Tetris: Proactive Container Scheduling for Long-Term Load Balancing in Shared Clusters

Fei Xu[®], Member, IEEE, Xiyue Shen, Shuohao Lin, Li Chen[®], Member, IEEE, Zhi Zhou[®], Member, IEEE, Fen Xiao, and 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 among 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 proactively 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.

Manuscript received 5 September 2023; revised 17 July 2024; accepted 5 August 2024. Date of publication 13 August 2024; date of current version 9 October 2024. This work was supported in part by the NSFC under Grant 62372184, in part by the Science and Technology Commission of Shanghai Municipality under Grant 22DZ2229004, and a grant from the Tencent Corporation. The work of Li Chen was supported in part by the NSF under Grant OIA-2019511 and Grant OIA-2327452. The work of Zhi Zhou was supported in part by the National Key Research & Development (R&D) Plan under Grant 2022YFB4501703. The work of Fangming Liu was supported in part by the Major Key Project of PCL under Grant PCL2022A05. (Corresponding author: Fei Xu.)

Fei Xu, Xiyue Shen, and Shuohao Lin are with the Shanghai Key Laboratory of Multidimensional Information Processing, School of Computer Science and Technology, East China Normal University, Shanghai 200062, China (e-mail: fxu@cs.ecnu.edu.cn).

Li Chen is with the School of Computing and Informatics, University of Louisiana at Lafayette, 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, Guangzhou 510006, China (e-mail: zhouzhi9@mail.sysu.edu.cn).

Fen Xiao is with Tencent Incorporation, Shenzhen 518054, China (e-mail: cherryxiao@tencent.com).

Fangming Liu is with the Peng Cheng Laboratory, Huazhong University of Science and Technology, Wuhan 430074, China (e-mail: fangminghk@gmail.com).

This article has supplementary downloadable material available at $\frac{1}{100}$ https://doi.org/10.1109/TSC.2024.3442544, provided by the authors.

Digital Object Identifier 10.1109/TSC.2024.3442544

Index Terms—Container scheduling, load balancing, long-running containerized workloads, migration cost.

I. INTRODUCTION

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

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 among servers over a short time window) in shared production clusters. As evidenced by our motivational analysis (in Section II-B) on both Alibaba cluster traces v2018 and v2022 [9], more than half of container migrations are invalid migrations as the time window size is set as 3 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 Section II-B reveal that 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 longrunning workloads, which is actually the *root cause of invalid migrations* as discussed in Section II-B. There have also been

1939-1374 © 2024 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

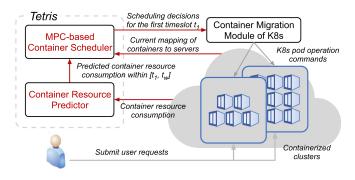


Fig. 1. Overview of *Tetris*. It comprises two pieces of modules including a *Container Resource Predictor* and an *MPC-based Container Scheduler*. By *jointly* optimizing the cluster load balancing degree and container migration cost, *Tetris* decides the appropriate container scheduling plans over a time window $[t_1, t_W]$, which are performed by the *Container Migration Module of K8s* in the containerized cluster.

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 schedule*) over the entire time period (i.e., *infinite future*). The containers to be scheduled (i.e., *which container to schedule*) 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 Section II-C. As a result, scant 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* shown in Fig. 1, a model predictive control [15] (MPC)-based container scheduling strategy to proactively make scheduling 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 Section V-B) 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 (Section III). 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

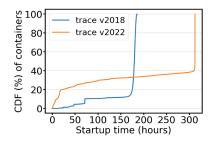


Fig. 2. Container uptime in Alibaba cluster traces v2018 and v2022 [9]. More than 80% of containers last for over 175 hours in the v2018 trace, while around 60% of containers run over 312 hours in the v2022 trace.

optimization problem based on MPC, by jointly considering the cluster load imbalance degree and migration cost of containers.

Description Second, we design Tetris to achieve long-term scheduling optimization for cluster load balancing (Section IV). Specifically, the Container Resource Predictor in Tetris first accurately predicts the container resource consumption (Section IV-A). The MPC-based Container Scheduler then leverages the Monte Carlo method [16] to judiciously identify the container scheduling decisions for each timeslot, by minimizing our formulated cost function of container scheduling (Section IV-B). To particularly reduce the complexity of Tetris container scheduling strategy, we calculate the thresholds 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 (Section V). Compared with the conventional scheduling strategy (i.e., Sandpiper [8]) and the state-of-the-art RL-based method (i.e., Metis⁺, a modified version of Metis [3]), *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.

II. BACKGROUND AND MOTIVATION

A. Long-Running Containerized Workloads

Large-scale shared production clusters often host many long-running containers in response to latency-critical user requests [17]. Taking Alibaba as an example, its online services are mainly long-running and *stateful* containerized applications, such as online shopping, database, and machine learning [18]. As shown in Fig. 2, more than half of the containers are executed over 175 hours in the Alibaba production clusters. Additionally, the stateful (including partially stateful) workloads account for over 50% in the latest measurement study in Microsoft production clusters [19]. Due to the unpredictability of user requests, such long-running workloads place diverse demands on multidimensional container resources including CPU, GPU, memory, network and disk I/O, which are characterized by dynamic variability and uncertainty in the resource consumption [20].

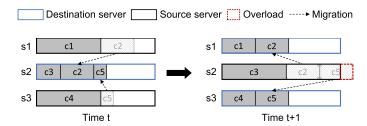


Fig. 3. Illustration of invalid migration. Containers c2 and c5 are migrated back and forth among servers s1, s2, and s3.

Accordingly, the co-location of containerized workloads can lead to contention for shared resources, which can easily cause performance interference and resource wastage in the shared clusters.

Based on the above, enabling adequate (re)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 *rescheduling*¹ (i.e., container migrations [21]). Currently, CRIU (Checkpoint/Restore In Userspace)² supports the *pre-copy*-based *live* migration of containers [22]. To simplify the container migration process, we *simulate* it with the following three steps: i.e., stopping containerized services on the source server first, then iteratively replicating the memory dirty pages to restore services, and finally restarting the services on the destination server. In general, the container scheduling inevitably brings a noticeable amount of migration cost, which requires to be considered during the scheduling of containers, especially for *stateful* workloads.

B. 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. Such a *short-term* optimization of container scheduling policy can easily 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. 3 as an example. The cluster scheduler (e.g., the default K8s scheduler) first migrates container c2 and c5 from server s1 and s3, respectively, to server s2, in order to obtain a *temporary* benefits of load balancing at time t. Unfortunately, server s2 becomes overloaded as the resource consumption of three containers (i.e., c3, c2, c5) dramatically increases, which triggers two container migrations (i.e., container c2, c5) from server s2 to s1 and s3, respectively. In such a case, naive scheduling policies are likely to make long-running containers migrated back and forth among servers, thereby causing heavy and unnecessary migration cost to containers. Accordingly, we formally define *invalid migrations* as in Definition 1.

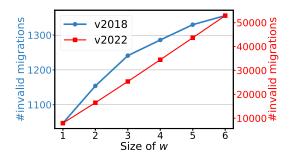


Fig. 4. Number of invalid migrations within a time window size varying from 1 to 6 hours for 67,437 containers in Alibaba cluster trace v2018 and 676,938 containers in Alibaba cluster trace v2022.

Definition 1: If a container is migrated back and forth among servers over a 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 both Alibaba cluster traces v2018 (i.e., an 8-day period) and v2022 (i.e., a 13-day period) [9]. As shown in Fig. 4, the number of invalid migrations increases as the scheduling time window increases. Specifically, the cumulative number of invalid migrations exceeds 1,200 and 30,000 for the v2018 and v2022 traces, respectively, as w is set as 3 hours. The number reaches over 50,000, accounting for over 75% of container migrations for the v2022 trace when w increases to 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 (partially stateful workloads). Our experiment results reveal that 61.7% of the containers are affected by invalid migrations. 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 future resource consumption of containers.

C. An Illustrative Example

To avoid invalid migrations, we design *Tetris* in Section IV, a simple yet effective container scheduling strategy that leverages the MPC approach to proactively 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 [23]. To illustrate how Tetris works, we show a motivation example in Fig. 5 by comparing Tetris with an RL-based scheduling method (i.e., Metis⁺, a modified version of Metis [3] which will be introduced in Section V-A). Though Metis⁺ 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 s2 at time t+2 and trigger the migration of container c3 from server s2 to s1. 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

¹Descheduler for Kubernetes: https://github.com/kubernetes-sigs/descheduler

²CRIU: https://criu.org/Main_Page

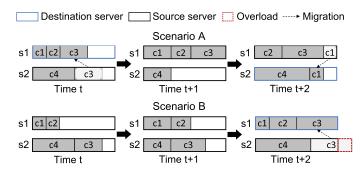


Fig. 5. An illustrative example of MPC-based container scheduling (our proposed *Tetris* in scenario A), as compared with RL-based container scheduling (Metis⁺ [3] in scenario B).

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 Section V-D.

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.

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

A. 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 $t = \{t_0, t_1, \ldots, t_W\}$ 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.

Cluster load imbalance degree: 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

 ${\it TABLE I} \\ {\it Key Notations in Our Discrete-Time Dynamic Model At Time } t \\$

Notation	Definition				
\mathcal{M}, \mathcal{N}	Sets of servers and containers				
\overline{W}	Time window size				
$C_b(t)$	Load imbalance degree of the cluster				
$C_m(t)$	Container migration cost of the cluster				
$cpu^k(t)$	CPU consumption of a container k				
$mem^k(t)$	Memory consumption of a container k				
$\overline{CPU(t)}$	Average CPU consumption of a server				
$\overline{MEM(t)}$	Average memory consumption of a server				
α	Normalized model coefficient of migration cost to load imbalance degree				
β	Normalized model coefficient of memory resources to CPU resources				
γ , δ	Linear model coefficients of container migration cost in terms of the memory size				
$x_i^k(t)$	Indicator if a container k is on a server i				
$m^k(t)$	Indicator whether a container k is migrated				

imbalance degrees of each resource of servers, which is given by

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

$$(1)$$

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

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

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.

Container migration cost: We use the sum of unit cost of the migrated containers $C^k_{mig}(t)$ to denote the migration cost $C_m(t)$, which is formulated as

$$C_m(t) = \sum_{k \in \mathcal{N}} C_{mig}^k(t) \cdot m^k(t), \tag{2}$$

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. As elaborated in Section II-A, the migration cost of containers is *roughly linear* to their memory footprint $mem^k(t)$ [24]. Accordingly, $C_{mig}^k(t)$ can be given by

$$C_{mig}^{k}(t) = \delta + \gamma \cdot mem^{k}(t), \tag{3}$$

where γ and δ are linear model coefficients, which are empirical and constant values. We elaborate the configuration process of model coefficients in Appendix C, available online.

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,

$$C(t) = C_b(t) + \alpha \cdot C_m(t), \tag{4}$$

where $\alpha \in [0,1]$ denotes the normalized parameter of migration cost $C_m(t)$ to cluster load imbalance degree $C_b(t)$, which can practically be obtained by the trace analysis. By incorporating $\overline{CPU(t)}$ and $\overline{MEM(t)}$ into (1), and (3) and $m^k(t)$ into (2), we yield $C_b(t)$ and $C_m(t)$, respectively, which are formulated as

$$C_b(t) = \frac{2}{|\mathcal{M}| - 1} \sum_{i \in \mathcal{M}} \sum_{k,l \in \mathcal{N}} x_i^k(t) x_i^l(t) \cdot \left(cpu^k(t) cpu^l(t) \right)$$
(5)

$$+\beta \cdot mem^k(t)mem^l(t) + C_1,$$

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

(6)

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, available online. By analyzing the formulations above, (5) 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. Equation (6) implies that the migration cost across the cluster can be determined by the negative sum of migration costs of the migrated containers.

B. 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 (7)). We assume the scheduling within the time window $t \in [t_1, t_W]$, and our optimization problem can be formulated as

$$\min_{x_i^k(t)} \sum_{t=t_1}^{t_W} C(t) \tag{7}$$

s.t.
$$\sum_{i \in M} x_i^k(t) = 1, \quad \forall k \in \mathcal{N}$$
 (8)

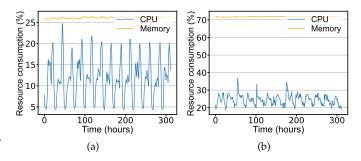


Fig. 6. Resource consumption of containers with respect to the physical machine resources over time: (a) a periodic container (id = 404), and (b) an aperiodic container (id = 22) in Alibaba cluster trace v2022.

$$\sum_{k \in \mathcal{N}} x_i^k(t) \cdot cpu^k(t) \le CPU_i^{cap}(t), \quad \forall i \in \mathcal{M}$$
 (9)

$$\sum_{k \in \mathcal{N}} x_i^k(t) \cdot mem^k(t) \le MEM_i^{cap}(t), \forall i \in \mathcal{M}$$
 (10)

where the hard constraints in our optimization problem are as follows. Constraint (8) limits each container to be hosted only on one server per time to avoid scheduling conflicts. Constraints (9)–(10) ensure the total CPU and memory resource consumption on each server cannot exceed the server resource capacity (i.e., $CPU_i^{cap}(t)$ and $MEM_i^{cap}(t)$).

Problem analysis: For each timeslot, the problem defined in (7) can be easily reduced to a multiprocessor scheduling problem (MSP) [25]. 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 [26]. Moreover, (7) 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.

IV. TETRIS DESIGN

Based on our discrete-time dynamic model defined in Section III, this section designs *Tetris*, including a container resource predictor (Section IV-A) and an MPC-based container scheduler (Section IV-B), with the aim of jointly optimizing the cluster load balancing degree and container migration cost over a certain time window.

A. 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 [27]. Specifically, we first obtain the *frequency* of container resource consumption values and calculate its corresponding *periodicity*. With such a value, we then slice the container resource consumption into several segments. We finally calculate the Pearson correlation coefficients between every

two consecutive segments and decide whether such coefficients exceed a given periodicity threshold thr_{period} (e.g., 0.85). If the correlation coefficients exceed thr_{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. Specifically, we adopt ARIMA [28] to predict the resource consumption of periodic containers. 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. In addition, we leverage a simple yet effective LSTM [29] model to predict the resource consumption of aperiodic containers, rather than a more sophisticated attention-based model [30]. We adopt a 4-layer stacked LSTM model with 50 neurons per layer and add a Dropout layer between every two LSTM layers to reduce overfitting. To alleviate the impact of error propagation, we directly conduct the multistep-ahead predictions and set the number of prediction steps as W. In particular, we use the historical container resource consumption to train the shared ARIMA and LSTM models for newly-launched containers. To improve the prediction accuracy, we further leverage incremental learning [31] to update a personalized model for each container periodically.

B. MPC-Based Container Scheduling

We design *Tetris* in Algorithm 1 by leveraging the MPC approach to make container scheduling decisions over the time window 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 Section III-B to obtain $\mathbf{X}_t, \forall t \in [t_1, t_W]$, at the current time t_0 . Then, we proactively 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 to maintain the load balancing of the containerized cluster.

In more detail, we solve the scheduling optimization problem using the Monte Carlo method [16] as discussed in Section III-B. 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 Algorithm 1 including server classification (lines 3 to 9) and container scheduling (lines 10 to 18), 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., (5)) defined in our dynamic model, we can obtain the simplified load imbalance degree (i.e., $load_i(t)$) for each server i given

Algorithm 1: *Tetris:* Proactive MPC-Based Container Scheduling Algorithm for Long-Term Load Balancing.

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.

```
1 for all z \in [1, Z] do
       for all t \in [t_1, t_W] do
            // Phase I: server classification
           Initialize: the load threshold of migration
3
             source servers load_{src}(t) \leftarrow Eq. (12) and that
             of destination servers load_{dest}(t) \leftarrow \text{Eq. (13)},
             the set of migration source/destination
             servers \mathcal{M}_{src} \leftarrow \phi, \mathcal{M}_{dest} \leftarrow \phi;
           for each server i \in \mathcal{M} do
4
                Calculate load_i(t) \leftarrow Eq. (11);
 5
                if load_i(t) > load_{src}(t) then
 6
                 Put server i in the set \mathcal{M}_{src};
                if load_i(t) < load_{dest}(t) then
 8
                 Put server i in the set \mathcal{M}_{dest};
            // Phase II: container scheduling
           for all k \in [1, K] do
10
                Obtain the set of migrated containers
11
                 \mathcal{N}_{mig} \leftarrow \text{sampling containers from } \mathcal{M}_{src};
                for each container c \in \mathcal{N}_{mig} do
12
                    for each server i \in \mathcal{M}_{dest} do
13
                         Calculate the container scheduling
14
                          cost C(t) \leftarrow Eq. (4), if container c
                          is migrated to server i;
                    Update X_t by migrating container c to
15
                      the server i with the smallest C(t);
                if X_t satisfies the Constraints (8) – (10) then
16
                 break;
                Update \mathbf{X}_t \leftarrow \mathbf{X}_{t-1}; // \mathbf{X}_t cannot
                    satisfy scheduling constraints
       Calculate the overall container scheduling cost
19
        over the time window W as cost_z \leftarrow Eq. (7)
        with \mathbf{X}_t for each sample z;
```

20 Obtain $\mathbf{X}_t^{min} \leftarrow \mathbf{X}_t$ with the minimum $cost_z$ among all Z samples;

21 **return:** Container scheduling decisions for the first timeslot (*i.e.*, $\mathbf{X}_{t_1}^{min}$).

by

$$load_{i}(t) = \sum_{k,l \in \mathcal{N}} x_{i}^{k}(t)x_{i}^{l}(t) \cdot \left(cpu^{k}(t)cpu^{l}(t) + \beta \cdot mem^{k}(t)mem^{l}(t)\right). \tag{11}$$

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 $load_i(t)$ of each server i can vary with the hosted containers and their resource consumption, we obtain the load

threshold of migration source servers $load_{src}(t)$ and that of destination servers $load_{dest}(t)$ as in Theorem 1. In particular, Tetris can tolerate moderate prediction errors of container resource consumption for classifying migration source and destination servers.

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, we formulate the two load thresholds $load_{src}(t)$ and $load_{dest}(t)$ as

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

$$load_{dest}(t) = \frac{1}{2} \left(\overline{CPU(t)}^2 + \beta \cdot \overline{MEM(t)}^2\right)$$

$$-\sum_{k=1}^{n(t)} \left(cpu^k(t)^2 + \beta \cdot mem^k(t)^2\right).$$
(12)

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 kth 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, available online. \Box

Phase II. container scheduling: To make Algorithm 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 extensive 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 (line 11). For each container to be migrated, we further select the migration destination server in \mathcal{M}_{dest} with the smallest container scheduling cost and then update \mathbf{X}_t (lines 12 to 15). Considering the randomness of sampling, the obtained container scheduling decisions \mathbf{X}_t can violate constraints (8) – (10). 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 16 to 17).

Finally, we iteratively calculate the overall container scheduling cost over the time window W for each sample z. We choose the container scheduling decisions \mathbf{X}_t^{min} with the minimized scheduling 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 19 to 21).

Determining parameters in Tetris: We obtain three key parameters (i.e., v, W, Z) in Algorithm 1 as below. First, we 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 configure 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., 20%). W_2 is the window size with the fastest growing number of invalid migrations (in Fig. 4). 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 Algorithm 1, the time complexity of *Tetris* is in the order of $\mathcal{O}(ZWK \cdot |\mathcal{N}||\mathcal{M}|)$. 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* can be practically acceptable, which will be validated in Section V-D.

Prototype implementation: We implement a prototype of *Tetris* based on K8s, with over 1,000 lines of Python and Linux Shell codes which are publicly available on GitHub.³ In particular, we implement the container migration module of *Tetris* by converting container scheduling decisions into a series of K8s pod operations (i.e., pod deletion and creation commands executed on migration source and destination servers). We plan to integrate the container live migration of CRIU into *Tetris* as our future work.

V. 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 real-world Alibaba production cluster traces v2018 and v2022. We focus on answering the questions listed below.

- Accuracy: Can Tetris accurately predict the resource consumption of long-running containers? (Section V-B)
- Effectiveness: Can our Tetris scheduling strategy jointly optimize the cluster load balancing degree and migration cost for long-running containers without incurring SLO violations? (Section V-C)
- Overhead: How much runtime overhead does Tetris practically bring? (Section V-D)

A. Experimental Setup

Cluster configurations and workloads: 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. We select three representative containerized workloads including Apache Tomcat server⁴, Redis⁵, and

(13)

³https://github.com/icloud-ecnu/Tetris

⁴https://tomcat.apache.org/

⁵https://redis.io/

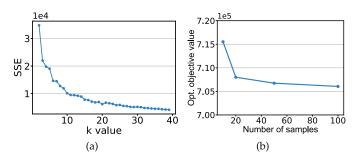


Fig. 7. (a) Sum of the squared errors (SSE) of k-means clustering on CPU and memory resource consumption of containers with different k values, and (b) the optimization objective value decreasing slowly from 20 samples in the Monte Carlo method.

ResNet50 [32] training. We use Apache-benchmark⁶ and redisbenchmark⁷ as the time-varying user requests for Tomcat Web server and Redis, respectively. We adopt CIFAR-10⁸ as the dataset for ResNet50. We *randomly* deploy the three workloads on the 60 containers in our cluster. To obtain complementary insights, we build a trace-driven simulator based on a discrete-event simulation framework [33]. We adopt Alibaba cluster traces v2018 and v2022 [9] for the prediction of container resource consumption and large-scale simulations.

Tetris parameters and model coefficients: We first configure four Tetris parameters. 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. The number of trials K and the sampling ratio v of migration containers are empirically set as 10 and 40%, respectively. Second, we empirically set four model coefficients including α , β , γ , and δ as 0.004, 0.0025, 1, and 10, respectively. The detailed configuration process can be found in Appendix C, available online.

Baselines and metrics: We evaluate Tetris against the conventional Sandpiper [8] and the state-of-the-art Metis⁺ (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 container index volume calculated by the multi-dimensional resource consumption [8]. Metis⁺ uses the method in *Tetris* to select migration source servers and migrated containers and leverages the negative value of (4) as the reward of RL to make container scheduling (i.e., where to schedule) decisions. In particular, we use the mean absolute percentage error (MAPE) [34] to evaluate the efficacy of the *predictor* module of *Tetris*. Meanwhile, we adopt the load imbalance degree defined in (1) and the migration cost defined in (2) 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

TABLE II

NUMBER OF CONTAINERS WITH PERIODIC OR APERIODIC CPU/MEMORY
RESOURCE CONSUMPTION IN ALIBABA CLUSTER TRACE V2018 AND V2022

	Container type	CPU	Memory
v2018	#periodic containers #aperiodic containers	49,962 17,475	40, 501 26, 936
v2022	#periodic containers #aperiodic containers	416, 905 260, 033	151, 147 525, 791

scheduling performance with error bars of standard deviation by repeating experiments three times.

B. Validating Container Resource Predictor in Tetris

We first classify the resource consumption of containers as periodic or aperiodic using the method elaborated in Section IV-A. As shown in Table II, the periodic CPU and memory resource consumption accounts for 61.6%-74.1% and 22.3%-60.1% of containers, respectively, while the proportion of aperiodic CPU and memory resource consumption are only 25.9%-38.4% and 39.9%-77.7%, respectively, for both Alibaba cluster traces v2018 and v2022.

We next examine the prediction accuracy of container resource consumption. We use 70% of the trace data to train the ARIMA and LSTM models in *Tetris*. As shown in Fig. 8, we *first* observe that ARIMA is more accurate than LSTM in predicting the resource consumption of periodic containers. This is because the moving average and autoregression in ARIMA accurately captures the periodic resource consumption of containers, with a small prediction error (i.e., less than 15%). Second, LSTM achieves a more accurate prediction error (i.e., 0.9% - 13.9%) for the resource consumption of aperiodic containers compared with ARIMA. This is mainly because ARIMA can easily be affected by abnormal spikes in aperiodic resource consumption, while LSTM can learn such large resource fluctuations through model training. Interestingly, LSTM achieves around 10.4% - 23.9%of prediction error for periodic containers as depicted in Fig. 8(b) and (d), simply because LSTM contains nonlinear activation functions and thus overfits the periodic data [35]. Third, the prediction error of container resource consumption can significantly increase from 0.9% to 23.9% with the time window size ranging from 1 to 6. That is the reason why *Tetris* controls the time window size W within 3 to maintain an acceptable prediction error (i.e., within 20%) of container resource consumption.

We finally examine the prediction overhead of *Tetris*. Specifically, the average prediction time of ARIMA and LSTM are 0.25 seconds and 0.36 seconds, respectively, for each container. Though *Tetris* requires a shared model for all containers and a personalized model for each container, the average storage footprint of an ARIMA model and an LSTM model are just 556KB and 908KB, respectively. Such time and space overhead above for the prediction of container resource consumption is practically acceptable.

⁶https://httpd.apache.org/docs/2.4/programs/ab.html

⁷https://redis.io/topics/benchmarks

⁸https://www.cs.toronto.edu/~kriz/cifar.html

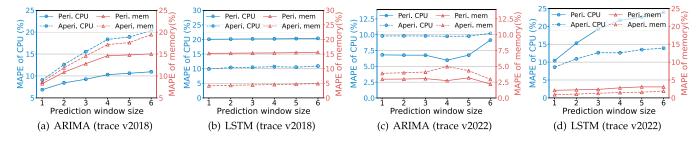


Fig. 8. Average prediction accuracy of CPU and memory resource consumption of 67,437 containers and 676,938 containers in Alibaba cluster traces v2018 and v2022, respectively, by varying the prediction window size from 1 to 6.

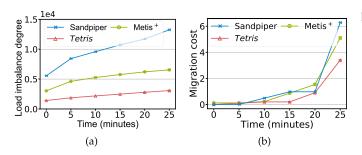


Fig. 9. Performance comparison of *Tetris* with Sandpiper and Metis⁺ strategies under K8s-based prototype experiments, in terms of the cumulative values of (a) cluster load imbalance degree and (b) migration cost of containers over time (i.e., 5 timeslots with 5 minutes each).

C. Effectiveness of Tetris

Prototype experiments on Amazon EC2: We first evaluate the effectiveness of *Tetris* in a 60-container K8s cluster by setting the timeslot⁹ 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⁺, respectively. Moreover, the migration cost with *Tetris* is 7.8% - 79.5% and 10.3% - 77.0% lower than that with Sandpiper and Metis⁺, 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⁺ 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* achieves the most load balancing than the other two strategies at both time 0 and time 1, by proactively migrating out four small containers (i.e., c3-c6) at time 0. In comparison, Sandpiper causes an invalid migration of container c1. As server s1 is overloaded at time 0 by setting the overload threshold as 85%, Sandpiper chooses to migrate

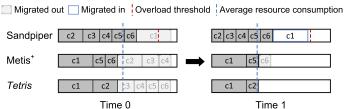


Fig. 10. Comparison of container scheduling decisions on server s1 of the 60-container K8s cluster in the timeslot [0,1] under *Tetris*, Sandpiper, and Metis⁺ strategies.

TABLE III
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
Sandpiper	0	0	14	22	26	61
Metis ⁺	0	0	0	8	8	20
Tetris	0	0	0	0	0	0

container c1 (i.e., the container with the maximum volume) to server s2. It then migrates container c1 back to server s1 as server s2 is overloaded at time 1, resulting in poor load balancing and large migration cost. As for Metis⁺, it migrates out containers (i.e., c2-c4) and container c6 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⁺ 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 III, *Tetris* causes zero SLO violations, while Sandpiper and Metis⁺ 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 (9)–(10)) on the CPU and memory resource capacities of servers, guaranteeing the workload SLO 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

⁹The timeslot can be determined by the cluster administrator according to the cluster requirements.

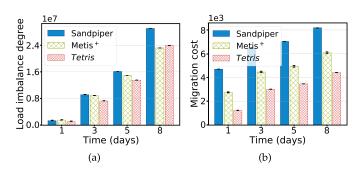


Fig. 11. Performance comparison of *Tetris* with Sandpiper and Metis⁺ strategies under simulations on Alibaba cluster trace v2018, in terms of the cumulative values of (a) cluster load imbalance degree and (b) migration cost of containers over 8 days (the timeslot is 1 h).

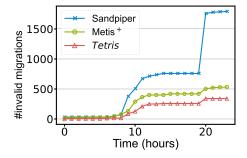


Fig. 12. Cumulative number of invalid migrations of *Tetris* compared with that of Sandpiper and Metis⁺ strategies under simulations on a 24-hour Alibaba cluster trace v2018.

server, causing unexpected SLO violations to containers. The RL method in Metis⁺ is not designed to guarantee the scheduling constraints, so that 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 tracedriven simulations, by setting the timeslot as 1 h. 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⁺. Additionally, *Tetris* achieves the lowest migration cost, which is reduced by 27.4% - 74.1% than the other two strategies. Such results above 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 many invalid migrations and a high growth rate of migration cost as depicted in Fig. 11(b). Though Metis⁺ can optimize the load imbalance degree and migration cost over the infinite future, it achieves a worse load balancing degree in the early stage (i.e., the first 5 days) and a better load 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 a moderate number of invalid migrations in the early stages.

We further examine whether *Tetris* can alleviate invalid migrations. As shown in Fig. 12, *Tetris* reduces the number of

TABLE IV
CUMULATIVE NUMBER OF SLO VIOLATIONS OVER TIME IN A SIMULATED 67,437-CONTAINER CLUSTER DRIVEN BY ALIBABA CLUSTER TRACE V2018

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

TABLE V PERFORMANCE COMPARISON OF TETRIS WITH VARIOUS PARAMETER CONFIGURATIONS (I.E., [W,Z]) ON A 24-HOUR ALIBABA CLUSTER TRACE V2018, BY JOINTLY CONSIDERING CLUSTER LOAD IMBALANCE DEGREE $(C_b(t))$ AND MIGRATION COST OF CONTAINERS $(C_m(t))$

[W,Z]	[2, 20]	[4, 20]	[6, 20]	[2, 10]	[2, 30]
$C_b(t)$	59.8	60.8	63.1	62.4	60.4
$C_m(t)$	1.6	2.0	3.5	2.8	1.7

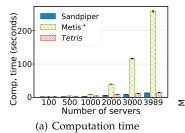
[2, 20] particularly denotes the default Tetris.

invalid migrations by 59.4% – 81.3% and 32.4% – 78.9% as compared with Sandpiper and Metis⁺, respectively. *Tetris* and Metis⁺ 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⁺ 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 the resource consumption of containers varies greatly which overloads a number of servers at time 20.

Similar to our prototype experiments, we proceed to examine whether *Tetris* can guarantee the SLO of workloads in large-scale simulations. As shown in Table IV, we observe that *Tetris* does not cause any SLO violations for 8 days, while both Sandpiper and Metis⁺ 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. Our simulation results above are consistent with the prototype experiments, demonstrating the effectiveness of *Tetris* in guaranteeing workload SLOs.

We finally conduct a sensitivity analysis of two key *Tetris* parameters (i.e., the window size W and the number of samples Z) in Algorithm 1. As shown in Table V, the default *Tetris* achieves the *lowest* container scheduling performance among various parameter configurations. Specifically, the performance metrics (i.e., $C_b(t)$ and $C_m(t)$) both increase *moderately* as W varies from 2 to 6, simply because the prediction error of container resource consumption increases as shown in Fig. 8. Additionally, the performance of container scheduling first decreases and then stabilizes as Z increases from 10 to 30, which is consistent with our configuration analysis of Z in Fig. 7.

17 containers.



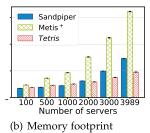


Fig. 13. Runtime overhead of *Tetris* compared with that of Sandpiper and Metis⁺ strategies for each timeslot (i.e., 1 h) on a 24-hour Alibaba cluster trace v2018, by varying the number of servers from 100 to 3,989 with each hosting

TABLE VI CUMULATIVE COMPUTATION OVERHEAD (IN MINUTES) OF *TETRIS*, SANDPIPER AND METIS⁺ SCHEDULING ALGORITHMS EXECUTED ON A 8-DAY ALIABABA CLUSTER TRACE V2018

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

D. 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⁺. The rationale is that the search space of Metis⁺ contains all possible mappings of containers to servers, while Sandpiper and *Tetris* only consider the overloaded servers (containers) and the containers on the migration source and destination servers, respectively. Though the time complexity of Tetris is in the order of $\mathcal{O}(|\mathcal{N}| \cdot |\mathcal{M}|)$ (as discussed in Section IV-B), the *quadratic term ratio* (obtained using polynomial regression) of *Tetris* is only 1.9e - 07, which is much smaller than that of Metis⁺ (i.e., 2.4e - 05) and Sandpiper (i.e., 5.5e - 07). Similarly, *Tetris* consumes a much smaller amount of memory than Metis⁺ and Sandpiper as shown in Fig. 13(b). This is because Metis⁺ 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 VI. The computation time of *Tetris* is slightly longer than that of Sandpiper (i.e., 0.7 – 4.8 minutes), while much shorter than that of Metis⁺ (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⁺ greedily selects the largest reward action across all state spaces in each timeslot. Also, it requires online training to update the RL model. In contrast, *Tetris* periodically triggers the execution of Algorithm 1 in each timeslot. Such time overhead is much

smaller than Metis⁺ per timeslot as depicted in Fig. 13(a). As a result, the runtime overhead of *Tetris* is practically acceptable.

VI. RELATED WORK

Short-term optimization of workload scheduling: There have been many works devoted to scheduling VMs or containers by considering the short-term (e.g., the current or upcoming timeslot) benefits. To predict the resource consumption of workloads, Dabbagh et al. [36] leverage the Wiener filter method for achieving VM consolidations [37], while Madu [30] designs a more complicated attention-based machine learning model to reduce workload latency. Bayesian optimization [38] and predictionbased vertical auto-scaling [39] have been proposed to improve cluster utilization and workload performance. Medea [2] further considers several constraints like affinity and cardinality of containers. To particularly optimize the initial placement of containers, Lv et al. [24] propose a communication-aware worst-fit decreasing algorithm and select the best server for satisfying the network SLO in K8s [40]. To achieve cluster load balancing, Sandpiper [8] performs the worst fitting algorithm using the server load index based on cluster resources. ATCM [41] schedules containers among the cloud servers to minimize the migration cost while maintaining the degree of load balance in a short time scale. Nevertheless, such short-term optimizations of workload scheduling above can inevitably cause invalid migrations as illustrated in Section II-B. 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 initial placement plans for long-running containers, by modeling the reward as an indicator of container performance and constraint violations. To improve cluster utilization Mondal et al. [4] schedule time-varying workloads to servers using deep reinforcement learning (DRL). To consider the long-term optimization of workload scheduling, Megh [42] leverages an online RL-based VM migration method to minimize energy consumption and avoid SLO violations. A-SARSA [43] adopts an RL-based container horizontal scaling approach to adjust the number of actions corresponding to each state to avoid repeated scheduling. However, RL-based scheduling methods can cause unexpected SLO violations by considering the long-term optimization over the infinite future, as evidenced by Section II-C and Section V-C. Additionally, RL methods have high time complexity and memory consumption even after the dimensionality reduction, which requires efforts to be deployed in large-scale clusters as shown in Section V-D. 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.

VII. CONCLUSION AND FUTURE WORK

This paper presents the design and implementation of *Tetris*, a proactive MPC-based container scheduling strategy for

achieving long-term 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 container to schedule* and *where to schedule*) 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 traces v2018 and v2022. 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 GPU, network and disk I/O resources. We also plan to incorporate CRIU into *Tetris* to support container live migrations and examine the effectiveness and scalability of *Tetris* in large-scale production clusters.

REFERENCES

- M. Brooker, M. Danilov, C. Greenwood, and P. Piwonka, "On-demand container loading in AWS lambda," in *Proc. USENIX Annu. Tech. Conf.*, 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. 13th EuroSys Conf.*, 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. IEEE Int. Conf. High Perform. Comput., Netw., Storage Anal.*, 2020, pp. 1–17.
- [4] S. S. Mondal, N. Sheoran, and S. Mitra, "Scheduling of time-varying workloads using reinforcement learning," in *Proc. AAAI Conf. Artif. Intell.*, 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.* ACM Symp. Cloud Comput., 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. 10th Eur. Conf. Comput. Syst.*, 2015, pp. 1–17.
- [7] D. Bernstein, "Containers and cloud: From LXC to docker to kubernetes," IEEE Cloud Comput., vol. 1, no. 3, pp. 81–84, Sep. 2014.
- [8] T. Wood, P. Shenoy, A. Venkataramani, and M. Yousif, "Sandpiper: Black-box and gray-box resource management for virtual machines," *Comput. Netw.*, vol. 53, no. 17, pp. 2923–2938, 2009.
- [9] Alibaba, "Alibaba cluster trace program," 2024. [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 Gener. Comput. Syst.*, 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 Comput.: Pract. Experience*, vol. 29, no. 12, 2017, Art. no. e4123.
- [12] H. Qiu et al., "AWARE: Automate workload autoscaling with reinforcement learning in production cloud systems," in *Proc. USENIX Annu. Techn. Conf.*, 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. Amer. Control Conf.*, 2016, pp. 1987–1992.
- [14] M. Gaggero and L. Caviglione, "Model predictive control for energy-efficient, quality-aware, and secure virtual machine placement," *IEEE Trans. Automat. Sci. Eng.*, vol. 16, no. 1, pp. 420–432, Jan. 2019.
- [15] C. E. Garcia, D. M. Prett, 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," *J. Amer. Statist. Assoc.*, vol. 44, no. 247, pp. 335–341, 1949.

- [17] H. Wu et al., "Aladdin: Optimized maximum flow management for shared production clusters," in *Proc. IEEE Int. Parallel Distrib. Process. Symp.*, 2019, pp. 696–707.
- [18] J. Guo et al., "Who limits the resource efficiency of my datacenter: An analysis of alibaba datacenter traces," in *Proc. IEEE/ACM 27th Int. Symp. Qual. Serv.*, 2019, pp. 1–10.
- [19] A. Parayil et al., "Towards cloud efficiency with large-scale workload characterization," 2024, arXiv:2405.07250.
- [20] M. Xu et al., ": Deep neural network based multivariate workload prediction in cloud computing environments," ACM Trans. Internet Technol., vol. 22, no. 3, pp. 1–24, 2022.
- [21] S. Deng et al., "Cloud-native computing: A survey from the perspective of services," in *Proc. IEEE*, vol. 112, no. 1, pp. 12–46, Jan. 2024.
- [22] J. Thiede, N. Bashir, D. Irwin, and P. Shenoy, "Carbon containers: A system-level facility for managing application-level carbon emissions," in *Proc. ACM Symp. Cloud Comput.*, 2023, pp. 17–31.
- [23] Y. Zhang, L. Jiao, J. Yan, and X. Lin, "Dynamic service placement for virtual reality group gaming on mobile edge cloudlets," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 8, pp. 1881–1897, Aug. 2019.
- [24] L. Lv et al., "Communication-aware container placement and reassignment in large-scale internet data centers," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 3, pp. 540–555, Mar. 2019.
- [25] E. S. Hou, N. Ansari, and H. Ren, "A genetic algorithm for multiprocessor scheduling," *IEEE Trans. Parallel Distrib. Syst.*, vol. 5, no. 2, pp. 113–120, Feb. 1994.
- [26] M. R. Garey and D. S. Johnson, Computers and Intractability: A Guide to the Theory of NP-Hardness. San Fransisco: W. H. Freeman, 1979.
- [27] G. Christofidi, K. Papaioannou, and T. D. Doudali, "Is machine learning necessary for cloud resource usage forecasting?," in *Proc. ACM Symp. Cloud Comput.*, 2023, pp. 544–554.
- [28] G. E. Box, G. M. Jenkins, G. C. Reinsel, and G. M. Ljung, *Time Series Analysis: Forecasting and Control*. Hoboken, NJ, USA: Wiley, 2015.
- [29] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [30] S. Luo et al., "The power of prediction: Microservice auto scaling via workload learning," in *Proc. 13th Symp. Cloud Comput.*, 2022, pp. 355– 369.
- [31] Z. Li and D. Hoiem, "Learning without forgetting," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 12, pp. 2935–2947, Dec. 2018.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [33] F. Li and B. Hu, "DeepJS: Job scheduling based on deep reinforcement learning in cloud data center," in *Proc. 4th Int. Conf. Big Data Comput.*, 2019, pp. 48–53.
- [34] 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.
- [35] 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.
- [36] M. Dabbagh, B. Hamdaoui, M. Guizani, and A. Rayes, "An energy-efficient VM prediction and migration framework for overcommitted clouds," *IEEE Trans. Cloud Comput.*, vol. 6, no. 4, pp. 955–966, Dec. 2016.
- [37] Z. Á. Mann, "Allocation of virtual machines in cloud data centers-a survey of problem models and optimization algorithms," ACM Comput. Surv., vol. 48, no. 1, pp. 1–34, 2015.
- [38] C. Luo et al., "Intelligent virtual machine provisioning in cloud computing," in *Proc. Int. Joint Conf. Artif. Intell.*, 2021, pp. 1495–1502.
- [39] J. Zhu et al., "QoS-Aware co-scheduling for distributed long-running applications on shared clusters," *IEEE Trans. Parallel Distrib. Syst.*, vol. 33, no. 12, pp. 4818–4834, Dec. 2022.
- [40] E. Kim, K. Lee, and C. Yoo, "Network SLO-Aware container scheduling in kubernetes," *J. Supercomputing*, vol. 79, no. 10, pp. 11478–11494, 2023.
- [41] 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 Gener. Comput. Syst.*, vol. 135, pp. 303–314, 2022.
- [42] D. Basu, X. Wang, Y. Hong, H. Chen, and S. Bressan, "Learn-as-you-go with megh: Efficient live migration of virtual machines," *IEEE Trans. Parallel Distrib. Syst.*, vol. 30, no. 8, pp. 1786–1801, Aug. 2019.
- [43] 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. IEEE Int. Conf. Web Serv.*, 2020, pp. 489–497.



Fei Xu (Member, IEEE) received the BS, ME, and PhD degrees from the Huazhong University of Science and Technology (HUST), Wuhan, China, in 2007, 2009, and 2014, respectively. He is currently a full professor with the School of Computer Science and Technology, East China Normal University, Shanghai, China. He received Outstanding Doctoral Dissertation Award in Hubei province, China, and ACM Wuhan & Hubei Computer Society Doctoral Dissertation Award, in 2015, and a recipient of the Best Student Paper Award of CCFSys 2023. His re-

search interests include cloud computing and datacenter, and distributed machine learning systems.



Zhi Zhou (Member, IEEE) received the BS, ME, and PhD degrees from the School of Computer Science and Technology, Huazhong University of Science and Technology (HUST), Wuhan, China, in 2012, 2014, and 2017, respectively. He is currently an associate professor in the School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou, China. In 2016, he was a visiting scholar with the 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.



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.



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 with Tencent, Shenzhen, China. Her research interests focus on NoSQL and distributed storage systems.



Shuohao 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.



Fangming Liu (Senior Member, IEEE) received the BEng degree from the Tsinghua University, Beijing, and the PhD 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.



Li Chen (Member, IEEE) received the BEng degree from the Department of Computer Science and Technology, Huazhong University of Science and Technology, China, in 2012, 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, Uni-

versity of Louisiana with Lafayette, Lafayette, USA. Her research interests include Big Data analytics systems, cloud computing, datacenter networking, and resource allocation.