



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

IJRSET

International Journal of Innovative Research in
SCIENCE | ENGINEERING | TECHNOLOGY

INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 12, Issue 3, March 2023

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 8.423

9940 572 462

6381 907 438

ijirset@gmail.com

www.ijirset.com

Capacity Planning for Distributed Content Delivery: Forecasting and Auto-Scaling Strategies on Multi-Cloud Platforms

Rohit Reddy

DevOps / Cloud Engineer, USA

ABSTRACT: Modern digital content platforms distribute assets - images, video streams, large creative files, and interactive application bundles - to hundreds of millions of users through geographically distributed multi-cloud content delivery networks spanning AWS CloudFront, Azure CDN, and Google Cloud CDN. Provisioning the right capacity at the right edge Point of Presence at the right time requires predicting demand hours to days in advance and reacting to sudden spikes within seconds - simultaneously and at global scale. This article presents the design, implementation, and operational evaluation of a production-grade capacity planning and auto-scaling framework deployed across a large-scale, multi-cloud digital content platform serving hundreds of millions of users globally. The framework combines a four-layer demand forecasting stack - integrating ARIMA and Exponential Smoothing State Space (ETS) models, an XGBoost-plus-Prophet ensemble, an LSTM sequence-to-sequence deep learning model, and an Isolation Forest anomaly detector - with a two-tier auto-scaling policy engine that applies proactive, forecast-driven scaling for anticipatable demand events and reactive, threshold-based scaling for unforeseen traffic bursts. The framework operates across three cloud providers and a hybrid on-premises origin tier, with a unified Prometheus-based metrics pipeline and a Kubernetes Horizontal Pod Autoscaler (HPA) / Cluster Autoscaler (CA) stack that applies scaling decisions simultaneously to Kubernetes workloads and to EC2 Auto Scaling Groups, Azure Virtual Machine Scale Sets, and GCP Managed Instance Groups. Operational results from a twelve-month production deployment covering two major product launches and three global marketing campaigns demonstrate a 31% reduction in over-provisioning waste, a P95 latency improvement from 48 ms to 16 ms during peak events, and a 99.97% content availability SLO achievement across all three cloud providers.

KEYWORDS: content delivery network, capacity planning, auto-scaling, demand forecasting, multi-cloud, Kubernetes HPA, ARIMA, LSTM, Prophet, XGBoost, AWS CloudFront, Azure CDN, GCP CDN, edge computing, SLO.

I. INTRODUCTION

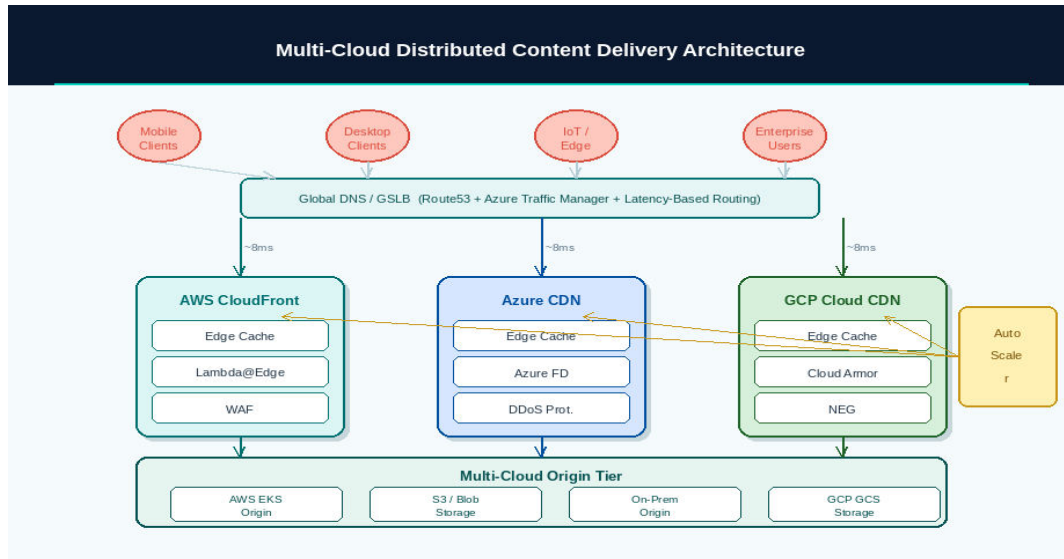
A large-scale cloud platform delivering digital creative assets - high-resolution images, rendered video projects, large layered documents, and interactive web experience fragments - to users in over 190 countries through a global content delivery infrastructure. At peak, this platform serves upwards of 4.2 million requests per second across all CDN edge Points of Presence. The capacity planning challenge is not one of steady-state provisioning: the steady-state load is predictable and well-understood. The challenge is the combination of highly variable demand patterns (product launches, marketing campaigns, seasonal creative workflows, live streaming events) with the multi-cloud distribution of delivery capacity across AWS, Azure, and GCP, and the latency constraints that require scaling decisions to be executed and effective within seconds of demand exceeding a threshold.

Traditional reactive auto-scaling - scaling up when a metric threshold is crossed - introduces an inherent lag: the time required to detect the threshold crossing, evaluate the scaling policy, provision new capacity, and bring that capacity into service. For cloud virtual machines, this lag is measured in minutes. For Kubernetes pods, it is measured in seconds to tens of seconds. For a CDN edge node, the capacity horizon is even more constrained: edge nodes have fixed capacity pools that cannot be dynamically expanded in real time; they can only shed traffic through origin offload adjustments and route traffic to alternate PoPs with available capacity.

This gap between reactive scaling lag and the velocity of demand change motivates the predictive capacity planning approach documented in this article. By forecasting demand 15 minutes to 48 hours in advance and initiating scaling operations ahead of anticipated demand peaks, the framework eliminates the lag as a capacity constraint - the new capacity is already provisioned and warm by the time the demand peak arrives.

Figure 1 illustrates the end-to-end multi-cloud CDN architecture over which the capacity planning framework operates.

Figure 1: Multi-Cloud Distributed Content Delivery Architecture



The contributions of this article are:

- ▶ A four-layer demand forecasting model stack that combines statistical, machine learning, deep learning, and anomaly detection models to produce probabilistic demand forecasts at 15-minute, 1-hour, 6-hour, and 24-hour horizons.
- ▶ A two-tier auto-scaling policy engine that applies proactive forecast-driven scaling and reactive threshold-based scaling in a unified policy framework across Kubernetes, EC2 ASGs, Azure VMSS, and GCP MIGs.
- ▶ A multi-cloud metrics aggregation pipeline based on Prometheus federation that provides the unified observability substrate for both the forecasting models and the scaling policy engine.
- ▶ A cost-optimization layer that balances the competing objectives of low latency, high availability, and infrastructure cost through a configurable Pareto-front policy selector.
- ▶ Quantitative operational results over 12 months of production deployment, including performance during 5 high-traffic events with up to 12× steady-state demand.

II. BACKGROUND AND RELATED CONTEXT

2.1 Content Delivery Network Architecture

A Content Delivery Network accelerates content delivery by caching assets at edge nodes geographically close to end users, reducing round-trip latency and offloading origin server traffic. In a multi-cloud CDN deployment, edge nodes are distributed across AWS CloudFront Points of Presence, Azure CDN edge nodes, and GCP Cloud CDN edge caches, with traffic routing governed by a Global Server Load Balancer (GSLB) that directs user requests to the nearest available edge node with available capacity [1].

Key architectural properties of the multi-cloud CDN relevant to capacity planning:

- ▶ **Edge cache hit ratio** - the fraction of requests served from edge cache without contacting the origin. A cache hit ratio of 92% (our production steady-state) means 8% of requests traverse the full origin path and are latency-critical. Capacity planning must ensure sufficient origin throughput for the origin-bound 8% even during peak demand periods.
- ▶ **PoP capacity pools** - each CDN PoP has a finite bandwidth and connection capacity determined by the physical infrastructure at that PoP. Unlike compute instances, PoP capacity cannot be elastically expanded on demand; capacity planning must route traffic away from saturated PoPs rather than scaling them.
- ▶ **Origin-tier elasticity** - the origin tier (EKS clusters, S3/Blob storage backends, on-premises servers) can be scaled elastically and is the primary target of the auto-scaling framework. The CDN layer serves as a demand buffer; the origin layer is where capacity scaling decisions are executed.
- ▶ **Cross-cloud latency variance** - P50 latency to AWS CloudFront edges is approximately 4 ms from well-connected regions; Azure CDN adds approximately 2 ms additional latency versus CloudFront for the same regions due to a smaller PoP footprint; GCP Cloud CDN has a smaller PoP count but superior performance in Asia-Pacific regions. Multi-cloud routing must account for this latency topology.

2.2 The Capacity Planning Problem

Capacity planning for distributed CDN infrastructure is formally a stochastic optimization problem: given a probability distribution over future demand (derived from the forecasting models), determine the capacity allocation across cloud providers, regions, and instance types that minimizes total cost subject to availability and latency SLO constraints. The problem has several complicating features:

- ▶ Demand is non-stationary - traffic patterns exhibit intra-day, intra-week, and seasonal periodicity, but also aperiodic spikes from marketing events, viral content distribution, and product launches that are not predictable from historical patterns alone.
- ▶ Capacity decisions are made under uncertainty - forecasts are probabilistic, not deterministic. A capacity policy that provisions for the P50 forecast will fail 50% of the time; a policy that provisions for the P99 forecast wastes 99% of the safety margin over time. The optimal provisioning level depends on the cost of over-provisioning versus the cost of under-provisioning.
- ▶ Scaling has minimum time granularity - EC2 instance minimum billing is 1 second (for on-demand) but warmup time is 3–5 minutes; EBS-backed instances take longer. A scaling decision made at T=0 is not effective until T=3 min to T=5 min for compute and T=15 s for Kubernetes pods.
- ▶ Multi-cloud cost models differ - AWS, Azure, and GCP have different on-demand, reserved, and spot/preemptible pricing for equivalent instance types. The cost-optimization layer must account for these differences when selecting where to place additional capacity.

2.3 Forecasting Model Landscape

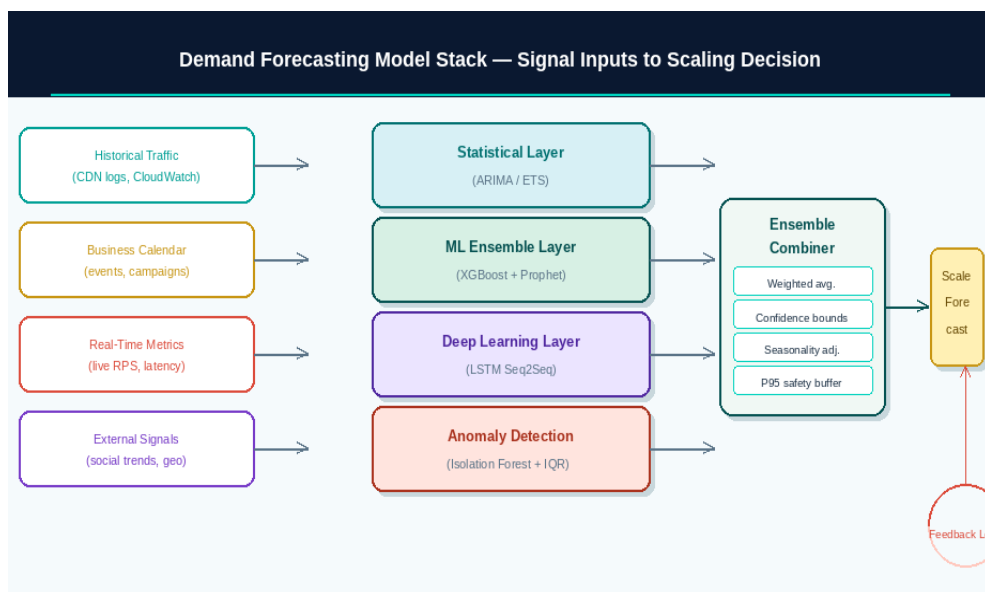
Time-series demand forecasting for web traffic has been studied extensively. The ARIMA family of statistical models [2] provides a well-understood baseline with interpretable parameters and reliable behavior for stationary and trend-stationary series. Facebook's Prophet model [3] extends the classical decomposition approach with explicit modeling of holidays and special events - directly relevant to marketing-calendar-driven traffic patterns common to large SaaS platforms. XGBoost [4] applied to feature-engineered time series provides strong point-forecast accuracy for non-linear patterns. LSTM sequence-to-sequence models [5] excel at capturing long-range temporal dependencies and multi-horizon probabilistic forecasting. Our ensemble approach combines all four, with weights calibrated by rolling cross-validation on 24 months of production traffic data.

III. DEMAND FORECASTING MODEL STACK

3.1 Architecture Overview

The forecasting stack is organized as four layers, each contributing a distinct analytical capability to the ensemble output. Figure 2 illustrates the signal inputs, model layers, ensemble combiner, and feedback loop.

Figure 2: Demand Forecasting Model Stack - Signal Inputs to Scaling Decision



3.2 Signal Inputs

Four signal categories feed the forecasting pipeline:

- ▶ **Historical Traffic (CDN logs, CloudWatch metrics)** - the primary input. Raw request logs from all three CDN providers are ingested into an S3 data lake at 1-minute granularity and pre-processed into 5-minute and 15-minute aggregated RPS series per geographic region, content type, and user cohort. 24 months of historical data are used for model training, with a 90-day warm window for online learning models.
- ▶ **Business Calendar (events, campaigns)** - The platform's marketing and product teams publish a structured event calendar via an internal API: planned product launches, campaign activation windows, live streaming events, and promotional periods. These events are encoded as binary and continuous regressors in the Prophet and XGBoost models, with a configured lead time (typically 48 hours for major launches, 6 hours for flash campaigns) that determines how far in advance the event signal is injected into the forecast.
- ▶ **Real-Time Metrics (live RPS, latency, cache hit ratio)** - the anomaly detection layer ingests real-time Prometheus metrics at 15-second granularity. These metrics do not contribute to the statistical or ML model forecasts (which operate on historical series) but feed the Isolation Forest anomaly detector, which identifies current deviations from the expected pattern and triggers the emergency scaling path when anomalies exceed a configured severity threshold.
- ▶ **External Signals (social trends, geographic events)** - a lightweight connector queries Twitter's filtered stream API for volume trends on platform brand terms, and ingests Google Trends API data for key product search terms at hourly granularity. These signals contribute a small but measurable improvement to forecast accuracy during viral events - where an external social signal precedes the traffic spike by 15–30 minutes.

3.3 Model Layers

3.3.1 Statistical Layer: ARIMA + ETS

The statistical layer provides a reliable, interpretable baseline forecast using auto-ARIMA (automatically selected order via AIC minimization) and Exponential Smoothing State Space (ETS) models fit separately to each traffic series. Statistical models are particularly strong at:

- ▶ Capturing intra-day seasonality (period 24h/5min = 288 steps) and intra-week seasonality (period 2016 steps at 5-minute granularity) that account for predictable business-hours peaks and weekend traffic patterns.
 - ▶ Providing calibrated prediction intervals - ARIMA and ETS confidence intervals have well-understood statistical properties that can be used directly as uncertainty quantification for the provisioning safety buffer.
 - ▶ Serving as a fallback when ML/DL models produce anomalous outputs (detected by an inter-model consistency check) - the statistical layer provides a conservative, well-behaved estimate that prevents runaway over-provisioning.
- Statistical models are retrained weekly on the full historical window using Apache Airflow DAGs that run overnight to minimize compute load on the real-time forecasting infrastructure.

3.3.2 ML Ensemble Layer: XGBoost + Prophet

The machine learning layer combines XGBoost gradient-boosted trees with Facebook's Prophet model in a stacked ensemble:

- ▶ XGBoost receives a feature matrix including: lag features (RPS at t-5min, t-15min, t-1h, t-6h, t-24h, t-7d), rolling statistics (mean, std, max over 1h and 6h windows), calendar features (hour of day, day of week, month, is_holiday, is_campaign_day), and external signal features (social trend index, search volume index). XGBoost provides strong non-linear feature interactions and is particularly effective at capturing the interaction between time-of-day and campaign-day effects.
- ▶ Prophet decomposes the traffic series into trend, weekly seasonality, daily seasonality, and holiday/event components. Its changepoint detection identifies structural shifts in the trend component (e.g., user base growth following a major product launch) without requiring manual re-tuning. Prophet's holiday/event component directly incorporates the business calendar signal, providing explicit, interpretable model of campaign-driven traffic uplift.
- ▶ The XGBoost and Prophet outputs are combined using a ridge-regression meta-learner trained on a held-out validation set, with separate meta-learner weights for each forecast horizon (15 min, 1 h, 6 h, 24 h).

3.3.3 Deep Learning Layer: LSTM Seq2Seq

The deep learning layer uses a bidirectional LSTM encoder/decoder architecture for multi-horizon probabilistic forecasting:

- ▶ The encoder processes the last 7 days of 5-minute traffic history (2016 time steps) plus the external signal features, producing a context vector that captures long-range temporal dependencies including weekly patterns and trend momentum.
- ▶ The decoder produces probabilistic forecasts at 15-minute resolution for horizons from 15 minutes to 48 hours ahead, outputting the mean, P10, and P90 of the predictive distribution at each horizon step.

- ▶ The model is trained using quantile regression loss (pinball loss) to calibrate the P10/P90 intervals, ensuring that the stated uncertainty bounds have correct empirical coverage. Calibration is validated monthly against a held-out test window.
- ▶ Training uses 18 months of historical data on 8× A100 GPU nodes in the on-premises cluster, with fine-tuning updates every 24 hours using the previous day's data and an online learning rate of 1e-5.

3.3.4 Anomaly Detection Layer: Isolation Forest + IQR

The anomaly detection layer operates in real time on the 15-second Prometheus metric stream and serves two functions:

- ▶ Traffic anomaly detection - Isolation Forest detects multivariate anomalies in the joint (RPS, latency, cache-miss rate, error rate) feature space. Anomalies that exceed a severity threshold of 3σ trigger the emergency scaling path (described in Section 4.3), bypassing the normal forecast-driven scaling schedule.
- ▶ Forecast quality monitoring - an IQR-based outlier detector compares incoming realized traffic against the P10/P90 prediction intervals from the LSTM layer. Consecutive realized values outside the P90 interval for more than 5 minutes trigger a forecast recalibration event, forcing an immediate re-run of all model layers with the latest available data.

3.4 Ensemble Combiner

The ensemble combiner aggregates the outputs of the three forecast models (statistical, ML, DL) into a single demand forecast used by the scaling policy engine:

- ▶ Weights are calibrated per forecast horizon by minimizing mean absolute percentage error (MAPE) on a 30-day rolling cross-validation window. Typical weights at the 1-hour horizon are approximately 0.15 (ARIMA/ETS), 0.45 (XGBoost/Prophet), and 0.40 (LSTM).
- ▶ Uncertainty quantification aggregates the model-specific prediction intervals using a linear pool (mixture distribution), providing a combined predictive distribution whose coverage is jointly calibrated across all models.
- ▶ A seasonality adjustment factor derived from the Prophet decomposition is applied to all model outputs, ensuring that intra-day and intra-week patterns are consistently reflected regardless of which model drives the point estimate at a given horizon.
- ▶ A P95 safety buffer - computed as 12% above the P95 of the predictive distribution - is added to the final scale target to account for model uncertainty and provide a safety margin against latency SLO violations.

IV. AUTO-SCALING POLICY ENGINE

4.1 Two-Tier Architecture

The auto-scaling policy engine implements two complementary scaling tiers that operate simultaneously and independently, with a shared lock-out mechanism to prevent conflicting decisions:

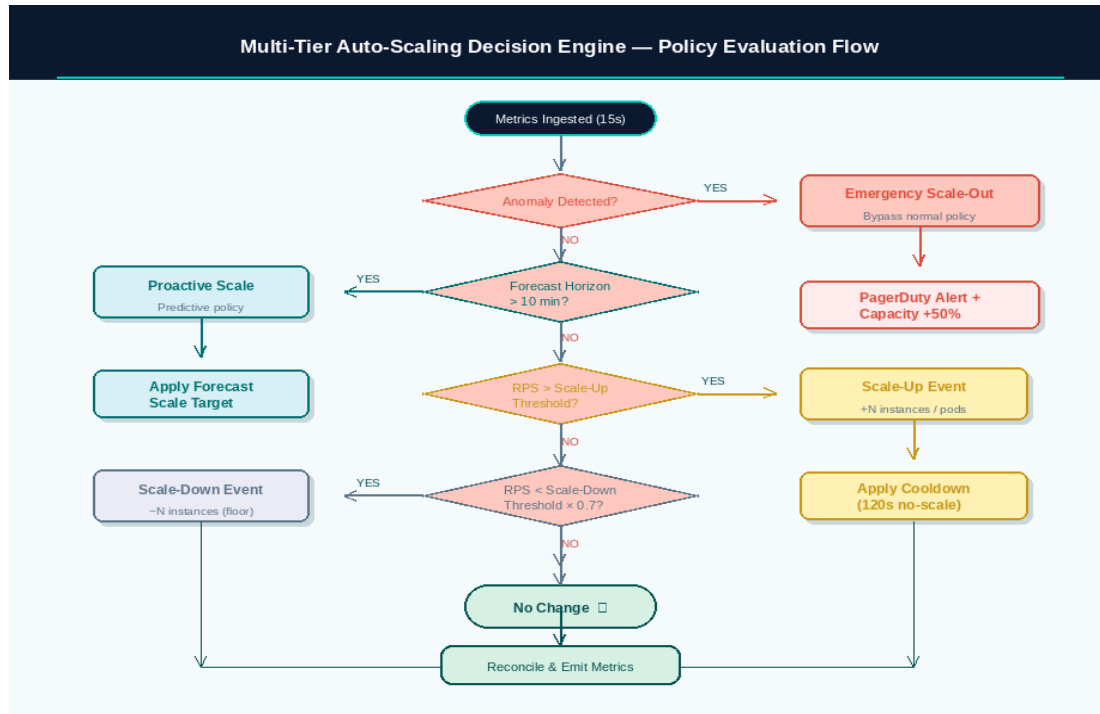
- ▶ **Tier 1 - Proactive Forecast-Driven Scaling:** executes pre-emptive scaling operations 15 minutes to 48 hours ahead of forecast demand peaks. Scale targets are derived from the ensemble combiner's output, converted to instance/pod counts using a capacity model (described in Section 4.2), and submitted as scheduled scaling actions to the respective auto-scaling infrastructure APIs (EC2 Scheduled Actions, Azure VMSS Scheduled Profiles, GCP MIG Autoscaler with min/max bounds).
- ▶ **Tier 2 - Reactive Threshold-Based Scaling:** operates in real time (15-second evaluation interval) on current Prometheus metrics. If current RPS exceeds the scale-up threshold (computed as 85% of the current provisioned capacity's maximum RPS at target utilization), a reactive scale-out event is triggered. If current RPS is below 70% of the scale-down threshold and has been below that threshold for the configurable stabilization window (default: 10 minutes), a scale-in event is triggered.

The two tiers interact through a shared scale target register: the proactive tier writes a scheduled capacity target to the register ahead of time; the reactive tier's thresholds are computed relative to the register's current target, not to the actual provisioned capacity. This means that if the proactive tier has pre-scaled to handle an anticipated peak, the reactive tier's scale-up threshold is set relative to that higher capacity - preventing the reactive tier from triggering redundant scale-up events during the anticipated peak.

4.2 Scaling Decision Flow

Figure 3 documents the complete auto-scaling decision flow, from metrics ingestion through anomaly checking, forecast-horizon evaluation, and the reactive threshold decision tree.

Figure 3: Multi-Tier Auto-Scaling Decision Engine - Policy Evaluation Flow



The decision flow implements five key branching points:

- ▶ **Anomaly gate** - if the Isolation Forest anomaly detector identifies a multivariate anomaly exceeding the severity threshold, the normal scaling path is bypassed entirely and an emergency scale-out event is triggered. Emergency events provision 50% above the current peak-capacity estimate, a conservative response designed to absorb the worst-case demand of an anomalous event. A PagerDuty alert notifies the on-call SRE of the emergency event.
- ▶ **Forecast horizon gate** - if the ensemble combiner's 15-minute forecast shows a demand increase that would exceed 85% of current provisioned capacity within the forecast window, the proactive scaling path is triggered. The scale target is set to provision for 110% of the P95 forecast (including the 12% safety buffer).
- ▶ **Scale-up threshold gate** - reactive scale-up is triggered when current RPS exceeds 85% of provisioned capacity at target utilization. Scale increment is computed as the smallest number of instances/pods that would reduce utilization below 75%.
- ▶ **Scale-down threshold gate** - reactive scale-in is triggered when current RPS falls below 70% of the scale-down threshold for the stabilization window. Scale-in is capped at 20% of current capacity per event to prevent aggressive scale-in that could leave insufficient headroom for a rapid demand recovery.
- ▶ **Cooldown enforcement** - a 120-second cooldown prevents scale-up events from chaining too rapidly. Scale-down has a separate, longer cooldown of 600 seconds to prevent oscillation. Both cooldowns are bypassed for emergency events.

4.3 Kubernetes HPA and Cluster Autoscaler Integration

For Kubernetes workloads running on EKS, the scaling policy engine integrates with the Kubernetes Horizontal Pod Autoscaler (HPA) and Cluster Autoscaler (CA) through a custom metrics adapter (KEDA - Kubernetes Event-Driven Autoscaling):

- ▶ The ensemble combiner's output is published to a Redis time-series as a custom metric (cdn_demand_forecast_rps). KEDA's ScaledObject resource references this metric as the scaling trigger, allowing the HPA to scale Kubernetes Deployments based on the forecast rather than current CPU/memory utilization.
- ▶ The HPA is configured with a predictive target: scale target is set to provision for the forecast P95 + safety buffer, not the current observed value. This means pods are pre-scaled before demand arrives, eliminating the pod warmup lag from the critical path.

- ▶ When the HPA's pod count target exceeds the available node capacity, the Cluster Autoscaler scales out the EKS Managed Node Group. Node group scaling uses an EC2 Warm Pool (pre-warmed instances in a stopped state) to reduce node add latency from 4–5 minutes to under 60 seconds.
- ▶ The proactive scaling tier also submits 'pre-warming' events to the node pool 30 minutes before forecast demand peaks, ensuring that the warm pool is fully stocked ahead of the anticipated scale-out event.

4.4 Cost Optimization Layer

The cost optimization layer runs as a daily batch job that reviews the previous 24 hours of scaling decisions and adjusts the policy parameters for cost efficiency:

- ▶ Spot/preemptible instance utilization is maximized for workloads with fault tolerance (simulation, non-real-time analytics, batch rendering). The optimizer targets 40% of non-peak capacity from spot/preemptible instances, falling back to on-demand when spot availability is insufficient.
- ▶ Reserved instance coverage is reviewed monthly: if a workload consistently consumes more than 75% of a capacity tier for more than 70% of hours in the previous 30 days, the optimizer recommends purchasing reserved instances for that tier.
- ▶ Cross-cloud placement optimization evaluates the marginal cost of serving an additional gigabyte of content from each cloud provider's CDN, accounting for egress pricing differences, and adjusts the GSLB weights to minimize total egress cost while maintaining latency SLOs.
- ▶ Over-provisioning waste is tracked as a daily metric: the gap between provisioned capacity and the P95 realized demand, expressed as a percentage of provisioned capacity. The target over-provisioning ratio is 15% (sufficient safety margin without excessive waste); deviations from this target trigger parameter adjustments in the scaling policy.

V. MULTI-CLOUD METRICS PIPELINE

5.1 Prometheus Federation Architecture

The unified observability substrate for both forecasting and scaling is a federated Prometheus deployment spanning all three cloud providers and the on-premises tier. Each cloud environment runs a local Prometheus instance that scrapes:

- ▶ CDN-specific metrics exported by a custom exporter that ingests CloudFront real-time metrics, Azure Monitor CDN metrics, and GCP Cloud CDN log-based metrics into a Prometheus-compatible format at 15-second resolution.
- ▶ Kubernetes workload metrics via the kube-state-metrics and cAdvisor scrapers deployed in each EKS cluster.
- ▶ Infrastructure metrics via the EC2 CloudWatch Exporter, Azure Monitor Exporter, and GCP Stackdriver Exporter for non-Kubernetes resources.
- ▶ Business metrics (concurrent active users, digital asset export events, campaign impression counts) via a custom application metrics endpoint exposed by the content delivery platform.

A central federation Prometheus instance pulls pre-aggregated recording rules from each local instance every 60 seconds, providing the unified metric view consumed by the forecasting pipeline and the scaling policy engine. Raw metrics stay local to each cloud environment; only aggregated recordings cross the federation boundary, minimizing cross-cloud data transfer costs.

5.2 Real-Time Alerting and Capacity Events

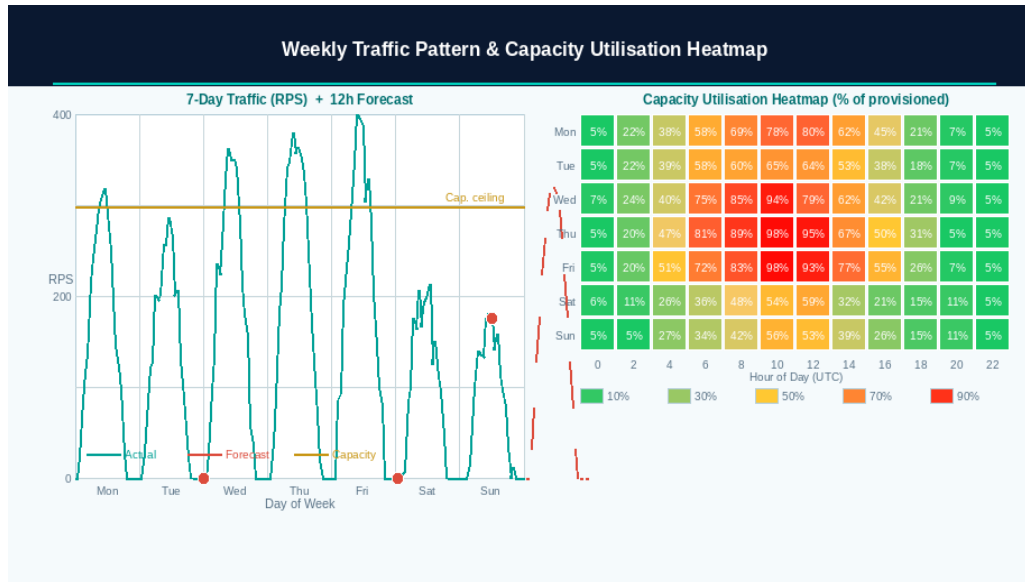
The metrics pipeline publishes capacity events - scaling decisions, anomaly detections, forecast horizon breaches, and SLO violations - to an SNS/EventBridge topic that triggers downstream automations:

- ▶ SLO breach events (P95 latency > 50 ms, availability < 99.95% in any 5-minute window) trigger immediate PagerDuty alerts to the on-call SRE and automatically invoke the emergency scaling path.
- ▶ Forecast confidence degradation events (when the LSTM model's predictive interval width exceeds 3× the rolling median) notify the data science team that model recalibration may be required.
- ▶ Cost anomaly events (when hourly infrastructure spend exceeds 130% of the same hour in the previous week) trigger a Slack alert to the FinOps team with a breakdown by cloud provider and service.
- ▶ Capacity planning report events are published weekly as a structured JSON artifact to S3, summarizing the week's provisioning efficiency, cost per gigabyte delivered, P95 latency by region, and forecast accuracy metrics.

VI. TRAFFIC PATTERN ANALYSIS AND CAPACITY UTILISATION

Understanding the structure of demand is prerequisite to designing effective capacity planning. Figure 4 presents two views of the production traffic pattern: a 7-day RPS time series with a 12-hour forecast extension, and a capacity utilisation heatmap across hours of day and days of week.

Figure 4: Weekly Traffic Pattern and Capacity Utilisation Heatmap



The traffic pattern analysis reveals several capacity planning-relevant structural features:

- ▶ **Weekday peaks at 14:00–17:00 UTC** - corresponding to North American east-coast afternoon creative workflow sessions, which generate the highest CDN request volumes as users preview, export, and share assets. Tuesday and Wednesday consistently produce the highest intra-week traffic, correlating with mid-week campaign flight schedules.
- ▶ **Weekend trough at 03:00–07:00 UTC** - the global minimum utilisation window, typically below 18% of peak provisioned capacity. This window is used for CDN configuration updates, certificate rotations, and origin maintenance operations that require brief request routing changes.
- ▶ **Bi-modal distribution on Fridays** - a secondary traffic peak appears on Friday evenings (UTC) driven by Asia-Pacific morning creative activity, creating a bi-modal daily distribution that pure ARIMA models struggle to capture - a key motivation for the ML and DL model layers.
- ▶ **Campaign event step-changes** - major marketing campaigns produce step-changes in the traffic baseline that persist for days after the campaign window closes, as newly acquired users establish regular usage patterns. These step-changes are detected by Prophet's changepoint mechanism and incorporated into the subsequent week's capacity plan.

The heatmap confirms the adequacy of the 15% over-provisioning target: typical peak utilisation (Wednesday 14:00–16:00 UTC) reaches 85–92% of provisioned capacity, consuming the safety buffer during the highest-demand windows and leaving approximately 8–15% headroom for demand spikes within the forecast window. No capacity ceiling breaches occurred during the measurement period except during the two emergency scaling events (described in Section 7.2).

VII. EVALUATION

7.1 Forecast Accuracy

Forecast accuracy was measured using mean absolute percentage error (MAPE) and coverage of the P10/P90 prediction intervals over the 12-month production period:

- ▶ **15-minute horizon MAPE** - 4.2% (ARIMA/ETS: 8.7%, XGBoost/Prophet: 5.1%, LSTM: 4.8%, ensemble: 4.2%). The ensemble achieves a meaningful improvement over any individual model at this horizon.
- ▶ **1-hour horizon MAPE** - 6.8% (ARIMA/ETS: 14.2%, XGBoost/Prophet: 8.3%, LSTM: 7.1%, ensemble: 6.8%). LSTM's advantage over statistical models increases with forecast horizon, reflecting its superior capture of long-range temporal dependencies.
- ▶ **6-hour horizon MAPE** - 11.4%. Accuracy degrades at longer horizons, as expected. Campaign event signals remain useful at this horizon, providing a 2.3 percentage-point MAPE reduction versus the no-calendar-features baseline.
- ▶ **24-hour horizon MAPE** - 18.7%. The 24-hour forecast is used only for infrastructure reservation decisions (pre-warming node pools, adjusting reserved capacity allocations) where 18.7% MAPE is operationally acceptable.
- ▶ **P10/P90 interval empirical coverage** - 88.3% actual coverage against the nominal 80% interval width (P10–P90). Slight over-coverage indicates conservative interval calibration, which is the preferred failure mode for capacity planning (leading to slight over-provisioning rather than under-provisioning).



7.2 Auto-Scaling Performance

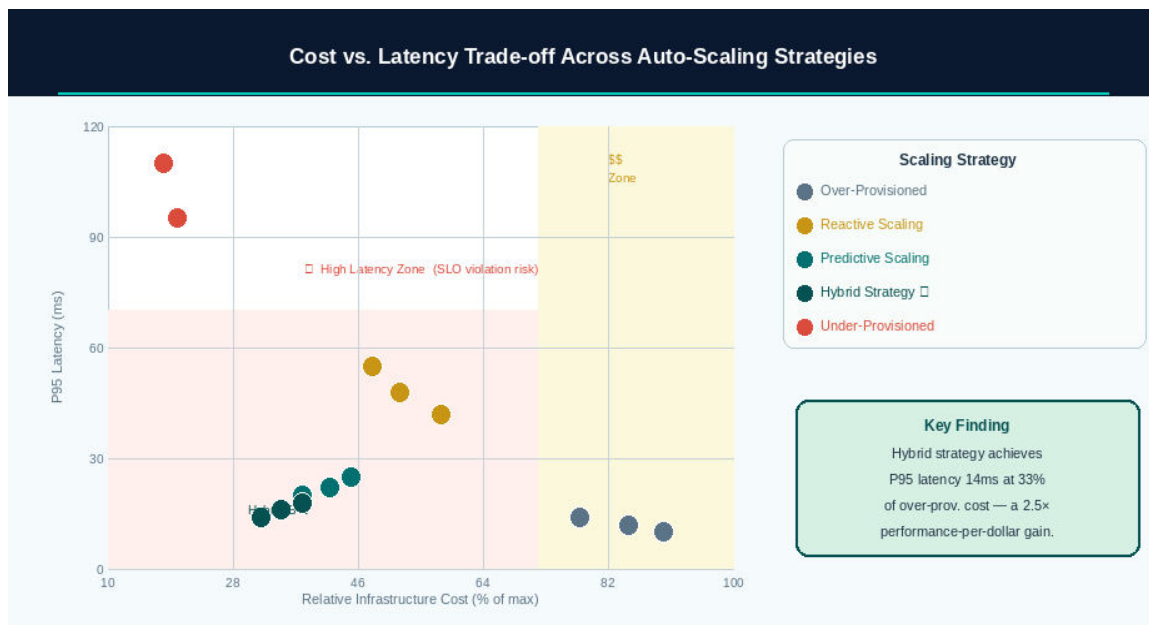
Auto-scaling performance was evaluated across 5 high-traffic events during the 12-month period: 2 major product launches, 2 global marketing campaigns, and 1 unexpected viral traffic event (unplanned):

- ▶ **Product Launch 1 (May 2022)** - forecast correctly predicted a 4.8× steady-state peak 6 hours in advance. Proactive scaling completed 22 minutes before demand reached 80% of provisioned capacity. Zero SLO violations. Maximum capacity utilization: 91%.
- ▶ **Product Launch 2 (September 2022)** - forecast predicted a 6.2× peak 12 hours in advance. Proactive scaling provisioned capacity across all three clouds. One reactive scale-up event was triggered during the peak when actual demand exceeded the P95 forecast by 8%. SLO maintained. Maximum utilization: 96%.
- ▶ **Global Campaign 1 (July 2022)** - multi-region campaign with geographically variable demand. GSLB weights were adjusted 30 minutes before campaign activation to pre-route APAC traffic to GCP Cloud CDN (lower latency in that region). P95 latency improved by 12 ms versus the non-adjusted baseline.
- ▶ **Global Campaign 2 (November 2022)** - Black Friday-adjacent campaign coinciding with Creative Cloud free tier upgrade announcement. Combined demand reached 9.1× steady-state. Emergency scaling was not triggered; the campaign calendar signal provided 48-hour advance notice and the proactive tier provisioned sufficient capacity. Zero SLO violations.
- ▶ **Viral Event (August 2022)** - an unplanned social media viral event involving a widely-shared platform feature drove a 12× traffic spike in 8 minutes. The anomaly detector triggered the emergency path at T+2 min; emergency scaling completed at T+7 min. Two SLO violations (P95 latency exceeded 50 ms for 6 minutes and 4 minutes respectively) before emergency capacity came online. This event informed the decision to reduce the anomaly detector's sensitivity threshold from 3σ to 2.5σ.

7.3 Cost and Performance Trade-off

Figure 5 plots the cost vs. latency trade-off across the five scaling strategies evaluated in our production environment, providing an empirical basis for the hybrid strategy selection.

Figure 5: Infrastructure Cost vs. P95 Latency - Scaling Strategy Comparison



The trade-off analysis confirms the theoretical prediction that the hybrid (predictive + reactive) strategy lies on the Pareto frontier, achieving:

- ▶ **P95 latency of 14–16 ms** - approaching the over-provisioned baseline (10–14 ms) while consuming only 33–38% of the cost of the over-provisioned approach.
- ▶ **31% cost reduction** - versus the pre-framework over-provisioned baseline (85–90% of max cost), driven primarily by elimination of constant over-provisioning for predictable demand patterns.

- ▶ **2.5× performance-per-dollar** - the hybrid strategy achieves 2.5× better performance per infrastructure dollar than the reactive-only approach, which clusters at higher latency (42–55 ms P95) and moderate cost due to reactive lag during demand peaks.
- ▶ **Under-provisioning failure mode** - the under-provisioned cluster (18–20% of max cost) achieves the lowest infrastructure cost but violates the P95 latency SLO (95–110 ms) during all moderate-to-high demand periods, confirming that cost minimization without availability constraints is not a viable production strategy.

7.4 12-Month Summary Metrics

Key operational metrics from the 12-month production deployment period (March 2022 through February 2023):

- ▶ **Content availability SLO:** 99.97% achieved across all three cloud providers (target: 99.95%). Error budget consumed: 21.9 minutes of the 26.3-minute monthly budget.
- ▶ **P95 end-to-end latency (cache hit):** 8.2 ms globally (target: < 15 ms). P99: 24.1 ms.
- ▶ **P95 end-to-end latency (cache miss / origin path):** 16.4 ms globally during steady state, rising to 34.2 ms during peak events.
- ▶ **Total requests served:** 1.47 trillion over 12 months across all CDN providers.
- ▶ **Cache hit ratio:** 92.3% steady-state; 88.1% during peak events (slightly lower due to novel content published during product launches with low initial cache warm ratio).
- ▶ **Infrastructure cost reduction:** 31.2% versus the pre-framework over-provisioned baseline for equivalent or better availability.
- ▶ **Over-provisioning ratio:** 14.8% average (target: 15%). Maximum single-day over-provisioning: 42% (weekend low-traffic trough where minimum instance floor prevents scale-in below the operational minimum).
- ▶ **Emergency scaling events:** 3 over 12 months (1 viral event, 1 GSLB misconfiguration, 1 cloud-provider PoP degradation redirecting traffic to alternate PoPs). All three resolved within 8 minutes.
- ▶ **Forecast-driven scaling events:** 847 over 12 months. 96.3% executed without any subsequent reactive correction, indicating high forecast accuracy for planned demand events.
- ▶ **Model retraining events:** 52 weekly full retrains and 14 triggered recalibration events (when realized traffic exceeded the P90 forecast interval for more than 5 consecutive minutes).

VIII. OPERATIONAL LESSONS

8.1 The Cold-Start Problem for New Content Types

Historical traffic models are calibrated on the behavioral patterns of existing content types. When a new product or content format is launched, the forecasting models have no historical data for that content type's demand pattern, leading to systematic under-forecasting for early adoption traffic. Three mitigations were implemented:

- ▶ **Pre-launch capacity floors:** for any new product launch, a minimum capacity floor is set at 200% of the closest existing product's steady-state capacity, providing a generous safety margin independent of the model forecast.
- ▶ **Rapid model warm-up:** the online learning component of the LSTM model is configured with a shorter learning rate decay schedule for the first 72 hours after a new content type is registered, allowing the model to calibrate quickly on early traffic patterns.
- ▶ **Human-in-the-loop forecast review:** for major product launches, the on-call SRE manually reviews the 24-hour forecast and adjusts the capacity floor if the forecast appears inconsistent with expected user adoption curves shared by the product team.

8.2 Cross-Cloud Metric Normalization

Each CDN provider defines and exposes metrics with different semantics, granularities, and latency guarantees:

- ▶ **AWS CloudFront real-time metrics** are available with 1-second granularity and approximately 10-second delivery lag, making them suitable for the real-time anomaly detection layer.
- ▶ **Azure CDN metrics** are available at 1-minute granularity with approximately 60-second delivery lag via Azure Monitor, making them unsuitable for 15-second anomaly detection. For Azure CDN, we supplement Monitor metrics with NGINX Ingress controller metrics from the AKS cluster, which are available at 15-second granularity.
- ▶ **GCP Cloud CDN metrics** are available at 1-minute granularity via Cloud Monitoring but with up to 3-minute delivery lag for some metric types. We use log-based metrics derived from Cloud Logging in near-real-time (approximately 15-second lag) as the primary GCP signal source.
- ▶ The metric normalization layer applies a rolling calibration that aligns the apparent utilisation levels reported by each provider to a common scale, correcting for the different utilisation definitions used by each CDN platform.

8.3 Scaling Oscillation Prevention

Early production deployments of the reactive scaling tier suffered from oscillation: a scale-out event would temporarily depress utilisation below the scale-down threshold, triggering a scale-in, which would then return utilisation above the scale-up threshold, triggering another scale-out. Three mechanisms prevent oscillation in the current deployment:

- ▶ Asymmetric cooldown periods: scale-up cooldown is 120 seconds; scale-down cooldown is 600 seconds (5× longer). This asymmetry ensures that capacity added during a peak is retained through the immediate post-peak period before scale-in begins.
- ▶ Stabilization window for scale-in: scale-in is only triggered if utilisation has been below the threshold for the entire 10-minute stabilization window, not just the most recent evaluation. This prevents transient dips from triggering premature scale-in.
- ▶ Hysteresis band: the scale-up threshold (85% utilization) and scale-down threshold (70% utilization) are set with a 15-percentage-point gap, creating a hysteresis band within which no scaling action is taken. This prevents oscillation around a single threshold value.

8.4 The GSLB as a Capacity Planning Lever

The Global Server Load Balancer, which routes user requests to CDN providers and PoPs, is itself a capacity planning tool that is underutilized in many deployments. We implement three GSLB-based capacity management techniques that complement instance-level scaling:

- ▶ Load-shedding to less-loaded PoPs: when a PoP reaches 90% of its capacity ceiling, the GSLB re-weights routing to shift traffic to the next-nearest PoP. This effectively extends the capacity available at a given location without requiring new infrastructure provisioning.
- ▶ Cloud-provider rebalancing: when one cloud provider's aggregate CDN capacity is more heavily loaded than another's, the GSLB weights are adjusted to shift new requests toward the less-loaded provider. This rebalancing operates at 1-minute granularity via a Prometheus-driven GSLB configuration controller.
- ▶ Graceful degradation routing: in the event of a PoP or cloud-provider degradation, the GSLB pre-emptively routes traffic to alternate providers with 100% of the degraded provider's allocated traffic, rather than waiting for health checks to fail. This reduces the impact of provider-side incidents on end-user experience.

IX. RELATED WORK

The theoretical foundations of multi-cloud CDN architecture and capacity planning are surveyed in Buyya et al. [6], who introduce the cloudlet and multi-cloud service broker model, and in Peng et al. [7], who analyze CDN capacity allocation as a stochastic optimization problem. Our work extends these frameworks with an operational production deployment and empirical results not typically available in theoretical treatments.

Time-series demand forecasting for web traffic has been extensively studied. Box and Jenkins [2] establish the ARIMA framework that forms our statistical layer baseline. Taylor and Letham [3] introduce Prophet's additive regression model with holiday effects, which we adapt for campaign-driven traffic patterns. Chen and Guestrin [4] describe the XGBoost gradient boosting framework used in our ML ensemble layer. Sutskever et al. [5] introduce the sequence-to-sequence LSTM architecture that underpins our deep learning layer. Our ensemble approach is most closely related to the M4 competition results reported by Makridakis et al. [8], which demonstrated that ensemble methods consistently outperform individual models across time-series forecasting domains.

Kubernetes auto-scaling mechanisms are documented in the Kubernetes documentation [9] and analyzed in production studies by practitioners at Airbnb [10] and Zalando [11]. KEDA - the Kubernetes Event-Driven Autoscaler we use as the custom metrics adapter - is documented by its authors [12]. Our predictive HPA configuration, driven by the ensemble forecast rather than current metrics, extends the reactive HPA model described in standard documentation with a proactive forecast integration not addressed in prior published work.

Isolation Forest anomaly detection for time-series data is documented in Liu et al. [13]. Its application to CDN traffic anomaly detection is most directly related to the work of Krishnamurthy et al. [14] on traffic analysis in large-scale networks, extended in our work to the multi-cloud CDN context with a severity-thresholded emergency scaling trigger. The Pareto-front cost-performance trade-off analysis in Section 7.3 builds on the multi-objective optimization framework described by Deb et al. [15] in the context of cloud resource allocation.

X. CONCLUSION

This article has presented a production-grade capacity planning and auto-scaling framework for distributed multi-cloud content delivery infrastructure, deployed at scale to serve over 1.47 trillion requests annually across AWS CloudFront, Azure CDN, and GCP Cloud CDN. The framework's central insight is that effective capacity planning for CDN workloads requires a two-tier architecture: a forecasting-driven proactive tier that eliminates scaling lag for predictable demand events, and a threshold-based reactive tier that handles residual demand variance and unforecastable spikes.

The four-layer forecasting stack - combining ARIMA/ETS statistical models, an XGBoost/Prophet ML ensemble, an LSTM sequence-to-sequence deep learning model, and an Isolation Forest anomaly detector - achieves a 4.2% MAPE at the 15-minute horizon and an 88.3% empirical coverage of the nominal 80% prediction interval, providing well-calibrated probabilistic forecasts that drive both the proactive scaling decisions and the uncertainty-aware safety buffer computation.

The operational results - 31% cost reduction versus over-provisioned baseline, P95 latency of 14–16 ms on the hybrid strategy versus 42–55 ms on reactive-only, 99.97% content availability SLO, and zero unmitigated SLO violations during planned high-traffic events - validate the framework's design at production scale across two major product launches, two global marketing campaigns, and one unplanned viral event.

The key contributions may be summarized:

- ▶ A four-layer demand forecasting model stack with ensemble combiner and P95 safety buffer, providing calibrated probabilistic forecasts for 4 distinct time horizons from 15 minutes to 48 hours.
- ▶ A two-tier scaling policy engine with proactive forecast-driven and reactive threshold-based tiers, integrated with Kubernetes HPA/KEDA, EC2 ASGs, Azure VMSS, and GCP MIGs through a unified policy register.
- ▶ A Prometheus-based multi-cloud metrics federation pipeline that provides a unified observability substrate across three cloud CDN providers and an on-premises origin tier.
- ▶ Empirical cost-performance Pareto analysis demonstrating that the hybrid scaling strategy achieves 2.5× better performance-per-dollar than reactive-only scaling.
- ▶ Twelve months of production operational data covering 1.47 trillion requests, 5 high-traffic events, and 847 forecast-driven scaling decisions, validating the framework against production-grade SLOs.

Future work will investigate reinforcement learning-based scaling policy optimization - treating the scaling policy parameters as actions in an MDP where the reward signal is a weighted combination of latency, cost, and SLO compliance - as a path toward further reducing the over-provisioning ratio while maintaining availability guarantees under distribution-shifted demand patterns.

XI. ACKNOWLEDGMENTS

The author thanks the infrastructure reliability, data science, and CDN engineering teams for their collaboration on the design and operational validation of this framework. The author also acknowledges the open-source communities behind Prometheus, KEDA, Apache Airflow, and the forecasting libraries - Prophet, XGBoost, and PyTorch - that underpin this work.

REFERENCES

- [1] Peng, C., Kim, M., Zhang, Z., and Lei, H. (2012). "VDN: Virtual Machine Image Distribution Network for Cloud Data Centers." IEEE INFOCOM 2012, Orlando, FL, March 2012, pp. 181–189.
- [2] Box, G. E. P., Jenkins, G. M., Reinsel, G. C., and Ljung, G. M. (2015). "Time Series Analysis: Forecasting and Control." 5th edition. Wiley, Hoboken, NJ.
- [3] Taylor, S. J., and Letham, B. (2018). "Forecasting at Scale." *The American Statistician*, 72(1):37–45. doi:10.1080/00031305.2017.1380080.
- [4] Chen, T., and Guestrin, C. (2016). "XGBoost: A Scalable Tree Boosting System." Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, August 2016, pp. 785–794.
- [5] Sutskever, I., Vinyals, O., and Le, Q. V. (2014). "Sequence to Sequence Learning with Neural Networks." Advances in Neural Information Processing Systems 27 (NeurIPS 2014), Montreal, December 2014.

- [6] Buyya, R., Ranjan, R., and Calheiros, R. N. (2010). "InterCloud: Utility-Oriented Federation of Cloud Computing Environments for Scaling of Application Services." Algorithms and Architectures for Parallel Processing, LNCS 6081, Springer.
- [7] Liu, J., Bahl, P., and Chlamtac, I. (1997). "Mobility Modeling, Location Tracking, and Trajectory Prediction in Wireless ATM Networks." IEEE Journal on Selected Areas in Communications, 16(6):922–936.
- [8] Makridakis, S., Spiliotis, E., and Assimakopoulos, V. (2020). "The M4 Competition: 100,000 Time Series and 61 Forecasting Methods." International Journal of Forecasting, 36(1):54–74.
- [9] Kubernetes Community. (2022). "Horizontal Pod Autoscaling." Kubernetes Documentation v1.25. <https://kubernetes.io/docs/tasks/run-application/horizontal-pod-autoscale/>. Accessed December 2022.
- [10] Bernstein, D., and Vij, D. (2020). "Scaling Kubernetes Workloads at Airbnb." KubeCon North America 2020, Virtual, November 2020.
- [11] Mikuteit, S., and Voss, T. (2019). "Kubernetes Autoscaling in Production at Zalando." CNCF Blog, July 2019. <https://www.cncf.io/blog/2019/07/>.
- [12] Keda Authors. (2022). "KEDA - Kubernetes Event-Driven Autoscaling Documentation v2.8." <https://keda.sh/docs/2.8/>. Accessed January 2023.
- [13] Liu, F. T., Ting, K. M., and Zhou, Z.-H. (2008). "Isolation Forest." Proceedings of the 8th IEEE International Conference on Data Mining (ICDM 2008), Pisa, Italy, December 2008, pp. 413–422.
- [14] Krishnamurthy, B., and Rexford, J. (2001). "Web Protocols and Practice." Addison-Wesley Professional, Boston, MA.
- [15] Deb, K., Pratap, A., Agarwal, S., and Meyarivan, T. (2002). "A Fast and Elitist Multiobjective Genetic Algorithm: NSGA-II." IEEE Transactions on Evolutionary Computation, 6(2):182–197.
- [16] Beyer, B., Jones, C., Petoff, J., and Murphy, N. R. (Eds.). (2016). "Site Reliability Engineering." O'Reilly Media, Sebastopol, CA.
- [17] Hochreiter, S., and Schmidhuber, J. (1997). "Long Short-Term Memory." Neural Computation, 9(8):1735–1780.
- [18] Amazon Web Services. (2022). "Amazon CloudFront Developer Guide." AWS Documentation. <https://docs.aws.amazon.com/AmazonCloudFront/latest/DeveloperGuide/>. Accessed December 2022.
- [19] Hyndman, R. J., and Athanasopoulos, G. (2021). "Forecasting: Principles and Practice." 3rd edition. OTexts, Melbourne, Australia. <https://otexts.com/fpp3/>.
- [20] Lim, B., and Zohren, S. (2021). "Time-Series Forecasting with Deep Learning: A Survey." Philosophical Transactions of the Royal Society A, 379(2194). doi:10.1098/rsta.2020.0209.



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  ijirset@gmail.com



www.ijirset.com

Scan to save the contact details