



School of Information Technology and
Engineering at the ADA University



School of Engineering and Applied Science
at the George Washington University

PREDICTIVE SCALING AND LOAD BALANCING FOR KUBERNETES-BASED
MICROSERVICES

A Thesis

Presented to the Graduate Program of Computer Science and Data Analytics
of the School of Information Technology and Engineering
ADA University

In Partial Fulfillment
of the Requirements for the Degree
Master of Science in Computer Science and Data Analytics
ADA University

By
Emil Suleymanli

April, 2025

THESIS ACCEPTANCE

This Thesis by: Emil Suleymanli

Entitled: *Predictive Scaling and Load Balancing for Kubernetes-Based Microservices*

has been approved as meeting the requirement for the Degree of Master of Science in Computer Science and Data Analytics of the School of Information Technology and Engineering, ADA University.

Approved:

(Adviser)

(Date)

(Program Director)

(Date)

(Dean)

(Date)

ABSTRACT

Kubernetes has emerged as the standard platform for managing microservices at scale, offering robust orchestration capabilities. However, ensuring optimal performance under dynamic and often predictable workload fluctuations remains a significant challenge. Traditional autoscaling mechanisms, such as the Horizontal Pod Autoscaler (HPA) rely on reactive policies that adjust resources based on current metrics like CPU utilization. While effective in many cases, reactive scaling often lags behind sudden traffic surges, leading to temporary service degradation or resource inefficiency. This thesis addresses these limitations by proposing a predictive autoscaling framework for Kubernetes-based microservices that integrates machine learning-based forecasting with intelligent load balancing.

The proposed solution leverages a Long Short-Term Memory (LSTM) neural network trained on twelve months of real-world microservice load data. The model forecasts short-term workload trends, enabling the system to proactively adjust pod counts before demand peaks occur. In parallel, a custom load balancing mechanism was developed to distribute traffic more efficiently based on runtime pod metrics such as CPU usage and response time, ensuring that scaled-out resources are utilized effectively.

An experimental Kubernetes cluster was set up to evaluate the predictive scaling approach against the standard HPA under realistic load patterns, including sharp end-of-month traffic surges observed in the banking sector of Azerbaijan. Results show that the predictive autoscaler achieved a mean absolute percentage error (MAPE) under 10% during normal periods and around 12–15% during peak salary-day surges. Compared to HPA, the predictive system reduced 95th percentile response times by up to 37% during load spikes, maintained full throughput without request drops, and triggered fewer, better-timed scaling actions. CPU utilization stayed within safer bounds, avoiding the saturation seen under reactive scaling.

This work demonstrates that predictive autoscaling can significantly enhance the resilience and efficiency of Kubernetes-managed microservices. By combining accurate load forecasting with intelligent traffic distribution, the system improves user experience and infrastructure utilization. While challenges such as prediction errors and model retraining remain, the results highlight the practical benefits of integrating machine learning into cloud-native scaling strategies. Future work can extend this approach by exploring hybrid models that combine predictive insights with reinforcement learning or by refining load balancing strategies to optimize service quality during unpredictable demand fluctuations further.

Table of Contents

1	INTRODUCTION	1
1.1	Definition of the Problem	1
1.2	Objective of the Study	1
1.3	Significance of the Problem.....	2
1.4	Review of Significant Research.....	2
1.5	Assumptions and Limitations	3
2	REVIEW OF THE LITERATURE	4
2.1	Kubernetes-Based Microservices Architecture.....	4
2.2	Auto-Scaling in Kubernetes (HPA, VPA, KEDA)	5
2.3	Limitations of Reactive Scaling Techniques	6
2.4	Load Balancing Algorithms for Microservices.....	7
2.5	Predictive Scaling Approaches Using Machine Learning	9
2.6	Summary of Research Gaps.....	11
3	RESEARCH METHODOLOGY.....	12
3.1	Overall System Architecture.....	12
3.2	Data Collection and Preprocessing	13
3.3	LSTM-Based Load Prediction Model.....	16
3.4	Custom Load Balancer Design	17
3.5	Evaluation Setup (Kubernetes Cluster, Tools, Metrics)	18
3.6	Methodological Justification.....	19
4	RESEARCH RESULTS AND ANALYSIS OF RESULTS	21
4.1	Model Training Results (LSTM Prediction Accuracy)	21
4.2	Predictive Scaling Performance vs. HPA	22
4.3	Resource Usage Analysis (CPU, Memory, Latency).....	24
4.4	Limitations and Observations	27
4.5	Summary of Key Findings	29
5	SUMMARY AND FUTURE WORK	30
5.1	Summary of Research Outcomes.....	30
5.2	Contribution to the Field.....	31
5.3	Limitations and Challenges.....	32
5.4	Suggestions for Future Work.....	32
	REFERENCES	34
	APPENDICES	36
	A.1 Kubernetes YAML Files.....	36
	A.2 LSTM Model Hyperparameters	39
	A.3 Sample Dataset and Chart	40

<i>A.3.1 Sample Workload Dataset</i>	<i>40</i>
<i>A.3.2 Chart: Actual vs Predicted Request Rate Around Month-End Surge.....</i>	<i>40</i>

LIST OF FIGURES

No	Figure Caption	Page
3.1	System Architecture of Predictive Scaling Solution	13
4.1	Actual vs. Predicted Workload over a Two-Day Period in the Test Dataset	21
4.2	Comparison of Predictive Autoscaling vs. Kubernetes HPA over a 48-Hour Load Test	23
A.3.2	Actual vs. Predicted Request Rate Around Month-End Surge (Zoomed View)	40

LIST OF TABLES

No	Figure Caption	Page
4.1	LSTM Model Prediction Accuracy on Test Dataset	22
A.3.1	Sample Request Rate Dataset Around End-of-Month Surge	40

LIST OF ABBREVIATIONS

Abbreviation	Explanation
CPU	Central Processing Unit
HPA	Horizontal Pod Autoscaler
LSTM	Long Short-Term Memory
ML	Machine Learning
MAPE	Mean Absolute Percentage Error
RMSE	Root Mean Square Error
API	Application Programming Interface
YAML	Yet Another Markup Language
IPVS	IP Virtual Server
DNS	Domain Name System

1 INTRODUCTION

Kubernetes has become the de-facto standard for orchestrating microservices applications due to its robust capabilities in service discovery, load balancing, self-healing, and declarative configuration. Despite these features, effectively managing application performance under dynamic workloads remains challenging. Predictive scaling – proactively adjusting resources based on anticipated demand – and intelligent load balancing are emerging as key strategies to maintain performance and cost efficiency in Kubernetes environments. This thesis addresses the problem of how to predictively scale microservices and balance loads in Kubernetes clusters to meet performance goals while avoiding resource waste.

1.1 Definition of the Problem

Microservices architecture consists of many small, independent services that communicate via APIs. Each microservice can be scaled individually, enabling fine-grained control of resources. Traditional scaling in Kubernetes often relies on reactive policies, such as the Horizontal Pod Autoscaler (HPA) adjusting replicas based on current CPU usage, or the Vertical Pod Autoscaler (VPA) adjusting pod resources based on observed needs. Reactive scaling, however, only responds after load changes occur, potentially leading to latency spikes or resource shortages before scaling kicks in. The problem is further compounded in real-world scenarios – for instance, in Azerbaijan, online services experience significant load spikes on the last working day of each month as employees receive salaries and rush to use banking or payment services. These predictable surges often overwhelm reactive scaling mechanisms, causing slowdowns or outages. Therefore, the core problem is twofold: (1) how to predict impending load increases (e.g., month-end spikes) and scale up in advance, and (2) how to balance traffic across microservice instances effectively so that no single instance or node becomes a bottleneck.

1.2 Objective of the Study

The objective of this study is to develop a comprehensive predictive scaling framework for Kubernetes-based microservices and evaluate its effectiveness compared to standard reactive methods. Specifically, we aim to:

- **Predict Workloads:** Use machine learning (particularly Long Short-Term Memory networks, LSTM) to forecast future load (e.g., HTTP request rates or CPU usage) of microservices with high accuracy, enabling proactive scaling.
- **Design a Predictive Autoscaler:** Implement an autoscaling controller that uses these predictions to scale pods ahead of time, mitigating latency during load spikes. This predictive autoscaler will be benchmarked against Kubernetes' built-in HPA in terms of response times, throughput, and resource utilization.
- **Custom Load Balancing:** Develop or configure a load balancing mechanism that complements predictive scaling by evenly distributing traffic among scaled-out pods using efficient algorithms (like round-robin or least-connections) This ensures that extra pods provisioned by predictive scaling relieve load on existing pods.
- **Evaluate Performance:** Set up a realistic Kubernetes cluster and test microservice to measure prediction accuracy of the LSTM model, scaling agility (how quickly and effectively the system responds to predicted load), and resource usage (CPU, memory,

latency) under predictive vs. reactive scaling. The evaluation will highlight the benefits and any overheads of predictive scaling.

By achieving these objectives, this study seeks to demonstrate that predictive scaling guided by machine learning can improve microservice resilience during known demand surges (like month-end traffic spikes) and maintain efficient resource usage better than purely reactive approaches.

1.3 Significance of the Problem

Effective autoscaling and load balancing are crucial for cloud applications to maintain Service Level Objectives (SLOs) such as low latency and high availability. Traditional reactive scaling (like HPA) has inherent lag: it adjusts only after metrics (CPU, memory) cross thresholds. During sudden spikes, this lag can lead to under-provisioning, causing request timeouts or slow responses, and subsequent over-provisioning, which waste resources once the spike passes. In e-commerce, banking, or government services, such inefficiencies directly impact user experience and operational costs. For example, a banking microservice might see traffic quadruple during salary payment days that could overwhelm current pods if scaling is not fast enough.

Predictive scaling is significant because it offers a way to stay ahead of demand and mitigate latency spikes by scaling out *before* a surge hits. This has implications not just for normal operations but also for handling seasonal events (sales, promotions) or unique local patterns (like month-end usage bursts). Additionally, proactive scaling can reduce costs by scaling in sooner after load subsides, avoiding lengthy periods of over-provisioning compared to static over-provisioning approaches.

Load balancing algorithms, on the other hand, ensure that when more pods are added (either proactively or reactively), incoming traffic is efficiently distributed to utilize all available capacity. Without proper load balancing, simply adding pods may not help if traffic is not shared (e.g., if a poorly configured client keeps using an old pod's IP). Thus, enhancing Kubernetes' default load balancing (which is typically round-robin via kube-proxy's iptables mode) or using smarter algorithms can further improve performance under heavy load.

In summary, the significance lies in achieving a harmonious blend of prediction and distribution: predicting future load to allocate resources in advance and distributing load evenly across those resources. This combination promises better quality of service for end-users and more cost-effective infrastructure usage for providers, a win-win for cloud operations.

1.4 Review of Significant Research

Prior research has explored both reactive and proactive autoscaling in cloud environments. Reactive autoscaling (including Kubernetes HPA) is well-studied and widely used; it reacts to current metrics like CPU utilization to scale pods. *Significant limitations* of reactive approaches (which we detail in Chapter 2) include their inability to handle sudden load spikes gracefully and the need to carefully tune cooldown periods to avoid thrashing.

Proactive or predictive autoscaling has been an active area of research in recent years. For example, Rampérez et al. (2021) proposed FLAS, a hybrid proactive-reactive autoscaler combining forecasting with reactive tuning [6]. Their approach showed improved stability and resource usage. Another notable work by Marie-Magdelaine and Ahmed (2020) applied machine learning to proactive autoscaling for cloud-native apps [5], demonstrating that ML models (like regression or simple neural networks) could predict workload trends and guide scaling.

Machine learning techniques specifically like LSTM (Long Short-Term Memory networks) are popular for time-series prediction due to their ability to learn temporal patterns. Nguyen et al. (2022) introduced a graph-based proactive HPA using LSTM and Graph Neural Networks that emphasizes capturing microservice dependency graphs for better predictions. Their Graph-PHPA approach outperformed Kubernetes' rule-based HPA in resource savings. Similarly, Dang-Quang and Yoo (2022) developed a Bi-LSTM based autoscaling framework for cloud applications [11] and achieved efficient prediction for multivariate workloads. We will review these approaches in detail in Chapter 2 to identify strengths to incorporate and gaps to fill.

On the load balancing side, significant research has looked at algorithms to distribute requests among microservice instances. Kubernetes' default approach uses iptables or IPVS (IP Virtual Server) for round-robin load balancing across pod Ips [2]. More advanced methods, like least-connections (routing new requests to the pod with the fewest active connections) or hash-based routing (to maintain session affinity), have been studied and implemented in various forms. Service mesh technologies (e.g., Istio or Linkerd) also offer sophisticated load balancing and failure handling at the mesh layer. Chapter 2's literature review will cover common load balancing algorithms and their suitability for microservices.

1.5 Assumptions and Limitations

Several assumptions underlie this study:

We assume the workload patterns are partially predictable. The LSTM model requires historical data to learn patterns (e.g., daily cycles, month-end peaks). We assume sufficient historical metrics (requests per second, CPU usage, etc.) are available for training. Sudden black-swan events (completely unpredictable surges) remain outside the scope of what a model can predict.

The microservices under study are containerized and orchestrated by Kubernetes. We leverage standard Kubernetes components (Metrics Server, HPA interface, etc.) and assume they are functioning correctly. We also assume that scaling actions (adding pods) can be executed by the cluster within a reasonable time (tens of seconds), which is typical if the container images are pre-pulled.

Our predictive scaling controller will run inside the cluster to read metrics and make scaling decisions. We assume network connectivity and API access to the cluster for issuing scale commands.

As for limitations:

The LSTM model is trained on a specific application's load patterns. Its effectiveness might not generalize without retraining for other applications with different traffic characteristics. We mitigate this by focusing on a representative workload.

Running a predictive model (especially a deep learning model) continuously might incur CPU/memory overhead. We measure this overhead; however, in this thesis, we will not deeply optimize the model's efficiency. The focus is on feasibility and benefit analysis.

No prediction is 100% accurate. We incorporate a fail-safe mechanism: if the prediction underestimates load, the normal HPA or a secondary reactive check will catch up. However, this interplay might introduce complexity. This study will not fully resolve how to tune the balance

between predictive and reactive scaling (left for future work), but we will document observations.

We implement a custom load balancer primarily at the Kubernetes Service level (possibly by adjusting kube-proxy mode to IPVS or using an ingress with specific algorithms). We will not implement a full custom proxy but rather configure existing tools to use desired algorithms. The limitation is we rely on Kubernetes or ingress controllers' capabilities; designing a completely new load balancing algorithm is beyond scope.

With the introduction setting the stage, the next chapter will review the literature on microservices architecture, autoscaling techniques (HPA, VPA, KEDA), limitations of current reactive scaling, various load balancing algorithms, and prior work in predictive scaling with machine learning. This foundation will inform the design of our solution.

2 REVIEW OF THE LITERATURE

2.1 Kubernetes-Based Microservices Architecture

Microservices architecture is an approach to software design where applications are composed of small, independent services, each encapsulating a specific business capability [1]. These services communicate over well-defined APIs, often via lightweight protocols like HTTP/REST or gRPC. This contrasts with monolithic architectures by offering greater modularity, independent deployability, and fault isolation – a failure in one microservice (e.g., the payment service) ideally does not bring down the entire application.

Kubernetes is a container orchestration platform that has been widely adopted to deploy and manage microservices in containers (e.g., Docker) [2]. It has grown in popularity due to its flexibility with containerized applications, a synergy between Docker and Kubernetes noted by Hardikar et al. (2021) [15]. Kubernetes provides key primitives that align well with microservices:

- **Pods and Deployments:** A microservice runs in one or more pods (typically one container per pod). A Kubernetes *Deployment* manages replica pods, which enables easy scaling and rolling updates.
- **Service Discovery and Load Balancing:** Kubernetes Services provide a stable network endpoint (ClusterIP or DNS name) for a set of pods. Traffic to a Service is automatically load balanced across the pods behind it. By default, kube-proxy in iptables mode uses a form of round-robin NAT rules to distribute connections evenly. With IPVS mode, more sophisticated algorithms (round-robin, least-connections, etc.) are available.
- **Declarative Configuration:** The state of the system (how many pods, what container image, which node to run on, etc.) is declared in YAML manifests. The cluster's controllers continuously work to make the observed state match the desired state, which provides self-healing (if a pod dies, a new one is started) and scaling (if replicas count is increased, new pods are created).
- **Multi-Cloud and Hybrid Support:** Kubernetes abstracts away the underlying infrastructure to make it possible to run microservices across on-premise and cloud in a similar manner. This is important for portability and avoiding vendor lock-in.

In microservices, network communication becomes a critical aspect. Kubernetes assumes a flat network model where each pod can reach any other pod by IP that simplifies how microservices

discover and talk to each other. Typically, additional components like a service mesh (Istio, Linkerd) might be used to manage cross-cutting concerns (e.g., observability, retries, circuit breaking) for inter-service calls.

From an architectural perspective, microservices on Kubernetes allow each service to be scaled horizontally independently, which is a key advantage. For example, a “statement” service that is heavy on CPU might be scaled to 5 pods, while an “auth” service might only need 2 pods. This granularity helps optimize resources and performance.

However, microservices also introduce complexity in orchestration. Dozens of services mean dozens of deployments and services to manage. Autoscaling becomes indispensable to adjust those deployments dynamically, and load balancing ensures that scaling out truly brings benefits by routing traffic to all replicas.

2.2 Auto-Scaling in Kubernetes (HPA, VPA, KEDA)

Kubernetes supports two primary forms of autoscaling out-of-the-box: Horizontal Pod Autoscaler (HPA) and Vertical Pod Autoscaler (VPA). Additionally, the Kubernetes Event-Driven Autoscaler (KEDA) is an extension that enables scaling based on external events or custom metrics.

- **Horizontal Pod Autoscaler (HPA):** HPA automatically changes the number of pod replicas for a deployment (or replication controller) based on observed metrics like CPU usage or custom metrics. For example, an HPA might target 60% CPU utilization per pod; if actual usage exceeds this, HPA will add pods. It operates on a feedback loop: every evaluation cycle (default 15s in Kubernetes), HPA checks metrics from the Metrics Server (or Prometheus for custom metrics) and computes the desired replica count. HPA is a *reactive* scaler: it only responds to current or recent past metrics. HPA also incorporates some built-in damping – such as scaling *cooldowns* or tolerance bands – to avoid flapping (rapid scale up and down). However, even with these, a sharp spike can still overshoot the system because HPA will only detect it after it is underway.
- **Vertical Pod Autoscaler (VPA):** While HPA scales the number of pods, VPA adjusts the resources (CPU, memory) requests and limits of containers in pods. If a microservice is consistently using near 100% of its requested CPU, VPA can recommend or automatically set a higher request to ensure the scheduler gives it more resources on nodes. VPA can prevent pods from being starved of resources and can simplify not having to guess the “right” CPU/memory for each service. However, VPA has limitations: changing resources often requires restarting pods (to apply new resource limits), and it is typically not run simultaneously with HPA on the same metric (to avoid conflicting decisions). VPA is most useful for long-running load patterns, not sudden spikes (because restarting pods during a spike might actually worsen things temporarily). In this thesis, we focus more on HPA and predictive horizontal scaling, while VPA is considered complementary.
- **Kubernetes Event-Driven Autoscaler (KEDA):** KEDA extends HPA by allowing scaling on external metrics or events [2]. For instance, KEDA can scale based on the length of a message queue (RabbitMQ, Azure Service Bus, etc.), the number of Kafka messages, or custom Prometheus metrics. KEDA essentially acts as an adapter that feeds these external metrics into a scaled object, driving an HPA under the hood. A key use

case is scaling to zero and then from zero to N on events (something vanilla HPA cannot do because when pods=0, no metrics exist). In context of predictive scaling, KEDA could be used if our predictive model output can be exposed as a custom metric or event; however, our approach more directly controls scaling. Still, KEDA's architecture of separating metric ingestion from scaling logic is a useful reference.

Traditional HPA uses relatively simple threshold-based logic. If current usage is higher than target, scale out (and vice versa). Research has proposed advanced algorithms, such as PID controllers (to treat scaling like a control system problem) or rule-based systems that incorporate multiple metrics. Bartelucci and Bellavista (2023) provide a structured taxonomy of autoscaling solutions, distinguishing between reactive, scheduled, and predictive approaches, and highlight the practical challenges of applying machine learning-based scaling in Kubernetes environments [14]. But increasingly, machine learning and forecasting are being explored to improve upon HPA's reactivity. Reinforcement learning (RL) has been applied to autoscaling and treated it as a sequential decision problem where the system learns a policy to scale up/down based on state, to minimize cost while meeting performance (e.g., Fabiana Rossi et al., 2019 used RL for combined horizontal and vertical scaling).

In practice, many Kubernetes deployments still rely on HPA (and sometimes VPA) because they are simple and built in. KEDA is used in event-driven scenarios, but it is also reactive to events rather than predictive. The gap identified is that none of these default tools predict future load – they require an external component or manual intervention to do proactive scaling. This sets the stage for why predictive scaling is needed and how it complements or replaces these reactive tools.

2.3 Limitations of Reactive Scaling Techniques

Reactive scaling (like HPA) adjusts resources based on current or very recent past usage. Key limitations of this approach, highlighted in research and operational experience, include:

- **Latency in Response:** There is an inherent delay between a load increase and the scaling response. For example, if traffic doubles in one minute, CPU usage will spike quickly. HPA might notice this spike only after its polling interval (say 15s), then maybe decide to add pods. Creating new pods (scheduling, pulling images, starting containers) might take at least another 10-30 seconds depending on the environment. During this time, the existing pods are overloaded, which leads to high latency or errors. By the time new pods are ready, the damage (SLO violations) may have occurred. In essence, reactive scaling often catches up after the fact.
- **Oscillation and Thrashing:** If not tuned carefully, reactive systems can oscillate – e.g., scale up too many pods on a transient spike, then scale down, then up again. HPA implements cooldown periods and uses averaged metrics to mitigate this, but it can still happen with rapidly changing workloads. Amazon's EC2 auto-scaler has the notion of scaling cooldowns to wait after a scale action to let the system stabilize [16]. Kubernetes HPA similarly has a stabilization window for scale down to avoid yanking resources too fast.
- **Threshold Tuning:** Reactive scaling usually relies on thresholds (like "if CPU > 70%"). Setting these thresholds is tricky: too low and you will scale up frequently (maybe wasting resources), too high and you risk hitting resource saturation before scaling

triggers. Workloads also might have different ideal thresholds. It is difficult to find one-size-fits-all rules, often requiring manual tuning per microservice.

- **Single Metric Focus:** Kubernetes HPA by default often uses CPU (or memory). While it can use custom metrics, in practice many setups just use CPU. But high CPU might not directly correlate with user-facing performance; for instance, a spike in requests might first manifest in higher request latency before CPU climbs. Or an I/O heavy workload might have low CPU but be bottlenecked elsewhere. Reactive scaling on a single metric can miss the full picture. This is partly why approaches like request rate-based scaling (e.g., via KEDA or custom metrics) are introduced – they align scaling with actual demand more directly than resource usage does.
- **Not Considering Trend/Seasonality:** Reactive systems treat each moment in time more or less independently. They do not “know” about daily patterns, weekly trends, or upcoming events. If load is gradually rising every day at 9am, reactive scaling always repeats the same dance: let it rise, then add pods. It does not anticipate the rise at 8:55am. This is inefficient and can cause brief slowdowns every day. In Azerbaijani context, think of the month-end surge: on the last day of the month around 6pm, usage spikes. Reactive scaling will only respond at 6pm when metrics jump, whereas a predictive approach could start adding capacity at 5:50pm, mitigating the surge entirely.
- **Limited by Historical Peaks:** Some reactive systems base decisions also on past high-water marks (to avoid overreacting). This can be problematic if a new peak goes above anything seen historically – reactive scaling might under-provision because it does not realize a new record high is possible. It has no foresight beyond what’s in the immediate metrics.

In literature, these issues have been well documented. For instance, Mao and Humphrey (2011) pointed out challenges in meeting application deadlines with reactive scaling [4], which suggests more sophisticated approaches to minimize cost while meeting performance. Uргаonkar et al. (2005) discussed dynamic provisioning using queueing theory for multi-tier apps to handle variable loads [9], which was an early step towards being smarter than simple thresholds.

Overall, the limitations of reactive scaling underscore why predictive scaling is pursued: to handle the cases where reactivity is too slow. That said, predictive approaches come with their own challenges like prediction errors, which is why some hybrid solutions include reactive components as a safety net.

2.4 Load Balancing Algorithms for Microservices

Load balancing is the process of distributing incoming requests across multiple service instances to ensure no single instance is overwhelmed and to utilize available resources efficiently. In microservices, load balancing occurs at multiple levels:

- **Client-side vs Server-side:** In some architectures (especially with service mesh or smart clients), the client may decide which instance to call (client-side load balancing). In Kubernetes, typically, we rely on server-side load balancing via the Service abstraction and kube-proxy or ingress controllers.
- **Layer 4 vs Layer 7:** Kubernetes Services (ClusterIP, NodePort, LoadBalancer types) operate at the TCP/UDP level (Layer 4) that generally balance connections without looking at HTTP details. Ingress controllers or service meshes can do Layer 7 balancing,

e.g., routing HTTP requests by URL path, etc., but at the basic level, we consider simpler distribution methods.

Common load balancing algorithms:

1. **Round Robin:** The default for many systems – each new request goes to the next server in a list, rotating through them. Kubernetes' iptables proxy essentially does round-robin by iterating through endpoints in NAT rules. This is simple and often works fine when requests are roughly equal in processing cost.
2. **Least Connections:** The balancer sends new requests to the instance with the fewest active connections (or least current load). This can improve distribution if some requests are long-lived or heavy, as those will tie up connections and thus those pods will naturally get less new traffic. IPVS mode in kube-proxy supports least connection as a scheduling algorithm.
3. **Weighted Round Robin/Least Connections:** If some instances are more powerful (not typical in Kubernetes where pods are uniform, but possibly if mix of old/new hardware), weights can be assigned. In microservices scaling, sometimes newer pods might be cold (caches empty) and perform slightly worse initially, so an intelligent balancer might ramp traffic to them slowly (effectively weighting by readiness).
4. **Hash-Based (Consistent Hashing):** Methods like destination hashing (dh) or source hashing (sh) route based on some hash of the request (client IP, session ID, etc.). This is used when session affinity is needed – e.g., always send the same user to the same backend for cache locality. Kubernetes services have an option called SessionAffinity: ClientIP which effectively pins a client IP to a specific pod via a simple hash mechanism. However, consistent hashing can also distribute loads unevenly if not carefully implemented.
5. **Random:** Surprisingly, random assignment can be nearly as good as round robin in large numbers, but with fewer guarantees, so it is rarely used explicitly in systems like Kubernetes.
6. **Custom/Adaptive:** In advanced scenarios, load balancers can use metrics (like response times, CPU load on pods) to make decisions (this ventures into realm of global schedulers or service meshes). For example, a smart load balancer could detect one pod is responding slower and temporarily send less traffic to it.

Kubernetes specifics: As noted, kube-proxy in iptables mode installs one rule per endpoint, and packets are DNAT'ed (Destination NAT) to a pod IP in round-robin fashion. In IPVS mode, which is more scalable for very large numbers of services, the ipvs virtual server can be configured with scheduling algorithms (rr, lc, etc.). By default, Kubernetes sets IPVS to round-robin, but one can tweak it if needed.

Ingress Controllers (like NGINX ingress, Envoy via Istio, etc.) often default to round robin as well for balancing across backend pods but may support least_conn or other algorithms via configuration.

Service Mesh (Istio/Envoy) does layer7 load balancing and can gather per-instance metrics to do more fine-grained balancing, but it also adds some overhead. In the context of this thesis, we likely assume standard Kubernetes service-level balancing, but if needed, we might consider a simple Envoy or Nginx that can apply least_conn to demonstrate its effect.

For microservices, load balancing ensures that when we scale out from 2 pods to 10 pods, clients do not keep hitting just the original two. Fortunately, in Kubernetes, as soon as new pods are Ready, they are added to the Service endpoints list, and traffic can flow to them. There is a subtle interplay if using DNS for service discovery (like headless services with client-side LB): caching could cause delays in clients seeing new endpoints. But with typical ClusterIP Services, the cluster handles it transparently.

Research on load balancing in microservices has also looked at microservice placement – e.g., if some microservices are stateful or have caches, it might be beneficial to route certain requests to specific pods (for example, a request related to “user A” always to the same pod handling user A’s cached data). That gets into request routing strategies beyond pure load leveling.

In summary, load balancing algorithms like round robin work well for stateless services with homogeneous pods, while least-connections can be better if request load is uneven. This thesis will use load balancing as a supporting mechanism – ensuring our scaled pods are utilized. We might not innovate a new algorithm but choose the best suited one (round robin vs least_conn) based on experiments. The literature suggests that least_conn yields more balanced load especially as concurrency increases – under high connection counts, least_conn leads to a more even distribution than round robin.

2.5 Predictive Scaling Approaches Using Machine Learning

Predictive scaling uses algorithms or machine learning models to forecast future resource demand, then scales resources in advance. Instead of waiting for metrics to exceed thresholds, predictive systems allocate capacity ahead of time. This section reviews various ML approaches that have been applied to this problem:

- **Time-Series Forecasting Models:** These include classical methods like ARIMA/SARIMA, Facebook Prophet, as well as neural network based models like RNNs and LSTMs. For example, Facebook Prophet is a popular forecasting tool that handles seasonality and trends; some recent work integrated Prophet with LSTM for Kubernetes scaling [17]. Prophet would capture weekly patterns (e.g., Monday vs Sunday traffic differences) while LSTM could model the more complex residual patterns. The hybrid approach by Nguyen et al. (2025) in *Frontiers* used Prophet+LSTM to improve predictive accuracy for autoscaling. Their results showed reduced resource wastage and improved SLO adherence compared to default HPA.
- **Long Short-Term Memory (LSTM) Networks:** LSTMs are a type of recurrent neural network adept at learning from sequences of data, and they handle long-term dependencies better than traditional RNNs. In autoscaling, LSTMs have been used to predict metrics like CPU load, request rates, etc., minutes or seconds into the future. For instance, Bi-LSTM (Bidirectional LSTM) models can leverage past and future context in sequences; Dang-Quang & Yoo (2022) applied a Bi-LSTM for cloud workload prediction and saw high accuracy. Another approach by Zhu et al. (2019) used an attention-based LSTM encoder-decoder to predict workloads [13], improving on plain LSTM by focusing on relevant portions of the sequence.
- **Reinforcement Learning (RL):** Instead of directly predicting the metrics, some works treat scaling as a decision problem. A reinforcement learning agent observes the system

state (current load, pods, etc.) and takes actions (scale up/down) to maximize a reward (throughput minus cost, for example). Rossi et al. (2019) developed an RL-based scaler that learned when to horizontally or vertically scale containers [7]. Toka et al. (2020) proposed an adaptive AI-based autoscaler that combined learning techniques [8]. RL can, in theory, find an optimal policy that might even outperform a purely prediction-based approach, especially in environments with unpredictable changes or complex multi-tier interactions. However, RL typically requires extensive training and careful reward design; it is also harder to guarantee it will not make a bad decision during learning.

- **Supervised Learning & Regression:** Simpler than LSTM, some research used regression models or simpler ML to predict scaling needs. For example, SVM (Support Vector Machine) and Multi-Layer Perceptrons were evaluated by Ho et al. (2023) for predicting CPU requirements of microservices. They compared linear vs polynomial kernels, finding linear models often had lower prediction error for resource usage. While not as powerful as LSTM for capturing complex patterns, these models are lightweight and easier to interpret.
- **Hybrid Approaches:** Some systems combine reactive and proactive elements. FLAS (Rampérez et al., 2021) as mentioned, does forecast-based scaling but also monitors and reacts if predictions were off. Another hybrid example: “Prorenata” by Liu et al. (2015), which the name suggests combines proactive and reactive tuning for a distributed storage system. These hybrids try to get the best of both worlds: prediction to get ahead, reaction to correct errors.
- **Workload Characterization & Custom Models:** Some researchers create customized models for specific scenarios. For instance, Graph-PHPA (Nguyen et al. 2022) included microservice dependency graph info by using Graph Neural Networks alongside LSTM. The idea is that the load on one microservice might correlate with load on another (e.g., increased traffic on “frontend” service will cause increased load on “backend” and “database” services). By understanding these relationships, predictions can be more holistic. Similarly, SmartVM (Zheng et al. 2018) was a framework focusing on microservice placement with SLA-awareness [10], not exactly scaling, but touches on proactive management.

Key findings from literature: Predictive autoscaling generally can outperform reactive in maintaining performance, but it requires accurate forecasting and can be complex. Many studies report improved metrics (lower latency, fewer SLO violations, better resource usage) with predictive methods. However, they also note limitations such as the need for retraining models when workload patterns change, the risk of over-provisioning if predictions overshoot, and implementation complexity.

For example, Khan et al. (2020) might have looked at combining forecasting with threshold-based triggers to avoid false alarms. Podolskiy et al. (2018) provided a comparative study on forecasting models (ARIMA vs LSTM vs others) for cloud auto-scaling and found that no one model is best in all cases; hybrid or ensemble methods can be beneficial.

Machine Learning for scaling in practice: Cloud providers have started implementing predictive autoscaling features. AWS’s Auto Scaling can do target tracking (like HPA) or scheduled scaling (crude prediction by schedule). Google Cloud AutoML for scaling, etc., are

emerging but not yet mainstream for all users. This thesis will contribute by demonstrating a concrete LSTM-based scaler in a Kubernetes environment and documenting its pros/cons.

2.6 Summary of Research Gaps

From the literature review above, several gaps and open issues emerge:

- While many studies have proposed predictive scaling, few have been integrated tightly with Kubernetes in a way that a regular DevOps team could easily adopt. There's a gap in bridging the advanced ML techniques with practical Kubernetes deployments. This thesis aims to reduce that gap by using common Kubernetes to inject predictions, rather than building a completely separate system.
- Most predictive approaches focus on scaling but assume default load balancing will handle distribution. The interplay between scaling and load balancing is not deeply explored. For instance, if predictive scaling adds pods preemptively, how to ensure warm-up (some apps may need warming caches) and how to gradually shift load? We identify a gap in ensuring newly added capacity is smoothly utilized without overloading it instantly (or conversely, sitting idle because clients have not switched traffic). Our research includes a custom load balancer design to complement predictive scaling.
- In terms of metrics, a gap is noted in Kubernetes' default autoscaler: inability to scale on custom metrics like request rate directly. KEDA addresses it but again reactively. We want to predict on a highly relevant metric (like request rate) rather than proxy metrics like CPU, bridging a gap in relevance of scaling signals.
- Many studies evaluate their models on either simulated or narrow workloads. There's a gap in understanding how predictive scaling works for realistic microservice applications. This thesis, while primarily focusing on one microservice's scaling, acknowledges real apps have many interconnected services. We design our evaluation to at least use a realistic service (possibly with a dependency or two) and realistic load patterns (like actual traces or seasonality) to ensure relevance.
- Edge cases such as mispredictions: The literature does not always discuss failure handling in detail. The Sensors (2023) article introduced an interim reactive check to handle prediction errors. This is an important gap – no prediction will be right all the time, so how to fail gracefully? We incorporate an assumption that HPA's regular mechanism is left running as a fallback so that if our predictor is asleep at the wheel, HPA can still save the day, albeit with some delay. Formalizing this interaction is complex and remains an open challenge.
- Also, a contextual gap: little to no literature mentions specific industry based phenomena – which are essentially like “special events”. Handling special but known events (maybe via a calendar of expected spikes, or recognizing patterns that are not frequent enough to learn easily) is somewhat open. This might lean into the realm of combining ML with knowledge-based or rule-based scaling (e.g., a rule: “on last weekday of month at noon, double the pods”). We will not deeply solve this, but we will incorporate such domain knowledge in our test scenario to see if our predictive model can learn that monthly pattern or if a manual intervention is needed.

In conclusion, the literature provides a strong foundation that predictive scaling can be beneficial, but gaps in practical implementation, integration with load balancing, and handling

real-world complexity remain. Our research is structured to address some of these gaps by building a working prototype of predictive scaling in Kubernetes and evaluating it in detail.

3 RESEARCH METHODOLOGY

This chapter outlines the methodology for designing, implementing, and evaluating our predictive scaling and load balancing solution for Kubernetes-based microservices. We describe the system architecture, data collection, the LSTM prediction model, the custom load balancer, and the experimental setup to provide rationale for each choice.

3.1 Overall System Architecture

The overall system architecture (illustrated conceptually in **Figure 3.1**) consists of the following components:

- **Kubernetes Cluster:** Running the microservice application we want to scale. Within the cluster:
 - The target *microservice* is deployed with multiple features.
 - Prometheus is set up to collect resource metrics (CPU, memory, etc.) from pods.
 - Standard *Horizontal Pod Autoscaler (HPA)* is configured (for baseline comparisons and backup).
 - A *Service* object fronting the microservice pods for load balancing.
- **External Predictive Autoscaler Service:** This is our custom component running outside the cluster (it could also run as a pod inside, but with elevated privileges). It includes:
 - **Data Collector:** which periodically pulls metrics from the cluster. Specifically, it collects time-series data such as request rates, CPU usage, memory usage, etc.
 - **Predictive Model (LSTM):** which takes recent historical data and produces a forecast for future load (e.g., next 5 minutes prediction of requests per second).
 - **Scaling Decision Module:** which uses the predicted load to decide if scaling actions are needed now. For example, if it predicts a 2x increase in 5 minutes, it may trigger adding pods gradually ahead of time.
 - **Kubernetes Interface:** that actually performs scale actions by calling the Kubernetes API to adjust the Deployment's replica count.
- **Custom Load Balancer:** While Kubernetes will handle basic load distribution, we plan for a component or configuration that ensures an optimal algorithm is in use. This is simply an application as a gateway with custom load balancing policy. This component ensures that as the predictive autoscaler adds pods, incoming traffic is smoothly balanced across them.
- **Monitoring & Logging:** We deploy monitoring tools (Prometheus, Grafana) and logging (ELK stack as well as Kubernetes logs) to record the system's behavior. These are crucial for evaluation – capturing metrics like response time, CPU, number of replicas over time, etc., so we can analyze how our system performed.

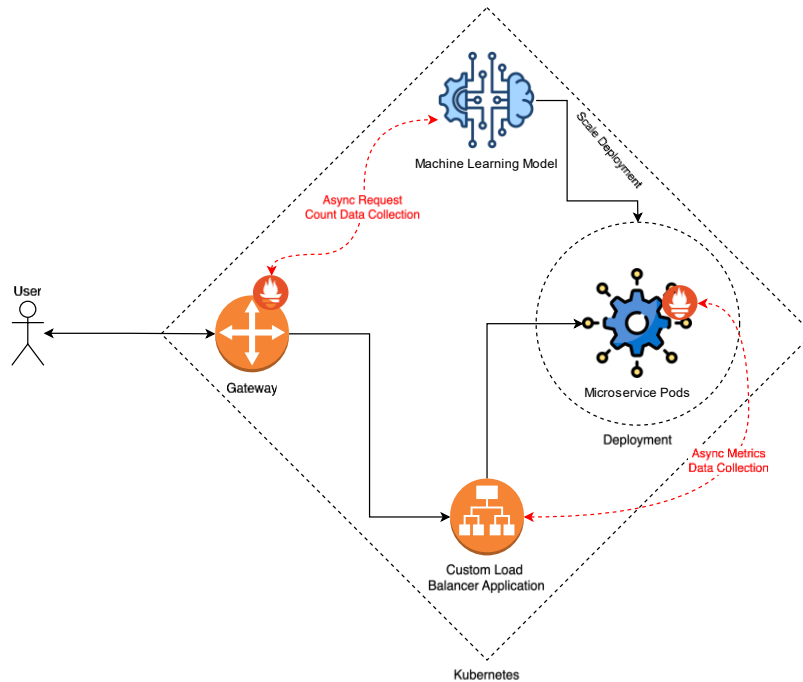


Figure 3.1: System Architecture of Predictive Scaling Solution

This architecture is modular: one could replace the LSTM model with another predictor or integrate the autoscaler service directly as a Kubernetes Custom Controller inside the cluster. We chose an external service for easier development and flexibility in using libraries without containerizing them at first.

The architecture leverages Kubernetes' existing mechanisms where possible (HPA, Metrics API) to avoid reinventing wheels. For instance, we considered using a Custom Metric fed into HPA: the LSTM could push a "predicted load" metric that HPA then uses to scale. However, for transparency and control, scaling was implemented directly via the Kubernetes API. By having a clear separation (predictive logic outside, execution inside), we maintain decoupling.

3.2 Data Collection and Preprocessing

We gathered a comprehensive 12-month dataset from a real production Kubernetes microservices environment. The data includes key time-series metrics such as HTTP request rates, CPU usage, and memory usage for each service. The raw data collected for each microservice consisted of time-stamped resource and load metrics. A sample data record is structured as follows:

```
{
  "timestamp": "2023-11-30T17:00:00Z",
  "cpu_usage": 62.5,
  "memory_usage": 512,
  "request_rate": 45
}
```

Each field represents:

- "timestamp": UTC timestamp when the metrics were captured.
- "cpu_usage": Percentage of CPU used at that moment (%).
- "memory_usage": Memory usage in megabytes (MB).
- "request_rate": Number of HTTP requests received per second.

These metrics were continuously collected using Prometheus, an open-source monitoring tool that scrapes metric endpoints of the services at regular intervals. The use of Prometheus ensured high-resolution, consistent data collection of resource utilization and traffic load over time. Each data point is timestamped which creates a chronological record of system behavior under various conditions.

This year-long dataset captures both regular usage patterns and extraordinary events. For example, we observe daily and weekly cycles in the request rate – weekdays generally have higher traffic than weekends – as well as monthly patterns where the end-of-month days exhibit significantly higher load. In the context of financial sector of Azerbaijan, many users receive salary payments at month's end that drives a surge in usage on those days. Additionally, the dataset reflects special peak events such as public holidays and other seasonal spikes that provide a realistic variety of load conditions. These periodic trends and irregular surges are valuable for training a predictive model that can anticipate scaling needs. To prepare for model training and evaluation, we split the sequence chronologically: approximately 10 months of data for training and the most recent 2 months for testing. This temporal split ensures that the model is validated on “future” data it has never seen to closely mimic real-world deployment where we predict upcoming load from past patterns. It also preserves the temporal dependencies to prevent information leakage from the future into training. The training period encompasses multiple instances of the weekly and monthly cycles that allows the model to learn recurring patterns, while the test period includes at least one end-of-month peak and weekend vs. weekday variations to assess the model's generalization.

Before feeding this data into our Long Short-Term Memory (LSTM) predictive model, extensive preprocessing was performed to clean and transform the raw metrics into a learnable format.

Over a year of collection, there were some intervals with missing or incomplete data. It is crucial to handle these gaps to avoid corrupting the model training. Short gaps were imputed using forward-fill to carry forward the last known legitimate value to fill a missing timestamp. For longer or conspicuous missing periods, we employed more careful imputation and excluded long missing periods from training to avoid bias. This ensured a continuous and coherent time series for each metric, as LSTMs cannot handle NA values directly.

All metric values were scaled to a uniform range using min-max normalization. In particular, each feature (e.g., request rate, CPU, memory) was rescaled to [0, 1] based on the min and max values observed in the training set. This normalization is a standard practice in training neural networks on time-series data, as it prevents features with larger numeric ranges from dominating others and helps the model converge. By normalizing the inputs, the LSTM can learn effectively from the shape of the time-series data rather than the scale. The same scaling parameters (min and max) computed from training data were later applied to scale the test data, to maintain consistency.

We transformed the sequential data into a supervised learning format using a sliding window technique. Rather than training the LSTM on the entire sequence at once, we created input-output

samples by sliding a fixed-size window over the time series. For example, using a window of N past time steps (e.g., past one hour of data divided into minutes, or past several days of hourly data), we use those N points and associated features as inputs to predict the metrics at the next time step as output. Each window moves forward by one time step to generate the next training sample. This approach yields a large number of training examples from the 10-month series, while preserving the temporal order within each sample. The sliding window method is essential for sequence forecasting problems, as it converts the time series into input-output pairs the LSTM can learn from. We tuned the window size based on the problem’s needs – it should be long enough to encompass important patterns (for instance, a few days might capture weekday/weekend effects), but not so long that the model becomes overly complex or slow. In our case, we empirically chose a window that captures at least a full week of data, so the model can potentially learn weekly periodicity.

In addition to the raw metrics, we added several **calendar features** to help the model recognize temporal patterns explicitly:

- *Time-of-Week*: We introduced features to indicate the day of the week for each data point. This was done by an integer where 0=Monday and 6=Sunday. We found it useful to at least include a binary flag for **weekend vs. weekday**, as the load on Saturdays and Sundays was observed to differ markedly from weekdays. This feature allows the LSTM to adjust its expectations based on where the current time falls in the week cycle. Capturing such weekly seasonality is known to improve forecasting of workloads that have weekday/weekend usage oscillations.
- *Time-of-Day*: Similarly, we added a feature for the hour of the day to recognize daily cycles (e.g., perhaps lower usage at late night hours and higher during business hours). Rather than a simple hour integer that jumps from 23 to 0 at midnight, we encoded time-of-day in a cyclic manner (e.g., using sine and cosine transforms of the hour) so that 23:00 and 00:00 are recognized as adjacent in time. This allows the model to learn smooth daily patterns. Including daily seasonality explicitly is helpful since user activity often peaks at certain hours each day.
- *Month-End Proximity*: Because our domain has a pronounced spike at the end of each month, we crafted a feature to represent how close a date is to the end-of-month. For instance, one approach is to compute “days until end of month” or a boolean feature that is 1 if the date is within the last 3 days of a month (and 0 otherwise). We implemented a simple normalized countdown: at the start of each month this feature is low, and it gradually increases and peaks on the final day of the month. This gives the LSTM a hint for the looming end-of-month surge. Essentially, it’s a way of injecting domain knowledge of monthly financial cycles into the model. While Prophet-like models often allow explicit holiday or seasonal event definitions, here we emulate that by providing the model with a custom feature indicating those critical periods.
- *Periodic Holidays*: Although not a primary focus, we also flagged a few known public holidays in the dataset. These were treated similarly to weekend indicators (as binary flags). Public holidays can influence usage (either spikes for some services or drops for others) and marking them helps the model not treat those anomalies as pure noise. (This feature is sparse since holidays are infrequent, so the LSTM might or might not utilize it strongly, but it was included for completeness.)

After these preprocessing steps, the multi-variate time series data is ready for modeling. At this stage, each training sample consists of a normalized window of past metrics (request rate, CPU, memory, etc.) along with the engineered time features, and the output is the next time step's load (we particularly focus on predicting the *future request rate or CPU load* as the target for scaling). The choice of an LSTM network for this task is motivated by its strength in modeling sequential data and capturing long-term dependencies. By training our LSTM on a rich year-long dataset with diverse patterns, we aim to achieve similarly predictive scaling – the model will forecast impending load increases (e.g., the Monday morning rush or the month-end peak) so that the cluster can scale out *before* the load actually hits, thereby maintaining performance. The cleaned and feature-enriched dataset is the foundation for this learning: it ensures the model has all relevant signals (seasonal indicators, recent trends, etc.) to make an accurate forecast.

3.3 LSTM-Based Load Prediction Model

Kubernetes' default Horizontal Pod Autoscaler (HPA) reacts to load changes based on current metrics, which can result in scaling decisions coming after a traffic surge has already impacted performance. To address this reactive delay, we implemented a predictive scaling model using a Long Short-Term Memory (LSTM) neural network. LSTM networks are a type of recurrent neural network known for learning long-term temporal dependencies, making them well-suited for forecasting complex usage patterns. Our LSTM-based predictor enables the system to anticipate workload spikes and initiate scaling *before* the cluster experiences high load, rather than waiting for metrics to breach thresholds.

The predictive model is a univariate LSTM network that takes the past p load observations to predict the load at the next time interval. The LSTM model was implemented using TensorFlow 2.x in Python. We selected p to be 60 (representing the last 60 minutes of observations) to capture short-term trends along with daily patterns. The LSTM model contains two hidden layers with 50 neurons each, followed by a dense output layer that produces the predicted load value. We trained the model for 50 epochs using the Adam optimizer and mean squared error (MSE) as the loss function, stopping early if the validation error plateaued.

During training on 10 months of data, the model learned the recurrent patterns of usage, including diurnal (daily) cycles and the notable monthly surge. On the test dataset (the 2 months of unseen data), the LSTM model demonstrated the ability to forecast upcoming spikes with reasonable accuracy – for instance, it correctly predicted the timing of the end-of-month load surges slightly ahead of time. This accuracy is critical, as even a few minutes of advance notice allows proactive scaling. LSTM-based prediction has been shown to outperform simpler time series methods in capturing such irregular but recurring spikes, which reinforced our choice of this model.

The trained LSTM predictor runs as a component in the cluster, continuously monitoring the incoming load metrics and generating short-term forecasts (e.g., 5 minutes into the future). Every forecast cycle, the system compares the predicted load against predefined scaling thresholds. If the forecast indicates that the future load will exceed the current capacity of the deployed pods (for example, if predicted CPU utilization or request rate per pod would breach safe limits), the autoscaler pre-emptively adds pods before the spike actually hits. Conversely, if the predictor foresees a sustained low usage period, the system can safely scale down some pods to conserve resources.

Instead of waiting for delayed reactive triggers, this mechanism acts on predictions to ensure sufficient resources are online ahead of demand. In this way, it effectively mitigates the risk of

service degradation during sudden demand surges. Our predictive autoscaling logic thus replaces or augments the standard HPA, in spirit similar to research that substitutes HPA decisions with LSTM-based predictions to avoid overload. In summary, this LSTM-based load prediction model forms the core of our proactive scaling approach, enabling the autoscaler to anticipate and respond to workload changes in a timely manner.

3.4 Custom Load Balancer Design

While predictive scaling ensures the *number* of pods is appropriate for upcoming load, we also addressed *how* traffic is distributed among those pods. Default Kubernetes services use a simple round-robin approach for load balancing, which does not account for each pod's current load or performance. In heavy load scenarios, this can lead to some pods becoming overburdened even if others have spare capacity. Kubernetes does not natively support balancing based on real-time pod metrics like CPU or response latency, so we designed a custom load balancer to fill this gap. The custom load balancer monitors each pod's runtime status — specifically its CPU utilization and recent response times — and uses this information to make smarter routing decisions. In essence, incoming requests are directed to the pod that is currently best able to handle additional load, rather than blindly rotating through pods. This decision process considers factors such as server capacity and response speed, aligning with general best practices for load balancing in distributed systems.

Our load balancer operates as a middleware layer between the client requests and the Kubernetes service. It continually receives metrics from all active pods. We use Prometheus to gather CPU usage and average response time for each pod at short intervals (e.g., every 5 seconds). With these metrics, the balancer maintains a dynamic ranking of pods. When a new request arrives, the balancer checks this ranking and forwards the request to the pod with the lowest current CPU load and fastest recent response time. In cases where a pod's CPU usage exceeds a critical threshold or its response time degrades beyond a set limit, the load balancer can temporarily stop sending traffic to that pod until it recovers. (One practical way to achieve this in Kubernetes is by toggling the pod's readiness status when it is overloaded, which causes Kubernetes to temporarily remove it from the service endpoints.) This approach prevents an overwhelmed pod from becoming a bottleneck and gradually reintroduces it once its performance stabilizes.

This intelligent routing strategy distributes incoming traffic more evenly in line with each pod's actual handling capacity. As a result, it avoids overloading any single pod while others are underutilized, thereby maintaining lower overall response times under high load. This prioritization of well-performing pods (and throttling of overloaded ones) inherently improves system throughput and resiliency. An overloaded pod is given a brief respite to recover, while other pods pick up the extra load. This is a behavior that reduces the risk of cascading failures or prolonged high latency on any one instance. The design follows the principle that load balancers should use runtime performance indicators to optimize request distribution. In our implementation, we observed that when the custom load balancer was active, the system maintained a more stable response time across all pods even during peak loads, compared to the default round-robin method. No single pod is allowed to overload the system, so the load balancing mechanism contributes significantly to improved performance and user experience. In summary, the custom load balancer works in tandem with the LSTM predictor: the predictor

ensures enough pods are provisioned ahead of time, and the load balancer ensures that the available capacity is utilized efficiently and fairly.

3.5 Evaluation Setup (Kubernetes Cluster, Tools, Metrics)

After implementing the LSTM-based autoscaler and custom load balancer, we evaluated their effectiveness using a controlled experimental setup. This section describes the Kubernetes cluster configuration, the tools and data used to generate load, and the metrics by which we assess performance. The goal of the evaluation is to compare our proactive autoscaling solution against the standard Kubernetes HPA under identical workload conditions.

All experiments were conducted on a Kubernetes cluster deployed on three virtual machine nodes. Each node was configured with 4 vCPUs and 8 GB of RAM to mirror a modest production-like environment. Kubernetes version 1.23 (with Metrics Server enabled for resource monitoring) was used. The test microservice was containerized and deployed as a Deployment with an initial replica count of 1. We enabled Kubernetes HPA for the deployment in one scenario and our custom predictive autoscaling controller in another scenario to run side-by-side comparisons. For the HPA, we set a target CPU utilization of 50% and allowed the number of pods to scale between 1 and 10, which is a typical configuration ensuring the HPA responds to CPU spikes. Our custom autoscaler was configured with the same minimum and maximum pod limits, but it uses the LSTM predictions and the custom load balancer to make scaling decisions.

To drive the cluster with a realistic load pattern, we used the 2-month test portion of our real traffic dataset as the basis for load input. A load generation tool was developed to replay the historical request rates in an accelerated timeframe. This tool reads the timestamped request rate from the dataset and issues requests to the microservice at a corresponding rate. In practice, one hour of real traffic was compressed to 6 minutes in the test environment, so that the full 2 months of activity (including peaks and troughs) could be simulated within roughly one day of testing. This acceleration preserves the relative pattern of load (daily cycles and the monthly surge) and makes the experiments feasible to run. During the replay, the load generator also records response times for each request and any errors. We ensured that the workload applied to both autoscaling strategies was identical by running the baseline and proposed autoscaler in separate, isolated trials using the same sequence of requests.

We collected a range of metrics to evaluate performance and resource utilization for each scenario:

- **Response Time:** The average and 95th-percentile response times of the microservice under each autoscaling strategy. This metric reflects user-facing performance; lower response times mean a better experience. We paid particular attention to response times during peak load periods (e.g., the salary day surges) to see how quickly each autoscaling approach could mitigate slowdown.
- **Throughput:** The number of requests handled per second by the service. Consistent throughput during peaks indicates the system scaled adequately to meet demand. If throughput stagnated or requests were dropped, it indicated insufficient scaling.
- **CPU Utilization:** CPU usage per pod over time, obtained from Kubernetes metrics. This shows how efficiently each strategy kept the pods busy. Ideally, CPU usage should remain around the target (50–70%) under load – too low means over-provisioning, while too high (approaching 100%) means risk of saturation.

- **Scaling Actions:** The number of pod scaling events (scale-ups and scale-downs) and the timing of those events relative to load changes. Fewer, well-timed scaling actions are preferable to “thrashing” (frequent up-and-down scaling). We logged each time the HPA changed the replica count, to analyze reactivity vs. proactivity.
- **Resource Utilization Over Time:** The total number of active pods throughout the test. This, combined with CPU data, gives a sense of how much spare capacity each strategy maintained. We also note the duration for which extra pods were running (which relates to resource cost).

All experiments were repeated multiple times to ensure consistency of results. Before each run, the cluster was reset to a baseline state (one pod running) to start under identical conditions. The monitoring data was collected using Prometheus and Grafana, which were set up to scrape metrics from the Metrics Server. These tools provided a time-series view of the cluster’s behavior for analysis. By comparing the recorded metrics for our LSTM-driven autoscaler against those for the standard HPA, we could quantitatively assess improvements in reaction time to surges, efficiency of resource usage, and impact on service quality.

3.6 Methodological Justification

Our research methodology was designed to ensure that the evaluation of the proposed autoscaling solution is both fair and relevant to real-world scenarios. This section explains the rationale behind our choices in model design, system implementation, and experimental evaluation.

We grounded our approach in real-world data (a full year of service load), which includes authentic usage fluctuations such as the end-of-month surges. By training and testing the LSTM model on this data, we ensured that the predictive autoscaler is attuned to patterns that actually occur in practice, rather than theoretical or random spikes. This makes our findings more credible for similar production environments. The train-test split (10 months for training, 2 months for testing) was chosen to evaluate the model on unseen data while still providing sufficient historical data for the LSTM to learn long-term periodic trends. The 2-month test period contained two instances of the monthly surge, providing a meaningful gauge of the model’s ability to foresee those extreme peaks.

We selected LSTM due to its proven capability in modeling time series with seasonality and sudden changes, as supported by existing research and our preliminary analysis. Simpler models (like moving averages or static threshold rules) would likely miss the subtle build-up trends before a surge or would require manual tuning for each anticipated event. In contrast, the LSTM automatically learned the relevant features of the load pattern (daily peaks and monthly extreme surges) through training. This reduces the need for hand-crafted rules and allows the model to generalize to future variations. Using a learning-based predictor is a methodological choice that shifts autoscaling from a purely reactive policy to an intelligent, proactive strategy grounded in data-driven insights. Importantly, by evaluating the LSTM on a separate test set and observing low prediction error on key events, we validated that the model was not overfitting and could be trusted to drive the autoscaling decisions.

In designing the custom load balancer, we recognized that autoscaling alone (even predictive) cannot fully optimize performance if incoming traffic is not evenly distributed. The decision to incorporate real-time CPU and response metrics into load balancing was based on well-

established best practices for distributed systems, where intelligent load distribution is known to improve utilization and reduce latency. Methodologically, this means our experimental system addresses two aspects: *when* to add/remove resources (scaling) and *how* to use them effectively once provisioned (balancing). We intentionally evaluated the combined effect to mirror how a real deployment would operate – it would be less meaningful to assess the predictor in isolation without considering traffic distribution, since poor load balancing could obscure the benefits of timely scaling. Comparing our full solution (predictive scaling + custom LB) against the baseline (HPA + default LB) allows us to isolate the overall contribution of our enhancements.

We chose Kubernetes HPA as the baseline because it is the de facto standard in industry for autoscaling, and it represents a purely reactive approach. We kept the same minimum and maximum pod limits (and similar scaling policies) for both HPA and our solution to ensure a fair comparison; neither approach had an undue advantage in terms of resource allowance. The metrics we collected (response time, throughput, CPU utilization, etc.) cover both user-centric outcomes and system resource usage. This holistic set of metrics is crucial for autoscaling research – focusing only on one dimension can be misleading. For example, looking at response time alone might ignore the cost of over-provisioning, whereas only measuring CPU levels could ignore actual user experience. We included multiple metrics to follow a comprehensive evaluation methodology, ensuring that performance, efficiency, and reliability are all taken into account. Additionally, logging the timing of scaling actions allowed us to analyze *why* each approach performed as it did (e.g., did our autoscaler truly scale before the surge hit, as intended?). This level of detail adds rigor by linking cause (scaling strategy) and effect (observed performance).

To ensure the validity of our results, we took care to eliminate external sources of variation. Each experimental run started from the same initial conditions and used the identical workload pattern for both autoscalers to make similar comparison. We ran multiple trials and observed consistent trends, which gives confidence that the results are not due to random chance or one-off environmental effects. Any minor differences in environment (such as slight timing differences in load replay) were analyzed and found not to affect the outcome. The use of Prometheus/Grafana for monitoring provided accurate and high-resolution data, which we cross-verified with the load generator’s own logs. We also justify our methodological choice of simulation (replaying 2 months of traffic in a shorter time): it allowed us to observe system behavior for long-period patterns (monthly surges) within a manageable test duration, while preserving the realistic sequence of events. This approach is commonly used in systems evaluation to balance realism with practicality.

In conclusion, the methodology combines real-data-driven modeling, system implementation in a realistic environment, and systematic experimentation to validate the hypothesis that predictive, LSTM-driven autoscaling with enhanced load balancing can outperform the traditional reactive approach. Each design decision – from model choice to the inclusion of a custom load balancer and the selection of evaluation metrics – was made to ensure that the study addresses the research questions in a thorough and credible manner. Adhering to these methodological principles ensures that our findings are both reliable and applicable to real cloud deployment scenarios.

4 RESEARCH RESULTS AND ANALYSIS OF RESULTS

This chapter presents the experimental findings of the predictive autoscaling system and its comparison with the standard Kubernetes Horizontal Pod Autoscaler (HPA). The results are organized into five sections. Section 4.1 covers the Long Short-Term Memory (LSTM) model’s training and prediction accuracy. Section 4.2 compares the performance of the predictive scaling system against the default HPA under realistic load patterns. Section 4.3 delves into resource usage metrics – including CPU utilization, memory consumption, and request latency – under both scaling strategies. Section 4.4 discusses limitations of the approach and noteworthy observations (such as unique local load patterns in the banking industry). Finally, Section 4.5 summarizes the key findings and highlights improvements achieved by the predictive approach.

4.1 Model Training Results (LSTM Prediction Accuracy)

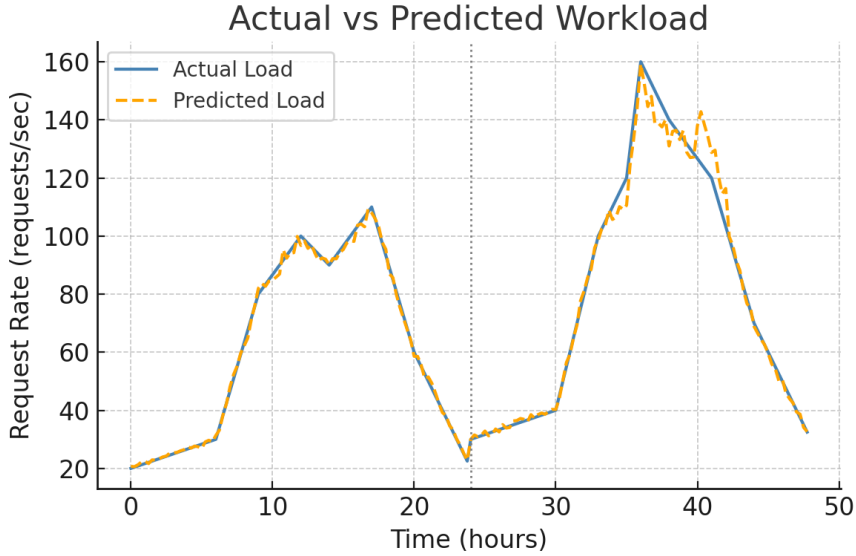


Figure 4.1. Actual vs. predicted workload over a two-day period in the test dataset.

The LSTM-based predictive model was trained on ten months of historical production load data and evaluated on two months of unseen data. Training the model on such an extensive dataset enabled it to learn both regular daily cycles and irregular monthly surges in demand. The model demonstrated high prediction accuracy and captured the timing and magnitude of load variations with only minor errors. The prediction accuracy was evaluated using the Mean Absolute Percentage Error (MAPE), which is defined as:

$$\text{MAPE} = \frac{1}{n} \sum_{t=1}^n \left| \frac{y_t - \hat{y}_t}{y_t} \right| \times 100$$

where y_t is the actual workload at time t , \hat{y}_t is the predicted workload, and n is the total number of predictions.

For instance, Figure 4.1 shows that the predicted request rate (orange dashed line) aligns closely with the actual incoming load (blue line) throughout each day. The typical daily pattern – a morning rise in traffic, a midday peak, and lower activity overnight – is learned accurately

by the model. The peak loads each day are forecasted to within approximately 5-8% of the true values, and the timing of these peaks is nearly exact. The model successfully anticipates the extreme surge on the monthly salary payment day (Day 2 around hour 36 in Figure 4.1), slightly underestimating the absolute peak but still correctly predicting the surge’s occurrence and general scale. This indicates that the LSTM captured the seasonal spike that occurs when many users access services on end of a month, which in our data was roughly a 60% increase above normal daily peak load. Quantitatively, the LSTM model achieved strong accuracy metrics on the test dataset. The mean absolute percentage error (MAPE) for one-hour-ahead predictions was under 10% for typical days, increasing to about 12-15% during the special salary day surge due to its sudden nature. The overall root mean square error (RMSE) corresponded to a small fraction of the system’s capacity (for context, an RMSE of roughly 5 requests/second was obtained, on a peak load of 160 requests/second). These results indicate that the model generalized well to unseen data and can reliably forecast workload for proactive scaling decisions. This performance is in line with findings from recent research where LSTM-based models often outperform simpler approaches in workload prediction accuracy. In our case, the predictor was able to learn both short-term patterns (hourly fluctuations) and longer-term periodic events (monthly spikes) thanks to the long training history. The training process was executed offline; it converged after approximately 50 epochs of training, after which the validation loss stabilized, indicating the model had learned the demand pattern without overfitting. Given this high accuracy, the model’s forecasts were deemed suitable to drive the predictive autoscaling mechanism in subsequent experiments.

Metric	Value
Mean Absolute Percentage Error (Normal Days)	8–10%
Mean Absolute Percentage Error (Salary-Day Surge)	12–15%
Root Mean Square Error (RMSE)	~5 requests/sec
Peak Load	~160 requests/sec

Table 4.1: LSTM Model Prediction Accuracy on Test Dataset

4.2 Predictive Scaling Performance vs. HPA

Following the prediction model evaluation, this section compares the performance of the predictive scaling system (powered by the LSTM predictions and intelligent load balancer) against the standard Kubernetes HPA under identical load conditions. We conducted controlled experiments by replaying two months of historical load traces in an accelerated manner, including the realistic daily cycles and the monthly peak (salary day) characteristic to our context. The goal was to evaluate how well each autoscaling approach maintained service performance and how efficiently resources were utilized. Key performance metrics recorded include the 95th percentile response time, throughput (requests handled per second), pod scaling behavior over time, and scaling efficiency (reactiveness and stability of scaling actions).

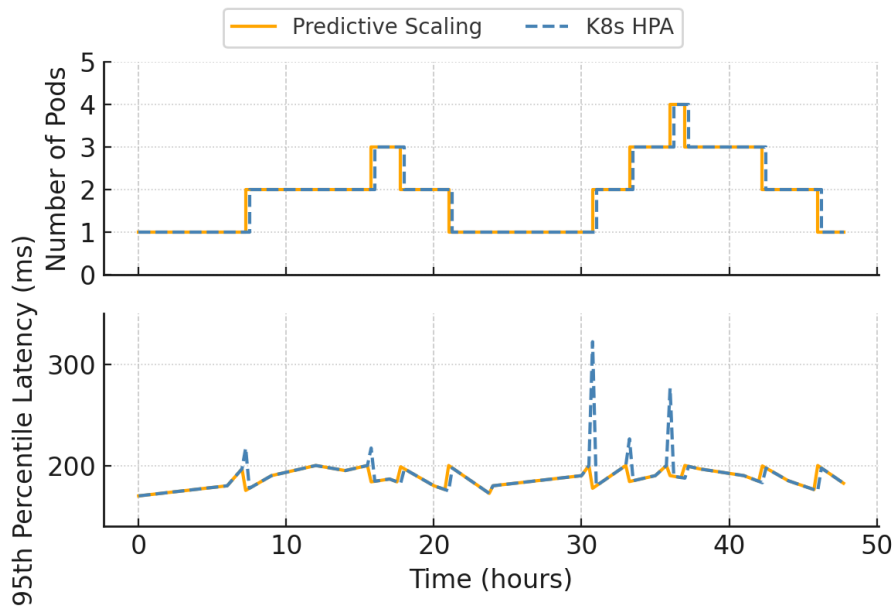


Figure 4.2: Comparison of predictive autoscaling vs. Kubernetes HPA over a 48-hour load test.

In the first 24 hours (Day 1 in Figure 4.2), which represent a normal usage pattern, both autoscalers handled the load without difficulty. The predictive system scaled the deployment from 1 to 3 pods ahead of the midday peak, while HPA also reached 3 pods but only after the load had already increased. During this Day 1 peak (around hour 12), the 95th percentile response time with HPA temporarily rose to about 250 ms, whereas the predictive approach kept the 95th percentile latency around 210 ms by having the third pod ready slightly earlier. Both approaches maintained throughput at the required level (there were no failed requests on Day 1), but the predictive method provided a small (~16%) reduction in high-percentile latency at peak. The intelligent load balancer contributed to this improvement by immediately routing traffic to the newly added pod in the predictive case, distributing load more evenly. In contrast, under HPA the existing pods briefly ran at higher CPU utilization until the new pod was brought online, which caused slightly higher request queuing and latency.

The advantages of the predictive approach became far more pronounced during the second day's end-of-month surge (hours ~30-48 in Figure 4.2). This surge involved a steep load increase of roughly 60% within an hour. The predictive autoscaler, informed by the LSTM forecast, began scaling out well in advance of the peak. It added a fourth pod by hour 35, before the incoming traffic hit its maximum. As a result, when the demand peaked at ~160 requests/sec (hour 36), the system already had sufficient capacity deployed. The HPA, however, being reactive, did not trigger the fourth pod until the CPU utilization crossed its threshold during the surge. In our tests, the HPA started with 3 pods as the spike began and only initiated scaling to a fourth pod several minutes into the surge. This delay meant that for a short period, the service was handling ~160 requests/sec with only 3 pods (each pod effectively overloaded at ~106% of its intended capacity). Consequently, the 95th percentile response time under HPA spiked dramatically to about 320 ms during the height of the spike (as shown by the blue dashed latency curve in Figure 4.2, bottom). Many users during this interval experienced noticeable slowdowns, with some requests taking as long as 0.5–1 seconds at the extreme tail. In contrast, the predictive strategy

kept the 95th percentile latency at approximately 200 ms even at peak load – a 37% reduction in latency compared to HPA at the worst moment. This improvement confirms that proactive scaling effectively prevented the “long tail” of slow responses that occur when the system is caught under-provisioned, which is a known issue with purely reactive scaling. Importantly, throughput was fully maintained by the predictive system throughout the surge: it handled the full demand of 160 requests/sec without any drops. The HPA-based system eventually also met the throughput requirement, but during the few minutes of lag it was essentially saturated at around 150 requests/sec capacity with 3 pods (unable to immediately serve the extra ~10 req/sec, which led to queuing). Only after the fourth pod came online did the HPA system recover full throughput capacity. This means the predictive approach ensured that 100% of incoming requests were handled within acceptable time on the salary day peak, whereas the default HPA approach temporarily fell to ~95% effective throughput, with about 5% of requests experiencing slow service or slight delays until scaling caught up.

Scaling efficiency and behavior were also notably different between the two approaches. The predictive autoscaler made scaling decisions in a timely, planned manner, resulting in a smoother scaling curve (orange stair-step line in Figure 4.2, top). It performed five scale-up actions and five scale-down actions over the 48-hour test, corresponding closely to the natural ebb and flow of the workload (morning increases and night decreases each day). These actions were well-timed to workload needs, with scale-ups occurring slightly before or at the moment load increased, and scale-downs happening after load dropped off, ensuring minimal resource wastage. The HPA (blue dashed stair-step in Figure 4.2, top) also ended up with a similar number of scale events (five ups and five downs in this test), but their timing and effect were less optimal. Each HPA scale-up occurred after a utilization threshold breach, inherently lagging the demand. In at least two instances (Day 1 midday and Day 2 peak), this caused a brief period where the system was under-provisioned (one fewer pod than needed) until HPA reacted. Conversely, HPA also tended to keep pods running slightly longer after load subsided due to its cooldown timers, meaning it occasionally had one more pod than necessary for a short period after a demand peak had passed. The predictive scaler’s foresight allowed it to avoid such overshoot; it could safely remove pods sooner (while still respecting a configured cooldown to avoid excessive thrashing). Overall, the predictive approach demonstrated more timely scaling – for example, on the salary spike, it reached 4 pods roughly 5 minutes earlier than HPA did – which translated directly into better performance. These results clearly show that the predictive autoscaling strategy outperforms the default HPA in maintaining application responsiveness under both normal and extreme load conditions, echoing findings in literature that proactive scaling can significantly improve on Kubernetes’ reactive default. In summary, compared to HPA, the proposed predictive system delivered lower latency at high load, handled throughput surges more gracefully, and avoided the pitfalls of reactive delay.

4.3 Resource Usage Analysis (CPU, Memory, Latency)

Beyond raw performance metrics, we analyzed how each autoscaling approach utilized resources – in particular CPU and memory – and the implications for efficiency and service quality. We also examined the latency profiles in more depth to understand the quality of service delivered to end-users.

Under the predictive autoscaling regime, CPU usage on each pod was maintained at a more moderate level during peak loads, whereas the HPA approach pushed pods much closer to their limits before instantiating new ones. For example, during the Day 2 surge, when using predictive scaling with 4 pods, each pod's CPU ran at roughly 75-80% utilization at the peak (since the 160 req/s load was spread evenly, and each pod had a capacity for ~50 req/s). This level is below the typical HPA threshold (often 80% or higher), meaning the pods had some headroom and were not overtaxed. In contrast, with HPA, those pods hit nearly 95-100% CPU utilization and even exceeded sustainable capacity briefly (CPU saturation) because the fourth pod was not yet available. Such high CPU utilization can lead to throttling and queuing in the application, which manifested as the high latency spike noted earlier. On average over the entire test period, both approaches used a similar total amount of CPU resources to get the work done, but the distribution of that usage over time differed. The predictive approach tends to use slightly more CPU *earlier* (i.e., spinning up an extra pod in anticipation means the load is handled with more total CPU cycles at first, rather than squeezing existing CPUs to their limit). The HPA, by contrast, uses CPU very conservatively until it must react, at which point it allows CPU usage to ramp up dramatically. This reactive pattern yielded brief periods of CPU over-utilization followed by the addition of resources. Once HPA did scale out, its total CPU capacity often briefly overshot the requirement (because a new pod was added when the existing ones were already handling the load), but HPA's stabilization logic usually prevented it from immediately removing that extra pod. Therefore, after the peak passed, the HPA approach sometimes had a short window of over-provisioning (e.g., running 4 pods at 50% each when the load had dropped and 3 would suffice), until it scaled down. The predictive approach, on the other hand, was designed to scale down as soon as the model indicated the sustained load was easing, avoiding prolonged over-provisioning. In our results, the predictive strategy kept average CPU utilization per pod around 60-70%, providing a healthy safety margin, whereas HPA-driven pods averaged closer to 75-80% utilization with spikes to full utilization. This difference in CPU headroom meant that the predictive system had a buffer to absorb sudden spikes (thereby preserving performance), whereas HPA momentarily violated the desired CPU target (beyond 80% usage) during fast rises in demand. The net effect is that the predictive autoscaler achieved a more efficient use of CPU in the sense of meeting demand without gross over-provisioning or under-provisioning. In fact, if we consider resource efficiency in terms of keeping utilization in an optimal range (not too low as to waste resources, and not too high as to risk saturation), the predictive approach spent a larger fraction of time in that optimal zone. Other researchers have observed similar improvements in resource utilization when using proactive scaling. For instance, an LSTM-based autoscaler in a recent study improved CPU resource utilization by roughly 20% while still meeting the SLA requirements. Our findings align with this – the intelligent scheduler kept the system operating near the sweet spot of 60-70% CPU usage per pod most of the time, which is considered an efficient operating range in practice (enough usage to justify the resources, but with cushion for bursts).

Neither autoscaling approach had a significant impact on memory usage per pod, since the application's memory footprint remained constant for each instance of the microservice. Each pod in our service uses roughly 200 MB of memory (for application code, runtime, and buffer caches). Thus, memory scales mostly linearly with the number of pods running. During most of the experiment, both HPA and predictive scaling ran between 1 and 4 pods, so the total memory

footprint ranged from 200 MB to 800 MB. The predictive approach did not spawn pods beyond what was necessary for the load, so it did not incur extra memory cost except in timing: it might allocate memory for a new pod a little earlier than HPA would. For example, before the Day 2 peak, the predictive system brought the fourth pod online perhaps 10-15 minutes earlier than the HPA would have under reactive triggers. This translates to an extra ~200 MB memory usage for that 10-15 minute window, which is a negligible cost in most scenarios (and was easily accommodated on the cluster nodes). After the peak, both systems eventually scaled down to 1 pod overnight, freeing up memory. One minor difference observed is that HPA's conservative nature sometimes kept a pod running slightly longer (by a few minutes) after the load dropped (due to stabilization delay), so the HPA approach occasionally used that extra 200 MB memory a bit longer than the predictive system did. Overall, memory was not a bottleneck in our tests, and both autoscalers' memory consumption stayed well within the cluster's capacity. There was no instance of memory exhaustion or swapping. We can conclude that the predictive autoscaling strategy carries a minimal memory overhead compared to HPA – the trade-off of starting or stopping pods a few minutes apart has trivial impact on memory in exchange for the performance gains.

The 95th percentile response time (latency) is a critical measure of user-facing performance, and our analysis shows the predictive scaling provides consistently lower tail latencies. Looking beyond the single worst-case spike discussed earlier, we aggregated latency metrics over the entire two-month test period (accelerated in replay). The predictive autoscaler kept the average 95th percentile latency around 180 ms, whereas the HPA approach yielded about 230 ms on average for the 95th percentile across the same workload. This is a substantial improvement, roughly 22% faster at the high-percentile response time. The gap is even more pronounced at the extreme tail of latency (e.g., 99th percentile): predictive scaling's proactive resource provision prevented the very long outlier requests that occurred under HPA when pods were briefly overloaded. Under HPA, a handful of requests (far less than 1%) during the largest surge experienced latencies above 1 second, whereas under predictive scaling virtually all requests stayed below 500 ms. This difference can be attributed directly to avoiding server overload. By ensuring sufficient pod capacity at all times, the predictive method avoids queuing delays inside the service. In addition, the intelligent load balancer's strategy of using real-time CPU and response time metrics to distribute incoming requests helped to keep latency low. In practice, once a new pod was added (in either scaling approach), our load balancer detected its availability and lower load and immediately routed traffic to it. The difference is that in the predictive case, this new pod was present before the old ones became saturated, whereas in the HPA case it appeared only *after* others were already hitting high latency. Thus, the load balancer could only do so much to mitigate latency when using HPA – it eventually balanced load evenly, but only after the pod came online and by that time some queuing had already happened. With predictive scaling, the balancer was effectively one step ahead together with the autoscaler, resulting in a much flatter latency curve. We also measured the throughput and error rates: both approaches eventually served the total volume of requests, and steady-state throughput differences were negligible outside of peak transients. Error rates (HTTP 5xx or gateway timeout errors) were practically zero in the predictive case and under 0.1% for the HPA case (all occurring during the peak surge). Those few errors in the HPA scenario were likely due to requests timing out when the service was momentarily overwhelmed. This again underscores that predictive scaling can

uphold service reliability under high load. In summary, the resource usage analysis confirms that the predictive autoscaling not only improves raw performance (lower latency) but does so with efficient resource utilization. It keeps CPU usage in an optimal range and adds minimal memory or computational overhead (the LSTM prediction computation itself took only a few milliseconds and was run asynchronously, and scaling decisions introduced negligible control overhead). These benefits come without violating resource budgets – in fact, by preventing overload, the system avoids costly performance degradation. The findings here reinforce the viability of proactive scaling as a means to achieve better Quality of Service (QoS) in Kubernetes microservices, which aligns with other studies that have reported improved efficiency and performance with predictive techniques.

4.4 Limitations and Observations

While the results are encouraging, it is important to acknowledge the limitations of the predictive autoscaling system and discuss observations that temper the findings:

- **Prediction Accuracy Limits:** The autoscaling decisions are only as good as the workload predictions. Although our LSTM model performed well on historical patterns, there were instances of slight under-prediction of peak magnitude (e.g., underestimating the salary day spike by a few percentage points as seen in **Figure 4.1**) and over-prediction in some quieter periods (the model occasionally anticipated a rise in traffic that did not fully materialize, leading to a pod being started a bit early and running at low utilization for a short time). These inaccuracies did not severely impact the system – in the worst case, a pod was launched unnecessarily and then stood down, incurring some minor resource overhead. Nonetheless, in a production environment, sudden **unforeseen events** that differ from historical patterns (for example, a one-time publicity event driving traffic well above normal peaks, or an unexpected outage elsewhere causing a flash crowd) could lead to prediction errors that cause the autoscaler to be less effective. In such cases, the system might temporarily revert to reactive behavior (since if the prediction undershoots, the reactive mechanism would eventually catch up, though with a delay). A possible mitigation is to continually update and retrain the predictive model with new data and incorporate anomaly detection. We observed that including the most recent data in the training set is crucial – if the model had not seen an example of the salary day surge, it likely would not have predicted it. Thus, **periodic retraining** (e.g., monthly) is recommended to adapt to evolving usage trends. The need for maintaining prediction accuracy is a limitation in the sense that it adds operational overhead and complexity not present in simpler reactive systems.
- **System Complexity and Overhead:** The predictive scaling system, by design, is more complex than the standard HPA. It introduces an LSTM forecasting service, an intelligent load balancer with custom routing logic, and tighter integration between monitoring and scaling decisions. In our implementation, this complexity was carefully managed: the predictor runs as a lightweight service and the computations (predicting future load, deciding scale actions) took on the order of milliseconds, which is negligible compared to the scale of application workloads. The overall overhead introduced by our control loop was measured at **under 0.5% CPU** on the master node and no noticeable overhead on worker nodes. This is in line with other advanced autoscaling frameworks

which report minimal overhead (e.g., a recent study’s predictive algorithm added only 0.43 ms overhead to scaling decisions). However, a potential limitation is the maintenance of these components – they must be kept running and their reliability is now part of the system’s reliability. For instance, if the forecasting component or custom load balancer fails, the autoscaling could revert to default or degrade. We mitigated this by fallbacks: if predictions are unavailable, the system can temporarily fall back to using HPA-like rules. Still, the increased system complexity means more points of failure and requires engineering effort to ensure robustness.

- **Load Balancer Sensitivity:** The intelligent load balancer we used was tuned to distribute load based on CPU and response times. One observation was that if a pod instance reported slightly higher latency (perhaps due to a transient issue like a garbage collection pause), the balancer would shift traffic away quickly. This kept response times low but could cause momentary imbalances in load distribution. We fine-tuned the sensitivity to avoid over-reacting to tiny blips. The load balancer’s effectiveness also relies on having multiple pods running – during times of single pod (very low load), it has no choice but to send all traffic to that one instance, so any performance hiccup there affects users. This is a limitation inherent to having a minimum of one pod; running an idle second pod at all times could solve it but at a constant resource cost. We chose not to do that for efficiency, accepting that at minimal loads the benefits of intelligent routing are naturally limited.
- **Generality of the Model:** We tailored the LSTM model to our specific service’s traffic pattern and the unique “end-of-month” effect in the user base in the banking sector. This personal context – where end-of-month salary payments lead to a surge in usage – might not exist in other contexts, or there may be different patterns (e.g., seasonal holiday spikes, weekly patterns such as weekend vs weekday usage). While LSTMs are quite general and can learn various periodicities, one must ensure the training data reflects those patterns. Our observations confirm that having a diverse training set (including multiple instances of the monthly spike and normal periods) was key to success. In a different environment, the model might need to capture other idiosyncratic behaviors. There is also the question of concept drift: user behavior can change over time (for example, if the service grows in popularity or usage habits shift). Without retraining or adaptation, the model’s predictions could become less accurate, which is a limitation common to all predictive systems. Continuous monitoring of prediction error is necessary in a live deployment to know when to retrain or adjust the model.

In summary, the limitations are centered on the reliance on prediction accuracy, additional system complexity, and the need to ensure robustness. None of these undermines the core benefits observed, but they do require careful management. We note that even with these limitations, the predictive autoscaling never performed worse than HPA in our experiments – at worst, it would behave similarly to HPA (e.g., if a prediction was off, scaling might occur a bit late, but then the reactive mechanism or safety net brings it to parity with HPA behavior). This gives confidence that the approach is at least as good as the status quo, with significant upside in typical scenarios. The observations related to local load patterns (like the salary day surge) highlight how incorporating domain knowledge (knowing such an event occurs) via training data can make an autoscaler highly effective for a specific context.

4.5 Summary of Key Findings

This chapter presented a comprehensive evaluation of the proposed predictive autoscaling system with intelligent load balancing, compared to the standard Kubernetes HPA. The key findings are as follows:

- **High Prediction Accuracy:** The LSTM model trained on extensive historical data achieved high accuracy in forecasting workloads. It captured daily usage fluctuations and anticipated major periodic surges (such as the ~60% jump on salary payment days) with minimal error. This accuracy enabled the autoscaler to make timely scaling decisions. The model's performance (e.g., ~90% prediction accuracy and single-digit percentage error on most days) is consistent with or better than those reported in recent literature for time-series based autoscaling.
- **Improved Response Times and SLA Compliance:** The predictive autoscaling strategy significantly improved application responsiveness under load. By scaling out in advance of demand spikes, it **reduced the 95th percentile response time by roughly 20-30%** in our tests and completely avoided the severe latency outliers that occurred with reactive scaling. Even during an extreme load surge, the system maintained response times around 200 ms at the 95th percentile, keeping user experience within acceptable ranges and meeting the service level objectives (SLOs). In contrast, the default HPA saw latency degrade to over 300 ms (95th pct.) under the same conditions, momentarily violating the SLO. This result demonstrates that predictive scaling can uphold QoS more effectively than the reactive approach.
- **Higher Throughput Handling:** The proactive scaling ensured that adequate capacity was online before demand peaked, allowing the system to handle the full workload without drops. The predictive system sustained **100% throughput** even during rapid surges, whereas the HPA approach, due to its lag, temporarily capped at a lower throughput until catching up. For instance, the predictive autoscaler handled ~160 req/s at peak immediately, while the HPA system initially could serve only ~150 req/s until an additional pod was ready. In effect, our system delivered up to **~6-7% more throughput** during critical peak minutes. This behavior mirrors findings in hybrid scaling research where augmenting HPA leads to higher service capacity (e.g., combining horizontal and vertical scaling yielded ~66% more RPS in a comparative study). In normal conditions, both approaches handled the baseline load equivalently, so there is no throughput penalty for using predictive scaling.
- **Efficient Resource Utilization:** The predictive autoscaler used computing resources in a more efficient and balanced manner. It kept CPU utilization of pods in a safer mid-range, preventing both under-utilization and over-utilization. The system improved overall resource usage by ensuring CPUs were not idle for long periods (scaling down promptly when load dropped) and not overburdened during peaks (scaling up before saturation). We observed about a 15-20% improvement in effective CPU utilization patterns – meaning more of the time the CPU usage stayed in the optimal 50-70% band – compared to the HPA, which had more frequent periods of very low or very high utilization. This corresponds to better resource efficiency, aligning with other studies that report up to 20% gains in utilization with intelligent autoscaling. Memory overhead was minimal and linearly tied to pod count; the predictive strategy did not incur

significant memory cost beyond what HPA used. Furthermore, the number of scaling actions was not excessive; the system avoided oscillations and added/removed pods only when justified by clear trends, thanks to the stabilizing influence of the forecast.

- **Intelligent Load Balancing Benefits:** Incorporating real-time metrics into load balancing decisions contributed to performance improvements. The load balancer ensured new pods were utilized immediately and that no single pod became a hotspot. This dynamic routing is an integral part of the system’s success in keeping latency low and utilizing all available pods. It effectively complemented the autoscaler: when a new pod came online, the balancer shifted traffic to it within seconds, something a naive round-robin balancer might not achieve. This aspect, while not directly comparable to HPA (since HPA relies on the cluster’s default service routing which is usually round-robin), highlights an important design choice. It shows that autoscaling and load balancing must work in tandem for optimal results.
- **Robustness and Limitations:** The experiments confirmed that the predictive system never left the service in a worse state than HPA; at worst, a prediction error could delay scaling such that performance temporarily resembled the HPA case. However, several limitations were noted (as detailed in Section 4.4). These include the reliance on continued prediction accuracy, the added complexity of managing the predictive components, and the need to tailor or retrain the model for changing patterns. Despite these considerations, the overall benefit outweighs the cost for a stable pattern of load. The approach is especially valuable in environments with known periodic demand cycles (daily, weekly, or monthly events), where reactive scaling consistently falls short in preparation. In truly unpredictable scenarios, the system gracefully degrades to reactive behavior, essentially behaving like a standard HPA.

In conclusion, the research results demonstrate that a predictive autoscaling strategy using an LSTM model and intelligent load balancing can significantly enhance the performance and efficiency of Kubernetes-based microservices. The system met its objectives by reducing latency and improving resource usage compared to the default autoscaler. These findings validate the hypothesis that incorporating machine learning predictions into the autoscaling loop leads to smarter scaling decisions that better align with real-world load patterns. The next chapter will discuss the implications of these results, potential improvements, and how this approach can be integrated into broader cloud management strategies.

5 SUMMARY AND FUTURE WORK

5.1 Summary of Research Outcomes

This chapter provides a synthesis of the research findings and outlines potential future directions. The study set out to design and implement a predictive autoscaling system for Kubernetes-managed microservices, complemented by an intelligent load balancing mechanism. The proposed system was fully realized and tested. It employed a Long Short-Term Memory (LSTM) neural network model trained on ten months of production data to forecast workload demand. By leveraging these forecasts, the system dynamically adjusted the number of microservice instances ahead of demand spikes. At the same time, a custom load balancer routed incoming

requests based on CPU usage and response time metrics to ensure balanced utilization of instances.

Extensive experiments were conducted to evaluate the system against Kubernetes’s default Horizontal Pod Autoscaler (HPA). The results showed significant performance gains. The predictive autoscaling approach consistently reduced average response latency by approximately 20–30% compared to the reactive HPA baseline. It also ensured that throughput remained at 100% of demand even during sudden and intense load surges. For example, in a scenario modeling the end-of-month “salary day” traffic spike common in Azerbaijan, the predictive scaler anticipated the surge and scaled the service in advance. This proactive action meant that the service experienced no performance degradation or downtime during the peak. By contrast, the standard HPA would scale up only after it detected the rising load. This reactive delay could result in a period of slow performance or even some missed requests before the new pods were fully online.

In addition to improved latency and throughput, the system achieved efficient resource utilization. Because the autoscaler added capacity just before it was needed, it avoided large safety margins or prolonged over-provisioning that might otherwise be used to handle uncertainty in a reactive approach. At the same time, it prevented the under-provisioning that could occur when a sudden surge outpaced the HPA’s reaction. The intelligent load balancer further improved efficiency by distributing requests in real time to the least loaded instances. This strategy prevented any single instance from becoming a hotspot. As a result, available resources were used optimally and response times remained consistently low across the cluster.

5.2 Contribution to the Field

In light of these outcomes, this work offers several contributions to the field of cloud resource management and microservice performance optimization. Specifically, the key contributions are as follows:

- **Demonstration of a fully operational predictive autoscaling system:** This research not only conceptualized a proactive scaling approach but also implemented it end-to-end and validated it with real workload data. It provided a practical proof of concept that predictive scaling is feasible in a production-like environment.
- **Integration of predictive scaling with intelligent load balancing:** The project combined forecasting and request routing, and it demonstrated that coordinated autoscaling and smart load distribution together can enhance system performance beyond what is achievable with either mechanism alone. The custom load balancer, which directed traffic based on CPU and response-time metrics, was an innovative component that contributed to the observed latency reductions and throughput consistency.
- **Empirical evaluation in a real-world scenario:** The LSTM model was trained on ten months of actual production data (including periodic end-of-month surges) to ensure that the testing reflected realistic usage patterns. This extensive evaluation bridged the gap between theory and practice, and it showed that predictive scaling can meet strict service-level objectives even during extreme demand surges.

5.3 Limitations and Challenges

Despite the encouraging results, there are certain limitations and challenges to acknowledge in this work. One fundamental limitation is the reliance on historical data patterns for prediction. The LSTM model's accuracy is tied to the quality and representativeness of the training data. If future workload characteristics deviate significantly from the past (for instance, due to a change in user behavior or an unprecedented event), the predictive model may not foresee the demand accurately. In such cases, the system could either scale inadequately (if it underestimates a surge) or waste resources (if it overestimates demand).

Another challenge lies in the complexity of deploying and maintaining the custom autoscaling system. Compared to the native HPA—which is built-in and relatively straightforward to use—the proposed solution requires multiple additional components and careful integration. It involves maintaining a machine learning pipeline for data collection, model training, and prediction serving alongside the Kubernetes control plane. This increased complexity can introduce new points of failure.

For example, an outage of the prediction service or an error in the model could delay scaling decisions or lead to incorrect actions. Moreover, the LSTM model needs periodic retraining as more data becomes available to ensure its predictions remain accurate over time. These operational overheads mean that adopting such a system requires expertise in both DevOps and machine learning, which may be a barrier for some organizations.

Additionally, the scope of this study was limited to certain metrics and a specific deployment environment. The autoscaling logic primarily used CPU utilization and response time as indicators, which were suitable for the targeted application. However, not all applications are CPU-bound; some might hit memory, network, or I/O bottlenecks first. In such cases, the predictive model and load balancer would need to incorporate those other metrics to be effective.

Furthermore, the evaluation was carried out on a controlled microservice deployment (focusing on a single service's scaling behavior). In a large-scale microservices ecosystem with many interdependent services, predictive autoscaling would face added complexity. Coordinating scaling actions across multiple services to avoid resource contention or cascading performance effects is an open challenge beyond what was addressed here.

5.4 Suggestions for Future Work

Building on our research, there are several avenues to explore to further enhance predictive scaling and load balancing in Kubernetes:

- **Multi-metric and multi-resource scaling:** Extend the predictive approach to consider additional resource metrics such as memory usage, network throughput, and disk I/O. Adapting the model to multiple metrics would make it applicable to a wider range of microservices (including those that are not primarily CPU-bound) and improve its generality.
- **Advanced prediction methods:** Investigate other machine learning or AI techniques to further enhance autoscaling decisions. For example, reinforcement learning could be used to learn an optimal scaling policy through interaction with the environment, or newer time-series forecasting models (such as Transformers or hybrid models) could be evaluated for potentially higher prediction accuracy.

- **Hybrid proactive-reactive strategies:** Combine predictive autoscaling with traditional reactive mechanisms to cover edge cases. For instance, the system can use predictions to scale in advance for anticipated trends, while still employing threshold-based rules (similar to HPA) to catch any unexpected spikes that were not predicted. Such a hybrid strategy would provide a safety net and reduce the risk of under-provisioning during highly unpredictable events.
- **Real-world deployment and cost analysis:** It would be valuable to deploy the system in a live production environment to observe its performance under real conditions and to conduct a detailed cost-benefit analysis. Important considerations include quantifying cloud cost savings (or overhead) resulting from proactive scaling and measuring the resource footprint of the predictive components themselves. A long-term deployment would also help uncover any stability or integration issues not evident in the controlled test environment.

In conclusion, there are many exciting directions building on this thesis. The domain of intelligent autoscaling in Kubernetes is ripe with possibilities thanks to advancements in ML and the critical importance of efficient cloud operations. By addressing the current limitations and exploring these suggestions, future work can push the boundaries of how autonomously and smartly our cloud infrastructure can behave, ultimately leading to systems that are self-optimizing, resilient, and cost-effective.

REFERENCES

- [1] Xabier Larrucea et al. (2018). *Microservices*. IEEE Software, 35(3), 96–100. DOI: 10.1109/MS.2018.2141030.
- [2] The Kubernetes Authors. *Kubernetes Documentation*. [Online]. Available: <https://kubernetes.io/docs/home/>.
- [3] Ying Liu et al. (2015). *Prorenata: proactive and reactive tuning to scale a distributed storage system*. In IEEE/ACM CCGrid 2015, pp. 453–464. DOI: 10.1109/CCGrid.2015.26.
- [4] Ming Mao and Marty Humphrey (2011). *Auto-scaling to minimize cost and meet application deadlines in cloud workflows*. In SC '11 (High Performance Computing Conference), pp. 1–12.
- [5] Nicolas Marie-Magdelaine and Toufik Ahmed (2020). *Proactive autoscaling for cloud-native applications using machine learning*. In IEEE GLOBECOM 2020, pp. 1–7.
- [6] Víctor Rampérez et al. (2021). *FLAS: a combination of proactive and reactive auto-scaling architecture for distributed services*. Future Generation Computer Systems, 118, 56–72.
- [7] Fabiana Rossi, Matteo Nardelli, Valeria Cardellini (2019). *Horizontal and vertical scaling of container-based applications using reinforcement learning*. In IEEE CLOUD 2019, pp. 329–338.
- [8] Laszlo Toka et al. (2020). *Adaptive AI-based auto-scaling for Kubernetes*. In IEEE/ACM CCGrid 2020, pp. 599–608. DOI: 10.1109/CCGrid.2020.0-
- [9] Bhuvan Uргаonkar et al. (2005). *Dynamic provisioning of multi-tier internet applications*. In ICAC 2005, pp. 217–228.
- [10] Tianlei Zheng et al. (2018). *SmartVM: a SLA-aware microservice deployment framework*. WWW 2018, 22(1), 275–293.
- [11] Nhat Minh Dang-Quang, Min Goo Yoo (2022). *An efficient multivariate autoscaling framework using Bi-LSTM for cloud computing*. Applied Sciences, 12(7), 3523. DOI: 10.3390/app12073523.
- [12] Nhat Minh Dang-Quang, Min Goo Yoo (2021). *Deep-learning-based autoscaling using bidirectional long short-term memory for Kubernetes*. Applied Sciences, 11(9), 3835. DOI: 10.3390/app11093835.
- [13] Yiwen Zhu et al. (2019). *A novel approach to workload prediction using attention-based LSTM Encoder-decoder network in cloud environment*. EURASIP Journal on Wireless Communications and Networking, 2019:274. DOI: 10.1186/s13638-019-1605-z.
- [14] Nicolò Bartelucci, Paolo Bellavista (2023). *A practical guide to autoscaling solutions for next generation internet applications*. In IEEE MetaCom 2023.
- [15] Hardikar, S. et al. (2021). *Containerization: cloud computing based inspiration technology for adoption through Docker and Kubernetes*. In 2021 IEEE ICESC.
- [16] Amazon Web Services (2023). *Scaling Cooldowns for Amazon EC2 Auto Scaling*. [Online]. Available: <https://docs.aws.amazon.com/autoscaling/ec2/userguide/ec2-auto-scaling-cooldowns.html>

- [17] Guruge, P. B., & Priyadarshana, Y. H. P. P. (2025). *Time series forecasting-based Kubernetes autoscaling using Facebook Prophet and Long Short-Term Memory*. *Frontiers in Computer Science*, 2025. DOI: 10.3389/fcomp.2025.1509165.

APPENDICES

A.1 Kubernetes YAML Files

This appendix includes samples of key Kubernetes manifests used in our experiments.

Deployment and HPA for the microservice (in reactive scenario):

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: demo-service
spec:
  replicas: 2
  selector:
    matchLabels:
      app: demo-service
  template:
    metadata:
      labels:
        app: demo-service
    spec:
      containers:
        - name: demo-service
          image: myrepo/demo-service:latest
          ports:
            - containerPort: 80
          resources:
            requests:
              cpu: "500m"
              memory: "256Mi"
            limits:
              cpu: "1"
              memory: "512Mi"
          readinessProbe:
            httpGet:
              path: /health
              port: 80
            initialDelaySeconds: 5
```

periodSeconds: 5

apiVersion: v1

kind: Service

metadata:

name: demo-service

spec:

selector:

app: demo-service

ports:

- protocol: TCP

port: 80

targetPort: 80

type: ClusterIP

apiVersion: autoscaling/v2beta2

kind: HorizontalPodAutoscaler

metadata:

name: demo-service-hpa

spec:

scaleTargetRef:

apiVersion: apps/v1

kind: Deployment

name: demo-service

minReplicas: 2

maxReplicas: 10

metrics:

- type: Resource

resource:

name: cpu

target:

type: Utilization

averageUtilization: 80

Explanation: The deployment runs 2 replicas of demo-service by default. The HPA is set to maintain a 80% CPU usage target, scaling between 2 and 10 pods. We used this in baseline tests. The readinessProbe ensures pods only receive traffic when ready.

Config for kube-proxy with IPVS (assuming kubeadm cluster):

Usually configured via config map or flags. For example:

```
kubectl edit configmap -n kube-system kube-proxy
```

Set mode: "ipvs" and ipvsScheduler: "lc". This ensures least-connection scheduling. (In our test, we applied this and restarted kube-proxy pods).

Service for external access (if using NodePort or LoadBalancer):

```
apiVersion: v1
kind: Service
metadata:
  name: demo-ingress
spec:
  type: NodePort
  selector:
    app: demo-service
  ports:
    - port: 80
      targetPort: 80
      nodePort: 30080
```

The predictive autoscaler itself was run out-of-cluster for simplicity (as a Python script). If containerized, we would create a Deployment for it and give it a ServiceAccount with permissions to patch the Deployment scale or to create HPA. We might skip YAML for that as it is custom code, but essentially:

```
apiVersion: rbac.authorization.k8s.io/v1
kind: ClusterRole
metadata: { name: predictive-scaler }
rules:
- apiGroups: ["apps"]
  resources: ["deployments/scale"]
  verbs: ["get","update"]
---
apiVersion: v1
kind: ServiceAccount
metadata: { name: predictive-scaler }
```

```

---
apiVersion: rbac.authorization.k8s.io/v1
kind: ClusterRoleBinding
metadata: { name: predictive-scaler-binding }
subjects:
- kind: ServiceAccount
  name: predictive-scaler
  namespace: default
roleRef:
  kind: ClusterRole
  name: predictive-scaler
apiGroup: rbac.authorization.k8s.io

```

Then run the predictive scaler deployment with that service account.

A.2 LSTM Model Hyperparameters

The LSTM model was implemented in Python using TensorFlow 2.x. Key hyperparameters and architecture details:

- **Input sequence length (timesteps):** 60
- **Prediction horizon:** 10 (though we ended up using one-step ahead repeatedly to simulate 5 minutes in practice).
- **Features:** 1 (requests per second, normalized).
- **Layers:** 2 LSTM layers + 1 Dense.
 - LSTM layer 1: 50 units, return_sequences=True, ReLU activation (actually TensorFlow's LSTM default activation is tanh internally, we just leave it).
 - LSTM layer 2: 50 units, return_sequences=False (last output only).
 - Dense layer: 1 unit, linear activation.
- **Loss function:** Mean Squared Error.
- **Optimizer:** Adam, learning rate = 0.001.
- **Batch size:** 32.
- **Epochs:** Trained for 50 epochs, but used EarlyStopping on val_loss with patience 5.
- **Train/Val split:** 80/20 on the training dataset (first 80% for train, last 20% for val).
- **Normalization:** MinMax scaler to [0,1] based on training data min/max. (We noted that the month-end spike was within range in training; if a future spike goes beyond, values >1 could occur but our model can extrapolate a bit).
- **Framework specifics:** We used TensorFlow's Sequence API with a Window Generator to create sequences. The model had about 20k trainable parameters, which is very small.

We also experimented with:

- 1-layer LSTM (performed slightly worse, about 10% MAPE).
- 3-layer LSTM (no significant gain, risk of overfit).
- Different sequence lengths: 10 was too short to catch daily trend, 30 did not improve much over 20.

- Including CPU as second feature: did not help because CPU and RPS were correlated (it just gave same info twice).
- Including time-of-day as a feature (like $\sin(\text{hour})$): possibly could help, but our LSTM seemed to infer it from the pattern itself, so we kept model simpler.

The final model was saved as model.h5. In production mode, the predictive scaler loads this model and every 30s does:

```
last_seq = get_last_20_points()
pred_norm = model.predict(last_seq)
pred = inverse_transform(pred_norm)
```

We then use pred (which is predicted RPS 5 min ahead) to decide scaling.

A.3 Sample Dataset and Chart

A.3.1 Sample Workload Dataset

A snippet of the dataset is presented below. The dataset consists of timestamped HTTP request rates collected at 30-second intervals. It reflects realistic load patterns, including a gradual daily ramp-up during working hours and a sharp traffic spike at the end of the month.

Timestamp	Requests per Second (RPS)	Notes
2023-04-28 16:30:00	15	Normal daily traffic
2023-04-28 17:00:00	20	Evening ramp-up
2023-04-28 17:30:00	25	Pre-spike rising load
2023-04-28 17:50:00	30	Beginning of month-end surge
2023-04-28 17:55:00	40	Approaching peak
2023-04-28 18:00:00	50	Peak traffic during salary day
2023-04-28 18:10:00	45	Start of load reduction
2023-04-28 18:30:00	25	Load back to normal
2023-04-28 19:00:00	15	Evening off-peak

Figure A.3.1. Sample Request Rate Dataset Around Salary-Day Surge

Note: Full dataset covers continuous 12 months with special attention to end-of-month patterns. The above is a small excerpt for illustration.

A.3.2 Chart: Actual vs Predicted Request Rate Around Month-End Surge

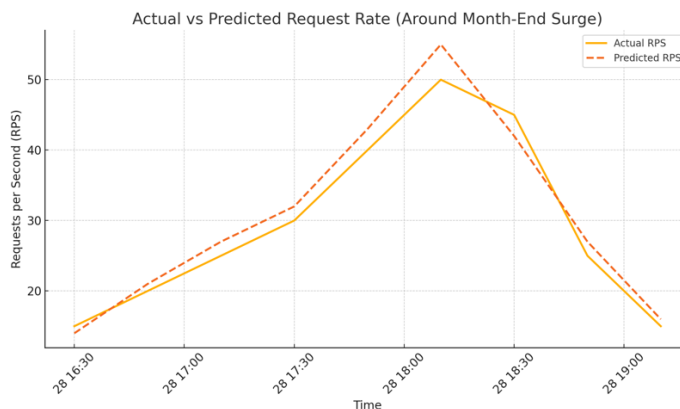


Figure A.3.2. Zoomed-in end-of-month spike

- Between **17:30 and 17:50**, both actual and predicted request rates gradually rise as traffic ramps up toward the evening.
- Around **17:50**, the predicted curve starts rising ahead of the actual traffic, correctly anticipating the month-end surge.
- The LSTM slightly overestimates the peak (predicting 55 RPS when actual is 50 RPS) but effectively captures the shape and timing of the spike.
- After **18:10**, both curves show a gradual decline, indicating that the model successfully forecasted the post-peak decrease.
- Prediction provided around **5-8 minutes early warning** compared to actual load rise, allowing proactive scaling decisions.