

When Less is More: Core-Restricted Container Provisioning for Serverless Computing

Gaetano Somma^{*}, Constantine Ayimba^{‡§}, Paolo Casari[‡], Simon Pietro Romano^{*}, Vincenzo Mancuso[‡]

^{*}University of Naples “Federico II,” Italy

[§]University Carlos III of Madrid, Spain

[‡]IMDEA Networks Institute, Madrid, Spain

[‡]DISI – University of Trento, Italy

Abstract—Cloud applications are exposed to workloads whose intensity can change unpredictably over time. Hence, the ability to quickly scale the amount of computing resources provisioned to applications is essential to minimize costs while providing reliable services. In this context, *containers* are deemed to be a promising technology to enable fast elasticity in resource allocation schemes.

In this paper, we propose and experimentally test an efficient container-based cloud computing provisioning system. First, we address the container deployment problem and discuss how to manage container provisioning and scaling. Second, we devise a resource management mechanism leveraging on both admission control and auto-scaling techniques. We propose to drive auto-scaling decisions through a Q-Learning algorithm, which is agnostic to the specific computing environment, and proceeds based only on the load of the physical processors assigned to a container. We evaluate our solution in two experimental setups, and show that it yields significant advantages when compared to popular container managers such as Kubernetes.

Index Terms—Autoscaling, Provisioning, Q-Learning, Container, Docker, Kubernetes

I. INTRODUCTION

The adoption of containers is growing at a rapid pace. This is mainly due to their comparatively simpler management with respect to virtual machines (VMs), and to the efficient resource utilization they enable. Moreover, the performance of containers is comparable to that of native computing platforms in terms of throughput and CPU utilization whereas VMs typically impose non-negligible overhead [1], [2]. Characteristics such as quick deployment, short boot time [2], easy network management and the use of layered, small-sized container images have encouraged container adoption in global systems like the Google Cloud Platform and Amazon Web Services (AWS) [3].

The features mentioned above make it possible to achieve rapid elasticity (the ability to quickly scale the amount of allocated resources according to workload intensity) using containers. Managing container deployments in cloud environments is still an open issue that mainly involves provisioning and scaling strategies.

The performance of a container depends not only on how many CPU threads it uses but also on which core of the CPU these threads are executed. It is also contingent on the level of contention by other applications on the same CPU. In order to make the cost calculation transparent to the user, such inter-dependencies must be avoided and the relationship between the quality of service delivered and the amount of used resources

needs to be clarified. Existing container provisioning platforms delegate the CPU scheduling to the operating system. Given that current operating systems are not hyper-threading aware, interference among running containers cannot be completely avoided, and load is typically distributed sub-optimally across CPU threads. This results in unpredictable service times.

Containers are currently managed using software platforms such as HPA by Kubernetes or Docker Swarm, which offer reliability by default [4]. However, some limitations exist regarding performance guarantees and adaptability to rapid changes in the operating environment. For instance, currently adopted solutions favor over-provisioning policies over rapid elasticity involving system adaptation. Capacity allocations are statically sized to serve peak loads, so resources remain underutilized most of the time. We however argue that over-provisioning is neither efficient nor strictly needed. We show that automatic scaling is a better option provided that containers can be mapped onto hardware resources to avoid resource access conflicts.

The objective of our work is two-fold: (*i*) to design mechanisms that make access to computing resources simple and effective for container provisioning engines and, (*ii*) to validate an adaptive scaling strategy based on reinforcement learning, which optimizes throughput and costs without sacrificing the application’s response time.

We implement *core pinning* to ensure that a CPU core is reserved for a container, thereby forestalling contention side-effects due to hyper-threading. This approach also simplifies pricing models by making the cost of a container proportional to that of a CPU core. With regard to provisioning, we show that scaling the number of containers allotted to an application yields better performance than scaling the amount of resources allotted to a single container. Scaling the number of containers makes service time predictable, and thus provides the technical basis for, e.g., the stipulation of *Service Level Agreements* (SLAs) between service providers and customers.

Moreover, we implement an automatic scaling system that predicts the required amount of resources and proactively makes scaling decisions in order to maximize the application workload processing throughput and minimize the infrastructure allocation costs. Our automatic scaling subsystem is a Q-Learning agent. Since it is based on a model-free reinforcement learning technique, this agent can learn the operating environment autonomously and adapt its scaling

policies without manual intervention.

The rest of the paper is structured as follows. In Section II, we discuss the challenges involved in container provisioning and propose our self-scaling provisioning solution in Section III. We describe our experimental testbed in Section IV. We present and analyze our experimental results in Section V. In Section VI, we review the literature on container provisioning and draw final conclusions in Section VII.

II. CONTAINER PROVISIONING

In this section we explain how to avoid interference (i.e., contention in accessing shared computing resources) among running containers and explain how we operate to make service time predictable.

A Linux container is a group of isolated processes running on the host machine without any resource virtualization. A container can be granted an arbitrary amount of resources on the host machine: the amount of actually available resources depends both on the host capacity and on the resources allotted to other containers. For this reason, containers running on the same host will interfere with each other.

The use of hyper-threaded CPUs results in additional interference. Each CPU comprises multiple cores, each of which can run two threads. These threads share the hardware for the execution phase. Hyper-threading leads to a performance improvement for each core since it minimizes the impact of cache-miss interruptions. Unfortunately, such an architecture may also yield unpredictable performance, depending on which thread is used to run a process [5], [6]. In fact, as different threads in the same core share part of the architecture [7], execution performance is affected by other processes in the same core.

The solution to these problems is the use of resource limitation in conjunction with core pinning, as demonstrated in [8]. Concretely, this means that we dedicate one or more (entire) cores of a CPU to a given container for its exclusive use. This configuration curtails any interference with other processes and eliminates the interference related to L1 caching mechanisms since each process stably runs on the same core. CPU core pinning also leads to efficient resource isolation [9], improved throughput and improved power utilization [10].

In Linux, both resource limitation and core pinning can be achieved by leveraging the `cgroup` feature. Using this feature requires the specification of the threads to be used for each container. To avoid disparity in hyper-threaded architectures, we select threads belonging to the same core. For the Linux distribution we use in our test-bed, this involves consulting the `cpuinfo` file.

Core pinning also clearly delineates the number of resources used as containers are mapped onto a known number of allocated cores. Therefore, following [11], the cost of running a container over k time intervals can be computed as:

$$\text{Cost}(k) = \alpha \sum_{n=1}^k C_{n,n-1} \cdot (t_n - t_{n-1}), \quad (1)$$

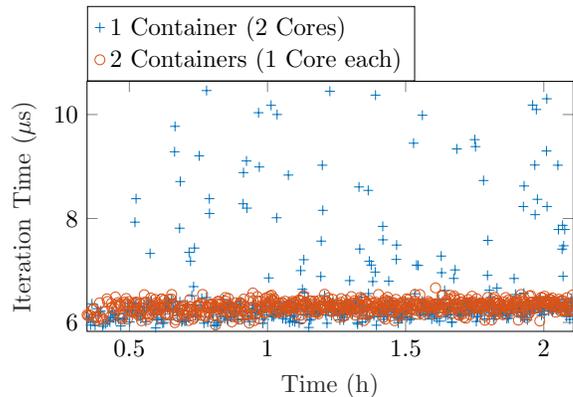


Fig. 1. Vertical vs. horizontal scaling: two containers running on distinct cores provide more predictable performance than one container running on two cores.

where $C_{n,n-1}$ is the number of cores dedicated to the container during the interval $[t_{n-1}, t_n]$, and α is a configurable cost scaling parameter chosen by the Cloud provider.

As regards scaling, two strategies are possible: vertical scaling and horizontal scaling. The former entails the addition or removal of cores from a running container while the latter involves instantiating new containers or decommissioning active ones without adjusting their resource allocations. If vertical scaling is used such that multiple cores are assigned to a container, even with CPU core pinning on hyper-threaded cores, Linux’s completely fair scheduler (CFS) might assign all processes to some cores leaving the rest in an idle state. This is because these schedulers are hyper-threading unaware and may introduce interference among competing processes on the same container. Horizontal scaling, instead, allows fine-grained resource management and prevents intra-container interference. This effect is exemplified in Fig. 1 where the time per iteration of the double 256-bit bitcoin hash computation is shown for multiple such requests over a time period. The response times exhibited by two containers running on independent cores is markedly predictable compared to one container with two cores.

In the next section we employ core pinning to spawn containers using different cores. In particular we allocate one core for each container.

III. AUTOMATIC PROVISIONING SYSTEM

Here we exploit the container-provisioning approach to build a system capable of optimizing resource utilization. Our system scales the number of allocated containers to align with the varying demand in order to minimize costs while maintaining a high level of service. As a result of employing core pinning as explained in Section II, scaling decisions do not affect response time. A schematic diagram of our system is shown in Fig. 2. It comprises three components: Load Balancer (LB), Admission Controller (AC) and Auto-Scaler (AS). As suggested in [12], the LB component directs an admitted request to the active container reporting the most recent and lowest utilization rate. The AC component leverages CPU utilization statistics from the `cgroup` file-system to decide

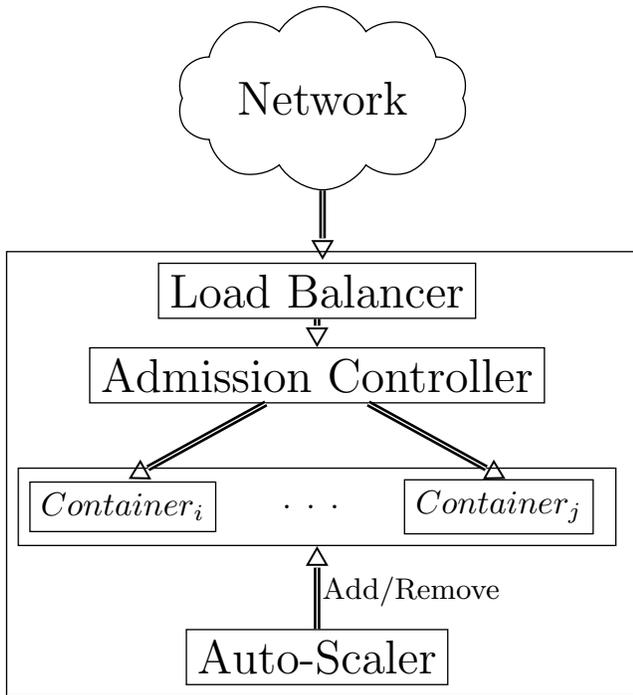


Fig. 2. Proposed container auto-scaling architecture.

whether an incoming request should be handled. The AS spawns new containers or removes active ones as appropriate, depending on demand.

We now delve into our design of the AC and AS components. Our LB implementation is inherited from [12].

A. Admission controller

To remain profitable, cloud providers need to limit the amount of resources they allocate to each tenant while still offering the sufficient amount required to honor SLAs. When resources are constrained in this way, an AC is necessary to limit accepted requests, thus keeping the service time predictable and reducing SLA violations.

Considering the iterative double 256-bit bitcoin hashing algorithm as an exemplary cloud application, the behaviour of a container (with a dedicated CPU core) is shown in Fig. 3. Each request triggers a different number of iterations. The response time, normalized by the number of iterations, is shown to be independent of the particular request’s characteristics.

The plot in Fig. 3 also shows that the relationship between the amount of allotted resources and SLA terms of service (such as the minimum response time) is not necessarily linear. In particular, there exists a discrepancy between CPU utilization (as reported by the operating system) and the actual occupied capacity of the core due to the use of hyper-threading [5]. This disparity is because the operating system considers two threads of the same core as two independent cores. Bearing this in mind and using the operating system metrics, a container exhibits a tri-stable CPU behavior: 0% CPU resource utilization when idle, 50% while continuously busy on a single core, 100% when continuously hyper-

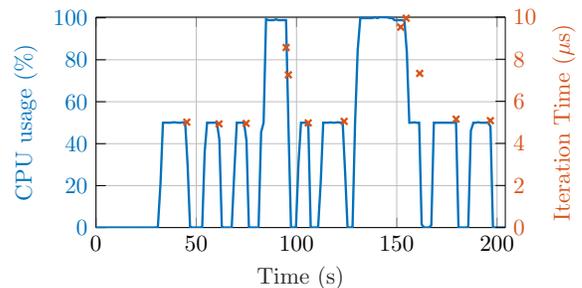


Fig. 3. Relationship between CPU utilization and service time. When the reported utilization is $\leq 50\%$, the iteration time is predictable.

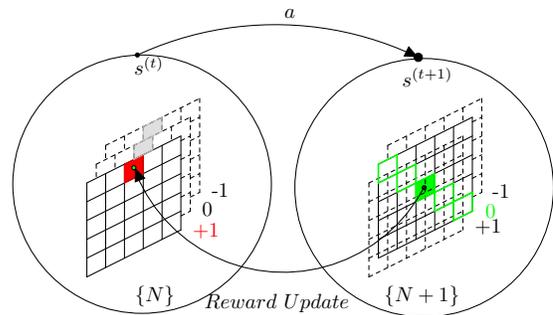


Fig. 4. Short-Term Memory Q-Learning mechanism used for the auto-scaler, cf. [13]. This schematic shows a scale out operation where in the epoch given by the interval $[t, t + 1)$, action a increases the number of containers from N to $N + 1$ and the state transitions from $s^{(t)}$ to $s^{(t+1)}$.

threading on the same physical core. The latter case points to saturation and unpredictable service times. The service time remains predictable only when a single thread (50% of a hyper-threaded core) is busy, as shown in Fig. 3. The above suggests that new requests should be admitted only when the container is not busy (0% utilization).

We finally remark that transients have a non-negligible impact on the correctness of admission decisions. Specifically, an admission error may occur in two cases: (i) a request was just assigned to the container, but the reported CPU usage value is still close to 0%, triggering the admission of an (otherwise undesirable) additional request; and (ii) the container just finished serving a request, but the reported CPU usage value is still close to 50%, triggering the rejection of a request that should have been admitted. These transients in reported CPU utilization are typically short-lived such that the inter-arrival time of requests, even at peak time, is much longer in comparison. However, in order to further reduce the likelihood of the first event, a new request is admitted if the reported value is $\leq 25\%$. In practice, this solution also reduces the likelihood of the second case.

B. Auto-scaler

The purpose of the AS is to allocate the minimum number of containers that is commensurate to the demand while still minimizing the number of dropped requests as reported by the empirical blocking probability. The scaling mechanism we

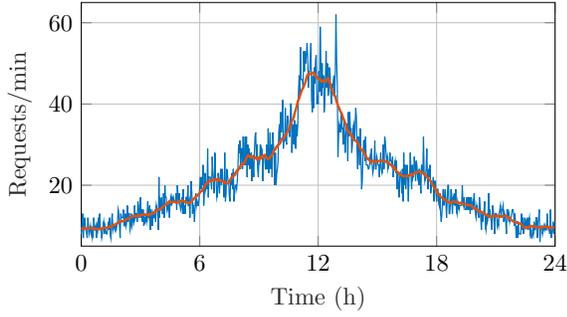


Fig. 5. Traffic profile observed during the Docker experiments. Traffic rates are taken over two-minute windows. The red line is the moving average over 30 samples.

employ is similar to the Q-Learning paradigm used in [13] as shown in Fig. 4.

We designate the permissible scaling actions as: $-n$ (remove n containers), $+n$ (add n containers) and 0 (maintain the existing number of containers). For $n = 1$, we can therefore describe the action space as follows:

$$A = \begin{cases} [-1; 0; +1] & \text{if } 1 < N_t < M; \\ [0; +1] & \text{if } N_t = 1; \\ [-1; 0] & \text{if } N_t = M; \end{cases} \quad (2)$$

where N_t is the current number of active containers and M is the maximum number of containers that can be allocated.

The state space is described by the triplet set of the number of containers, and the utilization in the previous and current epoch. We quantize the latter two components of the state space into nine levels, from 0% to 45% in steps of 5%. The last level encompasses the range from 45% to 100% utilization. The latter detail is required because the admission control function ensures that utilization never exceeds 50%. Fine-grained quantization of the state space is unnecessary as the action space is restricted to 3 discrete actions (cf. 2). Higher levels of quantization would increase the training time (owing to the curse of dimensionality) with little benefit to the quality of scaling policies learned.

The reward function (R_{sqr}), as presented in [13], consists of a penalty for blocking (R_{blk}) and another based on the number of containers provisioned (R_{res}):

$$\begin{aligned} R_{\text{sqr}} &= R_{\text{blk}} + R_{\text{res}} \\ R_{\text{blk}} &= \begin{cases} R_{\text{min}}, & \text{if } P \leq P_{\text{apt}} \\ \theta (P_{\text{apt}} - P), & \text{if } P > P_{\text{apt}}, \end{cases} \quad (3) \\ R_{\text{res}} &= \beta(1 - N_t) \end{aligned}$$

where P is the actual blocking rate, P_{apt} the acceptable level of outage as per the SLA, R_{min} is a small positive reward assigned to the agent when it keeps within the acceptable limits of service outage due to blocking, θ is the weight given to outage exceeding the acceptable level and β is the weighted cost incurred by the provider in deploying containers to handle client requests.

Compared to this approach given in [13], we simplify the mechanism in order to expedite learning. In particular, we set

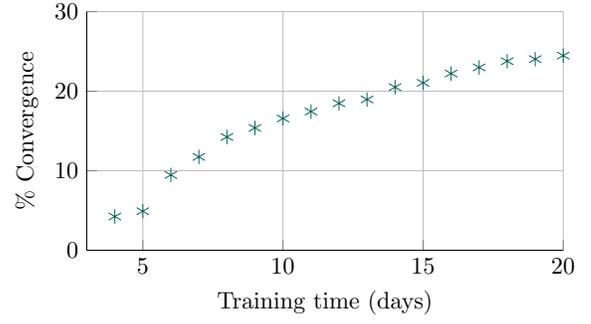


Fig. 6. Convergence evaluated at the end of a training day. High entropy in request inter-arrival times means that the scaler rarely visits the same states on subsequent days.

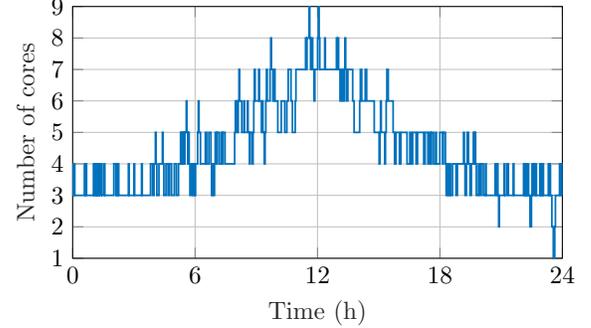


Fig. 7. Scaling decisions taken by Auto-Scaler algorithm during the Docker experiment at 20% of convergence (using $\beta = 0.02$).

$\theta = 1$ so that only β is adjusted. The ratio of the two influences the policies learned.

IV. EXPERIMENT SETUP

We implement our provisioning scheme using Docker containers on a Dell T640 server with 20 hyper-threaded cores running Ubuntu 18.04. In order to expedite learning, we choose $n = 1$ (cf. (2)) which effectively reduces the size of the action space. We implement CPU pinning and limit the maximum number of containers, M (cf. (2)), that can be provisioned to 9. This ensures that server capacity is never exceeded and that the host processes run on an independent core. The AC and LB functions are implemented as python applications and run on the host core.

The server is connected to client PCs in an isolated LAN via a high-speed switch. Bash scripts on the client PCs spawn requests to the server with varying frequency at different hours of the day to mimic peak and off-peak periods of typical real-world traffic profiles.

A 24-hour cycle is split into hourly periods as shown in Fig. 5. Each period has inter-arrival times following a discrete distribution $\lambda \sim U(0, \lambda_{\text{max}})$. By varying λ_{max} we create a suitable peak/off-peak profile. The use of a uniform distribution ensures high entropy in order to evaluate the robustness of the schemes in challenging conditions. As mentioned in Section III-A, we deploy the double 256-bit bitcoin hashing algorithm as our cloud application. Each admitted request triggers a different number of iterations, which makes it possible to mimic the diverse complexity of cloud applications.

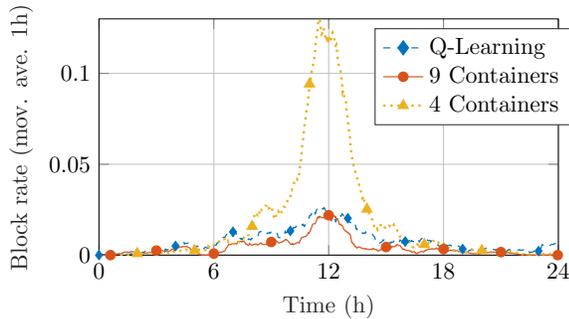


Fig. 8. The blocking rate observed over the time during Docker experiments in terms of rejected requests per second.

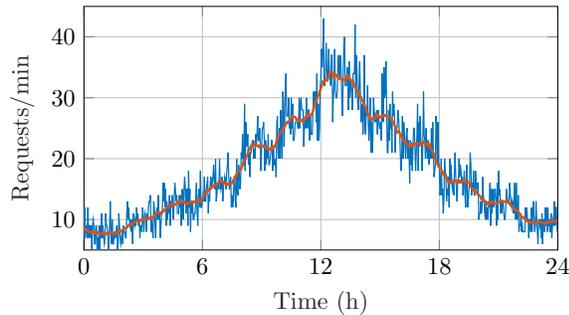


Fig. 10. Traffic profile observed during the Kubernetes experiments. Traffic rates are taken over two-minute windows. The red line is the moving average over 30 samples.

V. RESULTS AND DISCUSSION

We compare our provisioning scheme with the Google Horizontal Pod Autoscaler (HPA) for Kubernetes, a widely used container management tool. We also benchmark our scheme with static over-provisioning and under-provisioning. As comparative measures, we consider the following metrics:

- i Saved cost: cf. (1), the difference in cost between employing the maximum amount of resources throughout and using an auto-scaling algorithm to provision variable amounts of resources,
- ii Service time: the time taken to process a request normalized by the number of iterations triggered by it, measured at the server side in order to exclude network effects,
- iii Blocking rate: a measure of service availability defined as the percentage of dropped requests with respect to those received by the server.

We also introduce *convergence* as a key metric to measure the level of learning attained by our auto-scaler. This is the proportion of states visited by the scaler for which the policy is fully learned. Initially, a target number of statistically significant visits (30 in our case) is attributed to each state and serves as the baseline for decreasing ϵ , the probability of acting randomly. Each state starts off with $\epsilon = 1$ which is monotonically decayed as the number of visits to that given state increases and the agent acts less randomly in making scaling decisions. When a given state attains the specified number of visits, $\epsilon = 0$ and the scaler acts greedily (when in

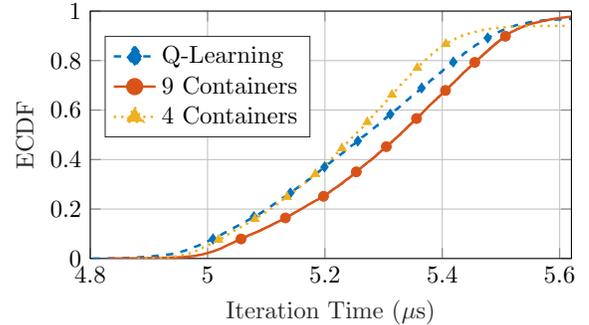


Fig. 9. Empirical CDF of the service time for Docker experiments.

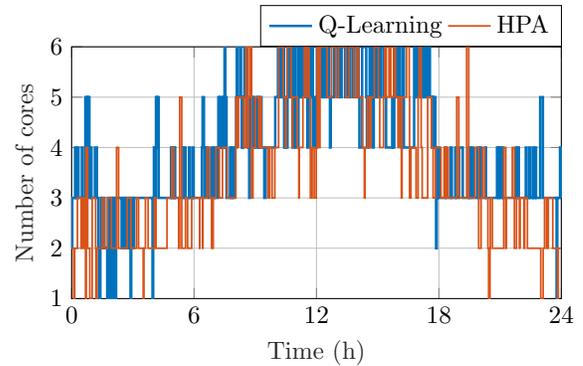
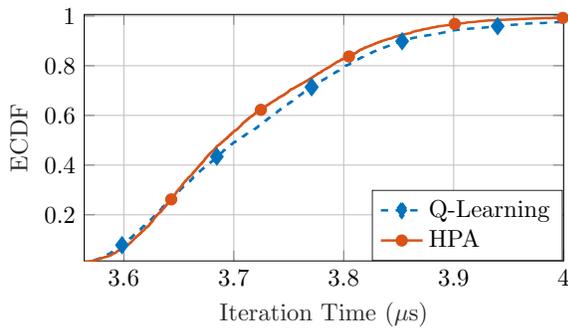


Fig. 11. Scaling decisions taken by the AS algorithm during the Kubernetes experiment at 25% of convergence (with $\theta = 1.0$ and $\beta = 0.02$).

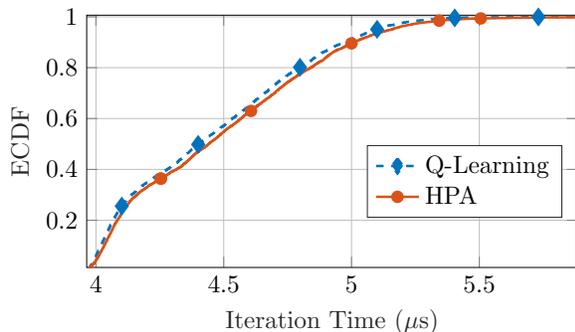
that state) according to the policy learned. In this way a good trade-off between exploration and exploitation is achieved. Owing to the uniform distribution in inter-arrival times in the various traffic periods the likelihood of the same state (the triplet of values defined in Section III-B) is low. As a consequence of the latter, the convergence rate is very gradual, as shown in Fig. 6.

1) *Docker experiments*: We initially test our scheme in a Docker environment. With reference to the offered load presented in Fig. 5, our scheme generates the scaling profile shown in Fig. 7 with 20% *convergence*. With the parameters we have adopted ($\theta = 1$ and $\beta = 0.02$), the system learns to act quite aggressively with respect to blocking and encourages the provisioning of additional containers even with modest rises in the traffic profile. A higher setting of β would result in a stiffer reaction of the scaler, and would likely result in higher blocking rates. With this setting, our solution achieves 51% saved cost.

The performance in terms of blocking rate is shown in Fig. 8. Our scaler achieves blocking rates that are comparable to the over-provisioned case with 9 containers throughout the 24 hours. It considerably outperforms the under-provisioned case in which only 4 containers are statically deployed and no adaptation is enforced over time. The saved cost for the under-provisioned case stands at 55% which is only marginally higher than that of our auto-scaler but with much poorer service availability, especially at peak traffic.



(a) Master node



(b) Slave node

Fig. 12. CDF of the service time for the Kubernetes experiments.

From Fig. 9, it is clear that the number of containers deployed has an effect on the service time. This is due to the fact that container processes share the L2 and L3 cache memory. The greater the number of active containers the more pronounced the impact. For this reason, the over-provisioned case with 9 containers exhibits slightly higher service times, whereas the under-provisioned case with 4 containers yields the lowest service time. Our scaling solution suffers a small deviation from the low service times for about half of the cases owing to the instances when it provisions more than the benchmark 4 containers at peak traffic. However for all cases considered, the maximum difference in service times is small with respect to the minimum values observed, i.e., the difference is less than $0.1 \mu s$, for more than 93% of the cases. This is because admission control ensures that the system only rarely reaches saturation.

2) *Kubernetes experiments*: We now compare the performance of our scaler against the commercial HPA for Kubernetes. We re-run the experiments with the traffic profile shown in Fig. 10 using two computers with different specifications. Our scaler autonomously learns the appropriate operating conditions for each computer to trigger the addition or removal of containers. HPA however requires that the threshold be set as an external input. Such a setting is often a trial and error process and is both application and configuration dependent. To obtain comparable results to our scaler, we set this threshold as 28%.

The comparison between the decisions of our scaler (at 25% convergence) and the ones made by HPA are shown in

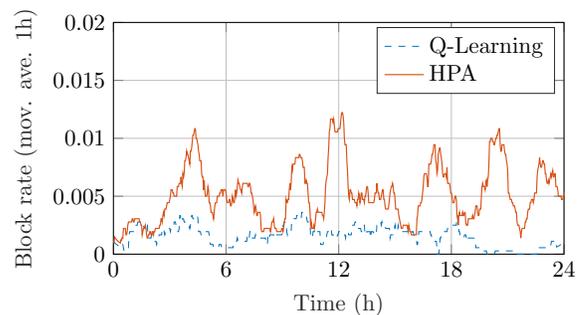


Fig. 13. Blocking rate observed over time during Kubernetes experiments in terms of rejected requests per second.

Fig. 11. Two different servers are used in these experiments to guarantee different service times. Even in this scenario, our provisioning scheme achieves predictable predictable service times, albeit different given the difference in compute power. In this case the less powerful compute engine (slave node) determines the service time benchmark.

While the two schemes save about the same costs and achieve similar results in terms of service time, as depicted in Fig. 12, their blocking performance differs. Our approach based on CPU core pinning and Q-Learning largely outperforms HPA in terms of blocking rate, as shown in Fig 13.

VI. RELATED WORK

The container provisioning problem is an evolving field of research given that this method of virtualization is quite recent compared to virtual machines (VMs). The authors of [14] leverage the maturity of VM provisioning schemes and propose an Integer Linear Programming (ILP) optimal mapping between containers and the available VMs.

In [15], the authors consider containers hosted within VMs and propose a system able to coordinate the vertical scaling of both. The authors of [11] propose ELASTICDOCKER, a vertical scaling system based on MAPE-K principles. Their scheme follows best-practices of setting thresholds to act as triggers for scaling operations. These works do not consider the effects of hyper-threading on vertical scaling and make assumptions based on queuing theory [16]. The prevalence of hyper-threaded CPUs therefore limits the performance of the proposed schemes [7].

In [17], Ye *et al.* propose a scheme which predicts the resource demands of an application and scales appropriately. An application-dependent, proactive resource provisioning scheme is proposed in [18]. The latter leverages the historical record of resource usage to trigger scaling actions.

The authors of [19] implement a horizontal scaler leveraging both reactive and proactive approaches. The amount of needed resources is calculated from a reactive term, based on a threshold, and a proactive term, based on traffic forecasting. The latter is performed with a simple ARMA model and may be ineffective for rapidly varying traffic.

In [20] Sangpetch *et al.* carry out a comparative study of three auto-scaling algorithms based on either Q-Learning,

artificial neural networks, or rules on thresholds. In this study, Q-Learning is found to achieve superior performance.

The preceding proposals do not consider CPU architecture and hence fail to capture the effect of hyper-threading on performance. Our approach takes this into account and curtails its pervasive effects by implementing CPU core pinning and by leveraging horizontal scaling. These approaches, coupled with our model free Q-Learning scaler, result in a robust and configuration-agnostic scheme that ensures predictable response times regardless of the application under consideration.

VII. CONCLUSIONS

We have presented a robust container provisioning system that leverages Q-Learning for autoscaling. We have demonstrated the consistent performance achieved by implementing CPU core pinning and horizontal scaling when compared to vertical scaling with hyper-threaded cores. We show that CPU core pinning simplifies the pricing models for cloud providers by facilitating an easy mapping between actual resources used and container resources assigned to tenants. Although our scheme curtails the application of hyper-threading and its advantages, the benefits of predictable and consistent high performance outweighs this disadvantage by far.

We have also demonstrated the superior performance of our scaling scheme when compared to the Horizontal Pod Autoscaler (HPA) for Kubernetes which is widely adopted in container provisioning platforms. Our Q-Learning scheme attains predictable response times in the face of highly dynamic traffic without the need for manual threshold setting. It is able to autonomously learn the appropriate scaling triggers without prior knowledge of the system configuration or the cloud application.

ACKNOWLEDGMENT

This work has been supported in part by the Spanish Ministry of Science, Innovation and Universities under Grant TIN2017-88749-R (DisCoEdge). The work of V. Mancuso was supported by the Spanish Ministry of Economy and Competitiveness through Ramon-y-Cajal under Grant RYC-2014-16285.

REFERENCES

- [1] W. Felten, A. Ferreira, R. Rajamony, and J. Rubio, "An updated performance comparison of virtual machines and linux containers," in *ISPASS*. IEEE, 2015, pp. 171–172.
- [2] K.-T. Seo, H.-S. Hwang, I.-Y. Moon, O.-Y. Kwon, and B.-J. Kim, "Performance comparison analysis of linux container and virtual machine for building cloud," *Adv. Sci. Tech. Lett.*, vol. 66, no. 105-111, p. 2, 2014.
- [3] D. Bernstein, "Containers and cloud: From lxc to docker to kubernetes," *IEEE Cloud Computing*, vol. 1, no. 3, pp. 81–84, 2014.
- [4] K. Hightower, B. Burns, and J. Beda, *Kubernetes: up and running: dive into the future of infrastructure*. O'Reilly Media, Inc., 2017.
- [5] M. Tegtmeier, "CPU utilization of multi-threaded architectures explained," 2015. [Online]. Available: <https://blogs.oracle.com/solaris/cpu-utilization-of-multi-threaded-architectures-explained-v2>
- [6] T. Leng, R. Ali, J. Hsieh, V. Mashayekhi, and R. Rooholamini, "An empirical study of hyper-threading in high performance computing clusters," *Linux HPC Revolution*, vol. 45, 2002.
- [7] Y. Ding, E. D. Bolker, and A. Kumar, "Performance implications of hyper-threading," in *Int. CMG Conference*, 2003, pp. 21–29.
- [8] C. D. Balaji Subramaniam. (2018) Future Highlight CPU Manager. [Online]. Available: <https://kubernetes.io/blog/2018/07/24/feature-highlight-cpu-manager/>
- [9] C. Prakash, P. Prashanth, U. Bellur, and P. Kulkarni, "Deterministic container resource management in derivative clouds," in *IC2E*. IEEE, 2018, pp. 79–89.
- [10] A. Podzimek, L. Bulej, L. Y. Chen, W. Binder, and P. Tuma, "Analyzing the impact of cpu pinning and partial cpu loads on performance and energy efficiency," in *CCGrid*. IEEE, 2015, pp. 1–10.
- [11] Y. Al-Dhuraibi, F. Paraiso, N. Djarallah, and P. Merle, "Autonomic vertical elasticity of docker containers with elasticdocker," in *CLOUD*. IEEE, 2017, pp. 472–479.
- [12] V. Cardellini, M. Colajanni, and P. S. Yu, "Dynamic load balancing on web-server systems," *IEEE Internet Computing*, vol. 3, no. 3, pp. 28–39, May 1999.
- [13] C. Ayimba, P. Casari, and V. Mancuso, "SQLR: Short-Term Memory Q-Learning for Elastic Provisioning," 2019, arXiv:1909.05772 [cs.NI].
- [14] M. Nardelli, C. Hochreiner, and S. Schulte, "Elastic provisioning of virtual machines for container deployment," in *SPEC*. ACM, 2017, pp. 5–10.
- [15] Y. al dhuraibi, Z. Faiez, N. Djarallah, and P. Merle, "Coordinating vertical elasticity of both containers and virtual machines," 03 2018.
- [16] W. Dawoud, I. Takouna, and C. Meinel, "Elastic virtual machine for fine-grained cloud resource provisioning," in *International Conference on Computing and Communication Systems*. Springer, 2011, pp. 11–25.
- [17] T. Ye, X. Guangtao, Q. Shiyong, and L. Minglu, "An auto-scaling framework for containerized elastic applications," in *BIGCOM*. IEEE, 2017, pp. 422–430.
- [18] H. Zhang, H. Ma, G. Fu, X. Yang, Z. Jiang, and Y. Gao, "Container based video surveillance cloud service with fine-grained resource provisioning," in *CLOUD*. IEEE, 2016, pp. 758–765.
- [19] C. Kan, "Docloud: An elastic cloud platform for web applications based on docker," in *ICACT*. IEEE, 2016, pp. 478–483.
- [20] A. Sangpetch, O. Sangpetch, N. Juangmarisakul, and S. Warodom, "Thoth: Automatic resource management with machine learning for container-based cloud platform," in *CLOSER*, 2017, pp. 75–83.