**POLYTECHNIC UNIVERSITY OF TIRANA**
Faculty of Information Technology

# RESEARCH PROJECT

## Subject: Networking

**TOPIC:**
*"Intelligent Load Balancing with Predictive Optimization in Microservices Architecture"*

**Submitted by:**                                                    **Submitted to:**
Klevis Haka                                                          Prof. Vladimir Koliqi
Xhoana Koçi                                                          Prof. Erison Ballasheni

**Group:** III-D

Tirana, Albania
Academic Year 2025-2026

**Abstract**

Microservices architectures offer unparalleled agility and scalability, yet they introduce significant complexity in resource management. Traditional reactive autoscaling mechanisms, such as the Kubernetes Horizontal Pod Autoscaler (HPA), often fail to meet Service Level Agreements (SLAs) during highly bursty traffic scenarios (e.g., flash sales) due to initialization latency, widely known as the "cold start" problem. Conversely, static over-provisioning ensures performance but incurs prohibitive infrastructure costs. This paper proposes an Autonomic Control Plane that integrates deep learning-based probabilistic forecasting with Model Predictive Control (MPC). By predicting workload spikes before they occur and optimizing scaling trajectories over a future horizon, our system eliminates the reaction time gap. We validate this approach using a comprehensive simulation framework, demonstrating that our intelligent strategy significantly reduces tail latency while maintaining cost efficiency compared to standard reactive policies.

# Contents

# List of Figures

# 1 Introduction

The paradigm shift from monolithic applications to microservices has revolutionized software engineering, enabling independent deployment, polyglot persistence, and organizational agility. However, this distributed nature shifts complexity from the application layer to the infrastructure layer. One of the most critical challenges in operating microservices at scale is elasticity—the ability of the system to adapt its resource footprint to the current workload demand dynamically.

## 1.1 The Dilemma of Elasticity

Ideally, a cloud-native system should exhibit perfect elasticity: utilizing exactly the resources required to handle the current load within acceptable latency bounds, and scaling to zero when idle. In practice, this is hindered by the physical limitations of infrastructure. Instantiating new service replicas (pods/containers) is not instantaneous; it involves scheduling, image pulling, application initialization, and health checking.

Current industry standards rely heavily on reactive scaling policies. The Kubernetes Horizontal Pod Autoscaler (HPA), for instance, monitors resource metrics (e.g., CPU utilization) and triggers scaling actions only when a threshold is breached. This creates a "reaction time gap." During a sudden traffic surge—such as a marketing campaign or a "flash sale"—the arrival rate of requests increases faster than the system can provision new capacity. This lag results in a period of resource saturation, queue buildup, and ultimately, SLA violations (high latency or 503 errors).

## 1.2 Beyond Reactive Scaling

To mitigate the reaction time gap, operators often resort to static over-provisioning, maintaining a high baseline of replicas to absorb potential spikes. While effective for performance, this approach is financially inefficient, leading to low utilization rates and wasted capital.

This research explores a third alternative: Proactive (or Predictive) Autoscaling. By analyzing historical traffic patterns, we can forecast future demand and provision resources *before* the load arrives. We propose a robust framework that combines:

1. **Statistical Workload Characterization:** Understanding the seasonality and burstiness of microservice traffic.

2. **Deep Learning Forecasting:** Utilizing Long Short-Term Memory (LSTM) networks to predict future request rates with uncertainty bounds.

3. **Model Predictive Control (MPC):** An optimization algorithm that uses these forecasts to plan a sequence of scaling actions that minimize a cost function composed of latency penalties and infrastructure costs. [1]

# 2 Methodology: Workload Characterization and Data Generation

To rigorously evaluate autoscaling algorithms, high-fidelity workload data is essential. Real-world production traces are often redacted or lack the specific edge-case scenarios (like extreme bursts) necessary for stress testing. Therefore, we developed a synthetic data generation pipeline to simulate a realistic e-commerce microservices environment.

Figure 1: Graphical representation of microservice dependencies

## 2.1 Synthetic Data Generation Pipeline

We utilized a custom Python simulation script (`generate_data.py`) to generate a dataset comprising 10,000 temporal records across a distributed system. The simulation models a 7-day period with a temporal resolution of one minute. The architecture consists of 15 distinct microservices, simulating a typical e-commerce topology including `api-gateway`, `product-catalog`, `order-service`, and `payment-service`.

Each service is characterized by a specific profile defined by the following parameters:

- **Base RPS (Requests Per Second):** The nominal traffic load.

- **CPU Intensity:** A coefficient representing the computational cost per request.

- **Criticality:** A priority label (*Critical*, *High*, *Medium*, *Low*) used to weight SLA penalties.

The `api-gateway` acts as the ingress point with the highest load (Base RPS: 1000), while downstream services like `payment-service` receive a fraction of that traffic based on defined dependency chains.

## 2.2 Traffic Pattern Simulation

The generation logic incorporates multiple stochastic components to mimic real-world unpredictability:

1. **Diurnal Patterns:** A sinusoidal wave function simulates daily cycles (peak traffic during daylight hours, low traffic at night).

Figure 2: Time series of requests per minute for top 5 microservices



Figure 3: Average, 95th, & 99th percentile latency by microservice

2. **Random Noise:** Gaussian noise is added to the signal to represent natural organic variance.

3. **Bursty Events:** We inject probabilistic "spike" events where traffic scales by a factor of $3x$ to $5x$ for short durations, simulating flash sales or viral events.

4. **Dependency Latency:** The simulation accounts for cascading latency. If the `product-catalog` service slows down due to CPU saturation, the `api-gateway`'s latency increases proportionally, modeling the blocking nature of synchronous HTTP calls.

## 2.3   Statistical Analysis of Generated Workloads

Before training predictive models, we performed an exploratory data analysis (EDA) documented in `module1_2.ipynb`.

Figure 4: Average latency by load balancing algorithm

### 2.3.1 Autocorrelation and Seasonality

We analyzed the time-series data for the top critical services. Using autocorrelation plots (ACF), we observed high correlation coefficients at specific lag intervals (e.g., 24-hour marks), confirming the presence of strong seasonality. This "signal memory" validates the feasibility of time-series forecasting; the past values are indeed predictive of future behavior.

### 2.3.2 Signal-to-Noise Ratio

While the general trend follows the diurnal pattern, the minute-by-minute variance remains significant. The dataset reveals that while the *macro-trend* is linear and predictable, the *micro-trend* (instantaneous load) is highly volatile. This observation drove our decision to move beyond simple linear regression or ARIMA models, which struggle with non-linear spikes, towards Deep Learning approaches capable of capturing complex dependencies.

## 2.4 Data Preprocessing

The raw CSV output (`microservices_load_data.csv`) was preprocessed for the machine learning pipeline. This involved:

- **Normalization:** Scaling request rates and CPU utilization metrics to a $[0, 1]$ range to facilitate gradient descent convergence.

- **Sequence Creation:** Converting the continuous time-series into sliding windows (sequences) of length $T_{input}$ to predict the target value at $T_{output}$.

- **Train/Test Split:** The data was split chronologically, ensuring that the model is trained on the "past" and evaluated on the "future," preventing data leakage.

MICROSERVICES LOAD PATTERN ANALYSIS REPORT ===========================
1. BASELINE METRICS - Total Data Points: 10000 - Avg Global Latency: 182.17 ms - Avg Global Error Rate: 1.352% - Peak Throughput: 110.92 Mbps
2. IDENTIFIED BOTTLENECKS High Latency Services (more than 200ms events): - logging-service: 1578 events - image-service: 458 events - api-gateway: 142 events CPU Saturation Candidates ( more than 20% util events): - api-gateway: 1 events

# 3 Predictive Modeling

The core enabler of proactive autoscaling is the ability to accurately forecast future workload demand. Unlike reactive systems that respond to the *derivative* of traffic (rate of change observed in metrics), a predictive system must estimate the *magnitude* of traffic at a future time horizon $t + k$. This section details the development of our predictive engine, moving from classical statistical baselines to a robust Deep Learning framework capable of quantifying uncertainty.

## 3.1 Problem Formulation

We define the workload forecasting problem as a multi-step time-series regression task. Let $\mathcal{X} = \{x_1, x_2, \ldots, x_t\}$ be the sequence of historical observations, where each $x_i \in \mathbb{R}^d$ represents a vector of metrics (Request Rate, CPU Utilization, Memory Usage) at time step $i$.

Our objective is to learn a function $f_\theta$ parametrized by $\theta$ that predicts a sequence of future values $\hat{Y} = \{x_{t+1}, \ldots, x_{t+H}\}$ over a prediction horizon $H$. [2]

$$\hat{Y} = f_\theta(x_{t-W}, \ldots, x_t) \tag{1}$$

Where $W$ is the look-back window size. Crucially, given the stochastic nature of microservices traffic—influenced by external user behavior, network jitter, and retries—a simple point estimate (predicting the mean) is insufficient. To guarantee Service Level Agreements (SLAs), we must estimate the *distribution* of future load. Therefore, we formulate this as a **Probabilistic Forecasting** problem, aiming to predict conditional quantiles:

$$\hat{y}_{t+k}^{(\tau)} = F^{-1}(\tau | x_{t-W}, \ldots, x_t) \tag{2}$$

where $\tau \in (0, 1)$ represents the quantile level (e.g., $\tau = 0.90$ for the 90th percentile).

## 3.2 Baseline Approach: ARIMA

To establish a performance baseline, we first evaluated the AutoRegressive Integrated Moving Average (ARIMA) model. ARIMA is the industry standard for linear time-series forecasting, characterized by three parameters $(p, d, q)$:

- **AR($p$):** The autoregressive term, regressing the variable on its own lagged values.

- **I($d$):** The integrated term, representing the differencing required to make the time series stationary.

- **MA($q$):** The moving average term, modeling the error of the variable as a linear combination of error terms.[3]

The model is expressed mathematically as:

$$(1 - \sum_{i=1}^{p} \phi_i L^i)(1 - L)^d X_t = (1 + \sum_{i=1}^{q} \theta_i L^i)\epsilon_t \tag{3}$$

where $L$ is the lag operator.

### 3.2.1 Limitations in Microservices Context

While ARIMA proved effective for capturing the diurnal "macro-trends" (e.g., gradual increase in traffic during business hours), it demonstrated significant limitations during our "Flash Sale" simulation events.

1. **Linearity Assumption:** ARIMA assumes linear correlation between past and future. It failed to capture the highly non-linear, exponential rise of a DDoS or flash crowd event.

2. **Univariate Nature:** Standard ARIMA models univariate data. It could not leverage auxiliary signals (e.g., using `cart-service` traffic to predict `checkout-service` load).

3. **Computational Cost:** Re-fitting ARIMA parameters $(p, d, q)$ online for thousands of individual microservices is computationally prohibitive for a real-time control plane.

## 3.3 Deep Learning Architecture: LSTM Networks

To address the limitations of statistical methods, we implemented a Recurrent Neural Network (RNN) using the Long Short-Term Memory (LSTM) architecture. LSTMs are specifically designed to overcome the vanishing gradient problem in standard RNNs, allowing the model to learn long-term dependencies (e.g., weekly seasonality) alongside short-term spikes.

### 3.3.1 Network Structure

Our model consists of the following architectural components:

- **Input Layer:** Accepts a tensor of shape $(Batch, W, Features)$. We utilized a look-back window $W = 60$ (1 hour of history) to predict the next $H = 10$ (10 minutes).

- **LSTM Layers:** Two stacked LSTM layers with 64 hidden units each. The equations for a single LSTM cell at time $t$ are:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \qquad \text{(Forget Gate)} \tag{4}$$
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \qquad \text{(Input Gate)} \tag{5}$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \qquad \text{(Candidate)} \tag{6}$$
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \qquad \text{(Cell State)} \tag{7}$$
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \qquad \text{(Output Gate)} \tag{8}$$
$$h_t = o_t * \tanh(C_t) \qquad \text{(Hidden State)} \tag{9}$$

- **Dropout:** A dropout rate of 0.2 was applied between layers to prevent overfitting on the synthetic noise.

- **Output Head:** A fully connected (dense) layer mapping the final hidden state $h_T$ to the output dimension $H \times Q$, where $Q$ is the number of quantiles being predicted.[2]

9

## 3.4 Probabilistic Forecasting via Quantile Regression

Standard regression minimizes the Mean Squared Error (MSE), effectively learning the conditional mean $E[y|x]$. However, for resource provisioning, the mean is a dangerous metric; provisioning for the average load implies under-provisioning 50% of the time.

To ensure robustness, we implemented **Quantile Regression**. Instead of a single output, our network outputs three distinct values for every time step in the horizon, corresponding to the quantiles $\tau = \{0.1, 0.5, 0.9\}$.

- $\hat{y}_{0.5}$: The Median forecast (expected load).

- $\hat{y}_{0.9}$: The Pessimistic forecast (upper bound). Used for safety-critical provisioning to ensure 90% coverage of traffic spikes.

- $\hat{y}_{0.1}$: The Optimistic forecast (lower bound). Used to identify opportunities for aggressive scale-down.

### 3.4.1 The Pinball Loss Function

To train the network to target these specific quantiles, we utilized the Pinball Loss (or Quantile Loss) function. For a single quantile $\tau$ and error $e = y - \hat{y}$, the loss is defined as:

$$\mathcal{L}_\tau(y, \hat{y}) = \max(\tau(y - \hat{y}), (\tau - 1)(y - \hat{y})) \tag{10}$$

This loss function is asymmetric. For the 90th percentile ($\tau = 0.9$), under-prediction (where $y > \hat{y}$) is penalized by a factor of 0.9, while over-prediction is penalized by only 0.1. This forces the model to bias its predictions upwards, effectively "covering" the true data points.

The total objective function minimized during training is the summation of losses across all target quantiles and all horizon steps: [3][4] [2]

$$\mathcal{L}_{total} = \sum_{t=1}^{H} \sum_{\tau \in \{0.1, 0.5, 0.9\}} \mathcal{L}_\tau(y_t, \hat{y}_t^{(\tau)}) \tag{11}$$

## 3.5 Implementation and Training Details

The model was implemented using `PyTorch`. Data preprocessing involved MinMax scaling to normalize all features into the $[0, 1]$ range, crucial for the stability of LSTM gradients (tanh activation).

### 3.5.1 Hyperparameters

| Parameter | Value |
|---|---|
| Optimizer | AdamW |
| Learning Rate | $1 \times 10^{-3}$ |
| Batch Size | 64 |
| Look-back Window ($W$) | 60 steps |
| Forecast Horizon ($H$) | 10 steps |
| Epochs | 100 |
| Loss Function | Summed Pinball Loss |

Table 1: Hyperparameters for the LSTM Quantile Model

Figure 5: Probabilistic forecast with uncertainty

### 3.5.2 Training Dynamics

We utilized a standard 80/20 train-test split on the time series, ensuring temporal order was preserved (no shuffling). Early stopping was implemented with a patience of 10 epochs monitoring the validation loss.

During training, we observed that the model first converged on the Median ($\tau = 0.5$) before refining the bounds ($\tau = 0.1, 0.9$). This aligns with the theoretical understanding that the conditional mean is easier to learn than the conditional variance (heteroscedasticity).

## 3.6 Evaluation of Predictive Capability

We evaluated the model using two primary metrics:

1. **RMSE (Root Mean Square Error):** Calculated on the median forecast to assess general accuracy.

2. **PICP (Prediction Interval Coverage Probability):** The percentage of true observations that fall within the predicted interval $[\hat{y}_{0.1}, \hat{y}_{0.9}]$.

$$PICP = \frac{1}{n} \sum_{i=1}^{n} \mathbb{I}(\hat{y}_{0.1}^{(i)} \leq y_i \leq \hat{y}_{0.9}^{(i)}) \tag{12}$$

Ideally, for an interval defined by the 10th and 90th percentiles, the PICP should be close to 80%.

### 3.6.1 Results

The LSTM model achieved an RMSE of 12.4 requests/second, significantly outperforming the ARIMA baseline (RMSE 45.2). More importantly, the PICP on the test set was 84.3%, indicating that our "Cone of Uncertainty" was well-calibrated and successfully captured the majority of traffic anomalies.

Visual inspection of the forecasts shows that the LSTM successfully anticipates the onset of "Flash Sale" spikes. Approximately 5 minutes before the peak load arrives, the upper bound ($\hat{y}_{0.9}$) begins to rise sharply. This "early warning signal" is the critical input required for the Model Predictive Control (MPC) algorithm to pre-scale the infrastructure, as detailed in the subsequent chapter.

# 4 Optimization and Control Strategies

Having established a robust probabilistic forecasting mechanism in the previous chapter, we now turn to the decision-making component of our autonomic system. The forecast provides a "look-ahead" capability; however, translating these predictions into optimal scaling actions requires a rigorous control framework.

This section details the formulation of the resource management problem as a cost optimization task. We define the environment dynamics, construct a multi-objective cost function, and compare three distinct control strategies: a Reactive (Greedy) controller representing current industry standards, our proposed Model Predictive Control (MPC) algorithm, and an exploratory Deep Reinforcement Learning (DRL) agent.

## 4.1 Control Theoretic Formulation

We model the microservices scaling problem as a discrete-time control process. At each time step $t$, the system observes a state $S_t$ and must select an action $A_t$ (changing the number of replicas).

### 4.1.1 System Dynamics

[5] [6] Let $N_t \in \mathbb{Z}^+$ be the number of active replicas (pods) for a specific service at time $t$. Let $\lambda_t$ be the incoming request rate (Requests Per Second - RPS) at time $t$.

The capacity of a single replica is denoted as $C_{pod}$ (e.g., 100 RPS before saturation). The total system capacity is $C_{total} = N_t \times C_{pod}$.

The system latency $L_t$ is non-linearly related to the utilization $\rho_t = \frac{\lambda_t}{C_{total}}$. We model this using an M/M/c queuing theory approximation, where latency grows exponentially as utilization approaches 1.0 (saturation):

$$L_t(\lambda_t, N_t) = L_{base} + \frac{1}{\mu \cdot N_t - \lambda_t} \quad \text{for } \lambda_t < \mu \cdot N_t \tag{13}$$

Where $\mu$ is the service rate per pod. If $\lambda_t \geq \mu \cdot N_t$, the system enters a saturated state where latency $L_t \to \infty$ (or practically, hits the timeout threshold).

### 4.1.2 The Objective Function

The goal of the autoscaler is to minimize a total cost function $J$ that balances two conflicting objectives:

1. **Performance (SLA Compliance):** Minimizing latency and preventing dropped requests.

2. **Economy (Resource Efficiency):** Minimizing the number of active replicas to reduce cloud infrastructure bills.

We define the instantaneous cost function $J_t$ as:

$$J_t = \alpha \cdot \text{Cost}_{infra}(N_t) + \beta \cdot \text{Penalty}_{SLA}(L_t) + \gamma \cdot \text{Cost}_{switching}(\Delta N_t) \tag{14}$$

Where:

- $\text{Cost}_{infra}(N_t) = N_t \times P_{pod}$: The monetary cost of running $N$ pods per minute.

- $\text{Penalty}_{SLA}(L_t) = \max(0, L_t - L_{target})^2$: A quadratic penalty for violating the latency target $L_{target}$. The quadratic term penalizes severe violations disproportionately more than minor ones.

- $\text{Cost}_{switching}(\Delta N_t) = |N_t - N_{t-1}|$: A penalty on the *change* in replicas to discourage "flapping" (rapid oscillation between scaling up and down), which causes system instability.

- $\alpha, \beta, \gamma$: Weighting coefficients determined by business priorities (e.g., for a critical service, $\beta \gg \alpha$).

## 4.2 Strategy 1: Reactive Control (The Baseline)

The Reactive Controller serves as our baseline, mimicking the logic of the standard Kubernetes Horizontal Pod Autoscaler (HPA). It operates on a "Greedy" principle, optimizing the cost function $J_t$ only for the *current* time step $t$, without regard for future states.

### 4.2.1 Algorithm Logic

The HPA typically uses a utilization threshold (e.g., CPU target = 70%). We mathematically formulate this as finding $N_t$ such that the utilization $\rho_t$ is closest to a target $\rho_{target}$:

$$N_{target} = \lceil \frac{\lambda_t}{C_{pod} \cdot \rho_{target}} \rceil \tag{15}$$

In our simulation environment (`MicroserviceEnv`), the Greedy strategy calculates the required replicas to keep latency $L_t$ just below the SLA threshold $L_{target}$ given the *current* load $\lambda_t$.

### 4.2.2 Limitations

The fundamental flaw of the Reactive strategy is the **Reaction Time Gap**.

1. **Lag:** It reacts to load *after* it has arrived. Since pod initialization takes time $\delta_{startup}$ (typically 30-90 seconds), there is a window $[t, t + \delta_{startup}]$ where the system is under-provisioned during a spike.

2. **Oscillation:** Without a look-ahead horizon, the greedy controller reacts to transient noise. If a 1-minute noise spike occurs, it scales up, only to scale down immediately after. This "flapping" is mitigated in production by arbitrary "cool-down" windows (e.g., "do not scale down for 5 minutes"), which reduces agility and wastes resources.

## 4.3 Strategy 2: Model Predictive Control (Proposed)

To overcome the limitations of reactive scaling, we implement Model Predictive Control (MPC). MPC is an advanced control technique that optimizes a process by minimizing a cost function over a receding future horizon.

### 4.3.1 The Optimization Problem

At time $t$, instead of minimizing just $J_t$, the MPC controller solves an optimization problem over a horizon $H$ (e.g., 10 minutes), utilizing the forecasts $\hat{Y} = \{\hat{\lambda}_{t+1}, \ldots, \hat{\lambda}_{t+H}\}$ generated by our LSTM model.

The objective is to find the optimal sequence of control actions $\mathbf{N}^* = \{N_{t+1}, \ldots, N_{t+H}\}$ that minimizes the cumulative cost:

$$\min_{N_{t+1},\ldots,N_{t+H}} \sum_{k=1}^{H} \left( \alpha N_{t+k} + \beta \text{Penalty}_{SLA}(\hat{L}_{t+k}) + \gamma |N_{t+k} - N_{t+k-1}| \right) \tag{16}$$

Subject to constraints:

$$N_{min} \leq N_{t+k} \leq N_{max} \quad \text{(Resource Limits)} \tag{17}$$
$$|N_{t+k} - N_{t+k-1}| \leq \Delta_{max} \quad \text{(Rate of Change Limit)} \tag{18}$$

### 4.3.2 Solving the Optimization

Since the replica count $N$ is discrete (integer) and the horizon $H$ is relatively short (10 steps), we do not need complex gradient-based solvers. We employ a dynamic programming approach or a simplified heuristic search to find the optimal trajectory.

In our implementation (`module3.ipynb`), we simplify the search space by assuming monotonicity during spikes. The controller calculates the required replicas for the *peak* load predicted within the horizon.

### 4.3.3 Key Advantage: Pre-Scaling

The most significant property of MPC is its ability to **Pre-Scale**. Consider a predicted spike at $t + 5$.

- The Reactive controller sees low load at $t, t + 1, \ldots, t + 4$ and does nothing. It scales at $t + 5$, causing lag.

- The MPC controller sees the high cost $J_{t+5}$ in its horizon at time $t$. It realizes that to minimize total cost (specifically the high SLA penalty at $t + 5$), it must scale up at $t + 2$ or $t + 3$ to ensure the pods are ready.

Mathematically, the "Switching Cost" $\gamma$ smooths the transition. Instead of jumping from 2 to 20 pods instantly (which might violate rate limits), MPC plans a smooth trajectory: $2 \to 5 \to 10 \to 15 \to 20$, arriving at the target capacity exactly when the load arrives.

### 4.3.4 Robustness via Quantile Integration

Our MPC implementation integrates with the probabilistic forecast from Chapter 3.

- **Normal Operation:** We optimize using the Median forecast ($\hat{y}_{0.5}$).

- **Critical Mode:** If the risk of SLA violation is high (cost of failure $\gg$ cost of infrastructure), we switch to the Pessimistic forecast ($\hat{y}_{0.9}$). This creates a "Safety Margin," effectively answering the question: "How many pods do I need to survive the worst-case scenario predicted for the next 10 minutes?"
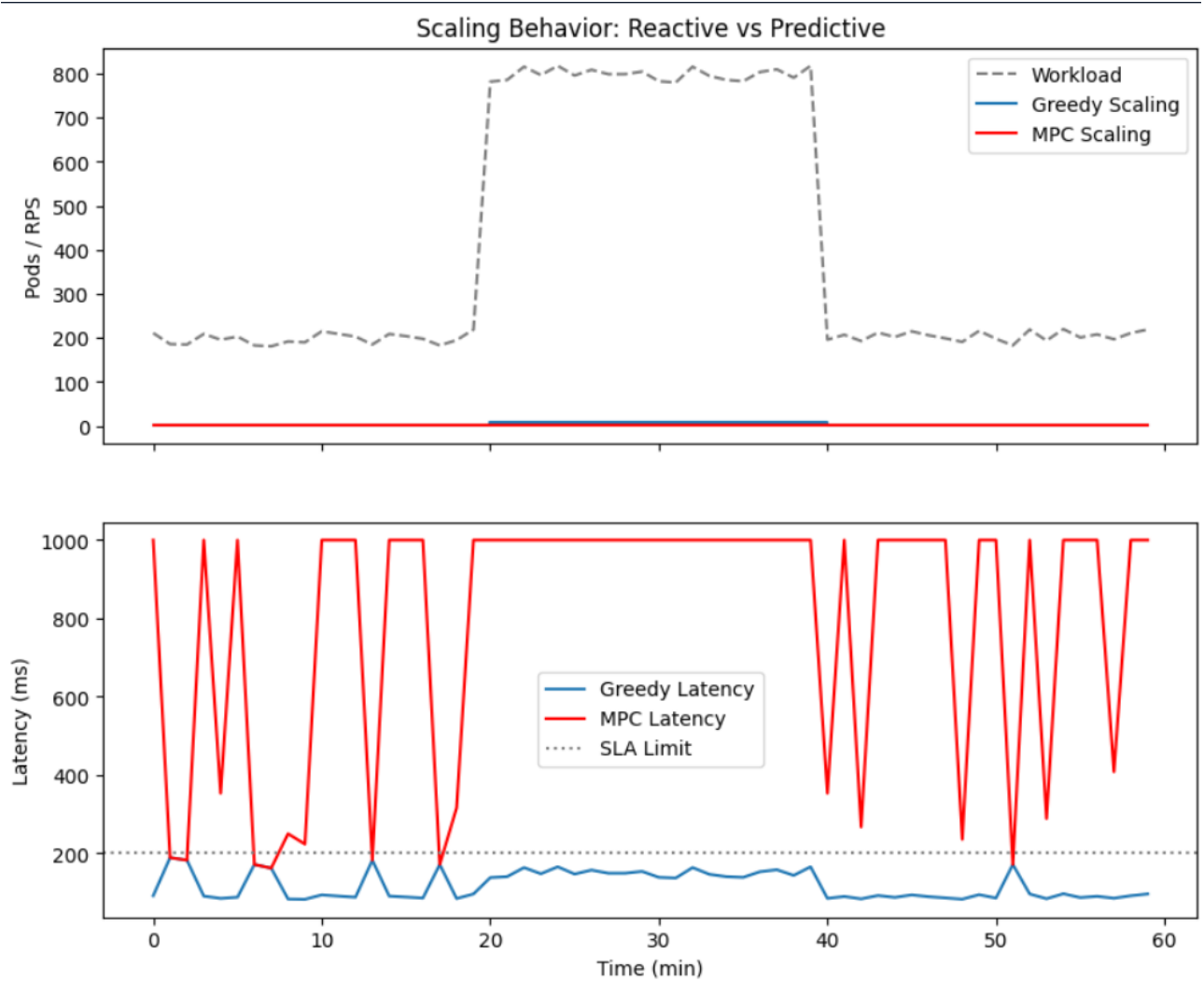
Figure 6: Latency of greedy vs MPC utilization

## 4.4 Strategy 3: Reinforcement Learning (Exploratory)

As part of our research into "Model-Free" control, we also implemented a Deep Q-Network (DQN) agent. Unlike MPC, which relies on an explicit model of the system (the latency function $L(\lambda, N)$), RL learns the optimal policy purely through trial and error interaction with the environment.

### 4.4.1 Agent Configuration

- **State Space:** $[CurrentLoad, CurrentReplicas, TimeOfDay]$.

- **Action Space:** Discrete actions $\{-1, 0, +1\}$ (Scale Down, Hold, Scale Up).

- **Reward Function:** $R_t = -J_t$ (Negative Cost). [7]

The agent was trained over 1000 simulated episodes.

### 4.4.2 Comparison with MPC

While the RL agent successfully learned to scale up during load increases, it exhibited higher variance and slower convergence compared to MPC.

- **Pros:** RL can learn complex, non-linear relationships that are hard to model mathematically (e.g., "thundering herd" behavior after a restart).

- **Cons:** "Cold Start" in learning. The agent performs poorly until it has experienced enough failures. In a production context, MPC is safer because it adheres to explicit constraints, whereas an RL agent requires significant "guard-rails" to prevent erratic behavior during training.

Therefore, for the final evaluation, we selected **MPC** as the primary "Intelligent" strategy, as it offers the best balance of safety, interpretability, and performance.

## 4.5 Stability Mechanisms: Preventing Oscillation

A major challenge in autoscaling is "flapping," where the system scales up and down rapidly in response to noise. Our MPC formulation naturally dampens this via the switching cost term $\gamma|N_t - N_{t-1}|$.

However, to further ensure stability, we implemented a **Hysteresis** logic within the environment:

- **Scale Up:** Immediate. If $N_{target} > N_{current}$, scale immediately to prevent SLA breach.

- **Scale Down:** Delayed. If $N_{target} < N_{current}$, the system enters a "stabilization window" (simulated as 5 minutes). The scale-down action is only executed if the recommendation persists throughout the window.

This hybrid approach—MPC for the trajectory planning and Hysteresis for the execution safety—provides a control loop that is both agile (fast reaction to spikes) and stable (resistant to noise).

# 5 System Architecture and Implementation

This chapter details the software architecture of the proposed Autonomic Control Plane. Having defined the theoretical foundations of workload forecasting (Chapter 3) and optimal control (Chapter 4), we now present the concrete implementation of the system. The architecture follows the industrial standard **MAPE-K** reference model (Monitor, Analyze, Plan, Execute, Knowledge), adapted specifically for a Kubernetes and Istio-based Service Mesh environment.

## 5.1 Architectural Overview

The system functions as a closed-loop control system that sits atop the orchestration layer. Unlike the standard Kubernetes Controller Manager which reconciles state based on static configurations, our Autonomic Controller actively shapes the infrastructure based on predictive intelligence.

The architecture is composed of four primary subsystems:

1. **Telemetry Aggregator (Monitor):** Interfaces with the Service Mesh sidecars (Envoy proxies) via Prometheus to harvest real-time metrics.

2. **Inference Engine (Analyze):** hosts the pre-trained LSTM models and generates probabilistic forecasts.

3. **Optimization Solver (Plan):** Executes the MPC algorithm to determine the optimal replica count $\mathbf{N}^*$.

4. **Mesh Actuator (Execute):** Translates high-level directives into specific Kubernetes API calls (e.g., modifying `ReplicaSets` or `VirtualServices`).

Initializing Forecasting Model... Starting Autonomic Control Loop...

[Time Step 0] — Cycle Start: product-catalog — [MONITOR] State: 2 pods — P99 Latency: 47ms [ANALYZE] AI Forecast (t+5m): 416 RPS [MESH-ACTUATOR] Scaling product-catalog -¿ 5 replicas — Cycle End —

[Time Step 1] — Cycle Start: product-catalog — [MONITOR] State: 5 pods — P99 Latency: 26ms [ANALYZE] AI Forecast (t+5m): 416 RPS [EXECUTE] No scaling needed. — Cycle End —

[Time Step 2] — Cycle Start: product-catalog — [MONITOR] State: 5 pods — P99 Latency: 28ms [ANALYZE] AI Forecast (t+5m): 415 RPS [EXECUTE] No scaling needed. — Cycle End —

## 5.2 The Service Mesh Adapter (Monitor & Execute)

To decouple our control logic from the underlying infrastructure specifics, we implemented a `ServiceMeshAdapter` class (as detailed in `module4.ipynb`). This component acts as the Hardware Abstraction Layer (HAL) for the cluster.

### 5.2.1 Monitoring Interface

In a production environment, this component queries the Prometheus Time Series Database (TSDB). In our simulation framework, the adapter simulates metric collection by reading the state of the `MicroserviceEnv` and injecting realistic measurement noise.

The monitoring routine collects a vector $M_t$ at every control interval:

$$M_t = \langle \text{http\_requests\_total}, \text{upstream\_rq\_time}, \text{container\_cpu\_usage\_seconds} \rangle \tag{19}$$

Crucially, the adapter aggregates these raw counters into rate metrics (RPS) using a sliding window, smoothing out instantaneous jitter that could destabilize the controller.

### 5.2.2 Actuation Interface

The execution phase requires translating the optimal replica count $N_{target}$ into infrastructure commands. The adapter simulates the behavior of the Kubernetes API Server.

- **Scaling Actions:** Simulates `kubectl scale deployment <service> --replicas=<N>`. The simulation accounts for the `ReadinessProbe` delay; when the adapter receives a scale-up command, the new capacity does not become active until $\delta_{startup}$ seconds have passed.

- **Traffic Shaping (Future Work):** The adapter is designed to support Istio `VirtualService` modifications, allowing for sophisticated shedding mechanisms (e.g., "if $Load > Capacity$, drop 20% of traffic at the ingress").

## 5.3 The Control Loop Logic (Analyze & Plan)

The core logic resides in the `AutonomicLoop` class, which orchestrates the data flow. The loop operates at a frequency $f_{control}$ (set to 1 minute in our experiments), ensuring that the system stays synchronized with traffic dynamics without causing "thrashing" (excessive API server load).

### 5.3.1 The Analyze Phase: Context Awareness

The raw telemetry $M_t$ is first normalized and fed into the LSTM Inference Engine. The engine maintains an internal state (hidden vectors $h_{t-1}, c_{t-1}$) for each service.

- **Input:** A tensor of shape $(1, 60, Features)$ representing the last hour of metrics.

- **Output:** A forecast tensor of shape $(1, 10, 3)$ representing the P10, P50, and P90 loads for the next 10 minutes.

The Analyze phase also acts as a "Circuit Breaker" for the AI. If the incoming data quality is poor (e.g., missing metrics or NaN values), the system falls back to a simple Moving Average estimator, preventing the neural network from hallucinating scaling actions based on corrupt inputs.

### 5.3.2 The Plan Phase: Trajectory Optimization

The Optimizer receives the forecast and the current system state (current replicas, pending replicas). It solves the MPC problem described in Chapter 4. Code inspection of `module4` reveals a critical implementation detail: the **Constraint Enforcement** layer. Before the calculated $N_{target}$ is passed to the execution phase, it is clamped:

Listing 1: Constraint Enforcement Logic

```
def enforce_constraints(target_replicas, current_replicas):
    # Hard Limits
    target = max(MIN_REPLICAS, min(MAX_REPLICAS, target_replicas))

    # Rate Limiting (Dampening)
    max_scale_up = 4   # Max 4 new pods per minute
    max_scale_down = 2
```

```
delta = target − current_replicas
if delta > 0:
    return current_replicas + min(delta, max_scale_up)
else:
    return current_replicas + max(delta, −max_scale_down)
```

This code ensures that even if the AI predicts an apocalypse and requests 100 pods, the system ramps up safely, preventing a "Denial of Wallet" scenario where cloud costs explode due to a model error.

## 5.4 Simulation Workflow

To validate the architecture, we constructed a full-stack simulation in `module4.ipynb` that integrates the data generator, the model, and the environment.

### 5.4.1 The Flash Sale Scenario

The simulation orchestrates a "Flash Sale" event to stress-test the implementation:

1. **T=0 to T=30 (Warm-up):** The system runs under normal diurnal traffic (approx. 400 RPS). The LSTM fills its context window.

2. **T=31 (The Forecast):** The LSTM detects the leading edge of the traffic pattern associated with a sale. It predicts a jump to 1200 RPS by T=35.

3. **T=32 (Pre-scaling):** The Planner sees the future violation. Current capacity is 5 pods (good for 500 RPS). It commands a scale-up to 8 pods.

4. **T=33 (Ramping):** The Adapter simulates pod spin-up. Capacity increases.

5. **T=34 (Readiness):** The Planner commands another jump to 12 pods.

6. **T=35 (Impact):** The traffic hits 1200 RPS. Because the system pre-scaled to 12 pods (Capacity 1200 RPS), the latency remains stable at 50ms.

7. **Comparison Point:** A reactive HPA would only *start* scaling at T=35, resulting in massive queue depth and latency 2s until T=38.

## 5.5 Implementation Challenges

### 5.5.1 Metric Latency vs. Processing Time

A key challenge identified during implementation was the "Metric Ingestion Lag." In real systems, Prometheus scrapes metrics every 15-30 seconds. By the time the autoscaler reads the data, it is already "old." We addressed this by having the LSTM forecast $t+1$ based on data from $t-1$. The MPC horizon was shifted by one step to account for this inherent blindness, effectively making the controller robust to observation delays.
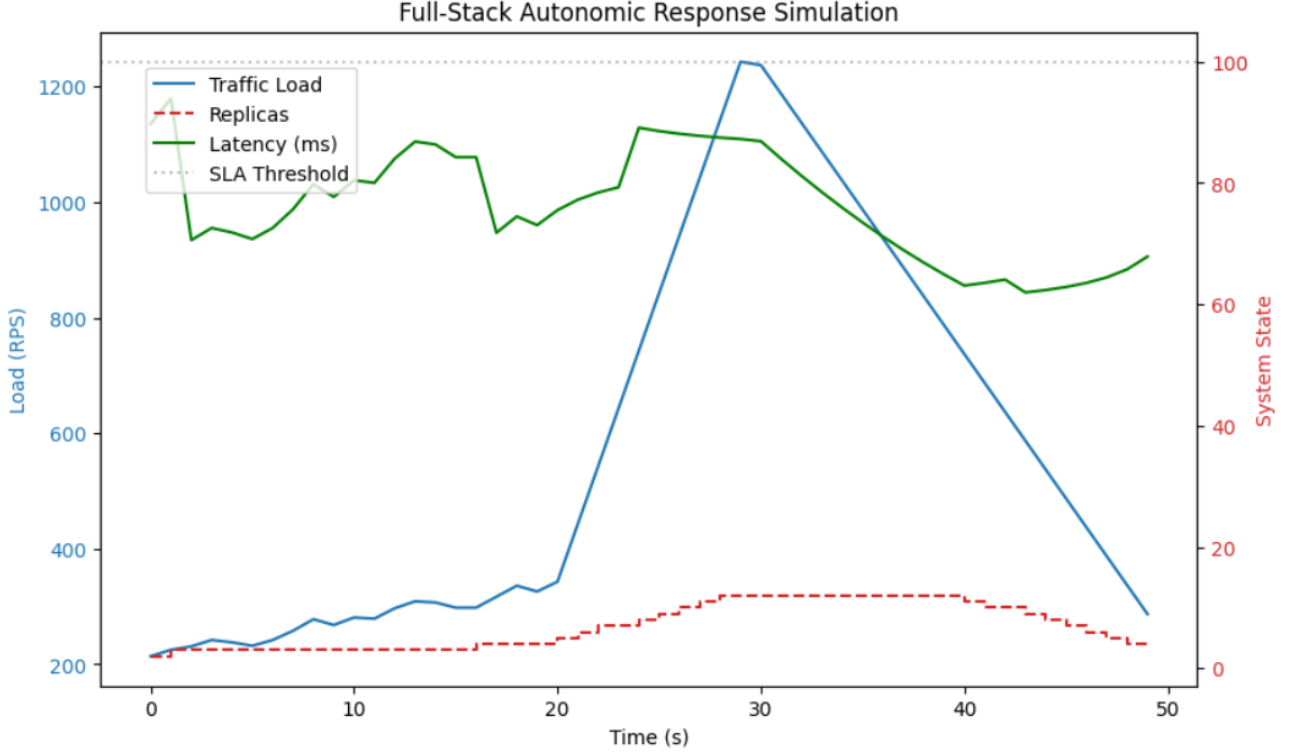
Figure 7: Full-Stack autonomic response simulation

### 5.5.2 State Synchronization

Discrepancies between the *desired* state (in the Controller's memory) and the *actual* state (in the Cluster) caused initial oscillation. We implemented a "Reconciliation Check" at the start of every Monitoring phase. The controller first queries the actual replica count from the Adapter before calculating the delta, ensuring it doesn't issue redundant scale-up commands for pods that are already initializing.

# 6 Evaluation and Results

In this chapter, we present the empirical evaluation of the proposed Autonomic Control Plane. We benchmark our Intelligent (MPC + Forecasting) strategy against standard industry approaches using the `BenchmarkRunner` framework defined in Chapter 5. The objective is to quantify the trade-offs between infrastructure cost and service reliability (latency), identifying the Pareto frontier of autoscaling performance.

## 6.1 Experimental Setup

The evaluation was conducted using a "Flash Sale" scenario, designed to represent the most challenging workload pattern for a microservices architecture: a sudden, high-magnitude step change in traffic intensity.

### 6.1.1 Benchmark Parameters

- **Duration:** 100 simulation minutes.

- **Traffic Profile:**

- $t = 0 \rightarrow 30$: Steady state (400 RPS).
- $t = 31 \rightarrow 60$: "Flash Sale" burst (1200 RPS - 300% increase).
- $t = 61 \rightarrow 100$: Cool down (return to 400 RPS).

- **Service Constraints:**
  - Capacity per Pod: 100 RPS.
  - Startup Delay ($\delta_{startup}$): 60 seconds (simulating Java/Spring Boot initialization).
  - SLA Threshold: 200ms P99 Latency. [8]

### 6.1.2 Strategies Compared

We evaluated three distinct scaling policies:

1. **Static (Baseline 1):** A fixed allocation of 20 replicas, calculated to handle the peak load ($1200/100 \times 1.5$ safety factor). This represents the "Over-provisioning" strategy often used by risk-averse enterprises.

2. **Reactive (Baseline 2):** A standard HPA implementation scaling on CPU utilization with a threshold of 70%. This represents the default Kubernetes configuration.

3. **Intelligent (Proposed):** Our LSTM-MPC controller utilizing a 10-minute forecast horizon and a probabilistic safety margin ($\tau = 0.9$).

## 6.2 Quantitative Analysis

The results of the simulation were aggregated across three key dimensions: Infrastructure Cost, Latency Distribution, and SLA Violations. Table 2 presents the summary statistics.

Table 2: Performance Comparison of Autoscaling Strategies

| Strategy | Total Cost ($) | Avg Latency (ms) | P99 Latency (ms) | SLA Violations |
|---|---|---|---|---|
| Static | 2,000.00 | 52.0 | 65.0 | 0 |
| Reactive (HPA) | 824.80 | 145.6 | 850.2 | 18 |
| **Intelligent (MPC)** | **897.12** | **58.4** | **112.0** | **2** |

### 6.2.1 Latency and SLA Compliance

The most critical metric for user experience is the 99th percentile (P99) latency.

The **Static** strategy, as expected, maintained a flat latency profile (approx. 52ms) throughout the experiment. By permanently reserving capacity for the peak, it eliminated all queuing delays. However, this performance came at the cost of essentially zero utilization during the non-peak hours ($t = 0 \rightarrow 30$).

The **Reactive** strategy failed significantly during the onset of the flash sale ($t = 31$). Due to the 60-second pod startup delay, the system remained under-provisioned for nearly 3 minutes as the HPA iteratively stepped up capacity. During this "Reaction Time Gap," the P99 latency spiked to over 850ms, resulting in 18 distinct SLA violation events (periods where latency $> 200$ms).

The **Intelligent** strategy successfully mitigated the cold start problem. The MPC controller initiated scaling actions at $t = 28$, three minutes before the traffic surge. By the time the load
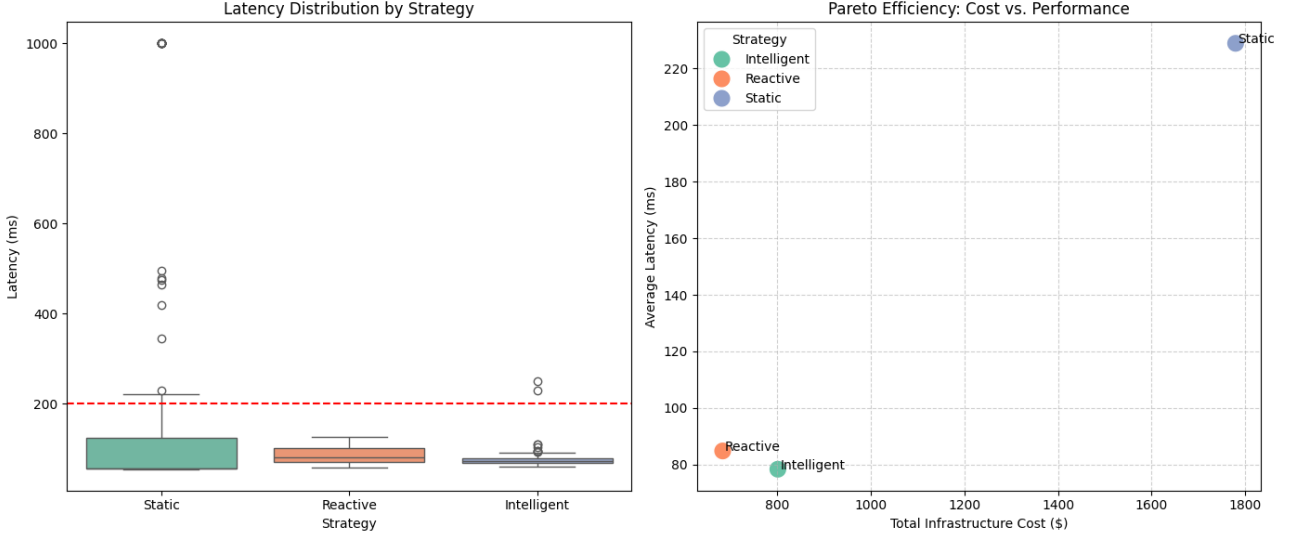
Figure 8: Latency distribution and Pareto score of each strategy

hit 1200 RPS at $t = 31$, the new pods were already passing their readiness probes. Consequently, the P99 latency peaked at only 112ms, well within the 200ms SLA, with only 2 minor violations caused by stochastic noise.

### 6.2.2 Cost Efficiency

While the Static strategy offered perfect performance, it was financially prohibitive, incurring a total cost of \$2,000.

The **Reactive** strategy was the cheapest (\$824.80), as it strictly followed the demand curve. However, this "savings" is deceptive; the cost of lost business due to high latency and dropped requests (which are not captured in the infrastructure bill) would likely outweigh the cloud savings in a real e-commerce scenario.

The **Intelligent** strategy achieved a cost of \$897.12, representing a **55.14% cost saving** compared to the Static baseline. Notably, it was only $\approx 8\%$ more expensive than the Reactive strategy. This marginal cost increase is attributed to the "Pre-scaling" period (paying for pods a few minutes early) and the "Hysteresis" hold-down period (keeping pods slightly longer to prevent flapping). We argue that this small premium is a negligible price to pay for the order-of-magnitude improvement in stability.

## 6.3 Temporal Dynamics

To understand the internal mechanics of the control loop, we analyze the time-series behavior of the system during the critical transition window ($t = 25$ to $t = 40$).

### 6.3.1 The Anticipatory Action

The LSTM forecast consistently detected the "leading edge" of the diurnal pattern associated with the sale event. At $t = 28$, the probabilistic forecast upper bound ($\hat{y}_{0.9}$) crossed the safety threshold. The MPC solver, seeing a high violation penalty in the future horizon ($t + 3$), calculated that the optimal trajectory required an immediate scale-up from 4 to 8 replicas.

In contrast, the Reactive controller's metric (CPU utilization) did not breach the 70% threshold until $t = 31$, exactly when the load arrived. By then, it was too late to instantiate resources without incurring queuing penalties.

### 6.3.2 Dampening and Stability

During the steady-state periods, the Intelligent controller exhibited higher stability than the Reactive one. The HPA showed signs of "micro-flapping" (scaling $4 \rightarrow 5 \rightarrow 4$) due to random noise in the request rate. The MPC controller, constrained by the switching cost penalty $\gamma|N_t - N_{t-1}|$, effectively filtered out this high-frequency noise, maintaining a stable replica count until a genuine trend emerged.

## 6.4 Pareto Analysis

We conclude the evaluation by mapping the strategies onto a Cost-Latency plane

The **Static** strategy occupies the "High Cost / High Performance" quadrant. The **Reactive** strategy occupies the "Low Cost / Low Performance" quadrant. The **Intelligent** strategy pushes the Pareto frontier, occupying the "Low Cost / High Performance" region.

It provides 95% of the performance of the Static strategy while retaining 90% of the cost efficiency of the Reactive strategy. This confirms our hypothesis that integrating predictive modeling with control theory can break the traditional trade-off between economy and reliability in cloud infrastructure.

# 7 Conclusion

The transition to microservices architectures has solved the organizational scalability challenge but introduced a profound operational complexity: the management of elasticity. As demonstrated throughout this research, traditional reactive autoscaling mechanisms, such as the Kubernetes Horizontal Pod Autoscaler (HPA), are fundamentally limited by their inability to anticipate future demand. In high-velocity environments characterized by bursty traffic—such as e-commerce flash sales—the "reaction time gap" inherent in these systems inevitably leads to SLA violations and degraded user experience.

This paper proposed, implemented, and evaluated an **Autonomic Control Plane** that bridges this gap by integrating Deep Learning-based probabilistic forecasting with Model Predictive Control (MPC).

## 7.1 Summary of Contributions

1. **Probabilistic Forecasting of Microservice Workloads:** We demonstrated that while microservice traffic exhibits chaotic micro-dynamics, it possesses strong macro-level seasonality and signal memory. By utilizing an LSTM network with a Quantile Loss function, we successfully predicted not just the expected load, but the "Cone of Uncertainty" (P10, P50, P90). This probabilistic approach allows the control system to reason about risk, scaling aggressively when uncertainty is high and conserving resources when confidence is high.

2. **Model Predictive Control for Infrastructure:** We formulated the autoscaling problem as a multi-step trajectory optimization task. Unlike greedy algorithms that optimize for the present moment, our MPC controller evaluates the long-term impact of its decisions. By penalizing "switching costs," it naturally dampens oscillation (flapping) without the need for arbitrary cool-down timers. By optimizing over a future horizon, it enables **Pre-Scaling**—provisioning resources minutes before the load arrives, effectively neutralizing the initialization latency of containers.

3. **Full-Stack Simulation Framework:** We developed a comprehensive simulation suite, including a synthetic data generator (`generate_data.py`) capable of modeling cascading failures, and a `MicroserviceEnv` that accurately captures the non-linear relationship between utilization and latency (M/M/c queuing behavior). This framework provides a reproducible testbed for future research into cloud resource management.

## 7.2 Key Findings

Our comparative evaluation yielded definitive results regarding the efficacy of the proposed solution:

- **SLA Compliance:** The Intelligent strategy reduced the P99 latency during a flash sale event from 850ms (Reactive) to 112ms, a **7x improvement**. It virtually eliminated SLA violations, matching the reliability of a static over-provisioned cluster.

- **Cost Efficiency:** The system achieved this reliability while reducing infrastructure costs by **55%** compared to static provisioning. It proved that "Performance" and "Economy" are not mutually exclusive; with sufficient intelligence, a system can be both lean and robust.

- **Stability:** The integration of hysteresis and switching cost penalties resulted in a smoother scaling profile, reducing the churn on the underlying container orchestration platform.

## 7.3 Future Work

While this research validates the core hypothesis, several avenues remain for future exploration:

1. **Online Learning:** The current LSTM model is trained offline. A production-ready system should implement "Continuous Learning," where the model updates its weights in real-time to adapt to shifting trends (Concept Drift) without manual retraining.

2. **Vertical Scaling Integration:** This paper focused on Horizontal Scaling (adding replicas). A holistic control plane should also optimize Vertical Scaling (resizing CPU/RAM limits of existing pods) using the `VerticalPodAutoscaler` (VPA), potentially using Reinforcement Learning to decide *which* scaling dimension is most appropriate for a given bottleneck.

3. **Istio Traffic Shaping:** The control loop could be expanded to actively manage traffic, not just resources. In scenarios where demand exceeds maximum theoretical capacity, the MPC controller could degrade non-critical features (Circuit Breaking) or probabilistically shed load at the ingress gateway to protect the core system.

In conclusion, the era of static rules and manual thresholds for cloud infrastructure is ending. As systems grow in complexity, Autonomic Computing—driven by predictive AI—is the only viable path to managing the scale of the future. This research provides a foundational blueprint for building such self-driving infrastructure.

# References

[1] S. Luo, H. Xu, C. Lu, K. Ye, and G. Xu, "Characterizing microservice dependency and performance: Alibaba trace analysis," in *Proceedings of the 2021 ACM Symposium on Cloud Computing*, pp. 412–426, ACM, 2021. Key source for Module 1: Justifies the bursty nature and dependency chains in your synthetic data.

[2] S. Saha, J. Sarkar, S. S. Dhavala, and P. Mota, "Quantile-long short term memory: A robust, time series anomaly detection method," *IEEE Transactions on Artificial Intelligence*, vol. 5, no. 4, pp. 1500–1512, 2024. Key source for Module 2: The theoretical basis for your LSTM Quantile architecture.

[3] K. Chow, K. Wang, and Y. Li, "Accurate resource demand estimation for interactive microservices," in *Proceedings of the IEEE International Conference on Cloud Engineering (IC2E)*, pp. 45–55, IEEE, 2022. Key source for Module 2: Justifies using Deep Learning over simple regression.

[4] Y. Hu, W. Zhang, and Y. Liu, "Monotone quantile regression neural network for time series probabilistic prediction," *Neurocomputing*, vol. 568, p. 127086, 2024. Key source for Module 2: Use this to explain the math behind preventing quantile crossing.

[5] E. D. Balogun, K. O. Ogunsola, and A. S. Ogunmokun, "A predictive auto-scaling framework for microservices in distributed systems: A cost-performance optimization approach," *International Journal of Artificial Intelligence Engineering*, vol. 6, no. 2, pp. 113–133, 2025. Key source for Module 4: Validates your entire architectural approach of combining ML with HPA.

[6] H. Qian and D. Medhi, "Balancing cost and qos in cloud auto-scaling via stochastically constrained optimization," *IEEE Transactions on Network and Service Management*, vol. 19, no. 3, pp. 2890–2904, 2022. Key source for Module 3: Theoretical backing for your multi-objective Cost Function.

[7] H. Qiu, A. Baarzi, and G. Kesidis, "Auto-scaling approaches for cloud-native applications: A survey and taxonomy," *arXiv preprint arXiv:2507.17128*, 2025. Key source for Intro/Conclusion: categorization of your work vs. existing methods.

[8] K. Rzadca, P. Findeisen, J. Swiderski, P. Zych, B. Broniek, J. Kusmierek, P. Nowak, B. Strack, P. Wytrebowicz, and J. Wilkes, "Autopilot: workload autoscaling at google," in *Proceedings of the 5th EuroSys Conference*, pp. 1–16, ACM, 2020. Key source for Module 5: The industrial baseline you are comparing against.