

Technical University of Munich

Department of Informatics

Bachelor's Thesis in Informatics

# **Polynomial Time Competitive Repartitioning of Dynamic Graphs**

**Tobias Forner**

Technical University of Munich

Department of Informatics

Bachelor's Thesis in Informatics

# Polynomial Time Competitive Repartitioning of Dynamic Graphs

## Kompetitive Repartitionierung Dynamischer Graphen in polynomieller Zeit

**Tobias Forner**

Supervisor:	Prof. Dr. Harald Räcke
Advisors:	Prof. Dr. Harald Räcke Univ.-Prof. Dr. Stefan Schmid
Submission Date	16.03.2020

I confirm that this bachelor's thesis is my own work and I have documented all sources and material used.

---

Date

---

Tobias Forner

## Abstract

This thesis studies the Dynamic Balanced Graph Partitioning problem, an online graph partitioning variant. The task is to maintain a mapping of communication nodes to servers while requests are revealed according to an unknown pattern. As the requests appear they need to be served, either remotely between nodes on different servers at cost one or locally within a server at cost zero. Before a request has to be served there is an option to change the mapping of nodes to servers, i.e. to migrate, at cost  $\alpha$  per node move.

We discuss different algorithmic approaches and the respective implications for the competitive analysis.

Our main contribution is a polynomial-time deterministic algorithm for this problem which achieves a competitive ratio of  $O(2/\epsilon \cdot k \log k)$ . We analyze and empirically evaluate this algorithm, also by comparison with another algorithm variant and with existing implementations.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Contributions . . . . .	1
1.2	Organization . . . . .	2
<b>2</b>	<b>General Definitions and Notation</b>	<b>2</b>
<b>3</b>	<b>Problem Definition</b>	<b>3</b>
<b>4</b>	<b>Related Work</b>	<b>4</b>
4.1	Dynamic Balanced RePartitioning . . . . .	4
4.2	Restricted Variants of Balanced RePartitioning . . . . .	4
4.3	Clustering . . . . .	5
4.4	Online Paging . . . . .	5
4.5	Static Balanced Graph Partitioning . . . . .	6
4.6	Fast Algorithms and Heuristics . . . . .	6
<b>5</b>	<b>Algorithmic Ideas</b>	<b>6</b>
5.1	CREP-CORE . . . . .	8
5.2	CREP-ADJ . . . . .	9
<b>6</b>	<b>Competitive Analysis of CREP-ADJ</b>	<b>11</b>
6.1	Algorithm Definitions . . . . .	12
6.2	Structural Properties . . . . .	12
6.3	Proof Overview . . . . .	14
6.4	Upper Bound On CREP-ADJ . . . . .	14
6.5	Lower Bound on OPT . . . . .	16
6.6	Competitive Ratio . . . . .	20
<b>7</b>	<b>From Algorithmic Ideas to Polynomial Time Implementation</b>	<b>21</b>
7.1	Algorithm Explanations . . . . .	21
7.2	Algorithm Pseudocode . . . . .	22
7.3	On the Decomposition of a Subgraph . . . . .	25
7.4	Running Time . . . . .	26
<b>8</b>	<b>Evaluation</b>	<b>27</b>
8.1	Reference Algorithms . . . . .	27
8.2	Input Data . . . . .	27
8.3	Comparison of CREP-ADJ and CREP-CORE . . . . .	28
8.4	On the Influence of the Connectivity Threshold . . . . .	29
8.5	Results of ADAPT and STATIC . . . . .	30
<b>9</b>	<b>Conclusion</b>	<b>32</b>
<b>10</b>	<b>Future Work</b>	<b>32</b>
	<b>References</b>	<b>34</b>

# 1 Introduction

With the continually increasing importance of distributed and cloud computation it is of great importance to minimize the communication overhead due to the increased network load. One method for achieving such improvements is to relocate communication partners (e.g. virtual machines) that communicate or interact frequently on the same servers. This way expensive and slow inter-server communication can be minimized by handling this communication within a server, without generating network load.

As communication patterns may change over time it seems especially promising to develop algorithms which adapt this distribution of communication partners over time, i.e. in a dynamic fashion.

The Dynamic Balanced Graph Partitioning (DBGP) problem ([2, 5]) which we study in this thesis provides a suitable model for these situations. Here algorithms are tasked to provide a mapping of communication nodes (virtual machines) to servers of fixed equal size and to dynamically adapt this mapping as new communication requests are revealed in an online manner, i.e. without prior knowledge of future requests. The goal is then to provide algorithms that maintain such a mapping while minimizing communication and migration costs.

In our model, a communication request can be served remotely, i.e. between nodes mapped to different servers, at cost one or locally if both nodes are located on the same server at no cost. Before the cost for the current request is paid, an algorithm has the option to migrate nodes at a cost of  $\alpha$  for each node move.

This problem has a variety of interesting challenges and properties ([5]):

- **Rent or buy** ([10]): should a request be served remotely at a small cost (rent) or should the communicating nodes be colocated at a large cost upfront (buy)?
- **How should nodes be colocated?** If there are nodes which communicate frequently it seems advantageous to colocate them. Then the question is how to define such a threshold and also which nodes should be moved in order to migrate as little as possible.
- **Should nodes be evicted?** If nodes are to be colocated there might not be enough space left in a suitable target server. In this case it needs to be decided which other nodes should be evicted from the target server.

## 1.1 Contributions

We study the Dynamic Balanced Graph Partitioning problem and show a competitive ratio of  $O(2/\epsilon \cdot k \log k)$  for our component-based deterministic algorithm CREP in the variant CREP-ADJ which runs in polynomial time. We also discuss challenges in the competitive analysis of this algorithm and of another algorithm, CREP-CORE, for which we have not shown a competitive ratio. Both approaches are adaptations of prior work ([2, 5]). Previous algorithms had super-polynomial running time.

Our connectivity-based approach to the definition of components allows to achieve a polynomial running time.

We discuss the details of our implementation which uses a decomposition tree data structure in order to reduce the update times when recomputing the components upon receiving new requests.

We evaluate our implementation and compare it with CREP-CORE that is also based on components, but for which we have not shown a competitive ratio. Furthermore we compare it with two existing implementations in the METIS and ParMETIS framework and we explore first ideas of adjusting parameters to improve the results of the presented algorithms even further.

## 1.2 Organization

We first introduce general definitions and notation in Section 2 which we will use throughout this thesis. We then formally define the Dynamic Balanced Graph Partitioning (DBGP) problem in Section 3. In Section 4 we discuss previous work related to the DBGP problem. After that we describe two algorithmic ideas, denoted by CREP-ADJ and CREP-CORE, which are inspired by prior work ([2, 5]) and the respective implications and challenges for their competitive analysis. In Section 6 we show the competitive ratio of  $O(2/\epsilon \cdot k \log k)$  for the algorithm CREP-ADJ. Then we discuss our implementation of CREP-ADJ and show that the algorithm can be implemented in polynomial time in Section 7. After that we empirically evaluate our algorithms by comparing their results with those of existing implementations in Section 8. We summarize the results of this thesis in Section 9. Finally we describe avenues for future work in Section 10.

## 2 General Definitions and Notation

In this section we introduce several definitions and notations that will be used throughout the thesis.

First we define a graph  $G = (V, E, [w])$  where  $V$  is a set of vertices and  $E \subseteq V \times V$  is a set of edges.  $w : E \rightarrow \mathbb{N}$  assigns each edge an (integer) weight. Note that we always assume that a graph is undirected, i.e. edges are of the form  $\{u, v\}$  where  $u \neq v$  and  $u, v \in V$ . A graph is only unweighted or directed if it is explicitly mentioned or clear from the context.

Given a graph  $G = (V, E, w)$  we define an (*edge*) *cut* of  $G$  as a pair of two disjoint subsets  $X, Y$  of  $V$  such that  $X \cup Y = V$ . The value of this cut is then the sum of the weight of edges between nodes from  $X$  and  $Y$ , i.e.  $\sum_{e=\{u,v\} \in E: u \in X, v \in Y} w(e)$  is the value of the cut  $(X, Y)$ . Note that such a cut can also be defined by the set of the edges connecting  $X$  and  $Y$  that are cut. We call a cut a *minimum (edge) cut* of  $G$  if it is one of the cuts with minimum value.

We use this definition of a minimum cut to introduce the notion of *connectivity*. The connectivity of a graph  $G$  is equal to the value of a minimum edge cut of  $G$ . This definition will be used in order to define the communication components our algorithm maintains as these are subsets of  $V$  which induce subgraphs of high connectivity. We explain the concept of components in greater detail in Section 5.

Furthermore we define the term  $(s, t)$ -*cut* as a cut  $(X, Y)$  for which  $s \in X$  and  $t \in Y$ , i.e. a  $(s, t)$ -cut separates the nodes  $s$  and  $t$  in  $G$ . Then a *minimum  $(s, t)$ -cut*

is a  $(s, t)$ -cut of minimum value. Note that a minimum  $(s, t)$ -cut is not necessarily a minimum cut.

Finally we call an injective function  $m : X \rightarrow Y$  a *mapping of  $X$  to  $Y$* . We use this terminology for example when we talk about the assignment of nodes to servers.

### 3 Problem Definition

Avin et al. ([5]) first introduced the Dynamic Balanced Graph Partitioning problem under the name of *Balanced RePartitioning* (BRP) problem, but changed the name to the former in a more recent version of their paper ([2]).

The task is to maintain a partitioning of a dynamic graph consisting of  $n = k \cdot l$  nodes that communicate with each other into  $l$  parts, each of size  $k$  while minimizing both the cost due to communication and due to node migrations defined as follows.

The communication cost is zero if both nodes are located on the same server at the time the request needs to be served and it is normalized to one if they are mapped to different servers. An algorithm may perform node migrations in order to change the mapping of nodes to servers prior to serving the communication request at time  $t$ . Each move of a vertex incurs cost  $\alpha > 1$ .

More formally we are given  $l$  servers  $V_0, \dots, V_{l-1}$ , each with capacity  $k$  and an initial perfect mapping of  $n = k \cdot l$  nodes to the  $l$  servers, i.e. each server is assigned exactly  $k$  nodes. An input sequence  $\sigma = (u_1, v_1), (u_2, v_2), \dots, (u_i, v_i), \dots$  describes the sequence of communication requests: the pair  $(u_t, v_t)$  represents a communication request between the nodes  $u_t$  and  $v_t$  arriving at time  $t$ . At time  $t$  the algorithm is allowed to perform node migrations at a cost of  $\alpha > 1$  per move. After the migration step, the algorithm pays cost 1 if  $u_t$  and  $v_t$  are mapped to different servers and does not pay any cost otherwise. Note that an algorithm may also choose to perform no migrations at all.

We are in the realm of competitive analysis and as a result we compare an online algorithm ONL to the optimal offline algorithm OPT. ONL only learns of the requests in the input sequence  $\sigma$  as they happen and as a result only knows about the partial sequence  $(u_1, v_1), \dots, (u_t, v_t)$  at time  $t$  whereas OPT has perfect knowledge of the complete sequence  $\sigma$  at all times.

The goal is to design an online algorithm ONL with a good competitive ratio with regard to OPT defined as follows.

An online algorithm ONL is  $\rho$ -competitive if there exists a constant  $\beta$  such that

$$\text{ONL}(\sigma) \leq \rho \cdot \text{OPT}(\sigma) + \beta \forall \sigma$$

where  $\text{ONL}(\sigma)$  and  $\text{OPT}(\sigma)$  denote the cost of serving input sequence  $\sigma$  of ONL and OPT respectively.

We allow the online algorithm to use larger capacities per server. In this case we speak of an *augmentation* of  $\delta$  in the case where the online algorithm is allowed to assign  $\delta \cdot n/k$  nodes to each server where  $\delta > 1$ . This augmented online algorithm is then compared with the optimal offline algorithm OPT which is not allowed to use any augmentation.



## 4 Related Work

With the ever increasing importance of distributed computation systems also come new challenges. One such challenge is an increase in network traffic. This increased network load may impact the performance of cloud computing networks significantly and should hence be minimized. For example, studies ([12, 13]) have shown that reconfigurable datacenter networks can provide performance similar to regular datacenter networks while having significantly lower cost in some cases.

This suggests that there is great potential in exploring methods of designing and controlling dynamic networks in order to reduce the overall network load.

One avenue of reducing this network load and hence also increasing the performance of cloud applications is to relocate the different communication partners inside the network dynamically in order to reduce the volume of inter-server communication.

### 4.1 Dynamic Balanced RePartitioning

Avin et al. ([2]) initiated the study of the Dynamic Balanced Graph Partitioning (DBGP) problem that is the topic of this thesis. They propose a deterministic algorithm for the DBGP problem with augmentation  $2 + \epsilon$  for any  $\epsilon > 1/k$  that takes exponential time.

They also show a lower bound of  $k - 1$  for the competitive ratio of any online algorithm for the Dynamic Balanced Graph Partitioning problem on two clusters via a reduction to online paging.

Furthermore they show that no  $\delta$ -augmented deterministic online algorithm can achieve a competitive ratio smaller than  $k$  for any augmentation  $\delta < l$ .

### 4.2 Restricted Variants of Balanced RePartitioning

Restricted variants of the Balanced RePartitioning problem have also been studied. Here one assumes certain restrictions of the input sequence  $\sigma$  and then studies online algorithms for these cases.

Avin, Cohen, Parham and Schmid ([3]) study one such case: the authors assume that an adversary provides requests according to a fixed distribution of which the optimal algorithm OPT has knowledge while an online algorithm that is compared with OPT has not. Further they restrict the communication pattern to form a ring-like pattern, i.e. for the case of  $n$  nodes  $0, \dots, n - 1$  only requests  $r$  of the form  $r = \{i \bmod n, (i + 1) \bmod n\}$  are allowed. For this case they present a competitive online algorithm which achieves a competitive ratio of  $O(\log n)$  with high probability.

Henzinger, Neumann and Schmid ([15]) study a special *learning variant* of the Dynamic Balanced Graph Partitioning problem specified above. In this version it is assumed that the input sequence  $\sigma$  eventually reveals a perfect balanced partitioning of the  $n$  nodes into  $l$  parts of size  $k$  such that the edge cut is zero. In this case the communication patterns reveal connected components of the communication graph of which each forms one of the partitions. Algorithms are tasked to *learn* this partition and to eventually collocate nodes according to the partition while minimizing communication and migration costs.

The authors of [15] present an algorithm for the case where the number of servers is  $l = 2$  that achieves a competitive ratio of  $O((\log n)/\epsilon)$  with augmentation  $\epsilon$ , i.e. each server has capacity  $(1 + \epsilon)n/2$  for  $\epsilon \in (0, 1)$ .

For the general case of  $l$  servers of capacity  $(1 + \epsilon)n/l$  the authors construct an exponential-time algorithm that achieves a competitive ratio of  $O((l \log n \log l)/\epsilon)$  for  $\epsilon \in (0, 1/2)$  and also provide a distributed version. Additionally the authors describe a polynomial-time  $O((l^2 \log n \log l)/\epsilon^2)$ -competitive algorithm for the case with general  $l$ , servers of capacity  $(1 + \epsilon)n/l$  and  $\epsilon \in (0, 1/2)$ .

It is important to stress that the assumption that the requests reveal a perfect partitioning of the communication nodes is not applicable for most practical applications and thus it is important to study the general DBGP problem without restricting  $\sigma$ .

### 4.3 Clustering

Clustering is the process of generating subsets of elements with high similarity ([14]). Desirable properties are *homogeneity*, i.e. high similarity within elements of a cluster and that the similarity is low between elements of different clusters, a property called *separation*. Hartuv and Shamir ([14]) approach this problem by constructing a similarity graph and then separating the nodes via repeated computations of minimum edge cuts. This approach is quite similar to the one we use in order to determine sub-graphs of high connectivity of the communication graph induced by the requests.

Clustering has also been studied within a variety of other different contexts and approaches from data mining to image segmentation ([6, 25, 19]).

However, we consider an online problem, i.e. algorithms need to react dynamically to changes in the graph and need to maintain their data structures and adapt accordingly whereas clustering considers complete data sets which are static.

### 4.4 Online Paging

In the Online Paging problem ([11], [9]) one is given a scenario with a fast cache of  $k$  pages and  $n - k$  pages in slow memory. Pages are requested in an online manner, i.e. without prior knowledge of future requests. If a requested page is in the cache at the time of the request it can be served without cost. If it is in slow memory however, then a *page fault* occurs and the requested page needs to be moved into the cache. If the cache is full then a page from the cache needs to be evicted, i.e. moved to the slow memory in order to make space for the requested one. The goal is to design algorithms which minimize the number of such page faults.

The Dynamic Balanced Graph Partitioning problem can be seen as a generalization of Online Paging. In fact, Avin et al. ([2]) have shown a reduction of Online Paging to this problem.

However, the standard version of Online Paging has no equivalent to the option of serving a request remotely as is possible in the Dynamic Balanced Graph Partitioning problem. The variant *with bypassing* allows an algorithm to access pages in slow memory without moving them into the cache, thus providing such an equivalent. It

is worth stressing however that in our problem requests involve two nodes while in Online Paging the nodes themselves are requested.

## 4.5 Static Balanced Graph Partitioning

The Static Balanced Graph Partitioning problem is the static offline variant of the topic of this thesis. In this version an algorithm may not perform any migrations, but has perfect knowledge of the request sequence  $\sigma$  and then needs to provide a perfectly balanced partitioning of the  $n = k \cdot l$  nodes into  $l$  sets of equal size  $k$  that minimizes cost, i.e. the weight of edges between the servers. This scenario can be modelled as a graph partitioning problem where the weight of an edge corresponds to the number of requests between its end points in the input sequence  $\sigma$ .

An algorithm then has to provide a partition of the nodes into sets of exactly  $k$  nodes each while minimizing the total edge weights between partitions, i.e. an algorithm needs to minimize the edge cut of the graph.

This problem is NP-complete ([1]) and for the case where  $l \geq 3$ , Andreev and Räcke ([1]) have shown that there is no polynomial time approximation algorithm which guarantees a finite approximation factor unless  $P=NP$ .

## 4.6 Fast Algorithms and Heuristics

There are several algorithms and frameworks for graph partitioning problems. Usually these frameworks employ heuristics in order to achieve their results. The most successful such heuristic is *Multilevel Graph Partitioning* ([7]). This method consists of three phases. Initially the graph is repeatedly coarsened into a hierarchy of smaller graphs in such a way that cuts in the coarse graphs also correspond to cuts in the finer graphs. On the coarsest level a (potentially expensive) algorithm is used in order to compute an initial partition. This partitioning is then transferred to the finer graphs. In this process one usually uses other local heuristics in order to improve the partition quality even further with every step.

METIS ([16, 17]) and Jostle ([23, 24]) are examples of libraries that utilize this multilevel approach. We choose METIS as a reference for our empirical evaluation.

## 5 Algorithmic Ideas

In this section we describe two different solution approaches to the Dynamic Balanced Graph Partitioning problem. We call these approaches CREP-CORE and CREP-ADJ. Both approaches are inspired by the versions of the CREP algorithm in [2] and [5]. We first describe the general approach that is common to both algorithms and then address the specific differences and analysis ideas.

Both methods share a similar concept at their core: a second-order partitioning of the communication nodes into *communication components* which represent node-induced sub-graphs of the original communication graph given by the requests from the input sequence  $\sigma$ . As more requests from  $\sigma$  are revealed to the algorithms they merge the corresponding components once they are suitably connected and relocate

the nodes of the new component in such a way that all the nodes of a component are always located on the same server.

More formally, initially each node forms a singleton component, but as the input sequence  $\sigma$  is revealed to the algorithms new communication patterns unfold. The algorithm keeps track of these patterns by maintaining a graph in which the nodes represent the actual communication nodes and the weighted edges represent the number of communication requests between nodes that were part of different components at the time of the request, i.e. for edge  $e = \{u, v\}$ ,  $w(e)$  represents the number of paid communication requests between  $u$  and  $v$ . We say that a communication request between nodes  $u$  and  $v$  is *paid* if the nodes are located on different servers at the time of the request.

Both algorithms merge a set  $S$  of components into a new component  $C$  if the connectivity of the component graph induced by the components in  $S$  is at least  $\alpha$ . After each edge insertion the algorithm checks whether there exists a new component set  $S$  with  $|S| > 1$  which fulfils this requirement.

If after any request and the insertion of the resulting edge the algorithm discovers a new subset  $S$  of nodes whose induced subgraph has connectivity at least  $\alpha$  and which is of cardinality at most  $k$  it merges the components that form this set into one new component and collocates all the nodes in the resulting set on a single server. The algorithm reserves additional space  $\min\{\lfloor \epsilon \cdot |C| \rfloor, k - |C|\}$  for each component on the server it is currently located on. Note that the additional reservation may be zero for components smaller than  $1/\epsilon$ . This reservation guarantees that nodes are not migrated too often for the analysis to work. Note that this also limits the total space a component can use to a maximum of  $k$ . This makes sense as a component whose size exceeds  $k$  is never merged but is deleted instead and hence there would be no benefit to a component taking space more than  $k$ .

To this end the algorithms keep track of the reservations for each component.

Both algorithms use augmentation  $2 + \epsilon$  in order to guarantee that the collocation of such component sets of at most  $k$  individual communication nodes is always possible without moving a node not in  $C$ . This guarantees by an averaging argument that there is always at least one cluster with capacity at least  $k$  which a newly merged component can be moved to.

If the subset has cardinality greater than  $k$  the resulting component is deleted. The definition of this deletion process is the main difference between our algorithms which we discuss in the following subsections. The common part of both algorithms is also summarized in the form of pseudocode in Algorithm 1. Note that the subroutine *delete*( $Y$ ) of a component set  $Y$  is different for each of the algorithms. In the pseudocode description we denote the reservation of a component  $C$  by  $res(C)$  and the current server it is mapped to by  $serv(C)$ . The free capacity of a server  $i$  is denoted by  $cap(i)$ .

We also describe the particular challenges each approach entails when it comes to the competitive analysis in their respective subsections.

These approaches are fairly similar to the algorithms defined in previous work ([2], [5]). The main differentiating factor is that we merge once a component set reaches connectivity  $\alpha$  while prior approaches do so once the component set reaches a certain density threshold. More specifically they merge a component set  $S$  once it

---

**Algorithm 1** CREP

---

Initialize an empty graph on  $n$  nodes  
turn each of the  $n$  nodes into a singleton component  
**for all**  $r = \{u, v\} \in \sigma$  **do**  
  **if**  $\text{comp}(v) \neq \text{comp}(u)$  **then**  
     $w(\{u, v\}) \leftarrow w(\{u, v\}) + 1$   
  **end if**  
  **if**  $\exists$  component set  $X$  with connectivity at least  $\alpha$  and  $|X| > 1$  and  $\text{nodes}(X) \leq k$  **then**  
     $\text{mergeAndRes}(X)$   
  **end if**  
  **if**  $\exists$  component set  $Y$  with connectivity at least  $\alpha$  and  $\text{nodes}(Y) > k$  **then**  
     $\text{delete}(Y)$  //to be specified later  
  **end if**  
**end for**

---

---

**Algorithm 2**  $\text{mergeAndRes}(X)$ 

---

**for all**  $C \in X$  **do**  
   $\text{cap}(\text{serv}(C)) \leftarrow \text{cap}(\text{serv}(C)) + \text{res}(C)$   
**end for**  
 $N \leftarrow \text{collocate}(X)$  //moves all components from  $X$  to the same server as described  
  
// $N$  contains the newly created component  
**if**  $|N| > 2/\epsilon$  **then**  
  //reserve additional space  
   $\text{res}(N) \leftarrow \min\{\lfloor \epsilon \cdot |N| \rfloor, k - |N|\}$   
   $\text{cap}(\text{serv}(N)) \leftarrow \text{cap}(\text{serv}(N)) - \text{res}(N)$   
**end if**

---

fulfills  $w(S) \geq (|S| - 1) \cdot \alpha$  where  $w(S)$  denotes the cumulative weight of the edges between nodes contained in the components of  $S$ . The similarities are especially apparent when comparing the respective lemmas and properties that are used in order to bound the weight between partitions of mergeable component sets. These are Lemma 4.3 in [2], Property 3 in [5] and Lemma 6 in this thesis.

Our connectivity-based approach however allows us to achieve a polynomial run-time.

Now we describe the differences in the deletion steps of CREP-CORE and CREP-ADJ and present the implications for their respective competitive analysis.

## 5.1 CREP-CORE

We address CREP-CORE first. In this version edges are reset which are *contained* in the deleted component, i.e. all edges  $e = \{u, v\}$  are reset to zero if both  $u$  and  $v$  were contained in component  $C$  at the time of its deletion. This approach resembles the one suggested by Avin et al. ([2]) and is also written in pseudocode in Algorithm 3.

The idea of the analysis is then to relate the cost of CREP-CORE with the cost of

---

**Algorithm 3** delete( $Y$ ) of CREP-CORE

---

```
for all  $e = \{u, v\} \in E$  do
  if  $u \in Y$  and  $v \in Y$  then
     $w(e) \leftarrow 0$ 
  end if
end for
for all  $C \in Y$  do
   $cap(serv(C)) \leftarrow cap(serv(C)) + res(C)$ 
   $res(C) \leftarrow 0$ 
end for
```

---

OPT by considering the respective costs due to requests from a deleted component  $C$  in the solution of CREP-CORE as these are of high connectivity and are impossible for OPT to collocate on one server as each deleted component contains more than  $k$  nodes. Then one could sum these costs over all such deleted components in order to bound the total cost. For this approach to work, however, one would have to find a way to cleanly separate requests belonging to one deleted component from those belonging to another in order to establish a lower bound on the cost of OPT.

In this case the edges which are adjacent to a deleted component  $C$  remain even after the deletion of  $C$  and may then contribute to the creation of a new component  $D$ . It is now very challenging to attribute any significant cost to OPT for these requests.

Figure 1 shows an example sequence of requests for which this approach does not work. The diagram shows horizontal lines, each representing one of the vertices. A vertical line represents a communication request between its end points. For example the sequence shown contains a communication request between nodes 1 and 2 at time  $t = 5$ . Now consider the case where  $\alpha = 3$  and  $k = 3$ .

In this sequence the first two requests between nodes 0 and 4 happen without leading to a merge. The following 12 requests lead to a merge of a new component  $C$  consisting of the nodes 1, 2, 3 and 4 which gets deleted by CREP-CORE at time  $t = 14$ . Note now that the edges corresponding to the requests between nodes 0 and 4 are still present even after this deletion. Finally the request at time  $t = 15$  leads to a merge of nodes 0 and 4.

One can see that such cases may also happen at a much larger scale, where almost all requests happen much earlier than the time at which an actual merge of the nodes involved in the requests happens.

These problems suggest that the approach of accounting cost for each component deletion in order to achieve a competitive ratio may be very challenging or even impossible. It remains to be seen whether there is another analysis approach for this case that may be more promising.

## 5.2 CREP-ADJ

The second algorithm, CREP-ADJ, resets all the edges contained in the deleted component  $C$  but also resets the weights of edges *adjacent* to  $C$ , i.e. all edges  $e = \{u, v\}$  are reset to zero if  $u$  or  $v$  were contained in component  $C$  at the time of

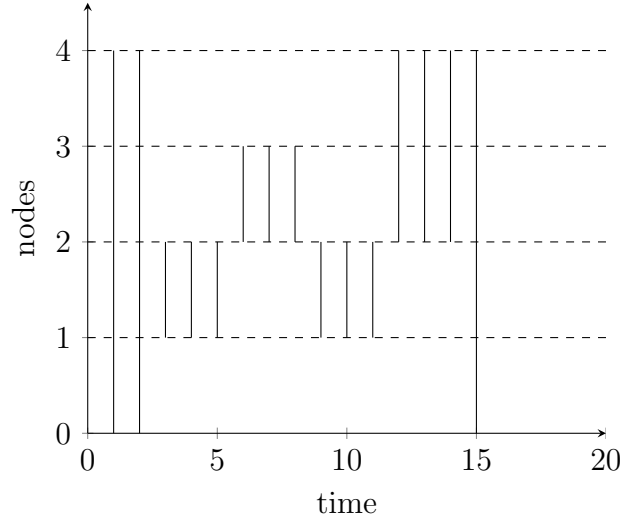


Figure 1: illustration of the analysis problem with the approach CREP-CORE

---

**Algorithm 4** delete( $Y$ ) of CREP-ADJ

---

```

for all  $e = \{u, v\} \in E$  do
  if  $u \in Y$  or  $v \in Y$  then
     $w(e) \leftarrow 0$ 
  end if
end for
for all  $C \in Y$  do
   $cap(serv(C)) \leftarrow cap(serv(C)) + res(C)$ 
   $res(C) \leftarrow 0$ 
end for

```

---

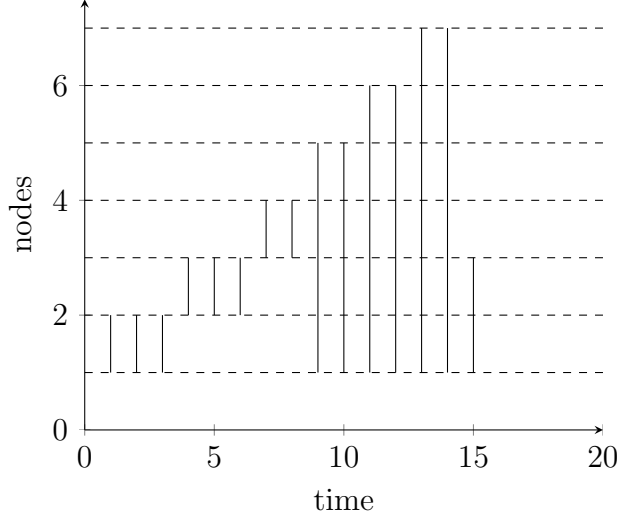


Figure 2: illustration for the CREP-ADJ approach

its deletion. This second version is very similar to the algorithm proposed by Avin et al. in [5]. The deletion method is also described in pseudocode in Algorithm 4.

The idea for the analysis is once again to relate the cost of both OPT and CREP-ADJ to the deleted components in the solution of CREP-ADJ.

The fact that CREP-ADJ also resets adjacent edges means that we can uniquely identify requests with the deleted component whose deletion led to the reset of the corresponding edge weights to zero. This means that we do not have to face the problems we discussed in the previous section on CREP-CORE. However the question remains whether the deletion of adjacent edges is valid, i.e. whether we can still preserve a good competitive ratio for this approach.

Figure 2 illustrates an example of this challenge mentioned above for the CREP-ADJ approach. In the illustration it is assumed that  $\alpha = 3$  and  $k = 3$ . In this case the input sequence leads CREP-ADJ to first merge nodes 1 and 2 at time  $t = 3$  and then to add node 3 to the resulting component at time  $t = 6$ . The next two requests do not quite lead to the merge of the node 4 with the component. Instead a series of requests follows where node 1 communicates with other nodes that are outside of its component without any merges. Note that this sequence can be extended until node 1 has communicated with every node except from the nodes 1, 2, 3, 4. Finally the first 4 nodes are merged at time  $t = 15$  at which point the resulting component is deleted as well as *all* edges in the sequence. This means that the amount of adjacent edges that are reset may greatly exceed the amount of edges inside a deleted component.

We show in our analysis in Section 6 that this approach is valid and allows for a competitive ratio of  $O(2/\epsilon \cdot k \log k)$ .

## 6 Competitive Analysis of CREP-ADJ

For the sake of this analysis we examine the CREP-ADJ algorithm from Section 5 and name the algorithm CREP from now on.



We analyse the competitive ratio of CREP with augmentation  $(2 + \epsilon)$  and show in Theorem 15 that CREP is  $O(2/\epsilon \cdot k \log k)$ -competitive.

In Theorem 17 in Section 7.4 we show that this algorithm can be implemented in polynomial time.

## 6.1 Algorithm Definitions

We begin our analysis by introducing two general definitions that we will use throughout the analysis.

**Definition 1.** Define for any subset  $S$  of components  $w(S)$  as the total weight of all edges between nodes of  $S$ .

**Definition 2.** Let a set of components of size at least 2 and of connectivity  $\alpha$  be a *mergeable* component set.

## 6.2 Structural Properties

These properties are adapted from previous work ([2, 5]) to use the connectivity-based approach which generally simplifies them, but guarantees slightly less minimum edge weight within mergeable component sets.

**Definition 3.** An  $\alpha$ -connected component is a maximal set of vertices that is  $\alpha$ -connected.

The following lemma states that there is always at most one mergeable component set after the insertion of a new edge which CREP then merges. The earliest point in time a new mergeable component set can emerge is after the next edge is inserted. This lemma is very similar to Lemma 4.1 from [2] in essence. The only difference for CREP lies in the definition of a mergeable component set which we have defined via a connectivity measure while the authors of [2] used density.

**Lemma 4.** At any time  $t$  after CREP performed its merge and delete actions, all subsets  $S$  of components with  $|S| > 1$  have connectivity less than  $\alpha$ , i.e. there exist no mergeable component sets after CREP performed its merges and deletions.

*Proof.* We prove the lemma by an induction on steps. The lemma holds trivially at time 0.

Now assume that at some time  $t > 0$  the lemma does not hold, i.e. there is a subset  $S$  of components with connectivity at least  $\alpha$  and  $|S| > 1$ . We may assume that  $t$  is the earliest time for which  $S$  has connectivity  $\alpha$ .

Then the incrementation of the weight of some edge  $e$  at time  $t$  raised the connectivity of  $S$ , but  $S$  was not merged into a new  $\alpha$ -connected component  $C$ . If no new component was created at time  $t$  we arrive at a contradiction as CREP always merges if there exists a mergeable component set.

Now assume that a component  $C$  was created at time  $t$ . This means that  $C$  must also contain the endpoints of  $e$ . But then the conjunction of  $C$  and  $S$  forms an even larger subset of components with connectivity at least  $\alpha$  which is a contradiction to

the maximality of  $C$  and  $S$ . □

The following lemma is adapted for our connectivity-based approach from Corollary 4.2 in [2].

**Lemma 5.** Fix any time  $t$  and consider weights right after they were updated by CREP but before any merge or delete actions. Then all subsets  $S$  of components with  $|S| > 1$  have connectivity at most  $\alpha$  and a mergeable component set  $S$  has connectivity exactly  $\alpha$ .

*Proof.* This lemma follows directly from Lemma 4 as connectivities can only increase by at most 1 at each time  $t$  and Lemma 4 guarantees that mergeable component sets are merged by CREP directly after they emerge before a new request is revealed. □

The following two lemmas combined give us a result similar to Lemma 4.3 in [2]: bounds on the edge weight that is cut when partitioning a mergeable component set, i.e. a set of components of connectivity at least  $\alpha$ .

We start by establishing a lower bound on this edge weight in the following lemma.

**Lemma 6.** Given a mergeable set of components  $S$  and a partition of  $S$  into  $g > 1$  parts  $S_1, \dots, S_g$ . Then the weight between the parts of the partition is at least  $g/2 \cdot \alpha$ .

*Proof.* We construct a graph  $G$  with the different parts  $S_i, i \in \{1, \dots, g\}$  of the partition as nodes. Note that this graph is  $\alpha$ -connected. We insert an edge for each edge between the parts  $S_i$ . Now consider the sum of the weighted degrees of all such nodes  $S_i$  in the constructed graph:

$$\sum_{i \in \{1, \dots, g\}} \deg_G(S_i) = 2 \sum_{e \in G} w(e)$$

The equality follows as the left sum counts each edge twice, once for each endpoint. Now consider the fact that each node  $S_i$  must have degree at least  $\alpha$  with respect to the edges in  $G$  because  $G$  is  $\alpha$ -connected. Hence

$$2 \sum_{e \in G} w(e) = \sum_{i \in \{1, \dots, g\}} \deg_G(S_i) \leq \sum_{i \in \{1, \dots, g\}} \alpha = g \cdot \alpha$$

which gives us that  $\sum_{e \in G} w(e) \geq g/2 \cdot \alpha$ . □

In the following lemma we establish the upper bound on the cut edge weight when partitioning a mergeable set of components  $S$  into  $g \geq 2$  parts.

**Lemma 7.** Given a mergeable set of components  $S$  and a partitioning of  $S$  into  $g \geq 2$  parts  $S_1, \dots, S_g$ . The weight between the parts  $S_i$  is at most  $(g - 1) \cdot \alpha$  during the execution of CREP (in the version CREP-ADJ).

*Proof.* Similarly to before we construct a graph  $G = (V, E)$  with the different parts  $S_i, i \in \{1, \dots, g\}$  of the partition as nodes and we insert an edge for each edge between the parts  $S_i$ . Note again that this graph is  $\alpha$ -connected. We iteratively

partition  $G$  into subsets via minimum cuts with regard to edge weight, i.e. we consider a minimum edge cut of  $G$  which partitions the nodes of  $G$  into the subsets  $V_1$  and  $V_2$ . We continue to iteratively partition the resulting sets until all sets contain only one node of  $G$  each. As this required at most  $|V| - 1$  cuts of value at most  $\alpha$  and  $|V| = g$  by definition of  $G$  the lemma follows.  $\square$

### 6.3 Proof Overview

Now that we have defined some general notions and fundamental properties we give an overview over the general approach we use in order to achieve an upper bound on CREP-ADJ, a lower bound on OPT and finally the competitive ratio of CREP-ADJ.

As we have argued above, the definition of CREP-ADJ allows us to argue about each deleted component and the requests deleted during its deletion separately.

Furthermore we know that a component deleted by CREP-ADJ is both at least  $\alpha$ -connected as well as larger than  $k$ , i.e. its nodes must be spread across several servers in the solution of OPT. This allows us to assign significant cost to OPT.

Additionally we also know by the definition of CREP that the edges that are reset because they are adjacent to  $C$  can not be part of any component. We use this fact in Lemma 12 in order to show that the amount of such requests is at most  $k \cdot \alpha$  for each migration performed by OPT.

These are the main concepts we use in order to bound the cost of CREP-ADJ and to assign cost to OPT in the following sections. In Theorem 15 we finally show the competitive ratio of CREP-ADJ.

### 6.4 Upper Bound On CREP-ADJ

We start the analysis of the upper bound of CREP by introducing several notions that we will use throughout the analysis.

We define the set  $\text{DEL}(\sigma)$  as the set of components that were deleted by CREP during its execution given the input sequence  $\sigma$ .

We define the following notions for a deleted component  $C \in \text{DEL}(\sigma)$ .

Let  $\text{EPOCH}(C)$  denote the (node, time) pairs of nodes in  $C$  starting at the time after the time  $\tau(\text{node})$  when  $\text{node}$  was last turned into a singleton component, i.e.

$$\text{EPOCH}(C) = \bigcup_{n \in \text{nodes}(C)} \{n\} \times \{\tau(n) + 1, \dots, \tau(C)\}.$$

This definition of an epoch is very similar to the one from [5] and [2]. Note that for  $C \in \text{DEL}(\sigma)$ ,  $\tau(C)$  denotes both the time of the creation as well as the time of deletion of  $C$ . We can use this definition of a component epoch  $\text{EPOCH}(C)$  to uniquely assign each node to a deleted component  $C$  at each point in time  $t$  (except for nodes in components that persist until the end of sequence  $\sigma$ ).

We assign all requests to  $\text{EPOCH}(C)$  whose corresponding requests are deleted because of the deletion of component  $C$  and call the set of those requests  $\text{REQ}(C)$ . We split the requests from  $\text{REQ}(C)$  into two sets:  $\text{CORE}(C)$  and  $\text{HALO}(C)$ .  $\text{CORE}(C)$  contains all requests for which both nodes have already been assigned to  $C$  at the

time of the request, i.e.

$$\text{CORE}(C) = \{r = \{u, v\} \in \sigma \mid (u, \text{TIME}(r)) \in \text{EPOCH}(C) \text{ and } (v, \text{TIME}(r)) \in \text{EPOCH}(C)\}.$$

These are the requests that led to the creation of component  $C$  by increasing the connectivity within the corresponding subgraph.

We define  $\text{HALO}(C)$  as the set of all requests from  $\text{REQ}(C)$  for which exactly one end point was associated with  $C$  at the time of the request. Note that this means that  $\text{HALO}(C) = \text{REQ}(C) \setminus \text{CORE}(C)$ .

These definitions allow us to differentiate between the highly-connected sub-graph induced by the nodes of  $C$  which are connected by requests from  $\text{CORE}(C)$  and the edges leaving  $C$  from  $\text{HALO}(C)$  which are relatively less dense as CREP has not merged any outer node with the component.

We start the analysis by bounding the communication cost of CREP that is due to serving requests from  $\text{CORE}(C)$  for  $C \in \text{DEL}(\sigma)$ .

**Lemma 8.** With augmentation  $2 + \epsilon$ , CREP pays at most communication cost  $|C| \cdot \alpha$  for requests in  $\text{CORE}(C)$  where  $C \in \text{DEL}(\sigma)$ .

*Proof.* First note that due to Lemma 4 CREP merges mergeable component sets as soon as they emerge. Whenever CREP performs a merge of a mergeable component set  $S$ , Lemma 7 states that there was at most total edge weight  $(|S| - 1) \cdot \alpha$  between the merged components, i.e.  $w(S) \leq (|S| - 1) \cdot \alpha$ . Each such merge decreases the number of components that need to be merged in order to form component  $C$  by  $|S| - 1$ . Hence CREP has paid at most  $|C| \cdot \alpha$  communication cost for requests in  $\text{CORE}(C)$ .  $\square$

We define  $\text{FIN-WEIGHTS}(\sigma)$  as the total amount of edge weight between the components  $\text{FIN-COMPS}(\sigma)$  which are present after the execution of CREP given input sequence  $\sigma$ .

Together with the fact that CREP pays for all requests in  $\text{HALO}(C)$  for deleted components  $C$  we use these definitions as well as the previous lemma to bound the total communication cost of CREP in the following lemma.

**Lemma 9.** The cost of serving communication requests that CREP has to pay, denoted by  $\text{CREP}^{\text{req}}(\sigma)$  given input sequence  $\sigma$  is bounded by

$$\text{CREP}^{\text{req}}(\sigma) \leq \sum_{C \in \text{DEL}(\sigma)} (|C| \cdot \alpha + |\text{HALO}(C)|) + \sum_{C \in \text{FIN-COMPS}(\sigma)} |C| \cdot \alpha + \text{FIN-WEIGHTS}(\sigma).$$

*Proof.* The number of communication requests that led to the creation of a component  $C$  is bounded by  $|C| \cdot \alpha$  due to Lemma 7. If component  $C$  was deleted by CREP then also the edge weights corresponding to requests from  $\text{HALO}(C)$  were reset to zero. All other edge weights were not changed. The remaining communication requests that have not been accounted for so far have either led to the creation of component  $C \in \text{FIN-COMPS}(\sigma)$  and are hence also bounded by  $|C| \cdot \alpha$  or have not led CREP to any merge and are hence contained in  $\text{FIN-WEIGHTS}(\sigma)$ . This concludes the proof.  $\square$

We continue our analysis by bounding the migration cost of CREP in the following lemma.

**Lemma 10.** With augmentation  $2 + \epsilon$ , CREP pays at most migration costs of

$$\text{CREP}^{mig}(\sigma) \leq \sum_{C \in \text{DEL}(\sigma) \cup \text{FIN-COMPS}(\sigma)} |C| \cdot ((2/\epsilon + 1) + \log k) \cdot \alpha.$$

*Proof.* First note that CREP only performs migrations when it merges components. We fix a component  $C \in \text{DEL}(\sigma) \cup \text{FIN-COMPS}(\sigma)$  and bound the number of times each node of  $C$  is moved as CREP processes the requests that led to the creation of  $C$ .

As CREP only reserves additional space  $\lfloor \epsilon \cdot |B| \rfloor$  for each component  $B$  and only moves component  $B$  when a merge results in a component of size more than  $(1 + \epsilon) \cdot |B|$  each node of  $C$  is moved at most  $(2/\epsilon + 1) + \log k$  times. Summing over all nodes in  $C$  that were actually moved by CREP bounds the number of migrations by  $|C| \cdot ((2/\epsilon + 1) + \log k)$  as components get deleted without migrations once they contain more than  $k$  nodes. This leads to the desired bound on the migration costs as each node migration incurs cost  $\alpha$  to CREP.  $\square$

Finally we summarize our results from Lemma 9 and Lemma 10 in the following lemma in order to obtain the final upper bound on the cost of CREP.

**Lemma 11.** With augmentation  $2 + \epsilon$ , CREP pays at most total cost

$$2 \cdot \sum_{C \in \text{COMPS}(\sigma)} |C| \cdot ((2/\epsilon + 1) + \log k) \cdot \alpha + \sum_{C \in \text{DEL}(\sigma)} |\text{HALO}(C)| + \text{FIN-WEIGHTS}(\sigma).$$

where  $\text{COMPS}(\sigma) = \text{DEL}(\sigma) \cup \text{FIN-COMPS}(\sigma)$ .

*Proof.* We sum the results from Lemma 9 and Lemma 10 to obtain the lemma:

$$\begin{aligned} \text{CREP}(\sigma) &\leq \text{CREP}^{req} + \text{CREP}^{mig} \\ &\leq \sum_{C \in \text{DEL}(\sigma)} (|C| \cdot \alpha + |\text{HALO}(C)|) + \sum_{C \in \text{FIN-COMPS}(\sigma)} |C| \cdot \alpha + \text{FIN-WEIGHTS}(\sigma) \\ &\quad + \sum_{C \in \text{COMPS}(\sigma)} |C| \cdot ((2/\epsilon + 1) + \log k) \cdot \alpha \\ &\leq 2 \cdot \sum_{C \in \text{COMPS}(\sigma)} |C| \cdot ((2/\epsilon + 1) + \log k) \cdot \alpha + \sum_{C \in \text{DEL}(\sigma)} |\text{HALO}(C)| \\ &\quad + \text{FIN-WEIGHTS}(\sigma). \end{aligned}$$

$\square$

## 6.5 Lower Bound on $\text{OPT}$

In this section we bound the cost on  $\text{OPT}$  by assigning cost to  $\text{OPT}$  based on the size of the components  $C$  that CREP deletes and the associated adjacent edges  $\text{HALO}(C)$

which CREP resets to zero during the deletion of  $C$ . In order to achieve this we introduce some additional notions.

First we define the term *offline interval* of a node  $v$  to be the time between two migrations of  $v$  in the solution of OPT. More specifically an offline interval of node  $v$  either starts at time zero (if it is the first offline interval of  $v$ ) or after a migration of  $v$  and ends with the next migration of node  $v$  that OPT performs.

Furthermore we say that an offline interval is contained in the epoch  $\text{EPOCH}(C)$  of a component  $C \in \text{DEL}(\sigma)$  if it ends before the time  $\tau(C)$ . Note that  $\tau(C)$  is both the time of the creation of  $C$  in the solution of CREP and the time of its deletion as  $C \in \text{DEL}(\sigma)$ .

We assign a request  $r$  involving the node  $v$  to an offline interval of  $v$  if it is both the first offline interval of one of the end points of  $r$  that ends and if the offline interval ends before the deletion of the edge representing  $r$  due to a component deletion.

The requests from  $\mathcal{H} = \bigcup_{C \in \text{DEL}(\sigma)} \text{HALO}(C)$  that are not assigned to any offline interval are then those which are deleted due to the deletion of a component that took place before the corresponding offline interval ended.

Let  $P$  denote the set of edges from  $\bigcup_{C \in \text{DEL}(\sigma)} \text{HALO}(C)$  that both CREP and OPT pay for and let  $I$  denote the set of requests we have assigned to offline intervals.

These definitions are illustrated in Figure 3. Note that we only show some requests explicitly for the sake of readability. The grey horizontal lines represent the nodes at each time  $t$ . The red outline surrounds the (node,time) pairs of  $\text{EPOCH}(C)$ . Blue dots mark migrations of the corresponding node performed by OPT while red dots mark deletions of the component the node was assigned to at that time. The dashed vertical lines in black mark requests that are assigned to another component because it is deleted before component  $C$ . The dotted green line is a request from  $\text{HALO}(C)$  assigned to the offline interval of node 5 between the two blue dots. The regular green lines are assigned to an offline interval which is not contained in  $\text{EPOCH}(C)$ . We define this concept more formally at a later point in the analysis. The lines in magenta are sample requests from  $\text{CORE}(C)$ .

We start by bounding the total edge weight (the total number of requests) we assign to any one offline interval when limiting ourselves to requests from  $\mathcal{H}$  which CREP pays for but OPT does not. We denote the set of these requests by  $N$ , i.e.  $N = \mathcal{H} \setminus P$ . Note that  $\mathcal{H}$  only contains requests which CREP paid for due to the definition of  $\text{HALO}(C)$ .

**Lemma 12.** We assign at most  $k \cdot \alpha$  requests from  $N$  to any one offline interval.

*Proof.* We fix an arbitrary offline interval of node  $v$ . Observe that none of the nodes involved in the assigned requests are moved by OPT during the offline interval, hence all the requests in question involve only nodes that OPT has placed on the same server as  $v$  during the offline interval.

The number of such nodes is hence limited by the server capacity  $k$ . As we only examine requests from  $\mathcal{H}$  we know that none of these requests have led CREP to perform any merges, hence there were at most  $\alpha$  requests between  $v$  and any one of

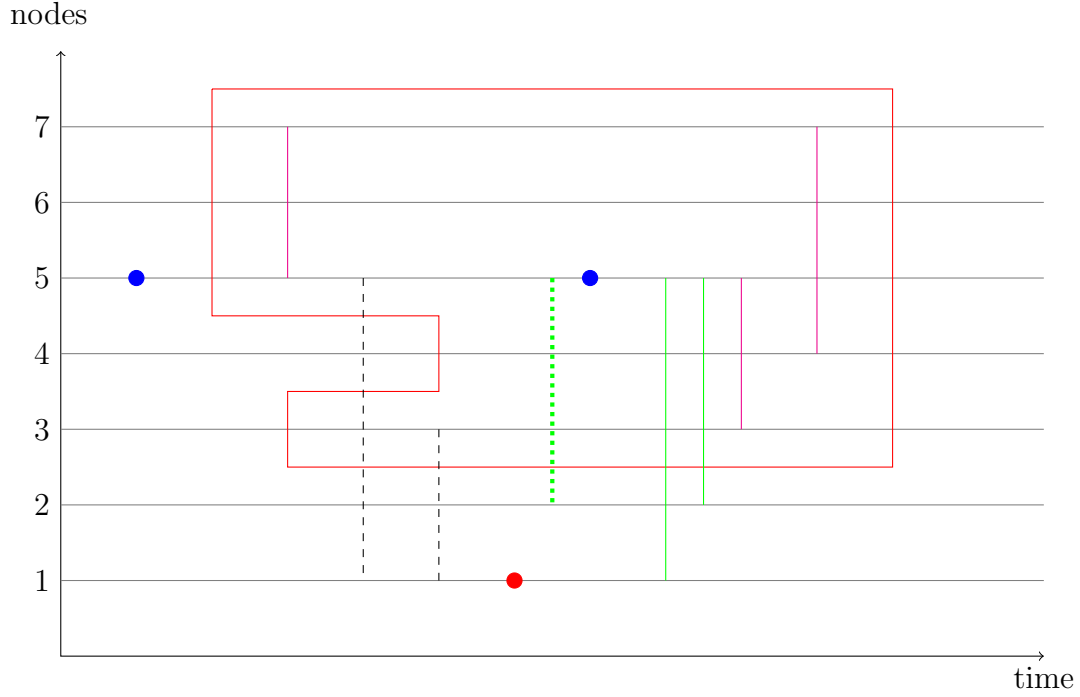


Figure 3: illustration of definitions used in the analysis

the other nodes on its server during the offline interval. This bounds the number of requests assigned to the offline interval by  $k \cdot \alpha$ .  $\square$

Let  $R(C)$  denote the set of requests from  $\text{HALO}(C)$  that were not assigned to any offline interval for a deleted component  $C \in \text{DEL}(\sigma)$ .

We say that a migration of node  $v$  at time  $t$  in the solution of  $\text{OPT}$  is *contained* in  $\text{EPOCH}(C)$  if  $(v, t) \in \text{EPOCH}(C)$ .

Let  $\text{OPT-MIG}(C)$  denote the cost of  $\text{OPT}$  due to migrations of nodes from component  $C$  that are contained in  $\text{EPOCH}(C)$  and let  $\text{OPT-REQ}(C)$  denote the cost of  $\text{OPT}$  due to serving requests from  $\text{CORE}(C)$ .

We show the following lower bound on the cost of  $\text{OPT}$  for migrations from  $\text{OPT-MIG}(C)$  and requests from  $\text{OPT-REQ}(C)$  for all deleted components  $C$ .

**Lemma 13.**

$$\sum_{C \in \text{DEL}(\sigma)} \text{OPT-MIG}(C) + \text{OPT-REQ}(C) \geq |C|/k \cdot \alpha + |R(C)|/k$$

*Proof.* For the following part of the proof we fix an arbitrary component  $C \in \text{DEL}(\sigma)$ . Note that the nodes involved in requests from  $R(C)$  were not moved by  $\text{OPT}$  during the processing of requests from  $R(C)$  until the time of deletion of  $C$  as otherwise they would be assigned to an offline interval.

The number of nodes contained in  $C$  or connected to  $C$  via edges representing requests from  $R$  is at least  $|C| + |R(C)|/\alpha$  since requests from  $R(C)$  have not led  $\text{CREP}$  to perform any migrations. Because of this fact  $\text{OPT}$  must have placed those nodes on at least  $\frac{|C| + |R(C)|/\alpha}{k}$  different servers. As  $\text{OPT}$  does not pay for any requests

from  $R$  it follows that OPT must have placed the nodes from  $C$  in  $\frac{|C|+|R(C)|/\alpha}{k}$  different servers.

We first examine the case in which OPT does not move any nodes from  $C$  during  $\text{EPOCH}(C)$ . In this case OPT must partition a graph containing the nodes from  $C$  which are connected via edges representing the requests from  $\text{CORE}(C)$ . As stated earlier OPT placed those nodes in  $\frac{|C|+|R(C)|/\alpha}{k}$  different servers at time  $\tau(C)$ . As CREP merged component  $C$  this graph is  $\alpha$ -connected and hence Lemma 6 gives that OPT has to cut at least edges of total weight  $\frac{|C|+|R(C)|/\alpha}{k} \cdot \alpha = |C|/k \cdot \alpha + |R(C)|/k$ .

For the more general case in which OPT may perform node migrations during  $\text{EPOCH}(C)$  we adapt the graph construction from above as follows: we add a vertex representing each (node, time) pair from  $\text{EPOCH}(C)$ . We connect each (node, time) pair  $p$  with edges of weight  $\alpha$  to the pairs of the same node that represent the time step directly before and directly after  $p$  (if they exist in the graph). These edges represent the fact that OPT may choose to migrate a node between any two time steps in  $\text{EPOCH}(C)$ . Additionally we add an edge of weight one for each request  $r = \{u, v\}$  from  $\text{CORE}(C)$  by connecting the nodes in the graph that represent the pairs  $(u, t)$  and  $(v, t)$ , respectively. OPT once again has to partition this graph into  $\frac{|C|+|R(C)|/\alpha}{k}$  parts.

Note that we only added edges of weight  $\alpha$  to the graph and hence this graph is also  $\alpha$ -connected. We conclude that once again OPT has to cut edges of weight at least  $\frac{|C|+|R(C)|/\alpha}{k} \cdot \alpha = |C|/k \cdot \alpha + |R(C)|/k$ .

In both cases only edges representing either requests from  $\text{OPT-REQ}(C)$  or migrations from  $\text{OPT-MIG}(C)$  were cut.

As the sets  $\text{core}(C)$ ,  $R(C)$ ,  $\text{core}(D)$  and  $R(D)$  are disjoint for two different components  $C, D \in \text{DEL}(\sigma)$  per their definition we conclude that

$$\sum_{C \in \text{DEL}(\sigma)} \text{OPT-MIG}(C) + \text{OPT-REQ}(C) \geq |C|/k \cdot \alpha + |R(C)|/k.$$

□

In the following lemma we combine the results of the previous lemmas in order to bound the cost of OPT given input sequence  $\sigma$ , denoted by  $\text{OPT}(\sigma)$ .

**Lemma 14.** The cost of the solution of OPT given input sequence  $\sigma$  is bounded by

$$\text{OPT}(\sigma) \geq 1/2 \cdot \sum_{C \in \text{DEL}(\sigma)} |C|/k \cdot \alpha + |\text{HALO}(C)|/k.$$

*Proof.* We combine the results from Lemma 12 and Lemma 13. Note that the cost from Lemma 13 may contain migration costs. In this case the corresponding migrations represent the end of an offline interval. We denote the number of offline intervals by  $o$ . This gives us that

$$2\text{OPT}(\sigma) \geq \sum_{C \in \text{DEL}(\sigma)} \text{OPT-MIG}(C) + \text{OPT-REQ}(C) + o \cdot \alpha + |P|$$



as we account for each migration at most twice.

Consider that due to Lemma 12 we have the inequality  $o \geq |N|/k$ . We repeat that  $\mathcal{H} = \bigcup_{C \in \text{DEL}(\sigma)} \text{HALO}(C)$ . Note that  $N$  is the subset of requests of  $\mathcal{H}$  for which  $\text{OPT}$  does not pay while  $P$  is the subset of  $\mathcal{H}$   $\text{OPT}$  pays for. It follows that the disjoint union of  $N$  and  $P$  is  $\mathcal{H}$ . Hence we obtain

$$\begin{aligned} 2\text{OPT}(\sigma) &\geq \sum_{C \in \text{DEL}(\sigma)} \text{OPT-MIG}(C) + \text{OPT-REQ}(C) + o \cdot \alpha + |P| \\ &\geq \sum_{C \in \text{DEL}(\sigma)} |C|/k \cdot \alpha + |R(C)|/k + (|N| + |P|)/k \\ &\geq \sum_{C \in \text{DEL}(\sigma)} |C|/k \cdot \alpha + |\text{HALO}(C)|/k. \end{aligned}$$

This gives us the lemma.  $\square$

## 6.6 Competitive Ratio

In this section we combine the results of Lemma 11 and Lemma 14 to obtain the following theorem which gives us the desired competitive ratio.

**Theorem 15.** With augmentation  $(2 + \epsilon)$  the competitive ratio of  $\text{CREP}$  is in  $O(2/\epsilon \cdot k \log k)$ .

*Proof.* We arbitrarily fix an input sequence  $\sigma$  and use our previous results to bound the competitive ratio of  $\text{CREP}$ . We define  $\text{COMPS}(\sigma) := \text{DEL}(\sigma) \cup \text{FIN-COMPS}(\sigma)$  in order to improve readability. Let  $P$  denote the set of edges from  $\bigcup_{C \in \text{DEL}(\sigma)} \text{HALO}(C)$  that both  $\text{CREP}$  and  $\text{OPT}$  pay for.

$$\begin{aligned} &\frac{\text{CREP}(\sigma) - \text{FIN-WEIGHTS}(\sigma)}{\text{OPT}(\sigma)} \\ &\leq \frac{2 \cdot \sum_{C \in \text{COMPS}(\sigma)} |C| \cdot ((2/\epsilon + 1) + \log k) \cdot \alpha + \sum_{C \in \text{DEL}(\sigma)} |\text{HALO}(C)|}{1/2 \cdot \sum_{C \in \text{DEL}(\sigma)} |C|/k \cdot \alpha + |\text{HALO}(C)|/k + |P|} \\ &\leq k \log k \frac{2 \cdot \sum_{C \in \text{DEL}(\sigma)} |C| \cdot (2/\epsilon + 1) \cdot \alpha + \sum_{C \in \text{DEL}(\sigma)} |\text{HALO}(C)|}{1/2 \sum_{C \in \text{DEL}(\sigma)} |C| \cdot \alpha/2 + |\text{HALO}(C)|} + \beta \\ &= O(2/\epsilon \cdot k \log k) + \beta \end{aligned}$$

where

$$\beta = \sum_{C \in \text{FIN-COMPS}(\sigma)} |C| \cdot ((2/\epsilon + 1) + \log k) \cdot \alpha$$

Let  $\beta' = \beta + \text{FIN-WEIGHTS}(\sigma)$ . Then it follows that

$$\frac{\text{CREP}(\sigma)}{\text{OPT}(\sigma)} \leq O(2/\epsilon \cdot k \log k) + \beta'.$$

To obtain the bound on  $\beta'$  we observe that the components in  $\text{FIN-COMPS}(\sigma)$  each are of size at most  $k$  since they were not deleted by CREP. This allows us to derive to bound  $\sum_{C \in \text{FIN-COMPS}(\sigma)} |C| \cdot ((2/\epsilon + 1) + \log k) \leq l \cdot k \cdot ((2/\epsilon + 1) + \log k)$ . Since at the end of the execution of CREP there can be at most  $k \cdot l$  components, Lemma 7 allows us to bound  $\text{FIN-WEIGHTS}(\sigma)$  by  $k \cdot l \cdot \alpha$ . Hence we conclude that  $\beta' \leq l \cdot k \cdot ((2/\epsilon + 1) + \log k) \cdot \alpha + k \cdot l \cdot \alpha \in O(2/\epsilon \cdot k \log k)$ .  $\square$

So far we have shown that CREP is competitive. In Section 7.4 we will show that our version CREP-ADJ of CREP can be implemented in polynomial time.

## 7 From Algorithmic Ideas to Polynomial Time Implementation

So far we have only shown that the algorithm CREP-ADJ described in Section 5 is competitive with competitive ratio  $O(2/\epsilon \cdot k \log k)$ . In this section we discuss the implementation and running time and conclude in Theorem 17 in Section 7.4 that it is indeed polynomial.

We first describe the ideas behind our implementation of CREP and then illustrate them by discussing pseudocode of the implemented algorithm.

### 7.1 Algorithm Explanations

In this section we describe our implementation of CREP with augmentation  $2 + \epsilon$  in greater detail.

In order to limit the section of the graph  $G$  maintained by CREP that needs to be updated upon a new request between nodes of different components we maintain a decomposition tree defined as follows: the root represents the whole graph and is assigned the connectivity of the whole graph. Given a node  $v$  in the tree that represents a subgraph  $G'$  of  $G$ , we decompose  $G'$  into subgraphs whose connectivity is strictly larger than that of  $G'$  and add children to  $v$  for each such subgraph. We do not decompose sub-graphs of connectivity at least  $\alpha$  any further as we only need to identify whether a new subgraph of connectivity at least  $\alpha$  was created by the insertion of the most recent request. Additionally we keep track of the connectivity of each such subgraph.

Figure 5 illustrates this decomposition for the graph shown in Figure 4. In the decomposition tree we have labelled each node with the corresponding subset of vertices and the connectivity of the graph induced by these vertices.

If a new request is revealed to CREP then we only need to update the smallest subtree of the decomposition tree which still contains both end points of the request. For example in the case of a new request between 0 and 1 we only need to recompute the decomposition of the sub-graph induced by the vertex set  $\{0, 1, 2\}$  in the example.

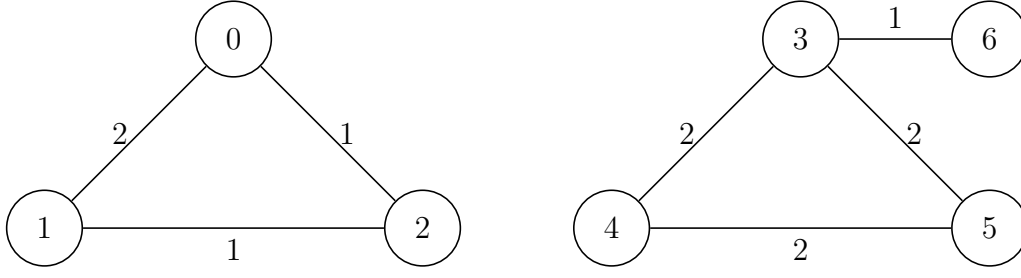


Figure 4: example graph

This is correct because we can view each decomposition of a subgraph  $G'$  into smaller graphs of a higher connectivity as a set of cuts that separates the nodes of  $G'$ . Inserting a new edge within a subgraph  $G'$  may only increase the value of the cuts which result in the decomposition of  $G'$ , but do not affect cuts separating  $G'$  itself from other subgraphs.

If a new request led to the creation of a new component this means that two old components that were at least  $\alpha$ -connected were merged and hence the number of leaves in the decomposition tree decreased. If this is the case then the algorithm checks whether the new component contains more than  $k$  nodes. In this case the component is deleted and split into singleton components, each containing one node from the deleted component.

Upon such a component deletion the edges inside of and adjacent to the component are deleted, i.e. their weight is reset to zero. This means that the decomposition tree needs to be recomputed in order to reflect this change.

If however the resulting component  $C$  contains at most  $k$  nodes the algorithm tries to collocate the nodes of the component while minimizing migration costs, i.e. looking for a cluster which contains as many nodes of the newly merged component as possible but which also has enough free capacity for the remaining nodes to be moved there and for additional reservation  $\min\{\lfloor \epsilon \cdot |C| \rfloor, k - |C|\}$ .

Both the decision whether to delete as well as possible node migrations are handled in the subroutine *updateMapping* which is given the newly computed  $\alpha$ -connected components as input.

In the next section we provide more detailed pseudocode describing our implementation.

## 7.2 Algorithm Pseudocode

Algorithm 5 is the main function that is called upon each new request. It checks whether the new request is between different  $\alpha$ -connected components. If this is not the case it determines that this request can not change the decomposition and returns.

Otherwise the weight of the corresponding edge is increased and other routines are called that update the decomposition based on this new edge.

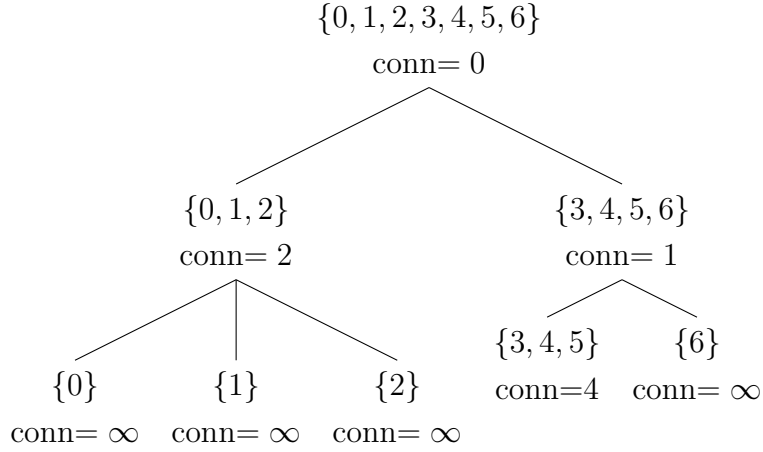


Figure 5: decomposition tree for the graph from Figure 4 for  $\alpha = 4$

---

**Algorithm 5** insertAndUpdate( $a, b$ )

---

```

if comp[ $a$ ] == comp[ $b$ ] then
    return
end if
addEdge( $a, b$ )
updateDecomposition( $a, b$ )
 $del \leftarrow$  updateMapping( $\alpha$ ConnectedComponents)
delComponents( $del$ )

```

---

Algorithm 6 first determines the smallest sub-graph in the decomposition tree that contains both end points of the request and decomposes this sub-graph. For this decomposition step we use the algorithm proposed by Chang et al. ([8]). We explain this algorithm in greater detail in Section 7.3.

Afterwards the routine *updateMapping* is called which compares the number of components to the number of components before the arrival of the request. Only if the number of components has decreased it checks for the new component as otherwise there was no component merge.

If there was a merge then the routine examines the size of the newly created (merged) component and decides whether to delete or to collocate based on the logic described in the previous section.

If a deletion has to be performed then this step is done in the routine *delComponents* (Algorithm 7) which resets all edge weights both of edges between nodes of the component as well as all adjacent edges and finally starts the decomposition of the whole graph in order to arrive at a new decomposition that follows the definition from the previous section.

In the case of a collocation the nodes are moved to a cluster that has enough space while updating the reservations and cluster capacities accordingly.

---

**Algorithm 6** updateDecomposition( $a, b$ )

---

```
 $q \leftarrow \text{findSmallestSubgraph}(a, b)$ 
while  $q$  not empty do
   $\text{current} \leftarrow q.\text{popFront}()$ 
  if  $\text{res.connectivity} == \alpha$  then
    continue
  end if
   $\text{res} \leftarrow \text{decompose}(\text{current}, \text{current.connectivity}+1)$  //decomposition based on
  ( $s, t$ )-cuts
   $\text{current.connectivity} \leftarrow$  value of smallest encountered cut
  if  $\text{current.connectivity} \geq \alpha$  then
    continue
  end if
   $\text{childrenQueue} \leftarrow \text{res}$ 
  //make sure that only subgraphs with higher connectivity are added as children

  while  $\text{childrenQueue}$  not empty do
     $c \leftarrow \text{childrenQueue.pop}()$ 
     $cRes \leftarrow \text{decompose}(c, \text{current.connectivity}+1)$ 
     $c.\text{connectivity} \leftarrow$  value of smallest encountered cut
    if decompose returned only one graph then
       $\text{current.children.add}(cRes)$ 
      if  $cRes$  has connectivity smaller than  $\alpha$  then
         $q.\text{push}(cRes)$ 
      end if
    else
       $\text{childrenQueue.add}(cRes)$ 
    end if
  end while
end while
```

---

---

**Algorithm 7** delComponents( $del$ )

---

```
 $\text{delAllEdges}(del)$ 
 $\text{root.connectivity} = 0$ 
 $\text{root.children} = \{\}$ 
 $\text{updateDecomposition}(0, 1)$ 
```

---

### 7.3 On the Decomposition of a Subgraph

In this section we describe our algorithm for the decomposition of a given subgraph represented by a node in the decomposition tree. Specifically we describe the general idea behind our implementation of the subroutine *decompose*(*treeNode*, *c*) which we have based on the algorithm described by Chang et al. ([8]). In this algorithm *c* denotes the connectivity that is the basis of this decomposition step.

Given a node *v* of the decomposition tree, first a *partition graph* is constructed which is a graph consisting of the nodes in the subgraph represented by *v* and the edges which are between the nodes of the subgraph. This partition graph also supports merges and cuts of its nodes. More specifically the partition graph is initialized as a graph  $P = (V, E)$  with  $V = \text{nodes}(v)$  and  $E = \{e = \{u, w\} \in E' | u \in \text{nodes}(v) \text{ and } w \in \text{nodes}(v)\}$  where  $E'$  represents the set of edges in the graph maintained by CREP. Additionally we maintain a mapping *M* which assigns each node from *V* a set of the nodes in the subgraph represented by the decomposition tree node *v*. Initially *M* assigns each node in *V* the subset containing only the node itself.

We now run a *maximum adjacency search* algorithm, sometimes also called *maximum cardinality search* algorithm, ([22]) in order to obtain an arbitrary minimum (*s*, *t*)-cut of the graph.

The maximum adjacency search algorithm is defined as follows: We start with an empty list *L* to which we add an arbitrary node of *P*. We then continually add the most tightly connected node from *V* to *L*, i.e. the node which is connected to the nodes in *L* via edges of the most total weight. Stoer and Wagner ([22]) have shown that the edges between the last two nodes *s* and *t* added to *L* form a minimum (*s*, *t*)-cut.

We use the value of this cut in order to decide whether to merge the nodes *s* and *t* or whether to separate them. If the cut has value less than *c* we separate the nodes, otherwise we perform a merge.

Here the separation of the nodes *s* and *t* means that we remove all edges in the cut from the edges *E* of the partition graph *P*.

In the case of a merge we combine the nodes *s* and *t* and merge the outgoing edges, i.e. we replace the set of nodes *V* of *P* by the set  $V' = V \setminus \{s, t\} \cup \{v'\}$ . The edges *E* of *P* are modified by removing all edges adjacent to *s* and *t* and adding an edge  $e' = \{v', u\}$  of weight  $w(\{s, u\}) + w(\{t, u\})$  where  $w(e)$  denotes the weight of edge *e* if it exists and is equal to zero otherwise. Furthermore we adjust the mapping *M* by setting  $M(v') = M(s) \cup M(t)$ .

We continually run this algorithm until *P* contains no edges, i.e. until  $E = \emptyset$ . The sets of nodes mapped to each of the nodes of *P* by *M* now represent candidate subgraphs for the decomposition. Note though that we have only cut and merged according to minimum (*s*, *t*)-cuts and not according to minimum cuts. This means that the specific sequence in which we have performed the cuts may influence the result, e.g. if we merged based on a minimum (*s*, *t*)-cut which is not a minimum cut.

This can be remedied by repeating the procedure on the resulting subgraphs until it returns a subgraph of only one node, i.e. until no separation step is performed during the decomposition. This is due to the fact that this procedure always cuts a

subgraph of connectivity less than  $c$  at least once, as Chang et al. have shown (see Cutability Property in [8]).

In order to speed up this computation we use the heap data structure proposed and analyzed by Chang et al. ([8]).

Thus we conclude that this procedure correctly decomposes a given subgraph as Chang et al. have also stated in Theorem 3.1 in [8].

In the following section we show that this decomposition can be implemented in polynomial time.

## 7.4 Running Time

In this section we discuss the polynomial running time of an update of the decomposition tree after a new request has been received in our implementation of CREP-ADJ.

The main bottleneck of the algorithm lies in the decomposition updates. The following lemma shows a polynomial running time for these updates. The other parts of the algorithm can be implemented in polynomial time already.

**Lemma 16.** The subroutine *decompose* which decomposes a subgraph can be implemented in polynomial time  $O(\alpha|V|^2|E|)$ .

*Proof.* The worst case is given when the whole tree has to be recomputed. We first discuss the time complexity of decomposing a single tree node  $v$ . Let the corresponding subgraph be denoted by  $G_v = (V_v, E_v)$ . Since each iteration of the subroutine *decompose* performs at least one cut as long as the connectivity of the given graph is smaller than the current threshold  $c$  we conclude that after at most  $|V_v|$  iterations of *decompose* a correct decomposition is found.

Each step of *decompose* can be performed in  $O(|V_v| \cdot |E_v|)$  as the maximum adjacency search algorithm finds an arbitrary minimum  $(s, t)$ -cut in time  $O(|E_v|)$  as shown in theorem 4.1 in [8] and as there are at most  $|V_v|$  minimum  $(s, t)$ -cuts computed for each invocation of *decompose*.

Hence the complexity of decomposing the subgraph represented by a tree node  $v$  is in  $O(|V_v|^2|E_v|)$ . Let  $C_v$  denote the time needed for the decomposition of the subgraph represented by decomposition tree node  $v$ .

We now sum this complexity over the nodes for each connectivity level of the decomposition tree. To this end let  $level(i)$  denote all nodes in the decomposition tree which are of connectivity exactly  $i$ .

$$\sum_{i=0}^{\alpha} \sum_{v \in level(i)} C_v \leq \sum_{i=0}^{\alpha} O(|V|^2|E|) \in O(\alpha|V|^2|E|).$$

We conclude our analysis of the time complexity by observing the polynomial-time complexity of  $O(\alpha|V|^2|E|)$ .  $\square$

We conclude our analysis of the running time by observing that the remaining subroutines can be implemented in polynomial time which results in the following theorem on the running time.

**Theorem 17.** The algorithm CREP-ADJ can be implemented in polynomial time.

*Proof.* The theorem follows from the Lemma 16 and the fact that the remaining subroutines for edge insertions and component deletions can hence also be implemented in polynomial time. Note in particular that there can be at most one component deletion during the update step after any given request.  $\square$

## 8 Evaluation

In this section we evaluate the quality of the results and the performance of our algorithm implementation by comparing it with several algorithms described in Section 8.1 on input data sets described in Section 8.2.

### 8.1 Reference Algorithms

We first present the algorithms we compare in order to evaluate our results.

We consider both algorithmic approaches based on maintaining a second-order partition of the nodes into components we discussed in Section 5, i.e. CREP-ADJ where on a component deletion all edges inside of and adjacent to the component are deleted and CREP-CORE which only deletes internal edges. Both algorithms start with randomly initialized mappings of nodes to servers.

We have shown the competitive ratio of CREP-ADJ in Theorem 15 and its polynomial running time in the previous section.

We compare these algorithms both with the static graph partitioning algorithm METIS\_PartGraphRecursive implemented in the METIS framework ([16, 17]) which we will refer to from now on as STATIC and the adaptive/dynamic algorithm ParMETIS\_V3\_AdaptiveRepart ([18, 21, 20]) available in the ParMetis framework, referred to as ADAPT. Both frameworks are known to produce very good results and to be very fast.

We run ADAPT in an online manner while maintaining a communication graph by inserting an edge or incrementing a corresponding edge weight for each request, i.e. ADAPT is called after each request is inserted into the graph.

For STATIC we first record the communication graph and give it as input to STATIC. We then compute the cost of STATIC due to migrations from the random initial mapping of nodes to servers and due to remote requests, i.e. the total weight of edges between nodes of the communication graph which STATIC mapped to different servers.

### 8.2 Input Data

As input data we use several HPC traces (Mocfe, NeckBone, MultiGrid), the nature of the data is described in more detail by Avin, Ghobadi, Griner and Schmid in [4]. For the sake of readability we use the abbreviations NB for NeckBone and MG for MultiGrid.

All data sets contain 1024 different communication nodes and are limited to the first 300 000 requests. The value of  $\alpha$  is set to 6 and the algorithm was tasked to



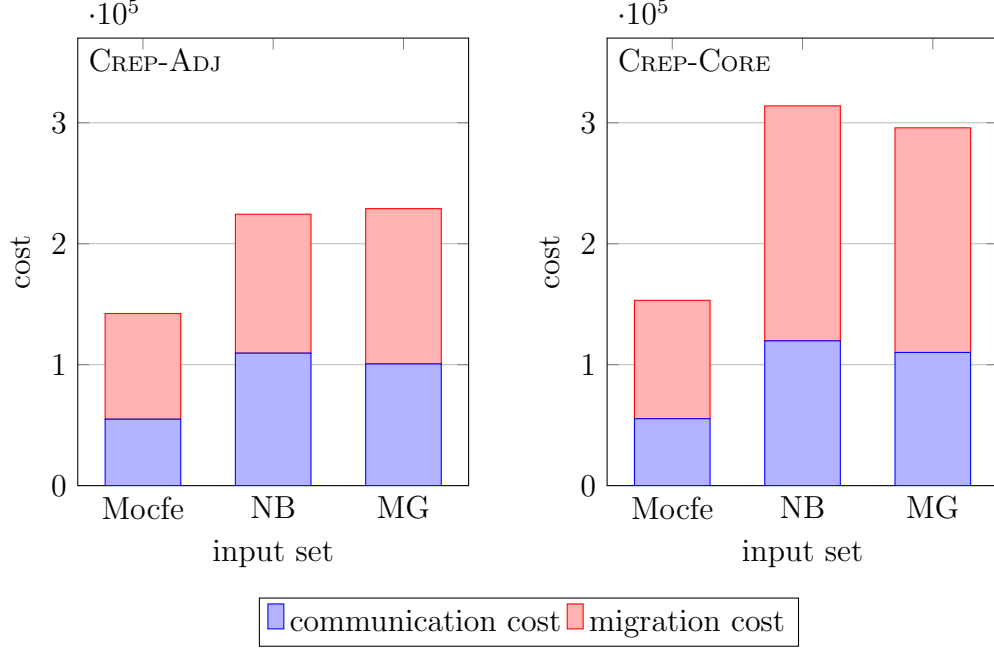


Figure 6: comparison of the total cost of CREP-ADJ (left) and CREP-CORE (right)

partition the nodes into 32 clusters of size 32 each, i.e.  $k = l = 32$ . The dynamic algorithms are allowed to use augmentation with a factor of 2.1, i.e. for the dynamic algorithms the maximum cluster capacities are  $\lfloor 32 \cdot 2.1 \rfloor = 67$ .

### 8.3 Comparison of CREP-ADJ and CREP-CORE

First we compare the results of CREP-ADJ and CREP-CORE. Figure 6 shows the resulting total cost of both algorithms on the different input sets as well as how this cost is split between communication and migration costs. One can see that CREP-ADJ always produces better results, especially for the second and third input set.

We now discuss the differences of the two algorithms in terms of communication and migration cost. The figures illustrate that CREP-ADJ pays only very slightly less for communication requests than CREP-CORE whereas the former pays significantly reduced migration cost. This suggests that the deletion of adjacent edges indeed improves the quality of the results by reducing the number of migrations.

Finally Figure 7 compares the running times of both algorithms. One can see that CREP-ADJ performs drastically better than CREP-CORE in this regard. This may be due to the fact that CREP-ADJ deletes edges more frequently and thus needs to take less edges into account when updating the decomposition tree.

Also note that both algorithms generally pay more for migrations than for communication. This suggests that there may be room for fine-tuning these algorithms in order to achieve a more balanced distribution of the cost. One such adjustment is investigated in the next section, namely whether we can improve the results by changing the algorithms in such a way that they only merge once a connectivity of  $2 \cdot \alpha$  is reached.

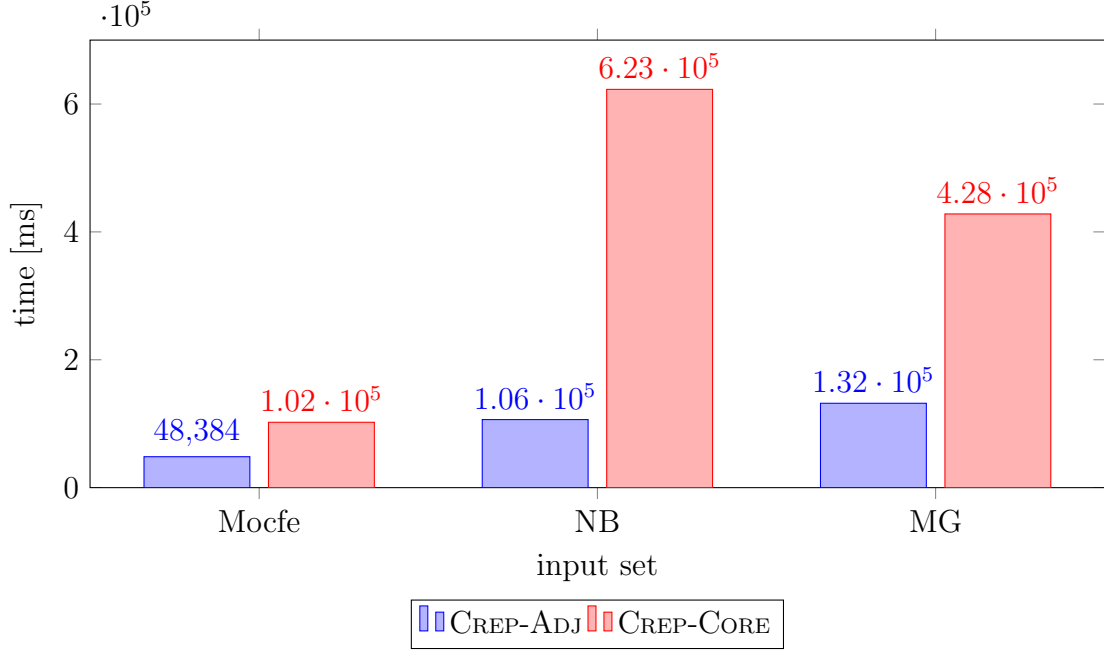


Figure 7: comparison of running time of CREP-ADJ and CREP-CORE

## 8.4 On the Influence of the Connectivity Threshold

In the previous section we observed that both CREP-CORE and CREP-ADJ have higher migration cost than communication cost. This leads to the question whether we can find adjustments that allow us to achieve a better balance of the different costs. In this section we investigate the influence of the connectivity threshold which determines when components are merged. Up until this point it was set to  $\alpha$ . As we show there lies potential in adjusting this threshold in order to improve the quality of the solution.

Namely we discuss the results of both algorithms for the case where they only merge components once a connectivity of at least  $2 \cdot \alpha$  is reached. Note that this only affects the analysis by constant factors, namely the statements from Lemma 6 and Lemma 7 are multiplied by a factor of two. This may only impact the other cost bounds by a factor of two which leaves our bound on the competitive ratio of CREP-ADJ unaffected. Also note that we do not change the cost of a node migration, this cost is still  $\alpha$ .

Figure 8 shows the total cost of CREP-ADJ and CREP-CORE for this scenario. One can see that the costs of both algorithms are reduced, but CREP-CORE has improved significantly more than CREP-ADJ.

When looking at the distribution of these costs to communication and migration costs one can see that now the communication costs of both algorithms are higher than the migration cost in all cases. For CREP-CORE this has improved the balance of both costs, but CREP-ADJ now has significantly higher communication costs than migration costs. This suggests that by adjusting the exact value of connectivity at which these algorithms perform their merges one might be able to further improve the quality of the results. It may even be possible to adjust this parameter dynamically based on some data gathered as requests are processed.

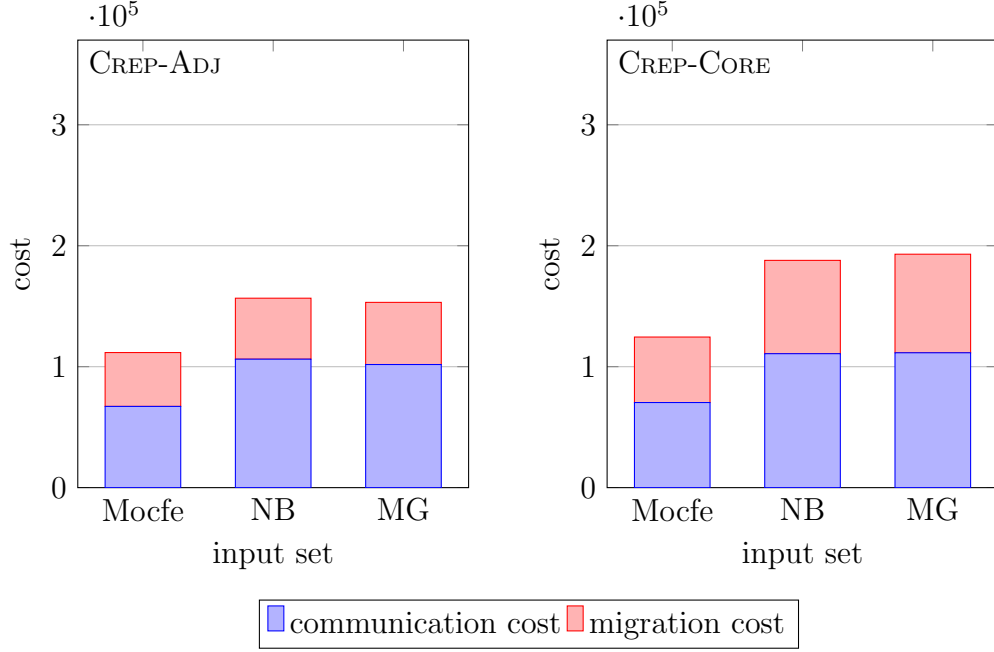


Figure 8: comparison of the total cost of CREP-ADJ (left) and CREP-CORE (right) in the case where merges are performed at connectivity  $2 \cdot \alpha$

These improvements come at a cost of running time as Figure 9 shows, almost doubling or tripling the times of the previous results in some cases. Especially the already worse run times of CREP-CORE are drastically increased.

## 8.5 Results of ADAPT and STATIC

Figure 10 shows the total cost of the solutions of ADAPT and STATIC. One can see that STATIC is able to produce results that far surpass ADAPT. In fact the results of STATIC are also a significant improvement over the results of CREP-ADJ and CREP-CORE presented before. But it is also important to note that STATIC is only useful in scenarios where the communication patterns and frequencies stay mostly the same over time. Otherwise the approach to record the communication graph and then use STATIC in order to compute a partitioning may lead to inconsistent results.

We stress that the implementations of STATIC and ADAPT rely on heuristics only and do not provide any guarantees while we have shown that our algorithm CREP-ADJ is a competitive algorithm with competitive ratio  $O(2/\epsilon \cdot k \log k)$ .

Finally Figure 11 illustrates the run times of ADAPT and STATIC. Note that static is only run once on the communication graph after all requests have been revealed which means that it is expected for STATIC to have the fastest running time. But also ADAPT achieves faster running times than our dynamic implementations.

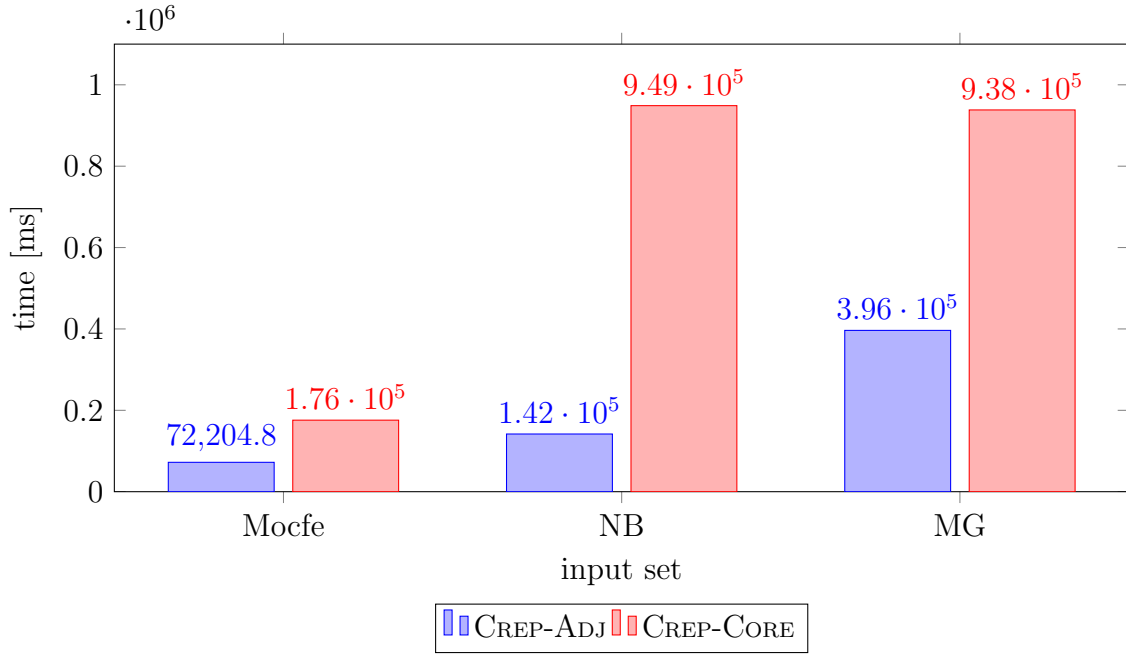


Figure 9: comparison of running time of CREP-ADJ and CREP-CORE in the case where merges are performed at connectivity  $2 \cdot \alpha$

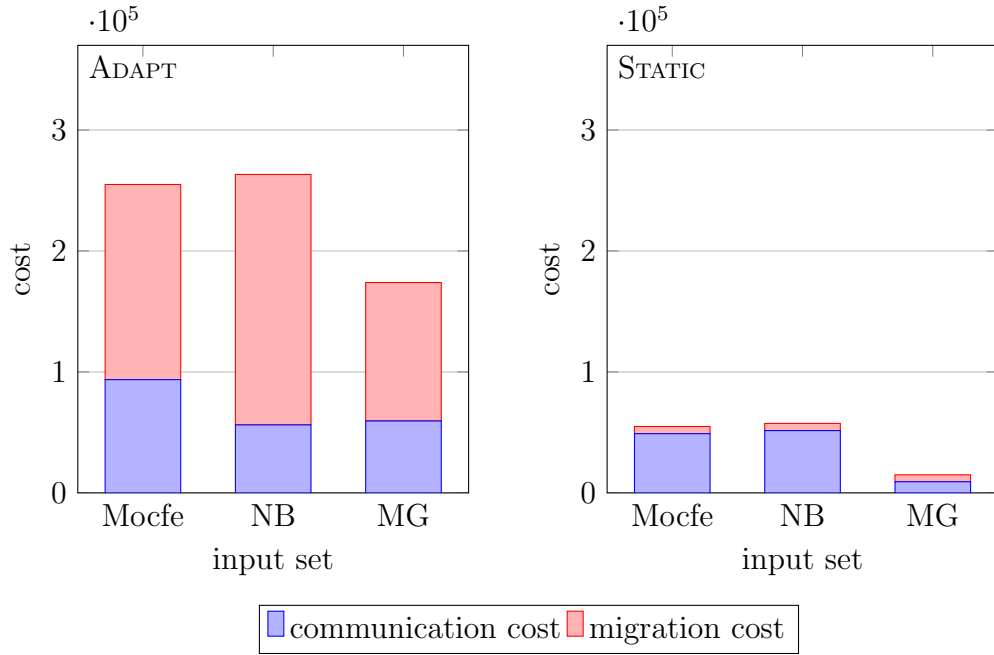


Figure 10: comparison of the total cost of ADAPT (left) and STATIC (right)

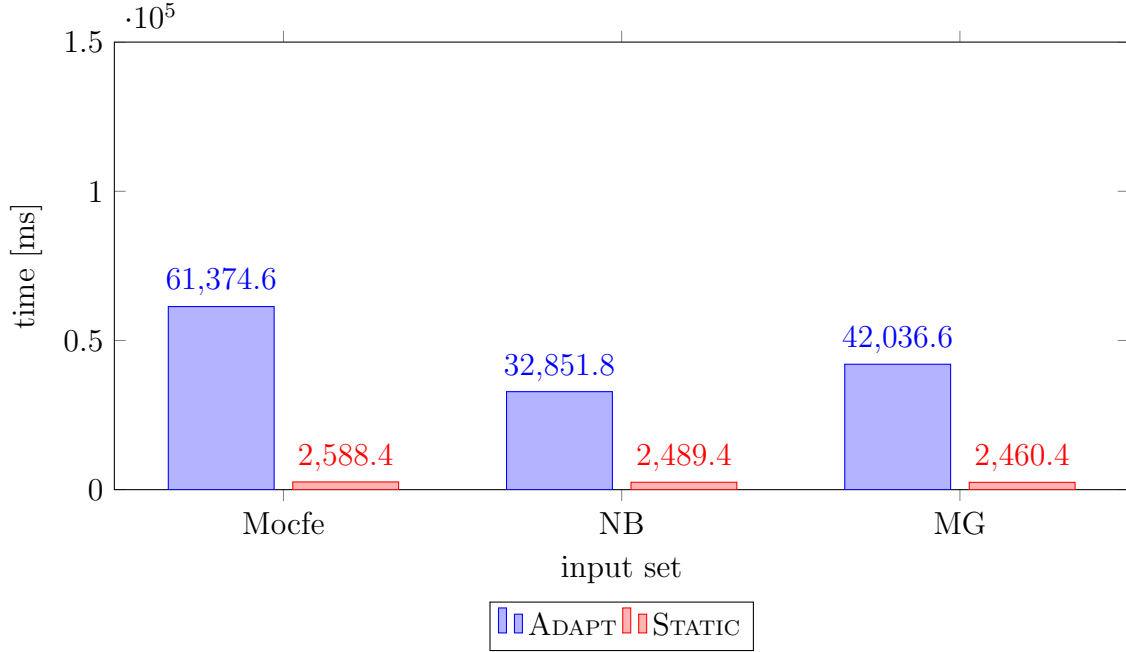


Figure 11: comparison of running time of ADAPT and STATIC

## 9 Conclusion

We have discussed different ideas for component-based algorithms for the Dynamic Balanced Graph Partitioning Problem as well as the corresponding implications and challenges for the competitive analysis.

We have shown that our algorithm CREP-ADJ has a competitive ratio of  $O(2/\epsilon \cdot k \log k)$  and we have presented our decomposition-based approach to its implementation.

Finally we have evaluated the results of our implementation by comparing it with three other algorithms: CREP-CORE which is a variation of CREP-ADJ which deletes fewer edges and ADAPT and STATIC which are available via METIS and ParMETIS frameworks, respectively. Furthermore we have explored the possibility of making slight adaptations of our algorithms in order to improve their results.

## 10 Future Work

Even though we have already shown very promising results as to the overall performance of our algorithms together with theoretical guarantees there are still areas in which our work could be expanded upon.

As discussed in Section 8.4 we see great potential in making adjustments to our algorithms, especially CREP-CORE, in order to improve their results even further while also preserving our result on the competitive ratio.

These kinds of adjustments may be achieved for example by changing certain algorithm parameters such as the connectivity threshold.

There may also be potential in adapting and improving our decomposition tree data structure in order to improve running times.

One such avenue for improvement may be to parallelize the computation of the decomposition of disjoint subgraphs. This could be both used to handle multiple requests in parallel if they only affect disjoint sub-graphs, i.e. their corresponding trees do not overlap in the decomposition tree, as well as to speed up the computation of the decomposition of the children trees during the update process after a new request has been received by CREP.

## References

- [1] Konstantin Andreev and Harald Räcke. Balanced Graph Partitioning. *Theory of Computing Systems*, 39(6):929–939, oct 2006.
- [2] Chen Avin, Marcin Bienkowski, Andreas Loukas, Maciej Pacut, and Stefan Schmid. Dynamic Balanced Graph Partitioning. *arXiv preprint arXiv:1511.02074v5*, 2015.
- [3] Chen Avin, Louis Cohen, Mahmoud Parham, and Stefan Schmid. Competitive clustering of stochastic communication patterns on a ring. *Computing*, 101(9):1369–1390, sep 2018.
- [4] Chen Avin, Manya Ghobadi, Chen Griner, and Stefan Schmid. Measuring the Complexity of Packet Traces. *arXiv:1905.08339v1*, 2019.
- [5] Chen Avin, Andreas Loukas, Maciej Pacut, and Stefan Schmid. Online Balanced Repartitioning. In *Lecture Notes in Computer Science*, pages 243–256. Springer Berlin Heidelberg, 2016.
- [6] Abba Chouni Benabdellah, Asmaa Benghabrit, and Imane Bouhaddou. A survey of clustering algorithms for an industrial context. *Procedia Computer Science*, 148:291–302, 2019.
- [7] Aydın Buluç, Henning Meyerhenke, Ilya Safro, Peter Sanders, and Christian Schulz. Recent Advances in Graph Partitioning. In *Algorithm Engineering*, pages 117–158. Springer International Publishing, 2016.
- [8] Lijun Chang, Jeffrey Xu Yu, Lu Qin, Xuemin Lin, Chengfei Liu, and Weifa Liang. Efficiently computing k-edge connected components via graph decomposition. In *Proceedings of the 2013 international conference on Management of data - SIGMOD '13*. ACM Press, 2013.
- [9] Leah Epstein, Csanád Imreh, Asaf Levin, and Judit Nagy-György. On variants of file caching. In *Automata, Languages and Programming*, pages 195–206. Springer Berlin Heidelberg, 2011.
- [10] Leah Epstein and Hanan Zebedat-Haider. Rent or buy problems with a fixed time horizon. *Theory of Computing Systems*, 56(2):309–329, jun 2014.
- [11] Amos Fiat, Richard Karp, Mike Luby, Lyle McGeoch, Daniel Sleator, and Neal E. Young. Competitive Paging Algorithms. *arXiv preprint cs/0205038*, 2002.
- [12] Monia Ghobadi, Daniel Kilper, Ratul Mahajan, Amar Phanishayee, Nikhil Devanur, Janardhan Kulkarni, Gireeja Ranade, Pierre-Alexandre Blanche, Houman Rastegarfar, and Madeleine Glick. ProjecToR: Agile Reconfigurable Data Center Interconnect. In *Proceedings of the 2016 conference on ACM SIGCOMM 2016 Conference - SIGCOMM '16*. ACM Press, 2016.

- [13] Navid Hamedazimi, Zafar Qazi, Himanshu Gupta, Vyas Sekar, Samir R. Das, Jon P. Longtin, Himanshu Shah, and Ashish Tanwer. FireFly. In *Proceedings of the 2014 ACM conference on SIGCOMM - SIGCOMM '14*. ACM Press, 2014.
- [14] Erez Hartuv and Ron Shamir. A clustering algorithm based on graph connectivity. *Information Processing Letters*, 76(4-6):175–181, dec 2000.
- [15] Monika Henzinger, Stefan Neumann, and Stefan Schmid. Efficient Distributed Workload (Re-)Embedding. *no idea*, 2019.
- [16] George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, jan 1998.
- [17] George Karypis and Vipin Kumar. Multilevelk-way partitioning scheme for irregular graphs. *Journal of Parallel and Distributed Computing*, 48(1):96–129, jan 1998.
- [18] George Karypis and Vipin Kumar. Parallel multilevel series k-way partitioning scheme for irregular graphs. *SIAM Review*, 41(2):278–300, jan 1999.
- [19] M. Pavan and M. Pelillo. A new graph-theoretic approach to clustering and segmentation. In *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.* IEEE Comput. Soc., 2003.
- [20] K. Schloegel, G. Karypis, and V. Kumar. A unified algorithm for load-balancing adaptive scientific simulations. In *ACM/IEEE SC 2000 Conference (SC'00)*. IEEE, 2000.
- [21] Kirk Schloegel, George Karypis, and Vipin Kumar. Multilevel diffusion schemes for repartitioning of adaptive meshes. *Journal of Parallel and Distributed Computing*, 47(2):109–124, dec 1997.
- [22] Mechthild Stoer and Frank Wagner. A simple min-cut algorithm. *Journal of the ACM*, 44(4):585–591, jul 1997.
- [23] C. Walshaw and M. Cross. Mesh partitioning: A multilevel balancing and refinement algorithm. *SIAM Journal on Scientific Computing*, 22(1):63–80, jan 2000.
- [24] Chris Walshaw and Mark Cross. JOSTLE: parallel multilevel graph-partitioning software—an overview. *Mesh partitioning techniques and domain decomposition techniques*, pages 27–58, 2007.
- [25] Z. Wu and R. Leahy. An optimal graph theoretic approach to data clustering: theory and its application to image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(11):1101–1113, 1993.