VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135













Journal Website: http://sciencebring.co m/index.php/iiasr

Copyright: Original content from this work may be used under the terms of the creative commons attributes 4.0 licence.



Adaptive Cloud Orchestration: Mitigating Cold-Start Latency and Optimizing Cost-Performance Trade-offs in Kubernetes via Reinforcement Learning and Ansible Integration

Submission Date: October 25, 2025, Accepted Date: November 10, 2025,

Published Date: November 27, 2025

Dr. A. Vance and J. Sterling

Department of Advanced Computing Systems, Institute of Cloud Architecture

ABSTRACT

The rapid proliferation of microservices architectures has established Kubernetes as the de facto standard for container orchestration. However, efficient auto-scaling remains a persistent challenge, particularly when balancing strict Service Level Agreements (SLAs) against the operational expenditures of cloud infrastructure. Traditional rule-based scaling mechanisms, such as the Horizontal Pod Autoscaler (HPA), often exhibit reactive latency, leading to "cold-start" delays during traffic bursts and resource overprovisioning during idle periods. This paper proposes a novel, hybrid orchestration framework that integrates Deep Reinforcement Learning (DRL) with Ansible-based configuration management to optimize the dynamic scaling of Azure PaaS environments. By treating the scaling problem as a Markov Decision Process (MDP), we develop a Q-learning agent capable of predicting workload fluctuations and preemptively provisioning resources. Furthermore, we leverage Ansible playbooks to parallelize node initialization, significantly reducing the initialization time of transient Virtual Machines (VMs). Our experimental results demonstrate that this approach reduces cold-start latency by approximately 40% and operational costs by 22% compared to standard reactive scaling methods, offering a robust solution for mixed interactive and batch workloads in enterprise environments.

KEYWORDS

Reinforcement Learning, Cloud Auto-scaling, Ansible, Cold-Start Latency, Cost Optimization, Container Orchestration.

INTRODUCTION

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











The The paradigm shift from monolithic software architectures to cloud-native microservices has fundamentally altered the landscape of application deployment and management. As organizations migrate critical workloads to the cloud, the ability to dynamically adjust resources in response to fluctuating demand—elasticity—has become a primary determinant of system reliability and costefficiency [1]. Kubernetes, originally developed by Google and now maintained by the Cloud Native Computing Foundation, has emerged as the dominant platform for automating the deployment, and management of containerized applications [4], [10]. Its declarative API and ecosystem allow for sophisticated robust orchestration strategies that were previously unattainable with traditional virtualization technologies.

However, despite the maturity of Kubernetes, the challenge of effective auto-scaling persists. The standard Horizontal Pod Autoscaler (HPA) typically relies on reactive metrics, such as CPU and memory utilization thresholds, to trigger scaling events. While effective for gradual traffic increases, this reactive approach is inherently flawed when dealing with sudden, high-velocity traffic spikes [30]. The delay between the detection of a metric threshold breach and the availability of a new ready-state container—often referred to as "coldstart latency"—can result in transient service unavailability, increased tail latency, and violations of Service Level Agreements (SLAs). Conversely, to mitigate these risks, operators often resort to overprovisioning, maintaining a buffer of idle resources that inflates operational costs without contributing to productive output [9].

The complexity of this problem is compounded in mixed-workload environments, where interactive user requests compete for resources with background batch processing tasks. In such scenarios, the scaling logic must distinguish latency-sensitive between and throughputsensitive workloads, a nuance often lost in simple threshold-based rules. Furthermore. the underlying infrastructure provisioning specifically, the time required to spin up new Virtual Machine (VM) nodes in a cloud environment—adds a significant bottleneck. While containers themselves start in seconds, the underlying infrastructure in a Platform-as-a-Service (PaaS) or Infrastructure-as-a-Service (IaaS) context may take minutes to become operational [25].

Recent advancements in Artificial Intelligence (AI) and machine learning offer promising avenues for these limitations. Reinforcement addressing Learning (RL), in particular, is well-suited for dynamic resource management problems where an agent learns optimal policies through interaction with an environment [26], [27]. By modeling the cloud cluster as a stochastic environment, an RL agent can learn to anticipate workload patterns and execute preemptive scaling actions, thereby smoothing out the volatility of demand. Simultaneously, the use of robust configuration management tools like Ansible has proven effective in streamlining the "end-to-end" provisioning process, reducing the time required to bring a new node from a raw state to a clusterready state [25].

This paper presents a comprehensive study on integrating these technologies. We propose a framework that combines a Q-learning-based

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











scheduler with Ansible automation to manage Kubernetes clusters on Azure PaaS. Our approach addresses two critical gaps in the current literature: the lack of predictive capability in standard Kubernetes scalers and the inefficiency of node initialization in dynamic cloud environments. By synthesizing insights from recent studies on Ansible-based scaling [25] and machine learningdriven orchestration [8], [30], we aim to demonstrate a scalable, cost-effective architecture for modern cloud applications.

The remainder of this paper is organized as follows: Section 2 details the methodology, including the mathematical formulation of the RL model and the system architecture. Section 3 presents the experimental results, comparing our proposed model against baseline methods. Section 4 discusses the implications, limitations, and potential future directions of this research.

2. Methodology

To address the dual challenges of latency minimization and cost optimization, we designed a closed-loop control system that overlays the standard Kubernetes control plane. This system consists of three primary components: the Monitoring Agent, the Reinforcement Learning (RL) Decision Engine, and the Ansible Execution Module.

2.1 System Architecture and Environment

The experimental environment was constructed using Azure Kubernetes Service (AKS), leveraging a hybrid node pool configuration. The primary node pool consisted of reserved instances to handle baseline loads, while a secondary pool of burstable, transient VMs (similar to AWS Spot Instances) was utilized to handle peak traffic. This setup mimics a real-world enterprise scenario where cost efficiency is prioritized for variable workloads [9], [25].

The Monitoring Agent is implemented as a Prometheus sidecar that aggregates cluster metrics at a 15-second granularity. Unlike standard HPA which primarily monitors CPU utilization, our agent captures a high-dimensional state vector including request arrival rates, queue depth, pod startup latency, and custom application-level metrics (e.g., transaction processing time). This data is fed into the RL Decision Engine.

The Ansible Execution Module serves as the actuator for the system. Upon receiving a scaling directive from the RL engine, it triggers specific Ansible playbooks designed to parallelize the node provisioning process. As noted by Donthi (2025), utilizing Ansible for end-to-end dynamic scaling allows for the pre-configuration of network interfaces and storage mounts, significantly reducing the "ready-time" of new nodes compared to standard cloud-init scripts [25].

2.2 Algorithmic **Formulation** of the Reinforcement Learning Scheduler

The core of our proposed solution is the decisionmaking engine, which we modeled as a Markov Decision Process (MDP). This allows the system to learn a policy \$\pi\$ that maps observed states to optimal scaling actions.

2.2.1 State Space (\$S\$)

The state of the cluster at time step \$t\$ is defined by a vector \$s_t \in S\$, representing the current resource utilization and workload characteristics. To capture the nuance of mixed workloads, we

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











expanded the state representation beyond simple CPU metrics. The state vector is defined as:

 $s_t = (u_{cpu}, u_{mem}, q_{len}, r_{in}, n_{pods}),$ n {nodes})\$\$

Where:

- \$u_{cpu}, u_{mem}\$: The aggregate CPU and memory utilization of the cluster (normalized between 0 and 1).
- \$q_{len}\$: The current length of the request queue, serving as a proxy for immediate saturation.
- \$r_{in}\$: The rate of incoming requests (requests per second), derived from the ingress controller.
- \$n_{pods}\$: The current number of active pods.
- \$n_{nodes}\$: The current number of active worker nodes.

This multi-dimensional state space allows the agent to distinguish between high-utilization scenarios caused by efficient processing versus those caused by bottlenecks [8].

2.2.2 Action Space (\$A\$)

The action space \$A\$ represents the possible control decisions available to the agent. To maintain stability and prevent oscillation (flapping), the action space is discrete:

\$\$A \{scale_out(k), $scale_in(k)$, $do\nothing\$

Where \$k\$ represents the magnitude of the scaling action (e.g., add 1 node, add 2 nodes). In our implementation, we limited \$k\$ to small integers (\$1, 2\$) to prevent aggressive over-scaling, relying on the high frequency of decision steps (every 30 seconds) to accumulate necessary resources.

2.2.3 Reward Function (\$R\$)

The reward function is the critical component that guides the agent's learning toward the desired trade-off between performance and cost. We formulated a composite reward function \$R(s_t, a t)\$ that penalizes both SLA violations and excessive resource usage.

 $R_t = \alpha \cdot R_{perf} + \beta \cdot Cdot$ $R_{cost} + \gamma \alpha \gamma Cdot R_{stable}$

Where:

\$R_{perf}\$: A negative reward (penalty) proportional to the number of failed requests or latency violations. If the 95th percentile latency \$L_{p95}\$ exceeds the SLA threshold \$L_{max}\$, the penalty increases exponentially:

$$R_{perf} = -e^{(L_{p95} - L_{max})}$$

\$R_{cost}\$: A negative reward proportional to the cost of active resources. This incentivizes the agent to scale down when resources are underutilized.

$$R_{\cos t} = -(C_{node} \times n_{nodes})$$

- \$R_{stable}\$: A penalty for frequent changes in the number of nodes, preventing "flapping" where the system adds and removes nodes rapidly, which incurs setup costs and instability [30].
- \$\alpha, \beta, \gamma\$: Weighting coefficients that determine the priority of the system. For this study, we prioritized performance (\$\alpha=0.6, \beta=0.3, \gamma=0.1\$).

2.2.4 Q-Learning Implementation

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











We employed a Q-learning algorithm, a model-free reinforcement learning technique, to estimate the value of taking action \$a\$ in state \$s\$. The O-value is updated iteratively using the Bellman equation:

 $\Q_{new}(s_t, a_t) \leq (1 - \epsilon)$ $Q(s_t, a_t) + \beta (cdot R_{t+1} + delta cdot$ $\max_{a} Q(s_{t+1}, a)]$ \$

Here, \$\eta\$ is the learning rate (set to 0.1), and \$\delta\$ is the discount factor (set to 0.95), emphasizing long-term rewards. We utilized an \$\epsilon\$-greedy strategy for exploration, where the agent chooses a random action with probability \$\epsilon\$ (initially 1.0, decaying to 0.05) to explore the state space, and the action with the highest Q-value otherwise [26].

2.3 Enhancing Responsiveness with Ansible

While the RL agent decides when and how much to scale, the speed of execution is governed by the infrastructure automation. Standard Azure Virtual Machine Scale Sets (VMSS) often rely on cloud-init scripts which execute sequentially at boot time. This introduces a "cold start" delay that can range from 3 to 5 minutes [25].

To mitigate this, we integrated Ansible Tower into the provisioning pipeline. Upon creating a new VM, the system triggers an Ansible callback. The Ansible control node then executes pre-compiled playbooks that perform necessary configurations (installing container runtimes, mounting volumes, configuring network routes) in parallel. This approach, as detailed in recent industry case studies, ensures that the node joins the Kubernetes cluster in a "Ready" state significantly faster than traditional methods [25], [28].

2.4 Workload Characterization and Traffic Generation

To train and validate the model, we required a realistic workload trace. We utilized the IMAPS email traffic dataset characteristics [29] to simulate bursty, interactive workload, superimposed with periodic batch processing jobs (e.g., scientific data analysis) [1].

The traffic generator was configured to produce:

- Diurnal Patterns: Sinusoidal traffic mimicking daily user activity.
- Poisson Bursts: Random spikes in request rates to test the system's responsiveness to unpredictable events.
- 3. Batch Jobs: Large, resource-intensive pods scheduled at irregular intervals, requiring the system to provision substantial capacity quickly.

This mixed workload is crucial for evaluating the "Federated" nature of modern cloud applications, where different types of jobs share the same underlying infrastructure [3], [27].

3. Results

The proposed framework (RL-Ansible) was evaluated against two baselines:

- Static Provisioning: A fixed cluster size calculated to handle peak load (the "safe" but expensive option).
- Standard K8s HPA: The default horizontal 2. pod autoscaler using CPU thresholds (target 70%).

3.1 Latency Reduction and SLA Compliance

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











The primary metric for performance was the 95th percentile tail latency (\$L_{p95}\$) of HTTP requests during burst periods.

Under the Standard HPA, we observed significant latency spikes during the onset of traffic bursts. The HPA's reactive nature meant that new pods were only requested after the CPU threshold was breached. Combined with the infrastructure spinup time, this resulted in a "saturation window" lasting approximately 4-6 minutes where the \$L {p95}\$ exceeded 2 seconds, violating the 500ms SLA.

In contrast, the RL-Ansible agent demonstrated predictive capability. By learning the temporal correlations in the traffic (e.g., a small increase in queue depth often precedes a spike), the agent initiated scaling actions before the saturation point. Furthermore, the Ansible-optimized provisioning reduced the node "ready time" from an average of 240 seconds to 145 seconds. Consequently, the \$L_{p95}\$ during bursts remained below 600ms for 92% of the experiment duration, a dramatic improvement over the HPA's 65% compliance rate.

3.2 Cost Efficiency Analysis

Cost efficiency was measured by calculating the total "resource-hours" billed by the cloud provider, adjusted for the utilization rate.

- Static Provisioning achieved 100% SLA compliance but incurred the highest cost, as resources remained idle during off-peak hours.
- Standard HPA reduced costs by roughly 30% compared to static provisioning but suffered from "thrashing"—frequently spinning up nodes that were only used for a few minutes before being

terminated, which is inefficient due to the minimum billing increments of many cloud providers.

RL-Ansible achieved a cost reduction of 22% compared to the Standard HPA (and nearly 50% compared to Static). This was achieved not just by scaling down, but by scaling smartly. The RL agent learned to tolerate brief, non-critical queue buildups rather than immediately provisioning expensive nodes for transient spikes, optimizing the integral of the cost function over time [9], [26].

3.3 Impact of Ansible on Cold-Start

We isolated the impact of the Ansible integration by running the RL agent with standard cloud-init scripts versus Ansible playbooks. The data showed that the Ansible-based approach reduced the variability (standard deviation) of node startup times. While cloud-init scripts occasionally hung due to repository timeouts or package dependencies, the Ansible playbooks, utilizing predownloaded artifacts and local mirrors, provided a consistent startup trajectory. This determinism is crucial for the RL agent, as it reduces the stochasticity of the environment, allowing the Qlearning algorithm to converge faster.

Expanded Analysis: Deep Dive into RL State Dynamics and Ansible Workflow

To fully appreciate the robustness of the proposed solution, it is necessary to expand on two critical areas: the nuanced behavior of the Reinforcement Learning agent during the "exploration vs. exploitation" phases, and the granular mechanics of the Ansible workflow that enables such rapid scaling. This section provides a detailed breakdown of these components, bridging the gap

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











between theoretical algorithmic design and practical infrastructure implementation.

3.3.1 **Dynamic Transition** State and **Convergence Analysis**

A common criticism of Reinforcement Learning in production systems is the "warm-up" period required for the agent to learn an effective policy. During the initial training phase, the agent operates with a high \$\epsilon\$ (exploration rate), taking random actions to map the state-action space. In a live production environment, this could be disastrous. To mitigate this, we employed a "shadow mode" training technique, utilizing historical trace data from the PhenoMeNal project [7] and synthesized IMAPS traffic [29] to pre-train the model before live deployment.

The analysis of the Q-table convergence reveals interesting behaviors regarding the agent's risk tolerance. Initially, the agent exhibited a high sensitivity to \$q_{len}\$ (queue length). Any nonzero queue depth triggered a scale-out action. However, as the training progressed, the agent "learned" the difference between a transient queue spike (lasting <30 seconds) and a sustained traffic surge. This is evident in the weight modifications within the Q-matrix.

Let us consider the state transition for a "moderate load" scenario, defined as CPU utilization \$u {cpu} \approx 60\%\$ and queue length \$q_{len} \approx 50\$.

In early epochs (0-100), the Q-value for a = 0 $scale_out(1)$ \$ was significantly higher than \$a = do\ nothing\$:

 $\$Q(s_{mod}),$ scale_out) $Q(s_{mod}),$ \gg do\ nothing)\$\$

This reflects a naive, reactive policy similar to HPA.

However, by epoch 5000, the relationship inverted for this specific state, provided that the rate of incoming requests \$r {in}\$ was stable. The agent learned that the cost penalty \$R {cost}\$ of adding a node outweighed the negligible performance gain \$R {perf}\$ for such a small queue, effectively learning to "absorb" small bursts.

This behavior is crucial for cost optimization in cloud environments utilizing spot instances or transient VMs [9]. Cloud providers often charge a minimum billing period (e.g., 60 seconds). Rapidly scaling up and down (thrashing) for micro-bursts accumulates costs without delivering value. The RL unlike static thresholds. inherently internalizes this cost structure through the reward function \$R_{cost}\$.

Furthermore, we observed that the agent developed a "pre-fetching" behavior during diurnal upticks. By recognizing the specific state vector associated with the onset of the "morning rush" (increasing \$r {in}\$ combined with specific timeof-day encoding), the agent began scaling out utilization metrics saturated. anticipatory action is the primary driver behind the 40% reduction in cold-start latency observed in the main results. It transforms the scaling paradigm from "reaction to failure" to "preparation for demand."

3.3.2 The Ansible "Fast-Path" Architecture

While the RL agent provides the intelligence, the Ansible integration provides the velocity. The standard Kubernetes node bootstrapping process involves several sequential steps:

VM Provisioning (Cloud Provider) 1.

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











- 2. OS Boot & Network Initialization
- 3. **Kubelet Installation & Configuration**
- 4. Container Runtime Interface (CRI) Setup (e.g., Docker/containerd)
- 5. Joining the Cluster
- 6. **Image Pulling**

In a standard setup, steps 3 through 6 are often handled by cloud-init or startup scripts that fetch packages from public repositories at runtime. This introduces two failure points: network latency (downloading packages) and dependency resolution errors.

Our Ansible approach replaces steps 3 through 6 with a deterministic, local-execution model. The detailed workflow is as follows:

- Golden Image Utilization: We utilize a "Golden Image" for the VMSS that has the heavy artifacts (Docker images, Kubelet binaries) prebaked. However, configuration must still happen dynamically to connect to the specific cluster.
- The Ansible Callback: As soon as the Azure VM has an IP address, a webhook triggers the Ansible Tower (or AWX) controller.
- Parallel Execution: Ansible executes a specialized playbook that performs the following actions in parallel (using async and poll directives):
- Task A: Updates the kubelet configuration with the specific dynamic tokens for the current cluster state.
- Task B: Mounts the high-performance storage drivers required for the data analysis workloads [1].

0 Task C: Pre-warms the node by loading common container images from the local disk into the CRI cache, eliminating the network overhead of docker pull for base layers.

The specific impact of this "Fast-Path" architecture was quantified by measuring the "Time-to-Ready" (TTR) metric.

 $$TTR = T \{ready\} - T \{request\}$

Standard Azure provisioning yielded a TTR with a mean of 240s and a standard deviation (\$\sigma\$) of 45s. The high variance was largely attributed to external repository speeds during package installation.

The Ansible-enhanced provisioning yielded a TTR with a mean of 145s and a \$\sigma\$ of only 8s. The low standard deviation is critical for the RL agent. In Control Theory, a controller acts most effectively when the actuator's lag is predictable. By stabilizing the TTR, Ansible allows the RL agent to plan more accurately, knowing that a scale out action will yield capacity in exactly ~2.5 minutes, rather than a variable window of 3-6 minutes.

3.3.3 Handling Heterogeneity in Federated Clusters

Modern cloud architectures often span multiple clusters or "Federations" [3], [27]. While our primary experiment focused on a single cluster, the logic extends to federated environments. In a federated setup, the Action Space \$A\$ of the RL agent expands to include a target cluster variable:

\$\$A_{fed} = \{scale_out(k, cluster_i), ...\}\$\$

Our expanded simulations suggest that the RL-Ansible approach is particularly beneficial here. Different clusters (e.g., On-prem vs. Azure) have

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











different cost profiles and TTR characteristics. A standard HPA cannot account for the fact that an on-prem node might be "free" but slow to boot, while an Azure node is "expensive" but fast. The RL agent, however, learns these latent variables.

In our simulation of a hybrid burst scenario, the agent prioritized utilizing on-prem resources for predictable base loads (minimizing \$R_{cost}\$) but aggressively switched to Azure scaling actions when the rate of change in demand (\$\Delta r_{in}\$) exceeded a certain threshold, recognizing that the "cost" of the slow on-prem boot time (SLA violation penalty) outweighed the financial cost of the cloud resources. This intelligent arbitration is a significant advancement over static priority rules typically used in federation controllers [3].

4. Discussion

The transition to cloud-native architectures offers immense potential for scalability, but this potential is often bottlenecked by the rigidity of traditional management tools. This study has demonstrated that the convergence of cognitive techniques (Reinforcement Learning) and modern automation (Ansible) creates a robust framework for nextgeneration cloud orchestration.

4.1 Interpretation of Findings

The superior performance of the RL-Ansible model validates the hypothesis that "awareness" of state history is essential for effective scaling. By encoding the queue depth and request rate trends into the state space, the agent transcends the limitations of reactive CPU thresholds. Furthermore, the cost savings (22%) highlight the economic viability of this approach. In large-scale enterprise clusters, where monthly bills can reach hundreds of thousands of dollars, such a

percentage represents a significant absolute saving, justifying the complexity of implementing an ML-based controller.

4.2 Limitations

Despite the promising results, several limitations must be acknowledged. First, the training time for the Q-learning agent is non-trivial. While we used historical traces for pre-training, a "cold" agent deployed into a new environment would perform poorly for several days until it explores the state space sufficiently. This limits the applicability of the solution for short-lived, ephemeral projects.

Second, the complexity of the Ansible integration introduces a maintenance overhead. The "Golden Images" and playbooks must be kept in sync with Kubernetes version updates. A mismatch between the Ansible configuration and the cluster version can lead to nodes failing to join, a failure mode that is less common in managed, provider-native scaling solutions.

Third, the reward function tuning is manual and sensitive. The balance between \$\alpha\$ (performance) and \$\beta\$ (cost) is subjective and business-dependent. An incorrectly weighted reward function could lead to massive cost overruns or severe service degradation.

4.3 Future Directions

Future research will focus on two main areas. First, we aim to explore Federated Reinforcement Learning, where agents in different clusters share "lessons learned" (gradients or policy updates) without sharing raw data. This would allow a new cluster to bootstrap its knowledge from an existing mature cluster, mitigating the cold-start training problem [27].

48

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











Second, we intend to investigate Deep Q-Networks (DQN) or Proximal Policy Optimization (PPO) to handle continuous action spaces. This would allow the agent to request precise amounts of CPU/RAM (vertical scaling) rather than just integer node counts (horizontal scaling), offering an even finer granularity of control [2], [8].

Finally, integrating predictive auto-scaling with energy-aware scheduling is a critical frontier. As data centers consume an increasing portion of global energy, adding carbon footprint metrics to the Reward Function could allow the agent to optimize not just for cost and speed, but for environmental sustainability.

Conclusion

As Kubernetes cements its place as the operating system of the cloud, the tools used to manage it must evolve. This paper presented a synergistic approach combining Reinforcement Learning and Ansible to solve the persistent problems of coldstart latency and cost inefficiency. By moving from static rules to dynamic, learned policies, and from sequential scripts to parallelized automation, we can achieve a cloud infrastructure that is not only elastic but truly adaptive. The results suggest that for complex, mixed-workload environments, such intelligent orchestration systems will soon become not just an optimization, but a necessity.

References

1. Sai Nikhil Donthi. (2025). Ansible-Based End-To-End Dynamic Scaling on Azure Paas for Refinery Turnarounds: Cold-Start Latency and Cost-Performance Trade-Offs. Frontiers in **Emerging Computer Science and Information** Technology, 2(11), 01-17.

- **2.** K. Hightower, B. Burns, and J. Beda. Kubernetes: Up and Running Dive into the Future of Infrastructure. O'Reilly Media, Inc., 1st edition, 2017.
- 3. S. Horovitz and Y. Arian. Efficient cloud autoscaling with sla objective using glearning. In 2018 IEEE 6th International Conference on Future Internet of Things and Cloud (FiCloud). pages 85-92. IEEE, 2018.
- 4. J. Huang, C. Xiao, and W. Wu. Rlsk: A job scheduler for federated kubernetes clusters based on reinforcement learning. In 2020 IEEE International Conference on Cloud Engineering (IC2E), pages 116–123. IEEE, 2020.
- 5. J. Humble and D. Farley. Continuous Delivery: Reliable Software Releases Through Build, Test, and Deployment Automation. Addison-Wesley Professional, 2010.
- and 6. M. Karamollahi C. Williamson. Characterization of IMAPS email traffic. In 27th IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems, MASCOTS 2019, Rennes, France, October 21-25, 2019, pages 214–220. IEEE Computer Society, 2019.
- 7. T. Lorido-Botran, J. Miguel-Alonso, and J. A. Lozano. A review of auto-scaling techniques for elastic applications in cloud environments. Journal of grid computing, 12(4):559-592, 2014.
- 8. OMG. Business Process Model and Notation (BPMN), 2.0. Version http://www.omg.org/spec/BPMN/2.0, January 2011.
- 9. Tesliuk, S. Bobkov, V. Ilyin, A. Novikov, A. Poyda, and V. Velikhov, "Kubernetes container orchestration as a framework for flexible and effective scientific data analysis," Kubernetes

VOLUME 05 ISSUE 11 Pages: 40-50

OCLC - 1368736135











- Container Orchestration as a Framework for Flexible and Effective Scientific Data Analysis, pp. 67–71, Dec. 2019.
- 10.F. Antonescu, P. Robinson, and T. Braun, "Dynamic topology orchestration for Cloud-Based distributed applications," Dynamic Topology Orchestration Distributed Cloud-Based Applications, pp. 116-123, Dec. 2012.
- 11.D. Kim, H. Muhammad, E. Kim, S. Helal, and C. Lee, "TOSCA-Based and Federation-Aware cloud orchestration for Kubernetes container platform," Applied Sciences, vol. 9, no. 1, p. 191, Jan. 2019.
- **12.** E. A. Brewer, "Kubernetes and the path to Cloud Native," Kubernetes and the Path to Cloud Native, Aug. 2015.
- 13.Q. Lei, W. Liao, Y. Jiang, M. Yang, and H. Li, "Performance and Scalability Testing Strategy based on KubeMark," Performance Scalability Testing Strategy Based on Kubemark, Apr. 2019.

- 14.Q. Li, G. Yin, T. Wang, and Y. Yu, "Building a Cloud-Ready program," Building a Cloud-Ready Program, pp. 159–164, Jun. 2018.
- 15.K. Peters et al., "PhenoMeNal: processing and analysis of metabolomics data in the cloud," GigaScience, vol. 8, no. 2, Dec. 2018.
- 16.L. Toka, G. Dobreff, B. Fodor, and B. Sonkoly, "Machine Learning-Based scaling Management Kubernetes clusters," edge IEEE Service Transactions on Network and Management, vol. 18, no. 1, pp. 958–972, Mar. 2021.
- 17.P. Ambati and D. Irwin, "Optimizing the cost of executing mixed interactive and batch workloads on transient VMs," Optimizing the Cost of Executing Mixed Interactive and Batch Workloads on Transient VMs, pp. 45-46, Jun. 2019.
- 18.B. Burns, B. Grant, D. Oppenheimer, E. Brewer, and J. Wilkes, "Borg, Omega, and Kubernetes," Communications of the ACM, vol. 59, no. 5, pp. 50-57, Apr. 2016.