

Optimal Online Balanced Graph Partitioning

Maciej Pacut¹, Mahmoud Parham^{2*} and Stefan Schmid^{1,2,3}

¹ TU Berlin, Germany.

² University of Vienna, Austria.

³ Fraunhofer SIT, Germany.

*Corresponding author(s). E-mail(s):

mahmoud.parham@univie.ac.at;

Contributing authors: maciej.pacut@inet.tu-berlin.de;

stefan.schmid@univie.ac.at;

Abstract

Distributed applications generate a significant amount of network traffic in datacenters. By collocating nodes (e.g., virtual machines) that communicate frequently so that they reside 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 online. This paper revisits the online balanced partitioning problem, introduced by Avin et al. at DISC 2016, that asks for an algorithm that strikes a trade-off between the benefits of collocation (i.e., reduced network traffic) and its costs (i.e., migrations). Our first contribution is a significantly improved deterministic lower bound of $\Omega(\mathbf{k} \cdot \ell)$ on the competitive ratio, where ℓ is the number of clusters and \mathbf{k} is the cluster size. The bound holds even for scenarios where the communication pattern can be perfectly partitioned so that all communications are internal to the clusters. We match this result with an asymptotically tight upper bound of $\mathcal{O}(\mathbf{k} \cdot \ell)$ for this scenario. For $\mathbf{k} = \mathbf{3}$, we contribute an asymptotically tight upper bound of $\Theta(\ell)$ for the case where the communication pattern can change arbitrarily over time. We improve the result for $\mathbf{k} = \mathbf{2}$ by providing a strictly $\mathbf{6}$ -competitive upper bound.

Keywords: Distributed applications, cloud computing, online algorithms, competitive analysis

1 Introduction

Data-centric applications, from distributed machine learning to distributed databases, generate a significant amount of network traffic in datacenters [1, 2]. The performance of these distributed applications often depends on the performance of the underlying communication networks [3, 4].

The virtualization of resources in datacenters introduces an intriguing opportunity to reduce the network traffic and improve performance. In particular, it becomes possible to adaptively *migrate* frequently communicating virtual machines (or containers) closer to each other, in a demand-aware manner. Such adaptive migrations, however, come at a cost (e.g., resource and time overhead) and introduce a tradeoff.

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. This is particularly challenging in a setting where the traffic can be bursty and change over time, in a hard to predict manner, as it is often the case in practice [5]. We are in the realm of *online algorithms and competitive analysis*, and ideally, the algorithm should perform closely to an optimal offline algorithm without requiring any information about future traffic demands.

1.1 Online Algorithms and Competitive Analysis

We measure the quality of the proposed solutions with competitive analysis [6], which suits well the networking problems that are online by their nature. The sequence of requests σ is revealed one-by-one, in an online fashion. Upon seeing a request, the algorithm must serve it without the knowledge of future requests.

We measure the performance of an online algorithm by comparing to the performance of an optimal offline algorithm. For a given sequence of requests σ , let $\text{ALG}(\sigma)$ be the cost incurred by a deterministic online algorithm ALG , and let $\text{OPT}(\sigma)$ be the cost incurred by an optimal offline algorithm OPT . In contrast to ALG that learns the requests one-by-one as it serves them, OPT has complete knowledge of the entire request sequence σ *ahead of time*. The goal is to design online algorithms that provide worst-case guarantees. In particular, ALG is said to be α -*competitive* if there is a constant β , such that for any input sequence σ it holds that

$$\text{ALG}(\sigma) \leq \alpha \cdot \text{OPT}(\sigma) + \beta.$$

Note that β cannot depend on input σ but can depend on other parameters of the problem, such as the number of nodes. The minimum α for which ALG is α -competitive is called the *competitive ratio* of ALG . We say that ALG is *strictly α -competitive* if additionally $\beta = 0$.

1.2 Model

The problem known as *online balanced graph repartitioning* was introduced by Avin et al. [7] at DISC 2016. A special variant of the general problem, called the *learning* model, was later introduced by Henzinger et al. [8]. We define our terminology and the problems considered in this paper as follows.

Preliminaries

We say that an assignment of nodes to clusters is a *partitioning* of nodes, and it represents the *configuration* of the algorithm. The reconfiguration or migration of nodes is a *repartitioning* operation. We say that a subset of nodes are *collocated* if they reside in the same cluster. A communication request between two nodes is *internal* if the nodes are collocated in the same cluster. Otherwise, they are in different clusters, and the request is *external*. Algorithms serve internal requests *locally* at cost 0 and external requests *remotely* at cost 1.

We refer to the graph structure of a request sequence as its *communication graph*, which contains an edge between pairs of nodes with at least one request in the sequence. A set of nodes belong to a *communicating (connected) component* if a path exists between them in the communication graph. A component is a *singleton* if it contains exactly one node, and we refer to it as an *isolated* node. The *origin* of a node is the cluster in which the node resides in the initial partitioning. For any subset of nodes C collocated in the initial partitioning, we let $I(C)$ denote their cluster of origin.

The General Partitioning Problem

We are given a set of n nodes, initially arbitrarily partitioned into $\ell \in \mathbb{N}^+$ clusters each of capacity $k = n/\ell$ nodes. The nodes interact in an online manner: we are given 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 fixed amount of data at time t . We sometimes refer to the nodes as virtual machines or processes, and we refer to the clusters as servers.

The cost of serving a request (u, v) depends on the relative positions of u and v : if they reside in the same cluster, the request costs 0, and it costs 1 otherwise. Before serving a request, an online algorithm may perform a *repartition* of nodes. That is, it may migrate any number of nodes to different clusters while ensuring the number of nodes in each cluster does not exceed k by more than ϵk for constant $\epsilon \geq 0$. We refer to the extra capacity as *resource augmentation*. We consider the problem without any augmentation for most of this paper, i.e., $\epsilon = 0$. The cost of migrating a single node is $\alpha \in \mathbb{N}^+$. The goal is to minimize the total cost of repartitions and to serve the requests. We use the terms “partition”, “partitioning”, and “configuration” interchangeably.

The Learning Variant

Consider a variant of the general partitioning problem, where each pair of nodes either never communicates, or they communicate indefinitely. Once

093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138

139 a communication request (u, v) arrives, any algorithm must keep u and v
 140 collocated for the rest of the input sequence. We assume that requests of σ
 141 constitute a communication graph that admits a *perfect partitioning*: a par-
 142 titioning that assigns exactly k nodes to each cluster, and no inter-cluster
 143 request ever occurs in this partitioning. Moreover, an optimal offline algorithm
 144 moves to this partitioning before serving the sequence, and it stays there per-
 145 manently. The goal of the algorithm is to *learn* the communication graph (as it
 146 is revealed one edge at a time) while serving requests without performing too
 147 many node migrations. In the learning model, any two communicating nodes
 148 must be collocated, and only the migration cost is relevant; for simplicity, we
 149 may scale it down to $\alpha = 1$ (our bounds hold for any $\alpha > 1$ as well).

150

151 **1.3 Related Work**

152

153 Closest to our work are those of Avin et al. [7] at DISC 2016 (on the gen-
 154 eral partitioning model) and Henzinger et al. (on the learning model) [8] at
 155 SIGMETRICS 2019 and SODA 2021 [9]. Recently, a polynomial-time online
 156 algorithm achieving the same competitive ratio as in [7] has been proposed
 157 by Forner et al. [10]. However, the focus of these papers is primarily on mod-
 158 els with resource augmentation: the online algorithm can use slightly larger
 159 clusters than the offline algorithm. Avin et al. showed that their lower bound
 160 $\Omega(k)$ holds even in a scenario with significant capacity augmentation, and they
 161 provided an algorithm with the competitive ratio $O(k \log k)$ using the $(2 + \epsilon)$ -
 162 augmented cluster capacity. Their ratio is independent of ℓ , which is impossible
 163 without significant resource augmentation.

164

165 In contrast, we study the non-augmented setting, where the nodes need to
 166 be perfectly balanced among the clusters. This assumption is more realistic,
 167 as it utilizes all the processing capacity of a datacenter instead of leaving some
 168 CPUs idle. This variant is significantly more challenging, as it is related to
 169 hard problems such as integer partitioning [11]. Not much is known about
 170 the setting without augmentation. For $k = 2$, Avin et al. [7] presented a 7-
 171 competitive algorithm with a substantial ($\Omega(\ell^2)$) additive constant. For $k > 3$,
 172 a $O(k^2 \cdot \ell^2)$ -competitive (phase-based) algorithm was given by [7]. Later, a more
 173 sophisticated analysis by Bienkowski et al. [12] improved the ratio of the same
 174 algorithm to $O(2^{O(k)} \cdot \ell)$, which is significant for its linear dependency on ℓ .
 175 The best known lower bound for the problem without augmentation is $\Omega(k)$ [7].
 176 Given our lower bound of $\Omega(k \cdot \ell)$ in this paper, the quest for closing the gap
 177 remains open. See Table 1 for known results.

178

179 The problem has also been studied in a weaker model where the adversary
 180 can only sample requests from a fixed distribution [13] over the edges of a ring
 181 communication graph.

182

183 The *static* offline version of the partitioning problem is known as the ℓ -
 184 *balanced graph partitioning problem*, where the entire communication graph
 is known in advance, and the task is to partition n nodes into ℓ clusters
 of capacity n/ℓ each, minimizing the number of inter-cluster edges. The prob-
 lem is NP-complete, and cannot even be approximated within any finite factor

Variant	Lower bound	Upper bound	
Learning, $k \geq 3$	$\Omega(k\ell), \epsilon = 0$ (Thm. 1)	$O(k\ell), \epsilon = 0$ (Thm. 3)	185
Learning, $k \geq 3$	$\Omega(\ell \log k), \epsilon \leq \frac{1}{32}$ [9]	$O(\ell \log k), \epsilon \in \Omega(1)$ [9]	186
Learning, $k \geq 3$	$\Omega(\ell), \epsilon < 1/3$ (Thm. 2)	$O(\log k), \epsilon > 1$ [9]	187
General, $k = 3$	$\Omega(\ell), \epsilon = 0$ (Thm. 4)	$O(\ell), \epsilon = 0$ (Thm. 6)	188
General, $k = 2$	$3, \epsilon = 0$ [7]	$6, \epsilon = 0$ (Thm. 7)	189
General, $k > 3$	$\Omega(k\ell), \epsilon = 0$ (Thm. 4)	$O(2^{O(k)}\ell), \epsilon = 0$ [12]	190

Table 1: Overview of our contributions and known results on the deterministic online partitioning problem.

unless P=NP [14]. The static variant where $\ell = 2$ corresponds to the minimum bisection problem, which is already NP-hard [15], and currently the best approximation ratio is $O(\log n)$ [16–21].

The studied problem is further related to some classic online problems. In particular, it is related to online paging [22–25], 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) [22]. 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* [26–28]. 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* endpoints of a communication request are subject to optimization. In this light, we can see our model as a “symmetric” version of online paging. The general problem additionally generalizes symmetric ski rental [29].

Graph partitioning and clustering problems are fundamental in computer science and arise in multiple contexts, see [30, 31].

1.4 Contributions

This paper presents several new results for the online graph partitioning problem. Table 1 provides an overview of our contributions compared to prior work.

For both the learning model and the general model, we obtain a lower bound of $\Omega(k \cdot \ell)$ on the competitive ratio of any online deterministic online algorithm (that also holds in the general partitioning model). This improves over the best known lower bound $\Omega(k)$ [7] that holds only in the general partitioning model. The generalized lower bound $\Omega(k \cdot \ell)$ for the learning model holds in the general model as well. We further adjust the lower bounds to show that the factor of $\Omega(\ell)$ is unavoidable even with a significant augmentation.

We complement these result with an asymptotically optimal, $O(k \cdot \ell)$ -competitive algorithm for the learning model. For the general partitioning

231 model and $k = 3$, we design an asymptotically optimal $O(\ell)$ -competitive
 232 algorithm (after the conference version of this paper, the upper bound was
 233 generalized to arbitrary ℓ by Bienkowski et. al. [12], yielding an $O(2^{O(k)} \cdot \ell)$ -
 234 competitive algorithm). We further present a strictly 6-competitive algorithm
 235 for $k = 2$ that improves upon the previous 7-competitive algorithm with
 236 $O(\alpha\ell^2)$ additive constant. All algorithms in this paper have a strict competitive
 237 ratio (cf. Section 1.1), which improves over previous results with substantial
 238 additive in terms of $(\alpha \cdot k \cdot \ell)^2$.

239

240 **2 The Learning Model**

241

242 We begin with the learning variant of online balanced graph partitioning
 243 problem. First, we show a surprisingly high lower bound of $\Omega(k \cdot \ell)$ against
 244 deterministic algorithms for $k \geq 3$. The lower bound holds also in the general
 245 partitioning model (see Theorem 4 for details). Second, we provide a determin-
 246 istic algorithm that asymptotically matches this lower bound for the learning
 247 model.

248

249 **2.1 Lower Bound**

250

251 *Overview of the construction.* At each time step, we issue a new request,
 252 depending on the configuration of the online algorithm. First, we issue requests
 253 between $k - 1$ nodes of some arbitrarily chosen cluster, and we refer to this
 254 communicating component as B . In any partitioning, the communicating com-
 255 ponent B is collocated with exactly one isolated node, called the *pivot*. Second,
 256 we issue a request between the pivot and an arbitrarily chosen node from
 257 a different cluster. Any algorithm must collocate these nodes (recall that we
 258 consider a learning variant), and it must place them in a different cluster than
 259 B 's (otherwise, its capacity would be violated). Note that after collocating
 260 them, another isolated node must take the place of the pivot. Third, we issue
 261 the request between the new pivot and an arbitrarily chosen node from the
 262 same cluster of origin as the pivot. We repeat the last step of the construction
 263 $\Theta(k \cdot \ell)$ times (exact condition to be determined), using the node isolated node
 264 collocated with B as a new pivot. The algorithm pays 1 per each such request.

265 We claim that this sequence is cheap to serve offline. Roughly, if an offline
 266 algorithm would reside in the initial configuration, only the requests between
 267 x_0 and y_0 would be external, and all the subsequent requests would be free.
 268 Consider an offline strategy that collocates x_0, y_0 by swapping them with some
 269 initially collocated nodes x^*, y^* that did not participate in any request. To
 270 make sure that such a pair always exists, we stop repeating the requests con-
 271 cerning the pivot while there are still two isolated nodes collocated on some
 272 server. We illustrate the construction in Figure 1.

273

274 **Theorem 1** *Any deterministic online algorithm for the learning model of online*
 275 *balanced graph partitioning and $k \geq 3$ has the competitive ratio of at least $(k - 2)(\ell -$*
 276 *$1)/2 - 2$.*

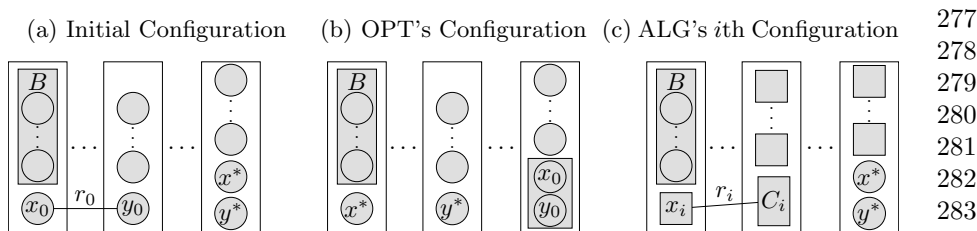


Fig. 1: Large rectangles represent clusters. Nodes are shown in gray circles, and gray rectangles represent components. Both ALG and OPT start in the configuration (a). OPT performs only two swaps and ends up at the configuration (b). At the beginning of the i th iteration, ALG is in the configuration (c) before evicting the i th pivot node x_i from the cluster of component B .

Proof Fix any online algorithm ALG. Initially, all nodes are isolated, i.e., each node is in a singleton communicating component. We issue requests one-by-one, in reaction to ALG's choices.

1. **Component B.** We issue requests among $k - 1$ nodes in an arbitrary cluster, and we refer to these nodes as a communicating component B . In any feasible partition, a single isolated node must be collocated with B (each cluster hosts exactly k nodes). We refer to the isolated node collocated with B at any time as the *pivot* node. Let x_0 denote the first pivot node.
2. **Request between nodes originating from different clusters.** We issue a request between x_0 and an arbitrarily chosen isolated node y_0 . This leads to the eviction of x_0 (otherwise, the algorithm incurs an arbitrarily large cost, while the optimal strategy is to collocate all communicating pairs). Hence, ALG must collocate this pair in a different cluster (cannot collocate it with B). In order to maintain a feasible partitioning of nodes after collocating $\{x_0, y_0\}$, ALG must replace x_0 with another isolated node, the new pivot node.
3. **Requests between nodes originating from the same clusters.** We continue to issue requests between the current pivot node and any node with the same origin as the pivot. Consider the i -th such request, and the isolated node x_i , collocated with B . Precisely, we issue a request between x_i and some node in C_i , where C_i is the largest communicating component s.t. $I(C_i) = I(x_i), C_i \neq \{x_0, y_0\}$. Then, ALG must collocate the communicating component $\{x_i\} \cup C_i$ in one cluster, and again, the algorithm replaces x_i with some isolated node x_{i+1} . We terminate the process once the number of remaining isolated nodes is smaller than $\ell + 3$. At each step i , the number of isolated nodes decreases either by one, or it decreases by two if C_i is a singleton. Therefore, once the process terminates, at least $\ell + 1$ nodes remain isolated.

To assure the correctness of this input sequence, we claim that the communicating components admit a feasible partition. Once we terminate, there are at least $\ell + 1$ isolated nodes left (the number of isolated nodes decreases at most by 2 at each

277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322

step). Therefore, two isolated nodes x^* and y^* exist with the same cluster of origin, $I(\{x^*\}) = I(\{y^*\})$. Consider a partition P^* , obtained from the initial partition after swapping x_0 and y_0 with x^* and y^* , respectively. In P^* , the pair $\{x_0, y_0\}$ resides in the cluster $I(\{x^*, y^*\})$. After the first request $\{x_0, y_0\}$, the requests are issued only between nodes originating from the same cluster, and all of these are collocated in P^* . This implies that no request is external in P^* , and we conclude that it is a feasible partition.

We bound the cost of ALG on the produced request sequence. At each request issued at step (3) of our construction, some communicating component grows, and ALG performs at least one swap. Let \mathcal{S} be the set of all communicating components created by issuing requests at step (3) of the construction. Each component $S \in \mathcal{S}$ grew exactly $|S| - 1$ times, each time joining an isolated node, and hence the number of nodes ALG swaps is at least

$$\text{ALG} \geq \sum_{S \in \mathcal{S}} (|S| - 1) = |\bigcup S| - |\mathcal{S}|.$$

In total, $\bigcup S$ contains all the $k \cdot \ell$ nodes excluding $k - 1$ nodes of B , the 2 nodes in $\{x_0, y_0\}$ and at most $\ell + 2$ singletons, which amounts to at least $k \cdot \ell - k - \ell - 3$ nodes. The set \mathcal{S} consists of at most $\ell - 1$ components, one per possible cluster of origin excluding B 's cluster of origin. Hence, the total number of swaps performed by ALG is

$$\text{ALG} \geq |\bigcup S| - |\mathcal{S}| = k \cdot \ell - k - 2\ell - 2 = (k - 2)(\ell - 1) - 4.$$

Finally, we bound offline algorithm's cost for the constructed sequence of requests. Consider an offline algorithm OPT that moves to P^* (described earlier in this proof) by performing only two node swaps. As argued earlier, no communicating component is split in P^* and OPT pays only for the two swaps. We combine the above arguments to conclude that the competitive ratio is bounded by $\text{ALG}/\text{OPT} \geq (k-2)(\ell-1)/2-2$. \square

We note that 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 3.4), a lower bound of 3 is known [7], and we provide a 6-competitive algorithm, see Section 3.4.

Later in Section 3.1, we elaborate on how to transform this construction to a lower bound for the general partitioning problem.

2.2 Lower Bound under Resource Augmentation

The majority of work on the online balanced partitioning problem so far [7–9] focuses on the scenario with resource augmentation, where the cluster capacity of an online algorithm is larger than that of the optimal offline algorithm to which we compare the performance. We say that an online algorithm uses augmentation $\epsilon > 1$ if each of its clusters has the capacity of $\epsilon \cdot k$ nodes. The number of all nodes remains $k \cdot \ell$.

By using a variant of the lower bound construction from Theorem 1, we show a lower bound of $\Omega(\ell)$ that holds even with significant resource augmentation.

Theorem 2 *With resource augmentation strictly smaller than $k/3$ (i.e., $\epsilon < 1/3$), the competitive ratio of any deterministic online algorithm for the learning model of online balanced graph partitioning is in $\Omega(\ell)$.*

Proof Fix k divisible by 3, and construct 3 communication components of size $k/3$ in each cluster. Consider any deterministic online algorithm with resource augmentation $\epsilon < 4/3$. Note that no more than 3 such communication components fit in one cluster. Then, apply the construction from the lower bound for $k = 3$ (Theorem 1), treating these communication components as individual nodes. The cost of any algorithm (including the optimal offline algorithm) scales up by $k/3$ due to the increased cost of moving entire components instead of individual nodes, and we conclude that the lower bound $\Omega(\ell)$ holds. \square

2.3 Upper Bound

We present an asymptotically optimal algorithm PPL (*Perfect Partition Learner*) for the learning model. The algorithm is a modification of the algorithm DET from [7]. The difference is in the choice of partition after a component merge. In DET, the choice of the partition was arbitrary, whereas our algorithm chooses an arbitrary partition closest to the initial partition that keeps all communicating nodes collocated. The algorithm PPL is listed in Algorithm 1.

To maintain the feasibility of the solution, the algorithm maintains *components* of communicating nodes. Initially, all nodes are in their own component, and upon a request between two nodes from different components, we *merge* the components. We maintain an invariant that the nodes of each communicating component are collocated. We say that a partition that collocates all nodes of all communicating components is a *communicating component respecting* partition.

Perfect Partition Learner

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 communication component-respecting partition that minimizes the distance to the initial partition P_I . We measure the distance between two configurations in the number of swaps to transform one configuration to the other. The distance to the initial partition is equivalent to the number of nodes that migrated from their cluster origin. Algorithm 1 describes the scheme of the algorithm.

Analysis

Fix a request sequence σ , the initial partition $P_I := \{I_1, \dots, I_\ell\}$ and an optimal offline strategy OPT with the final partition P_{OPT} . For each partition $P = \{C_1, \dots, C_\ell\}$ we define its *distance* from the initial partition as the number of nodes in P that do not reside in their initial cluster. Observe that at least

Algorithm 1 Perfect Partition Learner (PPL)

415 For each node v create a singleton component $C_v = \{v\}$ and add it to \mathcal{C} .
416 **for** each request $\sigma_t = \{u, v\}, 1 \leq t \leq N$ **do**
417 Let $C_1 \ni u$ and $C_2 \ni v$ be the components containing u and v ,
418 respectively.
419 **if** $C_1 \neq C_2$ **then**
420 Merge C_1 and C_2 into one component C'
421 and set $\mathcal{C} = (\mathcal{C} \setminus \{C_1, C_2\}) \cup \{C'\}$.
422 **if** C_1 and C_2 are not collocated **then**
423 Move to a component respecting partitioning, closest to P_I .
424 **end if**
425 **end if**
426 **end for**

429 $\Delta(P)$ node migrations are required in order to reach the partition P from P_I ,
430 and thus $\text{OPT}(\sigma) \geq \Delta(P_{\text{OPT}})$.
431

432 PPL replaces the current partition P with a perfect partition closest to P_I ,
433 thus it never moves to a partition that is more than $\Delta^* := \Delta(P_{\text{OPT}})$ migrations
434 away from P_I . Consequently, PPL never moves to a partition beyond the
435 distance Δ^* . We use this property to bound the cost of each repartitioning of
436 PPL.

437
438 **Property 1** Let P be any partitioning chosen by PPL at any time. Then, its
439 distance from the initial partitioning is $\Delta(P) \leq \Delta^*$.
440

441 **Lemma 1** The cost of each repartitioning by PPL is at most $2 \cdot \text{OPT}(\sigma)$, where
442 $\text{OPT}(\sigma)$ is the cost of an optimal offline algorithm for the request sequence σ .
443

444
445 *Proof* Let P_i denote the partition of PPL immediately after serving σ_i , the request
446 that arrives at time t . Consider the repartitioning that transforms P_{t-1} to P_t upon
447 the request σ_t . Let $M \subseteq V$ denote the set of nodes that migrate at t . Let M^-
448 and M^+ denote the subset of nodes that, respectively, enter or leave their initial
449 cluster during the repartitioning. In total, M^+ and M^- account for all migrations,
450 $M = M^+ \cup M^-$.

451 Since at least $|M^-|$ nodes are not in their initial cluster before the repartitioning
452 (i.e., in P_{t-1}), the distance from P_I before the repartitioning is $\Delta(P_{t-1}) \geq |M^-|$.
453 Analogously, the distance after the repartitioning is $\Delta(P_t) \geq |M^+|$. Thus, $|M| \leq$
454 $\Delta(P_{t-1}) + \Delta(P_t)$. By Property 1, $\Delta(P_{t-1}) \leq \Delta^*$ and $\Delta(P_t) \leq \Delta^*$. Since $\Delta^* \leq$
455 $\text{OPT}(\sigma)$, we conclude that the total cost of the algorithm is $|M| \leq 2 \cdot \text{OPT}(\sigma)$. \square

456
457 **Theorem 3** PPL is $(2(k-1) \cdot \ell)$ -competitive.
458
459
460

Proof Let $P_F := \{F_1, \dots, F_\ell\}$ be the partition of PPL after serving the sequence σ . At each request, the algorithm enumerates all communicating component-respecting ℓ -way partitions of components that are in the same (closest) distance to P_I . That is, once it reaches a partition P at distance $\Delta^* = \Delta(P)$, it does not move to a partition P' where $\Delta(P') > \Delta^*$ before it enumerates all partitions at distance Δ^* . Hence, P_F is at distance at most $\Delta^* = \text{OPT}(\sigma)$ from the initial partition.

We claim that PPL performs at most $(k-1) \cdot \ell$ repartitions while serving σ . Each component begins as a singleton, and with each request, the size of some component increases by one. Consequently, the number of repartitions of PPL is bounded by the number of times the components grow. Consider any cluster $F_i \in P_F$. Each cluster has the capacity k , thus the total number of times a component in F_i grows is at most $\sum_{C \in F_i} (|C| - 1) \leq k - 1$. Summing this bound over all ℓ clusters gives us at most $(k-1) \cdot \ell$ repartitions. By Lemma 1, each repartitioning costs at most $2 \cdot \text{OPT}(\sigma)$. The total cost of PPL is thus at most $2 \cdot (k-1) \cdot \ell \cdot \text{OPT}(\sigma)$, which implies the claim. \square

By Theorem 1, the lower bound for the competitive ratio of any deterministic algorithm is $\Omega(k \cdot \ell)$, and we conclude that PPL is asymptotically optimal.

Note on running time. Since the component sizes are in $O(n)$, computing a component-respecting partition for $\ell = 2$ is feasible in polynomial time using dynamic programming [13], but is strongly NP-hard for $\ell > 2$ [32]. However, we assume unlimited computational power and focus on competitiveness instead.

3 General Partitioning Model

In this section, we discuss the general online model where the request sequence can be arbitrary. First, in Section 3.1, we show a lower bound of $\Omega(k \cdot \ell)$ by generalizing the construction for the learning model from Section 2.1. Second, we generalize the lower bound for resource augmentation from Theorem 2. Third, in Section 3.3, we show an $O(k \cdot \ell)$ -competitive algorithm for $k = 3$. Finally, in Section 3.4, we show a strictly 6-competitive algorithm for $k = 2$.

3.1 Lower Bound

We generalize the lower bound construction of for the learning model in two dimensions. Recall that in the learning model, algorithms by assumption collocate any pair as soon as they communicate. Hence, we cannot apply the lower bound construction of Theorem 1 directly since algorithms may resist collocating such nodes. For the first dimension of our generalization, we impose a collocation by defining *ground sets* of nodes. Nodes that are assigned to the same ground set communicate as long as they are separated. For the second dimension of our generalization, we iterate the construction in batches. Repeating batches ensures that the input sequence with the claimed ratio can be arbitrarily long, and the algorithm cannot have a small competitive ratio with a large additive. After each batch, we ensure that the online algorithm resides in the same configuration as the optimal offline solution. A single batch

507 resembles the construction from Theorem 1, but in place of each request, we
 508 *reveal* a ground set: once the online algorithm splits a revealed set, we issue
 509 as many requests as it takes for it to collocate the split parts.

510 We start by revealing a ground set B of size $k - 1$ in an arbitrary cluster.
 511 Then, we reveal a cross-origin ground set with the pivot (the node collocated
 512 with B). Then, we repeatedly reveal ground sets of the current pivot and
 513 a node originating from the same cluster as the pivot. We repeat the last step
 514 of the construction $\Theta(k \cdot \ell)$ times (exact condition to be determined), using
 515 the isolated node collocated with B as a new pivot. The algorithm swaps at
 516 least one node at each such step.

517 We claim that this sequence is cheap to serve offline. Roughly, if an offline
 518 algorithm would reside in the initial configuration, only the requests between
 519 x_0 and y_0 would be external, and all the subsequent requests would be free.
 520 Consider an offline strategy that collocates x_0, y_0 by swapping them with some
 521 initially collocated nodes x^*, y^* that do not participate in any request. To make
 522 sure that such a pair always exists, we stop repeating the requests concerning
 523 the pivot while there are still two isolated nodes collocated on some server.

524 *Ground sets.* We construct our lower bound using *ground sets*. Instead of
 525 directly constructing the request sequence, we construct ground sets of nodes
 526 that start communicating if split by the online algorithm. This is possible since
 527 the algorithm is deterministic, and the adversary knows its configuration at
 528 any time. If the algorithm insists on keeping a ground set split, we continue to
 529 issue requests between non-located nodes of the ground set until the algo-
 530 rithm collocates them. Under such a request sequence, the algorithm must
 531 maintain a perfect partition of ground sets, as otherwise, it is not competi-
 532 tive. The ground sets are initially unknown to the algorithm, and the perfect
 533 partition is hidden from it. In contrast, the optimal offline algorithm knows
 534 the entire sequence in advance and may move to the perfect partition at the
 535 beginning.
 536

537
 538 **Theorem 4** *Any deterministic online algorithm for the general model of online*
 539 *balanced graph partitioning has the competitive ratio at least $(k - 2)(\ell - 1)/2 - 2$.*

540
 541
 542 *Proof* Fix any deterministic online algorithm ALG and any optimal offline algorithm
 543 OPT. We construct a sequence for the general model in *batches* of requests that can
 544 be repeated arbitrarily many times.

545 We construct the first batch by repeating the construction from Theorem 1,
 546 where in place of a request (u, v) , we merge the ground sets of u and v . This ensures
 547 that the algorithm collocates all nodes from each communicating component. If the
 548 algorithm does not collocate the nodes, it is not competitive as it pays an arbitrarily
 549 large cost for split ground sets. By the construction from Theorem 1, there exists
 549 a partitioning that collocates all the ground sets, and therefore any competitive
 550 algorithm eventually moves to such configuration.

551 After the batch finishes, we force the algorithm to move into the partitioning that
 552 is identical to OPT's partitioning. Let $\{C_1, C_2, \dots, C_\ell\}$ be OPT's configuration at

this point. We reveal additional ground sets C_i for $i \in [1, \ell]$ (i.e., containing all nodes that OPT has collocated). Note that these requests are free for OPT, and therefore OPT does not change its configuration. The batch ends once the algorithm moves to OPT's configuration. 553
554
555
556

Once the algorithm reaches OPT's configuration, we issue the next batch by repeating the construction. We may repeat issuing batches this way an arbitrary number of times. By applying similar reasoning to the proof of Theorem 1, in each batch the algorithm performs at least $(k \cdot \ell - k - 2\ell - 2)$ swaps, each for cost α . OPT performs at most two swaps in each batch. The competitive ratio from Theorem 1 holds for each batch separately, and therefore it holds for the entire sequence. \square 557
558
559
560
561
562

3.2 Lower Bound for Algorithms with Resource Augmentation 563 564 565

We generalize the lower bound for the learning model (Theorem 2) to show that the factor of $\Omega(\ell)$ is unavoidable even with significant augmentation. 566
567
568

Theorem 5 *With resource augmentation strictly smaller than $k/3$ (i.e., $\epsilon < 1/3$), the competitive ratio of any deterministic online algorithm for the general model of online balanced graph partitioning is in $\Omega(\ell)$.* 569
570
571
572
573

Proof Fix k divisible by 3, and construct 3 ground sets of size $k/3$ in each cluster. Consider any deterministic online algorithm with resource augmentation $1 + 1/3 - \epsilon$. Note that no more than 3 such ground sets fit in one cluster. Then, apply the construction from the lower bound (Theorem 4) for $k = 3$ using these communication components treating them as individual nodes. The cost of any algorithm (including the optimal offline algorithm) scales up by $k/3$ due to increased cost of moving entire ground sets instead of individual nodes, and we conclude that the lower bound $\Omega(\ell)$ holds. \square 574
575
576
577
578
579
580
581

3.3 Optimal Algorithm for Clusters of Size 3 582 583

For the setting with $k = 3$, Avin et al. [7] obtained a $O(\ell^2)$ -competitive algorithm. Their algorithm keeps track of external communication between nodes, and upon reaching a threshold α (the cost of migrating a node), we call the edge between them *saturated*. The algorithm keeps the invariant that the endpoint nodes of each saturated edge are collocated, and to this end, the algorithm tracks connected components consisting of nodes reachable via saturated edges. To collocate the nodes, the algorithm moves to an arbitrary partition where all nodes of all connected components are collocated. The algorithm operates in *phases*, and if no partition satisfying the invariant exists, it resets the counters for all pairs of nodes. 584
585
586
587
588
589
590
591
592
593

The algorithm ALG_3 presented in this section is a modified version of this algorithm. The difference is in the choice of partition after a component merge. Their algorithm chooses the partitioning arbitrarily, and our algorithm chooses the partition closest to the current partition (a repartition of minimum cost). 594
595
596
597
598

599 Our modification of the algorithm is straightforward, and our main con-
 600 tribution is an improved analysis of the algorithm. A straightforward analysis
 601 of our algorithm results in the bound $O(\ell^2)$, and we improve the analysis
 602 two aspects. First, we observe that the cost of each of $O(\ell)$ reconfigurations
 603 per phase is constant. Second, we deal with pairs of nodes that do not reach
 604 the threshold α (*unsaturated edges*). In the analysis of the algorithm from [7],
 605 these caused an additive $O(\alpha \cdot k^2 \cdot \ell^2)$. In our analysis, we observe that either the
 606 number of unsaturated edges is small, or any algorithm pays for a significant
 607 fraction of them. We bound the cost of the latter by estimating the capabil-
 608 ities of the optimal offline algorithm to *prepare* for an incoming sequence of
 609 requests.

610 *Saturated components.* We say that the pair (u, v) is *saturated* if the counter's
 611 value is α , and *unsaturated* otherwise (saturation of a pair leads to a merge
 612 action). We say that a partition that collocates all nodes of all saturated
 613 components is a *saturated component respecting* partition.
 614

615 *The algorithm ALG₃.* For each pair of nodes $\{x, y\}$, ALG₃ maintains a counter
 616 $C_{\{x,y\}}$ and increments it on every external request between x and y . Initially,
 617 each node is isolated (belongs to its own component). Once $C_{\{x,y\}} = \alpha$, ALG₃
 618 merges the components of u and v , and moves to the closest saturated com-
 619 ponent respecting partition. If no such partitioning exists, ALG₃ resets all
 620 components to singleton components, resets all counters to 0, and ends the
 621 phase.
 622

623 **Theorem 6** *ALG₃ is 60ℓ -competitive for $k = 3$.*
 624
 625

626 Before bounding the competitive ratio of ALG₃, we upper-bound the cost
 627 of a single repartition of ALG₃. In our analysis, we distinguish between three
 628 types of clusters: C_1, C_2 and C_3 . In a cluster of type C_i , the size of the largest
 629 component contained in this cluster is i .
 630

631 **Lemma 2** In a single repartition of ALG₃, it swaps at most two pairs of nodes.
 632
 633

634 *Proof* If no saturated component respecting partition exists after the merge of
 635 components, then ALG₃ resets all components, ends the phase, and performs no
 636 repartition. Thus it suffices to show that the merged component has a size at least 4
 637 to conclude that ALG₃ incurs no migration cost.

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

- 641 1. If either U or V is of type C_1 , then this cluster can fit the merged component,
 642 and the repartition is local within U and V , for the cost of at most 2 swaps.
 643
 644

2. 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 ALG_3 ends the phase with no repartition.
3. 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 ALG_3 incurs no cost. Otherwise, if one of u, v belongs to a component of size 2, then it suffices to swap 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 saturated component respecting partitioning exists. Otherwise, ALG_3 performs two swaps — it swaps the nodes u and v with any two nodes of any cluster of type C_1 .

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

Next, we observe that ALG_3 keeps saturated edges internal, and it increases counters only upon external communication, thus we upper-bound the ALG_3 's counter value for each unsaturated edge in any phase.

Observation 1 The external request counter for each unsaturated edge has a value at most $\alpha - 1$.

Now we are ready to bound the competitive ratio of ALG_3 .

Proof of Theorem 6 Fix a completed phase, and consider the state of ALG_3 's counters at the end of it (before the reset). By σ we denote the input sequence that arrived during the phase. We consider the incomplete phase later in this proof.

In our analysis, we focus on the requests that were external to ALG_3 at the moment of their arrival; these are the only requests that incur a cost for ALG_3 . We denote these external requests by σ_{cost} . We partition the sequence σ_{cost} into subsequences σ_I and σ_E . The sequence σ_I (inter-component requests) denotes the requests from σ_{cost} issued to pairs that belong to the same component of ALG_3 at the end of the phase. The sequence σ_E (extra-component requests) denotes the requests from σ_{cost} that do not appear in σ_I . Let $\text{ALG}_3(M)$ denote the cost of migrations performed by ALG_3 in this phase.

First, we bound the cost of ALG_3 in the phase. During the phase, ALG_3 performs at most 2ℓ component merge operations — exceeding this number would mean that a component of size 4 exists, and the phase would 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 α each, and the counter of the third, unsaturated pair is at most $\alpha - 1$ by Observation 1. The sum of counters inside C_1 is 0, and inside C_2 it is α . Summing over all ℓ clusters gives us $\text{ALG}_3(\sigma_I) \leq (3\alpha - 1) \cdot \ell \leq 3\alpha \cdot \ell$. Furthermore, ALG_3 paid for all requests from σ_E , and thus $\text{ALG}_3(\sigma_E) = |\sigma_E|$. In total, the cost of 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 the phase.

691 Second, we lower-bound the cost of the optimal offline solution. To this end, we
 692 fix any optimal offline algorithm OPT. By $\text{OPT}(\sigma_I)$ and $\text{OPT}(\sigma_E)$ we denote the
 693 cost of OPT on requests from sequences σ_I and σ_E , respectively. Note that these
 694 costs are defined with respect to components of ALG_3 in this phase. By $\text{OPT}(M)$ we
 695 denote the cost of migrations performed by OPT in this phase.

696 The cost of OPT is lower-bounded by the cost of serving σ_I and the cost of
 697 serving σ_E . While serving these requests, OPT may perform migrations, and we
 698 account for them in both parts: we separately bound OPT by $\text{OPT}(\sigma_I) + \text{OPT}(M)$
 699 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

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

703 First, we show $\text{OPT}(M) + \text{OPT}(\sigma_I) \geq \alpha$. Assume that OPT's partition is fixed
 704 throughout the phase (as otherwise OPT pays α for a migration). The phase ended
 705 when the components of ALG_3 could not be partitioned without splitting them.
 706 Hence, for every possible partition of OPT, there exists a non-located saturated
 pair, and OPT paid for α requests that saturated the pair.

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

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

719 Finally, to bound the competitive ratio, we transform the above inequality in
 720 the following way: $|\sigma_E| \leq 2(\text{OPT}(\sigma_E) + \text{OPT}(M)) + 4\alpha \cdot \ell$. For succinctness, let
 721 $\xi := \text{OPT}(\sigma_E) + \text{OPT}(M)$. Combining the bounds on the cost of ALG_3 and OPT
 722 during each finished phase, the competitive ratio is

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

726 It remains to consider the last, unfinished phase. First, consider the case where
 727 the unfinished phase is also the first one. Then, we cannot charge OPT due to the
 728 inability to partition the components. Instead, we use the fact that ALG_3 and OPT
 729 started with the same initial partition. If the input finished before the first α external
 730 requests, then ALG_3 is 1-competitive. If at least α external requests were issued, then
 731 OPT either paid α for serving them remotely or paid α for a migration. Charging
 732 this cost to OPT serves the purpose of charging α at the end of a finished phase, and
 733 thus we can apply the reasoning as for a finished phase. Second, consider the case,
 734 where there are at least two phases. Then we split the cost α of OPT accounted in the
 735 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 $\text{ALG}_3(\sigma)/\text{OPT}(\sigma) \leq 60\ell$. \square

Note on Arbitrary Capacity

The presented algorithm assumes the servers have capacity 3. The challenge in generalizing beyond this fixed capacity lies in bounding the cost of finding the closest saturated component respecting partition. For the capacity $k = 3$, we bound the cost of each reconfiguration by enumerating all cluster possibilities. We argue that the reconfiguration for $k = 3$ impacts only a constant number of clusters, and its cost is bounded only in terms of k , independently of ℓ . The major challenge at the time of our work was bounding the reconfiguration cost independently of ℓ , which was later addressed by Bienkowski et al. [12].

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. 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 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.

3.4 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 [7], in which the ratio comes with an additive factor $O(\alpha^2)$. We do not only improve upon their competitive ratio, but also we show that our ratio holds *strictly* (i.e., with no additive factor).

Our algorithm is slightly simpler than the one in [7], while our analysis is simpler and more concise. In the analysis of the algorithm, we propose a novel charging scheme for edges that share a vertex.

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

783 $C_{\{x,y\}} = \lambda$, it resets the counter $C_{\{x,y\}} := 0$ and collocates the two nodes by
 784 swapping one of them, say x , with the node collocated with y .

785

786 **Theorem 7** For $\lambda = \alpha$, the algorithm RM is strictly 6-competitive.

787

788

789

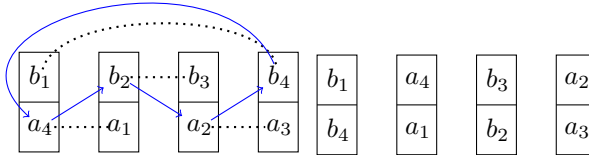
790

791

792

793

794



795 **Fig. 2:** Dashed lines represent external requests. An arrow from node x to
 796 node y indicates that x replaces y . In the left configuration, OPT collocates
 797 4 pairs by performing 4 migrations, which results in the right configuration.

798

799

800

801 *The Charging Scheme*

802

803 We charge both OPT and RM whenever RM collocates a pair. RM collocates
 804 a pair always with a swap which costs 2α , while OPT may save some costs by
 805 collocating multiple pairs at once. Thus it pays the price of only one migration
 806 per pair (see Figure 2). Therefore, OPT possibly collocates a pair by moving
 807 one node to the cluster of the other node paying only α , in contrast to the
 808 swapping cost 2α incurred by RM.

809 Consider two pairs that share the same node, i.e., *intersecting pairs*, and the
 810 set of requests that cause (first) collocations of these pairs. This set contains at
 811 least one request to each pair, and OPT must pay a non-zero cost over requests
 812 in this set, as it cannot have both pairs collocated at the same time. However,
 813 we can charge this cost to OPT only the first time RM collocates a pair and
 814 not at any later time when RM collocates it a second time. Otherwise, OPT
 815 is possibly charged for the same cost repeatedly. For this reason, we charge
 816 OPT a cost inflicted by a pair if and only if OPT incurs that cost after the
 817 last time RM separates the pair.

818 *Proof of Theorem 7* Fix an input sequence of requests $\sigma := \{\sigma_1, \dots, \sigma_m\}$. Assume
 819 that RM collocates a pair $\{u, v\}$ at time t . The value of $C_{\{u,v\}}$ at t , denoted $C_{\{u,v\}}^t$,
 820 reaches λ immediately before RM resets the counter. For any interval $[t_1, t_2]$, by
 821 $\sigma_{\{x,y\}}[t_1, t_2]$ we denote the set of all requests to a pair $\{x, y\}$ that arrive during
 822 $[t_1, t_2]$. We may use $\sigma_{\{x,y\}}$ whenever the interval $[t_1, t_2]$ is clear from the context.

823 If t is not the first time that RM collocates $\{u, v\}$ then let $0 < t' < t$ be the
 824 latest time when RM separates $\{u, v\}$ in order to collocate some intersecting pair
 825 $\{x, y\} \neq \{u, v\}$, $\{x, y\} \cap \{u, v\} \neq \emptyset$, e.g., $\{x, y\} = \{u, w\}$. Else, t is the first time that
 826 RM collocates $\{u, v\}$ and let $t' := 0$. Similarly, if $t' > 0$ is not the first time that RM
 827 collocates $\{u, w\}$ then let $0 < t'' < t'$ be the latest time before t' when RM separates
 828 $\{u, w\}$. 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 $T := \{\tau \in [1, m] \mid \exists \{x, y\} : C_{\{x, y\}}^\tau = \lambda\}$ 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 total cost incurred by 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 lead to the collocation of $\{u, v\}$ at $t \in T$. If t is the first time that RM collocates $\{u, v\}$, then OPT pays λ for serving requests in $\sigma_{\{u, v\}}[0, t]$ (remotely), or α 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 3). Then, $\text{OPT}(R_t) = \text{OPT}(\sigma_{\{u, w\}}) + \text{OPT}(\sigma_{\{u, v\}})$.

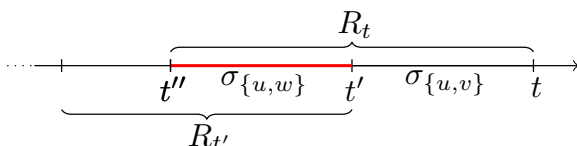


Fig. 3: Illustration of the timeline used in the proof of Theorem 7. The set R_t consists of requests to the pairs $\{u, w\}$ and $\{u, v\}$, which arrive in the interval $(t'', t]$. Similarly, $R_{t'}$ consists of requests to $\{u, w\}$ and some other pair irrelevant to the analysis. Hence, requests to $\{u, w\}$ are contained in both sets, and they arrive in $(t'', t']$ highlighted in red thick line.

If OPT has both pairs separated during their respective intervals, then clearly it pays 2λ in those intervals. Note that (trivially) OPT cannot have both pairs collocated at the same time. Else, OPT has one of the pairs, say $\{u, v\}$, already collocated prior its respective interval, $(t', t]$, and keeps it so during this interval. Then it pays zero for requests to this pair. Hence, it pays α for collocating the other pair, in this case $\{u, w\}$, or it pays λ for serving requests in $\sigma_{\{u, w\}}$. Note that OPT may deviate from the two cases by collocating the pair after serving some of its external requests. 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 it pays for 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, the total cost to OPT is

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

Finally, we bound the competitive ratio by aggregating the above bounds, obtaining

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

□

4 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.

Our algorithms allow for efficient distributed implementations. The algorithm PPL from Section 2.3 can be distributed similarly to the approach in [8]. The algorithm for $k = 2$ from Section 3.4 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 λ . Furthermore, we proposed an efficient distributed implementation of the algorithm for $k = 3$ in Section 3.3.

There remain several interesting avenues for future research. A general open direction concerns the study of the power of randomization in our setting. It would generally also be interesting to study the performance of our algorithms empirically, under realistic workloads, and engineer algorithms to speed up computations. But there are also more specific open questions. First open question regards the dependency on ℓ with resource augmentation. Our lower bounds from Theorems 2 and 5 shed light on the dependency on ℓ in the competitive ratio for the setting with augmentation. The algorithm CREP from [7] requires $(2 + \epsilon)$ -augmentation to guarantee the competitive ratio independent of ℓ . In contrast, our construction shows that the linear term ℓ is inevitable if the augmentation is smaller than $4/3$. This raises a question about the tradeoffs between augmentation and the dependency on ℓ for the competitive ratio. Another open question concerns the runtime. The algorithm Perfect Partition Learner runs in superpolynomial time, and the dominating term in

the runtime is due to finding the communicating component-respecting partition. We hence wonder if there exists a polynomial-time algorithm achieving an (asymptotically) optimal competitive ratio.

5 Acknowledgments

This project has received funding from the European Union's Horizon 2020 European Research Council (ERC), grant agreement No. 864228, 2020-2025.

References

- [1] Roy, A., Zeng, H., Bagga, J., Porter, G., Snoeren, A.C.: Inside the social network's (datacenter) network. In: Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication, pp. 123–137 (2015)
- [2] Singh, A., Ong, J., Agarwal, A., Anderson, G., Armistead, A., Bannon, R., Boving, S., Desai, G., Felderman, B., Germano, P., *et al.*: Jupiter rising: A decade of clos topologies and centralized control in google's datacenter network. ACM SIGCOMM Computer Communication review **45**(4), 183–197 (2015)
- [3] Mogul, J.C., Popa, L.: What we talk about when we talk about cloud network performance. ACM SIGCOMM Computer Communication Review **42**(5), 44–48 (2012)
- [4] Chowdhury, M., Zaharia, M., Ma, J., Jordan, M.I., Stoica, I.: Managing Data Transfers in Computer Clusters with Orchestra. ACM SIGCOMM **41**(4), 98–109 (2011)
- [5] Avin, C., Ghobadi, M., Griner, C., Schmid, S.: On the complexity of traffic traces and implications. In: Proc. ACM SIGMETRICS (2020)
- [6] Borodin, A., El-Yaniv, R.: Online Computation and Competitive Analysis, (1998)
- [7] Avin, C., Bienkowski, M., Loukas, A., Pacut, M., Schmid, S.: Dynamic balanced graph partitioning. SIAM J. Discret. Math. **34**(3), 1791–1812 (2020)
- [8] Henzinger, M., Neumann, S., Schmid, S.: Efficient distributed workload (re-)embedding. In: Proceedings of ACM SIGMETRICS / IFIP Performance 2019 (2019)
- [9] Henzinger, M., Neumann, S., Raecke, H., Schmid, S.: Tight bounds for online graph partitioning. In: Proc. ACM-SIAM Symposium on Discrete Algorithms (SODA) (2021)

- 967 [10] Forner, T., Raecke, H., Schmid, S.: Online balanced repartitioning of
968 dynamic communication patterns in polynomial time. In: Proc. SIAM
969 Symposium on Algorithmic Principles of Computer Systems (APOCS)
970 (2021)
971
- 972 [11] Andrews, G., Eriksson, K.: Integer Partitions. Cambridge University
973 Press
974
- 975 [12] Bienkowski, M., Böhm, M., Koutecký, M., Rothvoß, T., Sgall, J., Veselý,
976 P.: Improved analysis of online balanced clustering. Approximation and
977 Online Algorithms - International Workshop (WAOA) (2021)
978
- 979 [13] Avin, C., Cohen, L., Parham, M., Schmid, S.: Competitive clustering of
980 stochastic communication patterns on a ring. Computing **101**(9), 1369–
981 1390 (2019)
- 982 [14] Andreev, K., Räcke, H.: Balanced graph partitioning. Theory of Comput-
983 ing Systems **39**(6), 929–939 (2006)
984
- 985 [15] Garey, M.R., Johnson, D.S., Stockmeyer, L.J.: Some Simplified NP-
986 Complete Graph Problems **1**(3), 237–267 (1976)
987
- 988 [16] Saran, H., Vazirani, V.: Finding k cuts within twice the optimal. SIAM
989 Journal on Computing **24**(1), 101–108 (1995)
990
- 991 [17] Arora, S., Karger, D.R., Karpinski, M.: Polynomial time approximation
992 schemes for dense instances of NP-hard problems. Journal of Computer
993 and System Sciences **58**(1), 193–210 (1999)
994
- 995 [18] Feige, U., Krauthgamer, R., Nissim, K.: Approximating the minimum
996 bisection size (extended abstract). In: Proc. 32nd ACM Symposium on
997 Theory of Computing (STOC), pp. 530–536 (2000)
998
- 999 [19] Feige, U., Krauthgamer, R.: A polylogarithmic approximation of the
1000 minimum bisection. SIAM Journal on Computing **31**(4), 1090–1118
1001 (2002)
- 1002 [20] Krauthgamer, R., Feige, U.: A polylogarithmic approximation of the
1003 minimum bisection. SIAM Review **48**(1), 99–130 (2006)
1004
- 1005 [21] Räcke, H.: Optimal hierarchica decompositions for congestion mini-
1006 mization in networks. In: Proc. 40th ACM Symposium on Theory of
1007 Computing (STOC), pp. 255–264 (2008)
1008
- 1009 [22] Sleator, D., Tarjan, R.: Amortized efficiency of list update and paging
1010 rules. Communications of the ACM **28**(2), 202–208 (1985)
1011
- 1012 [23] Fiat, A., Karp, R.M., Luby, M., McGeoch, L.A., Sleator, D.D., Young,

- N.E.: Competitive paging algorithms. *Journal of Algorithms* **12**(4), 685–699 (1991) 1013
1014
1015
- [24] McGeoch, L., Sleator, D.: A strongly competitive randomized paging algorithm. *Algorithmica* **6**(6), 816–825 (1991) 1016
1017
1018
- [25] Achlioptas, D., Chrobak, M., Noga, J.: Competitive analysis of randomized paging algorithms. *Theoretical Computer Science* **234**(1–2), 203–218 (2000) 1019
1020
1021
1022
- [26] Epstein, L., Imreh, C., Levin, A., Nagy-György, J.: On variants of file caching. In: *Proc. 38th Int. Colloq. on Automata, Languages and Programming (ICALP)*, pp. 195–206 (2011) 1023
1024
1025
- [27] Epstein, L., Imreh, C., Levin, A., Nagy-György, J.: Online file caching with rejection penalties. *Algorithmica* **71**(2), 279–306 (2015) 1026
1027
1028
- [28] Irani, S.: Page replacement with multi-size pages and applications to web caching. *Algorithmica* **33**(3), 384–409 (2002) 1029
1030
1031
- [29] Karlin, A.R., Manasse, M.S., McGeoch, L.A., Owicki, S.S.: Competitive randomized algorithms for nonuniform problems. *Algorithmica* **11**(6), 542–571 (1994) 1032
1033
1034
1035
- [30] Stanton, I.: Streaming balanced graph partitioning algorithms for random graphs. In: *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms. SODA '14*, pp. 1287–1301 (2014) 1036
1037
1038
- [31] Alistarh, D., Iglesias, J., Vojnovic, M.: Streaming min-max hypergraph partitioning. In: *Advances in Neural Information Processing Systems 28*, pp. 1900–1908 (2015) 1039
1040
1041
1042
- [32] Garey, M.R., Johnson, D.S.: *Computers and Intractability; A Guide to the Theory of NP-Completeness*, (1990) 1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058