

Article

Not peer-reviewed version

Prompt–Centric Observability: Debugging and Securing Generative AI Pipelines in Enterprise Deployments

[Manaswini Bollikonda](#)*

Posted Date: 14 November 2025

doi: 10.20944/preprints202511.0008.v1

Keywords: observability; prompts; governance; LLM operations; security; RAG; telemetry; validation; enterprises



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Prompt—Centric Observability: Debugging and Securing Generative AI Pipelines in Enterprise Deployments

Manaswini Bollikonda

Independent Researcher, USA; manaswini.bollikonda@gmail.com

Abstract

Enterprises rarely fail because their language model is incapable; they fail because the prompt pathway is opaque. We introduce *prompt-centric observability* (PCO), an operational discipline that treats inputs, routing, retrieval, templating, safety overlays, generation, and validation as first-class, measurable surfaces. PCO emits low-cardinality signals (coverage, support, freshness, p95, spend), binds drafts to the exact evidence spans they used, and routes decisions through a governed gate (release, rewrite, redact, escalate). We present a compact architecture, a minimal telemetry hook, and actionable controls for entitlements, privacy, and auditability. The goal is not a new leaderboard but predictable behavior under change: when data, policies, or templates shift, the system responds within guardrails and every outcome is replayable.

Keywords: observability; prompts; governance; LLM operations; security; RAG; telemetry; validation; enterprises

1. Introduction

Generative systems are persuasive by design, which is precisely why observability is non-negotiable. In practice, most user-facing failures trace to the prompt pathway: a stale chunk slips into context, a template version drifts, or a safety overlay quietly stops applying. Traditional request/response logs cannot explain such incidents because they omit the mid-stream artifacts that shaped the answer.

We define prompt-centric observability (PCO) as the habit of capturing typed artifacts at each step—candidate lists, template versions, guard outcomes—and fusing them into shallow, stable metrics. The point is to make debugging a matter of comparing two prompt bundles rather than eyeballing token dumps. With typed links in place, “seems wrong” becomes a reproducible diff.

A second motivation is trust. Compliance reviewers and on-call engineers need to replay what the system saw, not merely what it said. If every draft is tied to specific chunk IDs, timestamps, and template versions, attribution becomes mechanical: who saw what, when, and why it was released. This supports both internal blameless postmortems and external audits without bespoke forensics.

We also care about speed. Observability has to be cheap enough for the hot path. PCO favors low-cardinality counters and bounded payloads: coverage (unique used vs. offered), support (used chunks with cited spans), freshness (age vs. policy), p95 latency, and token spend. These signals are sufficient to triage 80% of issues before deeper tracing is needed.

Beyond debugging, PCO functions as an organizational contract. Product, security, compliance, and SRE can agree on a small set of stable signals and actions, rather than negotiating over ad-hoc logs after an incident. This narrows ambiguity: if coverage or support drops below the agreed window, the gate must either rewrite or escalate—no exceptions hidden in one-off dashboards or tribal knowledge.

A useful mental model is “prompts as build artifacts.” Like binaries compiled from source, prompts are outputs assembled from retrieval, templates, and overlays. We therefore apply build-like guarantees: deterministic inputs, versioned specs, and repeatable outputs. When a regression appears, the question becomes “which artifact changed?” rather than “what did the model think?”

Finally, PCO reduces cognitive load for on-call engineers, organizational guesswork. When signals are versioned and visible, product and security can negotiate thresholds. Shadow policies can be evaluated safely, and promotions become data-backed rather than anecdotal. Most importantly, changes to data or policy become deliberate, reversible operations. Instead of scanning token walls, they compare compact records: candidate IDs, template version, guard outcomes, and the gate's decision trail. This shrinks time-to-reason (and incident MTTR) because the evidence that mattered is already summarized, and every field is designed to be actionable [1].

2. Background and Motivation

Early deployments treated prompts as throwaway strings, glued together from retrieval results and a few guardrails. As corpora grew and policies tightened, teams discovered that prompts are infrastructure: they deserve versioning, contracts, and health checks just like any API. This section motivates that shift and the operational wins it unlocks [2].

First, retrieval quality is unstable in the wild. Indexes age, embeddings drift, and near-duplicates proliferate. Without typed candidate sets and coverage metrics, regressions surface only as user complaints. PCO keeps retrieval explainable by logging which candidates were eligible, which ones were selected, and why.

Second, templates accumulate responsibilities: tone, disclaimers, formatting, and citation scaffolds. Minor edits can inadvertently remove a legal phrase or tighten token budgets enough to prune crucial evidence. Versioned template IDs and prompt quotas make such regressions visible and reversible.

Third, security is end-to-end. Entitlements from the source system must flow into retrieval filters and remain intact through prompting—never “mask after.” PCO binds answers to both evidence and permissions, so privilege violations are caught at the gate and can be audited later.

Fourth, governance is a control surface, not an afterthought. Release, rewrite, redact, and escalate are not vibes; they are functions of explicit signals: support, coverage, freshness, and policy risk. Because PCO keeps these numeric and typed, teams can tune thresholds without touching model code [3].

Lastly, the human layer benefits. When dashboards show coverage falling or freshness skewing, content owners can prioritize fixes; support teams can pinpoint brittle pages; legal can validate that disclaimers remain intact. PCO turns prompting from craft into engineering.

A recurring anti-pattern is “prompt sprawl,” where many teams copy and tweak templates without provenance. Over months, minor differences in disclaimers, tone, and quotas accumulate into inconsistent behavior and unpredictable risk. PCO curbs sprawl by pinning templates, exposing diffs, and tying every release to a specific version, which makes consolidation a mechanical task instead of a cross-team negotiation.

Another common failure is “mask after,” where sensitive spans are filtered post-generation. This fails both privacy and explainability: the model already saw the text, and the audit trail cannot prove that access controls were respected. By contrast, PCO pushes entitlements into retrieval and prompt assembly, ensuring sensitive content never enters the model's context [4].

Motivation also comes from the economics of scale. As traffic grows, small inefficiencies—overlong prompts, redundant evidence, heavy re-ranking—compound into material spend. Observability that tracks spend alongside support and coverage lets teams prioritize optimization where it preserves quality instead of blindly trimming tokens.

3. PCO Architecture

PCO organizes the pipeline into replaceable stages that emit typed artifacts: `InputRoute`, `CandidateSet`, `TemplateSpec`, `PromptBundle`, `AnswerTrace`, and `ValidationReport`. Each type is versioned so components evolve independently.

Figure 1 shows the single-column PCO flow. Blocks are stages; arrows are typed artifacts; feedback keeps templates and retrieval aligned with what validation and governance discover.

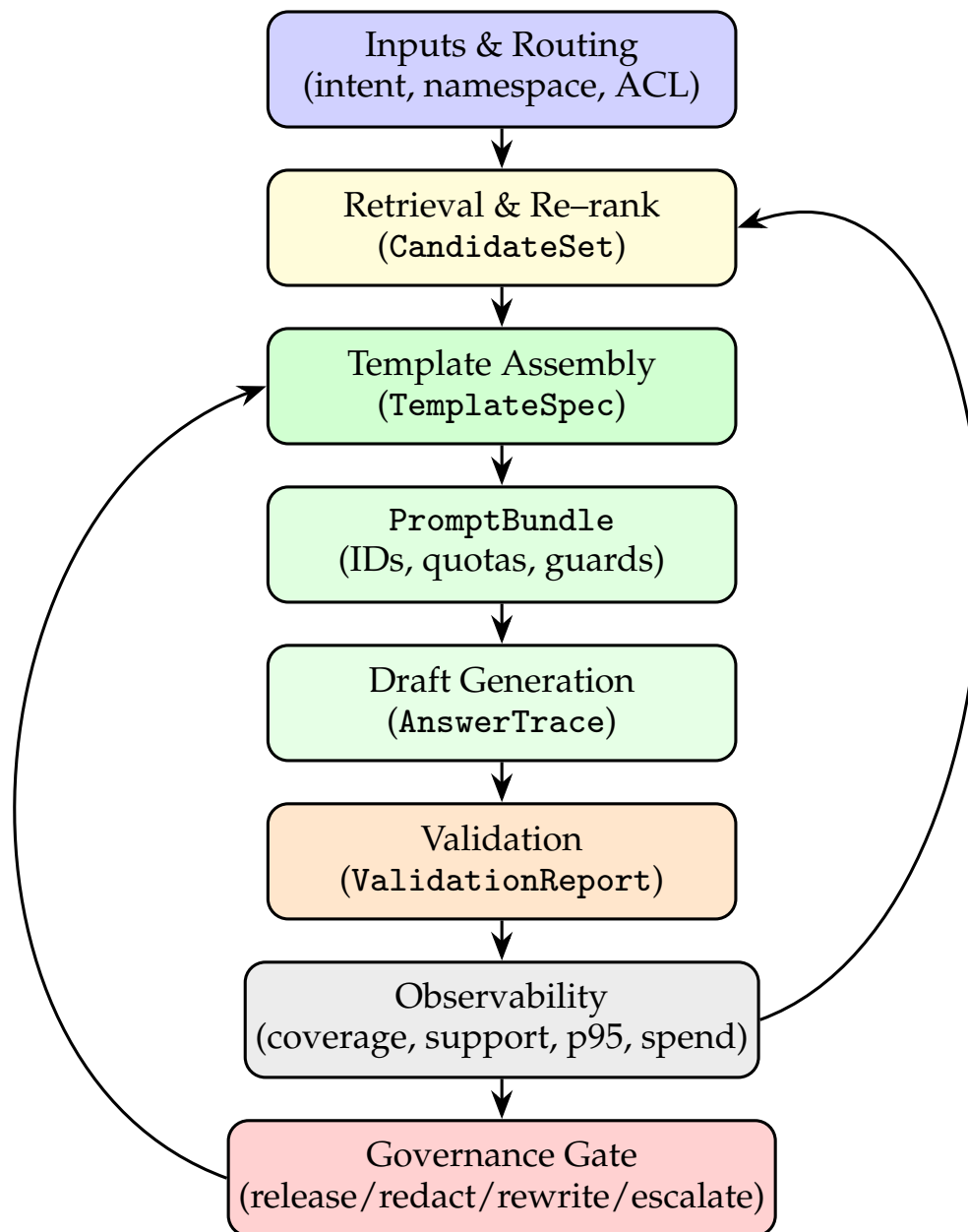


Figure 1. PCO pipeline: typed artifacts at each stage, shallow metrics for SLOs, and a governed decision ladder.

Inputs are normalized into intents with namespace and ACL context. Retrieval produces a `CandidateSet` that is explainable: each item carries base scores, freshness, and entitlement flags. Re-ranking adds diversity and trims near-duplicates so prompts cover distinct facets rather than repeating one page five ways.

Template assembly is treated as code. The `TemplateSpec` fixes instruction text, safety overlays, and quota budgets for evidence and exemplars. The resulting `PromptBundle` references chunk IDs (not just text), ensuring the validator can later check claims against specific spans.

Generation yields a draft answer plus an `AnswerTrace`: used chunk IDs, model version, token counts, and latency. The validator then computes `ValidationReport` fields: citation support, entity/number consistency, freshness vs. policy windows, and policy language checks [5].

Observability fuses these artifacts into low-cardinality signals for SLOs. Governance is a small decision ladder: release if support and coverage clear thresholds and risks are low; otherwise rewrite/redact/escalate with structured reasons. Because the gate reasons over typed signals, teams can change thresholds safely and audit decisions later.

Typed artifacts pay off only if they are stable under change. We adopt forward-compatible schemas with explicit “unknown” fields that can be ignored safely by older services, and we version each type independently (e.g., `TemplateSpec@v4`). This lets retrieval, prompting, validation, and governance evolve at different cadences without breaking traffic.

Multi-tenant deployments add constraints. Namespaces encapsulate policy (e.g., retention, freshness windows, PII rules) and route both retrieval and templates. The same pipeline can therefore serve HR, Legal, and Support with different knobs while sharing infrastructure. Crucially, namespace identity flows into every artifact so auditors can reconstruct which policy governed a given release [6].

Finally, we treat feedback as a first-class edge. Governance outcomes (rewrite, redact, escalate) are not terminal; they carry hints that feed back to retrieval and template assembly—e.g., “insufficient coverage on exceptions,” “missing policy disclaimer,” or “stale evidence.” Over time, these hints drive targeted fixes: re-ranker adjustments, new micro-templates, or content refactors.

4. Telemetry Signals and Trace Design

Telemetry has to be cheap, stable, and specific. We standardize on five signals: coverage (unique used/unique offered), support (fraction of used chunks with cited spans), freshness (age of oldest/newest cited spans vs. namespace policy), latency (p50/p95), and spend (tokens in/out). These cover most operational questions with minimal overhead.

Signals are emitted alongside IDs and versions: candidate chunk IDs, template version, model/runtime versions, and gate action. This joins observability with reproducibility: engineers can reconstruct a response by fetching those artifacts and replaying the decision ladder [7].

We avoid high-cardinality explosions. Free-form text and token dumps do not enter the hot path. Instead, traces store small arrays of IDs and integers. If deeper forensics are needed, background jobs can fetch bodies by ID under access controls.

Freshness and coverage deserve special handling. Freshness counts toward risk only when policy demands current guidance (e.g., public benefits, regulatory updates). Coverage penalizes prompt bundles that over-concentrate on one page; Max Marginal Relevance (MMR) in re-ranking helps keep it high.

Finally, consistency beats completeness. It is better to emit five stable signals per response than twenty noisy ones that change schema each quarter. Downstream dashboards and alerts remain simple; developers do not chase shifting field names [8].

Signal quality hinges on cardinality discipline. We avoid free-form labels and compress IDs with short, opaque keys. Aggregations run on rolling windows with fixed buckets so dashboards remain legible and alerts do not flap. When deeper analysis is needed, we join on IDs offline under stricter access controls rather than bloating hot-path traces.

Privacy is designed in. Traces never store raw content; they store references and small integers. For regulated namespaces, we hash IDs and rotate salts, then gate any body fetch behind entitlement checks. This preserves replayability without risking accidental data exposure in log pipelines [9].

SLOs become operational when they correlate with user outcomes. We therefore pair support/-coverage with simple product metrics—deflection rate, time to resolution, or edit-after-release. When a template change moves support but harms deflection, the traces make the trade-off explicit and reversible.

Table 1 defines the core signals. Each row is testable and cheap to compute; together they explain most release decisions without heavyweight tracing.

Listing 1 shows a minimal hook that binds prompts to IDs, computes cheap coverage/support, and emits one JSON event for dashboards and gates.

Table 1. Signal dictionary for prompt-centric observability

Signal	Definition	Why it matters	Mode
Coverage	Unique used / unique offered (IDs)	Avoids over-relying on one source; probes diversity	SLO
Support	% of used chunks with at least one cited span	Catches confident but unsupported claims	Gate
Freshness	Age of cited spans vs. namespace policy	Prevents stale answers in time-sensitive domains	Gate
Latency	p50/p95 end-to-end	Keeps UX responsive and budgets predictable	SLO
Spend	Tokens in/out per response	Controls cost regressions	SLO

Listing 1. Minimal prompt-centric telemetry hook

```

1 from typing import Dict, List
2 from time import time
3
4 def emit_prompt_event(prompt_bundle: Dict, draft: Dict, cited: List[Dict],
5                       t_start: float, t_end: float) -> Dict:
6     used = set(draft.get("used_chunk_ids", []))
7     offered = [e["chunk_id"] for e in prompt_bundle.get("evidence", [])]
8     cov = len(used) / max(1, len(set(offered)))
9     cited_ids = {c["chunk_id"] for c in cited}
10    sup = len(used & cited_ids) / max(1, len(used))
11    event = {
12        "ts": int(time()*1000),
13        "template": prompt_bundle.get("template", "unknown"),
14        "model": draft.get("model", "unknown"),
15        "coverage": round(cov, 3),
16        "support": round(sup, 3),
17        "latency_ms": int(1000*(t_end - t_start)),
18        "token_in": draft.get("tok_in", 0),
19        "token_out": draft.get("tok_out", 0),
20        "used_ids": list(used),
21        "offered_ids": offered
22    }
23    log_to_sink(event) # e.g., Kafka/OTLP
24    return event

```

5. Security and Governance Enforcement

Security failures in LLM pipelines are rarely model bugs; they are entitlement gaps, prompt leaks, or unredacted spans. PCO enforces *deny-by-default* at retrieval, carries ACLs into prompts, and blocks release when citations support is weak or freshness windows are violated [10].

Entitlements flow end-to-end. Indexes store ACLs; retrieval filters before scoring; prompts reference only entitled chunks; validators re-check that used IDs remain allowed at release time. This avoids “mask after” failures where sensitive text enters the model and reappears in paraphrase [11].

Governance translates signals into actions via a small decision ladder. If support and coverage exceed thresholds and policy risk is low, release; otherwise rewrite (tighten template), redact (remove spans), or escalate. Shadow thresholds let teams preview stricter policies without surprising users.

Redaction and masking are layered. Ingestion masks PII where possible; prompts apply overlays; the validator blocks or rewrites if sensitive terms appear. This redundancy keeps risk low even when one control blinks [12].

Finally, auditability is a feature. Every action carries reasons and artifact IDs so reviewers can replay the path. This reduces incident MTTR and builds trust with compliance teams.

Table 2 lists practical, low-cost controls. Each control is actionable and testable; together they reduce incident volume substantially.

Table 2. Prompt-centric security and governance controls

Control	What it enforces	Implementation hint	Mode
Entitlement filters	Only readable docs in prompts	Filter by ACL before scoring	Deny-by-default
Template pinning	Reproducible prompts	Include version in every trace	Versioned
Citation support	Claims backed by spans	Require span IDs for critical statements	Gate
Freshness window	No stale guidance	Boost/reject by doc age vs. policy	Gate
PII masking	No sensitive leakage	Mask in ingestion and overlays	Always on
Shadow thresholds	Safe tuning	Evaluate stricter rules without enforcement	Observe
Audit trail	Replayable incidents	Store inputs, IDs, scores, actions	Immutable

Threat modeling starts with the prompt surface. We assume attempts at injection (e.g., “ignore previous instructions”), context smuggling via retrieved snippets, and exfiltration through citations. PCO counters with layered controls: strict system prompts, retrieval filters that reject suspicious patterns, and validators that scan drafts for instruction-overrides before release.

False positives and negatives are managed deliberately. Over-blocking frustrates users; under-blocking risks harm. We log confusion sets—cases where reviewers disagree with the gate—and tune thresholds with shadow evaluation so changes do not surprise production. Over time, disagreement shrinks as rubrics and examples stabilize[13].

When violations occur, the runbook is procedural: freeze releases for the affected namespace, roll back template or index versions, quarantine offending artifacts, and replay golden prompts to verify the fix. Because every decision is tied to artifact versions, rollbacks are flips of pins, not emergency code pushes.

Figure 2 is the single-column decision ladder used in practice: signals feed thresholds; the gate emits one of four actions with reasons; rewrite/redact loop back into prompting.

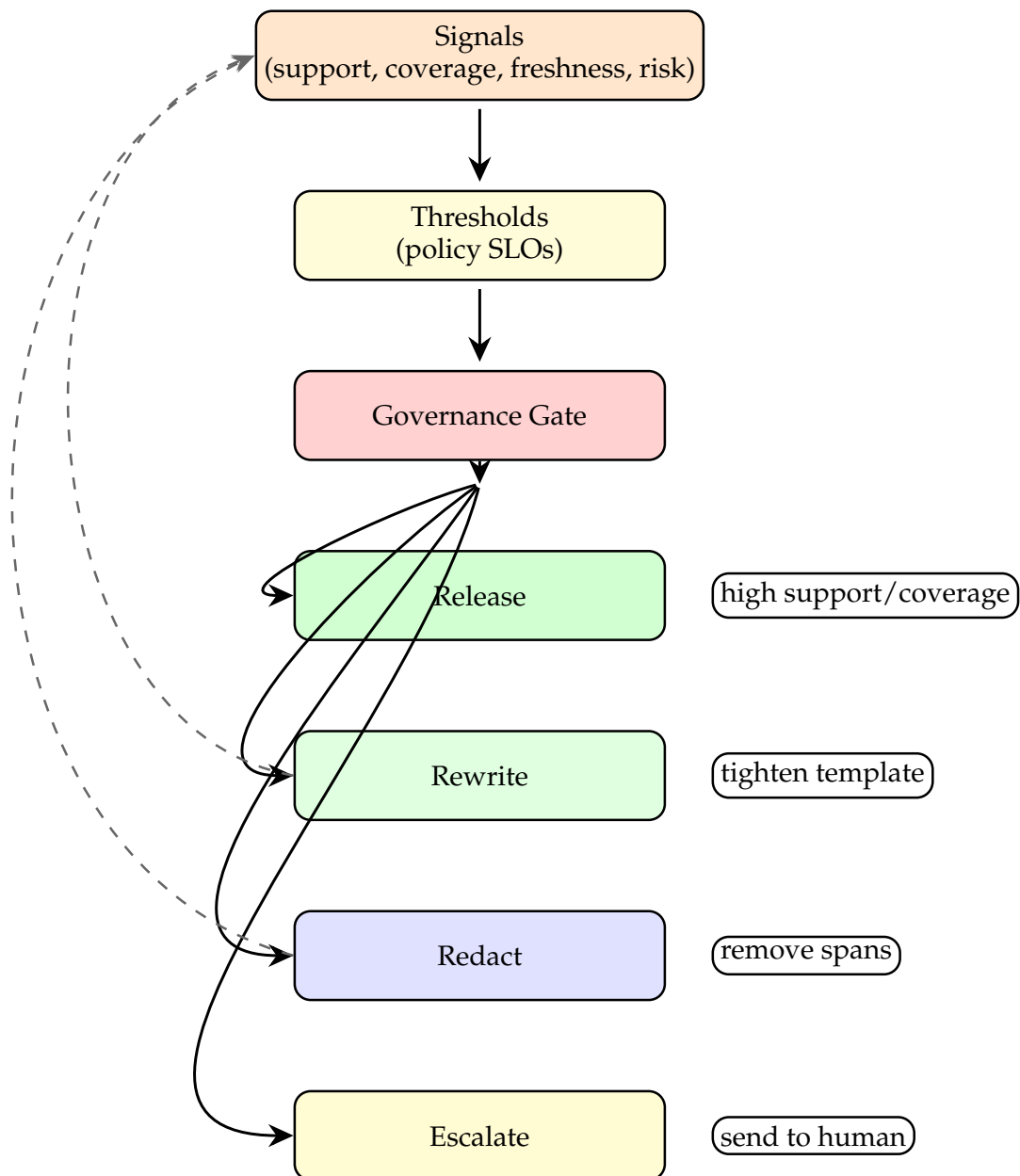


Figure 2. Governed decision ladder with distinct curved routes from the gate to each action (left entry) and dashed feedback arcs to the signals stage.

6. Evaluation Protocols and SLOs

PCO evaluation favors operational realism over synthetic perfection. We run three suites weekly: golden prompts (tricky real patterns), counterfactual pairs (date/name/term swaps), and freshness checks (time-boxed guidance). Each suite computes support, coverage, and abstention; regressions page the on-call.

Offline IR metrics remain useful for retrieval tuning, but they do not predict day-two reliability. We therefore track live SLOs: support ≥ 0.90 , coverage ≥ 0.70 , p95 ≤ 1200 ms, and spend per answer under a token budget. Shadow thresholds let us tighten policies gradually [14].

Sampling matters. We stratify by namespace, language, and user role to avoid blind spots. Underrepresented slices get overweighted in audits so improvements land where they matter most [15].

Finally, we surface trends, not snapshots. A single green week can hide drift; a rolling window reveals whether a fix stuck. Structured traces make it easy to attribute changes to a specific template edit, index rebuild, or validator update.

Good tests are specific and cheap. Golden prompts target brittle patterns (dates, units, disclaimers), counterfactuals probe sensitivity (old vs. new policy terms), and freshness checks guard recency. Each suite yields a compact scorecard that maps cleanly to gate thresholds, so policy changes can be rehearsed before enforcement [16].

We favor canaries over big-bang changes. New templates, validators, or retrieval weights first run on a small traffic slice with alerting on support, coverage, and p95. If stable, the slice grows; if not, rollbacks happen without user impact. This rhythm keeps improvements flowing without destabilizing service [17].

Human audits remain necessary. We sample disagreements between the gate and reviewers, record rationales, and periodically refactor rubrics to match policy intent. The objective is convergence: fewer surprises, tighter confidence bounds, and a smaller set of edge cases requiring human judgment.

Figure 3 shows SLO trends across a week. Labels are placed near ticks and margins are trimmed to keep the plot inside a single column.

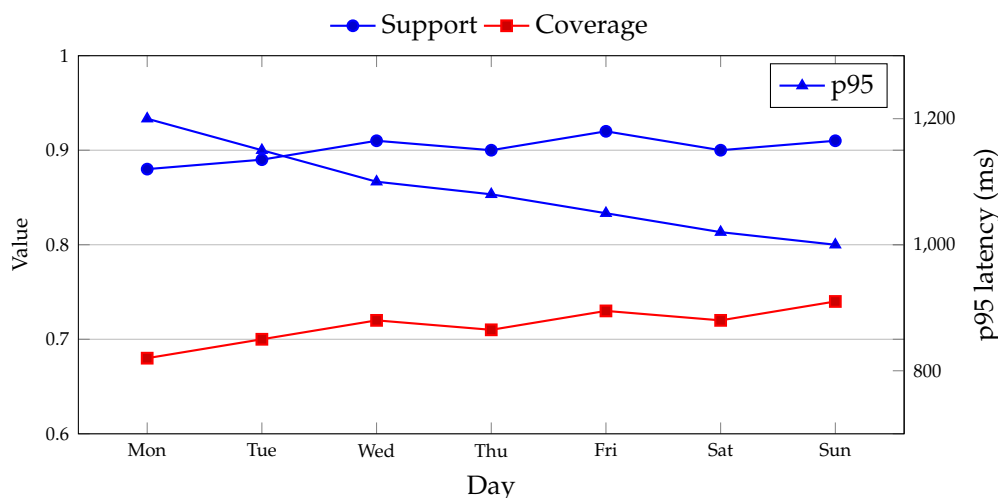


Figure 3. Weekly SLO trends: support and coverage (left axis), p95 latency (right axis).

7. Performance and Cost Engineering

Most budget sits in generation. We cap tokens, stream outputs, and prefer smaller models when evidence is strong. Retrieval and validation get bounded concurrency so a slow detector cannot dominate p95. Caches sit in front of embedding and cross-encoder calls; warmers run after major index rebuilds [18].

Adaptive prompting helps. If coverage/support clear thresholds early, we stop adding evidence; if they fall short, we expand the candidate set or switch to a stricter template. This keeps answers fast when the corpus is strong and careful when it is not.

We measure cost as tokens in/out per response and per namespace. Dashboards show spend outliers and cache hit ratios, guiding where to optimize first: overlong prompts, redundant evidence, or heavy re-ranking. Small, repeated trims often beat one big rewrite.

Finally, we practice safe rollouts. New embedding models, templates, or validator rules run in shadow mode, then under a small traffic slice with alerts on SLOs. Rollbacks are version pin flips, not emergency code pushes [19].

Costs drop fastest when prompts get smarter. We cap evidence by utility (score plus diversity) rather than a flat count, prefer structured snippets over full paragraphs, and stream answers so users see value while the validator finishes. Token budgets are explicit and enforced by the assembler, not left to chance.

Latency tails are managed with bounded concurrency and timeouts per stage. Retrieval, re-ranking, and validation each have budgeted slices; if a slice overruns, we fall back to cached results or

neutral signals rather than stalling the whole request. Dashboards break down p95 by stage so teams know where to invest [20].

Caching is surgical. We cache embedding lookups, cross-encoder scores for recent queries, and compiled templates by version. Invalidations happen on namespace events (index rebuild, template release) to avoid serving stale context. Warmers run after big changes so the first real users do not pay the cold-start tax.

8. Conclusion

This paper argued that reliability in generative AI is less about picking the “best” model and more about making the *prompt pathway* observable, governed, and reproducible. By treating inputs, routing, retrieval, template assembly, generation, validation, and release as first-class, typed interfaces, prompt-centric observability (PCO) turns opaque behavior into explainable operations. The result is not a new benchmark but a control surface: small, stable signals and explicit actions that keep systems within agreed risk bounds as data, policies, and usage evolve.

Our central contribution is a minimal but complete discipline for PCO: a typed artifact graph (CandidateSet, TemplateSpec, PromptBundle, AnswerTrace, ValidationReport); five low-cardinality signals (coverage, support, freshness, p95, spend) suitable for hot-path SLOs; and a governed decision ladder (release, rewrite, redact, escalate) that remains legible under audit. These pieces can be adopted incrementally and fit existing RAG and non-RAG pipelines without re-architecting the whole stack.

The practical payoff is faster, safer change management. When a template edit, index rebuild, or validator tweak shifts outcomes, traces make the cause visible and the rollback trivial—flip a version pin rather than ship a hotfix. Because controls are numeric and typed, product, security, and SRE align on thresholds instead of arguing from anecdotes. Shadow policies let teams rehearse stricter rules before enforcement, avoiding surprises in production.

PCO also reframes cost and latency as quality tools, not just constraints. Token budgets and coverage-aware evidence selection reduce spend while preserving support; bounded concurrency and per-stage timeouts keep tail latency predictable; caches and warmers are scoped by namespace and version so optimizations never leak stale context. These engineering moves are simple, but their consistency across the pipeline delivers compounding reliability.

There are limits. Signals are proxies for truth, reviewers disagree, and thresholds can be gamed. PCO does not remove these realities; it makes them measurable and correctable. Golden prompts, counterfactuals, and freshness checks provide ongoing pressure tests; disagreement sampling and rubric updates shrink ambiguity over time. Most importantly, every decision remains replayable—who saw what, which template and evidence shaped the draft, and why the gate acted as it did.

Prompt-centric observability reframes prompts from strings to infrastructure. By exporting typed artifacts and a handful of stable signals, teams can debug incidents quickly, enforce security consistently, and evolve policies without guessing. The approach favors predictability over heroics: when data or templates change, behavior shifts within guardrails, and every decision is replayable. This discipline turns persuasive systems into dependable services.

Looking ahead, the same contracts that stabilize today’s systems open space for tomorrow’s improvements: retrieval planning that asks for missing evidence, structured citations that hash spans for tamper-evident audit, learned validators that emit typed, checkable rationales, and policy compilers that turn plain-language rules into executable gates. The long-term goal is steady, predictable evolution—systems that change often and safely because the prompt pathway is treated as infrastructure, not as a string.

References

1. Baylor, D.; Brewster, E.; Chao, C.; Fiedel, N.; et al.. TFX: A TensorFlow-Based Production-Scale Machine Learning Platform. In Proceedings of the Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD). ACM, 2017, pp. 1387–1395. <https://doi.org/10.1145/3097983.3098021>.

2. Natarajan, G.N.; Veerapaneni, S.M.; Methuku, V.; Venkatesan, V.; Kanji, R.K. Federated AI for Surgical Robotics: Enhancing Precision, Privacy, and Real-Time Decision-Making in Smart Healthcare. In Proceedings of the Proceedings of the 2025 5th International Conference on Emerging Technologies in Healthcare (ICETH). IEEE, 2025.
3. Khattab, O.; Zaharia, M. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In Proceedings of the Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. ACM, 2020, pp. 39–48. <https://doi.org/10.1145/3397271.3401075>.
4. Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; Mittal, A. FEVER: A Large-scale Dataset for Fact Extraction and Verification. In Proceedings of the Proceedings of NAACL-HLT 2018. Association for Computational Linguistics, 2018, pp. 809–819. <https://doi.org/10.18653/v1/N18-1074>.
5. Petroni, F.; Piktus, A.; Fan, A. KILT: A Benchmark for Knowledge Intensive Language Tasks. In Proceedings of the Proceedings of NAACL-HLT 2021. Association for Computational Linguistics, 2021, pp. 2523–2544. <https://doi.org/10.18653/v1/2021.naacl-main.200>.
6. Karpukhin, V.; Oguz, B.; Min, S.; et al.. Dense Passage Retrieval for Open-Domain Question Answering. In Proceedings of the Proceedings of EMNLP 2020. Association for Computational Linguistics, 2020, pp. 6769–6781. <https://doi.org/10.18653/v1/2020.emnlp-main.550>.
7. Lin, J.; Ma, X.; Lin, S.C.; et al.. Pyserini: A Python Toolkit for Reproducible Research with Sparse and Dense Representations. In Proceedings of the Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval. ACM, 2021, pp. 2356–2362. <https://doi.org/10.1145/3404835.3463238>.
8. Lei, T.; Barzilay, R.; Jaakkola, T. Rationalizing Neural Predictions. In Proceedings of the Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (EMNLP). Association for Computational Linguistics, 2016, pp. 107–117. <https://doi.org/10.18653/v1/D16-1011>.
9. Gebru, T.; Morgenstern, J.; Vecchione, B.; Wortman Vaughan, J.; Wallach, H.; Daumé III, H.; Crawford, K. Datasheets for Datasets. *Communications of the ACM* **2021**, *64*, 86–92. <https://doi.org/10.1145/3458723>.
10. Pasam, V.R.; Devaraju, P.; Methuku, V.; Dharamshi, K.; Veerapaneni, S.M. Engineering Scalable AI Pipelines: A Cloud-Native Approach for Intelligent Transactional Systems. In Proceedings of the Proceedings of the 2025 International Conference on Intelligent Systems and Cloud Computing. IEEE, 2025.
11. Shokri, R.; Stronati, M.; Song, C.; Shmatikov, V. Membership Inference Attacks Against Machine Learning Models. In Proceedings of the 2017 IEEE Symposium on Security and Privacy (SP). IEEE, 2017, pp. 3–18. <https://doi.org/10.1109/SP.2017.41>.
12. Formal, T.; Lassance, C.; Piwowarski, B.; Clinchant, S. SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking. In Proceedings of the Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval. ACM, 2021, pp. 2288–2292. <https://doi.org/10.1145/3404835.3463098>.
13. Lundberg, S.M.; Lee, S.I. A Unified Approach to Interpreting Model Predictions. In Proceedings of the Proceedings of the 31st International Conference on Neural Information Processing Systems (NeurIPS). Curran Associates, Inc., 2017, pp. 4768–4777. <https://doi.org/10.5555/3295222.3295230>.
14. Veluguri, S.P. Deep PPG: Improving Heart Rate Estimates with Activity Prediction. In Proceedings of the Proceedings of the 2025 1st International Conference on Biomedical AI and Digital Health. IEEE, 2025.
15. Ribeiro, M.T.; Wu, T.; Guestrin, C.; Singh, S. Beyond Accuracy: Behavioral Testing of NLP Models with CheckList. In Proceedings of the Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL). Association for Computational Linguistics, 2020, pp. 4902–4912. <https://doi.org/10.18653/v1/2020.acl-main.442>.
16. Shahane, R.; Prakash, S. Quantum Machine Learning Opportunities for Scalable AI. *Journal of Validation Technology* **2025**, *28*, 75–89. <https://doi.org/10.1080/jvtnetwork.v28i1.131>.
17. Agarwal, S.; Peta, S.B. From Notes to Billing: Large Language Models in Revolutionizing Medical Documentation and Healthcare Administration. *Scholars Journal of Applied Medical Sciences* **2025**, *13*, 1558–1566. <https://doi.org/10.36347/sjams.2025.v13i08.005>.
18. Ribeiro, M.T.; Singh, S.; Guestrin, C. “Why Should I Trust You?”: Explaining the Predictions of Any Classifier. In Proceedings of the Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD). ACM, 2016, pp. 1135–1144. <https://doi.org/10.1145/2939672.2939778>.

19. Mitchell, M.; Wu, S.; Zaldivar, A.; Barnes, P.; Vasserman, L.; Hutchinson, B.; Spitzer, E.; Raji, I.D.; Gebru, T. Model Cards for Model Reporting. In Proceedings of the Proceedings of the 2019 ACM Conference on Fairness, Accountability, and Transparency (FAT*). ACM, 2019, pp. 220–229. <https://doi.org/10.1145/3287560.3287596>.
20. Abadi, M.; Chu, A.; Goodfellow, I.; et al.. Deep Learning with Differential Privacy. In Proceedings of the Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS). ACM, 2016, pp. 308–318. <https://doi.org/10.1145/2976749.2978318>.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.