



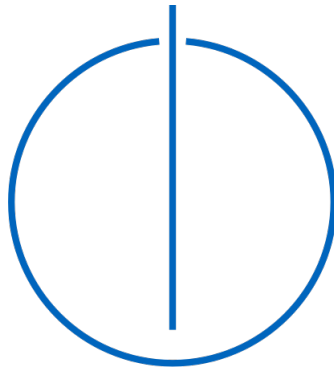
Technical University of Munich

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	15.03.2020

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

Date

Tobias Forner

Abstract

...

Contents

1	Introduction	1
2	Problem Definition	1
3	Related Work	1
3.1	Dynamic Balanced RePartitioning	1
3.2	Restricted Variant of Balanced RePartitioning	2
3.3	Clustering	2
3.4	Online Paging	2
3.5	Static Graph Partitioning	2
4	Algorithmic Ideas	2
4.1	CORE-DEL	3
4.2	ADJ-DEL	4
5	Problems With Analysis in Previous Work	4
6	Algorithm Overview	5
6.1	Overview	5
7	Analysis	6
7.1	Algorithm Definitions	6
7.2	Structural Properties	6
7.3	Upper Bound On CREP	7
7.4	Lower Bound on OPT	9
7.5	Competitive Ratio	11
8	Implementation Details	11
8.1	Algorithm Pseudocode	11
8.2	Algorithm Explanations	13
9	Evaluation	13
	References	17

1 Introduction

2 Problem Definition

The Dynamic Balanced Graph Partitioning problem was first introduced by Avin, Bienkowski, Loukas, Pacut and Schmid ([1]).

The task is to maintain a partitioning of a dynamic graph consisting of $n = k \cdot l$ nodes that communicate with each other into k parts, each of size l 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 . Such a move of one 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 σ 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 σ 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 ρ -competitive if there exists a constant β 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 σ of ONL and OPT respectively.

Often we allow the online algorithm to use larger capacities per server. In this case we speak of an *augmentation* of δ in the case where the online algorithm is allowed to assign $\delta \times n/k$ 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.

3 Related Work

3.1 Dynamic Balanced RePartitioning

Avin et al.[1] initiated the study of the online variant of the Balanced RePartitioning (BRP) problem that is the topic of this thesis. They propose a deterministic algorithm for the Dynamic Balanced Graph Partitioning problem with augmentation $2 + \epsilon$ for any $\epsilon > 1/k$. 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 δ -augmented deterministic online algorithm can achieve a competitive ratio smaller than k for any augmentation $\delta < l$.

3.2 Restricted Variant of Balanced RePartitioning

Restricted variants of the Balanced RePartitioning problem have also been studied. Here one assumes certain restrictions of the input sequence σ 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 ([6]) studies a special *learning variant* of the Dynamic Balanced Graph Partitioning problem specified above. In this version it is assumed that the input sequence σ 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 [6] 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 ϵ , 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 BRP problem without restricting σ .

3.3 Clustering

3.4 Online Paging

3.5 Static Graph Partitioning

4 Algorithmic Ideas

In this section we describe different solution approaches to the Dynamic Balanced Graph Partitioning problem. 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 σ . As more requests from σ 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. We first describe this general approach that is common to both algorithms and then address the specific differences and analysis ideas.

More formally, initially each node forms a singleton component, but as the input sequence σ 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 α . After each edge insertion the algorithm checks whether there exists a new component set S with $|S| > 1$ which fulfills 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 α 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. If the resulting component has size at least $2/\epsilon$ the algorithm reserves additional space $\min\{\epsilon \cdot |C|, k\}$, otherwise the reservation is zero.

The collocation of such component sets of at most k individual communication nodes is always possible without moving a node not in C due to the allowed augmentation of $2 + \epsilon$. This guarantees by an averaging argument that there is almost 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 differentiating factor between our algorithms which we discuss in the following subsections. We also describe the particular challenges each approach entails when it comes to the competitive analysis.

These approaches are fairly similar to the algorithms defined in previous work ([2], [1]). The main differentiating factor is that we merge once a component set reaches connectivity α while the approaches by the other authors do so once the component set reaches a certain density threshold. More specifically they merge a component set S once it 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 [1], property 3 in [2] and [Lemma 6](#).

4.1 CORE-DEL

The first one resets edges which are *contained* in the deleted component, i.e. all edges $e = \{u, v\}$ are reset to zero if u and v were contained in component C at the time of its deletion. This approach resembles the one suggested by Avin et al. ([1]).

4.2 ADJ-DEL

The second algorithm 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 its deletion. This second version is very similar to the algorithm proposed by Avin et al. in [2].

Algorithm 1 DynamicDecomp

```

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
    merge( $X$ )
  end if
  if  $\exists$  component set  $Y$  with connectivity at least  $\alpha$  and  $\text{nodes}(Y) > k$  then
    delete( $Y$ )
  end if
end for

```

The idea for the analysis of both approaches is to account cost to OPT for each component C that is deleted by one of the algorithms.

5 Problems With Analysis in Previous Work

In this section we point out problems in the analysis in two versions of [2] which share a similar approach to ours.

The problem in the analysis of [1] arises due to their usage of the concept of $F(c)$: This contains only edges incident to components from $S(c)$ that arrived after the involved components from $S(c)$ were created by CREP. In the analysis it is assumed that only those edges that are contained in $F(c)$ can contribute to the creation of component c . As we will show that is not the case and as a result it is very challenging to separate the costs of OPT that are due to the requests that actually led to the creation of c .

In order for the approach in [1] to work the different sets $F(c)$ that occur as CREP handles input σ need to form a partition of all requests such that each request can be mapped to one unique set $F(c)$. Only in this case it is possible to lower bound the cost of CREP via the requests in $\bigcup_{c \in \text{DEL}} F(c)$.

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 $t = 5$.

In this sequence the first two requests are not contained in $F(c)$ for any component c as nodes 3 and 4 were part of other components that were eventually deleted by CREP before finally being merged at time $t = 15$.

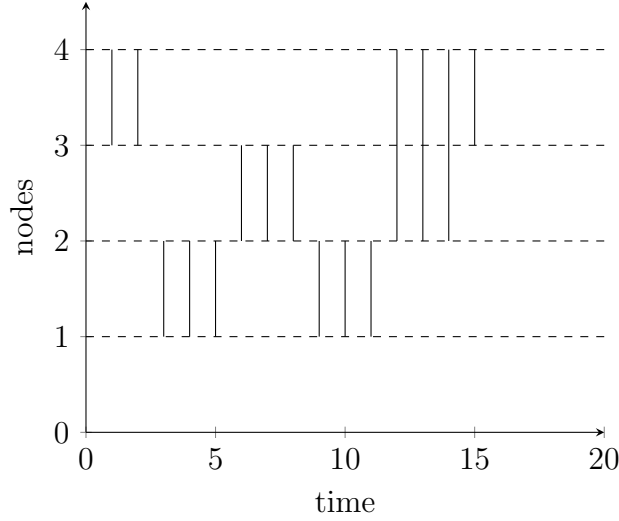


Figure 1: problem in the new version

In an older version ([2]) of the paper mentioned above we have also discovered some problems. In this version a similar component structure is maintained and merges are performed similarly to our approach as upon the deletion of a component the algorithm not only deletes inner edges but also those that leave the component, i.e. all edges $e = u, v$ are deleted where u or v are inside the component. However the additional deletions of edges that are leaving the component are not accounted for in the analysis. We show that this approach may lead to scenarios where the amount of edges that are deleted that leave the component may greatly exceed the amount of inner edges. We show in [Section 7](#) how the analysis can be adapted while preserving the competitive ratio.

[Figure 2](#) illustrates an example of the problems mentioned above for this older version of CREP. In the illustration it is assumed that $\alpha = 3$ and $k = 3$. In this case the input sequence leads CREP 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.

6 Algorithm Overview

In this section we give an overview of our algorithm CREP which builds upon the ideas of [2] and [1]. Implementation details will be discussed in greater detail in [Section 8](#).

6.1 Overview

The algorithm maintains a second-order partitioning of the nodes into *communication components* which are sets of nodes that communicated frequently. As more requests from the input sequence σ are revealed to the algorithm the components

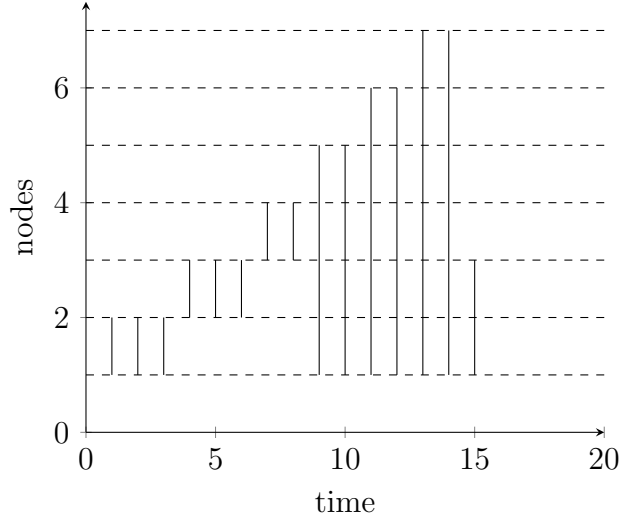


Figure 2: problem in the older version

grow in size until the algorithm discovers a component that is too large and hence decides to delete it.

7 Analysis

We analyse the competitive ratio of CREP with augmentation $(2 + \epsilon)$ and show in [Theorem 14](#) that CREP is $O(k \log k)$ -competitive.

7.1 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 .

Note that such an edge can only have positive weight if its endpoints are in different components.

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

7.2 Structural Properties

Note: These properties are changed to use the connectivity based approach which generally simplifies them, but guarantees slightly less minimum edge weight within mergeable component sets.

Definition 3. An α -connected component is a maximal set of vertices that is α -connected.

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 α , i.e. there exist no mergeable component sets after CREP performed its merges.

Proof. We proof 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 α and $|S| > 1$. We may assume that t is the earliest time for which S has connectivity α .

Then the incrementation of the weight of edge e at time t raised the connectivity of S , but S was not merged into a new α -connected component C . if no new component was created at time t then 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 α which is a contradiction to the maximality of C and S . \square

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 α and a mergeable component set S has connectivity exactly α .

Proof. This lemma follows directly from lemma 2 as connectivities can only increase by at most 1 at each time t . \square

Lemma 6. The weight between the components of a component subset S of connectivity α is at least $|S|/2 \cdot \alpha$.

Proof. Consider the sum of the weighted degrees of all components:

$$\sum_{c \in S} \deg_S(c) = 2 \sum_{e \in S} w(e)$$

The equality follows as the left sum counts each edge twice, once for each endpoint. Now consider the fact that each component must have degree at least α with respect to the edges in S as S has connectivity α and hence the lemma follows. \square

Lemma 7. The weight between the components of a component subset S of connectivity α is at most $(|S| - 1) \cdot \alpha$.

Proof. We iteratively partition S into subsets via minimum cuts with regard to edge weight, i.e. we consider a minimum edge cut of S which partitions S into the subsets S_1 and S_2 and iteratively partition the resulting sets until all sets contain only one component each. As this required at most $|S| - 1$ cuts of value at most α the lemma follows. \square

7.3 Upper Bound On CREP

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

We define the following notions for a component $C \in \text{DEL}(\sigma)$, i.e. the subgraph induced by the nodes of C has connectivity at least α and C consists of more than k nodes:

Let $\text{EPOCH}(C)$ denote the (node, time) pairs of nodes in C starting at the time after the time $\tau(\text{node})$ when node was last turned into a singleton component, i.e. $\text{EPOCH}(C) = \bigcup_{n \in \text{nodes}(C)} \{n\} \times \{\tau(n) + 1, \dots, \tau(C)\}$. 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 σ).

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: $\text{CORE}(C)$ 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 . $\text{HALO}(C)$ contains 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 $\text{HALO}(C) = \text{REQ}(C) \setminus \text{CORE}(C)$.

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. The lemma follows directly from Lemma 7 due to the fact that component sets of connectivity α get merged immediately by CREP as shown in Lemma 5. \square

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

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 σ is bounded by

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

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 requests from $\text{HALO}(C)$ were deleted. 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 $C \in \text{FINAL-COMPS}(\sigma)$ and are hence also bounded by $|C| \cdot \alpha$ or have not let CREP to any merge and are hence contained in $\text{FINAL-COMPS}(\sigma)$. 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

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

Proof. First note that CREP only performs migrations when it merges components. We fix a component $C \in \text{DEL}(\sigma) \cup \text{FINAL-COMPS}(\sigma)$ 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 components of size at least $2/\epsilon$ 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 α to CREP. \square

Finally we summarize our results in the following lemma.

Lemma 11. With augmentation $2 + \epsilon$, CREP pays at most total cost $2 \cdot \sum_{C \in \text{DEL}(\sigma) \cup \text{FINAL-COMPS}(\sigma)} |C| \cdot ((2/\epsilon + 1) + \log k) \cdot \alpha + \sum_{C \in \text{DEL}(\sigma)} |\text{HALO}(C)| + \text{FINAL-WEIGHTS}(\sigma)$.

Proof. We sum the results from Lemma 9 and Lemma 10:

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

\square

7.4 Lower Bound on OPT

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.

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 requests in question by N , i.e. those from $N = \mathcal{H} \setminus (P \cup OP)$ where OP denotes the requests from \mathcal{H} only OPT pays for.

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

Proof. We fix an arbitrary offline interval. 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 α requests between v and any one of the other nodes on its server. This bounds the number of requests assigned to the offline interval by $k \cdot \alpha$. \square

In the following lemma we combine this result with a bound on the cost of OPT due to requests we have not accounted for, namely those from $\text{CORE}(C)$ and those that are not contained in an offline interval because they get deleted by CREP due to a component deletion that takes place before their offline interval ends.

Lemma 13. The cost of the solution of OPT given input sequence σ 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. 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.

For the following part of the proof we fix an arbitrary component $C \in \text{DEL}(\sigma)$. Let R denote the set of requests from $\text{HALO}(C)$ that were not assigned to any offline interval. This means that the nodes involved in requests from R were not moved during the processing of requests from R until the time of deletion of C .

The number of nodes contained in C or connected to C via edges representing requests from R is at least $|C| + R/\alpha$ since requests from R have not led CREP to perform any migrations. Because of this fact OPT must have placed those nodes on at least $\frac{|C|+R/\alpha}{k}$ different servers. As OPT does not pay for any requests from R it follows that OPT must have placed the nodes from C in $\frac{|C|+R/\alpha}{k}$ 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 $\text{CORE}(C)$ into migrations. Because of this fact OPT must have placed those nodes on at least $\frac{|C|+R/\alpha}{k}$ parts. As CREP merged component C this graph is α -connected and hence Lemma 6 gives that OPT has to cut at least edges of total weight migrations. Because of this fact OPT must have placed those nodes on at least $\frac{|C|+R/\alpha}{k} \cdot \alpha = |C|/k \cdot \alpha + R/k$.

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 α 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 $\frac{|C|+R/\alpha}{k}$ parts.

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

Finally we need to account for the fact that the migrations of nodes from C that OPT performs also end offline intervals and might hence be accounted for twice in our analysis up to this point:

$$\begin{aligned} 2 \cdot \text{OPT}(\sigma) &\geq \sum_{(C \in \text{DEL}(\sigma))} |C|/k \cdot \alpha + R(C)/k + I/(k \cdot \alpha) \cdot \alpha + |P| + |OP| \\ &\geq \sum_{(C \in \text{DEL}(\sigma))} |C|/k \cdot \alpha + \text{HALO}(C)/k \end{aligned}$$

where the last equality follows from the fact that $\bigcup_{C \in \text{DEL}(\sigma)} R(C) \cup I \cup P \cup OP = \bigcup_{C \in \text{DEL}(\sigma)} \text{HALO}(C)$ as the different $R(C)$ as well as I, P and OP are disjoint. Hence the cost of OPT is at least $1/2 \cdot \sum_{(C \in \text{DEL}(\sigma))} |C|/k \cdot \alpha + \text{HALO}(C)/k$ as OPT pays for requests from P by the definition of P . \square

7.5 Competitive Ratio

In this section we combine the results of [Lemma 11](#) and [Lemma 13](#) to obtain the following theorem.

Theorem 14. With augmentation $(2 + \epsilon)$ the competitive ratio of CREP is in $O(k \log k)$.

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

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

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

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

8 Implementation Details

8.1 Algorithm Pseudocode

Algorithm 2 insertAndUpdate(a, b)

```
if comp[a] == comp[b] then
    return
end if
addEdge( $a, b$ )
updateDecomposition( $a, b$ )
del ← updateMapping(alphaConnectedComponents)
delComponents(del)
```

Algorithm 3 updateDecomposition(a, b)

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

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

Algorithm 4 delComponents(del)

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

8.2 Algorithm Explanations

- [Algorithm 2](#) calls the other routines as needed
- [Algorithm 3](#) starts at the smallest subgraph containing the nodes a and b in the decomposition tree and computes a new decomposition of the subgraph. Specifically it uses the decomposition approach from [5] to decompose one subgraph and then also computes the subgraphs with the next higher connectivity and recurses until the connectivity has reached α .
- `updateMapping` checks whether the `alphaConnectedComponents` were changed. If yes then it either collocates them if the resulting component is small enough or it adds the component to its return value. Then all the returned components are deleted, i.e. the edges connecting its nodes are deleted and the decomposition is recomputed
- this deletion is performed by [Algorithm 4](#)

9 Evaluation

In this section we evaluate the quality and performance of our algorithm implementation.

For this evaluation we compare our implementation described in [Section 8](#) to a static algorithm available via Metis (`METIS_PartGraphRecursive`) ([7, 8]) and an adaptive/dynamic algorithm (`ParMETIS_V3_AdaptiveRepart`) ([9, 11, 10]) implemented in the ParMetis framework. Both frameworks are known to produce very good results and to be very fast.

As input data we use several HPC traces, the nature of the data is described in more detail by Avin, Ghobadi, Griner and Schmid in [4].

All data sets contain 1024 different communication nodes and are limited to the first 300 000 requests. The value of α is set to be 6 and the algorithm was tasked to partition the nodes into 32 clusters of size 32 each. The dynamic algorithms were allowed to use augmentation with a factor of 2.1, i.e. for the dynamic algorithms the maximum cluster capacities were $\lfloor 32 \cdot 2.1 \rfloor = 67$.

To our knowledge it is not possible to specify hard limits for the capacities used by the static algorithm implemented in Metis and as a result there are some occurrences where the algorithm exceeds the capacities that are allowed by a small amount.

We will first discuss the overall results of our experiments, i.e. we will describe the quality of the results (see [Figure 3](#)) as well as the running time needed for each examined algorithm (see [Figure 4](#)).

The static algorithm is shown to give the best results in the shortest time, but the low running was also to be expected as it is only called once as opposed to the 300 000 times the other algorithms need to decide whether they want to change their partitioning. The static algorithm also has knowledge of all requests and as a result is able to produce the best results.

For data set A our decomposition algorithm beats the adaptive ParMetis algorithm by a significant amount of about one third of the cost of the latter while ParMetis is significantly faster. For data sets B and C ParMetis is shown to produce slightly better results within drastically less computation time. It is worth mentioning that

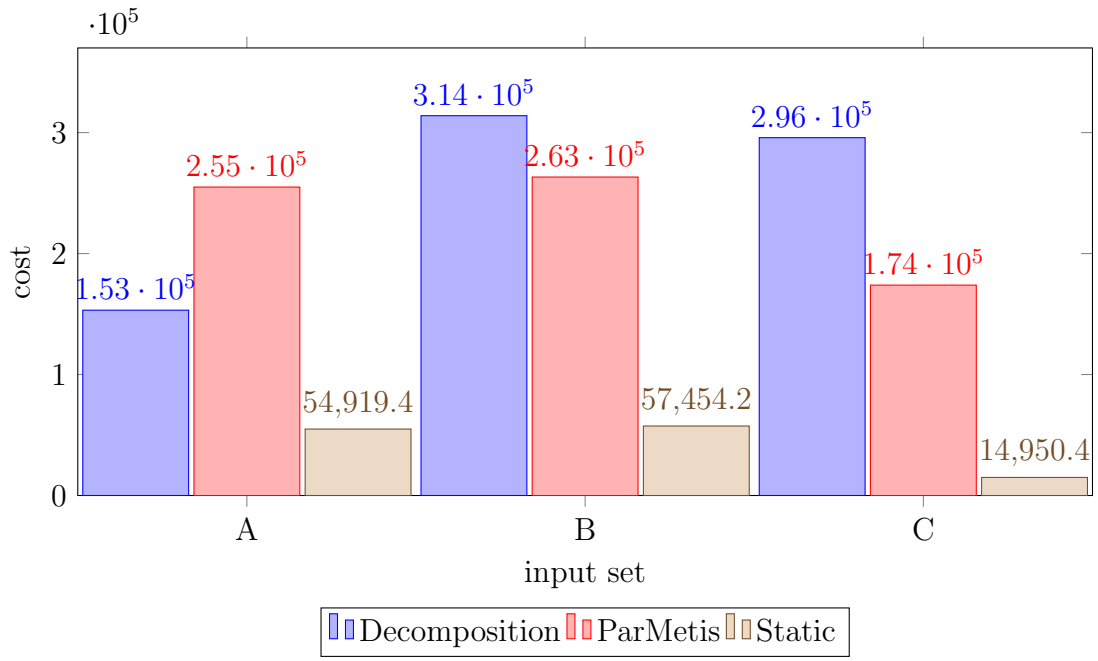


Figure 3: comparison of total cost

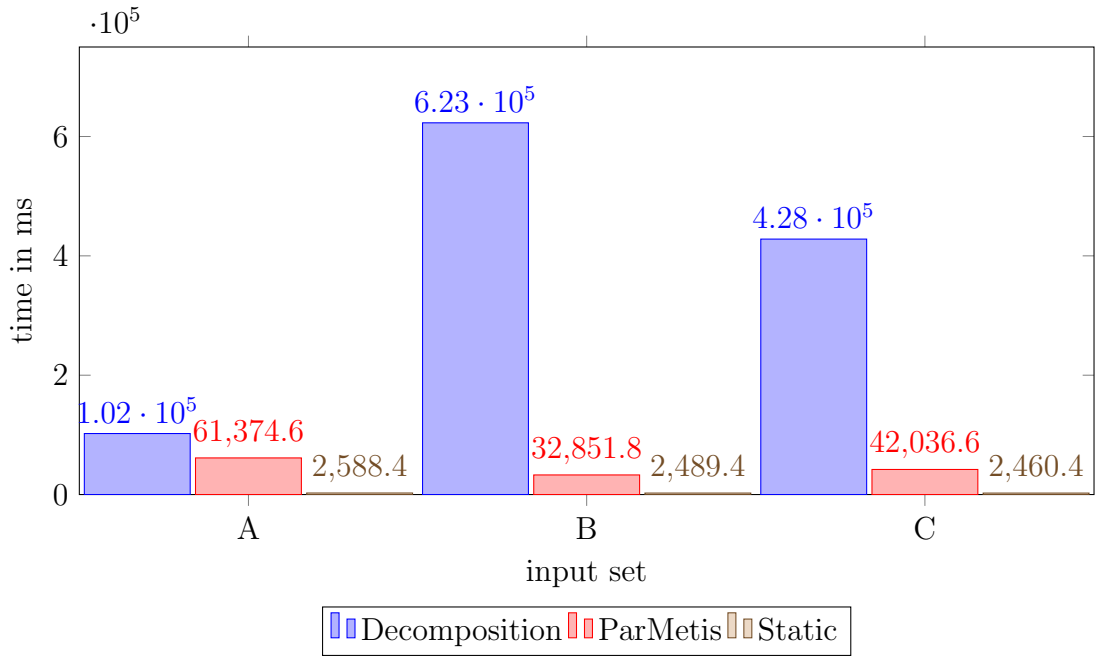


Figure 4: comparison of run time

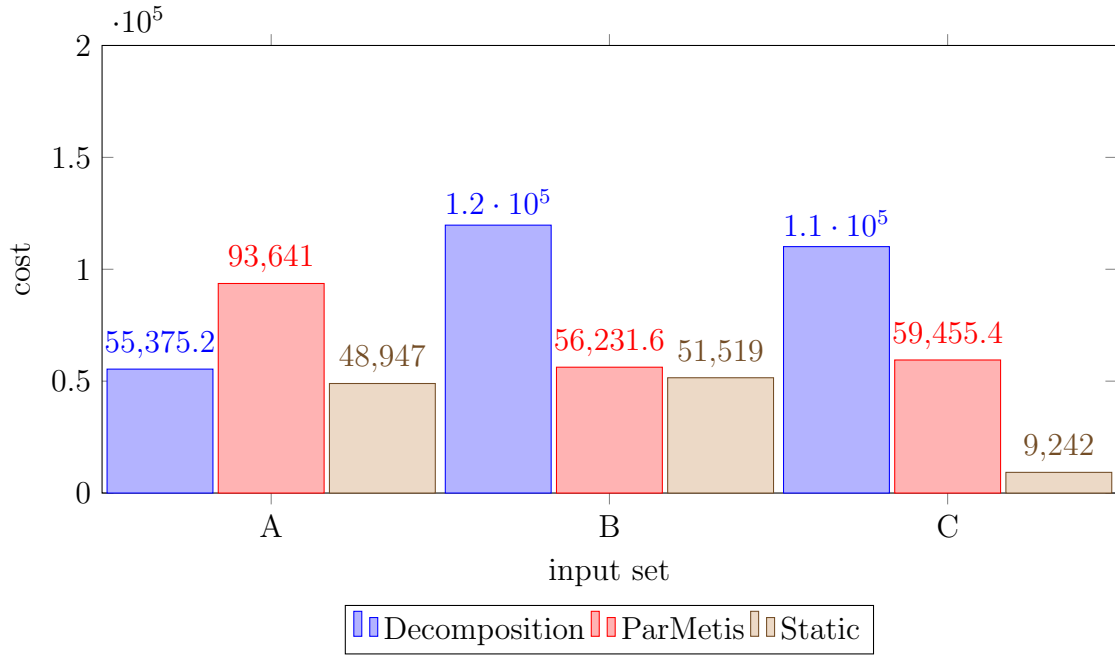


Figure 5: comparison of communication cost

ParMetis uses the ... graph description format that does not allow for easy adaptation on the fly and needed to be recomputed after every request during our tests. However, we chose not to include this in the running time calculations.

In the next section we discuss the general distribution of the total costs of each algorithm to communication (see Figure 5) and migration costs (Figure 6). Both ParMetis as well as our decomposition algorithm produce significantly more migration cost than communication cost while the static algorithm predominantly pays for communication. This shows that the dynamic algorithms tend to migrate too much while the static implementation is restricted to only migrate once to a static configuration it finds suitable and as a result has to pay more for communication. This also shows that there is potential to refine the dynamic implementations in such a way that they produce more balanced, and hopefully also less, cost overall.

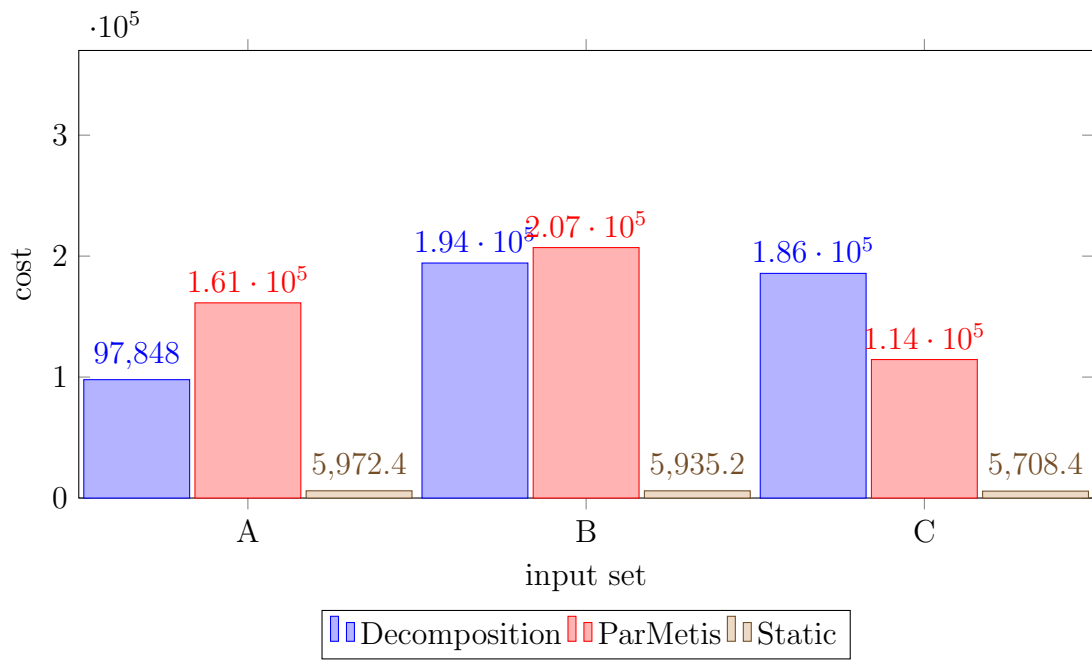


Figure 6: comparison of migration cost

References

- [1] Chen Avin, Marcin Bienkowski, Andreas Loukas, Maciej Pacut, and Stefan Schmid. Dynamic Balanced Graph Partitioning.
- [2] Chen Avin, Marcin Bienkowski, Andreas Loukas, Maciej Pacut, and Stefan Schmid. Dynamic Balanced Graph Partitioning. *no idea*, 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.
- [5] 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.
- [6] Monika Henzinger, Stefan Neumann, and Stefan Schmid. Efficient Distributed Workload (Re-)Embedding. *no idea*, 2019.
- [7] 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.
- [8] George Karypis and Vipin Kumar. Multilevelk-way partitioning scheme for irregular graphs. *Journal of Parallel and Distributed Computing*, 48(1):96–129, jan 1998.
- [9] George Karypis and Vipin Kumar. Parallel multilevel series k-way partitioning scheme for irregular graphs. *SIAM Review*, 41(2):278–300, jan 1999.
- [10] 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.
- [11] 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.