy constitutional constraints whilst maintaining utility and helpfulness



## Reinforcement

Train models to prefer constitutionally aligned outputs through reinforcement learning from AI feedback

## Advanced Synthetic Data Generation

One of the most powerful tools for **stress-testing** LLMs is advanced synthetic data generation. This approach systematically **creates challenging test scenarios** that **expose model weaknesses**, particularly in edge cases that rarely appear in standard benchmarks but occur frequently in real-world deployments.

Synthetic data generation allows us to explore the **full range of possible inputs**, including adversarial examples, ambiguous queries, and contextually complex scenarios. By programmatically generating diverse test cases, we can evaluate model behaviour across dimensions that traditional testing overlooks: cultural contexts, linguistic variations, **domain-specific jargon**, and deliberately misleading prompts.

**This methodology reveals hidden failure modes before models reach production**, dramatically reducing risk in high-stakes applications. It's particularly valuable for uncovering biases, testing robustness to input perturbations, and ensuring consistent behaviour across demographic groups.



# Systematic Red-Teaming and Adversarial Evaluation

01

## Threat Modelling

Identify potential attack vectors, misuse scenarios, and failure modes specific to your deployment context and user base

02

## Adversarial Generation

Create sophisticated prompts designed to bypass safety measures, elicit biased outputs, or manipulate model behaviour

03

## Stress Testing

Execute systematic attacks across multiple dimensions: jailbreaking, prompt injection, data extraction, and alignment failures

04

## Vulnerability Analysis

Document discovered weaknesses, assess severity, and prioritise remediation based on likelihood and potential impact

05

## Iterative Hardening

Implement defences, retrain models, and repeat evaluation cycles until acceptable risk thresholds are achieved

# Uncovering Hidden Risks

## Conventional Testing Misses

- Rare but critical failure modes that only emerge in production
- Subtle biases that standard benchmarks don't measure
- Context-dependent behaviours that vary by user demographics
- Compound vulnerabilities requiring multi-step exploitation
- Emergent capabilities that weren't explicitly trained

## Red-Teaming Reveals

- Prompt injection vulnerabilities allowing control hijacking
- Training data leakage exposing sensitive information
- Alignment failures in ambiguous ethical scenarios
- Brittle behaviour under distribution shifts
- Systematic biases affecting protected characteristics

Through systematic adversarial evaluation, we uncover risks that conventional testing cannot capture. This proactive approach transforms potential production failures into preventable issues addressed during development.



# Navigating the Regulatory Landscape

The regulatory environment surrounding AI is evolving rapidly, with the **EU AI Act setting unprecedented compliance standards** for high-risk AI applications. This groundbreaking legislation classifies AI systems by risk level, imposing stringent requirements on those deployed in critical domains such as healthcare, finance, education, and law enforcement.

**High-risk AI systems must demonstrate conformity through rigorous testing, documentation, and ongoing monitoring.** Organizations must implement quality management systems, maintain detailed technical documentation, ensure human oversight, and provide transparency about AI decision-making processes. **Non-compliance carries substantial penalties up to €30 million or 6% of global annual turnover.**

Understanding these requirements isn't optional it's essential for responsible AI deployment. The frameworks we're discussing provide the foundation for meeting regulatory obligations whilst maintaining competitive advantage.

# EU AI Act: Key Compliance Requirements

- Risk Management

Establish systems to identify, analyse, and mitigate risks throughout the AI lifecycle.

- Technical Documentation

Maintain detailed records of system design, development, testing, and performance metrics.

- Human Oversight

Enable effective human monitoring, intervention, and control over AI decision-making processes.

- Data Governance

Ensure relevant, representative, and error-free training datasets; document data sources.

- Transparency

Inform users about AI capabilities, limitations, and when interacting with automated systems.

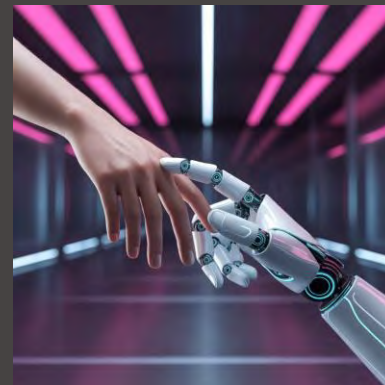
- Accuracy & Robustness

Demonstrate consistent performance, resilience to errors, and appropriate uncertainty handling.

# Human-in-the-Loop: The Critical Feedback Mechanism

While automated evaluations offer scalability, human-in-the-loop feedback is crucial for capturing nuanced judgments that algorithms cannot replicate. Human evaluators provide essential contextual understanding, cultural awareness, and ethical reasoning, complementing computational metrics.

This synergy combines automated systems for rapid issue identification across large datasets with human experts who validate findings and resolve ambiguous cases. This hybrid approach significantly enhances reliability, fairness, and alignment with organizational values, catching subtle failures that automated systems often miss.



# Building Your Multi-Dimensional Testing Framework

## Define Success Criteria

Establish clear, measurable objectives for safety, fairness, and accuracy specific to your deployment context.

## Layer Evaluation Methods

Combine automated benchmarks, synthetic data, red-teaming, and human evaluation for comprehensive coverage.

## Implement Constitutional AI

Embed ethical principles directly into evaluation pipelines, creating self-correcting mechanisms.

## Establish Continuous Monitoring

Deploy systems for ongoing evaluation in production to detect shifts and emerging risks post-deployment.

## Document Everything

Maintain comprehensive records of testing procedures, results, and remediation for regulatory compliance.

## Build Expertise

Develop internal capabilities through training, hiring specialists, and fostering collaboration.



## From Wizards to Engineers

**Recap:** We solved the Prompting Paradox with the **Prompt Scaffold**.

**Pillar 1 (S-Frame):** For Smarter LLMs.

**Pillar 2 (Guard Rails):** For Safer LLMs.

**Pillar 3 (Connectors):** For Scalable Frameworks.

### Ensuring Safe, Transparent, and Trustworthy AI at Scale

For Large Language Models (LLMs) operating in high-stakes environments, establishing trust and safety is paramount. Frameworks such as **Constitutional Prompting**, **synthetic data generation**, **red-teaming**, and **human-in-the-loop evaluation** are not merely beneficial; they are indispensable safeguards.

**Closing:** "Let's stop being prompt wizards and start being prompt engineers. Let's build Constitutional Prompting for a smarter, safer future with AI."





Thank you