

# Optimal Online Balanced Graph Partitioning

**Abstract**—Distributed applications generate a significant amount of network traffic. By collocating frequently communicating nodes (e.g., virtual machines) on the same clusters (e.g., server or rack), we can reduce the network load and improve application performance. However, the communication pattern of different applications is often unknown a priori and may change over time, hence it needs to be learned in an online manner. This paper revisits the online balanced partitioning problem that asks for an algorithm that strikes an optimal tradeoff between the benefits of collocation (i.e., lower network load) and its costs (i.e., migrations). Our first contribution is a significantly improved deterministic lower bound of  $\Omega(k \cdot \ell)$  on the competitive ratio, where  $\ell$  is the number of clusters and  $k$  is the cluster size, even for a scenario in which the communication pattern is static and can be perfectly partitioned; we also provide an asymptotically tight upper bound of  $O(k \cdot \ell)$  for this scenario. For  $k = 3$ , we contribute an asymptotically tight upper bound of  $\Theta(\ell)$  for the general model in which the communication pattern can change arbitrarily over time. We improve the result for  $k = 2$  by providing a strictly 6-competitive upper bound for the general model.

## I. INTRODUCTION

The popularity of data-centric, distributed applications has led to an explosive growth of network traffic, especially in data centers [1, 2]. The performance of these distributed applications often critically depends on the underlying network [3], and efficient operation of these networks is important. At the same time, distributed systems are often highly virtualized today, and provide interesting new opportunities for resource optimization. In particular, it has become possible to operate data centers in a more demand-aware manner: by dynamically migrating nodes (e.g., virtual machines) which communicate frequently topologically closer to each other, network traffic can be reduced significantly. However, migrations entail overhead and should be used moderately.

This paper studies the algorithmic problem underlying such demand-aware optimizations, aiming to strike a balance between the benefits of migrations (e.g., reduced network load) and their costs. In particular, we are interested in an online variant of the problem: since communication patterns can change over time, an online algorithm needs to react dynamically to new traffic patterns, and migrate nodes accordingly. Ideally, this algorithm should perform close to an optimal offline algorithm, without requiring any information about future traffic demands.

This problem is known as the online balanced graph repartitioning problem and was introduced by Avin et al. [4, 5] at DISC 2016. A special variant of the general problem has later been studied by Henzinger et al. [6] at SIGMETRICS 2019. We refer to the latter as the learning model.

## A. Model

We study two models in this paper: the *general partitioning* model, and its subproblem, the *learning* model. In both models, we assume that communication patterns are not known to our algorithms at the beginning. We measure the quality of presented algorithmic solutions by competitive analysis [7], which is well-suited for problems that are online by their nature. In the competitive analysis, the goal is to optimize the *competitive ratio* of a given online algorithm: the ratio of its cost to the cost of an optimal offline algorithm that knows the entire input sequence in advance.

**General partitioning model.** In the *online balanced graph partitioning* problem, we are given a set  $V$  of  $n$  nodes (e.g., virtual machines or processes), initially arbitrarily partitioned into  $\ell$  clusters (e.g., servers or entire racks), each of size  $k$ . The nodes interact using a sequence of pairwise communication requests  $\sigma = (u_1, v_1), (u_2, v_2), (u_3, v_3), \dots$ , where a pair  $(u_t, v_t)$  indicates that nodes  $u_t$  and  $v_t$  exchange a certain amount of data. Nodes in  $C \subset V$  are *collocated* if they reside in the same cluster.

An algorithm serves a communication request between two nodes either *locally* at cost 0 if they are collocated, or *remotely* at cost 1 if they are located in different clusters. We refer to these two types of requests as *internal* and *external* requests, respectively. We may refer to external requests as *inter-cluster* requests or edges interchangeably. Before serving a request, an online algorithm may perform a *repartition*, i.e., it may move (“migrate”) some nodes into clusters different from their current clusters, while respecting the capacity of every cluster. Afterward, the algorithm serves the request. The cost of migrating a node from one cluster to another is  $\alpha \in \mathbb{Z}^+$ . For any algorithm ALG, its cost, denoted by  $\text{ALG}(\sigma)$ , is the total cost of communications and the cost of migrations performed by ALG while serving the sequence  $\sigma$ .

**Learning model.** We further study a *learning* variant of online balanced graph partitioning, where the communication pattern is *static*: whether a pair of nodes ever communicate or not, is determined a priori and is unknown to algorithms, and such pairs communicate forever. As in Henzinger et al. [6], we assume that the communication graph admits a *perfect partition*, i.e., a partition in which no inter-cluster request ever occurs. Any algorithm must eventually collocate pairs of communicating nodes, as otherwise it cannot be competitive. The objective is to *learn* the communication graph while serving all requests, and without performing too many node migrations. In the learning model, for simplicity we assume that the migration cost is  $\alpha = 1$  (our bounds hold for any  $\alpha > 1$  as well).

## B. Related Work

The works closest to ours are by Avin et al. at DISC '16 (on the general partitioning model) [4, 5], by Henzinger et al. (on the learning model) [6] at SIGMETRICS '19 and SODA '20 [8]. However, the focus of these papers is primarily on models with resource augmentation: the online algorithm can use slightly larger clusters than the offline algorithm. Avin et al. actually showed that their lower bound  $\Omega(k)$  holds even in a scenario with significant resource augmentation, and they provided an algorithm with the competitive ratio  $O(k \log k)$  using the  $(2 + \epsilon)$ -augmented cluster capacity. Their ratio is independent of  $\ell$ , which is impossible without significant resource augmentation.

In contrast, we study the non-augmented setting, where the nodes need to be perfectly balanced among the clusters. This assumption is not only more realistic but also significantly more challenging, as it is related to hard problems such as integer partitioning [9]. In terms of results without augmentation, so far, it is only known that there exists an  $O(k^2 \cdot \ell^2)$ -competitive algorithm [4] (with a variant running in polynomial time [10]); the best known lower bound is significantly lower, namely  $\Omega(k)$ . For  $k = 2$ , Avin et al. [4] presented a 7-competitive algorithm with a substantial ( $\Omega(\ell^2)$ ) additive constant.

The problem has also been studied in a weaker model where the adversary can only sample requests from a fixed distribution [11].

The *static* offline version of the partitioning problem is known as the  $\ell$ -balanced graph partitioning problem, where the entire communication graph is known in advance, and the task is to partition  $n$  nodes into  $\ell$  clusters of capacity  $n/\ell$  each, minimizing the number of inter-cluster edges. The problem is NP-complete, and cannot even be approximated within any finite factor unless  $P=NP$  [12]. The static variant where  $\ell = 2$  corresponds to the minimum bisection problem, which is already NP-hard [13], and the currently best approximation ratio is  $O(\log n)$  [14, 15, 16, 17, 18, 19].

Our problem is further related to some classic online problems. In particular, it is related to online paging [20, 21, 22, 23], sometimes also referred to as online caching, where requests for data items (nodes) arrive over time and need to be served from a cache of finite capacity, and where the number of cache misses must be minimized. Classic problem variants usually boil down to finding a smart eviction strategy, such as Least Recently Used (LRU) [20]. In our setting, requests can be served remotely (i.e., without fetching the corresponding nodes to a single physical machine). In this light, our model is more reminiscent of caching models *with bypassing* [24, 25, 26]. A major difference between these problems is that in the caching problems, each request involves a single element of the universe, while in our model *both* end-points of a communication request are subject to optimization. In this light, we can see our model as a "symmetric" version of online paging.

Graph partitioning problems are fundamental in computer science, and arise in many different contexts [27, 28].

## C. Our Contributions

This paper presents several new results on the online graph partitioning problem without augmentation. For the learning model, we present a lower bound of  $\Omega(k \cdot \ell)$  on the competitive ratio of any online deterministic online algorithm (that holds also in the general partitioning model). The best known lower bound so far was  $\Omega(k)$  [4] that holds only in the general partitioning model. We complement this result with an asymptotically optimal,  $O(k \cdot \ell)$ -competitive algorithm for the learning model.

For the general partitioning model, we design an asymptotically optimal,  $\Theta(\ell)$ -competitive algorithm for  $k = 3$ , improving the best known upper bound so far  $O(\ell^2)$  [4]. We further present a strictly 6-competitive algorithm for  $k = 2$  that improves upon the previous 7-competitive algorithm with  $O(\alpha \ell^2)$  additive constant.

All algorithms in this paper have a strict competitive ratio (i.e., without an additive term). Table I provides an overview of our contributions compared to prior work.

Variant	Lower bound	Upper bound
$k = 2$	3 [4]	6 (§III-B)
$k = 3$	$\Omega(\ell)$ (§II-A)	$O(\ell)$ (§III-A)
$k > 3$	$\Omega(k \cdot \ell)$ (§II-A)	$O(k^2 \cdot \ell^2)$ [4]
Learning model	$\Omega(k \cdot \ell)$ (§II-A)	$O(k \cdot \ell)$ (§II-B)

TABLE I: Overview of known results and our contributions. The table summarizes the results for the general partitioning model, except for the last row that summarizes the results for the learning model for arbitrary  $k$  and  $\ell$ .

**Algorithmic techniques.** The variant for general  $k$  is still unresolved, hence it is vital to summarize algorithmic techniques used in this paper. A straightforward analysis of the algorithm from Section III-A results in the bound  $O(\ell^2)$ , and the analysis must be improved in two places. One must bound the cost of each of  $O(\ell)$  reconfigurations per phase by a constant, and show that the optimum offline algorithm must pay a significant cost for inter-cluster requests. We bound the cost of the latter by estimating the capabilities of the optimum offline algorithm to *prepare* for an incoming sequence of requests. Furthermore, in the analysis of the algorithm from Section III-B, we propose a novel charging scheme for edges that share a vertex.

## II. THE LEARNING MODEL

In this section, we consider the learning variant of online balanced graph partitioning problem. For this setting, we show a surprisingly high lower bound of  $\Omega(k \cdot \ell)$  for  $k \geq 3$ . The lower bound holds also in the general partitioning model (studied in Section III). At the end of this section, we discuss an asymptotically optimal upper bound for the learning variant.

### A. Lower Bound

We provide a lower bound  $\Omega(k \cdot \ell)$  for the competitive ratio of any deterministic online algorithm for the learning problem. Later, we elaborate on how to efficiently transform it to a lower bound for the general partitioning problem. The lower bound requires  $k \geq 3$ . In contrast, for  $k = 2$  the learning problem is trivial: immediate collocation of communicating pairs is 1-competitive. In contrast, the general partitioning problem for  $k = 2$  is non-trivial (see Section III-B).

Throughout this paper, we often refer to groups of communicating nodes. We use this concept slightly differently in the lower bound than the upper bounds. In our algorithms, we group nodes with a communication history into *components*. In this section, we group nodes that may ever communicate, into *ground sets*.

A *ground set* is a subset of nodes that are collocated in the same cluster under a given perfect partition unknown to algorithms. If an algorithm keeps a ground set split at any point, we (as adversary) issue as many requests as it takes between non-located parts of the ground set until the algorithm collocates them. Observe that under such input construction every algorithm is forced to maintain a perfect partition of ground sets, otherwise it is not competitive. Henceforth, we assume every deterministic algorithm collocates the endpoints of every request as soon as it arrives.

With each new inter-cluster request a part of some ground set is revealed. The algorithm recovers the whole ground set and eventually the (hidden) perfect partition gradually by unifying and collocating separated parts on each inter-cluster request. We say that a ground set is a *singleton* if it contains exactly one node, which we refer to as an *isolated node*.

We start by constructing a ground set called  $B$  of size  $k - 1$  on an arbitrarily chosen cluster. In any balanced partition, the ground set of size  $k - 1$  must be collocated with some isolated node. We issue requests between this isolated node and some node that was collocated with it in the initial partition but currently separated, forcing the algorithm to collocate the two nodes which means collocating the ground set  $B$  with another isolated node. By issuing requests in this manner repeatedly, almost every node gets collocated with  $B$ . We generate the sequence in a way that it admits a perfect partition and the optimal offline algorithm OPT reaches this partition by performing only two node exchanges ("swaps").

**Theorem 1.** *The competitive ratio of any deterministic online algorithm for the learning model of Online Balanced Graph Partitioning is at least  $(k - 2)(\ell - 1)/2 - 2$  for any  $k \geq 3$  and  $\ell \geq 2$ .*

*Proof:* Fix any online algorithm ALG. For a subset of nodes  $C$  that are collocated in the initial partition, let  $I(C)$  denote the cluster where  $C$  resides initially. We refer to  $I(C)$  as the cluster of *origin* when  $C$  is clear from the context. Initially, all nodes are isolated, i.e., each node is in a singleton ground set. First, we choose a cluster arbitrarily and create a ground set  $B$  of  $k - 1$  nodes in this cluster. Each cluster hosts exactly  $k$  nodes, and in any feasible partition, a single

isolated node must be collocated with  $B$ . At any time, we refer to the isolated node currently collocated with  $B$  as the *pivot node*. Let  $x_0$  denote the first pivot node.

Then, we join the pivot node to a larger ground set to force its eviction. Precisely, we create a ground set  $\{x_0, y_0\}$ , where  $y_0$  is an arbitrary isolated node. Since ALG does not have  $\{x_0, y_0\}$  collocated, we issue an external request to this pair so that ALG collocates it. ALG cannot collocate  $\{x_0, y_0\}$  with  $B$  (as  $B$ 's size is  $k - 1$ ), hence it collocates them in a different cluster. In order to preserve a feasible partition of nodes after collocating  $\{x_0, y_0\}$ , ALG must replace  $x_0$  with another isolated node that becomes the new pivot.

We proceed in similar steps by joining the current pivot node to a ground set of the same origin residing in a different cluster. Consider the step  $i$ , when the isolated node  $x_i$  is collocated with  $B$ . We issue a request between  $x_i$  and some node in  $C_i$ , where  $C_i$  is the largest ground set s.t.  $I(C_i) = I(x_i)$ ,  $C_i \neq \{x_0, y_0\}$ . Then ALG must collocate the new ground set  $\{x_i\} \cup C_i$  in one cluster. Any feasible partition replaces  $x_i$  with some isolated node  $x_{i+1}$ , as the new ground set  $\{x_i\} \cup C_i$  may not be ever split. We terminate the process once the number of remaining isolated nodes is less than  $\ell + 3$ . At each step  $i$ , the number of isolated nodes decreases either by one or by two if  $C_i$  is a singleton. Therefore, once the process terminates, in any case at least  $\ell + 1$  isolated nodes are left.

Next, we argue that a feasible partition exists when the process terminates. This implies that a feasible partition exists after any earlier step as well. Since there are at least  $\ell + 1$  isolated nodes left, there must be two isolated nodes  $x^*$  and  $y^*$ , with the same cluster of origin, i.e.,  $I(\{x^*\}) = I(\{y^*\})$ . Consider the partition  $P^*$  obtained from the initial partition after swapping  $x_0$  and  $y_0$  with  $x^*$  and  $y^*$  (respectively). In this partition, the ground set  $\{x_0, y_0\}$  is collocated in the cluster  $I(\{x^*, y^*\})$ . Note that after the first request  $\{x_0, y_0\}$ , we issue requests only between nodes that have the same cluster of origin and all these nodes are collocated in  $P^*$ . Therefore all ground sets constructed so far are collocated in  $P^*$ , and it is a feasible partition.

Consider nodes  $x^*$  and  $y^*$  and the partition  $P^*$  obtained previously. OPT moves to  $P^*$  by performing only two node swaps. Precisely, OPT collocates  $\{x_0, y_0\}$  by swapping them with  $x^*$  and  $y^*$ . No ground set is split in  $P^*$  and OPT pays only for the two swaps.

ALG performs at least one swap at each step  $i$ , and some ground set grows. Consider any ground set  $C^* \neq B$  after the termination. This ground set has grown exactly  $|C^*| - 1$  times until the termination. Let  $\mathcal{S}$  be the set of all ground sets after the process terminates. Thus,  $\mathcal{S}$  includes ground sets  $B$ ,  $\{x_0, y_0\}$ , and (up to)  $\ell + 2$  singleton ground sets. Among the remaining ground sets in  $\mathcal{S}$ , no two ground sets have the same origin. Otherwise, the smaller ground set is either a singleton, which contradicts the bound  $\ell + 2$  on the number of singletons, or we have joined nodes to it at some step, contradicting our choice of the largest  $C_i$  at step  $i$ . Hence, there are at most  $\ell - 1$  such ground sets, one per possible cluster of origin, excluding the cluster containing  $B$ . Therefore,  $|\mathcal{S}| \leq 1 + 1 + (\ell + 2) +$

$(\ell - 1) = 2\ell + 3$ . Note that among all non-singleton ground sets in  $\mathcal{S}$ , only  $B$  does not grow during the process. Thus, the total number of times that a ground set in  $\mathcal{S}$  has grown is

$$\begin{aligned} \sum_{C^* \in \mathcal{S}} (|C^*| - 1) - (k - 1) &= \sum_{C^* \in \mathcal{S}} |C^*| - \sum_{C^* \in \mathcal{S}} 1 - (k - 1) \\ &\geq k\ell - (2\ell + 3) - (k - 1) = (k - 2)(\ell - 1) - 4, \end{aligned}$$

which bounds the number of swaps performed by ALG. The competitive ratio is then  $\text{ALG}/\text{OPT} \geq ((k - 2)(\ell - 1) - 4)/2$ . ■

**Lower bound for the general problem.** For a lower bound  $\Omega(k \cdot \ell)$  for the general partitioning problem, we continue to issue requests to split nodes of a ground set until the algorithm collocates them. Note that the ground sets constructed in our lower bound can be perfectly partitioned into the clusters. Hence, the optimal algorithm moves to a perfect partition at the beginning (where requests incur the cost 0), and its cost is bounded. This means that the algorithm must eventually collocate all nodes of a ground set to be competitive. We reveal the next ground set only after the collocation, hence we can repeat the analysis of the algorithm for the learning problem. Finally, we note that the construction is oblivious to the choice of the reconfiguration cost  $\alpha$ : we compare the number of node exchanges of ALG and OPT.

**Resource augmentation.** The majority of work on the online balanced partitioning problem so far [4, 6] focuses on the scenario with resource augmentation, where the clusters of an online algorithm are larger than the clusters of the offline optimal algorithm that we compare the performance to. We can adjust our construction to show a lower bound of  $\Omega(\ell)$  for a setting with resource augmentation.

Consider a partitioning problem with resource augmentation  $1 + 1/3 - \epsilon$ . Fix  $k$  divisible by 3, and construct 3 ground sets of size  $k/3$  in each cluster. Note that no more than 3 such ground sets fit in one cluster. Then, apply the construction from the lower bound for  $k = 3$ , using these ground sets in the way we used individual nodes. The cost of any algorithm (including OPT) scales up by  $k/3$ , and the lower bound  $\Omega(\ell)$  holds.

Finally, we note the possibility of improvement. The algorithm CREP [4] requires  $(2 + \epsilon)$ -augmentation to guarantee the competitive ratio independent of  $\ell$ . In contrast, our construction shows that the linear term  $\ell$  is inevitable if the augmentation is smaller than  $1 + 1/3$ .

## B. Upper Bound

We present an asymptotically optimal algorithm for the learning problem. The algorithm collocates a pair as soon as they communicate and it never separates them. In order to preserve collocated pairs, we employ the concept of components, introduced by Avin et al. [4].

We maintain subsets of frequently communicating nodes as *components*. Initially, each node constitutes a single-node component which we refer to as a *singleton* component, and the node in such component is an *isolated* node. We define larger components in terms of smaller components. Concretely,

given a (sub)sequence of requests, two components  $C_1$  and  $C_2$  *merge* into one component as soon as for some pair of nodes  $v_1 \in C_1$  and  $v_2 \in C_2$ , the frequency of requests  $\{v_1, v_2\}$  reaches a certain threshold through the sequence. We keep all nodes of a component always collocated in the same cluster, i.e., when we move a node, we move the whole component that contains it. A partition that has every component collocated is a *component respecting* partition.

In addition, we maintain a balanced partition of our components as long as such partition exists, a reminiscent of partitioning given integers into sets of equal sum [13], which is an NP-hard problem. In contrast, our partition is time-varying: two components are merged into one component once they communicate, and we adjust the partition accordingly.

The algorithm in this section and the algorithm for  $k = 3$  (cf. Section III-A) are modified versions of the algorithm DET from [4]. The difference is in the choice of partition after a component merge. In DET, the partition was arbitrary. In the algorithm for the learning model, we choose a component respecting partition closest to the initial partition. In the algorithm from Section III-A, we choose the partition closest to the current partition (a repartition of minimum cost). Since the component sizes are in  $O(n)$ , computing a component respecting partition for  $\ell = 2$  is feasible in polynomial time using dynamic programming [11], but is strongly NP-hard for  $\ell > 2$  [29]. However, we assume unlimited computational power and focus on competitiveness instead.

**Perfect Partition Learner algorithm.** Now we describe the algorithm PPL. On each inter-cluster request  $\{u, v\}$ , PPL creates new components by merging the two components that contain nodes  $u$  and  $v$ . In order to collocate nodes of the new component, PPL moves to a component respecting partition that minimizes the distance to the initial partition  $P_I$ . The scheme of the algorithm can be found in the appendix (Algorithm 1).

Fix the initial partition  $P_I := \{I_1, \dots, I_\ell\}$  and OPT's final partition  $P_F := \{F_1, \dots, F_\ell\}$ . The *distance* of a partition  $P = \{C_1, \dots, C_\ell\}$  from the initial partition, defined as  $\Delta(P) := \sum_{j=1}^{\ell} |C_j \setminus I_j|$ , is the number of nodes in  $P$  that do not reside in their initial cluster. In other words, at least  $\Delta(P)$  node migrations are required in order to reach the partition  $P$  from  $P_I$ , and thus  $\text{OPT} \geq \Delta(P_F)$ .

PPL never moves to a partition that is more than  $\Delta^* := \Delta(P_F)$  migrations away from  $P_I$ . This invariant latter ensures us that PPL does not pay too much while recovering  $P_F$ . We emphasize that a *repartitioning* by PPL replaces the current partition  $P$  with a perfect partition closest to  $P_I$ . This way PPL never moves to a partition beyond the distance  $\Delta^*$ .

**Property 1.** *Let  $P$  be any partition chosen by PPL at any time. Then,  $\Delta(P) \leq \Delta^*$ .*

**Lemma 1.** *The cost of each repartitioning by PPL is  $2 \cdot \text{OPT}$ .*

*Proof:* Let  $P_i$  denote the partition of PPL immediately after serving  $\sigma_i$ . Consider the repartitioning that transforms  $P_{t-1}$  to  $P_t$  upon the request  $\sigma_t$ . Let  $M \subset V$  denote the set

of nodes that migrate during this process. Let  $M^-$  and  $M^+$  denote the subset of nodes that, respectively, enter or leave their initial cluster during the repartitioning. Then,  $M = M^+ \cup M^-$ . Since at least  $|M^-|$  nodes are not in their initial cluster before the repartitioning (i.e., in  $P_{t-1}$ ), the distance before the repartitioning is  $\Delta(P_{t-1}) \geq |M^-|$ . Analogously, the distance afterward is  $\Delta(P_t) \geq |M^+|$ . Thus,  $|M| \leq \Delta(P_{t-1}) + \Delta(P_t)$ . By Property 1,  $\Delta(P_{t-1}) \leq \Delta^*$  and  $\Delta(P_t) \leq \Delta^*$ . Since  $\Delta^* \leq \text{OPT}$ , we obtain  $|M| \leq 2 \cdot \text{OPT}$ . ■

**Theorem 2.** PPL reaches the final partition  $P_F$  and it is  $(2 \cdot (k-1) \cdot \ell)$ -competitive.

*Proof:* On each inter-cluster request, the algorithm enumerates all component respecting  $\ell$ -way partitions of components that are in the same (closest) distance to  $P_t$ . That is, once it reaches a partition  $P$  at distance  $\Delta^* = \Delta(P)$ , it does not move to a partition  $P'$ ,  $\Delta(P') > \Delta^*$ , before it enumerates all partitions at distance  $\Delta^*$ . Therefore, PPL eventually reaches the partition  $P_F$  at distance  $\Delta^* = \text{OPT}$ . With each distinct request, the size of some component increases by one. For any cluster  $F_i \in P_F$ , we have  $\sum_{C \in F_i} |C| = k$ . A component  $C \in F_i$ , initially begins as an isolated node and it grows by gaining  $|C| - 1$  more nodes. Hence, the total number of times a component in  $F_i$  grows is  $\sum_{C \in F_i} (|C| - 1) \leq k - 1$ . Therefore, there are at most  $(k-1) \cdot \ell$  distinct requests for which PPL performs a repartitioning and PPL performs at most  $(k-1) \cdot \ell$  repartitions. By Lemma 1, each repartitioning costs at most  $2 \cdot \text{OPT}$ . The total cost is thus at most  $2 \cdot \text{OPT} \cdot (k-1) \cdot \ell$ , which implies the competitive ratio. ■

### III. GENERAL PARTITIONING MODEL

Now we discuss the general online model where the request sequence can be arbitrary. In Section III-A, we show an  $O(k \cdot \ell)$ -competitive algorithm for  $k = 3$  using the classic *rent-or-buy* approach [30]. Prior to this section, we showed a lower bound of  $\Omega(k \cdot \ell)$  that holds for the general model (cf. Section II-A), hence the result from this section is asymptotically optimal. Furthermore, in Section III-B, we show a strictly 6-competitive algorithm for  $k = 2$ .

#### A. Optimal Algorithm for Clusters of Size 3

The algorithm analyzed in this section is a modified version of the algorithm DET proposed by Avin et al. [4], which for  $k = 3$  is  $O(\ell^2)$ -competitive. In our algorithm, we choose the closest partition after a component merge instead of an arbitrary one. This allows to bound the cost of repartition by a constant (Lemma 2).

This modification alone is insufficient to obtain  $O(\ell)$ -competitive algorithm, and the analysis must be further improved. In particular, pairs of nodes that did not reach the collocation threshold  $\alpha$  (called external requests) incur the cost  $O(\ell^2)$  for the algorithm in each phase. The novel part of the analysis lower-bounds the cost of OPT on external requests while considering its savings from migrations and possibly different configuration at the beginning of the phase. This way,

we show that OPT paid a significant portion of the algorithm's cost on external requests.

**Component-based algorithm.** The algorithm  $\text{ALG}_3$  partitions nodes into components, and initially, each node is isolated (belongs to its own component). For each pair of nodes  $\{x, y\}$ ,  $\text{ALG}_3$  maintains a counter  $C_{\{x, y\}}$  and increments it on every external request between  $x$  and  $y$ . Once  $C_{\{x, y\}} = \alpha$ ,  $\text{ALG}_3$  merges the components of  $u$  and  $v$ , and moves to the closest component respecting partitioning. If no such partitioning exists,  $\text{ALG}_3$  resets all components to singleton components, resets all counters to 0, and ends the phase.

**Theorem 3.**  $\text{ALG}_3$  is  $60\ell$ -competitive for  $k = 3$ .

Before bounding the competitive ratio of  $\text{ALG}_3$ , we upper-bound the cost of a single repartition of  $\text{ALG}_3$ . In our analysis, we distinguish among three types of clusters:  $C_1, C_2$  and  $C_3$ . In a cluster of type  $C_i$ , the size of the largest component contained in this cluster is  $i$ .

**Lemma 2.** In a single repartition of nodes (after a merge of components),  $\text{ALG}_3$  exchanges at most two pairs of nodes.

*Proof:* If no component respecting partition exists after the merge of components, then  $\text{ALG}_3$  resets all components, ends the phase, and performs no repartition. It suffices to show that the merged component has size at least 4 to conclude that  $\text{ALG}_3$  incurs no cost.

Consider a request between  $u$  and  $v$  that triggered the repartition and let  $U$  and  $V$  be their respective clusters. The request triggered the repartition, hence it was external and  $U \neq V$ . We consider cases based on the types of clusters  $U$  and  $V$ .

If either  $U$  or  $V$  is of type  $C_1$ , then this cluster can fit the merged component, and the repartition is local within  $U$  and  $V$ , for the cost of at most 2 swaps. If either  $U$  or  $V$  is of type  $C_3$ , a component of size 3 participates in a merge, and we have a component of size at least 4, and  $\text{ALG}_3$  ends the phase with no repartition.

It remains to consider the case where both  $U$  and  $V$  are of type  $C_2$ . If  $(u, v)$  both belong to components of size 2, then the merged component has size 4, and  $\text{ALG}_3$  incurs no cost. Otherwise, if one of  $u, v$  belongs to a component of size 2, then it suffices to exchange components of size 1 between  $U$  and  $V$ . Finally, if  $u$  and  $v$  belong to components of size 1, then we must place them in a cluster different from  $U$  and  $V$ . Note that if  $C_1$ -type cluster does not exist, then no component respecting partitioning exists. Otherwise,  $\text{ALG}_3$  performs two swaps — it exchanges the nodes  $u$  and  $v$  with any two nodes of any cluster of type  $C_1$ .

In each case, we showed that a component respecting partition is reachable in at most two swaps. ■

*Proof of Theorem 3:* Fix a completed phase, and consider the state of  $\text{ALG}_3$ 's counters at the end of it (before the reset). We consider the incomplete phase later in this proof.

$\text{ALG}_3$  is component respecting, hence it never increases any counter above  $\alpha$ . We say that the pair  $(u, v)$  is *saturated* if

the counter's value is  $\alpha$ , and *unsaturated* otherwise (saturation of a pair leads to a merge action). By  $\sigma$  we denote the input sequence that arrived during the phase. In our analysis, we focus on the requests that were external to  $\text{ALG}_3$  at the moment of their arrival; these are the only requests that incur a cost for  $\text{ALG}_3$ . We denote these external requests by  $\sigma_{\text{cost}}$ . We partition the sequence  $\sigma_{\text{cost}}$  into subsequences  $\sigma_I$  and  $\sigma_E$ . The sequence  $\sigma_I$  (inter-component requests) denotes the requests from  $\sigma_{\text{cost}}$  issued to pairs that belong to the same component of  $\text{ALG}_3$  at the end of the phase. The sequence  $\sigma_E$  (extra-component requests) denotes the requests from  $\sigma_{\text{cost}}$  that do not appear in  $\sigma_I$ .

Let  $\text{ALG}_3(M)$  denote the cost of migrations performed by  $\text{ALG}_3$  in this phase. During the phase,  $\text{ALG}_3$  performs at most  $2\ell$  component merge operations — exceeding this number would mean that a component of size 4 exists, and the phase should have ended already. We bound the cost of each repartition after a merge by Lemma 2, obtaining  $\text{ALG}_3(M) \leq 8\alpha \cdot \ell$ .

We bound  $\text{ALG}_3(\sigma_I)$  by summing the intra-component counters of each cluster at the end of the phase. The sum of intra-component counters in a cluster of type  $C_3$  is at most  $3\alpha - 1$ : two pairs of nodes from the component are saturated and its counter is  $\alpha$  each, and the counter of the third, unsaturated pair is at most  $\alpha - 1$ . The sum of counters inside  $C_1$  is 0, and inside  $C_2$  it is  $\alpha$ . Summing over all  $\ell$  clusters gives us  $\text{ALG}_3(\sigma_I) \leq (3\alpha - 1) \cdot \ell \leq 3\alpha \cdot \ell$ .

Furthermore,  $\text{ALG}_3$  paid for all requests from  $\sigma_E$ , and thus  $\text{ALG}_3(\sigma_E) = |\sigma_E|$ . In total, the cost of  $\text{ALG}_3$  is at most  $\text{ALG}_3(\sigma_I) + \text{ALG}_3(\sigma_E) + \text{ALG}_3(M) \leq 11\alpha \cdot \ell + |\sigma_E|$  during this phase.

Now we lower-bound the cost of the optimal offline solution. To this end, we fix any optimal offline algorithm  $\text{OPT}$ . By  $\text{OPT}(\sigma_I)$  and  $\text{OPT}(\sigma_E)$  we denote the cost of  $\text{OPT}$  on requests from sequences  $\sigma_I$  and  $\sigma_E$ , respectively. Note that these costs are defined with respect to components of  $\text{ALG}_3$  in this phase. By  $\text{OPT}(M)$  we denote the cost of migrations performed by  $\text{OPT}$  in this phase.

The cost of  $\text{OPT}$  is lower-bounded by the cost of serving  $\sigma_I$  and the cost of serving  $\sigma_E$ . While serving these requests,  $\text{OPT}$  may perform migrations, and we account for them in both parts: we separately bound  $\text{OPT}$  by  $\text{OPT}(\sigma_I) + \text{OPT}(M)$  and  $\text{OPT}(\sigma_E) + \text{OPT}(M)$ . Combining those bounds and using the relation between the maximum and the average, we obtain the bound

$$\begin{aligned} \text{OPT} &\geq \max\{\text{OPT}(\sigma_I) + \text{OPT}(M), \text{OPT}(\sigma_E) + \text{OPT}(M)\} \\ &\geq (\text{OPT}(\sigma_I) + \text{OPT}(M))/2 + (\text{OPT}(\sigma_E) + \text{OPT}(M))/2. \end{aligned}$$

First, we show  $\text{OPT}(M) + \text{OPT}(\sigma_I) \geq \alpha$ . Assume that  $\text{OPT}$ 's partition is fixed throughout the phase (as otherwise  $\text{OPT}$  pays  $\alpha$  for a migration). The phase ended when the components of  $\text{ALG}_3$  could not be partitioned without splitting them. Hence, for every possible partition of  $\text{OPT}$ , there exists a non-allocated saturated pair, and  $\text{OPT}$  paid for  $\alpha$  requests that saturated the pair.

Next, we bound  $\text{OPT}(\sigma_E) + \text{OPT}(M)$ . The sequence  $\sigma_E$  accounts only for unsaturated edges, thus there are at most  $\alpha - 1$  requests to each pair in  $\sigma_E$ .  $\text{OPT}$  may have at most  $3\ell$  pairs of nodes collocated in its clusters, and thus avoid paying for  $3\ell \cdot (\alpha - 1)$  requests from  $\sigma_E$ . Hence, at least  $\chi := |\sigma_E| - 3\ell \cdot (\alpha - 1)$  requests from  $\sigma_E$  are external requests with respect to  $\text{OPT}$ 's configuration at the beginning of the phase. Faced with these requests,  $\text{OPT}$  may serve them remotely or perform migrations to decrease its cost. By swapping a pair of nodes  $(u, v)$ ,  $\text{OPT}$  collocates  $u$  with two nodes  $u', u''$ , and  $v$  with two nodes  $v', v''$ . This may allow serving requests between  $(u, u')$ ,  $(u, u'')$ ,  $(v, v')$  and  $(v, v'')$  for free afterward. Hence, by performing a single swap that costs  $2\alpha$ ,  $\text{OPT}$  may avoid paying the remote serving costs for at most  $4(\alpha - 1)$  requests from  $\sigma_E$ . The total cost of  $\text{OPT}$  is then at least

$$\text{OPT}(\sigma_E) + \text{OPT}(M) \geq \chi \cdot \frac{2\alpha}{4(\alpha - 1)} \geq \frac{|\sigma_E|}{2} - 2\alpha \cdot \ell.$$

Finally, to bound the competitive ratio, we transform the above inequality in the following way:  $|\sigma_E| \leq 2(\text{OPT}(\sigma_E) + \text{OPT}(M)) + 4\alpha \cdot \ell$ . For succinctness, let  $\xi := \text{OPT}(\sigma_E) + \text{OPT}(M)$ . Combining the bounds on the cost of  $\text{ALG}_3$  and  $\text{OPT}$  during each finished phase, the competitive ratio is

$$\frac{\text{ALG}_3(\sigma)}{\text{OPT}(\sigma)} \leq \frac{11\alpha \cdot \ell + |\sigma_E|}{\alpha/2 + \xi/2} \leq \frac{30\alpha \cdot \ell + 4 \cdot \xi}{\alpha + \xi} \leq 30\ell.$$

It remains to consider the last, unfinished phase. First, consider the case where the unfinished phase is also the first one. Then, we cannot charge  $\text{OPT}$  due to the inability to partition the components. Instead, we use the fact that  $\text{ALG}_3$  and  $\text{OPT}$  started with the same initial partition. If the input finished before the first  $\alpha$  external requests, then  $\text{ALG}_3$  is 1-competitive. If at least  $\alpha$  external requests were issued, then  $\text{OPT}$  either paid  $\alpha$  for serving them remotely or paid  $\alpha$  for a migration. Charging this cost to  $\text{OPT}$  serves the purpose of charging  $\alpha$  at the end of a finished phase, and thus we can repeat the analysis of a finished phase. Second, consider the case, where there are at least two phases, then we split the cost  $\alpha$  charged in the penultimate phase into the last two phases, and we repeat the analysis of a finished phase. This way, the competitive ratio increases at most twofold in comparison to a finished phase, and the competitive ratio is  $\text{ALG}_3(\sigma)/\text{OPT}(\sigma) \leq 60\ell$ . ■

**Distributed implementation.** While we have described the algorithm globally so far, we note that it allows for efficient distributed implementations. The algorithm performs two types of operations that require communication with other clusters: a component merge, and a broadcast of the end of the phase. We say that a cluster containing 3 isolated nodes is *fresh*. A merge of two components may require finding a fresh cluster (for details see the proof of Lemma 2). In the following, we show how to efficiently find a fresh cluster in a distributed manner. To this end, we organize the clusters into an arbitrary rooted balanced binary tree, and we broadcast the root to each cluster. Each cluster maintains the counter of

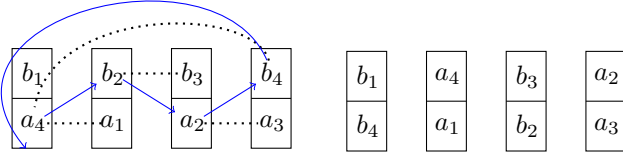


Fig. 1: Dashed lines represent requests. An arrow from node  $x$  to node  $y$  indicates that  $x$  replaces  $y$ . OPT collocates 4 pairs from the left partition, performing 4 migrations, resulting in the right partition.

fresh clusters in its subtree. To find a fresh cluster, we traverse an arbitrary path of non-zero counters from the root. Upon encountering a fresh cluster, we end the traversal and decrease the counters on the followed path by 1. Summarizing, ending the phase requires a single broadcast, and merging components has  $O(\log \ell)$  communication complexity.

### B. Improved Algorithm for Online Rematching

In this section, we present RM, an algorithm for clusters of capacity  $k = 2$ . We interpret a pair of nodes collocated in one cluster as a “matched” pair. Hence, the problem is an online variant of the maximal matching problem where a matched pair can separate in order to “rematch” with two other nodes. Rematching is necessary for maximizing intra-cluster communications, which is equivalent to minimizing inter-cluster communications. This is known as the *online rematching* problem and a non-strict 7-competitive algorithm is already given by [4], in which the ratio comes with an additive factor  $O(\alpha \ell^2)$ . We do not only improve upon their competitive ratio, but also show that our ratio holds *strictly* (i.e., with no additive factor). Our algorithm is slightly simpler than the one in [4], while our analysis is significantly simpler and more concise, thanks to the charging scheme we devise here.

**Algorithm ReMatch.** The algorithm ReMatch (RM) maintains a counter  $C_{\{x,y\}}$  for each pair of nodes  $\{x,y\}$  and increments it on every remote request between  $x$  and  $y$ . Once  $C_{\{x,y\}} = \lambda$ , it resets the counter  $C_{\{x,y\}} := 0$  and collocates the two nodes by swapping one of them, say  $x$ , with the node collocated with  $y$ .

**Theorem 4.** *For  $\lambda = \alpha$ , the algorithm RM is strictly 6-competitive.*

**The charging scheme.** We charge both OPT and RM whenever RM collocates a pair. RM collocates a pair always with a swap (that costs  $2\alpha$ ), while OPT may save some costs by collocating multiple pairs at once. Thus it pays the price of only one migration per pair (see Figure 1). For this reason, whenever OPT collocates a pair, we charge it only the cost  $\alpha$  of moving a single node to the other cluster (in contrast to the cost  $2\alpha$  incurred by RM).

Consider two pairs that share the same node, i.e. *intersecting pairs*, and the set of requests that cause (first) collocations of these pairs. This set contains at least one request to each pair and OPT must pay a non-zero cost over requests in this set because trivially it cannot have both pairs collocated at the same time. However, we can charge this cost to OPT only

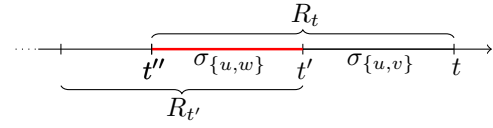


Fig. 2: Illustration of the timeline used in the proof of Theorem 4. Requests in  $R_t$  are only to  $\{u,w\}$  and  $\{u,v\}$ , which arrive during the interval  $(t'', t]$ . Similarly, requests in  $R_{t'}$  are to  $\{u,w\}$  and some other pair, irrelevant to the analysis. Hence, requests to  $\{u,w\}$  are included in at most two such sets, which are  $R_{t'}$  and  $R_t$ . This is because their intervals overlap on  $(t'', t']$ , shown with red thick line.

the first time RM collocates a pair, and not at any consequent time when RM collocates it a second time. Otherwise, OPT is possibly charged for the same cost repeatedly. For this reason, we charge OPT a cost inflicted by a pair if and only if OPT incurs that cost after the last time RM separates the pair.

*Proof of Theorem 4:* Fix an input sequence of requests  $\sigma := \{\sigma_1, \dots, \sigma_m\}$ . Assume that RM collocates a pair  $\{u,v\}$  at time  $t$ . The value of  $C_{\{u,v\}}$  at  $t$ , denoted  $C_{\{u,v\}}^t$ , reaches  $\lambda$  immediately before RM resets the counter. For any interval  $[t_1, t_2]$ , by  $\sigma_{\{x,y\}}[t_1, t_2]$  we denote the set of all requests to a pair  $\{x,y\}$  that arrive during  $[t_1, t_2]$ . We may use  $\sigma_{\{x,y\}}$  whenever the interval  $[t_1, t_2]$  is clear from the context.

If  $t$  is not the first time that RM collocates  $\{u,v\}$  then let  $0 < t' < t$  be the latest time when RM separates  $\{u,v\}$  in order to colocate some intersecting pair  $\{x,y\} \neq \{u,v\}$ ,  $\{x,y\} \cap \{u,v\} \neq \emptyset$ , e.g.,  $\{x,y\} = \{u,w\}$ . Else,  $t$  is the first time that RM collocates  $\{u,v\}$  and let  $t' := 0$ . Similarly, if  $t' > 0$  is not the first time that RM collocates  $\{u,w\}$  then let  $0 < t'' < t'$  be the latest time before  $t'$  when RM separates  $\{u,w\}$ . Else,  $t'$  is the first time that RM collocates  $\{u,w\}$  and we let  $t'' = 0$ .

First, we bound costs incurred by RM for requests that lead to the collocation of  $\{u,v\}$  at time  $t \in T$ , where  $T := \{i \in [1, m] \mid \exists \{x,y\} : C_{\{x,y\}}^i = \lambda\}$  is the set of times when RM performs a collocation. By definitions of  $t$  and  $t'$ , the overall cost of requests in  $\sigma_{\{u,v\}}$  incurred by RM, i.e., the total cost of remote serving and the moving cost is  $\lambda + 2\alpha$ . Next, we bound costs incurred by RM for requests that do not lead to collocations until the end of the sequence at  $t = m$ . Assume  $\{u,v\}$  is not collocated at  $t = m$  and  $0 < C_{\{u,v\}}^m < \lambda$ , which means RM pays  $C_{\{u,v\}}^m$  for requests in  $\sigma_{\{u,v\}}(t', m]$ . Then the overall cost of RM is  $\text{RM}(\sigma) = \sum_{t \in T} (\lambda + 2\alpha) + \sum_{\{u,v\}} C_{\{u,v\}}^m$ .

Next, we bound costs incurred by OPT for requests that trigger collocation of  $\{u,v\}$  at  $t \in T$ . If  $t$  is the first time that RM collocates  $\{u,v\}$ , then OPT pays  $\lambda$  for serving requests in  $\sigma_{\{u,v\}}[0, t]$  (remotely), or  $\alpha$  for collocating the pair and serving (some of) the requests with cost zero. Therefore in this case,  $\text{OPT}(\sigma_{\{u,v\}}(0, t]) \geq \min\{\lambda, \alpha\}$ . Otherwise, it is not the first collocation and consider times  $t'$  and  $t''$  as defined previously, and let  $R_t := \sigma_{\{u,w\}}(t'', t'] \cup \sigma_{\{u,v\}}(t', t]$ . We define  $R_{t'}$  for the collocation at  $t'$  analogously (see Figure 2). Then,  $\text{OPT}(R_t) = \text{OPT}(\sigma_{\{u,w\}}) + \text{OPT}(\sigma_{\{u,v\}})$ . If OPT has both pairs separated during their respective intervals, then



obviously it pays  $2\lambda$  during those intervals. Note that OPT cannot have both pairs collocated at the same time. Let us assume OPT has one of the pairs, e.g.  $\{u, v\}$ , collocated already prior its respective interval,  $(t', t]$ , and keeps it so during the interval. Then it pays zero while serving  $\sigma_{\{u, v\}}$ . Hence, it must pay  $\alpha$  for collocating the other pair, in this case  $\{u, w\}$ , or (resp., and) it pays (resp., up to)  $\lambda$  for serving (resp., some of) requests in  $\sigma_{\{u, w\}}$ . Therefore in any case,  $\text{OPT}(R_t) \geq \min\{\lambda, \alpha\} = \alpha$ .

It remains to bound the cost incurred by OPT due to requests to  $\{u, v\}$  that do not lead to its collocation until the end of the sequence at  $t = m$ . We bound the cost analogously to the case where RM collocates  $\{u, v\}$ . If  $\{u, v\}$  is not collocated in the initial matching and RM never collocates it, then  $C_{\{u, v\}}^m = |\sigma_{\{u, v\}}[1, m]|$ . OPT pays  $\text{OPT}(\sigma_{\{u, v\}}[1, m]) \geq \min\{\alpha, C_{\{u, v\}}^m\}$ , for collocating this pair or (and) paying for (resp. some of) requests in  $\sigma_{\{u, v\}}[1, m]$ . Else, either  $\{u, v\}$  is collocated in the initial matching or RM collocates it at some point. Then there exists an intersecting pair  $\{u, w\}$  that is collocated by RM at  $t' < m$ , separating  $\{u, v\}$ . We define times  $t'' < t' < m$  analogously to the former case. Let  $R_{\{u, v\}}^* := \sigma_{\{u, w\}}(t'', t') \cup \sigma_{\{u, v\}}(t', m)$ . Then, OPT must pay for collocating at least one pair or (and) serving requests to the other pair remotely. Thus,  $\text{OPT}(R_{\{u, v\}}^*) \geq \min\{C_{\{u, v\}}^m, \alpha\}$ .

Next, we sum up all costs incurred by OPT. By definitions of  $R_t$  and  $R_{\{u, v\}}^*$ , we have either  $R_{t'} \cap R_t = \sigma_{\{u, w\}}$  or  $R_{t'} \cap R_{\{u, v\}}^* = \sigma_{\{u, w\}}$ . This means,  $\text{OPT}(\sigma_{\{u, w\}})$  is counted at most twice in each of the expressions  $\text{OPT}(R_{t'}) + \text{OPT}(R_t)$  and  $\text{OPT}(R_{t'}) + \text{OPT}(R_{\{u, v\}}^*)$ . Hence, for all collocations performed by RM, and for final requests at  $t = m$ , OPT pays at least  $\frac{1}{2}(\sum_{t \in T} \text{OPT}(R_t) + \sum_{\{u, v\}} \text{OPT}(R_{\{u, v\}}^*))$ . Then, the total cost to OPT is

$$\begin{aligned} \text{OPT}(\sigma) &= \frac{1}{2} \left( \sum_{t \in T} \text{OPT}(R_t) + \sum_{\{u, v\}} \text{OPT}(R_{\{u, v\}}^*) \right) \\ &\geq \frac{1}{2} \left( \sum_{t \in T} \alpha + \sum_{\{u, v\}} C_{\{u, v\}}^m \right), \end{aligned}$$

and  $\text{RM}(\sigma)/\text{OPT}(\sigma) \leq$

$$2 \left( \sum_{t \in T} 3\alpha + \sum_{\{u, v\}} C_{\{u, v\}}^m \right) / \left( \sum_{t \in T} \alpha + \sum_{\{u, v\}} C_{\{u, v\}}^m \right) \leq 6.$$

#### IV. DISCUSSION AND FUTURE WORK

This paper revisited the online graph partitioning problem and presented several tight bounds for the important model where capacities cannot be exceeded, both for a general partitioning model and for a special learning model.

While our bounds are tight, there are several interesting avenues for future research. In particular, we have so far focused on deterministic algorithms, and it would be interesting to study the power of randomization in this context. On the practical side, it would also be interesting to study our algorithms empirically under realistic workloads.

Our algorithms allow for efficient distributed implementations. The algorithm PPL from Section II-B can be distributed similarly to the approach in [6]. The algorithm for  $k = 2$  from Section III-B performs only local communication for each request: counters are kept on the clusters and updated locally, and each migration is local within two clusters that reached the collocation threshold  $\lambda$ . Furthermore, we proposed an efficient distributed implementation of the algorithm for  $k = 3$  in Section III-A.

#### REFERENCES

- [1] A. Roy, H. Zeng, J. Bagga, G. Porter, and A. C. Snoeren, "Inside the social network's (datacenter) network," in *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*, 2015, pp. 123–137.
- [2] A. Singh, J. Ong, A. Agarwal, G. Anderson, A. Armistead, R. Bannan, S. Boving, G. Desai, B. Felderman, P. Germano *et al.*, "Jupiter rising: A decade of clos topologies and centralized control in google's datacenter network," *ACM SIGCOMM Computer Communication review*, vol. 45, no. 4, pp. 183–197, 2015.
- [3] J. C. Mogul and L. Popa, "What we talk about when we talk about cloud network performance," *ACM SIGCOMM Computer Communication Review*, vol. 42, no. 5, pp. 44–48, 2012.
- [4] C. Avin, A. Loukas, M. Pacut, and S. Schmid, "Online Balanced Repartitioning," *DISC*, pp. 243–256, 2016.
- [5] C. Avin, M. Bienkowski, A. Loukas, M. Pacut, and S. Schmid, "Dynamic Balanced Graph Partitioning," *arXiv e-prints*, p. arXiv:1511.02074, Nov. 2015.
- [6] M. Henzinger, S. Neumann, and S. Schmid, "Efficient distributed workload (re-)embedding," in *ACM SIGMETRICS / IFIP Performance 2019*, 2019.
- [7] A. Borodin and R. El-Yaniv, *Online Computation and Competitive Analysis*. Cambridge University Press, 1998.
- [8] M. Henzinger, S. Neumann, H. Raecke, and S. Schmid, "Tight bounds for online graph partitioning," in *Proc. ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2021.
- [9] G. Andrews and K. Eriksson, *Integer Partitions*. Cambridge University Press.
- [10] T. Forner, H. Raecke, and S. Schmid, "Online balanced repartitioning of dynamic communication patterns in polynomial time," in *Proc. SIAM Symposium on Algorithmic Principles of Computer Systems (APOCS)*, 2021.
- [11] C. Avin, L. Cohen, M. Parham, and S. Schmid, "Competitive clustering of stochastic communication patterns on a ring," *Computing*, vol. 101, no. 9, pp. 1369–1390, 2019.
- [12] K. Andreev and H. Räcke, "Balanced graph partitioning," *Theory of Computing Systems*, vol. 39, no. 6, pp. 929–939, 2006.
- [13] M. R. Garey, D. S. Johnson, and L. J. Stockmeyer, "Some Simplified NP-Complete Graph Problems," vol. 1, no. 3, pp. 237–267, 1976.
- [14] H. Saran and V. Vazirani, "Finding k cuts within twice the optimal," *SIAM Journal on Computing*, vol. 24, no. 1, pp. 101–108, 1995.
- [15] S. Arora, D. R. Karger, and M. Karpinski, "Polynomial time approximation schemes for dense instances of NP-hard problems," *Journal of Computer and System Sciences*, vol. 58, no. 1, pp. 193–210, 1999.
- [16] U. Feige, R. Krauthgamer, and K. Nissim, "Approximating the minimum bisection size (extended abstract)," in *Proc. 32nd ACM Symposium on Theory of Computing (STOC)*, 2000, pp. 530–536.
- [17] U. Feige and R. Krauthgamer, "A polylogarithmic approximation of the minimum bisection," *SIAM Journal on Computing*, vol. 31, no. 4, pp. 1090–1118, 2002.
- [18] R. Krauthgamer and U. Feige, "A polylogarithmic approximation of the minimum bisection," *SIAM Review*, vol. 48, no. 1, pp. 99–130, 2006.
- [19] H. Räcke, "Optimal hierarchical decompositions for congestion minimization in networks," in *Proc. 40th ACM Symposium on Theory of Computing (STOC)*, 2008, pp. 255–264.
- [20] D. Sleator and R. Tarjan, "Amortized efficiency of list update and paging rules," *Communications of the ACM*, vol. 28, no. 2, pp. 202–208, 1985.
- [21] A. Fiat, R. M. Karp, M. Luby, L. A. McGeoch, D. D. Sleator, and N. E. Young, "Competitive paging algorithms," *Journal of Algorithms*, vol. 12, no. 4, pp. 685–699, 1991.



- [22] L. McGeoch and D. Sleator, “A strongly competitive randomized paging algorithm,” *Algorithmica*, vol. 6, no. 6, pp. 816–825, 1991.
- [23] D. Achlioptas, M. Chrobak, and J. Noga, “Competitive analysis of randomized paging algorithms,” *Theoretical Computer Science*, vol. 234, no. 1–2, pp. 203–218, 2000.
- [24] L. Epstein, C. Imreh, A. Levin, and J. Nagy-György, “On variants of file caching,” in *Proc. 38th Int. Colloq. on Automata, Languages and Programming (ICALP)*, 2011, pp. 195–206.
- [25] L. Epstein, C. Imreh, A. Levin, and J. Nagy-György, “Online file caching with rejection penalties,” *Algorithmica*, vol. 71, no. 2, pp. 279–306, 2015.
- [26] S. Irani, “Page replacement with multi-size pages and applications to web caching,” *Algorithmica*, vol. 33, no. 3, pp. 384–409, 2002.
- [27] I. Stanton, “Streaming balanced graph partitioning algorithms for random graphs,” in *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, ser. SODA ’14, 2014, pp. 1287–1301.
- [28] D. Alistarh, J. Iglesias, and M. Vojnovic, “Streaming min-max hypergraph partitioning,” in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 1900–1908.
- [29] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, 1990.
- [30] A. Karlin, M. Manasse, and L. M. Karlin, “Competitive randomized algorithms for nonuniform problems,” *Algorithmica*, vol. 11, no. 6, pp. 542–571, 1994.

## APPENDIX

---

### Algorithm 1 Perfect Partition Learner (PPL)

---

```

For each node  $v$  create a singleton component  $C_v = \{v\}$ 
and add it to  $\mathcal{C}$ .
for each request  $\sigma_t = \{u, v\}, 1 \leq t \leq N$  do
    Let  $C_1 \ni u$  and  $C_2 \ni v$  be the components containing  $u$ 
    and  $v$ , respectively.
    if  $C_1 \neq C_2$  then
        Merge  $C_1$  and  $C_2$  into one component  $C'$  and  $\mathcal{C} =$ 
         $(\mathcal{C} \setminus \{C_1, C_2\}) \cup \{C'\}$ .
        if  $C_1$  and  $C_2$  are not collocated then
            Move to a partition closest to  $P_I$  and respecting all
            components in  $\mathcal{C}$ .
        end if
    end if
end for

```

---