ISSN: 3064-996X | Volume 2, Issue 4

Open Access | PP: 08-13

DOI: https://doi.org/10.70315/uloap.ulete.2025.0204002



Design and Implementation of a YAML-Driven Metrics Layer Framework for Standardized KPI Delivery in Microservices

Abhishek Anand

Data Engineer II - Analytics, Grubhub Holdings Inc.New York City, USA.

Abstract

The article reviews a declarative YAML layer framework implemented to standardize the delivery of key performance indicators within microservice architectures. Justification of work is laid down in the context of a highly rapid shift from monolithic systems to microservices, the widespread adoption of Kubernetes as an orchestrator, and increasing volumes of system telemetry, with fundamental challenges in unifying all these metrics and delivering them as understandable business KPIs. In the absence of a unified semantic vocabulary, the calculation of indicators across different services is performed using disparate filters, which prolongs incident-analysis processes and may result in financial losses reaching millions of dollars per hour. The objective of the study is to externalize the descriptions of source metrics and business formulas from microservice code into a declarative layer, formatted as YAML configurations, and to integrate this layer into a GitOps pipeline to ensure versioning, automatic propagation, and auditing of changes without requiring service restarts. The solution proposed here merges three logical entities: a light-sidecar adapter for Prometheus label normalization and data transfer over the Remote Write 2.0 specification; centralized YAML storage, registered and semantically version-controlled with pull requests; and a proxy processor that compiles declarative formulas into recording rules, performing calculation execution as well as aggregating series for publishing. Therefore, it is possible to say that a declarative framework reduces the time required to deploy a new metric to production from days to hours, minimizes the share of duplicate KPIs, improves SLA compliance regarding calculation latency and telemetry collection overhead, and reduces the need for code-based instrumentation. YAML specifications capture REQUEST (why/what) and encode ACCUT guardrails (what 'good' means) for KPI delivery. The scalable architecture, which separates hot and cold data-processing streams, employs semantic versioning and utilizes an isolated formula interpreter to ensure reliability, stack-storage independence, and compatibility with existing monitoring and tracing systems. This article will be valuable to distributed-systems architects, DevOps engineers, and observability researchers for designing manageable and scalable KPI-standardization platforms.

Keywords: YAML, Microservices, KPI, Declarative Framework, GitOps, Prometheus, Observability, Metrics, CI/CD, Semantic Versioning.

INTRODUCTION

The transition from monolithic systems to microservices began at major Internet companies in the early 2010s, when Netflix (a global video streaming platform), Amazon (an e-commerce and cloud computing giant), and eBay (an online auction and marketplace service) demonstrated that dividing a domain into small independent services accelerates feature releases and enhances fault tolerance. The standardization of containers, with the release of Docker 1.0 (a platform for building, packaging, and running applications inside lightweight, portable containers), followed by the emergence of Kubernetes (an open-source system for automating container deployment, scaling, and management), cemented this approach. By the end

of 2023, 66% of organizations were running Kubernetes in production, while another 18% were at the pilot stage, meaning that 84% of the market had effectively transitioned to cloud-native infrastructure [1]. Concurrently, a culture of system metrics emerged: Prometheus became the de facto standard for data collection, OpenTelemetry unified tracing, and business teams gained access to unprecedented volumes of telemetry that must be transformed into intelligible KPIs. Raw telemetry is published as Bronze, normalized and label-standardized to Silver, and business KPIs are materialized as Gold—keeping the metrics layer consistent with Medallion warehouse practice [9].

The increase in service count has sharply raised observability complexity. A typical product operates hundreds of

Citation: Abhishek Anand, "Design and Implementation of a YAML-Driven Metrics Layer Framework for Standardized KPI Delivery in Microservices", Universal Library of Engineering Technology, 2025; 2(4): 08-13. DOI: https://doi.org/10.70315/uloap.ulete.2025.0204002.

containers, each publishing dozens of metrics with varying naming schemes and labels. In the CNCF 2023 survey, over 90% of respondents reported using containers, with security being the top concern for 40% of companies, and monitoring and observability identified as rapidly growing challenges due to the scale of data [1]. In practice, this results in a single KPI being computed differently across services using varying filters, leading to divergent metrics and incident root-cause analyses that can take hours or even days.

An optimal response to these challenges is to externalize metric descriptions and business formulas from service code into a declarative layer, represented by YAML files, which are version-controlled alongside the infrastructure. The popularization of GitOps has reinforced this idea: according to a 2025 survey, 77% of organizations employ GitOps methodology in some form [2]. When KPIs are expressed in YAML and undergo the familiar pull-request process, the team obtains a unified metric vocabulary, a rigorous change history, and automated configuration propagation without repeated service deployments. This approach eliminates duplication, simplifies auditing, and shortens the time to introduce a new metric to production from days to hours—an outcome unattainable when metrics are calculated imperatively within each microservice.

MATERIALS AND METHODOLOGY

The study is based on the analysis of eight key sources, encompassing both quantitative data on microservice usage practices and technical specifications, as well as empirical evaluations. As starting points, the CNCF Annual Survey 2023 on Kubernetes prevalence and observability challenges [1], the Cloud Native Now 2025 survey on GitOps adoption rates [2], the New Relic report on IT-system downtime financial costs [3], and the DevOps.com 2024 survey on time spent on incident resolution [4] were employed. The empirical study by Hammad et al. provided data on tracing overhead in containerized code [5]. At the same time, the Grafana Labs Observability Survey 2024 facilitated an assessment of monitoring tool diversity and associated costs [6]. To understand GitOps practices in metric configuration management, the CNCF GitOps Microsurvey was consulted [7], and technical requirements for telemetry transmission were studied based on the Prometheus Remote Write 2.0 specification [8].

Methodologically, the work combines several approaches: comparative analysis of metric-collection and unification practices in microservice architectures using aggregated survey and report data [1, 2, 6]; systematic review of technical documentation and specifications to identify requirements for declarative metric description [8]; and performance analysis of tracing and monitoring tools based on empirical measurements [5]. To structure architectural requirements, a content analysis of reports on business costs of downtime and observability challenges was conducted [3, 4], which facilitated the formalization of target metrics and SLAs. Finally, the synthesis of GitOps implementation data

[7] informed the design of the CI/CD process for the YAML registry. As a result, the methodology integrates literature review, specification analysis, comparative practice analysis, and case-study content analysis to substantiate the choice of a declarative YAML layer in the proposed framework. Schema fields are selected to satisfy ACCUT (accuracy/correctness checks, completeness constraints, uniqueness keys, and timeliness SLAs). KPI entities are aligned with Kimball's 4-step modeling (process, grain, dimensions, facts) to keep business vocabulary stable [9].

RESULTS AND DISCUSSION

For business, key metrics are considered part of financial obligations: a survey of 1700 companies showed that highly critical outages last on average 77 hours per year, and each lost hour of iteration costs up to USD 1.9 million [3], engineers spend nearly one third of their workweek—approximately 30% of their time—on incident remediation [4]. Conducting root cause analysis (RCA) and post-incident reviews (37%), monitoring DevOps Research and Assessment (DORA) metrics (34%), monitoring golden signals (33%), and tracking, reporting, and influencing the mean time required to detect and resolve outages (33%) are the top observability best practices being followed, as shown in Figure 1.

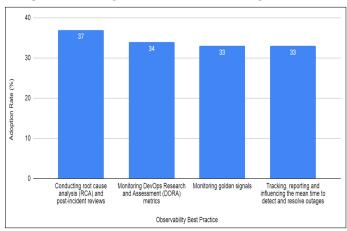


Fig. 1. Observability Practices Adoption and Impact [4]

Under these conditions, product owners demand that KPI data be updated with one-minute latency, support audit and versioning, and comply with the SLA. Discrepant KPIs arise when REQUEST is implicit. 'Reason/Entities/Units/Events' are externalized into YAML so that formulas are composed from the same contract across services [9]. Meanwhile, observability teams must provide a unified semantic metric vocabulary independent of the service implementation language, with the ability to trace each calculation formula back to its source metrics.

The operational environment imposes strict constraints: over 80% of organizations already run Kubernetes in production, with another 13% at the pilot stage, so most services operate in short-lived containers, scale horizontally, and communicate via a service mesh [2]. High label cardinality, regulatory restrictions on cross-region data transfer, and the need to process telemetry from a polyglot stack without

modifying business code further complicate observability load. Moreover, 45% of companies already use more than five monitoring tools, creating significant licensing and integration costs, and making it critical that a new platform operate on top of the existing landscape without an agent zoo [4].

Success criteria are quantified: the mean time to detect and resolve incidents must decrease, the proportion of duplicate KPIs must decline, and the time to deploy a new metric to production must not exceed two hours from the merge of its YAML configuration. The platform must ensure that aggregated KPI query latency is below a specified number of milliseconds at a given percentile, keep telemetry collection overhead within each node's CPU budget, and autoscale under load spikes without data loss. Meeting these conditions will demonstrate that the declarative layer truly addresses business objectives, fits technological constraints, and delivers measurable operational benefit.

In-service telemetry remains popular due to its ease of adoption; however, studies reveal its systemic costs. In a comparison of seventy APIs across two clouds, full-code tracing coverage reduced aggregate throughput by an average of 8.4%, with losses of up to 30% in some cases [5]. These percentages translate into direct costs when dozens of microservices operate under peak load; additionally, each team describes metrics differently, leading to discrepancies at the naming and labeling level. Any change to a business formula necessitates a new deployment, which slows KPI evolution and constrains analysts.

Second-generation observability platforms promise to relieve integration pain by offering a single pane of glass, yet in practice they exacerbate tooling fragmentation: the Grafana Observability Survey 2024 found that 70% of teams use at least four tools for metrics, logs, and tracing, and 61% of participants cite unpredictable costs and excessive bills as the primary operational challenge, as shown in Figure 2 [6]. RIPC guides platform design: Reduce duplicate collectors, pre-Index label spaces, Partition writes by KPI ID/tenant, Cache compiled rules [9].

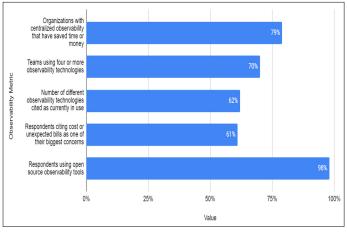


Fig. 2. Quantitative Assessment of Observability Tool Adoption, Diversity, and Cost Implications [6]

Each additional aggregator introduces its format, cardinality quotas, and pricing model, so computing a unified KPI often requires chains of import–export between systems, leading to version conflicts among libraries and exporters.

A declarative YAML layer eliminates both problem classes by externalizing KPI descriptions into a version-controlled repository and integrating them into a GitOps pipeline. A survey [7] revealed that 60% of respondents have employed GitOps for over a year, while another 31% started using it within the past year. Additionally, 67% of the remaining respondents plan to adopt it within the next year. Most organizations rely on one or two cloud platforms or a hybrid model, whereas purely on-premises or complex multi-cloud configurations are notably less common, as illustrated in Figure 3.

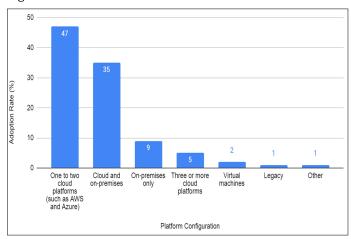


Fig. 3. Distribution of Organizational Infrastructure Platforms [7]

Thus, storing metrics in YAML naturally aligns with the established culture of declarative infrastructure: a change to a formula is recorded via a pull request, automatically subjected to tests and migrations, and then propagated to the proxy layer without requiring service restarts. This eliminates duplication, reduces operational expenditure, and limits the latency for deploying a new metric to the time needed for review rather than a complete deployment cycle, thereby directly supporting the defined success criteria. When a rule fails validation or exceeds cardinality budgets, the rule version is quarantined and the last good snapshot continues to be served [9].

The basic topology is built around three logical nodes: a lightweight adapter deployed as a sidecar to each microservice, a YAML registry that stores KPI configurations in Git, and a proxy processor responsible for computing and publishing aggregated series. Services continue to export telemetry in Prometheus format. The adapter normalizes labels, applies cardinality limiters, and, using the Remote Write 2.0 protocol, forwards samples losslessly to the receiver. The specification mandates control of protocol headers and acknowledgments, thus eliminating double reads and enhancing transmission reliability [8]. A Token Bucket is

Design and Implementation of a YAML-Driven Metrics Layer Framework for Standardized KPI Delivery in Microservices

applied per tenant at ingest; on overload/429s, limits are halved via AIMD and then slowly increased on recovery [9]. The YAML registry is hosted on a Git platform, where each modification undergoes the standard pull-request process and, thanks to the GitOps pipeline, automatically triggers a configuration update in the proxy.

Data flows along two independent channels. The hot path runs from adapters to the metric collector via Remote Write; the collector stores raw series in a scalable TSDB and immediately publishes them for PromQL queries. In parallel, upon startup, the proxy loads a snapshot of all YAML files, compiles formulas into recording rules, and, on schedule, recalculates KPIs—applying the same labels as the source metrics — so that consumers can switch to unified series without reworking their dashboards. The cold path for configuration changes originates from Git. After a merge, a bot performs schema validation and test computations on sample data, and upon success, delivers the new YAML package to the proxy, thereby preserving continuous service for historical version queries.

Scaling is achieved through strict node independence: sidecars scale horizontally in line with service replicas, while Kubernetes HPA manages proxies and collectors. Within the proxy, computation tasks are sharded by KPI identifier, allowing for linear throughput increases without requiring state coordination; an intermediate results cache reduces TSDB load during read spikes. Circuit-breaker policies define acceptable degradation windows: upon timeout, the proxy returns the last valid snapshot, and the adapter enters a drop-labels mode, discarding low-priority labels to cap cardinality. As a result, the system endures horizontal cluster expansion without component refactoring and continues to deliver KPIs.

The schema description begins with a minimal yet comprehensive file structure reflecting the domain entity of the key indicator. At the top level, it specifies a stable metric identifier, a version field, a list of source series, and a section for the calculation formula. For every KPI, the grain (e.g., service, region, minute) is declared and dimension keys are referenced; this STAR-compatible contract avoids silent semantic drift [9]. Each source metric records its name, a mandatory set of labels, and an aggregation type, guaranteeing result reproducibility regardless of the producing service's language or environment. The formula is written in a compact expression syntax closely resembling Python, enabling an analyst with basic arithmetic knowledge to define complex indicators without modifying microservice code.

The schema defines several source categories. The most common is Prometheus telemetry delivered via Remote Write. The second category comprises high-level series that are already aggregated in the time-series database and require only further normalization. The third includes events from logging systems, such as counters of successful

and failed transactions, which the proxy dynamically converts into time series. This source catalog spares teams from formatting data on the service side and simultaneously allows the proxy to optimize the pipeline for each load type.

Formula and aggregation descriptions follow a declarative approach: the KPI author operates on a list of reference tokens, such as metric0, metric1, and so on, while the compiler substitutes actual queries. Supported constructs include arithmetic, conditional expressions, and aggregate functions—such as total count and moving-window average. The constrained feature set prevents arbitrary code execution, thereby reducing the attack surface.

Versioning adheres to a semantic scheme: changes to source-metric descriptions increment the patch number, recomputation over historical intervals increments the minor number, and incompatible changes to aggregation type or measurement unit increment the major number. This rule simplifies automatic selection of the correct formula version when reading historical data, since the proxy always knows which parameters were valid at write time.

Processing begins at the validator stage, which checks JSON-schema compliance, tests the expression on a controlled value set, and rejects configurations that yield indivisible values or exceed allowable ranges. Static analysis further detects dead code branches and cyclic metric references. Next, the compiler transforms the declaration into PromQL or into an SQL dialect when the storage supports aggregated-table queries. The output is stored in a recording-rules catalog, which the executor polls at intervals defined by the update policy.

The executor performs computations, writes results back to the same time-series database, and stores intermediate values in a local cache. ACCUT telemetry (accuracy deltas, staleness) and stage-level metrics (latency p95) are emitted for each compute run, enabling rapid RCA similar to Spark UI heuristics [9]. The cache enables dashboards and alerts to function even during brief storage outages and dramatically reduces load when multiple queries hit the same KPI. Configuration reloads occur without service interruption. Git triggers a webhook, which delivers the update to the proxy. The proxy then loads the new file package, validates it locally, and atomically updates the rule-version reference. Consequently, consumers receive updated metrics within seconds, and computations initiated under the previous schema complete correctly, preventing partial data loss and ensuring SLA-level reliability. Bronze/Silver/Gold are versioned in lockstep with KPI YAML; rollbacks swap only the rule-version pointer, keeping raw series intact [9].

The parser is implemented using the Pydantic library, which converts the loaded YAML into strongly typed objects, verifying required fields, value ranges, and label-list integrity before the configuration enters the execution subsystem. This approach frees developers from manual validation and provides clear, human-readable error reports, simplifying

Design and Implementation of a YAML-Driven Metrics Layer Framework for Standardized KPI Delivery in Microservices

the repository review process. Upon successful validation, the object model is passed to the query generator. For each source metric, it matches the storage backend, applies filters and the aggregation type, and then constructs an expression in PromQL or the SQL dialect, as determined by the connected storage. As a result, analysts obtain a uniform query signature irrespective of the underlying time-series storage technology.

For reliable and secure formula evaluation, an isolated interpreter is employed. It is instantiated with an empty set of built-in functions, permitting only the arithmetic operations and aggregates required for business indicators. The code of each formula is compiled into an abstract syntax tree and then executed within this environment, where input/output operations and network access are unavailable, thereby virtually eliminating abuse. If computation completes successfully, the result is immediately sent to the cache, from which dashboards and alerting systems retrieve it. For external clients, a thin REST and gRPC layer has been developed, offering a list of available KPIs, their descriptions, and metadata, as well as serving precomputed series in formats compatible with popular visualization tools. This abstraction boundary conceals the internal pipeline details, enabling teams to adopt new storage backends or optimize algorithms without modifying client code.

The configuration delivery process to production is built around the familiar Git flow. Each change begins on a separate branch, undergoes automated review by code owners, and is then merged into the main development line via a pull request. Post-merge, a trigger launches a pipeline that repeats validation, runs a suite of tests on synthetic data, and deploys the proxy container in an isolated namespace, where analysts can verify how the new metric will appear on real dashboards. If the outcome satisfies stakeholders, the image is promoted to the production cluster. Should an error be discovered, reverting the last commit suffices to restore the previous rule set via the atomic configurationswitch mechanism. This cycle guarantees that changes never disrupt observability service operation, and the delivery time for new metrics is limited to the duration of review and automated tests.

Thus, the adoption of a declarative metrics layer based on YAML centralizes KPI descriptions, ensures uniform computation, and fully audits changes with minimal temporal overhead, while preserving scaling flexibility and independence from particular microservice implementations, reducing operational costs and the risk of metric discrepancies; all of this integrates seamlessly into existing GitOps processes and the established culture of cloud infrastructure.

CONCLUSION

The new setup showcases nice standardization and brings together key indicator calculations in lively small-service scenes. Moving metric details and business formulas out of service code into GitOps version control makes configuration management more centralized, provides a clear change history, and enables rule-based automatic deployment of new rules without requiring component restarts. This setup eliminates busy work, reduces metric drift, and reduces the time it takes to push a new metric to production from days to just hours.

The technical implementation, comprising three logical nodes—a lightweight sidecar adapter, a YAML configuration registry, and a proxy processor—has proven to be scalable and fault-tolerant. The hot telemetry stream ensures immediate availability of raw data in the TSDB, while the proxy layer periodically recalculates aggregated KPIs, preserving a unified semantics for labels and metrics. Caching mechanisms and circuit-breaker policies guarantee SLA compliance under peak loads, and the isolated formula interpreter reduces the attack surface and prevents arbitrary code execution.

Reduced time to detect and eliminate incidents due to unified metric vocabulary and fast root cause analysis;; percentage of duplicate KPIs reduced due to a single source of truth;; time to deploy a new metric to production does not exceed the established two-hour SLA has been the key quantitative success criteria. No additional agents are required, and compatibility with existing monitoring and tracing tools supports the adoption of the framework without disrupting the current landscape.

Therefore, the proposed YAML-oriented approach provides a resilient, flexible, and manageable platform for delivering standardized KPIs in microservice architectures, aligns with established GitOps workflows and cloud infrastructure practices, and lowers operational costs while enhancing the reliability of business metrics.

REFERENCES

- 1. CNCF, "CNCF Annual Survey 2023," *CNCF*, Apr. 09, 2024. https://www.cncf.io/reports/cncf-annual-survey-2023/ (accessed Jun. 26, 2025).
- 2. M. Vizard, "CNCF Survey Surfaces Steady Pace of Increased Cloud-Native Technology Adoption," *Cloud Native Now*, Apr. 04, 2025. https://cloudnativenow.com/topics/cloudnativedevelopment/cncf-survey-surfaces-steady-pace-of-increased-cloud-native-technology-adoption/(accessed Jun. 27, 2025).
- 3. New Relic, "New Relic Study Reveals IT Outages Cost Businesses Up to \$1.9 M Per Hour," *New Relic*, 2024. https://newrelic.com/press-release/20241022 (accessed Jun. 28, 2025).
- M. Vizard, "Survey: IT Teams Spend About a Third of Time Responding to Disruptions," *DevOps*, Nov. 2024. https://devops.com/survey-it-teams-spend-about-athird-of-time-responding-to-disruptions/ (accessed Jun. 29, 2025).

Design and Implementation of a YAML-Driven Metrics Layer Framework for Standardized KPI Delivery in Microservices

- 5. Y. Hammad, A. A.-S. Ahmad, and P. Andras, "An Empirical Study on the Performance Overhead of Code Instrumentation in Containerised Microservices," *Journal of Systems and Software*, vol. 230, pp. 112573–112573, Jul. 2025, doi: https://doi.org/10.1016/j.jss.2025.112573.
- 6. Grafana Labs, "Observability Survey Report 2024 key findings," *Grafana Labs*, 2024. https://grafana.com/observability-survey/2024/ (accessed Jul. 01, 2025).
- 7. G. Microsurvey, "Learning on the job as GitOps goes mainstream," *CNCF*. https://www.cncf.io/wp-content/uploads/2023/11/CNCF_GitOps-Microsurvey_Final.pdf (accessed Jul. 02, 2025).
- 8. Prometheus, "Prometheus Remote-Write 2.0 specification," *Prometheus*, 2024. https://prometheus.io/docs/specs/prw/remote_write_spec_2_0/ (accessed Jul. 03, 2025).
- 9. A. Anand, SELECT * FROM fact_DE. Independently published, 2025, p. 304. Accessed: Oct. 18, 2025. [Online]. Available: https://a.co/d/4qMytUP

Copyright: © 2025 The Author(s). This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.