






# Scheduling Multi-Server Jobs With Sublinear Regrets via Online Learning

Hailiang Zhao , Shuiguang Deng , *Senior Member, IEEE*, Zhengzhe Xiang , *Member, IEEE*, Xueqiang Yan, Jianwei Yin, Schahram Dustdar , *Fellow, IEEE*, and Albert Y. Zomaya , *Fellow, IEEE*

**Abstract**—Multi-server jobs that request multiple computing resources and hold onto them during their execution dominate modern computing clusters. When allocating the multi-type resources to several co-located multi-server jobs simultaneously in online settings, it is difficult to make the tradeoff between the parallel computation gain and the internal communication overhead, apart from the resource contention between jobs. To study the computation-communication tradeoff, we model the computation gain as the speedup on the job completion time when it is executed in parallelism on multiple computing instances, and fit it with utilities of different concavities. Meanwhile, we take the dominant communication overhead as the penalty to be subtracted. To achieve a better gain-overhead tradeoff, we formulate an cumulative reward maximization program and design an online algorithm, named OGASCHED, to schedule multi-server jobs. OGASCHED allocates the multi-type resources to each arrived job in the ascending direction of the reward gradients. It has several parallel sub-procedures to accelerate its computation, which greatly reduces the complexity. We proved that it has a sublinear regret with general concave rewards. We also conduct extensive trace-driven simulations to validate the performance of OGASCHED. The results demonstrate that OGASCHED outperforms widely used heuristics by 11.33%, 7.75%, 13.89%, and 13.44%, respectively.

**Index Terms**—Multi-server job, online gradient ascent, online scheduling, regret analysis.

## I. INTRODUCTION

IN TODAY'S computing clusters, whether in the cloud data centers or at the network edge, many jobs request multiple resources (CPUs, GPUs, etc.) simultaneously and hold onto them

Manuscript received 4 May 2023; revised 6 July 2023; accepted 4 August 2023. Date of publication 8 August 2023; date of current version 12 June 2024. This work was supported in part by the Key Research Project of Zhejiang Province under Grant 2022C01145 and in part by the National Science Foundation of China under Grants U20A20173 and 62125206. Recommended for acceptance by E. Damiani. (*Corresponding author: Shuiguang Deng.*)

Hailiang Zhao and Shuiguang Deng are with the Hainan Institute of Zhejiang University, Sanya 572025, China, and also with the College of Computer Science and Technology, Zhejiang University, Hangzhou 310027, China (e-mail: hliangzhao@zju.edu.cn; dengsg@zju.edu.cn).

Zhengzhe Xiang is with the School of Computer Science and Technology, Hangzhou City University, Hangzhou 310015, China (e-mail: xiangzz@zucc.edu.cn).

Xueqiang Yan is with Huawei Technologies Company Ltd, Shanghai 201206, China (e-mail: yanxueqiang1@huawei.com).

Jianwei Yin is with the College of Computer Science and Technology, Zhejiang University, Hangzhou 310027, China (e-mail: zjuyjw@zju.edu.cn).

Schahram Dustdar is with Distributed Systems Group, Technische Universität Wien, 1040 Vienna, Austria (e-mail: dustdar@dsg.tuwien.ac.at).

Albert Y. Zomaya is with the School of Computer Science, University of Sydney, Sydney, NSW 2006, Australia (e-mail: albert.zomaya@sydney.edu.au). Digital Object Identifier 10.1109/TSC.2023.3303344

during their executions. For example, graph computations [1], federated learning [2], distributed DNN model trainings [3], etc. In this paper, we refer to these jobs as *multi-server jobs* [4], [5]. Multi-server jobs of diverse resource requirements arrive at the cluster online, which puts great pressure to current resource allocation policies to achieve a high computation efficiency.

When allocating the multi-type resources to several co-located multi-server jobs simultaneously in online settings, it is difficult to make the tradeoff between the *parallel computation gain* and the *internal communication overhead*, apart from the resource contention between jobs. Here the parallel computation gain refers to the speedup on the job completion time when it is executed in parallelism on multiple computing instances, which could be modeled with a function of the allocated multi-type resources [6]. Correspondingly, the internal communication overhead refers to the cost caused by non-computation operations such as data synchronization, averaging, message passing, etc., between the distributed workers. To achieve better computation efficiency, we need to consider the following key challenges.

- *Resource contention with service locality*: With service locality, a multi-server job can only be processed by a subset of computing instances where the resource requirements, session affinity [7], and other obligatory constraints are satisfied. When several multi-server jobs arrive simultaneously, how to allocate the limited resources to them without degenerating the computation efficiency is challenging.
- *Unknown arrival patterns of jobs*: In real-life scenarios, the resource allocation should be made online without the knowledge of future job arrivals. The lack of information on the problem space could lead to a solution far from the global optimum.
- *The parallel computation gain does not increase in a linear rate with the quantity of allocated resources*: For instance, in distributed DNN model training or federated learning, adding workers (that request more resources) does not improve the training speed linearly [3], [6]. This is because the overhead of all-reduce operation between workers or the averaging of local gradients increase with the number of participated workers, especially when the workers are distributed in different machines and communicate with each other through network [8], [9], [10]. Compared with high-speed intra-node communication channels such as NVLink, the inter-node bandwidth through NIC is relatively much slower. Another example is graph computation. Without a well-designed graph partition policy, the

speedup of message-passing between graph nodes can be significantly slowed down [1], [11].

- *The type of resource which dominates the communication overhead varies to different job types:* For example, in graph computation jobs, the dominant communication overhead lies in the internal input-output data transferring between the interdependent CPU- and memory-intensive tasks [12]. However, the dominant overhead of the distributed training of DNNs lies in the data averaging and synchronizing between the GPU-intensive workers through network [13]. This variety greatly complicates the theoretical analysis for the gain-overhead tradeoff.

Despite the vast literature on the online resource allocation algorithms [3], [6], [13], [14], [15], [16], [17], [18], their model formulation and theoretical analysis which places emphasis on the gain-overhead tradeoff is limited. To fill the theoretical gap, in this paper, we propose an online scheduling algorithm, termed as OGASCHED, to *learn* to allocate multi-type resources to co-located multi-server jobs online to maximize the overall computation efficiency. We try to analyze the tradeoff in a generic way. The generality is embodied in the following points. First of all, different from the specific works on deep learning jobs [3], [6], [13] or query jobs [12], we allow different types of multi-server jobs to *co-locate* in the cluster which consists of heterogeneous computing resources. Different job types can have different resource requirements while different computing instances can be equipped with diverse quantities and types of resources. Second, we adopt general zero-startup non-decreasing utility functions to model the parallel computation gain in terms of the job completion time. Compared to existing literature, we allow the utilities to be diverse in their *level of concavity*. Specifically, we provide both analysis and experiments on linear, polynomial, logarithmic, and reciprocal utilities. Third, we makes no assumptions on the arrival patterns of multi-server jobs. OGASCHED requests no knowledge on the job arrival distributions but tries to learn them to make better scheduling decisions.

In our model formulation, the computation efficiency is modeled in the way of cumulative reward. Time is slotted, and the cumulative reward is obtained by summing up the reward in each time slot, where a single-time reward is a linear aggregation of each job's reward. Further, a job's reward at each time is designed as the achieved parallel computation gain aggregated over the allocated resources minus the penalty introduced by the dominant communication overhead. At each time, OGASCHED allocates resources to each arrived job in the direction that *makes the gradient of the reward increase*. OGASCHED is capable of handling high dimensional inputs in stochastic scenarios with unpredictable behaviors. We adopt regret, i.e., the gap on the cumulative reward between the proposed online algorithm and the offline optimum achieved by an oracle [19], to analyze the performance lower bound of OGASCHED. We prove that, OGASCHED has a State-of-the-Art (SOTA) regret, which is sub-linear with the time slot length and the number of job types. This work fulfills one of the key deficiencies of the past works in the modeling and analysis of the gain-overhead tradeoff for multi-server jobs. The contributions are summarized as follows.

- We systematically study the resource allocation of co-located multi-server jobs in terms of the tradeoff between the parallel computation gains and the internal communication overheads. Our study is general in scenario settings and it sufficiently takes the characters of the diminishing marginal effect of gains into consideration.
- We propose an algorithm, i.e., OGASCHED, to learn to strike a balanced computation-communication tradeoff. OGASCHED has no assumptions on the job arrival patterns. With a nice setup (defined in Section III-A), OGASCHED achieves a SOTA regret  $\mathcal{O}(\mathcal{H}_G \cdot \sqrt{T})$  for general concave non-linear rewards, where  $T$  is the time slot length, and  $\mathcal{H}_G$  (formally defined in (52)) is parameter that characterizes the bipartite graph model. OGASCHED is accelerated by well-designed parallel sub-procedures. The parallelism helps yield a complexity of  $\mathcal{O}(\log(K))$ , where  $K$  is the number of resource types.
- We conduct extensive trace-driven simulations to validate the performance of OGASCHED. The simulation results show that OGASCHED outperforms widely used heuristics including DRF [20], FAIRNESS, BINPACKING, and SPREADING by 11.33%, 7.75%, 13.89%, and 13.44%, respectively. We also provide large-scale validations.

The rest of this paper is organized as follows. We formulate the online scheduling problem for multi-server jobs in Section II. We then present the design details of OGASCHED with regret analysis and discuss its extensions in Section III. We demonstrate the experimental results in Section IV, and discuss related works in Section V. Finally, we conclude this paper in Section VI.

## II. BIPARTITE SCHEDULING WITH REGRETS

We consider a cluster of heterogenous computing instances serving several types of multi-server jobs. Here the computing instances can be VMs in clouds, or local servers at the network edge. The computing instances work collaboratively to provide resources to serve the considered jobs. Different computing instances are equipped with different types and quantities/specifications of resources, including CPU cores, memory, bandwidth, GPUs, etc. Jobs of different types can have different demands on them. Key notations used in this paper are summarized in Table I.

### A. Online Bipartite Scheduling

We use a bipartite graph  $\mathcal{G} = (\mathcal{L}, \mathcal{R}, \mathcal{E})$  to model the job-server constraints, as shown in Fig. 1. In graph  $\mathcal{G}$ ,  $\mathcal{L}$  is the set of job types and indexed by  $l$  while  $\mathcal{R}$  is the set of computing instances and indexed by  $r$ . The connections between the job types and the computing instances are recorded in  $\mathcal{E}$ . Because of the job-server constraints, type- $l$  job may only be served by a subset of  $\mathcal{R}$ . We denote the subset by

$$\mathcal{R}_l = \{r \in \mathcal{R} \mid (l, r) \in \mathcal{E}\}. \quad (1)$$

Similarly, we use

$$\mathcal{L}_r = \{l \in \mathcal{L} \mid (l, r) \in \mathcal{E}\}. \quad (2)$$

TABLE I  
SUMMARY OF KEY NOTATIONS

NOTATION	DESCRIPTION
$\mathcal{T}$	Time horizon of length $T$
$\mathcal{G} = (\mathcal{L}, \mathcal{R}, \mathcal{E})$	The bipartite graph
$l \in \mathcal{L}$	A job type (port)
$r \in \mathcal{R}$	A computing instance
$(l, r) \in \mathcal{E}$	The edge (channel) between $l$ and $r$
$\forall r : \mathcal{L}_r$	The set of job types connect to $r$
$\forall l : \mathcal{R}_l$	The set of computing instances connect to $l$
$\mathbf{x}(t)$	The job arrival status at time $t$
$\mathbf{y}(t)$	The scheduling decision at time $t$
$\mathcal{K}$	The set of different types of resources
$\forall l : \mathbf{a}_l$	Resource requirements of type- $l$ job
$\forall r, k : c_r^k$	The number of type- $k$ resources equipped by $r$
$q(\mathbf{x}(t), \mathbf{y}(t))$	The reward of time $t$
$\forall k : f_k(\cdot)$	Computation gain of type- $k$ devices
$\forall k : \beta_k \in [0, 1]$	Coefficient of type- $k$ communication overhead

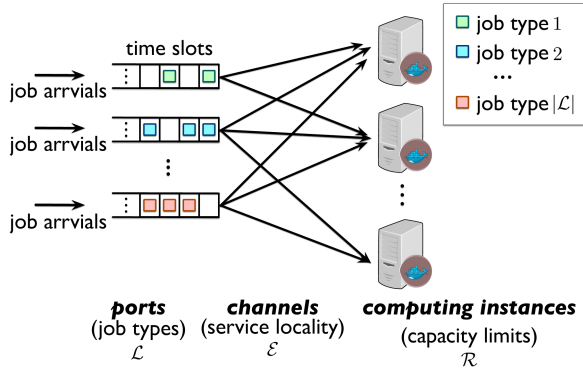


Fig. 1. The bipartite graph model for online job scheduling.

to represent the set of job types that connect to computing instance  $r$ . We designate each job type  $l \in \mathcal{L}$  as *port* and each connection  $(l, r) \in \mathcal{E}$  as *channel*.  $\mathcal{G}$  is called right  $d$ -regular iff the indegree of each right vertex is  $d$ , i.e.,  $\forall r \in \mathcal{R}, |\mathcal{L}_r| = d$ .

Time is discretized, and at each time  $t \in \mathcal{T} \triangleq \{1, \dots, T\}$ , from each port, at most one job yields. Let us denote by

$$\mathbf{x}(t) = [x_l(t)]_{l \in \mathcal{L}} \in \{0, 1\}^{|\mathcal{L}|}, \quad (3)$$

the job arrival status at time  $t$ . We do not make any assumption on the job arrival patterns or distributions. The cluster has  $K$  types of resources, and computing instance  $r$  has  $c_r^k$  type- $k$  resources, where  $k \in \mathcal{K} \triangleq \{1, 2, \dots, K\}$ . For each type- $l$  job, we denote its maximum requests on each resource by  $\mathbf{a}_l = [a_l^k]_{k \in \mathcal{K}} \in \mathbb{N}^{|\mathcal{K}|}$ . At time  $t$ , we use

$$\mathbf{y}(t) = \left[ y_{(l,r)}^k(t) \right]_{l \in \mathcal{L}, r \in \mathcal{R}_l, k \in \mathcal{K}} \in \mathbb{R}_{\geq 0}^{\sum_{l \in \mathcal{L}} |\mathcal{R}_l| \times K}, \quad (4)$$

to denote the scheduling decision. Here we allow  $y_{(l,r)}^k(t)$  to be fractional. Taking GPU as example, Machine-Learning-as-a-Service (MLaaS) platforms support GPU sharing in a space- and time-multiplexed manner by intercepting CUDA APIs [21], [22], [23].

The first constraint is that, through each channel, a job should not be allocated with resources more than it requires. Formally, we have

$$0 \leq y_{(l,r)}^k(t) \leq a_l^k, \forall l, r, k, t. \quad (5)$$

The second constraint  $\mathbf{y}(t)$  should satisfy is that, the resources allocated out from any computing instance  $r$  should not more than it has

$$\sum_{l \in \mathcal{L}_r} y_{(l,r)}^k(t) \leq c_r^k, \forall r, k, t. \quad (6)$$

We denote by  $\mathcal{Y} \triangleq \{\mathbf{y} \in \mathbb{R}^{\sum_{l \in \mathcal{L}} |\mathcal{R}_l| \times K} \mid (5) \text{ and } (6) \text{ hold}\}$  to represent the solution space from here on.

### B. Computation-Communication Tradeoff

The performance metric we use for online bipartite scheduling is designed as the gain obtained by the parallel computation through multi-type resources minus the penalty introduced by the dominant communication overheads. Specifically, we denote by  $q_l(\mathbf{x}(t), \mathbf{y}(t))$  the reward of port  $l$  at time  $t$ , and it is formulated as

$$q_l(\mathbf{x}(t), \mathbf{y}(t)) = x_l(t) \left[ \sum_{k \in \mathcal{K}} f_k \left( \sum_{r \in \mathcal{R}_l} y_{(l,r)}^k(t) \right) - \max_{k \in \mathcal{K}} \left\{ \beta_k \sum_{r \in \mathcal{R}_l} y_{(l,r)}^k(t) \right\} \right]. \quad (7)$$

In this formulation, the first part,  $\sum_{k \in \mathcal{K}} f_k(\sum_{r \in \mathcal{R}_l} y_{(l,r)}^k(t))$ , is the parallel computation gain, which is linearly aggregated over each type of resource, in proportional to each resource's weight. Jobs of different types can have different combinations of weights.  $f_k(\cdot)$  is the gain achieved by  $\sum_{r \in \mathcal{R}_l} y_{(l,r)}^k(t)$  type- $k$  resources collaboratively, where  $f_k(\cdot)$  is a zero-startup concave utility defined in  $\mathbb{R}_{\geq 0}$ . Note that  $\sum_{r \in \mathcal{R}_l} y_{(l,r)}^k(t)$  is the quota of the type- $k$  resources allocated to the type- $l$  job at  $t$ . As we have analyzed before,  $\{f_k(\cdot)\}_{k \in \mathcal{K}}$  are non-decreasing concave functions because the marginal effect of parallel computation decreases successively when increasing participated resources [3], [24]. We expect  $\{f_k(\cdot)\}_{k \in \mathcal{K}}$  to be *continuously differentiable* because it helps design a policy that yields a nice lower bound of the reward. The details will be demonstrated in Section III-C. If  $\{f_k(\cdot)\}_{k \in \mathcal{K}}$  are not differentiable everywhere, we can apply subgradient ascent-related techniques in the policy design. The second part in (7) is  $\max_{k \in \mathcal{K}} \{\beta_k \sum_{r \in \mathcal{R}_l} y_{(l,r)}^k(t)\}$ , which reflects the dominant weighted communication overheads over different types of resources. For example, in federated learning at the edge, the dominant communication overhead lies in the averaging and synchronizing of data between each edge server over the network [25]. Another example is graph computation, in which the job is organized into a direct acyclic graph (DAG), and the dominant communication overhead falls into the data & message passing between CPU- and memory-intensive tasks [12].  $\{\beta_k\}_{k \in \mathcal{K}}$  are the coefficients to balance the gain and the overhead. W.L.O.G., we set each  $\beta_k \in [0, 1]$ . Theoretically, the second part of (7) is a penalty, the minimization of which guides the



scheduling decisions to balance the communication overheads of different device types. Our reward design encourages each job to be served with the balance between the computation gain and the communication overhead being achieved.

### C. Regret Minimizing

Based on the above, we define the overall reward at time  $t$  as the linear aggregation over each port

$$q(\mathbf{x}(t), \mathbf{y}(t)) = \sum_{l \in \mathcal{L}} q_l(\mathbf{x}(t), \mathbf{y}(t)). \quad (8)$$

The cumulative reward of scheduling policy  $\pi$  over the time horizon  $\mathcal{T}$  is obtained by summing up the rewards obtained at each time until  $T$

$$Q^\pi(\{\mathbf{x}(t)\}_1^T, \{\mathbf{y}(t)\}_1^T) = \sum_{t \in \mathcal{T}} q(\mathbf{x}(t), \mathbf{y}(t)), \quad (9)$$

where the scheduling decisions  $\{\mathbf{y}(t)\}_1^T$  are made under the guidance of policy  $\pi$ . In the following, we just use  $Q$  and drop the superscript  $\pi$  for simplification.

We do not make any assumption on the distribution of the job arrival trajectory  $\{\mathbf{x}(t)\}_1^T$ . To obtain a non-trivial performance measure, we cast the multi-server bipartite scheduling problem into the framework of online learning, which prompts us to compare the performance of the online policy  $\pi$  with the best offline stationary policy  $\pi^*$  [26], [27]. Let us denote by  $\mathbf{y}^*$  the optimal offline stationary resource allocation decision guided by policy  $\pi^*$ , i.e.,

$$\mathbf{y}^* = \arg \sup_{\mathbf{y} \in \mathcal{Y}} Q(\{\mathbf{x}(t)\}_1^T, \mathbf{y}), \quad (10)$$

Physically,  $\mathbf{y}^*$  is the optimal stationary resource reservation decisions for each port whatever the actual job arrival status  $\mathbf{x}(t)$  is. Formally, we define the regret  $R_T^\pi(\{\mathbf{x}(t)\}_1^T)$  for the job arrival trajectory  $\{\mathbf{x}(t)\}_1^T$  as

$$R_T^\pi(\{\mathbf{x}(t)\}_1^T) \triangleq Q(\{\mathbf{x}(t)\}_1^T, \mathbf{y}^*) - Q(\{\mathbf{x}(t)\}_1^T, \{\mathbf{y}(t)\}_1^T).$$

The regret of policy  $\pi$  is further defined as the maximum regret achieved over every possible job arrival trajectory

$$R_T^\pi \triangleq \sup_{\{\mathbf{x}(t)\}_1^T} R_T^\pi(\{\mathbf{x}(t)\}_1^T). \quad (11)$$

Our goal is to find a policy  $\pi$ , under which a sequence of bipartite scheduling decisions  $\{\mathbf{y}(t)\}_1^T$  is yielded, to minimize  $R_T^\pi$ .

## III. ONLINE GRADIENT ASCENT

To minimize the regret  $R_T^\pi$ , we resort to an online variant of the gradient-based methods, online gradient ascent (OGA) [28]. A series of recent works have demonstrated that OGA achieves the best possible regret for online caching problems in different network settings when the rewards are linear [27], [29], [30], [31]. In this paper, we extend OGA to the online bipartite scheduling problem for multi-server jobs with non-linear rewards. Before presenting the design details, we first give some preliminary definitions and analysis.

### A. Preliminaries

*Definition 1. NICE SETUP:* If all the utilities  $\{f_k\}_{k \in \mathcal{K}}$  are (i) linearly separable over computing instances, i.e.,

$$f_k \left( \sum_{r \in \mathcal{R}_l} y_{(l,r)}^k \right) = \sum_{r \in \mathcal{R}_l} f_r^k \left( y_{(l,r)}^k \right), \quad (12)$$

and each concave utility  $f_r^k(\cdot)$  is (ii) continuously differentiable in  $\mathbb{R}_+$ , and (iii) there exist  $\varpi_r^k > 0$  such that

$$(f_r^k)'(0) \leq \varpi_r^k, \forall r, k, \quad (13)$$

we say this is a nice setup.

The following proposition demonstrates the property of the regret minimization problem, which will be used in the design and analysis of OGASCHED.

*Proposition 1. CONVEXITY:* (i) The feasible solution space  $\mathcal{Y}$  is convex. (ii) With a nice setup, at each time  $t$ , the single-slot reward function  $q(\mathbf{x}(t), \mathbf{y}(t))$  is a concave function of  $\mathbf{y}(t)$ .

*Proof:* In the following proof, we just drop  $(t)$  from  $\mathbf{x}(t)$  and  $\mathbf{y}(t)$  for simplification. Besides, we only prove the case that  $\mathcal{G}$  is right  $d$ -regular and  $d = |\mathcal{L}|$ . The left cases can be easily proved with the same techniques used in this proof.

We first prove (i). To do this, let us arrange the vector  $\mathbf{y}$  as

$$\mathbf{y} = \left[ \underbrace{\mathbf{y}^1}_{k=1}; \underbrace{\mathbf{y}^2}_{k=2}; \dots; \underbrace{\mathbf{y}^K}_{k=K} \right], \quad (14)$$

where  $\mathbf{y}^k \in \mathbb{R}^{(|\mathcal{L}| \times |\mathcal{R}|)}$  is arranged as

$$\mathbf{y}^k = \left[ \underbrace{y_{(1,1)}^k; \dots; y_{(1,|\mathcal{R}|)}^k}_{l=1}; \dots; \underbrace{y_{(|\mathcal{L}|,1)}^k; \dots; y_{(|\mathcal{L}|,|\mathcal{R}|)}^k}_{l=|\mathcal{L}|} \right]. \quad (15)$$

With this arrangement, the vector representation of (5) is

$$\mathbf{0} \leq \mathbf{y} \leq \mathbf{a}, \quad (16)$$

where  $\mathbf{a} = [\mathbf{a}^1; \dots; \mathbf{a}^K]$ , and

$$\mathbf{a}^k = \left[ \underbrace{a_1^k; \dots; a_{|\mathcal{R}|}^k}_{\text{of size } |\mathcal{R}|}; \dots; \underbrace{a_{|\mathcal{L}|}^k; \dots; a_{|\mathcal{L}|}^k}_{\text{of size } |\mathcal{R}|} \right]. \quad (17)$$

Similarly, we want to construct a matrix  $\mathbf{B}$  and a vector  $\mathbf{c}$  for the vector representation of (6).  $\forall k \in \mathcal{K}$ , we design  $\mathbf{B}' \in \mathbb{R}^{|\mathcal{R}| \times (|\mathcal{L}| \times |\mathcal{R}|)}$  as

$$[\mathbf{B}']_{ij} = \begin{cases} 1 & (j-i) \mid |\mathcal{R}| \\ 0 & \text{o.w.} \end{cases} \quad (18)$$

and  $\mathbf{c}^k = [c_1^k; \dots; c_{|\mathcal{R}|}^k]$ . Then, we have

$$\mathbf{B}' \mathbf{y}^k \leq \mathbf{c}^k, \forall k. \quad (19)$$

We can transform (19) into

$$\mathbf{B}' \mathbf{y}^k + \sum_{k' \neq k} \mathbf{O} \mathbf{y}^{k'} \leq \mathbf{c}^k, \forall k, \quad (20)$$

where  $\mathbf{O}$  is a zero matrix. As a result, (6) is equivalent to

$$\mathbf{B} \mathbf{y} \leq \mathbf{c}, \quad (21)$$

where  $\mathbf{c} = [\mathbf{c}^1; \dots; \mathbf{c}^K] \in \mathbb{R}^{(|\mathcal{R}| \times K)}$ , and

$$\mathbf{B} = \text{diag}(\mathbf{B}') \in \mathbb{R}^{(|\mathcal{R}| \times K) \times (|\mathcal{L}| \times |\mathcal{R}| \times K)}.$$

The above analysis leads to  $\mathcal{Y} = \{\mathbf{y} \mid \mathbf{0} \leq \mathbf{y} \leq \mathbf{a}, \mathbf{B}\mathbf{y} \leq \mathbf{c}\}$  being a polyhedron, which is well known to be a convex set.

We now prove (ii). Similarly, we try to find the vectorized representation of  $q(\mathbf{x}, \mathbf{y})$ . To do this, we define the operator  $\mathbf{f} : \mathbb{R}^{|\mathcal{L}| \times |\mathcal{R}| \times K} \rightarrow \mathbb{R}^{|\mathcal{L}| \times |\mathcal{R}| \times K}$  as

$$\mathbf{f} = [\mathbf{f}^1; \dots; \mathbf{f}^K], \quad (22)$$

where

$$\mathbf{f}^k = \left[ \underbrace{f_1^k; \dots; f_{|\mathcal{R}|}^k}_{\text{make } |\mathcal{L}| \text{ replicas}}; \dots; \underbrace{f_1^k; \dots; f_{|\mathcal{R}|}^k}_{l=|\mathcal{L}|} \right]. \quad (23)$$

Then, the first part of (7) can be transformed into

$$\begin{aligned} \sum_{l \in \mathcal{L}} x_l \sum_{k \in \mathcal{K}} f_k \left( \sum_{r \in \mathcal{R}_l} y_{(l,r)}^k \right) &= \sum_{l \in \mathcal{L}} \sum_{k \in \mathcal{K}} \sum_{r \in \mathcal{R}_l} x_l f_r^k \left( y_{(l,r)}^k \right) \\ &= \boldsymbol{\chi} \cdot \mathbf{f}(\mathbf{y}), \end{aligned} \quad (24)$$

where

$$\boldsymbol{\chi} = \left[ \underbrace{x_1; \dots; x_1}_{l=1, \forall r}; \dots; \underbrace{x_{|\mathcal{L}|}; \dots; x_{|\mathcal{L}|}}_{l=|\mathcal{L}|, \forall r}; \dots; \underbrace{x_1; \dots; x_{|\mathcal{L}|}}_{k=K} \right]. \quad (25)$$

For the second part of (7), we have

$$\sum_{l \in \mathcal{L}} x_l \max_{k \in \mathcal{K}} \left\{ \beta_k \sum_{r \in \mathcal{R}_l} y_{(l,r)}^k \right\} = \sum_{l \in \mathcal{L}} x_l \beta_{k^*} \sum_{r \in \mathcal{R}_l} y_{(l,r)}^{k^*}, \quad (26)$$

where

$$k^* = \arg \max_{k \in \mathcal{K}} \left\{ \beta_k \sum_{r \in \mathcal{R}_l} y_{(l,r)}^k \right\}. \quad (27)$$

Without loss of generality, we assume that  $k^* = 1$ . Then, the second part can be represented as  $\boldsymbol{\beta} \cdot \mathbf{y}$ , where

$$\boldsymbol{\beta} = \left[ \underbrace{x_1 \beta_1; \dots; x_1 \beta_1}_{l=1, \forall r}; \dots; \underbrace{x_{|\mathcal{L}|} \beta_1; \dots; x_{|\mathcal{L}|} \beta_1}_{l=|\mathcal{L}|, \forall r}; \underbrace{\mathbf{0}}_{\forall k \neq k^*} \right]. \quad (28)$$

The above analysis leads to

$$q(\mathbf{x}, \mathbf{y}) = \boldsymbol{\chi} \cdot \mathbf{f}(\mathbf{y}) - \boldsymbol{\beta} \cdot \mathbf{y}. \quad (29)$$

With the concavity of  $f_r^k(\cdot)$ , the result (ii) is immediate.  $\square$

As a result, the derivative of  $q(\cdot)$  at time  $t$  is

$$\frac{\partial q(\mathbf{x}(t), \mathbf{y}(t))}{\partial y_{(l,r)}^k(t)} = \begin{cases} x_l(t) \left( (f_r^k)' \left( y_{(l,r)}^k(t) \right) - \beta_k \right) & k = k^* \\ x_l(t) (f_r^k)' \left( y_{(l,r)}^k(t) \right) & \text{o.w.,} \end{cases} \quad (30)$$

where  $k^*$  is defined in (27)

### B. Online Gradient Ascent

In this section, we give the design details of the OGA-based bipartite scheduling policy.

**Definition 2. THE OGA POLICY:** For any feasible initial bipartite scheduling decision  $\mathbf{y}(1) \in \mathcal{Y}$ , at each time  $t \in \mathcal{T}$ , the OGA

policy gets  $\mathbf{y}(t+1)$  in the direction of ascending the gradient of  $q(\mathbf{x}(t), \mathbf{y}(t))$

$$\mathbf{y}(t+1) = \Pi_{\mathcal{Y}}(\mathbf{y}(t) + \eta_t \nabla q(\mathbf{x}(t), \mathbf{y}(t))), \quad (31)$$

where  $\eta_t$  is the step size, and

$$\Pi_{\mathcal{Y}}(\mathbf{z}) = \arg \min_{\hat{\mathbf{y}} \in \mathcal{Y}} \|\hat{\mathbf{y}} - \mathbf{z}\|_2^2, \quad (32)$$

is the euclidean projection of  $\mathbf{z}$  onto  $\mathcal{Y}$ .

To implement the projection (32) with low complexity, we propose OGASCHED, which is a combination of the OGA policy and the following fast projection technique. First, we introduce the Lagrangian of the projection (32) as

$$\begin{aligned} L(\hat{\mathbf{y}}, \boldsymbol{\rho}, \boldsymbol{\mu}, \boldsymbol{\lambda}) &= \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \sum_{k \in \mathcal{K}} \left( \hat{y}_{(l,r)}^k - z_{(l,r)}^k \right)^2 \\ &+ \sum_{r \in \mathcal{R}} \sum_{k \in \mathcal{K}} \rho_r^k \left( \sum_{l \in \mathcal{L}_r} \hat{y}_{(l,r)}^k - c_r^k \right) - \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \sum_{k \in \mathcal{K}} \lambda_{l,r}^k \hat{y}_{(l,r)}^k \\ &+ \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \sum_{k \in \mathcal{K}} \mu_{l,r}^k \left( \hat{y}_{(l,r)}^k - a_l^k \right), \end{aligned} \quad (33)$$

where  $\boldsymbol{\rho}$  is the dual variable for (6),  $\boldsymbol{\mu}$  is the dual variable for  $\mathbf{y}(t) \leq \mathbf{a}$ , and  $\boldsymbol{\lambda}$  is the dual variable for  $\mathbf{y}(t) \geq \mathbf{0}$ . Then, we can write the KKT conditions of the projection as

$$2(\hat{y}_{(l,r)}^k - z_{(l,r)}^k) + \rho_r^k - \lambda_{(l,r)}^k + \mu_{(l,r)}^k = 0 \quad (34)$$

$$\sum_{l \in \mathcal{L}_r} \hat{y}_{(l,r)}^k = c_r^k \ \& \ \rho_r^k > 0 \quad (35)$$

$$\hat{y}_{(l,r)}^k = a_l^k \ \& \ \mu_{l,r}^k > 0 \quad (36)$$

$$\hat{y}_{(l,r)}^k = 0 \ \& \ \lambda_{l,r}^k > 0, \quad (37)$$

for every  $l, r, k$ .

Our fast projection is implemented for each pair of  $(r, k)$  in parallel. Specifically, for each  $r \in \mathcal{R}$  and each  $k \in \mathcal{K}$ , we divide the ports  $l \in \mathcal{L}$  into three disjoint sets

$$\begin{cases} \mathcal{B}_{rk}^1 = \{l \in \mathcal{L}_r \mid \forall (l, r, k) : \hat{y}_{(l,r)}^k = a_l^k\} \\ \mathcal{B}_{rk}^2 = \{l \in \mathcal{L}_r \mid \forall (l, r, k) : \hat{y}_{(l,r)}^k = 0\} \\ \mathcal{B}_{rk}^3 = \{l \in \mathcal{L}_r \mid \forall (l, r, k) : 2(\hat{y}_{(l,r)}^k - z_{(l,r)}^k) + \rho_r^k = 0\}, \end{cases}$$

where

$$\rho_r^k = \frac{2}{|\mathcal{B}_{rk}^3|} \left( \sum_{l \in \mathcal{B}_{rk}^3} z_{(l,r)}^k - c_r^k + \sum_{l \in \mathcal{B}_{rk}^1} a_l^k \right), \forall r, k. \quad (38)$$

The fast projection works by solving the equation system (35) ~ (38) iteratively. Specifically, for each pair of  $(r, k)$ , we sort the elements of  $z_{(\cdot, r)}^k$  in descending order (step 7), and initialize  $\mathcal{B}_{rk}^1$  and  $\mathcal{B}_{rk}^2$  as  $\emptyset$  while initializing  $\mathcal{B}_{rk}^3$  as  $\mathcal{L}_r$  (step 10 and 12). Then, we repeat a loop, in which we calculate  $\rho_r^k$  with (38), and update the value of  $\hat{y}_{(l,r)}^k$  for each port  $l$  in  $\mathcal{B}_{rk}^3$  (step 25). Since the elements of  $z_{(\cdot, r)}^k$  are sorted from largest to smallest, if some  $\hat{y}_{(l,r)}^k < 0$ , we can derive that for all the  $l' \in \mathcal{S}_{rk} := \{l, \dots, |\mathcal{L}_r|\}$ , we have  $\hat{y}_{(l', r)}^k < 0$ . Thus, the resource allocation

**Algorithm 1: OGASCHED.**


---

**Input:** Graph  $\mathcal{G}$ , requirements  $\mathbf{a}$ , capacities  $\mathbf{c}$ , and the decay  $\lambda$

**Output:** Scheduling decisions  $\{\mathbf{y}(t)\}_{t \in \mathcal{T}}$

- 1 Initialize  $\mathbf{y}(1) \in \mathcal{Y}$  and  $\eta_0$
- 2 **for**  $t$  from 1 to  $T$  **do**
- 3     Observe the job arrival status  $\mathbf{x}(t)$
- 4     Calculate the gradient  $\nabla q(\mathbf{x}(t), \mathbf{y}(t))$  with (30)
- 5      $\mathbf{z}(t+1) \leftarrow \mathbf{y}(t) + \eta_t \nabla q(\mathbf{x}(t), \mathbf{y}(t))$
- 6     **foreach**  $(r, k)$  in  $\text{zip}(\mathcal{R}, \mathcal{K})$  **do in parallel**
- 7         Sort the elements of  $\mathbf{z}_{(:,r)}^k(t+1)$  in descending order
- 8          $\text{initialized} \leftarrow \text{False}$
- 9         **while**  $\text{True}$  **do**
- 10              $\mathcal{B}_{rk}^2 \leftarrow \emptyset$
- 11             **if not**  $\text{initialized}$  **then**
- 12                  $\mathcal{B}_{rk}^1 \leftarrow \emptyset, \mathcal{B}_{rk}^3 \leftarrow \mathcal{L}_r, \hat{\mathbf{y}} \leftarrow \mathbf{0}$
- 13                  $\text{initialized} \leftarrow \text{True}$
- 14             **else**
- 15                 **if**  $\hat{y}_{(1,r)}^k > a_1^k$  **then**
- 16                      $\mathcal{B}_{rk}^1 \leftarrow \{1\}, \mathcal{B}_{rk}^3 \leftarrow \mathcal{L}_r \setminus \{1\}$
- 17                     **else**
- 18                         **break**
- 19             **repeat**
- 20                 Calculate  $\rho_r^k$  with (39)
- 21                 **for**  $l \in \mathcal{L}_r$  **do**
- 22                     **if**  $l \in \mathcal{B}_{rk}^1$  **then**
- 23                          $\hat{y}_{(l,r)}^k \leftarrow a_l^k$
- 24                     **else if**  $l \in \mathcal{B}_{rk}^3$  **then**
- 25                          $\hat{y}_{(l,r)}^k \leftarrow z_{(l,r)}^k(t+1) - \rho_r^k/2$
- 26                         **if**  $\hat{y}_{(l,r)}^k < 0$  **then**
- 27                              $\mathcal{S}_{rk} \leftarrow \{l, l+1, \dots, |\mathcal{L}_r|\}$
- 28                             **break**
- 29                 // Update then re-calculate
- 30                  $\mathcal{B}_{rk}^2 \leftarrow \mathcal{B}_{rk}^2 \cup \mathcal{S}_{rk}, \mathcal{B}_{rk}^3 \leftarrow \mathcal{B}_{rk}^3 \setminus \mathcal{S}_{rk}$
- 31             **until**  $\mathcal{S}_{rk} = \emptyset$ ;
- 32      $\mathbf{y}(t+1) \leftarrow \hat{\mathbf{y}}$
- 33      $\eta_{t+1} \leftarrow \lambda \eta_t$  // Update learning rate
- 34 **return** the sequence of decisions  $\{\mathbf{y}(t)\}_{t \in \mathcal{T}}$

---

for all the ports in  $\mathcal{S}_{rk}$  is illegal, since  $\hat{y}_{(l,r)}^k \geq 0$  must hold. As a result, we update the sets  $\mathcal{B}_{rk}^2$  and  $\mathcal{B}_{rk}^3$ , and repeat the calculate loop again (step 29). The calculation loop stops when there are no illegal resource allocations, i.e.,  $\forall l \in \mathcal{L}_r$ , we have  $\hat{y}_{(l,r)}^k \geq 0$ . In other words,  $\mathcal{S}_{rk} = \emptyset$ . We call the calculation loop in step 18~step 30 the inner loop. The outer loop is the while loop defined in step 9. To exit the while loop, we need to guarantee that  $\hat{y}_{(1,r)}^k \leq a_1^k$ . Otherwise, the resource allocation is also illegal. Note that here we only need to check for  $l = 1$  since the elements in  $\mathbf{z}_{(:,r)}^k$  are sorted.

The number of projections is linearly proportional to the size of the solution's dimensions, i.e.,  $\sum_{l \in \mathcal{L}} |\mathcal{R}_l| \times K$ . Nevertheless,

as we have mentioned, we can do the projections for different combinations of  $r$  and  $k$  in parallel because they are not interwoven. Thus, the time complexity of the fast projection is of  $\mathcal{O}(|\mathcal{L}| \times \log(K \sum_{l \in \mathcal{L}} |\mathcal{R}_l|))$  in each time slot, where the  $\log(\cdot)$  operator comes from the sorting operation (step 7). The multiplier  $|\mathcal{L}|$  outside  $\log(\cdot)$  comes from the inner loop (step 19). In our experiments, the repeat loop's execution count is significantly less than the number of job types  $|\mathcal{L}|$ .

### C. Regret Analysis

In this section, we discuss the regret of OGASCHED. The main result is summarized in Theorem 1.

*Theorem 1: REGRET UPPER BOUND.* With a nice setup, the regret of OGASCHED is upper bounded by

$$R_T^{\text{OGASCHED}} \leq \sqrt{2T \sum_{k \in \mathcal{K}} \sum_{r \in \mathcal{R}} \bar{a}^k c_r^k} \times \sqrt{\sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} ((\beta^*)^2 + K(\varpi_r^*)^2)}, \quad (39)$$

where  $\bar{a}^k := \max_{l \in \mathcal{L}} a_l^k$ ,  $\beta^* := \max_{k \in \mathcal{K}} \beta_k$ , and  $\varpi_r^* := \max_{k \in \mathcal{K}} \varpi_r^k$ .

*Proof:* The result is based on the non-expansiveness property of euclidean projection and the concavity of  $\{f_r^k(\cdot)\}_{r,k}$ . Our proof has two parts. The first part gives the general form of the upper bound, which is similar to Theorem 2.13 in [32] and Theorem 3 in [29]. Meanwhile, the second part gives the specific upper bounds of involved variables.

At each time  $t > 1$ , for the  $\mathbf{y}(t)$  yielded by OGASCHED, we have

$$\begin{aligned} \|\mathbf{y}(t) - \mathbf{y}^*\|^2 &= \|\Pi_{\mathcal{Y}}(\mathbf{y}(t-1) + \eta_t \nabla q(t-1)) - \mathbf{y}^*\|^2 \\ &\stackrel{(i)}{\leq} \|\mathbf{y}(t-1) - \mathbf{y}^*\|^2 + \eta_t^2 \|\nabla q(t-1)\|^2 \\ &\quad + 2\eta_t \nabla q(\mathbf{y}(t-1))^T (\mathbf{y}(t-1) - \mathbf{y}^*), \end{aligned} \quad (40)$$

where  $\nabla q(\mathbf{y}(t-1))$  is a shorthand for  $\nabla q(\mathbf{x}(t-1), \mathbf{y}(t-1))$ . (i) is because the non-expansiveness property of the euclidean projection. By moving  $\|\mathbf{y}(t-1) - \mathbf{y}^*\|^2$  to the LHS of (40) and summing the inequality telescopically over  $\mathcal{T}$ , we have

$$\begin{aligned} &\sum_{t=2}^{T+1} \nabla q(\mathbf{y}(t-1))^T (\mathbf{y}^* - \mathbf{y}(t-1)) \\ &\stackrel{(i)}{\leq} \frac{\eta \sum_{t=1}^T \|\nabla q(\mathbf{y}(t))\|^2}{2} + \frac{\|\mathbf{y}(1) - \mathbf{y}^*\|^2 - \|\mathbf{y}(T) - \mathbf{y}^*\|^2}{2\eta} \\ &\stackrel{(ii)}{\leq} \frac{\eta T (\max \|\nabla q\|)^2}{2} + \frac{\text{diam}(\mathcal{Y})^2}{2\eta}. \end{aligned} \quad (41)$$

Inequality (i) is because  $\forall t \in \mathcal{T}$  we set  $\eta_t \equiv \eta$ . In (ii), we use the fact that  $\|\mathbf{y}(T) - \mathbf{y}^*\| \geq 0$ . In (40),  $\max \|\nabla q\|$  is the maximum euclidean norm of the gradient of  $q(\mathbf{x}(t), \mathbf{y}(t))$  over every possible  $\mathbf{y}(t)$ , and  $\text{diam}(\mathcal{Y})$  is the largest euclidean distance between any two elements of  $\mathcal{Y}$ . Because  $q(\cdot)$  is a concave function of

$\mathbf{y}(t)$ , we have

$$\begin{aligned} R_T^{\text{OGASCHED}} &= \sup_{\forall \{\mathbf{x}(t)\}_1^T} \sum_{t=1}^T (q(\mathbf{x}(t), \mathbf{y}^*) - q(\mathbf{x}(t), \mathbf{y}(t))) \\ &\leq \sup_{\forall \{\mathbf{x}(t)\}_1^T} \sum_{t=1}^T \nabla q(\mathbf{y}(t))^T (\mathbf{y}^* - \mathbf{y}(t)) \triangleright (41) \\ &\leq \frac{\text{diam}(\mathcal{Y})^2}{2\eta} + \frac{\eta T (\max \|\nabla q\|)^2}{2}. \end{aligned} \quad (42)$$

In the following, we give the upper bound of  $\max \|\nabla q\|$  and  $\text{diam}(\mathcal{Y})$ , respectively.

1) *The upper bound of  $\max \|\nabla q\|$ .* With the result of (30), we have

$$\begin{aligned} \|\nabla q\|^2 &= \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \left[ x_l(t)^2 \left( (f_r^{k^*})' (y_{(l,r)}^{k^*}(t)) - \beta_{k^*} \right)^2 \right] \\ &\quad + \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \sum_{k \neq k^*} x_l(t)^2 (f_r^k)' (y_{(l,r)}^k(t))^2 \\ &= \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} x_l(t)^2 \left[ \sum_{k \in \mathcal{K}} \left( (f_r^k)' (y_{(l,r)}^k(t)) \right)^2 \right. \\ &\quad \left. - 2\beta_{k^*} (f_r^{k^*})' (y_{(l,r)}^{k^*}(t)) \right] + \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} x_l(t)^2 \beta_{k^*}^2. \end{aligned} \quad (43)$$

where  $k^*$  is defined in (27). The second part of (43) can be upper bounded by

$$\sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} x_l(t)^2 \beta_{k^*}^2 \leq \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} (\beta^*)^2, \quad (44)$$

where  $\beta^* = \max_{k \in \mathcal{K}} \beta_k$ . If  $\mathcal{G}$  is right  $d$ -regular, the bound reduces to  $d|\mathcal{R}|(\beta^*)^2$ . For the first part of (43), we use  $(f_r^{k^*})'$  to replace  $(f_r^{k^*})'(y_{(l,r)}^{k^*}(t))$  for simplification. Then we have

$$\begin{aligned} &\sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} x_l(t)^2 \left[ \sum_{k \in \mathcal{K}} \left( (f_r^k)' \right)^2 - 2\beta_{k^*} (f_r^{k^*})' \right] \\ &\leq \underbrace{\sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \sum_{k \neq k^*} \left( (f_r^k)' \right)^2}_{\text{PART-A}} + \underbrace{\sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} (f_r^{k^*})' \left( (f_r^{k^*})' - 2\beta_{k^*} \right)}_{\text{PART-B}}. \end{aligned}$$

For PART-A we have

$$\text{PART-A} \leq (K-1) \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} (\varpi_r^*)^2, \quad (45)$$

where  $\varpi_r^* = \max_{k \in \mathcal{K}} \varpi_r^k$ . If  $\mathcal{G}$  is right  $d$ -regular, the bound reduces to  $d|\mathcal{R}|(K-1)(\text{varpi}_r^*)^2$ . To analyze the upper bound of PART-B, we need to partition the computing instances into two disjoint sets

$$\begin{aligned} \mathcal{R}_1 &= \{r \in \mathcal{R} : \varpi_r^{k^*} \leq 2\beta_{k^*}\} \\ \mathcal{R}_2 &= \{r \in \mathcal{R} : \varpi_r^{k^*} > 2\beta_{k^*}\}. \end{aligned}$$

For each  $r \in \mathcal{R}_1$ , the maximum of  $(f_r^{k^*})'((f_r^{k^*})' - 2\beta_{k^*})$  is 0 since  $(f_r^{k^*})' \geq 0$  holds. For each  $r \in \mathcal{R}_2$ , the maximum is

$(\varpi_r^{k^*})^2 - 2\beta_{k^*} \varpi_r^{k^*}$ . Thus,

$$\text{PART-B} \leq \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l \cap \mathcal{R}_2} \left( (\varpi_r^{k^*})^2 - 2\beta_{k^*} \varpi_r^{k^*} \right). \quad (46)$$

Recall that in (46)  $\mathcal{R}_l$  is the set of computing instances that connects to port  $l$ . Because  $\beta_k \in [0, 1]$  holds for each  $k \in \mathcal{K}$ ,  $\forall l \in \mathcal{L}, r \in \mathcal{R}_l \cap \mathcal{R}_2$ , we have

$$(\varpi_r^{k^*})^2 - 2\beta_{k^*} \varpi_r^{k^*} \leq (\varpi_r^*)^2 - 2\beta_{k^*} \varpi_r^* \leq (\varpi_r^*)^2, \quad (47)$$

Finally, we can get

$$\|\nabla q\|^2 \leq \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \left( (\beta^*)^2 + K(\varpi_r^*)^2 \right). \quad (48)$$

For the upper bound in (48), all the computing instances  $r \in \mathcal{R}_l$  fall into the set  $\mathcal{R}_2$ .

2) *The upper bound of  $\text{diam}(\mathcal{Y})$ .* By definition we have

$$\text{diam}(\mathcal{Y}) = \sup_{\mathbf{y}, \mathbf{z} \in \mathcal{Y}} \|\mathbf{y} - \mathbf{z}\|. \quad (49)$$

To find the upper bound of  $\|\mathbf{y} - \mathbf{z}\|$ , we can get

$$\begin{aligned} \|\mathbf{y} - \mathbf{z}\|^2 &= \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \sum_{k \in \mathcal{K}} \left( y_{(l,r)}^k - z_{(l,r)}^k \right)^2 \\ &\stackrel{(i)}{\leq} \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \sum_{k \in \mathcal{K}} |y_{(l,r)}^k - z_{(l,r)}^k| \cdot a_l^k \\ &\leq \sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \sum_{k \in \mathcal{K}} a_l^k \left( y_{(l,r)}^k + z_{(l,r)}^k \right) \\ &\leq \sum_{k \in \mathcal{K}} \bar{a}^k \sum_{r \in \mathcal{R}} \left( \sum_{l \in \mathcal{L}_r} y_{(l,r)}^k + \sum_{l \in \mathcal{L}_r} z_{(l,r)}^k \right) \\ &\stackrel{(ii)}{\leq} 2 \sum_{k \in \mathcal{K}} \bar{a}^k \sum_{r \in \mathcal{R}} c_r^k, \end{aligned} \quad (50)$$

where  $\bar{a}^k = \max_{l \in \mathcal{L}} a_l^k$ . In (50), (i) is because the constraint (5). In (ii), we use the capacity constraint (6). As a result, we have

$$\text{diam}(\mathcal{Y}) \leq \sqrt{2 \sum_{k \in \mathcal{K}} \bar{a}^k \sum_{r \in \mathcal{R}} c_r^k}. \quad (51)$$

Combing the result (48) and (51), and set  $\eta$  as  $\frac{\text{diam}(\mathcal{Y})}{\|\nabla q\| \sqrt{T}}$ , we finally get the result.

The theorem shows that the suboptimality gap between OGASCHED and the offline optimal is of  $\Theta(\mathcal{H}_{\mathcal{G}} \times \sqrt{T})$ , where

$$\mathcal{H}_{\mathcal{G}} := \sqrt{2 \sum_{k \in \mathcal{K}} \sum_{r \in \mathcal{R}} \bar{a}^k c_r^k} \times \sqrt{\sum_{l \in \mathcal{L}} \sum_{r \in \mathcal{R}_l} \left( (\beta^*)^2 + K(\varpi_r^*)^2 \right)}, \quad (52)$$

is a factor characterized the scale of the bipartite graph  $\mathcal{G}$ . In addition, we can find that the regret grows sublinearly with the number of job types  $|\mathcal{L}|$ . To the best of our knowledge, this is the best regret for the online bipartite scheduling problem with non-linear rewards. The proof also indicates that, to achieve a not-too-bad cumulative reward, at each time  $t$ , the learning rate



$\eta_t$  should be set as

$$\eta_t = \frac{\text{diam}(\mathcal{Y})}{\|\nabla q(\mathbf{x}(t), \mathbf{y}(t))\| \sqrt{T}}. \quad (53)$$

#### D. Extending to Multiple Job Arrivals

OGASCHED can be applied to the scenarios where multiple jobs are yielded from each port in each time slot. In this case, the job arrival status  $\mathbf{x}(t)$  is re-formulated as  $\mathbf{x}(t) = [x_l(t)]_{l \in \mathcal{L}} \in \mathbb{N}^{|\mathcal{L}|}$ , where  $x_l(t)$  indicates the number of jobs arrive at port  $l$  at time  $t$ . Further, the scheduling decisions at time  $t$  is re-formulated as

$$\mathbf{y}(t) = [y_{l,r}^{j,k}]_{l,j,r,k} \in \mathbb{R}^{\sum_{l \in \mathcal{L}} J_l \times |\mathcal{R}_l| \times K},$$

where  $J_l$  is the maximum number of the type- $l$  jobs arrive during each time slot, i.e.,  $J_l = \max_{t \in \mathcal{T}} x_l(t)$ . Correspondingly, the port- $l$  reward is re-formulated as

$$q_l(\mathbf{x}(t), \mathbf{y}(t)) = \sum_{j=1}^{J_l} \mathbb{1}\{j \leq x_l(t)\} \left[ \sum_{k \in \mathcal{K}} f_k \left( \sum_{r \in \mathcal{R}_l} y_{l,r}^{j,k}(t) \right) - \max_{k \in \mathcal{K}} \left\{ \beta_k \sum_{r \in \mathcal{R}_l} y_{l,r}^{j,k}(t) \right\} \right],$$

where  $\mathbb{1}\{p\}$  is the indicator function:  $\mathbb{1}\{p\}$  is 1 if the predicate  $p$  is true, otherwise 0. The new formulated problem can be solved by native OGASCHED after transformations.

#### E. Extending to Gang Scheduling

OGASCHED can be extended to the Gang Scheduling scenarios, where the scheduling decisions for the task instances of a job follows the ALL-OR-NOTHING property. In other words, only when *all* tasks<sup>1</sup> of a job are successfully scheduled, the job could be launched.

In the following, we show briefly how Gang Scheduling can be modeled. To start with, for each job type  $l \in \mathcal{L}$ , we denote the corresponding set of task components by  $\mathcal{Q}_l$  and indexed by  $q$ . Correspondingly, the job requests  $\mathbf{a}_l$  is redefined as

$$\mathbf{a}_l = [a_l^{q,k}]_{l,q,k} \in \mathbb{R}_{\geq 0}^{\sum_{l \in \mathcal{L}} |\mathcal{Q}_l| \times K}.$$

Similarly, we redefine the scheduling decisions at time  $t$  as

$$\mathbf{y}(t) = [y_{l,r}^{q,k}]_{l,q,r,k} \in \mathbb{R}_{\geq 0}^{\sum_{l \in \mathcal{L}} |\mathcal{Q}_l| \times |\mathcal{R}_l| \times K}.$$

As a result, the solution space  $\mathcal{Y}$  turns to

$$\mathcal{Y} = \left\{ y_{l,r}^{q,k} \mid \sum_{q \in \mathcal{Q}_l} \mathbb{1} \left\{ \sum_{r \in \mathcal{R}_l} \sum_{k \in \mathcal{K}} y_{l,r}^{q,k} > 0 \right\} \geq m_l(t), \forall l, \right. \\ \left. 0 \leq y_{l,r}^{q,k}(t) \leq a_l^{q,k}, \forall l, r, q, k, t, \right.$$

<sup>1</sup>In practice, not all tasks of a job need to be scheduled. In Kubernetes, the job submitter can specify the minimum number of tasks that must be scheduled successfully. In the following, we use  $m_l(t)$  to represent the minimum number of tasks that should be scheduled at time  $t$  of the type- $l$  job.

$$\left. \sum_{l \in \mathcal{L}_r} \sum_{q \in \mathcal{Q}_l} y_{l,r}^{q,k}(t) \leq c_r^k, \forall r, k, t \right\},$$

where in the first inequality,  $m_l(t)$  is the minimum number of task components that should be scheduled at time  $t$  of type- $l$  job. The port- $l$  reward at time  $t$  is re-formulated as

$$q_l(\mathbf{x}(t), \mathbf{y}(t)) = x_l(t) \left[ \sum_{k \in \mathcal{K}} f_k \left( \sum_{q \in \mathcal{Q}_l} \sum_{r \in \mathcal{R}_l} y_{l,r}^{q,k}(t) \right) - \max_{k \in \mathcal{K}} \left\{ \beta_k \sum_{q \in \mathcal{Q}_l} \sum_{r \in \mathcal{R}_l} y_{l,r}^{q,k}(t) \right\} \right].$$

The new formulated problem is more difficult because  $\mathcal{Y}$  is no longer a convex set and  $q_l(\mathbf{x}(t), \mathbf{y}(t))$  is not differentiable everywhere. Nevertheless, we can still develop a similar online scheduling algorithm with the subgradient ascent and mirror ascent related techniques which retains a sublinear regret. The design detail is omitted due to space limits.

## IV. EXPERIMENTAL RESULTS

In this section, we conduct extensive experiments to validate the performance of OGASCHED. Based on the Alibaba cluster trace datasets [33], we first examine the theoretically guaranteed superiority of OGASCHED against several baselines on the cumulative and average rewards. Then, we analyze the generality and robustness of it under different cluster settings. At last, we validate the efficacy of OGASCHED in large-scale scenarios. The trace-driven simulation is conducted on a server with 48 Intel Xeon Silver 4214 CPUs, 256 GB memory, and 2 Tesla P100 GPUs.

*Traces:* We hybrid the traces from cluster-trace-v2018 and cluster-trace-gpu-v2020 of the Alibaba Cluster Trace Program. Specifically, we leverage the specifications of the machines, the arrival patterns, and the resource requirements of different kinds of jobs to generate our simulation environment.

*Baselines:* The following widely used baselines are implemented to make comparisons with OGASCHED.

- DRF [20]. It is adopted by YARN [34] and Mesos [35]. In our scenario, DRF allocates resources to ports that yield jobs in the ascending order of their dominant resource shares. The dominant share  $s_l$  of port  $l$  is calculated as  $s_l = \max_{k \in \mathcal{K}} \{a_l^k / \sum_{r \in \mathcal{R}_l} c_r^k\}$ .
- FAIRNESS. We implement FAIRNESS in this way: at each time  $t$ , we allocate the type- $k$  resource of each node  $r$  to each port  $l$  that yield a job according to the job's share  $a_l^k / \sum_{l \in \mathcal{L}_r} a_l^k$ .
- BINPACKING. It is optional in Kubernetes with the name of MOSTALLOCATED strategy and supported in Volcano as a configurable plugin [36]. Specifically, it scores the computing instances based on the utilization of resources, favoring the ones with higher allocation.
- SPREADING. It is similar to BINPACKING in procedures but with an opposite favor. The nodes with lower utilizations of resources have higher scores.



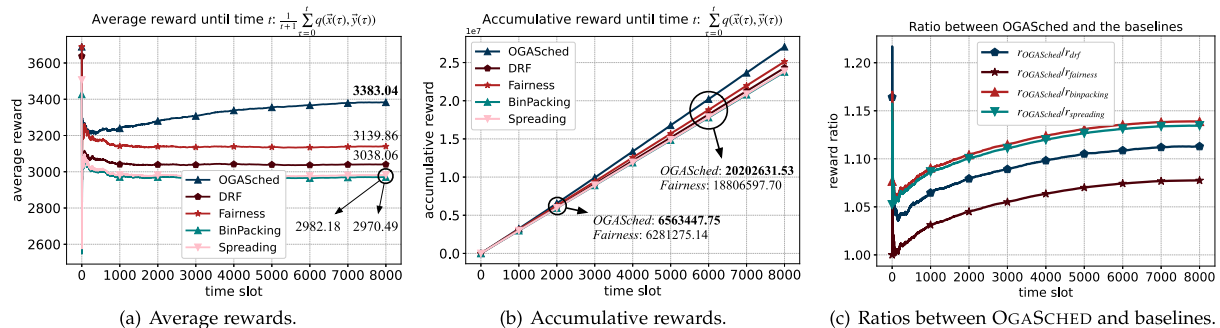


Fig. 2. Performance verification of OGASCHED. It takes one hour for OGASCHED to finish when  $T = 8000$ ,  $\beta \in [0.4, 0.6]$ , and contention level is 11.

TABLE II  
DEFAULT PARAMETER SETTINGS

PARAMETER	VALUE	PARAMETER	VALUE
job types num. $ \mathcal{L} $	10	node num. $ \mathcal{R} $	128
device type num. $K$	6	time slot num. $T$	2000
range of $\alpha$	[1.0, 1.5]	range of $\beta$	[0.3, 0.5]
initial learning rate $\eta_0$	25	decay $\lambda$	0.9999
job arrival prob. $\rho$	0.7	contention level	10

*Default Settings:* In default settings, our simulation environment has 128 computing instances, each equipped with 6 types of resources (CPUs, MEM, GPUs, NPUs, TPUs, and FPGAs), and 10 job types of different resource requirements. Large-scale validations will be demonstrated in Section IV-C. The computing instances and jobs are carefully selected from the trace to reflect heterogeneity. We support 4 types of utilities

$$f_r^k(y) = \begin{cases} \alpha y & \text{linear} \\ \alpha \ln(y+1) & \text{log} \\ \alpha^{-1} - (y+\alpha)^{-1} & \text{reciprocal} \\ \alpha \sqrt{y+1} - \alpha & \text{poly,} \end{cases} \quad (54)$$

The default settings of main parameters are listed in Table II. In this table, the initial learning rate and the decay are used to tune the learning rate at each time  $t$  around the value (53). Job arrival probability  $\rho$  is adopted to adjust the job arrival status with Bernoulli Distributions. This parameter is applied based on the actual arrival patterns from the trace to increase stochasticity. The contention level, located at the last cell of this table, is designed to tune the level of resource contention. The larger this value, the more fierce the contention. It is a multiplier to the resource requirements of jobs. The effect of it will be analyzed in detail in Section IV-B.

Note that in Section IV-A, the time slot length  $T$  is set as 8000. For the left experiments, the time slot length is 2,000, unless otherwise stated.

### A. Performance Verification

In this section, we compare the performance of OGASCHED with the baselines in terms of the achieved cumulative and average rewards.

In Fig. 2(a), the  $y$ -axis is the average reward until time  $t$ , i.e.,  $\frac{1}{t} \sum_{\tau=1}^t q(x(\tau), y(\tau))$ . Compared with the baselines, OGASCHED has a clear advantage on the performance (with the increases of 11.33%, 7.75%, 13.89%, and 13.44% compared with DRF, FAIRNESS, BINPACKING, and SPREADING, respectively). Besides, it shows that the performance of OGASCHED tends to increase as the length of the time horizon increases. The curve of OGASCHED starts steep and later flattens. The reason is that, as a learning-powered algorithm, OGASCHED learns the underlying distribution of job arrival patterns and it can make better decisions by adjusting the step directions. It is interesting to find that the rewards oscillate at the beginning time slots. One of the leading factors is that OGASCHED is not boosted with a well-designed initial solution. In our experiments, a 8000-time slot training only takes one hour. Thus, not surprisingly, the rewards achieved in the beginning can be easily surpassed when the time slot is sufficiently large.

It is not a surprise that FAIRNESS achieves the best among the baselines. FAIRNESS adopts a proportional allocation strategy and allocates resources to each non-empty port *without bias*, which increases the computation gains adequately. When the contention is not fierce while the communication overhead is low, the advantages of FAIRNESS will be more steady. By contrast, the advantages of BINPACKING and SPREADING are respectively high resource utilization and job isolation, which do not contribute to the reward directly.

Fig. 2(b) shows that the cumulative rewards achieved by all the five algorithms. In the beginning, FAIRNESS and DRF have the slight edge, benefiting by the propotional allocation idea. Nevertheless, as the time slot increases, OGASCHED is able to surpass them without difficulty. Fig. 2(c) demonstrates the ratio on the achieved average rewards between OGASCHED and the baselines. Similarly, the ratios oscillate at the beginning. After that, they increase steeply and later flattens.

The hyper-parameters of OGASCHED, especially the initial learning rate  $\eta_0$  and the decay, have a remarkable impact on its performance. From Fig. 4 we can find that, a wrong setting of these hyper-parameters could lead to a poor performance, even the decrease of the cumulative reward (which means, the reward is negative in many time slots). At the last of Section III-C, we claim that, to achieve an cumulative reward with a lower bound guarantee, at each time  $t$ , the learning rate should be set around (53). Note that in (53), the learning rate is encouraged to be

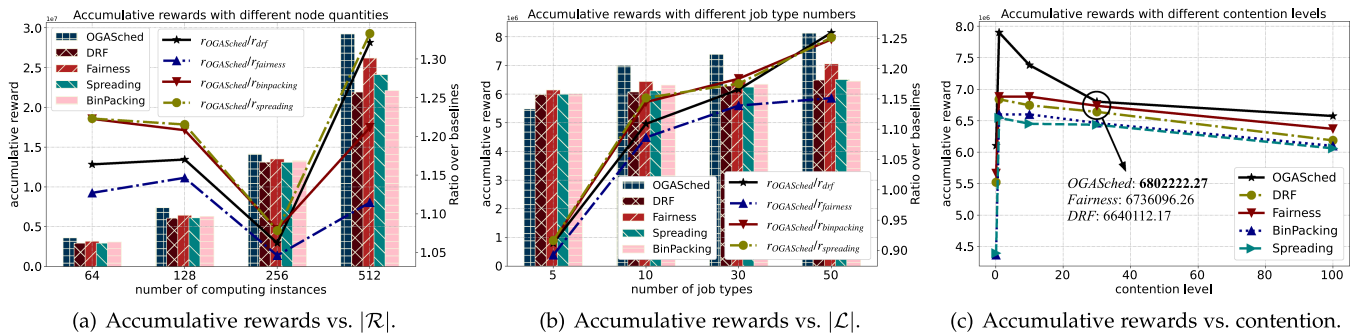


Fig. 3. Scalability verification of OGASCHED under different scales of the bipartite graph and the contention levels.

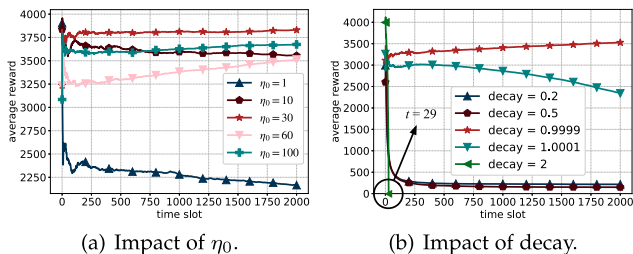


Fig. 4. The performance of OGASCHED with different hyper-parameters.

larger and larger as time moves, which is counterintuitive and it goes against the convergence to a local optimum. The curves in Fig. 4(b) also verify that, setting decay as 0.9999 is better than 1.0001. The best decay in practice does not follow the guidance of theory because the regret analysis only gives the *worst* case guarantee on the cumulative rewards. In our experiments, the best range for decay is [0.995, 0.9999].

### B. Scalability, Generality and Robustness Evaluations

In this section, we evaluate the performance of OGASCHED under different scales of scenario settings. Fig. 3(a) and (b) demonstrate the impact of the scale of the bipartite graph  $\mathcal{G}$ . In these two figures, the left  $y$ -axis is the cumulative reward while the right  $y$ -axis is the ratio  $r_a/r_b$ , where  $r_a$  is the cumulative reward achieved by OGASCHED, and  $r_b$  is the baselines'. First, we observe that, whatever the number of the computing instances is, OGASCHED takes the leading position. Besides, as  $|\mathcal{R}|$  increases, all the algorithms obtain a larger cumulative reward. The result is evident because a large cluster can provide sufficient resources, which leads to jobs being fully served. It is also worth noting that, when  $|\mathcal{R}|$  increases, the superiority of OGASCHED over the baselines first increases then decreases. It demonstrates that the resource contention is fierce when  $|\mathcal{R}| \in [128, 256]$ . In this case, it is necessary for OGASCHED to be trained with a larger time slot. Fig. 3(b) shows that the number of job types, i.e.,  $|\mathcal{L}|$ , has a weaker impact than  $|\mathcal{R}|$  to the performance of OGASCHED. The phenomenon verifies the conclusion we have concluded, i.e., the regret grows linearly with  $|\mathcal{R}|$ , but it is sublinear with  $|\mathcal{L}|$ .

Fig. 3(c) shows the impact of contention level. This parameter works as a multiplier to the resource requirements of jobs. We can observe that, when moving contention level from 0.1 to

1, all the achieved cumulative rewards increase. This is obvious because a larger resource requirement leads to a larger computation gain on the premise of low contention. However, increasing the multiplier further leads to the downgrade of performances and the reduction of the superiority of OGASCHED. Even so, OGASCHED always performs the best. Fig. 6 shows the average computation gain and communication overhead penalty of each time slot under different contention levels. We can find that the penalty increases with the contention level slowly.

Fig. 7 demonstrates the cumulative rewards with different utilities. Because of the diminishing marginal effect, the rewards with *poly*, *log*, and *reciprocal* utilities are significantly less than the rewards with *linear* utilities. Nevertheless, the diminishing marginal effect does not change the superiority of OGASCHED against the baselines. Even in the *all reciprocal* utility settings, for FAIRNESS, OGASCHED has its advantages.

In addition to the above evaluations, we also test the generality and robustness of OGASCHED under different settings of the following parameters: the time horizon length  $T$ , the job arrival probability  $\rho$ , and the dense of the bipartite graph. The graph dense is calculated as  $\sum_{r \in \mathcal{R}} |\mathcal{L}_r|/|\mathcal{R}|$ . The results are shown in Table III. The two largest values in each column of the table are made bold. Besides, for each parameter and each algorithm, the setting which leads to the largest reward is marked with a light-grey background. We summarize the key findings as follows.

- First, whatever the parameter settings, OGASCHED always performs the best, and its performance has a positive correlation with the time horizon length  $T$ . As we have analyzed, a large time horizon provides more chances for OGASCHED to learn the underlying distributions, thereby increasing the reward in the gradient ascent directions.
- Increasing the job arrival probability can lead to a high resource utilization, thereby increasing the rewards. However, a large job arrival probability also brings in a fierce resource contention. A direct consequence of it is that, for OGASCHED, many elements in the vector  $\mathbf{y}(t)$  fall into the interior of  $\mathcal{Y}$ , rather than the boundaries, thereby leading to a reward reduction. The phenomenon can be observed when moving  $\rho$  from 0.7 to 0.9.
- Graph dense has a similar effect on the reward to the job arrival probability. Nevertheless, the reasons behind are distinct. A larger graph dense increases the opportunities for a job to be served with a large possible parallelism,

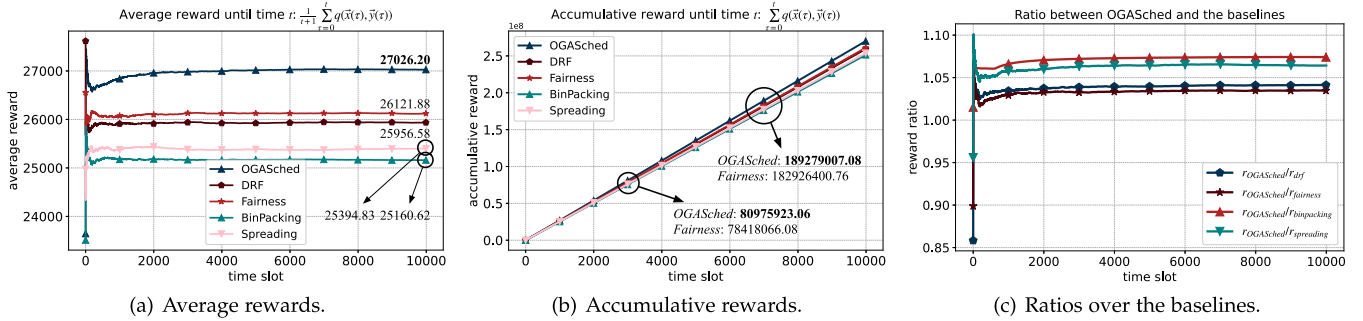


Fig. 5. Large-scale validations. It takes 15 hours for OGASCHED to complete when  $T = 10000$ ,  $\beta \in [0.01, 0.015]$ , and contention level is 5.

TABLE III  
GENERALITY AND ROBUSTNESS VALIDATION UNDER DIFFERENT SCENARIO SETTINGS

Avg. Reward	Time Horizon Length $T$				Job Arrival Probability $\rho$				Graph Dense		
	1000	2000	5000	10000	0.3	0.5	0.7	0.9	$\approx 2$	$\approx 2.5$	$\approx 3$
OGASCHED	2578.53	2886.33	2911.37	3104.98	1904.87	2154.18	3117.29	2938.22	2816.18	2904.51	3127.47
DRF	2422.47	2493.02	2449.23	2497.85	1364.53	2086.59	2503.01	2755.41	2417.08	2786.94	2795.42
FAIRNESS	2532.24	2582.80	2552.41	2436.22	1295.53	1997.19	2628.02	2873.84	2501.54	2857.60	2918.98
BINPACKING	2386.01	2449.15	2444.32	2365.13	1246.39	1897.79	2518.98	2740.19	2374.31	2757.71	2829.19
SPREADING	2382.01	2466.71	2436.60	2362.88	1250.67	1888.06	2519.37	2737.93	2382.87	2766.07	2836.37

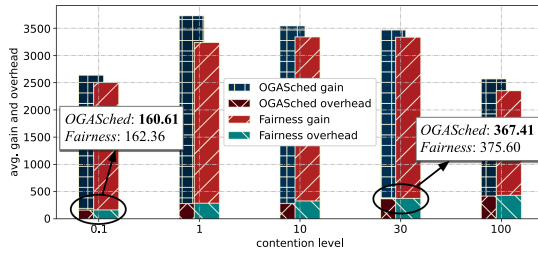


Fig. 6. Average computation gain and communication overhead of each time slot under different contention levels.

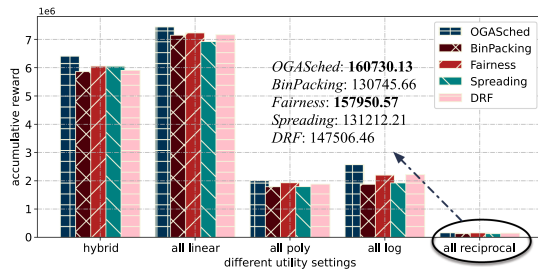


Fig. 7. Accumulative rewards with different utilities.

thereby increasing the computation gain. By contrast, the communication overhead has a slow rate of growth.

### C. Large-Scale Validations

To test the efficacy of OGASCHED in large-scale scenarios, we conduct the following experiments. In these experiments, the number of the job types is set as 100 while the quantity of the computing instances is 1,024 in default. The results in Fig. 5

show that the superiority of OGASCHED is preserved even in large-scale scenarios.

## V. RELATED WORKS

The design of online job scheduling algorithms that yield a nice theoretical bound is always the focus of attention from the research community. Existing online job scheduling algorithms can be organized into two categories.

In the first category, the online algorithms are elaborately designed for specific job types, such as DNN model training [3], [6], [6], [13], [18], [37], [38], big-data query & analytics [12], [39], multi-stage workflows [17], [40], [41], [42], etc. A typical work on DNN model training is [18], where the authors fully take the layered structure of DNNs into consideration and develop an efficient resource scheduling algorithm based on the sum-of-ratios multi-type-knapsack decomposition method. The authors further prove that the proposed algorithm has a SOTA approximation ratio within a polynomial running time. [13] is another work that fully explores the Bulk Synchronous Parallel (BSP) property of the DNN training jobs. The authors develop an algorithm which is  $\mathcal{O}(\ln |\mathcal{M}|)$ -approximate with high probability, where  $\mathcal{M}$  is the set of resources. These works are designed for specific job types, and they do not provide a general analysis of the gain-overhead tradeoff for multi-server jobs. This paper intends to fill the gap.

In the second category, the types of job are not specified, while the theoretical superiority is highlighted. The algorithms are designed with different theoretical basis, including online approximate algorithms [15], [43], [44], Online Convex Optimization (OCO) techniques [14], game-theoretical approaches [16], online learning and DRL-based algorithms [45], [46], etc. In these works, the performance of the proposed algorithms is



usually analyzed with approximate ratio, competitive ratio, Price of Anarchy (PoA), and regret. A typical recent work is [14]. The authors develop an algorithm whose dynamic regret is upper bounded by  $\mathcal{O}(\text{OPT}^{1-\beta})$ , where  $\beta \in [0, 1)$ . None of the existing works analyze the gain-overhead tradeoff and provide a regret of  $\mathcal{O}(\sqrt{|\mathcal{L}|T})$  as this paper demonstrates.

## VI. CONCLUSION

In this article, we study the online scheduling of multi-server jobs in terms of the gain-overhead tradeoff. The problem is formulated as an cumulative reward maximization program. The reward of scheduling a job is designed as the difference between the computation gain and the penalty on the dominant communication overhead. We propose an algorithm, i.e., OGASCHED, to learn the best possible scheduling decision in the ascending direction of the reward gradients. OGASCHED is the first algorithm that has a sublinear regret w.r.t. the number of job types and time slot length, which is a SOTA result for concave rewards. OGASCHED is well designed to be parallelized, which makes large-scale applications possible. The superiority of OGASCHED is also validated with extensive trace-driven simulations. Future extensions may include, i.e., more elaborate modeling and analysis of the intra-node and inter-node communication overheads.

## REFERENCES

- [1] M. Wu et al., "GraM: Scaling graph computation to the trillions," in *Proc. 6th ACM Symp. Cloud Comput.*, New York, NY, USA, 2015, pp. 408–421. [Online]. Available: <https://doi.org/10.1145/2806777.2806849>
- [2] L. Wang, W. Wang, and B. Li, "CMFL: Mitigating communication overhead for federated learning," in *Proc. IEEE 39th Int. Conf. Distrib. Comput. Syst.*, 2019, pp. 954–964.
- [3] Y. Peng, Y. Bao, Y. Chen, C. Wu, C. Meng, and W. Lin, "DL2: A deep learning-driven scheduler for deep learning clusters," *IEEE Trans. Parallel Distrib. Syst.*, vol. 32, no. 08, pp. 1947–1960, Aug. 2021.
- [4] W. Wang, Q. Xie, and M. Harchol-Balter, "Zero queueing for multi-server jobs," in *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 5, no. 1, Feb. 2021, Art. no. 7.
- [5] H. Zhao, S. Deng, F. Chen, J. Yin, S. Dustdar, and A. Y. Zomaya, "Learning to schedule multi-server jobs with fluctuated processing speeds," *IEEE Trans. Parallel Distrib. Syst.*, vol. 34, no. 1, pp. 234–245, Jan. 2023.
- [6] Y. Bao, Y. Peng, C. Wu, and Z. Li, "Online job scheduling in distributed machine learning clusters," in *Proc. IEEE Conf. Comput. Commun.*, 2018, pp. 495–503.
- [7] C. Carrión, "Kubernetes scheduling: Taxonomy, ongoing issues and challenges," *ACM Comput. Surv.*, vol. 55, 2022, Art. no. 138.
- [8] Y. Jiang, Y. Zhu, C. Lan, B. Yi, Y. Cui, and C. Guo, "A unified architecture for accelerating distributed DNN training in heterogeneous GPU/CPU clusters," in *Proc. 14th USENIX Symp. Operating Syst. Des. Implementation*, USENIX Association, 2020, pp. 463–479. [Online]. Available: <https://www.usenix.org/conference/osdi20/presentation/jiang>
- [9] L. Zheng et al., "Alpa: Automating inter- and intra-operator parallelism for distributed deep learning," in *Proc. 16th USENIX Symp. Operating Syst. Des. Implementation*, Carlsbad, CA: USENIX Association, 2022, pp. 559–578. [Online]. Available: <https://www.usenix.org/conference/osdi22/presentation/zheng-lianmin>
- [10] X. Jia et al., "Whale: Efficient giant model training over heterogeneous GPUs," in *Proc. USENIX Annu. Tech. Conf.*, Carlsbad, CA: USENIX Association, 2022, pp. 673–688. [Online]. Available: <https://www.usenix.org/conference/atc22/presentation/jia-xianyan>
- [11] V. Prabhakaran, M. Wu, X. Weng, F. McSherry, L. Zhou, and M. Haradasan, "Managing large graphs on multi-cores with graph awareness," in *Proc. USENIX Annu. Tech. Conf.*, Boston, MA: USENIX Association, 2012, pp. 41–52. [Online]. Available: <https://www.usenix.org/conference/atc12/technical-sessions/presentation/prabhakaran>
- [12] H. Mao, M. Schwarzkopf, S. B. Venkatakrisnan, Z. Meng, and M. Alizadeh, "Learning scheduling algorithms for data processing clusters," in *Proc. ACM Special Int. Group Data Commun.*, 2019, pp. 270–288.
- [13] Z. Han et al., "Scheduling placement-sensitive BSP jobs with inaccurate execution time estimation," in *Proc. IEEE Conf. Comput. Commun.*, 2020, pp. 1053–1062.
- [14] Y. Liu, H. Xu, and W. C. Lau, "Online job scheduling with resource packing on a cluster of heterogeneous servers," in *Proc. IEEE Conf. Comput. Commun.*, 2019, pp. 1441–1449.
- [15] K. Psychas and J. Ghaderi, "Scheduling jobs with random resource requirements in computing clusters," in *Proc. IEEE Conf. Comput. Commun.*, 2019, pp. 2269–2277.
- [16] R. Burra, C. Singh, and J. Kuri, "Service scheduling for bernoulli requests and quadratic cost," in *Proc. IEEE Conf. Comput. Commun.*, 2019, pp. 2584–2592.
- [17] B. Tian, C. Tian, H. Dai, and B. Wang, "Scheduling coflows of multi-stage jobs to minimize the total weighted job completion time," in *Proc. IEEE Conf. Comput. Commun.*, 2018, pp. 864–872.
- [18] M. Yu, C. Wu, B. Ji, and J. Liu, "A sum-of-ratios multi-dimensional-knapsack decomposition for DNN resource scheduling," in *Proc. IEEE Conf. Comput. Commun.*, 2021, pp. 1–10.
- [19] T. Anderson, *The Theory and Practice of Online Learning*. Athabasca, Canada: Athabasca University Press, 2008.
- [20] A. Ghodsi, M. Zaharia, B. Hindman, A. Konwinski, S. Shenker, and I. Stoica, "Dominant resource fairness: Fair allocation of multiple resource types," in *Proc. 8th USENIX Symp. Netw. Syst. Des. Implementation*, 2011, pp. 323–336.
- [21] F. Khorasani, H. Asghari Esfeden, A. Farmahini-Farahani, N. Jayasena, and V. Sarkar, "RegMutex: Inter-warp GPU register time-sharing," in *Proc. IEEE/ACM 45th Annu. Int. Symp. Comput. Architecture*, 2018, pp. 816–828.
- [22] T. N. Le, X. Sun, M. Chowdhury, and Z. Liu, "AlloX: Compute allocation in hybrid clusters," in *Proc. 15th Eur. Conf. Comput. Syst.*, 2020, Art. no. 31.
- [23] Q. Weng et al., "MLaaS in the wild: Workload analysis and scheduling in large-scale heterogeneous GPU clusters," in *Proc. 19th USENIX Symp. Netw. Syst. Des. Implementation*, Renton, WA, 2022, pp. 945–960.
- [24] J. F. Kurose and R. Simha, "A microeconomic approach to optimal resource allocation in distributed computer systems," *IEEE Trans. Comput.*, vol. 38, no. 5, pp. 705–717, May 1989.
- [25] G. Bao and P. Guo, "Federated learning in cloud-edge collaborative architecture: Key technologies, applications and challenges," *J. Cloud Comput.*, vol. 11, no. 1, 2022, Art. no. 94.
- [26] S. Shalev-Shwartz, "Online learning and online convex optimization," *Found. Trends Mach. Learn.*, vol. 4, no. 2, pp. 107–194, 2012.
- [27] R. Bhattacharjee, S. Banerjee, and A. Sinha, "Fundamental limits on the regret of online network-caching," in *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 4, no. 2, Jun. 2020, Art. no. 25.
- [28] Y. Ying and M. Pontil, "Online gradient descent learning algorithms," *Found. Comput. Math.*, vol. 8, no. 5, pp. 561–596, 2008.
- [29] G. S. Paschos, A. Destounis, L. Vigneri, and G. Iosifidis, "Learning to cache with no regrets," in *Proc. IEEE Conf. Comput. Commun.*, 2019, pp. 235–243.
- [30] G. S. Paschos, A. Destounis, and G. Iosifidis, "Online convex optimization for caching networks," *IEEE/ACM Trans. Netw.*, vol. 28, no. 2, pp. 625–638, Apr. 2020.
- [31] G. I. Ricardo, A. Tuhulokova, G. Neglia, and T. Spyropoulos, "Caching policies for delay minimization in small cell networks with coordinated multi-point joint transmissions," *IEEE/ACM Trans. Netw.*, vol. 29, no. 3, pp. 1105–1115, Jun. 2021.
- [32] F. Orabona, "A modern introduction to online learning," 2019, *arXiv:1912.13213*.
- [33] Alibaba cluster trace. 2022. [Online]. Available: <https://github.com/alibaba/clusterdata>
- [34] V. K. Vavilapalli et al., "Apache hadoop YARN: Yet another resource negotiator," in *Proc. 4th Annu. Symp. Cloud Comput.*, New York, NY, USA, 2013, Art. no. 5.
- [35] B. Hindman et al., "Mesos: A platform for fine-grained resource sharing in the data center," in *Proc. 8th USENIX Symp. Netw. Syst. Des. Implementation*, Boston, MA: USENIX Association, 2011, pp. 295–308.
- [36] volcano sh. "Volcano," 2022. [Online]. Available: <https://github.com/volcano-sh/volcano>
- [37] A. Qiao et al., "Pollux: Co-adaptive cluster scheduling for goodput-optimized deep learning," in *Proc. 15th USENIX Symp. Operating Syst. Des. Implementation*, 2021, pp. 1–18.
- [38] D. Narayanan, K. Santhanam, F. Kazhemiaka, A. Phanishayee, and M. Zaharia, "Heterogeneity-aware cluster scheduling policies for deep learning workloads," in *Proc. 14th USENIX Symp. Operating Syst. Des. Implementation*, 2020, pp. 481–498.
- [39] D. Cheng, X. Zhou, Y. Xu, L. Liu, and C. Jiang, "Deadline-aware MapReduce job scheduling with dynamic resource availability," *IEEE Trans. Parallel Distrib. Syst.*, vol. 31, no. 4, pp. 1101–1114, Apr. 2020.



- [40] Z. Hu, B. Li, C. Chen, and X. Ke, "FlowTime: Dynamic scheduling of deadline-aware workflows and ad-hoc jobs," in *Proc. IEEE 38th Int. Conf. Distrib. Comput. Syst.*, 2018, pp. 929–938.
- [41] D. Hu and B. Krishnamachari, "Throughput optimized scheduler for dispersed computing systems," in *Proc. IEEE 7th Int. Conf. Mobile Cloud Comput. Serv. Eng.*, 2019, pp. 76–84.
- [42] C.-S. Yang, R. Pedarsani, and A. S. Avestimehr, "Communication-aware scheduling of serial tasks for dispersed computing," *IEEE/ACM Trans. Netw.*, vol. 27, no. 4, pp. 1330–1343, Aug. 2019.
- [43] J. Meng, H. Tan, X.-Y. Li, Z. Han, and B. Li, "Online deadline-aware task dispatching and scheduling in edge computing," *IEEE Trans. Parallel Distrib. Syst.*, vol. 31, no. 6, pp. 1270–1286, Jun. 2020.
- [44] Z. Han et al., "OnDisc: Online latency-sensitive job dispatching and scheduling in heterogeneous edge-clouds," *IEEE/ACM Trans. Netw.*, vol. 27, no. 6, pp. 2472–2485, Dec. 2019.
- [45] S. Liang, Z. Yang, F. Jin, and Y. Chen, "Data centers job scheduling with deep reinforcement learning," in *Proc. Pacific-Asia Conf. Knowl. Discov. Data Mining*, Springer, 2020, pp. 906–917.
- [46] Y. Han, S. Shen, X. Wang, S. Wang, and V. C. Leung, "Tailored learning-based scheduling for kubernetes-oriented edge-cloud system," in *Proc. IEEE Conf. Comput. Commun.*, 2021, pp. 1–10.



**Hailiang Zhao** received the BS degree from the School of Computer Science and Technology, Wuhan University of Technology, Wuhan, China, in 2019. He is currently working toward the PhD degree with the College of Computer Science and Technology, Zhejiang University, Hangzhou, China. His research interests include cloud & edge computing, distributed systems and optimization algorithms. He has published several papers in flagship conferences and journals including IEEE ICWS 2019, *IEEE Transactions on Parallel and Distributed Systems*, *IEEE Transactions on Mobile Computing*, etc. He has been a recipient of the Best Student Paper Award of IEEE ICWS 2019. He is a reviewer for *IEEE Transactions on Services Computing* and *Internet of Things Journal*.



**Shuiguang Deng** (Senior Member, IEEE) received the BS and PhD degrees both in computer science from Zhejiang University, China, in 2002 and 2007, respectively. He is currently a full professor with the College of Computer Science and Technology, Zhejiang University, China. He previously worked with the Massachusetts Institute of Technology, in 2014 and Stanford University, in 2015 as a visiting scholar. His research interests include edge computing, service computing, cloud computing, and business process management. He serves for the journal *IEEE Trans. on Services Computing, Knowledge and Information Systems, Computing*, and *IET Cyber-Physical Systems: Theory & Applications* as an associate editor. Up to now, he has published more than 100 papers in journals and refereed conferences. In 2018, he was granted the Rising Star Award by IEEE TCSVC. He is a fellow of IET.



**Zhengzhe Xiang** received the BS and PhD degrees of computer science and technology from Zhejiang University, Hangzhou, China. He was previously a visiting student worked with the Karlstad University, Sweden, in 2018. He is currently a lecturer with Zhejiang University City College, Hangzhou, China. His research interests lie in the fields of service computing, cloud computing, and edge computing.



**Xueqiang Yan** is currently a technology expert with the Wireless Technology Lab, Huawei Technologies. He was a member of technical staff with Bell Labs from 2000 to 2004. From 2004 to 2016, he was the director of the Strategy Department, Alcatel-Lucent Shanghai Bell. His current research interests include wireless networking, the Internet of Things, edge AI, future mobile network architecture, network convergence, and evolution.



**Jianwei Yin** received the PhD degree in computer science from Zhejiang University (ZJU), in 2001. He was a visiting scholar with the Georgia Institute of Technology. He is currently a full professor with the College of Computer Science, ZJU. Up to now, he has published more than 100 papers in top international journals and conferences. His current research interests include service computing and business process management. He is an associate editor of the *IEEE Transactions on Services Computing*.



**Schahram Dustdar** (Fellow, IEEE) is a full professor of computer science (informatics) with a focus on Internet Technologies heading the Distributed Systems Group with the TU Wien. He is founding co-editor-in-chief of *ACM Transactions on Internet of Things (ACM TIoT)* as well as editor-in-chief of *Computing (Springer)*. He is an associate editor of *IEEE Transactions on Services Computing*, *IEEE Transactions on Cloud Computing*, *ACM Computing Surveys*, *ACM Transactions on the Web*, and *ACM Transactions on Internet Technology*, as well as on the editorial board of *IEEE Internet Computing* and *IEEE Computer*. He is recipient of multiple awards: TCI Distinguished Service Award (2021), IEEE TCSVC Outstanding Leadership Award (2018), IEEE TCSC Award for Excellence in Scalable Computing (2019), ACM distinguished scientist (2009), ACM distinguished speaker (2021), IBM Faculty Award (2012). He is an elected member of the Academia Europaea: The Academy of Europe, where he is chairman of the Informatics Section, an Asia-Pacific Artificial Intelligence Association (AAIA) president (2021) and fellow (2021). He is an EAI fellow (2021) and an I2CICC fellow (2021). He is a member of the 2022 IEEE Computer Society fellow Evaluating Committee (2022).



**Albert Y. Zomaya** (Fellow, IEEE) is the Peter Nicol Russell chair professor of computer science and director of the Centre for Distributed and High-Performance Computing with the University of Sydney. To date, he has published more than 600 scientific papers and articles and is (co-)author/editor of more than 30 books. A sought-after speaker, he has delivered more than 250 keynote addresses, invited seminars, and media briefings. His research interests span several areas in parallel and distributed computing and complex systems. He is currently the editor in chief of the *ACM Computing Surveys* and processed in the past as editor in chief of the *IEEE Transactions on Computers* (2010–2014) and the *IEEE Transactions on Sustainable Computing* (2016–2020). He is a decorated scholar with numerous accolades including Fellowship of the IEEE, the American Association for the Advancement of Science, and the Institution of Engineering and Technology (U.K.). Also, he is an elected fellow of the Royal Society of New South Wales and an elected foreign member of Academia Europaea. He is the recipient of the 1997 Edgeworth David Medal from the Royal Society of New South Wales for outstanding contributions to Australian Science, the IEEE Technical Committee on Parallel Processing Outstanding Service Award (2011), IEEE Technical Committee on Scalable Computing Medal for Excellence in Scalable Computing (2011), IEEE Computer Society Technical Achievement Award (2014), ACM MSWIM Reginald A. Fessenden Award (2017), the New South Wales Premier's Prize of Excellence in Engineering and Information and Communications Technology (2019), and the Research Innovation Award, IEEE Technical Committee on Cloud Computing (2021).