

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

Xiyue Shen^a, Fei Xu^{a,*}, Yifan Liu^b, Tao Song^b, Shuohao Lin^a, Li Chen^c, Fen Xiao^d

^a*Shanghai Key Laboratory of Multidimensional Information Processing, School of Computer Science and Technology, East China Normal University*

^b*Shanghai F7 Networks Technology Inc.*

^c*University of Louisiana at Lafayette*

^d*Tencent Inc.*

Abstract

Long-running containerized workloads (*e.g.*, online cloud services, machine learning) are increasingly prevailing in shared production clusters. As the request loads of such long-running containers typically show *time-varying* patterns, container scheduling is essential to improving the cluster resource utilization and workload performance. Existing cluster schedulers mainly focus on optimizing either the *short-term* benefits of cluster load balancing or the *initial placement* of long-running containers on servers. Such schedulers, however, would inevitably bring a noticeable number of *invalid migrations* (*i.e.*, containers are migrated back and forth between two servers over a short time window), leading to serious service level objective (SLO) violations. In this paper, we propose *Tetris*, a *model predictive control* (MPC)-based container scheduling strategy to judiciously migrate long-running workloads

*The corresponding author is Fei Xu.

Email addresses: `fxu@cs.ecnu.edu.cn` (Fei Xu), `yangjl@f7sys.com.cn` (Tao Song), `li.chen@louisiana.edu` (Li Chen), `cherryxiao@tencent.com` (Fen Xiao)

for cluster load balancing. Specifically, we first build a discrete-time dynamic model to formulate a *long-term* optimization problem of container scheduling. *Tetris* then solves such an optimization problem using two key components: (1) a container resource predictor, which leverages time-series analysis methods to accurately predict the container resource consumption; (2) an MPC-based container scheduler that jointly optimizes the load balancing and migration cost of long-running containers *over a certain sliding time window*. We implement a prototype of *Tetris* based on K8s, and evaluate *Tetris* with extensive prototype experiments on Amazon EC2 and large-scale simulations driven by Alibaba cluster trace v2018. Experiment results show that *Tetris* can improve the cluster load balancing degree by up to 77.8% without incurring any SLO violations, in comparison to the state-of-the-art container scheduling strategies, yet with acceptable runtime overhead.

Keywords:

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

9 unpredictable changes of container resource consumption can cause severe
10 contention of shared cluster resources. Such cluster load imbalance can re-
11 sult in many potential service level objective (SLO) violations to workloads.
12 Consequently, the mainstream cluster schedulers such as Borg [5] and Kuber-
13 netes [6] enable container scheduling to achieve load balancing [7] in shared
14 production clusters.

15 Though the existing container scheduling policies such as Sandpiper [7]
16 and Medea [2] perform well in achieving cluster load balancing, there still
17 exist a noticeable number of *invalid migrations* (*i.e.*, containers are migrated
18 back and forth between two servers over a short time window) in shared
19 production clusters. As evidenced by our motivational analysis (in Sec. 2.2)
20 on Alibaba cluster trace v2018 [8], the number of invalid migrations is over
21 1,200, accounting for more than 50% of container migrations within the win-
22 dow size of three timeslots. Such invalid migrations are likely to make the
23 workloads hosted on the migrated containers suffer from serious SLO viola-
24 tions, leading to unexpected performance interference to the containers that
25 are *co-located* on servers (*i.e.*, migration cost). Our motivation experiments
26 in Sec. 2.2 also reveal that the invalid migrations can significantly reduce
27 the number of requests processed by an Apache Tomcat Web server hosted
28 on a migrated container by up to 99.8%. Accordingly, it is critical for the
29 cluster scheduler to avoid invalid migrations during the process of container
30 scheduling, especially for long-running workloads.

31 Unfortunately, many research efforts have been devoted to making *short-*
32 *term* scheduling decisions based on the *current* cluster status for workload
33 consolidation [9] or load balancing [10]. Though such scheduling policies

34 can achieve short-term (*e.g.*, the current or upcoming timeslot) benefits,
 35 they are oblivious to the *time-varying* resource consumption of long-running
 36 workloads, which is actually the *root cause of invalid migrations* as discussed
 37 in Sec. 2.2. There have also been recent works on the *long-term* optimization
 38 of container scheduling, which are reinforcement learning (RL)-based [1, 11]
 39 or control-based [12, 13] scheduling policies, to *partially* tackle the issue of
 40 invalid migrations. Nevertheless, these techniques solely focus on optimizing
 41 the *initial placement* or resource auto-scaling of containerized workloads (*i.e.*,
 42 *where to migrate*) over the entire time period (*i.e.*, *infinite future*). The
 43 containers to be scheduled (*i.e.*, *which to migrate*) and the *migration cost* of
 44 containers have surprisingly received little attention. Accordingly, such *long-*
 45 *term* optimization techniques can still cause unexpected SLO violations, as
 46 evidenced in Sec. 2.3. As a result, there has been a paucity of research
 47 attention paid to developing container scheduling policies to *fully* deal with
 48 the invalid migrations of long-running workloads in containerized clusters.

49 To fill this gap, we propose *Tetris*, a model predictive control [14] (MPC)-
 50 based container scheduling strategy to judiciously make migration decisions
 51 for long-running containerized workloads. *Tetris* simply optimizes container
 52 scheduling *over a certain sliding time window* rather than *the infinite future*
 53 (as in RL-based method [3]) for two reasons: *first*, the prediction accuracy
 54 of container resource consumption dramatically decreases as the prediction
 55 window size increases. Our prediction results using time-series techniques (in
 56 Sec. 5.2) indicate that the prediction error can exceed 20% as the prediction
 57 window size reaches 6; *second*, blindly increasing the time window size can
 58 significantly increase the number of samples of *Tetris*, which requires notice-

59 able computation overhead to obtain container scheduling plans. To the best
60 of our knowledge, *Tetris* is the first attempt to achieve the *long-term* opti-
61 mization of container scheduling to *circumvent invalid migrations as many*
62 *as possible*, by jointly optimizing the cluster load balancing and container
63 migration cost over a certain sliding time window. Our main contributions
64 are summarized as follows.

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

72 ▷ *Second*, we design an MPC-based container scheduling strategy to
73 achieve *long-term* scheduling optimization for cluster load balancing (Sec. 4).
74 Specifically, *Tetris* first accurately predicts the container resource consump-
75 tion within a time window using time-series analysis methods. It then
76 leverages the Monte Carlo method [15] to judiciously identify the container
77 scheduling decision for each timeslot, by minimizing our formulated cost
78 function of container scheduling. In particular, we calculate the upper and
79 lower bounds of server load imbalance degree to classify the migration source
80 and destination servers, so that the complexity of *Tetris* container scheduling
81 strategy can be significantly reduced.

82 ▷ *Finally*, we implement a prototype of *Tetris*¹ based on K8s. We eval-

¹<https://github.com/icloud-ecnu/tetris>

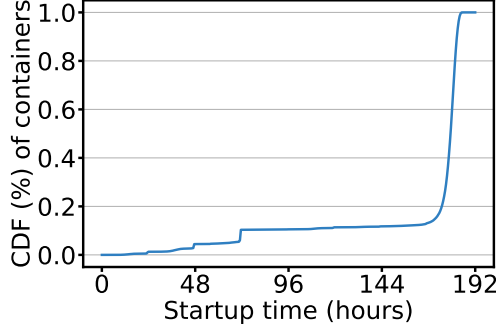


Figure 1: Container uptime in Alibaba cluster trace v2018.

83 uate the effectiveness and runtime overhead of *Tetris* with prototype exper-
 84 iments on 10 EC2 instances (*i.e.*, 60 containers) and large-scale simulations
 85 driven by Alibaba cluster trace v2018 (Sec. 5). Compared with the con-
 86 ventional scheduling strategy (*i.e.*, Sandpiper [7]) and the state-of-the-art
 87 RL-based method (*i.e.*, Metis⁺, a modified version of Metis [3]), *Tetris* is
 88 able to improve the cluster load balancing degree by up to 77.8% while cut-
 89 ting down the migration cost by up to 79.5%, yet with acceptable runtime
 90 overhead.

91 2. Background and Motivation

92 2.1. Long-running Containerized Workloads

93 Large-scale shared production clusters often host a number of long-running
 94 containers in response to latency-critical user requests [2]. A variety of long-
 95 running workloads including streaming, machine learning, and Web applica-
 96 tions run on them, in addition to storage services generated by multi-tenant
 97 NoSQL databases, *etc.* [16]. Taking Alibaba as an example, its online ser-
 98 vices are mainly user-interactive and long-running applications that run on
 99 containers, such as search engines and online shopping [17]. As depicted

100 in Fig. 1, over 80% of containers in the Alibaba production cluster last
 101 for more than 175 hours². Due to the unpredictability of user requests,
 102 such long-running containerized workloads have diverse demands on hard-
 103 ware resources including CPU, memory, network and disk I/O, which are
 104 generally dynamic variability and uncertainty in the amount of resource re-
 105 quests [18]. Furthermore, the co-location of containers can lead to contention
 106 for shared resources, which can easily cause performance interference and re-
 107 source wastage in the cluster.

108 Based on the above, enabling adequate scheduling of long-running con-
 109 tainerized workloads is essential to cluster load balancing. In shared pro-
 110 duction clusters, K8s achieves load balancing among servers by carrying out
 111 container migrations [6]. During the migration process, the container needs
 112 to stop services first, and then it is terminated on the source server and
 113 later restarted on the destination server, and finally it replicates the mem-
 114 ory/storage data to restore services. Accordingly, the container scheduling
 115 inevitably brings a noticeable amount of migration cost, which requires to
 116 be considered during the scheduling of containers.

117 2.2. Invalid Migrations of Long-Running Workloads

118 Load balancing is an effective way to alleviate resource contention in
 119 shared containerized clusters [7]. While many existing scheduling policies
 120 perform well in load balancing, they usually focus on the *short-term* benefits
 121 of container scheduling. As discussed above, the resource consumption of
 122 long-running containers are commonly *time-varying* (*i.e.*, fluctuating wildly

²The total time in the Alibaba cluster trace v2018 is 192 hours.

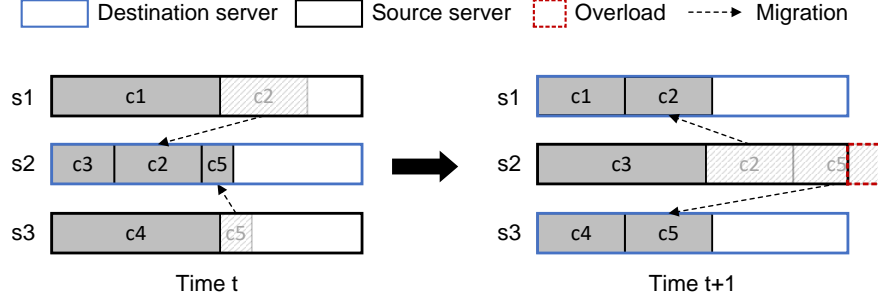


Figure 2: Illustration of invalid migration. Containers $c2$ and $c5$ are migrated back and forth among servers $s1$, $s2$, and $s3$.

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 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.

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 .

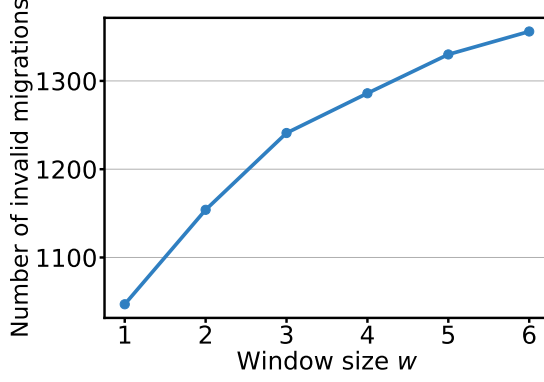


Figure 3: Number of invalid migrations within a time window size varying from 1 to 6.

Moreover, we validate the prevalence of invalid migrations using the Alibaba cluster trace v2018 with an 8-day period [8]. 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, and the number reaches around 1,400, accounting for 75.3% of container migrations when w is set as 6. To 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 future resource consumption of containers.

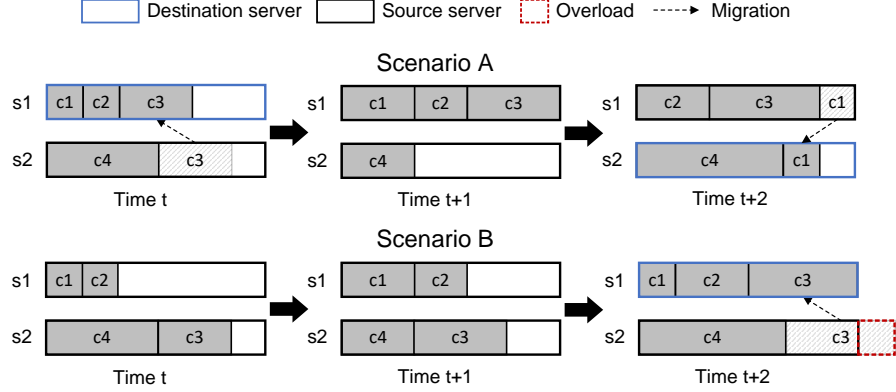


Figure 4: An illustrative example of MPC-based container scheduling (our proposed *Tetris* in scenario A), compared with RL-based container scheduling (Metis⁺ [3] in scenario B).

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 which allows the current timeslot to be optimized while taking *finite future timeslots* in account [19]. To illustrate how *Tetris* works, we show a motivation example in Fig. 4 by comparing *Tetris* with a RL-based scheduling method (*i.e.*, Metis⁺, a modified version of Metis [3]). 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 guarantees SLOs of all containerized workloads. Accordingly, the RL-based method can optimize the scheduling objective

171 *over the entire time period* at the cost of overloading a certain number of
 172 containers, while *Tetris* achieves the *long-term* optimization for container
 173 scheduling *within a certain sliding time window*. In addition, the RL-based
 174 method has the following problems such as requiring repeated training on
 175 a large amount of high-quality data samples and lacking interpretability as
 176 well as poor scalability [3], which will be further validated in Sec. 5.4.

177 **Summary:** Avoiding invalid migrations is critical to container schedul-
 178 ing, and such invalid migrations are mainly caused by the *short-term* (e.g.,
 179 the current or upcoming timeslot) optimization of container scheduling. Mean-
 180 while, greedily optimizing container scheduling over the entire time period
 181 (*i.e., infinite future*) are likely to cause unexpected SLO violations. Accord-
 182 ingly, there is a compelling need to design a *long-term* container scheduling
 183 strategy, by jointly optimizing the cluster load balancing and migration cost
 184 of containers within a certain sliding time window.

185 3. Model and Problem Formulation

186 In this section, we first build a discrete-time dynamic model to capture the
 187 load imbalance degree of clusters and migration cost of containers. Next, we
 188 formulate a container scheduling optimization problem based on our dynamic
 189 model. The key notations in our model are summarized in Table 1.

190 3.1. Discrete-Time Dynamic Model

191 We consider a containerized cluster including a set of servers \mathcal{M} hosting
 192 a set of containers \mathcal{N} . We assume the time t is discrete and slotted, where
 193 $t = t_0, t_1, \dots, t_W$ and W is the time window size. For each container $k \in \mathcal{N}$, its
 194 CPU and memory resource consumption at each time t is recorded as $cpu^k(t)$

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
$C_{mig}^k(t)$	Migration cost of container k at time t
$cpu^k(t)$	CPU consumption of container k at time t
$mem^k(t)$	Memory consumption of 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 server i at time t
$MEM_i^{cap}(t)$	Available capacity of memory resource of server i at time t
α	Normalized parameter of migration cost to load imbalance degree
β	Normalized parameter of memory resources to CPU resources
$x_i^k(t)$	Indicator whether container k is on server i at time t
$m^k(t)$	Indicator whether container k is migrated at time t

195 and $mem^k(t)$, respectively. For each server $i \in \mathcal{M}$, $CPU_i(t)$ and $MEM_i(t)$
196 represent its total CPU and total memory resource consumption at time t ,
197 respectively.

198 **Modeling cluster load imbalance degree at each timeslot.** We
199 use the variance of resource consumption of all cluster servers to represent
200 the load imbalance degree of the cluster. A larger degree value indicates a
201 severer load imbalance situation. By taking CPU and memory resources as
202 an example, the cluster load imbalance degree is represented by calculating
203 the weighted sum of the load imbalance degrees of each resource of servers,

204 which is given by

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

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

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

209 Then, we proceed to denote the average CPU resource consumption $\overline{CPU(t)}$
 210 and average memory resource consumption $\overline{MEM(t)}$ of a server in the cluster
 211 as below,

$$\overline{CPU(t)} = \frac{1}{|\mathcal{M}|} \sum_{k \in \mathcal{N}} cpu^k(t), \quad (3)$$

$$\overline{MEM(t)} = \frac{1}{|\mathcal{M}|} \sum_{k \in \mathcal{N}} mem^k(t). \quad (4)$$

212 **Modeling migration cost at each timeslot.** We use the sum of unit
 213 cost of the migrated containers to denote the migration cost. A larger number
 214 of migrated containers and a larger unit cost of each migrated container indi-
 215 cate a severer migration cost. Accordingly, the migration cost is formulated

$$C_m(t) = \sum_{k \in \mathcal{N}} C_{mig}^k(t) \cdot m^k(t), \quad (5)$$

217 where $m^k(t)$ denotes whether container k is migrated at time t , which is
 218 given by Eq. (6). $C_{mig}^k(t)$ is the unit migration cost of container k . As
 219 elaborated in Sec. 2.1, K8s [6] implements the container migration with three
 220 steps: service stop, container deletion and restart, and data replication. In
 221 particular, the data replication overhead is proportional to the container's
 222 memory footprint [20], while the overhead of the other two steps can be
 223 considered as a constant value δ . Hence, the unit migration cost of a container
 224 can be given by

$$m^k(t) = \sum_{i \in \mathcal{M}} x_i^k(t) \cdot (1 - x_i^k(t-1)), \quad (6)$$

$$C_{mig}^k(t) = \delta + \gamma \cdot mem^k(t), \quad (7)$$

225 where γ is a model coefficient obtained by workload profiling, and $mem^k(t)$
 226 is the memory footprint of container k at time t .

227 To sum up, we further formulate the *cost function* $C(t)$ of container
 228 *scheduling at each timeslot*. Based on the formulations above, we further
 229 combine the cluster load imbalance degree and the migration cost as below,

$$C(t) = C_b(t) + \alpha \cdot C_m(t), \quad (8)$$

230 where $\alpha \in [0, 1]$ denotes the normalized parameter of migration cost to cluster
 231 load imbalance degree. In more detail, we substitute Eq. (3) and Eq. (4) into
 232 Eq. (1) and substitute Eq. (6) and Eq. (7) into Eq. (5), yielding Eq. (9) and

Eq. (10), respectively, which are calculated as

$$C_b(t) = \frac{2}{|\mathcal{M}| - 1} \sum_{i \in \mathcal{M}} \sum_{k, l \in \mathcal{N}} x_i^k(t) x_i^l(t) \cdot (cpu^k(t) cpu^l(t) + \beta \cdot mem^k(t) mem^l(t)) + C_1, \quad (9)$$

$$C_m(t) = C_2 - \sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{M}} (\delta + \gamma \cdot mem^k(t)) \cdot x_i^k(t) x_i^k(t - 1), \quad (10)$$

where $C_1 = -\frac{|\mathcal{M}| \cdot (\overline{CPU}(t)^2 + \beta \cdot \overline{MEM}(t)^2) + \sum_{k \in \mathcal{N}} (cpu^k(t)^2 + \beta \cdot mem^k(t)^2)}{|\mathcal{M}| - 1}$ and $C_2 = \delta \cdot |\mathcal{N}| + \gamma \cdot \sum_{k \in \mathcal{N}} mem^k(t)$ are both constant values. The detailed mathematical derivations can be found in [Appendix .1](#). By analyzing Eq. (9), the cluster load imbalance degree is determined by the sum of the pairwise multiplications between the resource consumption of containers hosted on each server. By analyzing Eq. (10), the migration cost across the cluster is determined by the negative of the sum of migration costs of all 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. (11)). We assume the scheduling starts at time t_1 , and our optimization

248 problem can be formulated as

$$\min_{x_i^k(t)} \sum_{t=t_1}^{t_W} C(t) \quad (11)$$

$$\text{s.t.} \quad \sum_{i \in \mathcal{M}} x_i^k(t) = 1, \quad \forall k \in \mathcal{N}, t \in [t_1, t_W] \quad (12)$$

$$\sum_{k \in \mathcal{N}} x_i^k(t) \cdot \text{cpu}^k(t) \leq \text{CPU}_i^{\text{cap}}(t), \quad \forall i \in \mathcal{M}, t \in [t_1, t_W] \quad (13)$$

$$\sum_{k \in \mathcal{N}} x_i^k(t) \cdot \text{mem}^k(t) \leq \text{MEM}_i^{\text{cap}}(t). \quad \forall i \in \mathcal{M}, t \in [t_1, t_W] \quad (14)$$

249 The hard constraints in our optimization problem are described as follows:
 250 Constraint (12) limits each container to be hosted only on one server per time
 251 to avoid scheduling conflicts. Constraints (13)-(14) ensure the total resource
 252 consumption on each server cannot exceed the server resource capacity.

253 **Problem Analysis.** For each timeslot, the minimization problem de-
 254 fined in Eq. (11) can be easily reduced to a *multiprocessor scheduling problem*
 255 (*MSP*) [21]. MSP can be considered as finding a schedule for multiple tasks
 256 to be executed on a multiprocessor system at different execution time, so that
 257 the completion time can be minimized. Such a scheduling problem is known
 258 to be NP-hard [22]. Moreover, Eq. (8) is a multi-objective optimization prob-
 259 lem that jointly considers minimizing the migration cost and the cluster load
 260 imbalance degree, which further makes such a problem hard to solve. As a
 261 result, our optimization problem of Eq. (11) turns out to be NP-hard. In
 262 the upcoming section, we turn to leveraging the Monte Carlo method [15] to
 263 solve our long-term workload scheduling optimization problem.

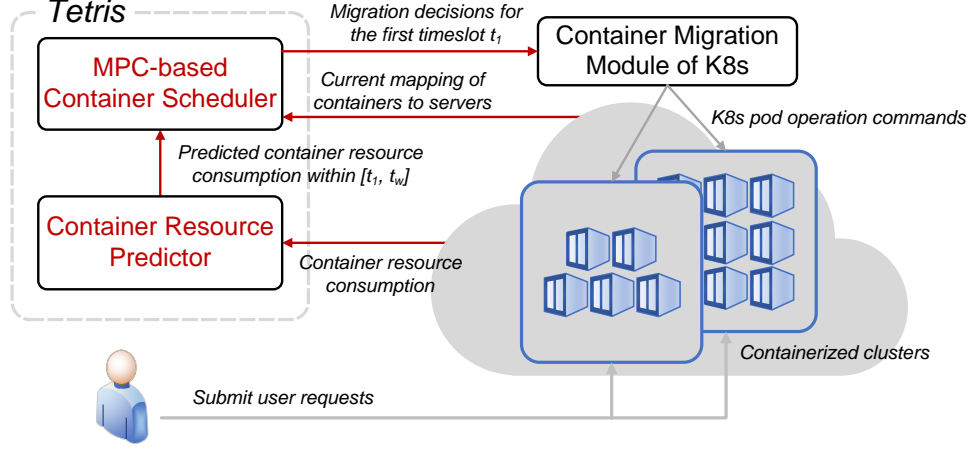


Figure 5: Overview of *Tetris* prototype in containerized clusters.

264 4. Tetris Design

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

270 4.1. Overview of *Tetris*

271 As illustrated in Fig. 5, *Tetris* comprises two pieces of modules including a
 272 *container resource predictor* (Sec. 4.2) and a *MPC-based container scheduler*
 273 (Sec. 4.3). After users submit workload requests to the containerized cluster,
 274 *Tetris* first leverages the *predictor* to estimate the resource consumption of
 275 containers over a given time window W , which is input to the *scheduler* for
 276 calculating the migration cost of containers. By *jointly* optimizing the cluster
 277 load balancing degree and container migration cost, the *scheduler* further
 278 decides the appropriate migration plans during the period of time window W ,

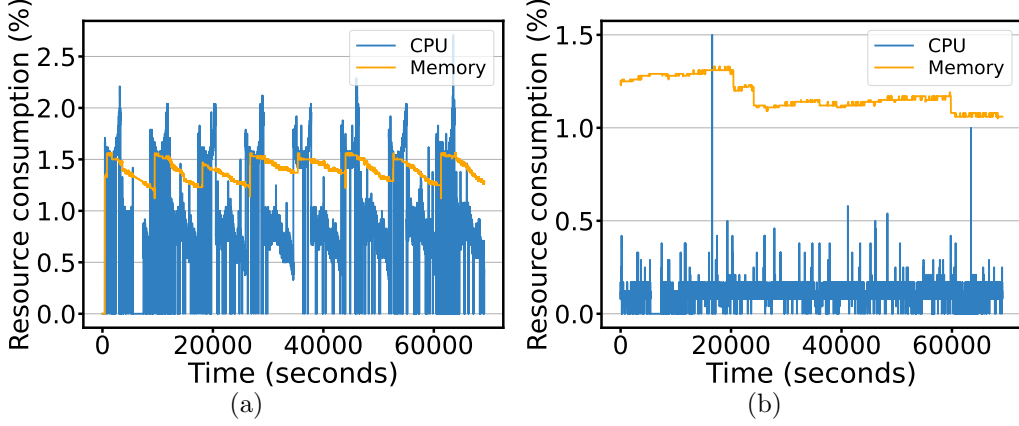


Figure 6: Resource consumption of containers over time: (a) a periodic container ($id = 71,417$), and (b) an aperiodic container ($id = 58$) in Alibaba cluster trace v2018.

279 which will be elaborated in Alg. 1. Finally, we implement a *container migration*
 280 *module of K8s*, which performs the appropriate migration decisions for
 281 the *first* timeslot in the containerized cluster, while discarding the migration
 282 decisions for the remaining timeslots. In particular, such a migration module
 283 can convert container scheduling decisions into a series of K8s pod³ operations
 284 (*i.e.*, pod deletion and creation commands executed on migration source and
 285 destination servers). The prototype of *Tetris* is implemented based on K8s,
 286 with over 1,000 lines of Python and Linux Shell codes which are publicly
 287 available on GitHub (<https://github.com/icloud-ecnu/tetris>).

288 4.2. Predicting Container Resource Consumption

289 We first classify workloads as *periodic* and *aperiodic* containers shown in
 290 Fig. 6, as most containerized workloads show a diurnal pattern [17]. Specifi-
 291 cally, we first obtain the *frequency* of container resource consumption values

³Note that each pod in our K8s cluster hosts one container in our experimental setup.

292 and calculate its corresponding *periodicity*. With such a periodicity value,
 293 we then slice the container resource consumption into a number of segments.
 294 We finally calculate the Pearson correlation coefficients between every two
 295 consecutive segments and decide whether such coefficients exceed a given pe-
 296 riodicity threshold thr_{period} (e.g., 0.85). If the coefficients exceed the thresh-
 297 old, we consider such a container as *periodic*. Otherwise, we consider the
 298 container as *aperiodic*.

299 Next, we proceed to predict the resource consumption for such two types
 300 of containers separately. In more detail, we first adopt ARIMA [23] to pre-
 301 dict the resource consumption of *periodic* containers. We use the *histori-*
 302 *cal* container resource consumption to pre-train a *shared* ARIMA model for
 303 newly-launched containers. To improve the prediction accuracy, we further
 304 leverage incremental learning [24] to train a *personalized* ARIMA model for
 305 each container that has been running for a certain period of time, and we
 306 can update the model periodically. In particular, we mainly obtain three key
 307 parameters (*i.e.*, the autoregressive process of order p , the moving average
 308 process of order q , and the difference order d) of the ARIMA model. Sec-
 309 ond, we leverage LSTM [25] to predict the resource consumption of *aperiodic*
 310 containers. We adopt a 4-layer stacked LSTM model with 50 neurons per
 311 layer to increase the depth of the neural network, and add a *Dropout* layer
 312 between every two LSTM layers to reduce overfitting. To keep the LSTM
 313 model simple and alleviate the impact of error propagation, we directly con-
 314 duct the multistep-ahead predictions and the number of prediction steps is
 315 set as W .

316 4.3. MPC-based Container Scheduling

317 We leverage the MPC approach to make container scheduling decisions
 318 over the time window W . Specifically, we use a $|\mathcal{M}| \times |\mathcal{N}|$ matrix \mathbf{X}_t of
 319 $x_i^k(t)$ ($i \in \mathcal{M}, k \in \mathcal{N}$) to denote the mapping of containers to servers at time
 320 $t \in [t_0, t_W]$. According to MPC, we first solve the scheduling optimization
 321 problem in Sec. 3.2 to obtain $\mathbf{X}_t, \forall t \in [t_1, t_W]$, at the current time t_0 . Then,
 322 we only perform the container scheduling decision \mathbf{X}_{t_1} for the first timeslot
 323 t_1 . After the containers are scheduled to migration destination servers at the
 324 “current” time t_1 , we continue solving the scheduling optimization problem
 325 and then perform the container scheduling decisions for the “first” timeslot
 326 t_2 . To maintain load balancing of the containerized cluster, *Tetris* can be
 327 *periodically* (e.g., for several minutes or hours) executed at each timeslot.

328 In more detail, we solve the scheduling optimization problem based on
 329 the Monte Carlo method [15] as discussed in Sec. 3.2. As shown in Alg 1,
 330 we take Z *sample* solutions in total, and each sample solution includes a set
 331 of container scheduling decisions $\mathbf{X}_t, \forall t \in [t_1, t_W]$ over the time window W
 332 (lines 1 to 2). To obtain the container scheduling decisions for each timeslot
 333 t , we design two phases in Alg. 1, *i.e.*, server classification (lines 3 to 9) and
 334 container scheduling (lines 10 to 19), which are elaborated as follows.

335 **Phase I: Server classification.** We classify the cluster servers into
 336 three categories, *i.e.*, migration source servers and destination servers as well
 337 as the remaining servers (lines 3 to 9). Based on the cluster load imbal-
 338 ance degree (*i.e.*, Eq. (9)) defined in our dynamic model, we can obtain the

Algorithm 1: *Tetris*: Long-term container scheduling algorithm for jointly optimizing load balancing and migration cost of containers.

Input: Current mapping \mathbf{X}_{t_0} of container set \mathcal{N} to server set \mathcal{M} , container CPU and memory resource consumption over a time window 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
5               $load_{src}(t) \leftarrow$  Eq. (16) and that for migration destination
6              server  $load_{dest}(t) \leftarrow$  Eq. (17);
7          for each server  $i \in \mathcal{M}$  do
8              Calculate the load imbalance degree  $load_i(t) \leftarrow$  Eq. (15);
9              if  $load_i(t) > load_{src}(t)$  then
10                 Put server  $i$  to migration source server set  $\mathcal{M}_{src}$ ;
11             if  $load_i(t) < load_{dest}(t)$  then
12                 Put server  $i$  to migration destination server set  $\mathcal{M}_{dest}$ ;
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  $C(t) \leftarrow$ 
20                         Eq. (8) if container  $c$  is migrated to server  $i$ ;
21                     Update  $\mathbf{X}_t$  by migrating container  $c$  to the server  $i$ 
22                         with the smallest  $C(t)$ ;
23                 if  $\mathbf{X}_t$  satisfies scheduling constraints (12) – (14) then
24                     break;
25             Update  $\mathbf{X}_t \leftarrow \mathbf{X}_{t-1}$ ; //  $\mathbf{X}_t$  cannot satisfy constraints
26         Calculate the container migration cost over the time window  $W$ 
27          $cost_z \leftarrow$  Eq. (11) of  $\mathbf{X}_t$  for each sample  $z$ ;
28     Obtain  $\mathbf{X}_t^{min} \leftarrow \mathbf{X}_t$  with the minimized  $cost_z$  among all  $Z$  samples;
29 return: Scheduling decisions for the first timeslot (i.e.,  $\mathbf{X}_{t_1}^{min}$ ).

```

339 simplified *load imbalance degree* (i.e., $load_i(t)$) for each server i given by

$$load_i(t) = \sum_{k,l \in \mathcal{N}} x_i^k(t) x_i^l(t) \cdot (cpu^k(t) cpu^l(t) + \beta \cdot mem^k(t) mem^l(t)). \quad (15)$$

340 Obviously, the CPU and memory resource consumption of each server are
 341 $\overline{CPU}(t)$ (i.e., Eq. (3)) and $\overline{MEM}(t)$ (i.e., Eq. (4)), respectively, when the
 342 cluster reaches the “ideal” *load-balanced state*. As the load imbalance degree
 343 of each server $load_i(t)$ can vary with the number of hosted containers and
 344 the resource consumption of containers, we obtain the thresholds of $load_i(t)$
 345 for migration source servers and destination servers as in Theorem 1.

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

$$load_{src}(t) = \left(\frac{1}{2} - \frac{1}{2|\mathcal{N}|} \right) \cdot \left(\overline{CPU}(t)^2 + \beta \cdot \overline{MEM}(t)^2 \right), \quad (16)$$

$$load_{dest}(t) = \frac{1}{2} \left(\overline{CPU}(t)^2 + \beta \cdot \overline{MEM}(t)^2 - \sum_{k=1}^{n(t)} (cpu^k(t)^2 + \beta \cdot mem^k(t)^2) \right). \quad (17)$$

350 *By sorting containers in a descending order with the CPU and memory re-*
 351 *source consumption, $cpu^k(t)$ and $mem^k(t)$ represent the CPU and memory*
 352 *consumption of the k -th container at time t , respectively. $n(t)$ denotes the*
 353 *minimum number of containers which satisfies $\sum_{k=1}^{n(t)} cpu^k(t) \leq \overline{CPU}(t)$ and*
 354 $\sum_{k=1}^{n(t)} mem^k(t) \leq \overline{MEM}(t)$.

Proof. The proof can be found in [Appendix .2](#). □

355 **Phase II: Container scheduling.** To make Alg. 1 practical, we sample
356 a set of containers to be migrated \mathcal{N}_{mig} from the set of migration source
357 servers \mathcal{M}_{src} . To achieve a more comprehensive coverage of possibilities
358 with a small number of samples, we adopt Latin Hypercube Sampling (LHS)
359 with a sampling rate of v per migration source server (lines 11 to 12). For
360 each container to be migrated, we further select the migration destination
361 server with the smallest container migration cost and then update \mathbf{X}_t (lines
362 13 to 16). Considering the randomness of sampling, the obtained container
363 scheduling decisions \mathbf{X}_t can violate constraints (12) – (14). If no violations
364 occur, \mathbf{X}_t is valid and we continue the container scheduling process for the
365 next timeslot $t + 1$. Otherwise, we require re-sampling \mathcal{N}_{mig} and obtain
366 another \mathbf{X}_t , until K trials are finished (line 10) or the scheduling constraints
367 are satisfied (lines 17 to 19).

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

373 **Determining parameters in *Tetris*.** We obtain three parameters (*i.e.*,
374 v , W , Z) in Alg. 1 using the methods below. *First* is to determine the sam-
375 pling rate v for migration containers. We empirically select the proportion
376 of the container resource consumption exceeding $\overline{CPU(t)}$ and $\overline{MEM(t)}$ in
377 migration source servers as the value of v . *Second* is to decide the time win-
378 dow size W . We obtain the maximum time window W_1 where the prediction
379 error of container resource consumption is within the required range (*e.g.*,

15%). Also, we obtain the window size W_2 with the fastest growing number of invalid migrations (in Fig. 3). We finally set W as the minimal value of W_1 and W_2 . *Third* is to determine 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 use the k clusters as the minimal value of Z .

Time complexity analysis. According to Alg. 1, the time complexity of *Tetris* is in the order of $\mathcal{O}(ZWK \cdot |\mathcal{N}||\mathcal{M}|)$. As the number of containers $|\mathcal{N}|$ and the number of servers $|\mathcal{M}|$ is *far larger* than the time window size W and the number of trials K , the complexity can be reduced to $\mathcal{O}(Z \cdot |\mathcal{N}| \cdot |\mathcal{M}|)$. 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)
- *Overhead*: How much runtime overhead does *Tetris* bring? (Sec. 5.4)

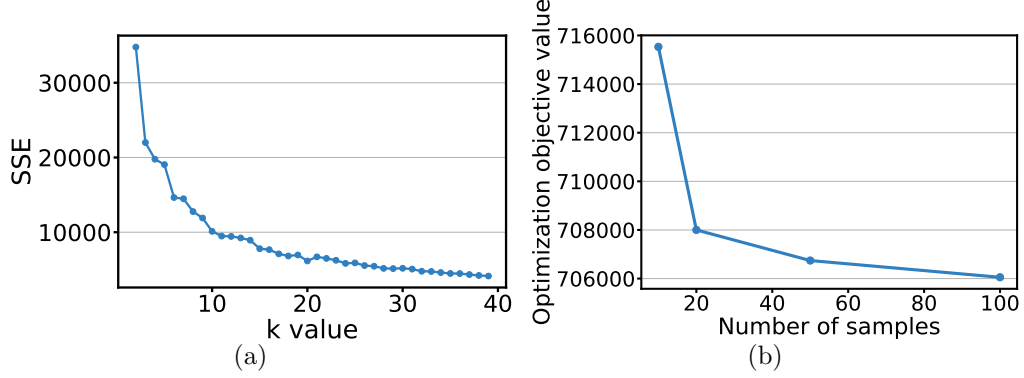


Figure 7: (a) Sum of the squared errors (SSE) of k -means clustering on CPU and memory resource consumption of containers with different k values, and (b) the optimization objective value converging at 20 samples in the Monte Carlo method.

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 [26] running on a commodity server equipped with 24 3.5 GHz Intel Core i9-10920X CPU cores and 32 GB memory. 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 trails K and the sampling ratio v of migration containers are empirically set as 10 and 40%, respectively.

418 **Workloads and datasets.** We employ three representative workloads
 419 in the containerized clusters, including Apache Tomcat server⁴, Redis⁵, and
 420 ResNet50 [27] training. We use Apache-benchmark⁶ and redis-benchmark⁷
 421 as the time-varying user requests for Tomcat Web server and Redis, respec-
 422 tively. We use CIFAR-10⁸ as the dataset for ResNet50. We *randomly* deploy
 423 the three workloads on the 60 containers in the K8s cluster. For the predic-
 424 tion of container resource consumption and large-scale simulations, we adopt
 425 Alibaba cluster trace v2018 [8] which contains around 4,000 machines over
 426 8 days. We obtain the CPU and memory resource consumption of 67,437
 427 containers using linear interpolation and zero-value padding. We set each
 428 timeslot as 1 hour and calculate the average resource consumption of con-
 429 tainers per hour.

430 **Baselines and metrics.** We evaluate *Tetris* against the conventional
 431 Sandpiper [7] and the state-of-the-art Metis⁺ (*i.e.*, a modified version of
 432 Metis [3]) scheduling algorithms. To achieve load balancing, Sandpiper per-
 433 forms a greedy worst-fit algorithm to schedule containers according to the
 434 container index *volume* calculated by the multi-dimensional resource con-
 435 sumption. Metis⁺ uses the method in *Tetris* to select migration source servers
 436 and migrated containers and leverages the negative value of Eq. (8) as the
 437 reward of RL to make container scheduling (*i.e.*, where to place) decisions.
 438 In particular, we use the *mean absolute percentage error (MAPE)* [28] to

⁴<https://tomcat.apache.org/>

⁵<https://redis.io/>

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

⁷<https://redis.io/topics/benchmarks>

⁸<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

439 evaluate the efficacy of the *predictor* module of *Tetris*. Meanwhile, we adopt
 440 the *load imbalance degree* defined in Eq. (1) and the *migration cost* defined
 441 in Eq. (5) as well as the number of *SLO violations* to evaluate the efficacy
 442 of the *scheduler* module of *Tetris*. Due to the randomness of sampling in
 443 Monte Carlo method and workload placement on containers, we illustrate
 444 the container scheduling performance with error bars of standard deviation
 445 by repeating experiments for three times.

446 5.2. Predicting Resource Consumption of Containers in Tetris

447 We first classify the resource consumption of containers as periodic or
 448 aperiodic using the method elaborated in Sec. 4.2. As shown in Table 2,
 449 the periodic CPU and memory resource consumption accounts for 74.1%
 450 and 60.1% of containers, respectively, while the proportion of aperiodic CPU
 451 and memory resource consumption are only 25.9% and 39.9%, respectively.
 452 Accordingly, most of the containers show a periodic resource consumption
 453 over time, which is consistent with a latest measurement study [17].

454 We next leverage ARIMA and LSTM to predict the periodic and aperi-
 455 odic container resource consumption. To train the prediction model, we use
 456 70% of the trace data as the training dataset. As shown in Fig. 8, we *first*
 457 observe that ARIMA is more accurate than LSTM for predicting the resource
 458 consumption of periodic containers. This is because the moving average and

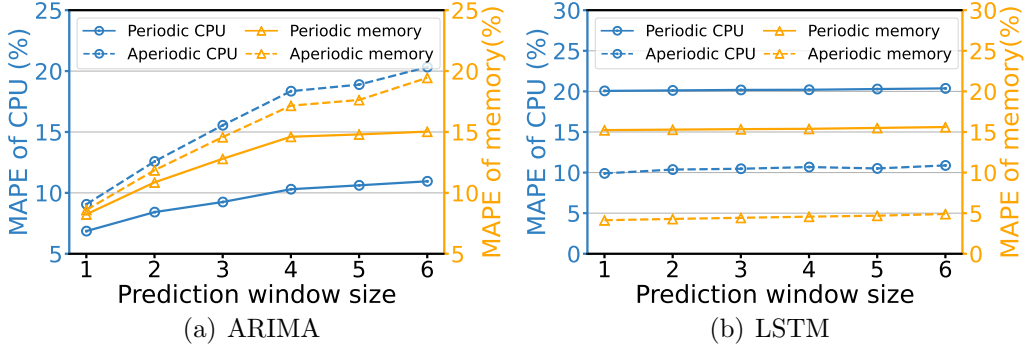


Figure 8: Average prediction accuracy of CPU and memory resource consumption of 67,437 containers in Alibaba cluster trace v2018, by varying the prediction window size from 1 to 6.

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 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 [29]. *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

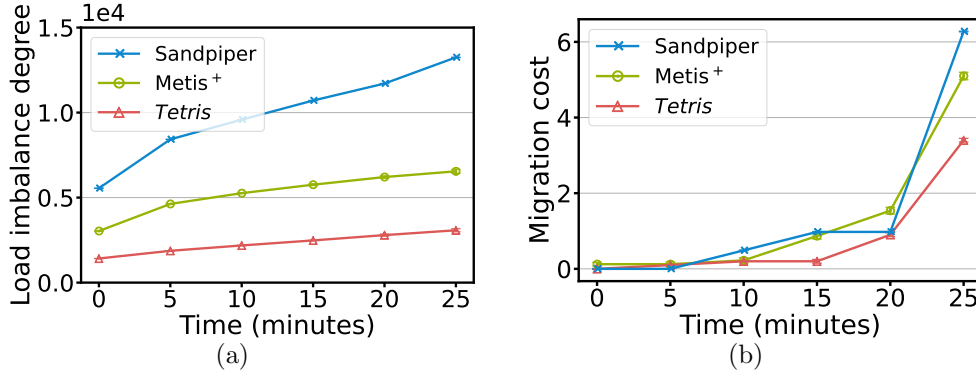


Figure 9: Performance comparison of *Tetris* with Sandpiper and Metis⁺ strategies under K8s-based prototype experiments, in terms of the cumulative values of (a) cluster load imbalance degree and (b) migration cost of containers over time (*i.e.*, 5 timeslots with 5 minutes each).

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 a 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 are both 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* 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

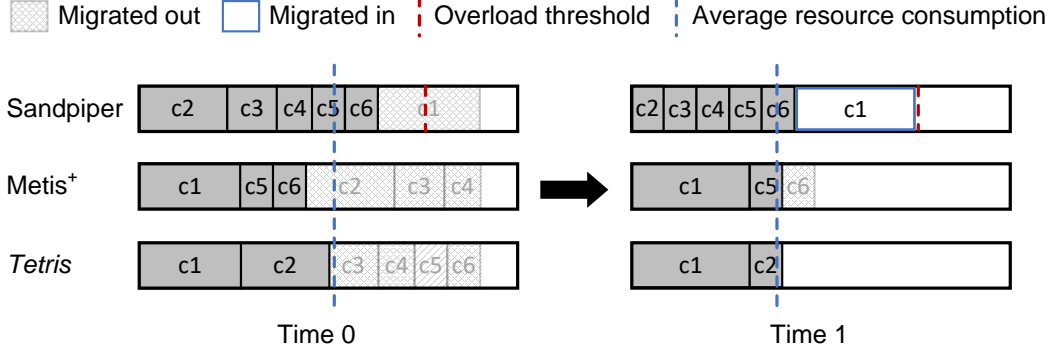


Figure 10: Comparison of container scheduling decisions on server s_1 of the 60-container K8s cluster in the timeslot $[0, 1]$ under *Tetris*, *Sandpiper*, and *Metis⁺* strategies.

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 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⁺*, it

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
Sandpiper	0	0	14	22	26	61
Metis ⁺	0	0	0	8	8	20
<i>Tetris</i>	0	0	0	0	0	0

migrates out containers $c2 - c4$ and $c6$ 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 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 3, *Tetris* causes zero SLO violations while Sandpiper and Metis⁺ can cause up to 35 and 12 SLO violations within one timeslot (*i.e.*, at time 5). The reason is that *Tetris* explicitly considers the hard constraints (*i.e.*, Constraints (13)–(14)) 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⁺ is not designed to fully guarantee Constraints (13)–(14). It allows the occurrence of SLO violations over a certain time period to optimize container

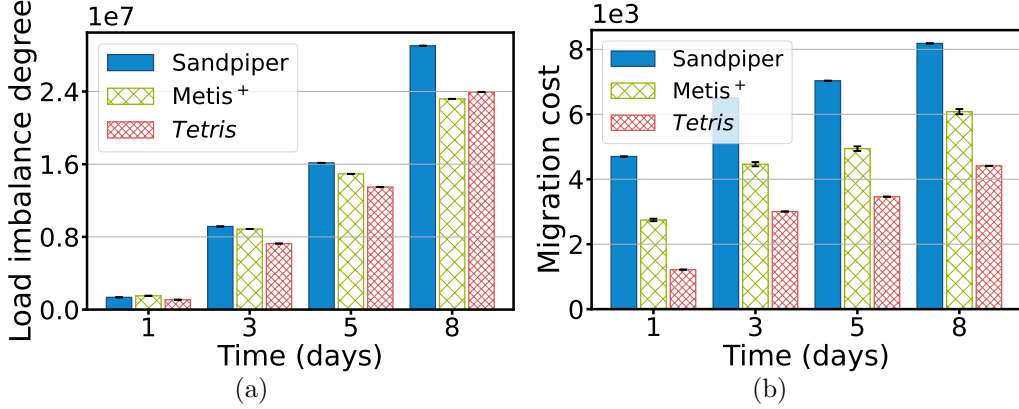


Figure 11: Performance comparison of *Tetris* with Sandpiper and Metis⁺ strategies on Alibaba cluster trace, in terms of the cumulative values of (a) cluster load imbalance degree and (b) migration cost of containers over 8 days (the timeslot is 1 hour).

526 scheduling (*i.e.*, maximize the action reward) over the infinite future.

527 **Large-scale simulations driven by Alibaba cluster trace.** We next
 528 evaluate the effectiveness of *Tetris* with real-world trace-driven simulations,
 529 by setting the *timeslot* as 1 hour. As shown in Fig. 11, *Tetris* can reduce
 530 the load imbalance degree of the cluster by 9.7% – 28.6% compared with
 531 Sandpiper and Metis⁺. It can also achieve the lowest migration cost, which
 532 is reduced by 27.4% – 74.1% than the other two strategies. The experimen-
 533 tal results are consistent with our prototype experiments. This is because
 534 *Tetris* co-optimizes the load balancing and migration cost to alleviate invalid
 535 migrations. In contrast, Sandpiper greedily alleviates the heavy server load
 536 using the worst-fit algorithm, resulting in a number of invalid migrations
 537 and a high growth rate of migration cost as depicted in Fig. 11(b). Though
 538 Metis⁺ can optimize the load imbalance degree and migration cost over the
 539 infinite future, it achieves worse load balance degree in the early stage (*i.e.*,
 540 the first 5 days) and better load balance degree in the late stage (*i.e.*, day

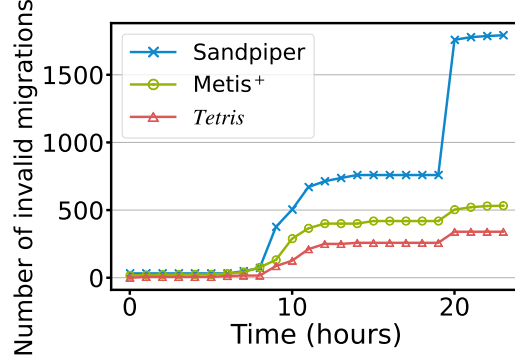


Figure 12: Cumulative number of invalid migrations of *Tetris* compared with that of Sandpiper and Metis⁺ strategies on a 24-hour Alibaba cluster trace.

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⁺, 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 and thus a

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 ⁺	0	220	220	918
<i>Tetris</i>	0	0	0	0

number of servers are overloaded 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⁺ 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 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* achieves 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 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⁺ (*i.e.*, $2.4e - 05$)

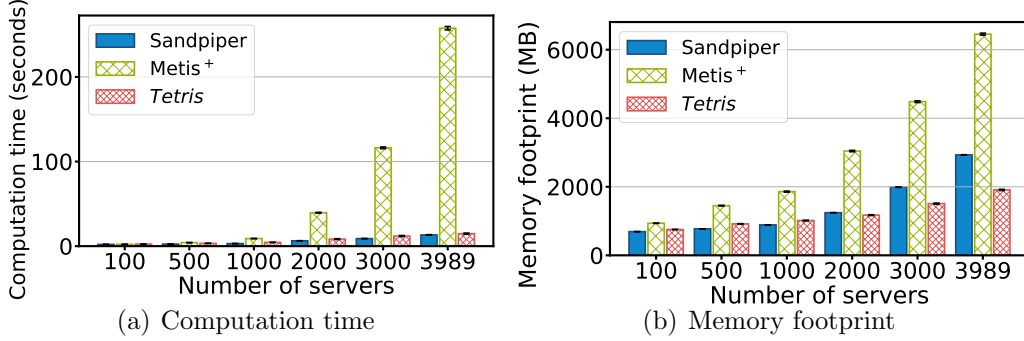


Figure 13: Runtime overhead of *Tetris* compared with that of Sandpiper and Metis⁺ strategies for each timeslot (*i.e.*, 1 hour) on a 24-hour Alibaba cluster trace, by varying the number of servers from 100 to 3,989 with each hosting 17 containers.

Table 5: Cumulative computation overhead (in minutes) of *Tetris*, Sandpiper and Metis⁺ scheduling algorithms executed on a 8-day Aliababa cluster trace.

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

576 and Sandpiper (*i.e.*, $5.5e - 07$). Similarly, *Tetris* consumes much smaller
577 amount of memory than Metis⁺ and Sandpiper as shown in Fig. 13(b). This
578 is because Metis⁺ stores the policy of RL which involves all states and actions
579 as well as the large parameters of neural networks. Sandpiper requires storing
580 the *volumes* of all containers and servers. In comparison, *Tetris* only needs
581 to calculate the load imbalance degree of all servers and the migration cost
582 of a subset of containers (*i.e.*, containers to be migrated).

583 In addition, we illustrate the computation overhead of three strategies
584 over various time scales (*i.e.*, from 1 to 8 days) in Table 5. The computation
585 time of *Tetris* is slightly longer than that of Sandpiper (*i.e.*, 0.7 – 4.8 min-
586 utes), while much shorter than that of Metis⁺ (*i.e.*, 81.6 – 780.4 minutes)

587 as the time scale increases. This is because Sandpiper triggers the fewest al-
 588 gorithm executions only when the resource hotspot occurs. Metis⁺ greedily
 589 selects the largest-reward action over all state spaces in each timeslot, and
 590 it also requires online training to update the RL model. In contrast, *Tetris*
 591 periodically triggers the execution of Alg. 1 in each timeslot and such time
 592 overhead is much smaller than Metis⁺ per timeslot as depicted in Fig. 13(a).
 593 As a result, the runtime overhead of *Tetris* is practically acceptable.

594 6. Related Work

595 **Short-term optimization of workload scheduling.** There have been
 596 a number of works devoted to scheduling VMs or containers by considering
 597 the *short-term* (*e.g.*, the current or upcoming timeslot) benefits. For exam-
 598 ple, in order to efficiently consolidate VMs on servers to save energy [30],
 599 Dabbagh et al. [31] minimize the energy consumption of migrations by lever-
 600 aging Wiener filter method to predict resource consumption. UAHS [32]
 601 leverages Bayesian optimization to maximize the CPU utilization of servers.
 602 To achieve load balancing of clusters, Sandpiper [7] performs the worst fitting
 603 algorithm using the server load index based on CPU, memory, and network
 604 resources. Lv et al. [20] propose a communication-aware worst-fit decreas-
 605 ing container placement algorithm to optimize the initial container place-
 606 ment. To particularly optimizing the scheduling of long-running workloads,
 607 Medea [2] explicitly considers several constraints like affinity and cardinal-
 608 ity of containers, while TOPOSCH [33] leverages a prediction-based vertical
 609 auto-scaling technique to avoid Quality of Service (QoS) violations and bal-
 610 ance the performance of batch jobs. The *short-term* optimization 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 [1]) 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 [34] leverages an online RL algorithm to perform VM migrations at each timeslot to minimize the energy consumption and avoid service level agreement (SLA) violations. A-SARSA [11] dynamically scales the number of instances based on RL methods to reduce SLA violations. To avoid repeated scheduling due to a fixed scaling upper limit, A-SARSA adjusts the number of actions corresponding to each state. 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.

636 7. Conclusion and Future Work

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

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

654 References

- 655 [1] S. Li, L. Wang, W. Wang, Y. Yu, B. Li, George: Learning to place long-
656 lived containers in large clusters with operation constraints, in: Proc. of
657 ACM SOCC, 2021, pp. 258–272.
- 658 [2] P. Garefalakis, K. Karanasos, P. Pietzuch, A. Suresh, S. Rao, Medea:

- scheduling of long running applications in shared production clusters,
in: Proc. of Eurosys, 2018, pp. 1–13.
- [3] L. Wang, Q. Weng, W. Wang, C. Chen, B. Li, Metis: learning to schedule long-running applications in shared container clusters at scale, in: Proc. of IEEE SC, 2020, pp. 1–17.
- [4] S. S. Mondal, N. Sheoran, S. Mitra, Scheduling of time-varying workloads using reinforcement learning, in: Proc. of AAAI, Vol. 35, 2021, pp. 9000–9008.
- [5] A. Verma, L. Pedrosa, M. Korupolu, D. Oppenheimer, E. Tune, J. Wilkes, Large-scale cluster management at google with borg, in: Proc. of Eurosys, 2015, pp. 1–17.
- [6] D. Bernstein, Containers and cloud: From lxc to docker to kubernetes, IEEE Cloud Computing 1 (3) (2014) 81–84.
- [7] T. Wood, P. Shenoy, A. Venkataramani, M. Yousif, Sandpiper: Black-box and gray-box resource management for virtual machines, Computer Networks 53 (17) (2009) 2923–2938.
- [8] Alibaba, [Alibaba Cluster Trace Program](#) (May 2022).
URL <https://github.com/alibaba/clusterdata/>
- [9] A. Beloglazov, J. Abawajy, R. Buyya, Energy-aware resource allocation heuristics for efficient management of data centers for cloud computing, Future Generation Computer Systems 28 (5) (2012) 755–768.

- 680 [10] M. Xu, W. Tian, R. Buyya, A survey on load balancing algorithms
681 for virtual machines placement in cloud computing, *Concurrency and*
682 *Computation: Practice and Experience* 29 (12) (2017) e4123.
- 683 [11] S. Zhang, T. Wu, M. Pan, C. Zhang, Y. Yu, A-sarsa: A predictive
684 container auto-scaling algorithm based on reinforcement learning, in:
685 *Proc. of IEEE ICWS*, 2020, pp. 489–497.
- 686 [12] M. Gaggero, L. Caviglione, Model predictive control for the placement
687 of virtual machines in cloud computing applications, in: *Proc. of ACC*,
688 2016, pp. 1987–1992.
- 689 [13] M. Gaggero, L. Caviglione, Model predictive control for energy-efficient,
690 quality-aware, and secure virtual machine placement, *IEEE Transac-*
691 *tions on Automation Science and Engineering* 16 (1) (2018) 420–432.
- 692 [14] C. E. Garcia, D. M. Prett, M. Morari, Model predictive control: Theory
693 and practice - a survey, *Automatica* 25 (3) (1989) 335–348.
- 694 [15] N. Metropolis, S. Ulam, The monte carlo method, *Journal of the Amer-*
695 *ican Statistical Association* 44 (247) (1949) 335–341.
- 696 [16] H. Wu, W. Zhang, Y. Xu, H. Xiang, T. Huang, H. Ding, Z. Zhang,
697 Aladdin: Optimized maximum flow management for shared production
698 clusters, in: *Proc. of IEEE IPDPS*, 2019, pp. 696–707.
- 699 [17] J. Guo, Z. Chang, S. Wang, H. Ding, Y. Feng, L. Mao, Y. Bao, Who
700 limits the resource efficiency of my datacenter: An analysis of alibaba
701 datacenter traces, in: *Proc. of IEEE/ACM IWQoS*, 2019, pp. 1–10.

- 702 [18] M. Xu, C. Song, H. Wu, S. S. Gill, K. Ye, C. Xu, esdnn: Deep neu-
703 ral network based multivariate workload prediction in cloud computing
704 environments, ACM Transactions on Internet Technology (2022).
- 705 [19] Y. Zhang, L. Jiao, J. Yan, X. Lin, Dynamic service placement for virtual
706 reality group gaming on mobile edge cloudlets, IEEE Journal on Selected
707 Areas in Communications 37 (8) (2019) 1881–1897.
- 708 [20] L. Lv, Y. Zhang, Y. Li, K. Xu, D. Wang, W. Wang, M. Li, X. Cao,
709 Q. Liang, Communication-aware container placement and reassignment
710 in large-scale internet data centers, IEEE Journal on Selected Areas in
711 Communications 37 (3) (2019) 540–555.
- 712 [21] E. S. Hou, N. Ansari, H. Ren, A genetic algorithm for multiprocessor
713 scheduling, IEEE Transactions on Parallel and Distributed Systems 5 (2)
714 (1994) 113–120.
- 715 [22] M. R. Garey, D. S. Johnson, Computers and intractability, Vol. 174,
716 1979.
- 717 [23] G. E. Box, G. M. Jenkins, G. C. Reinsel, G. M. Ljung, Time series
718 analysis: forecasting and control, 2015.
- 719 [24] Z. Li, D. Hoiem, Learning without forgetting, IEEE Transactions on
720 Pattern Analysis and Machine Intelligence 40 (12) (2017) 2935–2947.
- 721 [25] S. Hochreiter, J. Schmidhuber, Long short-term memory, Neural com-
722 putation 9 (8) (1997) 1735–1780.

- 723 [26] F. Li, B. Hu, Deepjs: Job scheduling based on deep reinforcement learn-
724 ing in cloud data center, in: Proc. of ICBDC, 2019, pp. 48–53.
- 725 [27] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image
726 recognition, in: Proc. of CVPR, 2016, pp. 770–778.
- 727 [28] B. L. Bowerman, R. T. O’Connell, A. B. Koehler, Forecasting, time
728 series, and regression: an applied approach, Vol. 4, 2005.
- 729 [29] D. Fan, H. Sun, J. Yao, K. Zhang, X. Yan, Z. Sun, Well production
730 forecasting based on arima-lstm model considering manual operations,
731 Energy 220 (2021) 1–13.
- 732 [30] Z. Á. Mann, Allocation of virtual machines in cloud data centers – a sur-
733 vey of problem models and optimization algorithms, ACM Computing
734 Surveys 48 (1) (2015) 1–34.
- 735 [31] M. Dabbagh, B. Hamdaoui, M. Guizani, A. Rayes, An energy-efficient
736 vm prediction and migration framework for overcommitted clouds, IEEE
737 Transactions on Cloud Computing 6 (4) (2016) 955–966.
- 738 [32] C. Luo, B. Qiao, X. Chen, P. Zhao, R. Yao, H. Zhang, W. Wu, A. Zhou,
739 Q. Lin, Intelligent virtual machine provisioning in cloud computing, in:
740 Proc. of IJCAI, 2021, pp. 1495–1502.
- 741 [33] R. Yang, J. Zhu, X. Sun, T. Wo, C. Hu, H. Peng, J. Xiao, A. Y. Zomaya,
742 J. Xu, Qos-aware co-scheduling for distributed long-running applications
743 on shared clusters, IEEE Transactions on Parallel and Distributed Sys-
744 tems (2022).

[34] D. Basu, X. Wang, Y. Hong, H. Chen, S. Bressan, Learn-as-you-go with
megh: Efficient live migration of virtual machines, IEEE Transactions
on Parallel and Distributed Systems 30 (8) (2019) 1786–1801.

Appendix .1. Mathematical Derivations of Eq. (9) and Eq. (10)

To obtain Eq. (9), the load imbalance degree of CPU resource part in
Eq. (1) can first be expanded as

$$\begin{aligned} \sum_{i \in \mathcal{M}} \left(\sum_{k \in \mathcal{N}} x_i^k(t) \cdot cpu^k(t) - \overline{CPU(t)} \right)^2 &= \sum_{i \in \mathcal{M}} \left(\sum_{k \in \mathcal{N}} x_i^k(t) \cdot cpu^k(t) \right)^2 \\ &\quad - 2 \cdot \overline{CPU(t)} \cdot \sum_{k \in \mathcal{N}} cpu^k(t) \\ &\quad + |\mathcal{M}| \cdot \overline{CPU(t)}^2, \end{aligned} \quad (.1)$$

where the last two terms of the expanded expression in Eq. (.1) are actually
constant for each timeslot t , which can be denoted by $-|\mathcal{M}| \cdot \overline{CPU(t)}^2$. The
load imbalance degree of CPU resource part can be mainly determined by
the first term of Eq. (.1), which is given by

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

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. (.2). The load imbalance degree of memory resource part in Eq. (1)
is similar to that of CPU resource part by substituting $mem^k(t)$ for $cpu^k(t)$.

759 We formulate $C_b(t)$ as

$$C_b(t) = \frac{2}{|\mathcal{M}| - 1} \sum_{i \in \mathcal{M}} \sum_{k, l \in \mathcal{N}} x_i^k(t) x_i^l(t) \cdot (cpu^k(t) cpu^l(t) + \beta \cdot mem^k(t) mem^l(t)) + C_1, \quad (.3)$$

760 where $C_1 = -\frac{|\mathcal{M}| \cdot (\overline{CPU(t)}^2 + \beta \cdot \overline{MEM(t)}^2) + \sum_{k \in \mathcal{N}} (cpu^k(t)^2 + \beta \cdot mem^k(t)^2)}{|\mathcal{M}| - 1}$ is a constant
 761 value. Accordingly, Eq. (.3) is actually the cluster load imbalance degree
 762 formulated in Eq. (9).

763 To obtain Eq. (10), we first incorporate Eq. (6) into Eq. (5), $C_m(t)$ can
 764 be represented by

$$C_m(t) = \sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{M}} C_{mig}^k(t) \cdot x_i^k(t) - \sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{M}} C_{mig}^k(t) \cdot x_i^k(t) x_i^k(t-1), \quad (.4)$$

765 where the first term is transformed into a constant $C_2 = \delta \cdot |\mathcal{N}| + \gamma \cdot$
 766 $\sum_{k \in \mathcal{N}} mem^k(t)$, which is related to the memory sum of all containers $mem^k(t)$, $\forall k \in$
 767 \mathcal{N} . Accordingly, by incorporating Eq. (7) into Eq. (.4), $C_m(t)$ is determined
 768 by the second term of Eq. (.4), which is given by

$$C_m(t) = C_2 - \sum_{k \in \mathcal{N}} \sum_{i \in \mathcal{M}} (\delta + \gamma \cdot mem^k(t)) \cdot x_i^k(t) x_i^k(t-1). \quad (.5)$$

769 As a result, Eq. (.5) is actually the migration cost across the cluster formu-
 770 lated in Eq. (10).

771 Appendix .2. Proof of Theorem 1

772 *Proof.* The load imbalance degree $load_i(t)$ of a server $i \in \mathcal{M}$ and Eq. (16)
 773 – (17) 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.*, $cpu^k(t), \forall k \in \mathcal{N}$). By assuming the cluster servers are in the “ideal” load-balanced state with the CPU resource consumption of $\overline{CPU}(t)$ 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 $mul_{n(t)}$ of $cpu^k(t)$ on a server in Eq. (15)) 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 $cpu^1(t), cpu^2(t), \dots, cpu^{n(t)}(t)$ and $\sum_{k=1}^{n(t)} cpu^k(t) = \overline{CPU}(t)$, where $n(t) \in [1, |\mathcal{N}|]$, the sum of pairwise multiplication $mul_{n(t)}$ of $cpu^k(t), \forall k \in [1, n(t)]$ becomes the maximum (*i.e.*, $\frac{\overline{CPU}(t)^2}{2}$), if and only if $cpu^1(t) = cpu^2(t) = \dots = cpu^{n(t)}(t)$ and $n(t) = |\mathcal{N}|$. Similarly, $mul_{n(t)}$ becomes the minimum (*i.e.*, $\frac{\overline{CPU}(t)^2}{2} - \frac{1}{2} \cdot \sum_{k=1}^{n(t)} cpu^k(t)^2$), if and only if $n(t)$ is the minimum value that satisfies $\sum_{k=1}^{n(t)} cpu^k(t) = \overline{CPU}(t)$.*

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

$$\begin{aligned}
 mul_{n(t)} &= \sum_{k=1}^{n(t)-1} \sum_{l=k+1}^{n(t)} (cpu^k(t) \cdot cpu^l(t)) \\
 &= \frac{1}{2} \left[\left(\sum_{k=1}^{n(t)} cpu^k(t) \right)^2 - \sum_{k=1}^{n(t)} cpu^k(t)^2 \right] \\
 &= \frac{\overline{CPU}(t)^2}{2} - \frac{\sum_{k=1}^{n(t)} cpu^k(t)^2}{2}.
 \end{aligned} \tag{.6}$$

The first term of Eq. (.6) is a constant value. By Cauchy-Buniakowsky-Schwarz Inequality, the second term $\frac{1}{2} \cdot \sum_{k=1}^{n(t)} cpu^k(t)^2$ becomes the minimum

791 (*i.e.*, $\frac{\overline{CPU(t)}^2}{2 \cdot n(t)}$) if and only if $cpu^k(t) = \frac{\overline{CPU(t)}}{n(t)}, \forall k \in [1, n(t)]$. Accordingly, we
 792 have

$$mul_{n(t)} \leq \left(\frac{1}{2} - \frac{1}{2 \cdot n(t)} \right) \cdot \overline{CPU(t)}^2, \quad (.7)$$

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

$$\max(mul_{n(t)}) = \left(\frac{1}{2} - \frac{1}{2 \cdot |\mathcal{N}|} \right) \cdot \overline{CPU(t)}^2. \quad (.8)$$

796 To calculate the minimum value of $mul_{n(t)}$, we have to find the maximum
 797 value of the second term $\frac{1}{2} \cdot \sum_{k=1}^{n(t)} cpu^k(t)^2$ of Eq. (.6). According to the perfect
 798 square tinomial, the square of the sum of a set of non-negative numbers is
 799 larger than or equal to the sum of squares of these numbers. Therefore, we
 800 find that $\frac{1}{2} \cdot \sum_{k=1}^{n(t)} cpu^k(t)^2$ can reach the maximum when $n(t)$ is the minimum
 801 number to satisfy $\sum_{k=1}^{n(t)} cpu^k(t) = \overline{CPU(t)}$. To find the minimum number of
 802 $n(t)$, we sort the CPU consumption of containers in a descending order and
 803 greedily select containers from the highest CPU consumption to the lowest.
 804 In such a situation above, Eq. (.6) can reach its minimum accordingly.

□

805 Similarly, we can calculate the maximum and minimum values of memory
 806 load imbalance degree of a server (*i.e.*, the sum of pairwise multiplication
 807 $mul_{n(t)}$ of $mem^k(t)$ on a server in Eq. (15)). To achieve that, we simply
 808 substitute $mem^k(t)$ for $cpu^k(t)$ and $\overline{MEM(t)}$ for $\overline{CPU(t)}$ in Lemma 1. As a
 809 result, we combine the CPU and memory resources together as in Eq. (15),
 810 and the maximum value of the load imbalance degree of a server in the
 811 “ideal” load-balanced state turns out to be Eq. (16). The minimum value

812 of the load imbalance degree of an “ideal” server is calculated as Eq. (17),
813 when such a server hosts the minimum number of containers that satisfies
814 $\sum_{k=1}^{n(t)} cpu^k(t) \leq \overline{CPU(t)}$ and $\sum_{k=1}^{n(t)} mem^k(t) \leq \overline{MEM(t)}$ (*i.e.*, by jointly
815 considering CPU and memory resources) in practice.

816 Finally, to achieve load balancing, we set the threshold of $load_i(t)$ for
817 migration source server (*i.e.*, $load_{src}(t)$) as the maximum value of the load
818 imbalance degree of an “ideal” server in Eq. (16). Meanwhile, we set the
819 threshold of $load_i(t)$ for migration destination server (*i.e.*, $load_{dest}(t)$) as the
820 minimum value of the load imbalance degree of an “ideal” server in Eq. (17).

□