

# Tetris: Scheduling Long-Running Workloads for Load Balancing in Shared Containerized Clusters

Fei Xu, *Member, IEEE*, Xiyue Shen, Shuohao Lin, Li Chen, *Member, IEEE*, Zhi Zhou, *Member, IEEE*,  
Fen Xiao, Fangming Liu, *Senior Member, IEEE*

**Abstract**—Long-running containerized workloads (e.g., machine learning), which typically show *time-varying* patterns, are increasingly prevailing in shared production clusters. To improve workload performance, current schedulers mainly focus on optimizing *short-term* benefits of cluster load balancing or *initial container placement* on servers. However, this would inevitably bring many *invalid migrations* (*i.e.*, containers are migrated back and forth between servers over a short time window), leading to significant service level objective (SLO) violations. This paper introduces *Tetris*, a *model predictive control* (MPC)-based container scheduling strategy to judiciously migrate long-running workloads for cluster load balancing. Specifically, we first build a discrete-time dynamic model for *long-term* optimization of container scheduling. To solve such an optimization problem, *Tetris* then employs two main components: (1) a container resource predictor, which leverages time-series analysis approaches to accurately predict the container resource consumption; (2) an MPC-based container scheduler that jointly optimizes the cluster load balancing and container migration cost over a certain sliding time window. We implement and open source a prototype of *Tetris* based on K8s. Extensive prototype experiments and trace-driven simulations demonstrate that *Tetris* can improve the cluster load balancing degree by up to 77.8% without incurring any SLO violations, compared to the state-of-the-art container scheduling strategies.

**Index Terms**—Long-running containerized workloads, load balancing, migration cost, container scheduling

## 1 INTRODUCTION

LARGE production clusters are commonly hosting various *long-running* containerized workloads, ranging from online Web services to streaming processing systems and machine learning applications [1]. Unlike traditional batched jobs that are typically executed within seconds to minutes, these long-running containerized workloads generally last for hours to months [2], and often have stringent service level objectives (SLOs) [3]. However, as the request loads of long-running workloads are *time-varying* [4], sudden and unpredictable changes in container resource consumption can cause severe contention of shared cluster resources. Such cluster load imbalance can result in many potential SLO violations to workloads [5]. Hence, the mainstream cluster schedulers such as Borg [6] and Kubernetes (K8s) [7] enable container scheduling to achieve load balancing [8] in shared production clusters.

- Fei Xu, Xiyue Shen, Shuohao Lin are with the Shanghai Key Laboratory of Multidimensional Information Processing, School of Computer Science and Technology, East China Normal University, 3663 N. Zhongshan Road, Shanghai 200062, China. Email: fxu@cs.ecnu.edu.cn.
- Li Chen is with the School of Computing and Informatics, University of Louisiana at Lafayette, 301 East Lewis Street, Lafayette, LA 70504, USA. E-mail: li.chen@louisiana.edu.
- Zhi Zhou is with the Guangdong Key Laboratory of Big Data Analysis and Processing, School of Computer Science and Engineering, Sun Yat-sen University, 132 E. Waihuan Road, Guangzhou 510006, China. E-mail: zhoushi9@mail.sysu.edu.cn.
- Fen Xiao is with Tencent Incorporation, 33 Haitian Second Road, Nanshan District, Shenzhen 518054, China. E-mail: cherryxiao@tencent.com.
- Fangming Liu is with Peng Cheng Laboratory, and Huazhong University of Science and Technology, China. E-mail: fangminghk@gmail.com.

Manuscript received September XX, 2023; revised December XX, 2023.

Though the existing container scheduling policies such as Sandpiper [8] and Medea [2] perform well in achieving cluster load balancing, there still exist a noticeable number of *invalid migrations* (*i.e.*, containers are migrated back and forth between servers over a short time window) in shared production clusters. As evidenced by our motivational analysis (in Sec. 2.2) on Alibaba cluster trace v2018 [9], the number of invalid migrations is over 1,200, accounting for more than 50% of container migrations within the time window size of three hours. Such invalid migrations are likely to make the workloads hosted on the migrated containers suffer from serious SLO violations, leading to unexpected performance interference to the containers that are *co-located* on servers (*i.e.*, migration cost). Meanwhile, our motivation experiments in Sec. 2.2 reveal that the invalid migrations can significantly reduce the number of requests processed by an Apache Tomcat Web server hosted on a migrated container by up to 99.8%. Accordingly, it is essential for the cluster scheduler to circumvent invalid migrations, especially for long-running workloads.

Unfortunately, many research efforts have been devoted to making *short-term* scheduling decisions based on the *current* cluster status for workload consolidation [10] or load balancing [11]. Though such scheduling policies can acquire short-term (*e.g.*, the current or upcoming timeslot) benefits, they are oblivious to the *time-varying* resource consumption of long-running workloads, which is actually the *root cause of invalid migrations* as discussed in Sec. 2.2. There have also been recent works on the *long-term* optimization of container scheduling, which are reinforcement learning (RL)-based [5], [12] or control-based [13], [14] scheduling

policies, to partially tackle the issue of invalid migrations. Nevertheless, these techniques solely focus on optimizing the *initial placement* or resource auto-scaling of containerized workloads (*i.e.*, where to migrate) over the entire time period (*i.e.*, infinite future). The containers to be scheduled (*i.e.*, which to migrate) and the *migration cost* of containers have surprisingly received little attention. Such *long-term* optimization techniques can still cause unexpected SLO violations, as evidenced in Sec. 2.3. As a result, little research attention has been paid to developing container scheduling policies to fully deal with the invalid migrations of long-running workloads in production containerized clusters.

To fill this gap, we propose *Tetris*, a model predictive control [15] (MPC)-based container scheduling strategy to judiciously make migration decisions for long-running containerized workloads. *Tetris* simply optimizes container scheduling over a certain sliding time window rather than the infinite future (as in RL-based method [3]) for two reasons: *first*, the prediction accuracy of container resource consumption dramatically decreases as the prediction window size increases. Our prediction results using time-series techniques (in Sec. 5.2) indicate that the prediction error can exceed 20% as the prediction window size reaches 6; *second*, blindly increasing the time window size can significantly increase the number of samples of *Tetris*, which requires noticeable computation overhead to obtain container scheduling plans. To the best of our knowledge, *Tetris* is the first attempt to achieve the *long-term* optimization of container scheduling to circumvent as many as possible invalid migrations, by jointly optimizing the cluster load balancing and container migration cost over a certain sliding time window. Our main contributions are summarized as follows.

▷ *First*, we build a discrete-time dynamic model for shared containerized clusters to capture the time-varying resource consumption of containers and the corresponding mapping of containers on servers (Sec. 3). Based on such a model, we further devise a cost function of container scheduling over a time window and formulate our *long-term* workload scheduling optimization problem based on MPC, by jointly considering the cluster load imbalance degree and migration cost of containers.

▷ *Second*, we design an MPC-based container scheduling strategy to achieve *long-term* scheduling optimization for cluster load balancing (Sec. 4). Specifically, *Tetris* first accurately predicts the container resource consumption within a time window using time-series analysis methods. It then leverages the Monte Carlo method [16] to judiciously identify the container scheduling decision for each timeslot, by minimizing our formulated cost function of container scheduling. To particularly reduce the complexity of *Tetris* container scheduling strategy, we calculate the upper and lower bounds of server load imbalance degree to classify the migration source and destination servers in *Tetris*.

▷ *Finally*, we implement a prototype of *Tetris* (<https://github.com/icloud-ecnu/Tetris>) based on K8s. We evaluate the effectiveness and runtime overhead of *Tetris* with prototype experiments on 10 EC2 instances (*i.e.*, 60 containers) and large-scale simulations driven by Alibaba cluster trace v2018 (Sec. 5). Compared with the conventional scheduling strategy (*i.e.*, Sandpiper [8]) and the state-of-the-art RL-based method (*i.e.*, Metis<sup>+</sup>, a modified version of Metis [3]),

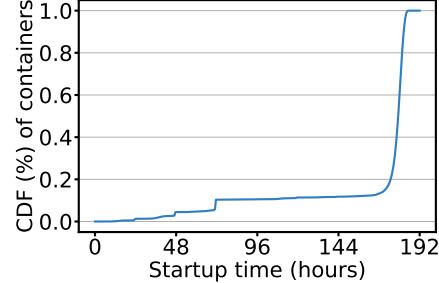


Fig. 1: Container uptime in Alibaba cluster trace v2018.

*Tetris* is able to improve the cluster load balancing degree by up to 77.8% while cutting down the migration cost by up to 79.5%, yet with acceptable runtime overhead.

## 2 BACKGROUND AND MOTIVATION

### 2.1 Long-running Containerized Workloads

Large-scale shared production clusters often host a number of long-running containers in response to latency-critical user requests [17]. Taking Alibaba as an example, its online services are mainly user-interactive and long-running applications that run on containers, such as search engines and online shopping [18]. As depicted in Fig. 1, over 80% of containers in the Alibaba production cluster last for more than 175 hours<sup>1</sup>. Due to the unpredictability of user requests, such long-running workloads have diverse demands on hardware resources including CPU, memory, network and disk I/O, which are generally dynamic variability and uncertainty in the amount of resource requests [19]. Furthermore, the co-location of containers can lead to contention for shared resources, which can easily cause performance interference and resource wastage in the cluster.

Based on the above, enabling adequate scheduling of long-running containerized workloads is essential to cluster load balancing. In shared production clusters, K8s achieves load balancing among servers by carrying out container migrations [7]. During the migration process, the container needs to stop services first, and then it is terminated on the source server and later restarted on the destination server, and finally it replicates the memory/storage data to restore services. Accordingly, the container scheduling inevitably brings a noticeable amount of migration cost, which requires to be considered during the scheduling of containers.

### 2.2 Invalid Migrations of Long-Running Workloads

While many existing scheduling policies (*e.g.*, Sandpiper [8]) perform well in load balancing in shared containerized clusters, they usually focus on the *short-term* benefits of container scheduling. As discussed above, the resource consumption of long-running containers is commonly *time-varying* (*i.e.*, fluctuating wildly over time). Such a *short-term* optimization of container scheduling policy can make *invalid migration* decisions for containers, resulting in unpredictable migration cost and potential SLO violations.

Specifically, we take a cluster of three servers hosting five containers illustrated in Fig. 2 as an example. The cluster

1. The total time of the Alibaba cluster trace v2018 is 192 hours.

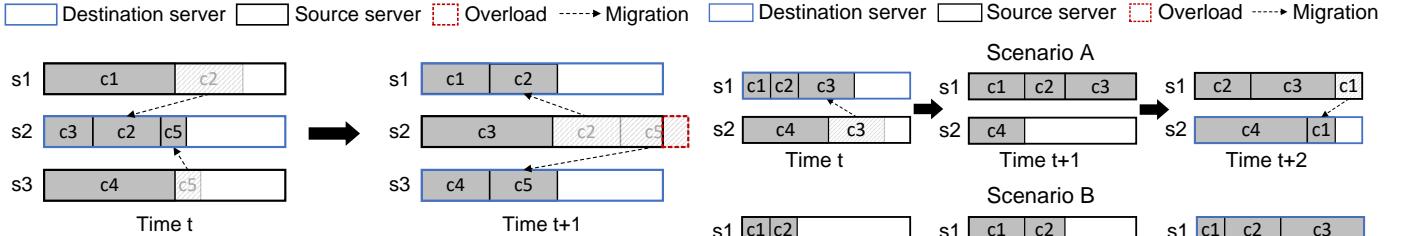


Fig. 2: Illustration of invalid migration. Containers  $c_2$  and  $c_5$  are migrated back and forth among servers  $s_1$ ,  $s_2$ , and  $s_3$ .

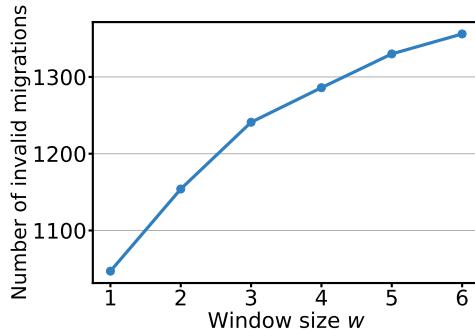


Fig. 3: Number of invalid migrations within a time window size varying from 1 to 6 hours.

scheduler (*e.g.*, the default K8s scheduler) first migrates container  $c_2$  and  $c_5$  from server  $s_1$  and  $s_3$ , respectively, to server  $s_2$ , in order to obtain a *temporary* benefits of load balancing at time  $t$ . Unfortunately, server  $s_2$  becomes overloaded as the resource consumption of three containers (*i.e.*,  $c_3$ ,  $c_2$ ,  $c_5$ ) dramatically increases, which triggers two container migrations (*i.e.*, container  $c_2$ ,  $c_5$ ) from server  $s_2$  to  $s_1$  and  $s_3$ , respectively. In such a case, naive scheduling policies are likely to make long-running containers migrated back and forth between servers, thereby causing heavy and unnecessary migration cost to containers. Accordingly, we formally define *invalid migrations* as in Definition 1.

**Definition 1.** If a container is migrated back and forth between two servers over a short time window  $[t, t + w]$ , we define such container migrations as *invalid migrations* within a time window size  $w$ .

Moreover, we validate the prevalence of invalid migrations with Sandpiper using the Alibaba cluster trace v2018 (*i.e.*, an 8-day period) [9]. As shown in Fig. 3, we observe that the number of invalid migrations increases as the scheduling time window increases. Specifically, the cumulative number of invalid migrations exceeds 1,200 (*i.e.*, around 50% of container migrations) as  $w$  is set as 3 hours, and the number reaches around 1,400, accounting for 75.3% of container migrations when  $w$  is set as 6 hours. To further illustrate the performance impact of invalid migrations, we conduct experiments on a cluster of 10 servers (*i.e.*, 60 containers) deployed with Apache Tomcat Web server. Our experiment results reveal that 61.7% of the containers are affected by invalid migrations. In addition, the number of requests processed by a migrated container can be reduced by up to 99.8% (74.0% on average), severely degrading the quality of long-running Web service. Therefore, it is essential to design a *long-term* container scheduling strategy to alleviate invalid migrations by explicitly considering the

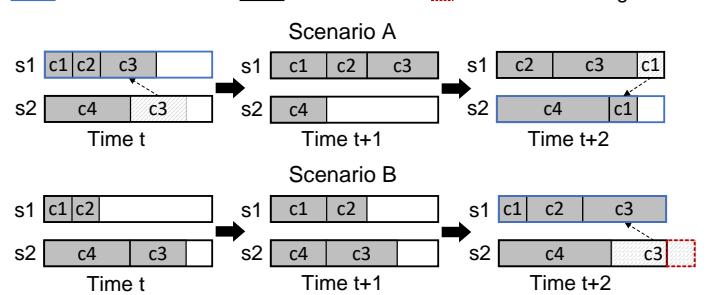


Fig. 4: An illustrative example of MPC-based container scheduling (our proposed *Tetris* in scenario A), as compared with RL-based container scheduling (Metis<sup>+</sup> [3] in scenario B).

future resource consumption of containers.

### 2.3 An Illustrative Example

To avoid invalid migrations, we design *Tetris* in Sec. 4, a *simple yet effective* container scheduling strategy that leverages the MPC approach to judiciously migrate long-running workloads for achieving the cluster load balancing. In particular, MPC adopts a compromise strategy that allows the current timeslot to be optimized while taking *finite future timeslots* into account [20]. To illustrate how *Tetris* works, we show a motivation example in Fig. 4 by comparing *Tetris* with an RL-based scheduling method (*i.e.*, Metis<sup>+</sup>, a modified version of Metis [3] which will be introduced in Sec. 5.1). Though Metis<sup>+</sup> greedily achieves the optimal load balancing degree and minimizes the migration cost over the *entire* time period of  $[t, t + 2]$ , it can still overload server  $s_2$  at time  $t + 2$  and trigger the migration of container  $c_3$  from server  $s_2$  to  $s_1$ . In contrast, our MPC approach (*i.e.*, *Tetris*, with the window size set as 2) jointly optimizes the load balancing degree and migration cost for two *sliding* timeslots (*i.e.*,  $[t, t + 1]$ , and  $[t + 1, t + 2]$ ) and guarantees SLOs of all containerized workloads. Accordingly, the RL-based method can optimize the scheduling objective *over the entire time period* at the cost of overloading a certain number of containers, while *Tetris* achieves the *long-term* optimization for container scheduling *within a certain sliding time window*. In addition, the RL-based method has the following problems such as requiring repeated training on a large amount of high-quality data samples and lacking interpretability as well as poor scalability [3], which will be further validated in Sec. 5.4.

**Summary.** Avoiding invalid migrations is critical to container scheduling, and such invalid migrations are mainly caused by the *short-term* (*e.g.*, the current or upcoming timeslot) optimization of container scheduling. Moreover, greedily optimizing container scheduling over the entire time period (*i.e.*, *infinite future*) is likely to cause unexpected SLO violations. Accordingly, there is a compelling need to design a *long-term* container scheduling strategy, by jointly optimizing the cluster load balancing and migration cost of containers within a certain sliding time window.

## 3 MODEL AND PROBLEM FORMULATION

In this section, we first build a discrete-time dynamic model to capture the load imbalance degree of clusters and the

TABLE 1: Key notations in our discrete-time dynamic model.

Notation	Definition
$\mathcal{M}, \mathcal{N}$	Sets of servers and containers
$W$	Time window size
$C_b(t), C_m(t)$	Cluster load imbalance degree and migration cost at time $t$
$C_{mig}^k(t)$	Unit migration cost of a container $k$ at time $t$
$cpu^k(t)$	CPU consumption of a container $k$ at time $t$
$mem^k(t)$	Memory consumption of a container $k$ at time $t$
$\overline{CPU}(t)$	Average CPU consumption of a server in the cluster at time $t$
$\overline{MEM}(t)$	Average memory consumption of a server in the cluster at time $t$
$CPU_i^{cap}(t)$	Available capacity of CPU resource of a server $i$ at time $t$
$MEM_i^{cap}(t)$	Available capacity of memory resource of a server $i$ at time $t$
$\alpha$	Normalized parameter of migration cost to load imbalance degree
$\beta$	Normalized parameter of memory resources to CPU resources
$x_i^k(t)$	Indicator whether a container $k$ is on a server $i$ at time $t$
$m^k(t)$	Indicator whether a container $k$ is migrated at time $t$

migration cost of containers. Next, we formulate a container scheduling optimization problem based on our dynamic model. The key notations in our model are summarized in Table 1.

### 3.1 Discrete-Time Dynamic Model

We consider a containerized cluster with a set of servers  $\mathcal{M}$  hosting a set of containers  $\mathcal{N}$ . We assume the time  $t$  is discrete and slotted, where  $t = t_0, t_1, \dots, t_W$  and  $W$  is the time window size. For each container  $k \in \mathcal{N}$ , its CPU and memory resource consumption at each time  $t$  is recorded as  $cpu^k(t)$  and  $mem^k(t)$ , respectively. For each server  $i \in \mathcal{M}$ ,  $CPU_i(t)$  and  $MEM_i(t)$  represent its total CPU and total memory resource consumption at time  $t$ , respectively.

**Modeling Cluster Load Imbalance Degree at Each Timeslot.** We use the variance of resource consumption of all cluster servers to represent the load imbalance degree of the cluster. A larger degree value indicates a severer load imbalance situation. By taking CPU and memory resources as an example, the cluster load imbalance degree is represented by calculating the weighted sum of the load imbalance degrees of each resource of servers, which is given by

$$C_b(t) = \frac{1}{|\mathcal{M}| - 1} \sum_{i \in \mathcal{M}} \left( \left( \sum_{k \in \mathcal{N}} x_i^k(t) \cdot cpu^k(t) - \overline{CPU}(t) \right)^2 + \beta \cdot \left( \sum_{k \in \mathcal{N}} x_i^k(t) \cdot mem^k(t) - \overline{MEM}(t) \right)^2 \right), \quad (1)$$

where  $\beta \in [0, 1]$  denotes the normalized parameter of memory resources to CPU resources. In practice,  $\beta$  can be empirically determined as the square of the ratio of CPU to memory resource consumption of workloads. Also,  $x_i^k(t)$  denotes whether the container  $k$  is hosted on the server  $i$ .

$$x_i^k(t) = \begin{cases} 1, & \text{if a container } k \text{ runs on a server } i, \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

Then, we proceed to denote the average CPU resource consumption  $\overline{CPU}(t) = \frac{1}{|\mathcal{M}|} \sum_{k \in \mathcal{N}} cpu^k(t)$  and average memory resource consumption  $\overline{MEM}(t) = \frac{1}{|\mathcal{M}|} \sum_{k \in \mathcal{N}} mem^k(t)$  of a server in the cluster.

**Modeling Migration Cost at Each Timeslot.** We use the sum of unit cost of the migrated containers to denote the migration cost  $C_m(t)$ , which is formulated as

$$C_m(t) = \sum_{k \in \mathcal{N}} C_{mig}^k(t) \cdot m^k(t), \quad (3)$$

where  $m^k(t) = \sum_{i \in \mathcal{M}} x_i^k(t) \cdot (1 - x_i^k(t-1))$  indicates whether a container  $k$  is migrated from a source server  $i$  at time  $t$ , and  $C_{mig}^k(t)$  is the unit migration cost of a container  $k$ . As elaborated in Sec. 2.1, K8s [7] implements the container migration with three steps: service stop, container deletion and restart, and data replication. The data replication overhead is proportional to the container's memory footprint  $mem^k(t)$  [21], while the overhead of the other two steps can be considered as a constant value  $\delta$ . Hence,  $C_{mig}^k(t)$  can be given by

$$C_{mig}^k(t) = \delta + \gamma \cdot mem^k(t), \quad (4)$$

where  $\gamma$  is a model coefficient obtained by workload profiling.

To sum up, we further formulate the *cost function*  $C(t)$  of *container scheduling at each timeslot*, by combining the cluster load imbalance degree  $C_b(t)$  and the migration cost  $C_m(t)$  as below,

$$C(t) = C_b(t) + \alpha \cdot C_m(t), \quad (5)$$

where  $\alpha \in [0, 1]$  denotes the normalized parameter of migration cost to cluster load imbalance degree. Specifically, we incorporate  $\overline{CPU}(t)$  and  $\overline{MEM}(t)$  into Eq. (1), and then incorporate Eq. (4) and  $m^k(t)$  into Eq. (3), yielding Eq. (6) and Eq. (7), respectively, which are formulated as

$$C_b(t) = \frac{2}{|\mathcal{M}| - 1} \sum_{i \in \mathcal{M}} \sum_{k, l \in \mathcal{N}} x_i^k(t) x_i^l(t) \cdot (cpu^k(t)cpu^l(t) + \beta \cdot mem^k(t)mem^l(t)) + C_1, \quad (6)$$

$$C_m(t) = C_2 - \sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{M}} (\delta + \gamma \cdot mem^k(t)) \cdot x_i^k(t) x_i^k(t-1). \quad (7)$$

Given a containerized cluster and workloads at time  $t$ ,  $C_1 = \frac{|\mathcal{M}| \cdot (\overline{CPU}(t)^2 + \beta \cdot \overline{MEM}(t)^2)}{|\mathcal{M}| - 1} - \frac{\sum_{k \in \mathcal{N}} (cpu^k(t)^2 + \beta \cdot mem^k(t)^2)}{|\mathcal{M}| - 1}$  and  $C_2 = \delta \cdot |\mathcal{N}| + \gamma \cdot \sum_{k \in \mathcal{N}} mem^k(t)$  are both *constant* values. The mathematical derivations can be found in Appendix A. By analyzing the formulations above, Eq. (6) indicates that the cluster load imbalance degree is determined by the sum of the pairwise multiplications between the resource consumption of containers hosted on each server. Eq. (7) implies that the migration cost across the cluster can be determined by the negative sum of migration costs of the migrated containers.

### 3.2 Workload Scheduling Optimization Problem over a Time Window

Based on our discrete-time dynamic model above, we proceed to formulate the *long-term* optimization problem of container scheduling based on MPC. At each time  $t$ , MPC leverages the predicted container resource consumption (*i.e.*,  $cpu^k(t)$ ,  $mem^k(t)$ ) to make scheduling decisions (*i.e.*, judiciously deciding  $x_i^k(t)$ ) to minimize the cost function  $C(t)$  over the time window (as in Eq. (8)). We assume the

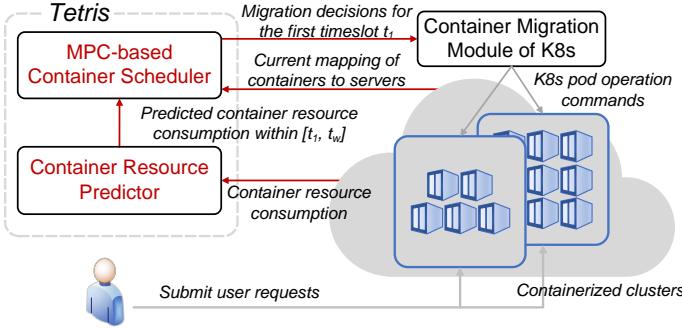


Fig. 5: Overview of *Tetris* prototype in containerized clusters.

scheduling within the time window  $t \in [t_1, t_W]$ , and our optimization problem can be formulated as

$$\min_{x_i^k(t)} \sum_{t=t_1}^{t_W} C(t) \quad (8)$$

$$\text{s.t. } \sum_{i \in \mathcal{M}} x_i^k(t) = 1, \quad \forall k \in \mathcal{N} \quad (9)$$

$$\sum_{k \in \mathcal{N}} x_i^k(t) \cdot \text{cpu}^k(t) \leq \text{CPU}_i^{\text{cap}}(t), \quad \forall i \in \mathcal{M} \quad (10)$$

$$\sum_{k \in \mathcal{N}} x_i^k(t) \cdot \text{mem}^k(t) \leq \text{MEM}_i^{\text{cap}}(t), \quad \forall i \in \mathcal{M} \quad (11)$$

where the hard constraints in our optimization problem are as follows. Constraint (9) limits each container to be hosted only on one server per time to avoid scheduling conflicts. Constraints (10)-(11) ensure the total resource consumption on each server cannot exceed the server resource capacity.

**Problem Analysis.** For each timeslot, the problem defined in Eq. (8) can be easily reduced to a *multiprocessor scheduling problem (MSP)* [22]. MSP is considered as finding a schedule for multiple tasks to be executed on a multiprocessor system at different timeslots, with the aim of minimizing the completion time. Such a scheduling problem is known to be NP-hard [23]. Moreover, Eq. (8) is a multi-objective optimization problem that jointly considers minimizing the migration cost and the cluster load imbalance degree, which makes it hard to solve. As a result, our optimization problem turns out to be harder than MSP. In the upcoming section, we turn to leveraging the Monte Carlo method [16] to solve our long-term workload scheduling optimization problem.

## 4 Tetris DESIGN

Based on our discrete-time dynamic model and optimization problem defined in Sec. 3, we proceed to design *Tetris*, a simple yet effective container scheduling strategy that leverages the MPC approach to optimize the cluster load balancing degree and container migration cost jointly over a certain time window.

### 4.1 Overview of *Tetris*

As illustrated in Fig. 5, *Tetris* comprises two pieces of modules including a *container resource predictor* (Sec. 4.2) and a *MPC-based container scheduler* (Sec. 4.3). After users submit workload requests to the containerized cluster, *Tetris* first leverages the *predictor* to estimate the resource consumption of containers over a given time window  $W$ , which is

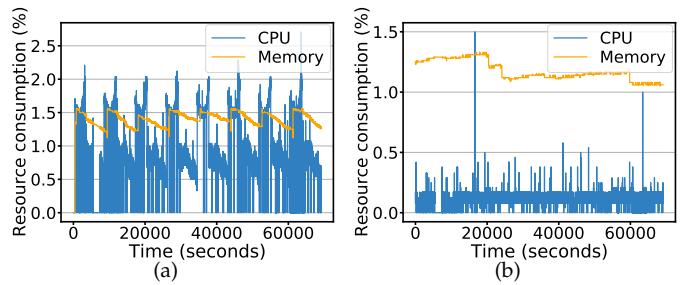


Fig. 6: Resource consumption of containers with respect to the physical machine resources over time: (a) a periodic container ( $id = 71,417$ ), and (b) an aperiodic container ( $id = 58$ ) in Alibaba cluster trace v2018.

input to the *scheduler* for calculating the migration cost of containers. By *jointly* optimizing the cluster load balancing degree and container migration cost, the *scheduler* further decides the appropriate migration plans during the period of time window  $W$ , which will be elaborated in Alg. 1. Finally, we implement a *container migration module of K8s*, which performs the appropriate migration decisions for the *first timeslot* in the containerized cluster, while discarding the migration decisions for the remaining timeslots. In particular, such a migration module can convert container scheduling decisions into a series of K8s pod operations (*i.e.*, pod deletion and creation commands executed on migration source and destination servers). The prototype of *Tetris* is implemented based on K8s, with over 1,000 lines of Python and Linux Shell codes which are publicly available on GitHub (<https://github.com/icloud-ecnu/Tetris>).

### 4.2 Predicting Container Resource Consumption

We first classify workloads as *periodic* and *aperiodic* containers shown in Fig. 6, as most containerized workloads show a diurnal pattern [18]. Specifically, we first obtain the *frequency* of container resource consumption values and calculate its corresponding *periodicity*. With such a periodicity value, we then slice the container resource consumption into a number of segments. We finally calculate the Pearson correlation coefficients between every two consecutive segments and decide whether such coefficients exceed a given periodicity threshold  $thr_{\text{period}}$  (*e.g.*, 0.85). If the correlation coefficients exceed  $thr_{\text{period}}$ , we consider such a container as *periodic*. Otherwise, we consider the container as *aperiodic*.

Next, we proceed to predict the resource consumption for such two types of containers separately. In more detail, we first adopt ARIMA [24] to predict the resource consumption of *periodic* containers. We use the *historical* container resource consumption to pre-train a *shared* ARIMA model for newly-launched containers. To improve the prediction accuracy, we further leverage incremental learning [25] to train a *personalized* ARIMA model for each container that has been running for a certain period of time, and we can update the model periodically. In particular, we mainly obtain three key parameters (*i.e.*, the autoregressive process of order  $p$ , the moving average process of order  $q$ , and the difference order  $d$ ) of the ARIMA model. Second, we leverage LSTM [26] to predict the resource consumption of *aperiodic* containers. We adopt a 4-layer stacked LSTM model with 50 neurons per layer to increase the depth of the neural

network, and add a *Dropout* layer between every two LSTM layers to reduce overfitting. To keep the LSTM model simple and alleviate the impact of error propagation, we directly conduct the multistep-ahead predictions, and the number of prediction steps is set as  $W$ .

**Algorithm 1:** *Tetris*: Long-term container scheduling algorithm.

---

**Input:** Current mapping  $\mathbf{X}_{t_0}$  of container set  $\mathcal{N}$  to server set  $\mathcal{M}$ , time window size  $W$ , number of samples  $Z$ , number of trials  $K$ .  
**Output:** Container scheduling decisions  $\mathbf{X}_{t_1}^{min}$  at time  $t_1$ .

```

1 for all  $z \in [1, Z]$  do
2   for all  $t \in [t_1, t_W]$  do
3     // server classification
4     Initialize: the  $load_i(t)$  threshold for
      migration source server  $load_{src}(t) \leftarrow$ 
      Eq. (13) and that for migration destination
      server  $load_{dest}(t) \leftarrow$  Eq. (14);
6     for each server  $i \in \mathcal{M}$  do
7       Calculate the load imbalance degree
         $load_i(t) \leftarrow$  Eq. (12);
8       if  $load_i(t) > load_{src}(t)$  then
9         Put server  $i$  to migration source server
          set  $\mathcal{M}_{src}$ ;
10      if  $load_i(t) < load_{dest}(t)$  then
11        Put server  $i$  to migration destination
          server set  $\mathcal{M}_{dest}$ ;
12
13    // container scheduling
14    for all  $k \in [1, K]$  do
15      for each server  $i \in \mathcal{M}_{src}$  do
16        Obtain the container set to be migrated
           $\mathcal{N}_{mig}$ ;
17        for each container  $c \in \mathcal{N}_{mig}$  do
18          for each server  $i \in \mathcal{M}_{dest}$  do
19            Calculate the container migration
              cost  $C(t) \leftarrow$  Eq. (5) if container  $c$ 
              is migrated to server  $i$ ;
20            Update  $\mathbf{X}_t$  by migrating container  $c$  to
              the server  $i$  with the smallest  $C(t)$ ;
21          if  $\mathbf{X}_t$  satisfies scheduling constraints (9) –
              (11) then
22            break;
23          Update  $\mathbf{X}_t \leftarrow \mathbf{X}_{t-1}$ ; //  $\mathbf{X}_t$  cannot
              satisfy constraints
24
25      Calculate the container migration cost over the
        time window  $W$  as  $cost_z \leftarrow$  Eq. (8) of  $\mathbf{X}_t$  for
        each sample  $z$ ;
26
27  Obtain  $\mathbf{X}_t^{min} \leftarrow \mathbf{X}_t$  with the minimized  $cost_z$ 
    among all  $Z$  samples;
28  return: Scheduling decisions for the first timeslot
    (i.e.,  $\mathbf{X}_{t_1}^{min}$ ).

```

---

### 4.3 MPC-based Container Scheduling

We design *Tetris* in Alg. 1 by leveraging the MPC approach to make container scheduling decisions over the time win-

dow  $W$ . Specifically, we use a  $|\mathcal{M}| \times |\mathcal{N}|$  matrix  $\mathbf{X}_t$  of  $x_i^k(t)$  ( $i \in \mathcal{M}, k \in \mathcal{N}$ ) to denote the mapping of containers to servers at time  $t \in [t_0, t_W]$ . According to MPC, we first solve the scheduling optimization problem in Sec. 3.2 to obtain  $\mathbf{X}_t, \forall t \in [t_1, t_W]$ , at the current time  $t_0$ . Then, we only perform the container scheduling decision  $\mathbf{X}_{t_1}$  for the first timeslot  $t_1$ . After the containers are scheduled to migration destination servers at the “current” time  $t_1$ , we continue solving the scheduling optimization problem and then perform the container scheduling decisions for the “first” timeslot  $t_2$ . Accordingly, *Tetris* can be *periodically* (e.g., for several minutes or hours) executed at each timeslot, so as to maintain load balancing of the containerized cluster.

In more detail, we solve the scheduling optimization problem based on the Monte Carlo method [16] as discussed in Sec. 3.2. We take  $Z$  *sample* solutions in total, and each sample solution includes a set of container scheduling decisions  $\mathbf{X}_t$  over the time window  $t \in [t_1, t_W]$  (lines 1 to 2). To obtain the container scheduling decisions for each timeslot  $t$ , we design two phases in Alg. 1 including *server classification* (lines 3 to 9) and *container scheduling* (lines 10 to 19), which are elaborated as follows.

**Phase I: Server Classification.** We classify the cluster servers into three categories, *i.e.*, migration source servers and destination servers as well as the remaining servers (lines 3 to 9). Based on the cluster load imbalance degree (*i.e.*, Eq. (6)) defined in our dynamic model, we can obtain the simplified *load imbalance degree* (*i.e.*,  $load_i(t)$ ) for each server  $i$  given by

$$load_i(t) = \sum_{k,l \in \mathcal{N}} x_i^k(t)x_i^l(t) \cdot (cpu^k(t)cpu^l(t) + \beta \cdot mem^k(t)mem^l(t)). \quad (12)$$

Obviously, the CPU and memory resource consumption of each server are  $\overline{CPU(t)}$  and  $\overline{MEM(t)}$ , respectively, when the cluster reaches the “ideal” *load-balanced state*. As the load imbalance degree of each server  $load_i(t)$  can vary with the number of hosted containers and the resource consumption of containers, we obtain the thresholds of  $load_i(t)$  for migration source servers and destination servers as in Theorem 1.

**Theorem 1.** *Given a server hosting a set of containers with the average CPU and memory resource consumption (*i.e.*,  $\overline{CPU(t)}$  and  $\overline{MEM(t)}$ ) at time  $t$ , the thresholds of  $load_i(t)$  for migration source server and destination server can be formulated as below:*

$$load_{src}(t) = \frac{|\mathcal{N}| - 1}{2|\mathcal{N}|} \cdot \left( \overline{CPU(t)}^2 + \beta \cdot \overline{MEM(t)}^2 \right), \quad (13)$$

$$load_{dest}(t) = \frac{1}{2} \left( \overline{CPU(t)}^2 + \beta \cdot \overline{MEM(t)}^2 - \sum_{k=1}^{n(t)} (cpu^k(t)^2 + \beta \cdot mem^k(t)^2) \right). \quad (14)$$

By sorting containers in descending order with the CPU and memory resource consumption,  $cpu^k(t)$  and  $mem^k(t)$  represent the CPU and memory consumption of the  $k$ -th container at time  $t$ , respectively.  $n(t)$  denotes the minimum number of containers which satisfies  $\sum_{k=1}^{n(t)} cpu^k(t) \leq \overline{CPU(t)}$  and  $\sum_{k=1}^{n(t)} mem^k(t) \leq \overline{MEM(t)}$ .

*Proof.* The proof can be found in Appendix B.  $\square$

**Phase II: Container Scheduling.** To make Alg. 1 practical, we sample a set of containers to be migrated  $\mathcal{N}_{mig}$  from the set of migration source servers  $\mathcal{M}_{src}$ . To achieve comprehensive coverage of possibilities with a small number of samples, we adopt Latin Hypercube Sampling (LHS) with a sampling rate of  $v$  per migration source server (lines 11 to 12). For each container to be migrated, we further select the migration destination server with the smallest container migration cost and then update  $\mathbf{X}_t$  (lines 13 to 16). Considering the randomness of sampling, the obtained container scheduling decisions  $\mathbf{X}_t$  can violate constraints (9) – (11). If no violations occur,  $\mathbf{X}_t$  is valid and we continue the container scheduling process for the next timeslot  $t + 1$ . Otherwise, we require re-sampling  $\mathcal{N}_{mig}$  and obtain another  $\mathbf{X}_t$ , until  $K$  trials are finished (line 10) or the scheduling constraints are satisfied (lines 17 to 19).

Finally, we iteratively calculate the container migration cost over the time window  $W$  for each sample. We choose the container scheduling decisions  $\mathbf{X}_t^{min}$  with the minimized migration cost among  $Z$  samples. According to MPC, we only perform the container scheduling decisions  $\mathbf{X}_{t_1}^{min}$  for the first timeslot  $t_1$  (lines 20 to 22).

**Determining Parameters in Tetris.** We obtain three key parameters (*i.e.*,  $v$ ,  $W$ ,  $Z$ ) in Alg. 1 as below. *First*, we empirically set the sampling rate  $v$  for migration containers as the minimum of two ratios (*i.e.*, one is the ratio of migration source servers with the CPU resource consumption exceeding  $\overline{CPU}(t)$ , and the other is the server ratio with the memory resource consumption exceeding  $\overline{MEM}(t)$ ). *Second*, we also empirically set the time window size  $W$  as the minimum of two values (*i.e.*,  $W_1$ ,  $W_2$ ).  $W_1$  is the maximum time window with the prediction error of container resource consumption below a threshold (*e.g.*, 15%).  $W_2$  is the window size with the fastest growing number of invalid migrations (in Fig. 3). *Third* is to set the number of samples  $Z$ . To uniformly sample different types of containers, we perform  $k$ -means clustering on the CPU and memory resource consumption of containers. We simply set the cluster number  $k$  as the lower bound of  $Z$ .

**Time Complexity Analysis.** According to Alg. 1, the time complexity of *Tetris* is in the order of  $\mathcal{O}(ZWK \cdot |\mathcal{N}| \cdot |\mathcal{M}|)$ . As the number of containers  $|\mathcal{N}|$  and the number of servers  $|\mathcal{M}|$  are both *far larger* than the time window size  $W$  and the number of trials  $K$ , the complexity can be reduced to  $\mathcal{O}(Z \cdot |\mathcal{N}| \cdot |\mathcal{M}|)$ . Accordingly, given a containerized cluster and workloads, the computation overhead of *Tetris* is practically acceptable, which will be validated in Sec. 5.4.

## 5 PERFORMANCE EVALUATION

In this section, we evaluate *Tetris* by conducting prototype experiments based on a 60-container K8s cluster in Amazon EC2 and complementary large-scale simulations driven by a real-world production cluster trace from Alibaba v2018. Our evaluation focuses on answering the following questions:

- *Accuracy*: Can *Tetris* accurately predict the resource consumption of long-running containers? (Sec. 5.2)
- *Effectiveness*: Can our *Tetris* strategy jointly optimize the cluster load balancing degree and migration cost for long-running containers without incurring SLO violations? (Sec. 5.3)

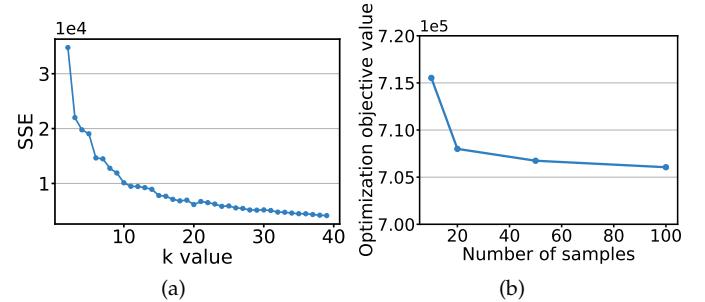


Fig. 7: (a) Sum of the squared errors (SSE) of  $k$ -means clustering on CPU and memory resource consumption of containers with different  $k$  values, and (b) the optimization objective value decreasing slowly from 20 samples in the Monte Carlo method.

- *Overhead*: How much runtime overhead does *Tetris* bring? (Sec. 5.4)

### 5.1 Experimental Setup

**Configurations of K8s Cluster and *Tetris* Parameters.** We carry out prototype experiments upon 10 m6a.large EC2 instances. Each instance is equipped with 2 vCPUs and 8 GB memory and initially hosting 6 containers. To obtain complementary insights, we also build a trace-driven simulator based on a discrete-event simulation framework [27]. In addition, we set several key parameters in *Tetris* dynamic model and scheduling strategy using the method in Sec. 4.3. The number of samples  $Z$  is set as 20, as the container resource consumption can be clustered with 20 types and the optimization value converges at 20 samples shown in Fig. 7. The time window size  $W$  is set as 2 based on the prediction accuracy of container resource consumption and the occurrence of invalid migrations in Fig. 3. The number of trials  $K$  and the sampling ratio  $v$  of migration containers are empirically set as 10 and 40%, respectively.

**Workloads and Datasets.** We employ three representative workloads in the containerized clusters, including Apache Tomcat server<sup>2</sup>, Redis<sup>3</sup>, and ResNet50 [28] training. We use Apache-benchmark<sup>4</sup> and redis-benchmark<sup>5</sup> as the time-varying user requests for Tomcat Web server and Redis, respectively. We adopt CIFAR-10<sup>6</sup> as the dataset for ResNet50. We randomly deploy the three workloads on the 60 containers in the K8s cluster. For the prediction of container resource consumption and large-scale simulations, we adopt Alibaba cluster trace v2018 [9] which contains around 4,000 machines over 8 days. We pre-process the CPU and memory resource consumption of 67,437 containers using linear interpolation and zero-value padding. We set each timeslot as 1 hour and calculate the average resource consumption of containers per hour.

**Baselines and Metrics.** We evaluate *Tetris* against the conventional Sandpiper [8] and the state-of-the-art Metis<sup>+</sup> (*i.e.*, a modified version of Metis [3]) scheduling algorithms. To achieve load balancing, Sandpiper performs a greedy worst-fit algorithm to schedule containers according to the

2. <https://tomcat.apache.org/>

3. <https://redis.io/>

4. <https://httpd.apache.org/docs/2.4/programs/ab.html>

5. <https://redis.io/topics/benchmarks>

6. <https://www.cs.toronto.edu/~kriz/cifar.html>

TABLE 2: Number of containers with periodic or aperiodic CPU/memory resource consumption.

Resource	CPU	Memory
Number of periodic containers	49,962	40,501
Number of aperiodic containers	17,475	26,936

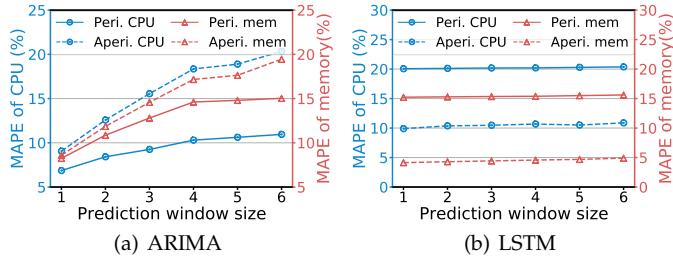


Fig. 8: Average prediction accuracy of CPU and memory resource consumption of 67,437 containers in Alibaba cluster trace v2018, by varying the prediction window size from 1 to 6.

container index *volume* calculated by the multi-dimensional resource consumption [8]. Metis<sup>+</sup> uses the method in *Tetris* to select migration source servers and migrated containers and leverages the negative value of Eq. (5) as the reward of RL to make container scheduling (*i.e.*, where to place) decisions. In particular, we use the *mean absolute percentage error* (MAPE) [29] to evaluate the efficacy of the *predictor* module of *Tetris*. Meanwhile, we adopt the *load imbalance degree* defined in Eq. (1) and the *migration cost* defined in Eq. (3) as well as the number of *SLO violations* to evaluate the efficacy of the *scheduler* module of *Tetris*. Due to the randomness of sampling in the Monte Carlo method and workload placement on containers, we illustrate the container scheduling performance with error bars of standard deviation by repeating experiments three times.

## 5.2 Predicting Resource Consumption of Containers in *Tetris*

We first classify the resource consumption of containers as periodic or aperiodic using the method elaborated in Sec. 4.2. As shown in Table 2, the periodic CPU and memory resource consumption accounts for 74.1% and 60.1% of containers, respectively, while the proportion of aperiodic CPU and memory resource consumption are only 25.9% and 39.9%, respectively. Accordingly, most of the containers show periodic resource consumption over time, which is consistent with the latest measurement study [18].

We next leverage ARIMA and LSTM to predict the periodic and aperiodic container resource consumption. To train the prediction model, we use 70% of the trace data as the training dataset. As shown in Fig. 8, we first observe that ARIMA is more accurate than LSTM for predicting the resource consumption of periodic containers. This is because the moving average and autoregression in ARIMA can accurately capture the periodic resource consumption of containers, with a moderate prediction error (*i.e.*, less than 15%). Second, LSTM achieves a more accurate prediction error (*i.e.*, 5% – 10%) in predicting the resource consumption of aperiodic containers compared with ARIMA as expected. This is mainly because ARIMA can easily be affected by abnormal spikes in aperiodic resource consumption, while

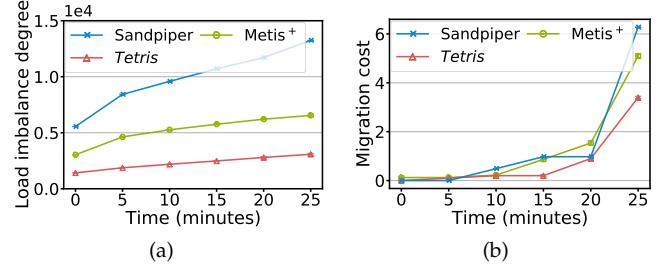


Fig. 9: Performance comparison of *Tetris* with Sandpiper and Metis<sup>+</sup> strategies under K8s-based prototype experiments, in terms of the cumulative values of (a) cluster load imbalance degree and (b) migration cost of containers over time (*i.e.*, 5 timeslots with 5 minutes each).

LSTM can learn such large resource fluctuations through model training. Interestingly, LSTM fails to predict the resource consumption of periodic containers depicted in Fig. 8(b) simply because LSTM contains nonlinear activation functions and thus overfits the periodic data, thereby causing around 15% – 20% of prediction error for periodic containers [30]. Third, the prediction error of container resource consumption with ARIMA can significantly increase from 7.8% to 20.1% with the time window size ranging from 1 to 6. That is the reason why *Tetris* controls the time window size to maintain an acceptable prediction error (*i.e.*, within 15%) of container resource consumption.

In addition, the average prediction time of ARIMA and LSTM are 0.25 seconds and 0.36 seconds, respectively, for each container. As our prediction method in Sec. 4.2 requires a shared model and a personalized model for each container, the average storage footprint of an ARIMA model and an LSTM model are just 556 KB and 908 KB, respectively. Accordingly, such time and space overhead for the prediction of container resource consumption is practically acceptable.

## 5.3 Effectiveness of *Tetris*

**Prototype Experiments on Amazon EC2.** We first evaluate the effectiveness of *Tetris* in a 60-container K8s cluster by setting the *timeslot*<sup>7</sup> as 5 minutes. As shown in Fig. 9, *Tetris* achieves a lower load imbalance degree by 74.4% – 77.8% and 53.0% – 59.6% compared with Sandpiper and Metis<sup>+</sup>, respectively. Moreover, the migration cost with *Tetris* is 7.8% – 79.5% and 10.3% – 77.0% lower than that with Sandpiper and Metis<sup>+</sup>, respectively. This is because (1) *Tetris* jointly optimizes the load imbalance degree and migration cost, while Sandpiper greedily migrates containers to alleviate the server hotspot without considering the negative impact of migration cost, which causes 37 invalid migrations in total. (2) Though Metis<sup>+</sup> leverages the RL method to explicitly consider the migration cost, it sacrifices the container performance for a certain number of timeslots to optimize container scheduling over the entire time period (*i.e.*, *infinite future*), bringing 12 invalid migrations for 5 timeslots. In contrast, *Tetris* continuously optimizes each timeslot (*i.e.*, 5 minutes) using a certain sliding time window.

To illustrate the effectiveness of *Tetris*, we further take a close look at the scheduling decisions made by the three strategies. As shown in Fig. 10, we observe that *Tetris*

<sup>7</sup> The timeslot can be determined by the cluster administrator according to the cluster requirements.

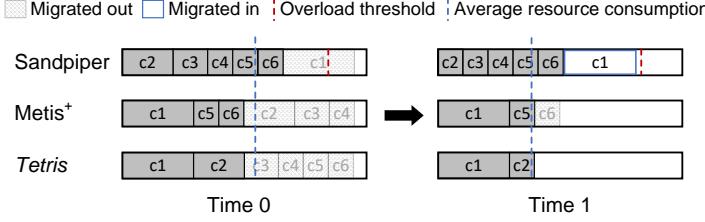


Fig. 10: Comparison of container scheduling decisions on server  $s_1$  of the 60-container K8s cluster in the timeslot  $[0, 1]$  under *Tetris*, *Sandpiper*, and *Metis<sup>+</sup>* strategies.

TABLE 3: Cumulative number of SLO violations over time in our prototype experiments on a 60-container K8s cluster.

Duration (timeslots)	0	1	2	3	4	5
<i>Sandpiper</i>	0	0	14	22	26	61
<i>Metis<sup>+</sup></i>	0	0	0	8	8	20
<i>Tetris</i>	0	0	0	0	0	0

achieves the most load balancing than the other two strategies at both time 0 and time 1, by judiciously migrating out four small containers (*i.e.*,  $c_3 - c_6$ ) at time 0. In comparison, *Sandpiper* causes an invalid migration of container  $c_1$ . As server  $s_1$  is overloaded at time 0 by setting the overload threshold as 85%, *Sandpiper* chooses to migrate container  $c_1$  (*i.e.*, the container with the maximum *volume*) to server  $s_2$ . It then migrates container  $c_1$  back to server  $s_1$  as server  $s_2$  is overloaded at time 1, resulting in poor load balancing and large migration cost. As for *Metis<sup>+</sup>*, it migrates out containers (*i.e.*,  $c_2 - c_4$ ) and container  $c_6$  at time 0 and time 1, respectively, because it adopts online training of the RL model, which has not been trained well enough over the initial time  $[0, 1]$ . In addition, *Metis<sup>+</sup>* mainly considers the optimization over the infinite future, which is likely to increase the load imbalance degree or migration cost over a certain time period (*e.g.*, time  $[0, 1]$ ).

We proceed to examine whether *Tetris* can guarantee the SLO of workloads. For simplicity, we consider the containers hosted on the *overloaded* servers as SLO violations. As shown in Table 3, *Tetris* causes zero SLO violations, while *Sandpiper* and *Metis<sup>+</sup>* can cause up to 35 and 12 SLO violations, respectively, within one timeslot (*i.e.*, at time 5). The reason is that *Tetris* explicitly considers the hard constraints (*i.e.*, Constraints (10)–(11)) on the CPU and memory resource capacities of servers, so as to guarantee the SLO of workloads for each timeslot. However, *Sandpiper* greedily migrates the container with the largest *volume* to the server with the lightest load. It does not check whether the server load exceeds the capacity on the migration destination server, causing unexpected SLO violations to containers. The RL method in *Metis<sup>+</sup>* is not designed to fully guarantee Constraints (10)–(11). It allows the occurrence of SLO violations over a certain time period to optimize container scheduling (*i.e.*, maximize the action reward) over the infinite future.

**Large-scale Simulations Driven by Alibaba Cluster Trace.** We next evaluate the effectiveness of *Tetris* with real-world trace-driven simulations, by setting the *timeslot* as 1 hour. As shown in Fig. 11, *Tetris* can reduce the load imbalance degree of the cluster by 9.7% – 28.6% compared with *Sandpiper* and *Metis<sup>+</sup>*. It can also achieve the lowest

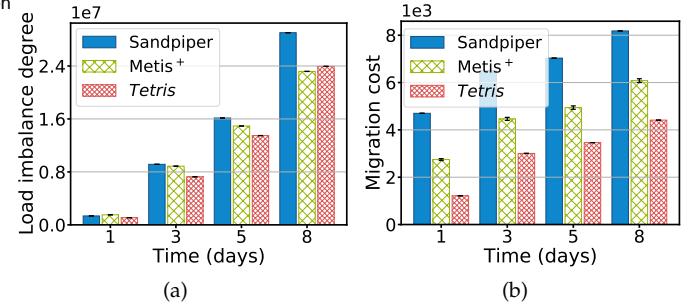


Fig. 11: Performance comparison of *Tetris* with *Sandpiper* and *Metis<sup>+</sup>* strategies under simulations on Alibaba cluster trace, in terms of the cumulative values of (a) cluster load imbalance degree and (b) migration cost of containers over 8 days (the timeslot is 1 hour).

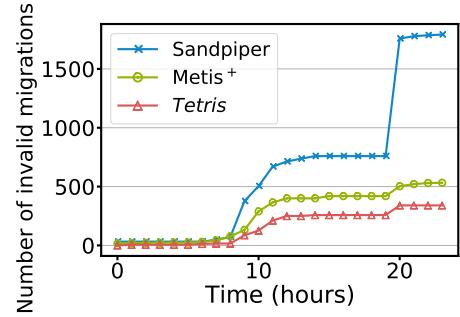


Fig. 12: Cumulative number of invalid migrations of *Tetris* compared with that of *Sandpiper* and *Metis<sup>+</sup>* strategies under simulations on a 24-hour Alibaba cluster trace.

migration cost, which is reduced by 27.4% – 74.1% than the other two strategies. The simulation results are consistent with our prototype experiments. This is because *Tetris* co-optimizes the load balancing and migration cost to alleviate invalid migrations. In contrast, *Sandpiper* greedily alleviates the heavy server load using the worst-fit algorithm, resulting in a number of invalid migrations and a high growth rate of migration cost as depicted in Fig. 11(b). Though *Metis<sup>+</sup>* can optimize the load imbalance degree and migration cost over the infinite future, it achieves a worse load balance degree in the early stage (*i.e.*, the first 5 days) and better load balance degree in the late stage (*i.e.*, day 8) compared with *Tetris* as shown in Fig. 11(a). This is because it requires continuous online training until the RL model converges (after day 5), which is likely to cause invalid migrations in the early stages.

To verify whether *Tetris* can alleviate invalid migrations, we record the cumulative number of invalid migrations under the three scheduling strategies in our simulation. As shown in Fig. 12, *Tetris* reduces the number of invalid migrations by 59.4% – 81.3% and 32.4% – 78.9% as compared with *Sandpiper* and *Metis<sup>+</sup>*, respectively. *Tetris* and *Metis<sup>+</sup>* can reduce the number of invalid migrations, simply because they both consider long-term co-optimization of load balancing and migration cost for container scheduling. In comparison, *Sandpiper* only considers short-term optimization of container scheduling, resulting in a significant number of invalid migrations. As *Metis<sup>+</sup>* achieves long-term container scheduling over the infinite future, which can cause more invalid migrations than *Tetris* over a certain number of timeslots. In particular, *Sandpiper* invokes a significant number of invalid migrations at time 20, because

TABLE 4: Cumulative number of SLO violations over time in a simulated 67,437-container cluster driven by Alibaba cluster trace.

Duration (days)	1	3	5	8
Sandpiper	881	3,512	12,248	14,403
Metis <sup>+</sup>	0	220	220	918
Tetris	0	0	0	0

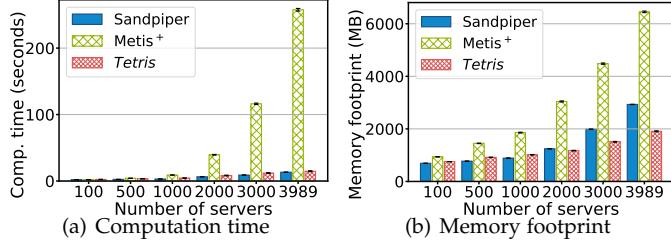


Fig. 13: Runtime overhead of *Tetris* compared with that of Sandpiper and Metis<sup>+</sup> strategies for each timeslot (*i.e.*, 1 hour) on a 24-hour Alibaba cluster trace, by varying the number of servers from 100 to 3,989 with each hosting 17 containers.

the resource consumption of containers varies greatly which overloads a number of servers at time 20.

Similar to our prototype experiments on Amazon EC2, we examine whether *Tetris* can guarantee the SLO of workloads in large-scale simulations. As shown in Table 4, we observe that *Tetris* does not cause any SLO violations for 8 days, while both Sandpiper and Metis<sup>+</sup> can cause up to 14,403 and 918 SLO violations, respectively, which account for 21.4% and 1.4% of the total number of containers. Accordingly, our simulation results are consistent with the prototype experiments.

#### 5.4 Runtime Overhead of *Tetris*

As shown in Fig. 13(a), we observe that the runtime overhead of the three strategies increases *quadratically* with the cluster scale, while Sandpiper and *Tetris* achieve much lower overhead than Metis<sup>+</sup>. The rationale is that the search space of Metis<sup>+</sup> contains all possible mappings of containers to servers, while Sandpiper and *Tetris* only consider a subset of servers and containers, which are the overloaded servers (containers) and the containers on the migration source and destination sets, respectively. Though the time complexity of *Tetris* is in the order of  $\mathcal{O}(|\mathcal{N}| \cdot |\mathcal{M}|)$  (as discussed in Sec. 4.3), the *quadratic term ratio* (obtained using polynomial regression) of *Tetris* is only  $1.9e - 07$ , which is much smaller than that of Metis<sup>+</sup> (*i.e.*,  $2.4e - 05$ ) and Sandpiper (*i.e.*,  $5.5e - 07$ ). Similarly, *Tetris* consumes a much smaller amount of memory than Metis<sup>+</sup> and Sandpiper as shown in Fig. 13(b). This is because Metis<sup>+</sup> stores the policy of RL which involves all states and actions as well as the large parameters of neural networks. Sandpiper requires storing the *volumes* of all containers and servers. In comparison, *Tetris* only needs to calculate the load imbalance degree of all servers and the migration cost of a subset of containers (*i.e.*, containers to be migrated).

In addition, we illustrate the computation overhead of three strategies over various time scales (*i.e.*, from 1 to 8 days) in Table 5. The computation time of *Tetris* is slightly longer than that of Sandpiper (*i.e.*, 0.7 – 4.8 minutes), while

TABLE 5: Cumulative computation overhead (in minutes) of *Tetris*, Sandpiper and Metis<sup>+</sup> scheduling algorithms executed on a 8-day Alibaba cluster trace.

Duration (days)	1	3	5	8
Sandpiper	4.5	17.6	28.5	46.7
Metis <sup>+</sup>	86.1	333.6	503.2	827.1
Tetris	6.5	19.8	29.2	51.5

much shorter than that of Metis<sup>+</sup> (*i.e.*, 79.6 – 775.6 minutes) as the time scale increases. This is because Sandpiper triggers the fewest algorithm executions only when the resource hotspot occurs. Metis<sup>+</sup> greedily selects the largest-reward action over all state spaces in each timeslot, and it also requires online training to update the RL model. In contrast, *Tetris* periodically triggers the execution of Alg. 1 in each timeslot, and such time overhead is much smaller than Metis<sup>+</sup> per timeslot as depicted in Fig. 13(a). As a result, the runtime overhead of *Tetris* is practically acceptable.

## 6 RELATED WORK

**Short-term optimization of workload scheduling.** There have been a number of works devoted to scheduling VMs or containers by considering the *short-term* (*e.g.*, the current or upcoming timeslot) benefits. For example, with the aim of achieving VM consolidations on servers to save energy [31], Dabbagh et al. [32] minimize the energy consumption of migrations by leveraging the Wiener filter method to predict resource consumption. UAHS [33] leverages Bayesian optimization to maximize the CPU utilization of servers. To achieve load balancing of clusters, Sandpiper [8] performs the worst fitting algorithm using the server load index based on CPU, memory, and network resources. Lv et al. [21] propose a communication-aware worst-fit decreasing container placement algorithm to optimize the initial container placement. To particularly optimize the scheduling of long-running workloads, Medea [2] explicitly considers several constraints like affinity and cardinality of containers, while TOPOSCH [34] leverages a prediction-based vertical auto-scaling technique to avoid Quality of Service (QoS) violations and balance the performance of batch jobs. A recent container scheduling method [35] is proposed to select a placement server for each container that is best suited for satisfying the network SLO in K8s. ATCM [36] migrates containers among the cloud servers to minimize the migration cost while maintaining the degree of load balance in a short time scale. Such *short-term* optimizations of workload scheduling above can cause invalid migrations as illustrated in Sec. 2.2. In contrast, *Tetris* leverages the MPC approach to achieve the *long-term* optimization of container scheduling and jointly optimize load balancing and migration cost for long-running workloads.

**RL-based workload scheduling.** Two recent works (*i.e.*, Metis [3], George [5]) design RL algorithms to obtain efficient container placement plans for long-running applications, by modeling the reward as an indicator of container performance and constraint violations. Mondal et al. [4] schedule time-varying workloads to servers based on deep reinforcement learning (DRL) to improve the cluster utilization. To explicitly consider the *long-term* optimization

of workload scheduling, Megh [37] leverages an online RL algorithm to perform VM migrations at each timeslot to minimize the energy consumption and avoid SLO violations. A-SARSA [38] dynamically scales the number of containers based on RL methods to reduce SLO violations. It further adjusts the number of actions corresponding to each state to avoid repeated scheduling. As evidenced by Sec. 2.3 and Sec. 5.3, RL-based scheduling methods can cause unexpected SLO violations by considering the long-term optimization over the infinite future. Also, as shown in Sec. 5.4, RL methods have high time complexity and memory consumption even after the dimensionality reduction, which requires efforts to be deployed in large-scale clusters. To reduce the computation complexity to minutes, *Tetris* leverages the MPC approach to obtain a sub-optimal solution over a certain and sliding time window, which is sufficient for our requirements of *long-term* container scheduling in production clusters.

## 7 CONCLUSION AND FUTURE WORK

This paper presents the design and implementation of *Tetris*, an MPC-based container scheduling strategy to optimize the migration of long-running workloads for achieving load balancing of clusters. By devising a discrete-time dynamic model of shared containerized clusters, *Tetris* judiciously identifies the container scheduling decisions (*i.e.*, *which and where* to migrate) for each timeslot using the Monte Carlo method, with the aim of jointly optimizing cluster load balancing and migration cost of containers. We implement a prototype of *Tetris* and conduct prototype experiments on Amazon EC2 and large-scale simulations driven by Alibaba cluster trace v2018. Our experiment results demonstrate that *Tetris* can improve the cluster load balancing degree by up to 77.8%, while reducing the migration cost by up to 79.5% with acceptable runtime overhead, compared with the state-of-the-art container scheduling strategies.

As our future work, we plan to extend our discrete-time dynamic model by incorporating more types of server resources such as network and disk I/O resources. We also plan to deploy *Tetris* in large-scale production clusters and examine the effectiveness and scalability of *Tetris* container scheduling strategy.

## REFERENCES

- [1] M. Brooker, M. Danilov, C. Greenwood, and P. Piwonka, "On-demand Container Loading in AWS Lambda," in *Proc. of USENIX ATC*, Jul. 2023, pp. 315–328.
- [2] P. Garefalakis, K. Karanasos, P. Pietzuch, A. Suresh, and S. Rao, "Medea: Scheduling of Long Running Applications in Shared Production Clusters," in *Proc. of Eurosys*, Apr. 2018, pp. 1–13.
- [3] L. Wang, Q. Weng, W. Wang, C. Chen, and B. Li, "Metis: Learning to Schedule Long-Running Applications in Shared Container Clusters at Scale," in *Proc. of IEEE SC*, Nov. 2020, pp. 1–17.
- [4] S. S. Mondal, N. Sheoran, and S. Mitra, "Scheduling of Time-Varying Workloads Using Reinforcement Learning," in *Proc. of AAAI*, vol. 35, no. 10, Feb. 2021, pp. 9000–9008.
- [5] S. Li, L. Wang, W. Wang, Y. Yu, and B. Li, "George: Learning to Place Long-Lived Containers in Large Clusters with Operation Constraints," in *Proc. of ACM SOCC*, Nov. 2021, pp. 258–272.
- [6] A. Verma, L. Pedrosa, M. Korupolu, D. Oppenheimer, E. Tune, and J. Wilkes, "Large-Scale Cluster Management at Google with Borg," in *Proc. of Eurosys*, Apr. 2015, pp. 1–17.
- [7] D. Bernstein, "Containers and Cloud: From LXC to Docker to Kubernetes," *IEEE Cloud Computing*, vol. 1, no. 3, pp. 81–84, 2014.
- [8] T. Wood, P. Shenoy, A. Venkataramani, and M. Yousif, "Sandpiper: Black-box and Gray-box Resource Management for Virtual Machines," *Computer Networks*, vol. 53, no. 17, pp. 2923–2938, 2009.
- [9] Alibaba, (2023, Sep.) Alibaba Cluster Trace Program. [Online]. Available: <https://github.com/alibaba/clusterdata/>
- [10] A. Beloglazov, J. Abawajy, and R. Buyya, "Energy-Aware Resource Allocation Heuristics for Efficient Management of Data Centers for Cloud Computing," *Future Generation Computer Systems*, vol. 28, no. 5, pp. 755–768, 2012.
- [11] M. Xu, W. Tian, and R. Buyya, "A Survey on Load Balancing Algorithms for Virtual Machines Placement in Cloud Computing," *Concurrency and Computation: Practice and Experience*, vol. 29, no. 12, p. e4123, 2017.
- [12] H. Qiu, W. Mao, C. Wang, H. Franke, A. Youssef, Z. T. Kalbarczyk, T. Başar, and R. K. Iyer, "AWARE: Automate Workload Autoscaling with Reinforcement Learning in Production Cloud Systems," in *Proc. of USENIX ATC*, Jul. 2023, pp. 387–402.
- [13] M. Gaggero and L. Caviglione, "Model Predictive Control for the Placement of Virtual Machines in Cloud Computing Applications," in *Proc. of ACC*, Jul. 2016, pp. 1987–1992.
- [14] ———, "Model Predictive Control for Energy-Efficient, Quality-Aware, and Secure Virtual Machine Placement," *IEEE Transactions on Automation Science and Engineering*, vol. 16, no. 1, pp. 420–432, 2018.
- [15] C. E. Garcia, D. M. Pretti, and M. Morari, "Model Predictive Control: Theory and Practice – A Survey," *Automatica*, vol. 25, no. 3, pp. 335–348, 1989.
- [16] N. Metropolis and S. Ulam, "The Monte Carlo Method," *Journal of the American Statistical Association*, vol. 44, no. 247, pp. 335–341, 1949.
- [17] H. Wu, W. Zhang, Y. Xu, H. Xiang, T. Huang, H. Ding, and Z. Zhang, "Aladdin: Optimized Maximum Flow Management for Shared Production Clusters," in *Proc. of IEEE IPDPS*, May 2019, pp. 696–707.
- [18] J. Guo, Z. Chang, S. Wang, H. Ding, Y. Feng, L. Mao, and Y. Bao, "Who Limits the Resource Efficiency of My Datacenter: An Analysis of Alibaba Datacenter Traces," in *Proc. of IEEE/ACM IWQoS*, Jun. 2019, pp. 1–10.
- [19] M. Xu, C. Song, H. Wu, S. S. Gill, K. Ye, and C. Xu, "esDNN: Deep Neural Network based Multivariate Workload Prediction in Cloud Computing Environments," *ACM Transactions on Internet Technology*, vol. 22, no. 3, pp. 1–24, 2022.
- [20] Y. Zhang, L. Jiao, J. Yan, and X. Lin, "Dynamic Service Placement for Virtual Reality Group Gaming on Mobile Edge Cloudlets," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 8, pp. 1881–1897, 2019.
- [21] L. Lv, Y. Zhang, Y. Li, K. Xu, D. Wang, W. Wang, M. Li, X. Cao, and Q. Liang, "Communication-Aware Container Placement and Reassignment in Large-Scale Internet Data Centers," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 3, pp. 540–555, 2019.
- [22] E. S. Hou, N. Ansari, and H. Ren, "A Genetic Algorithm for Multi-processor Scheduling," *IEEE Transactions on Parallel and Distributed Systems*, vol. 5, no. 2, pp. 113–120, 1994.
- [23] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Hardness*. San Fransisco: W. H. Freeman, 1979.
- [24] G. E. Box, G. M. Jenkins, G. C. Reinsel, and G. M. Ljung, *Time Series Analysis: Forecasting and Control*. John Wiley & Sons, 2015.
- [25] Z. Li and D. Hoiem, "Learning without Forgetting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 2935–2947, 2017.
- [26] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [27] F. Li and B. Hu, "Deepjs: Job Scheduling Based on Deep Reinforcement Learning in Cloud Data Center," in *Proc. of ACM ICBDC*, May 2019, pp. 48–53.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proc. of CVPR*, Jul. 2016, pp. 770–778.
- [29] A. De Myttenaere, B. Golden, B. Le Grand, and F. Rossi, "Mean Absolute Percentage Error for Regression Models," *Neurocomputing*, vol. 192, pp. 38–48, 2016.
- [30] D. Fan, H. Sun, J. Yao, K. Zhang, X. Yan, and Z. Sun, "Well Production Forecasting Based on ARIMA-LSTM Model Considering Manual Operations," *Energy*, vol. 220, pp. 1–13, 2021.
- [31] Z. Á. Mann, "Allocation of Virtual Machines in Cloud Data Centers—A Survey of Problem Models and Optimization Algorithms," *ACM Computing Surveys*, vol. 48, no. 1, pp. 1–34, 2015.

- [32] M. Dabbagh, B. Hamdaoui, M. Guizani, and A. Rayes, "An Energy-Efficient VM Prediction and Migration Framework for Overcommitted Clouds," *IEEE Transactions on Cloud Computing*, vol. 6, no. 4, pp. 955–966, 2016.
- [33] C. Luo, B. Qiao, X. Chen, P. Zhao, R. Yao, H. Zhang, W. Wu, A. Zhou, and Q. Lin, "Intelligent Virtual Machine Provisioning in Cloud Computing," in *Proc. of IJCAI*, Jan. 2021, pp. 1495–1502.
- [34] J. Zhu, R. Yang, X. Sun, T. Wo, C. Hu, H. Peng, J. Xiao, A. Y. Zomaya, and J. Xu, "QoS-Aware Co-Scheduling for Distributed Long-Running Applications on Shared Clusters," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 12, pp. 4818–4834, 2022.
- [35] E. Kim, K. Lee, and C. Yoo, "Network SLO-Aware Container Scheduling in Kubernetes," *The Journal of Supercomputing*, vol. 79, no. 10, pp. 11 478–11 494, 2023.
- [36] W. Zhang, L. Chen, J. Luo, and J. Liu, "A Two-Stage Container Management in the Cloud for Optimizing the Load Balancing and Migration Cost," *Future Generation Computer Systems*, vol. 135, pp. 303–314, 2022.
- [37] D. Basu, X. Wang, Y. Hong, H. Chen, and S. Bressan, "Learn-as-you-go with Megh: Efficient Live Migration of Virtual Machines," *IEEE Transactions on Parallel and Distributed Systems*, vol. 30, no. 8, pp. 1786–1801, 2019.
- [38] S. Zhang, T. Wu, M. Pan, C. Zhang, and Y. Yu, "A-SARSA: A Predictive Container Auto-Scaling Algorithm Based on Reinforcement Learning," in *Proc. of IEEE ICWS*, Oct. 2020, pp. 489–497.

**Fei Xu** received the B.S., M.E., and Ph.D. degrees in 2007, 2009, and 2014, respectively, all from the Huazhong University of Science and Technology (HUST), Wuhan, China. He received Outstanding Doctoral Dissertation Award in Hubei province, China, and ACM Wuhan & Hubei Computer Society Doctoral Dissertation Award in 2015. He is currently an associate professor with the School of Computer Science and Technology, East China Normal University, Shanghai, China. His research interests include datacenter, virtualization technology, and distributed systems.



**Xiyue Shen** received the master's degree with the School of Computer Science and Technology, East China Normal University, Shanghai, China. Her research interests focus on cloud computing and datacenter resource management.



**Shuhao Lin** is currently working toward the master's degree with the School of Computer Science and Technology, East China Normal University, Shanghai, China. His research interests focus on cloud computing and datacenter resource management.

**Li Chen** received the BEngr degree from the Department of Computer Science and Technology, Huazhong University of Science and Technology, China, in 2012 and the MAsc degree from the Department of Electrical and Computer Engineering, University of Toronto, in 2014 and the PhD degree in computer science and engineering from the Department of Electrical and Computer Engineering, University of Toronto, in 2018. She is currently an assistant professor with the Department of Computer Science, School of Computing and Informatics, University of Louisiana at Lafayette, Lafayette, USA. Her research interests include big data analytics systems, cloud computing, datacenter networking, and resource allocation.



**Zhi Zhou** received the B.S., M.E., and Ph.D. degrees in 2012, 2014, and 2017, respectively, all from the School of Computer Science and Technology at Huazhong University of Science and Technology (HUST), Wuhan, China. He is currently an associate professor in the School of Computer Science and Engineering at Sun Yat-sen University, Guangzhou, China. In 2016, he was a visiting scholar at University of Göttingen. He was nominated for the 2019 CCF Outstanding Doctoral Dissertation Award, the sole recipient of the 2018 ACM Wuhan & Hubei Computer Society Doctoral Dissertation Award, and a recipient of the Best Paper Award of IEEE UIC 2018. His research interests include edge computing, cloud computing, and distributed systems.



**Fen Xiao** received the master's degree from the School of Computer Science and Technology, Wuhan University, Wuhan, China, in 2009. She is currently a senior R&D engineer at Tencent, Shenzhen, China. Her research interests focus on NoSQL and distributed storage systems.



**Fangming Liu** (S'08, M'11, SM'16) received the B.Eng. degree from the Tsinghua University, Beijing, and the Ph.D. degree from the Hong Kong University of Science and Technology, Hong Kong. He is currently a Full Professor with the Huazhong University of Science and Technology, Wuhan, China. His research interests include cloud computing and edge computing, datacenter and green computing, SDN/NFV/5G and applied ML/AI. He received the National Natural Science Fund (NSFC) for Excellent Young Scholars, and the National Program Special Support for Top-Notch Young Professionals. He is a recipient of the Best Paper Award of IEEE/ACM IWQoS 2019, ACM e-Energy 2018 and IEEE GLOBECOM 2011, the First Class Prize of Natural Science of Ministry of Education in China, as well as the Second Class Prize of National Natural Science Award in China.

## APPENDIX

### .1 Mathematical Derivations of Eq. (6) and Eq. (7)

To obtain Eq. (6), the load imbalance degree of CPU resource part in Eq. (1) can first be expanded as

$$\begin{aligned} & \sum_{i \in \mathcal{M}} \left( \sum_{k \in \mathcal{N}} x_i^k(t) \cdot \text{cpu}^k(t) - \overline{\text{CPU}(t)} \right)^2 \\ &= \sum_{i \in \mathcal{M}} \left( \sum_{k \in \mathcal{N}} x_i^k(t) \cdot \text{cpu}^k(t) \right)^2 \\ &\quad - 2 \cdot \overline{\text{CPU}(t)} \cdot \sum_{k \in \mathcal{N}} \text{cpu}^k(t) + |\mathcal{M}| \cdot \overline{\text{CPU}(t)}^2, \end{aligned} \quad (15)$$

where the last two terms of the expanded expression in Eq. (15) are actually constant for each timeslot  $t$ , which can be denoted by  $-|\mathcal{M}| \cdot \overline{\text{CPU}(t)}^2$ . The load imbalance degree of CPU resource part can be mainly determined by the first term of Eq. (15), which is given by

$$\begin{aligned} & \sum_{i \in \mathcal{M}} \left( \sum_{k \in \mathcal{N}} x_i^k(t) \cdot \text{cpu}^k(t) \right)^2 \\ &= 2 \sum_{i \in \mathcal{M}} \sum_{k, l \in \mathcal{N}} x_i^k(t) \text{cpu}^k(t) \cdot x_i^l(t) \text{cpu}^l(t) - \sum_{k \in \mathcal{N}} \text{cpu}^k(t)^2, \end{aligned} \quad (16)$$

where the second term is also a constant at time  $t$ . Accordingly, the load imbalance degree of CPU resource part can be determined by the first term of Eq. (16). The load imbalance degree of memory resource part in Eq. (1) is similar to that of CPU resource part by substituting  $\text{mem}^k(t)$  for  $\text{cpu}^k(t)$ . We further formulate  $C_b(t)$  as

$$\begin{aligned} C_b(t) &= \frac{2}{|\mathcal{M}| - 1} \sum_{i \in \mathcal{M}} \sum_{k, l \in \mathcal{N}} x_i^k(t) x_i^l(t) \cdot (\text{cpu}^k(t) \text{cpu}^l(t) \\ &\quad + \beta \cdot \text{mem}^k(t) \text{mem}^l(t)) + C_1, \end{aligned} \quad (17)$$

where the parameter  $C_1 = -\frac{|\mathcal{M}| \cdot (\overline{\text{CPU}(t)}^2 + \beta \cdot \overline{\text{MEM}(t)}^2)}{|\mathcal{M}| - 1} - \frac{\sum_{k \in \mathcal{N}} (\text{cpu}^k(t)^2 + \beta \cdot \text{mem}^k(t)^2)}{|\mathcal{M}| - 1}$  is actually a constant value, simply because the parameters in  $C_1$  are all constant values at time  $t$ . Accordingly, Eq. (17) is actually the cluster load imbalance degree formulated in Eq. (6).

To obtain Eq. (7), we first incorporate  $m^k(t)$  into Eq. (3),  $C_m(t)$  can be represented by

$$\begin{aligned} C_m(t) &= \sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{M}} C_{mig}^k(t) \cdot x_i^k(t) \\ &\quad - \sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{M}} C_{mig}^k(t) \cdot x_i^k(t) x_i^k(t-1), \end{aligned} \quad (18)$$

where the first term is transformed into a constant  $C_2 = \delta \cdot |\mathcal{N}| + \gamma \cdot \sum_{k \in \mathcal{N}} \text{mem}^k(t)$ , which is related to the memory sum of all containers  $\text{mem}^k(t)$ ,  $\forall k \in \mathcal{N}$ . Accordingly, by incorporating Eq. (4) into Eq. (18),  $C_m(t)$  is determined by the second term of Eq. (18), which is given by

$$C_m(t) = C_2 - \sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{M}} (\delta + \gamma \cdot \text{mem}^k(t)) \cdot x_i^k(t) x_i^k(t-1). \quad (19)$$

As a result, Eq. (19) is actually the migration cost across the cluster formulated in Eq. (7).

### .2 Proof of Theorem 1

*Proof.* The load imbalance degree  $\text{load}_i(t)$  of a server  $i \in \mathcal{M}$  and Eq. (13) – (14) in Theorem 1 contain both CPU and memory resource consumption of containers. To simplify the proof, we first consider CPU resource of containers (*i.e.*,  $\text{cpu}^k(t)$ ,  $\forall k \in \mathcal{N}$ ). By assuming the cluster servers are in the “ideal” load-balanced state with the CPU resource consumption of  $\overline{\text{CPU}(t)}$  at time  $t$ , we can calculate the maximum and minimum values of CPU load imbalance degree of a server (*i.e.*, the sum of pairwise multiplication  $\text{mul}_{n(t)}$  of  $\text{cpu}^k(t)$  on a server in Eq. (12)) in Lemma 1 below.

**Lemma 1.** *Given a server in the “ideal” load-balanced state, which hosts a set of containers with the CPU consumption of  $\text{cpu}^1(t), \text{cpu}^2(t), \dots, \text{cpu}^{n(t)}(t)$  and  $\sum_{k=1}^{n(t)} \text{cpu}^k(t) = \overline{\text{CPU}(t)}$ , where  $n(t) \in [1, |\mathcal{N}|]$ , the sum of pairwise multiplication  $\text{mul}_{n(t)}$  of  $\text{cpu}^k(t)$ ,  $\forall k \in [1, n(t)]$  becomes the maximum (*i.e.*,  $\frac{\overline{\text{CPU}(t)}^2}{2}$ ), if and only if  $\text{cpu}^1(t) = \text{cpu}^2(t) = \dots = \text{cpu}^{n(t)}(t)$  and  $n(t) = |\mathcal{N}|$ . Similarly,  $\text{mul}_{n(t)}$  becomes the minimum (*i.e.*,  $\frac{\overline{\text{CPU}(t)}^2}{2} - \frac{1}{2} \cdot \sum_{k=1}^{n(t)} \text{cpu}^k(t)^2$ ), if and only if  $n(t)$  is the minimum value that satisfies  $\sum_{k=1}^{n(t)} \text{cpu}^k(t) = \overline{\text{CPU}(t)}$ .*

*Proof.* The sum of the pairwise multiplication of  $\text{cpu}^k(t)$ ,  $\forall k \in [1, n(t)]$  is calculated as

$$\begin{aligned} \text{mul}_{n(t)} &= \sum_{k=1}^{n(t)-1} \sum_{l=k+1}^{n(t)} (\text{cpu}^k(t) \cdot \text{cpu}^l(t)) \\ &= \frac{1}{2} \left[ \left( \sum_{k=1}^{n(t)} \text{cpu}^k(t) \right)^2 - \sum_{k=1}^{n(t)} \text{cpu}^k(t)^2 \right] \\ &= \frac{\overline{\text{CPU}(t)}^2}{2} - \frac{\sum_{k=1}^{n(t)} \text{cpu}^k(t)^2}{2}. \end{aligned} \quad (20)$$

The first term of Eq. (20) is a constant value. By Cauchy-Buniakowsky-Schwarz Inequality, the second term  $\frac{1}{2} \cdot \sum_{k=1}^{n(t)} \text{cpu}^k(t)^2$  becomes the minimum (*i.e.*,  $\frac{\overline{\text{CPU}(t)}^2}{2 \cdot n(t)}$ ) if and only if  $\text{cpu}^k(t) = \frac{\overline{\text{CPU}(t)}}{n(t)}$ ,  $\forall k \in [1, n(t)]$ . Accordingly, we have

$$\text{mul}_{n(t)} \leq \left( \frac{1}{2} - \frac{1}{2 \cdot n(t)} \right) \cdot \overline{\text{CPU}(t)}^2, \quad (21)$$

where  $\frac{1}{2} - \frac{1}{2 \cdot n(t)}$  is a monotonically increasing function with  $n(t) \in [1, |\mathcal{N}|]$ . As a result,  $\text{mul}_{n(t)}$  can reach the maximum when  $n(t) = |\mathcal{N}|$ , which is given by

$$\max(\text{mul}_{n(t)}) = \left( \frac{1}{2} - \frac{1}{2 \cdot |\mathcal{N}|} \right) \cdot \overline{\text{CPU}(t)}^2. \quad (22)$$

To calculate the minimum value of  $\text{mul}_{n(t)}$ , we have to find the maximum value of the second term  $\frac{1}{2} \cdot \sum_{k=1}^{n(t)} \text{cpu}^k(t)^2$  of Eq. (20). According to the perfect square trinomial, the square of the sum of a set of non-negative numbers is larger than or equal to the sum of squares of these numbers. Therefore, we find that  $\frac{1}{2} \cdot \sum_{k=1}^{n(t)} \text{cpu}^k(t)^2$  can reach the maximum when  $n(t)$  is the minimum number to satisfy  $\sum_{k=1}^{n(t)} \text{cpu}^k(t) = \overline{\text{CPU}(t)}$ . To find the minimum number of  $n(t)$ , we sort the CPU consumption of containers in descending order and greedily select containers from the

highest CPU consumption to the lowest. In such a situation above, Eq. (20) can reach its minimum accordingly.

□

Similarly, we can calculate the maximum and minimum values of memory load imbalance degree of a server (*i.e.*, the sum of pairwise multiplication  $mul_{n(t)}$  of  $mem^k(t)$  on a server in Eq. (12)). To achieve that, we simply substitute  $mem^k(t)$  for  $cpu^k(t)$  and  $MEM(t)$  for  $CPU(t)$  in Lemma 1. As a result, we combine the CPU and memory resources together as in Eq. (12), and the maximum value of the load imbalance degree of a server in the “ideal” load-balanced state turns out to be Eq. (13). The minimum value of the load imbalance degree of an “ideal” server is calculated as Eq. (14), when such a server hosts the minimum number of containers that satisfies  $\sum_{k=1}^{n(t)} cpu^k(t) \leq \overline{CPU(t)}$  and  $\sum_{k=1}^{n(t)} mem^k(t) \leq \overline{MEM(t)}$  (*i.e.*, by jointly considering CPU and memory resources) in practice.

Finally, to achieve load balancing, we set the threshold of  $load_i(t)$  for migration source server (*i.e.*,  $load_{src}(t)$ ) as the maximum value of the load imbalance degree of an “ideal” server in Eq. (13). Meanwhile, we set the threshold of  $load_i(t)$  for migration destination server (*i.e.*,  $load_{dest}(t)$ ) as the minimum value of the load imbalance degree of an “ideal” server in Eq. (14).

□