

Manuscript Number:

Title: Data Locality and Replica Aware Virtual Cluster Embeddings

Article Type: Regular Paper (10 - 40 pages)

Section/Category: A - Algorithms, automata, complexity and games

Keywords: Virtual Network Embeddings

Flow Algorithms

NP hardness

Corresponding Author: Prof. Stefan Schmid,

Corresponding Author's Institution: Telekom Innovation Labs

First Author: Carlo Fuerst

Order of Authors: Carlo Fuerst; Maciek Pacut; Stefan Schmid



11. DECEMBER 2015

AALBORG UNIVERSITY
DENMARK

ASSOCIATE PROFESSOR
DR. STEFAN SCHMID
DEPT. COMPUTER SCIENCE
SELMA LAGERLØFS VEJ 300
DK 9220 AALBORG
SCHMISTE@CS.AAU.DK

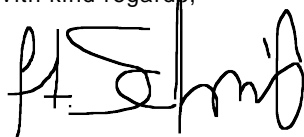
Elsevier TCS Submission

We are happy to submit our work on locality-aware cluster embeddings to TCS, and herewith confirm that the paper is not under submission at any other journal.

An early version of the paper appeared at the ICNP 2015 conference, and the paper includes a discussion of the differences.

We would like to thank the editor and the reviewers for their time and work.

With kind regards,

A handwritten signature in black ink, appearing to read 'Stefan Schmid'.

Stefan Schmid
+49 175 930 98 75
schmiste@gmail.com

Data Locality and Replica Aware Virtual Cluster Embeddings

Carlo Fuerst¹ Maciej Pacut² Stefan Schmid³
¹ TU Berlin, Germany ² University of Wroclaw, Poland
³ Aalborg University, Denmark

Abstract

Virtualized datacenters offer great flexibilities in terms of resource allocation. In particular, by decoupling applications from the constraints of the underlying infrastructure, virtualization supports an optimized mapping of virtual machines as well as their interconnecting network (the so-called *virtual cluster*) to their physical counterparts: a graph embedding problem.

However, existing virtual cluster embedding algorithms such as Oktopus, Proteus, and Kraken, often ignore a crucial dimension of the problem, namely *data locality*: the input to a cloud application such as MapReduce is typically stored in a distributed, and sometimes redundant, file system. Since moving data is costly, an embedding algorithm should be data locality aware, and allocate computational resources close to the data; in case of redundant storage, the algorithm should also optimize the *replica selection*.

This paper initiates the algorithmic study of data locality aware virtual cluster embeddings on datacenter topologies. We show that despite the multiple degrees of freedom in terms of embedding, replica selection and assignment, many problems can be solved efficiently. We also highlight the limitations of such optimizations, by presenting several NP-hardness proofs; interestingly, our hardness results also hold in uncapacitated networks of small diameter.

1. Introduction

Distributed cloud applications, such as batch-processing applications or scale-out databases, generate a significant amount of network traffic [26]. For instance, MapReduce consists of a network intensive shuffle phase, where data

is transferred from the mappers to the reducers. In order to ensure a predictable application performance, especially in shared cloud environments, it is important to provide isolation and bandwidth guarantees between the virtual machines [38], e.g., by making explicit network reservations [5]. Accordingly, modern batch-processing applications provide the abstraction of entire *virtual networks* [26], defining both the virtual machines as well as their interconnecting network. The most prominent virtual network abstraction is the *virtual cluster* [5, 16, 31, 35].

Virtualized datacenters offer great flexibilities on where these virtual networks can be instantiated or *embedded*. In order to maximize the resource utilization in the datacenter, it is in principle desirable to map the virtual machines of a given virtual network as close as possible in the underlying physical network, as this minimizes communication costs (respectively, bandwidth reservations) [5, 16, 31, 35].

However, existing systems often ignore a crucial dimension of the virtual network embedding problem: the fact that the input data for a cloud application, consisting of atomic *chunks*, is typically distributed across different servers and stored in a distributed file system [6, 18, 33]. In order to properly minimize communication costs, an embedding algorithm should hence also be *data locality aware* [4, 24, 37], and allocate (or *embed*) computational resources close to the to be processed data. Moreover, in case of redundant storage (batch processing applications often provide a 3-fold redundancy [33]), an algorithm should also be aware of, and exploit, *replica selection* flexibilities.

1.1. Our Contributions

This paper initiates the formal study of data-locality and replica aware virtual network embedding problems in datacenters. In particular, we decompose the general optimization problem into its fundamental aspects, such as assignment of chunks, replica selection, and flexible virtual machine placement, and answer questions such as:

1. Which chunks to assign to which virtual machine?
2. How to exploit redundancy and select good replicas?
3. How to efficiently embed virtual machines and their inter-connecting network?
4. Can the chunk assignment, replica selection and virtual machine embedding problems be jointly optimized, in polynomial time?

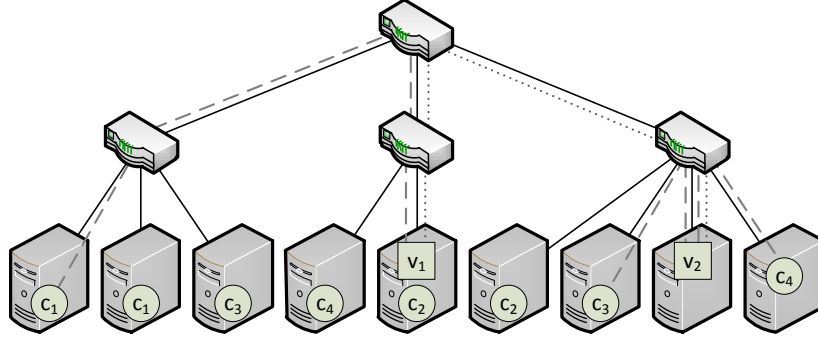


Figure 1: Overview: a 9-server datacenter storing $\tau = 4$ different chunk types $\{c_1, \dots, c_4\}$ (depicted as *circles*). The chunk replicas need to be selected and assigned to the two virtual machines v_1 and v_2 ; the virtual machines are depicted as *squares*, and the network connecting them to chunks (at bandwidth b_1) is *dashed*. In addition, the virtual machines are inter-connected among each other at bandwidth b_2 (*dotted*). The objective of the embedding algorithm is to minimize the overall bandwidth allocation (sum of *dashed* and *dotted* lines).

We draw a complete picture of the problem space: We show that even problem variants exhibiting multiple degrees of freedom in terms of replica selection and embedding, can be solved optimally in polynomial time, and we present several efficient algorithms accordingly. However, we also prove limitations in terms of computational tractability, by providing reductions from 3-D matching and Boolean satisfiability (SAT). Interestingly, while it is well-known that (unsplittable) multi-commodity flow problems are NP-hard in capacitated networks, our hardness results also hold in *uncapacitated* networks; moreover, we show that NP-hard problems already arise in small-diameter networks (as they are widely used today [2]).

1.2. Organization

Section 2 introduces our formal model in detail. Algorithms are presented in Section 3 and hardness results are presented in Section 4. Section 5 takes a deeper look at the NP-hardness variants. After discussing related work in Section 6, we conclude our work in Section 7.

2. Model

To get started, and before introducing our formal model and its constituting parts in detail, we will discuss the practical motivation. Figure 1 gives an overview of our model.

2.1. Background and Practical Motivation

Our model is motivated by batch-processing applications such as MapReduce. Such applications use multiple virtual machines to process data, often redundantly stored in a distributed file system implemented by multiple servers. [4, 10] The standard datacenter topologies today are (multi-rooted) fat-tree resp. *Clos* topologies [2, 21], hierarchical networks recursively made of sub-trees at each level; servers are located at the tree leaves. Given the amount of multiplexing over the mesh of links and the availability of multi-path routing protocol, e.g. ECMP, the redundant links can be considered as a single aggregate link for bandwidth reservations [5, 16, 31, 35].

During execution, batch-processing applications typically cycle through different phases, most prominently, a map phase and a reduce phase; between the two phases, a shuffling operation is performed, a phase where the results from the mappers are communicated to the reducers. Since the shuffling phase can constitute a non-negligible part of the overall runtime [7], and since concurrent network transmissions can introduce interference and performance unpredictability [38], it is important to provide explicit minimal bandwidth guarantees [26]. In particular, we model the virtual network connecting the virtual machines as a virtual cluster [5, 26, 35]; however, we extend this model with a notion of data-locality. In particular, we distinguish between the bandwidth needed between the assigned chunk and virtual machine (b_1) and the bandwidth needed between two virtual machines (b_2).

2.2. Formal Model

Let us now introduce our model more formally. The model combines three components: (1) the substrate network (the servers and the connecting physical network), (2) the input which needs to be processed (divided into data chunks), and (3) the virtual network (the virtual machines and the logical network connecting the machines to each other as well as to the chunks).

The Substrate Network. The substrate network (also known as the *host graph*) represents the physical resources: a set S of $n_S = |S|$ servers interconnected by a network consisting of a set R of routers (or switches) and a set E of (symmetric) links; we will often refer to the elements in $S \cup R$ as the *vertices*. We will assume that the inter-connecting network forms an (arbitrary, not necessarily balanced or regular) tree, where the servers are located at the tree leaves. Each server $s \in S$ can host a certain number of virtual machines (available server capacity $cap(s)$), and each link $e \in E$ has a certain bandwidth capacity $cap(e)$.

The Input Data. The to be processed data constitutes the input to the batch-processing application. The data is stored in a distributed manner; this spatial distribution is given and not subject to optimization. The input data consists of τ different *chunk types* $\{c_1, \dots, c_\tau\}$, where each chunk type c_i can have $r_i \geq 1$ instances (or replicas) $\{c_i^{(1)}, \dots, c_i^{(r_i)}\}$, stored at different servers. A single server may host multiple chunks. It is sufficient to process one replica, and we will sometimes refer to this replica as the *active* (or selected) replica.

The Virtual Network. The virtual network consists of a set V of $n_V = |V|$ virtual machines, henceforth often simply called *nodes*. Each node $v \in V$ can be placed (or, synonymously, *embedded*) on a server; this placement can be subject to optimization.

Depending on the available capacity $cap(s)$ of server s , multiple nodes may be hosted on s . We will denote the server s hosting node v by $\pi(v) = s$. Since these nodes process the input data, they need to be assigned and connected to the chunks. Concretely, for each chunk type c_i , exactly one replica $c_i^{(j)}$ must be processed by exactly one node v ; which replica $c_i^{(k)}$ is chosen is subject to optimization, and we will denote by μ the assignment of nodes to chunks.

In order to ensure a predictable application performance, both the connection to the chunks as well as the interconnection between the nodes may have to ensure certain minimal bandwidth guarantees; we will refer to the first type of virtual network as the (*chunk*) *access network*, and to the second type of virtual network as the (*node*) *inter-connect*; the latter is modeled as a complete network (a *clique*). Concretely, we assume that an active chunk is connected to its node at a minimal (guaranteed) bandwidth b_1 , and a node is connected to any other node at minimal (guaranteed) bandwidth b_2 .

2.3. Optimization Objective

Our goal is to develop algorithms which accept and embed a request whenever this is possible, and minimize the *resource footprint*: the amount of resources which have to be dedicated to a request, in order to realize its guarantees.

Formally, let $dist(v, c)$ denote the distance (in the underlying physical network T) between a node v and its assigned (active) chunk replica c , and let $dist(v_1, v_2)$ denote the distance between the two nodes v_1 and v_2 . We

define the *footprint* $F(v)$ of a node v as follows:

$$F(v) = \sum_{c \in \mu(v)} b_1 \cdot \text{dist}(v, c) + \underbrace{\frac{1}{2} \cdot \sum_{v' \in V \setminus \{v\}} b_2 \cdot \text{dist}(v, v')}_{\text{only for inter-connect}},$$

where $\mu(v)$ is the set of chunks assigned to v . Our goal is to minimize the overall footprint $F = \sum_{v \in V} F(v)$.

2.4. Problem Decomposition

In order to chart the landscape of the computational tractability and intractability of different problem variants, we decompose our problem into its fundamental aspects, namely replica selection (RS), multiple chunk assignment (MA), flexible node placement (FP), node interconnect (NI), and bandwidth constraints (BW), as described in the following. In this paper, we will consider all possible 32 problem variants, where each of these five aspects can either be enabled or disabled.

Replica Selection (RS). The first fundamental problem is replica selection: if the input data is stored redundantly, the algorithm has the freedom to choose a replica for each chunk type, and assign it to a virtual machine (i.e., *node*). In the following, we will refer to a scenario with redundant chunks by RS; in the RS-only scenario, the number of chunk types is equal to the number of nodes. Otherwise, we will add the +MA property discussed next.

Multiple Assignment (MA). If the number of chunk types τ is larger than the number of nodes, each node needs to be assigned multiple chunks. We will refer to such a scenario by MA. Since all nodes are identical and no additional information regarding the chunks is available at request time, we assume that each node will process an identical integer number of chunks $m = \tau/n_V$.

Flexible Placement (FP). While the nodes are placed a priori in some cases, the node placement (or synonymously: *embedding*) of nodes on physical servers can also be subject to optimization. We will refer to this degree of freedom by FP.

Node Interconnect (NI). We distinguish between scenarios where bandwidth needs to be reserved both from each node to its assigned chunks as well as to the other nodes (i.e., $b_1 > 0$ and $b_2 > 0$), and scenarios where only the (chunk) access network requires bandwidth reservation (i.e., $b_1 > 0$

and $b_2 = 0$). We will refer to the former scenario where bandwidth needs to be reserved also for the inter-connect, by NI. The node interconnect is modelled as a complete graph, to account for the all to all communication patterns of batch processing applications such as MapReduce.

Bandwidth Capacities (BW). We distinguish between an uncapacitated and a capacitated scenario where the links of the substrate network come with bandwidth constraints, and will refer to the bandwidth-constrained version by BW; the capacity of servers (the number of nodes which can be hosted concurrently) is always limited. Note that capacity constraints introduce infeasible problem instances, where it is impossible to allocate sufficient resources to satisfy an embedding request.

3. Polynomial-Time Algorithms

Despite the various degrees of freedom in terms of embedding and replica selection, we can solve many problem variants efficiently. This section introduces three general techniques, which can roughly be categorized into *flow* (Section 3.1), *matching* (Section 3.2) and *dynamic programming* (Section 3.3) approaches. First, let us make a simplifying observation:

Observation 1. *In problems without flexible placement (FP), the bandwidth required for the inter-connect network (NI) can be allocated upfront, as it does not depend on the replica selection and assignment. Accordingly, we can reduce problem variant RS + MA + NI + BW (as well as all its subproblems) to RS + MA + BW (resp. its subproblems).*

3.1. Flow Algorithms (RS + MA + NI + BW)

We first present an algorithm to solve the RS + MA + NI + BW problem. Recall that in this problem variant, we are given a set of redundant chunks (RS) and a set of nodes (the *nodes*) at fixed locations (no FP). The number of chunk types is larger than the number of nodes (MA), and each node needs to be connected to its selected chunks as well as to other nodes (NI), while respecting capacity constraints (BW). Our goal is to minimize the resource footprint F , consisting of the bandwidth reservations in the (chunk) access network and the (node) inter-connect. As we will see in the following, we can use a flow approach to solve this problem variant.

Construction of Artificial Graph. In order to solve the RS + MA + NI + BW problem, we first remove the NI property using Observation 1.

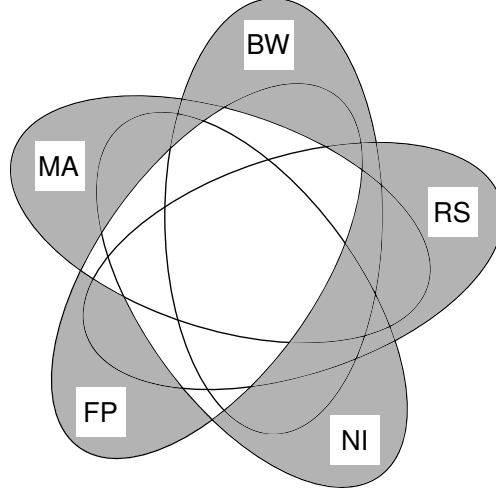


Figure 2: Variants solved by flow approach.

We then construct an artificial graph T^* , extending the substrate network T and normalizing bandwidth capacities, as follows. For T^* , we normalize the bandwidth of T to integer multiples of b_1 , i.e., for each link $e \in E(T)$, we set its new capacity in T^* to $\lfloor \text{cap}(e)/b_1 \rfloor$. After this normalization, we extend the topology T by introducing an artificial vertex for each chunk type. These artificial vertices are connected to each leaf (i.e., server) in T where a replica of the respective chunk type is located, connecting the replica of the respective chunk type by a link of capacity 1. In addition, we create a *super-source* s^+ , and connect it to each of the artificial chunk type vertices (with a link of capacity 1). Moreover, we create an artificial *super-sink* s^- and connect it to every leaf containing at least one node; the link capacity represents the number of nodes x hosted on this server, times the multi-assignment factor m . We additionally assign the following costs to edges of T^* : every edge of the original substrate network costs one unit, and all other artificial edges cost nothing.

A solution to the RS + MA + BW problem can now be computed from a solution to the *Min-Cost-Max-Flow* problem between super-source s^+ and super-sink s^- on the artificial graph T^* .

Example. Figure 3 shows an example of the extended substrate network T^* : The sink s^- is connected to the two leaves, which host the nodes. The artificial nodes are depicted below the leaves, are labeled with their respective chunk types (e.g., c_1), and are connected to the source s^+ as well as

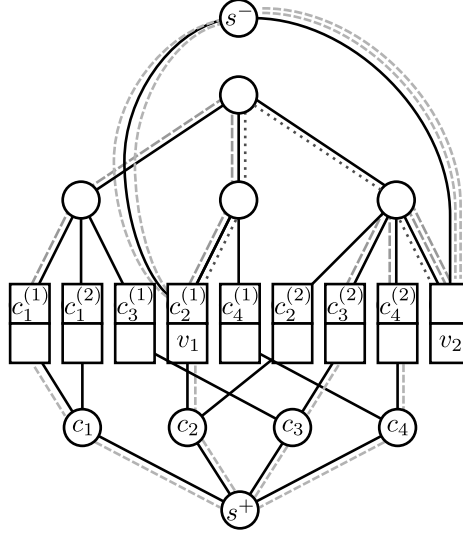


Figure 3: Example of flow construction: Problem instance with two nodes, four chunk types, and two replicas per type. The min-cost-max-flow is indicated by the dashed lines: each line represents one unit of flow.

to the leaves which contain replicas of their chunk type. The maximum flow with minimal costs is indicated by the dashed lines: each line represents one unit of flow. The dotted lines indicate links which have reduced capacity due to NI.

Algorithm. Our algorithm to solve RS + MA + NI + BW consists of three parts: *First*, we construct the normalized and extended graph T^* described above and compute a min-cost-max-flow solution, e.g., using [19, 34]. *Second*, we have to *round* the resulting, possibly fractional flow, to integer values. Due to the *integrality theorem* [1], there always exists an optimal integer solution on graphs with integer capacities. However, while algorithms like the successive shortest path algorithm [25] directly give us such an integral solution (in polynomial time), the fastest min-cost-max-flow algorithms (e.g., based on double-scaling methods [19] or minimum mean-cost cycle algorithms [34], may yield fractional solutions which need to be rounded to integral solutions (of the same cost). In order to compute integral solutions, we proceed as follows: we iteratively pick an arbitrary (loop-free) path currently having a fractional allocation of value f ($f > 0$), and distribute its flow f among all other fractional paths of the same length; due to the optimality of the fractional solution and due to the integrality theorem, such paths must always exist. After distributing this flow, the total allocation

on this path will be 0, and we have increased the number of integer paths by at least one. We proceed until we constructed the perfect matching. *Third*, given an integer min-cost-max-flow solution, we need to decompose the integer flow into the paths representing matched chunk-node pairs: The assignment can be obtained by decomposing the flow allocated in the original substrate network. In order to identify a matched chunk-node pair, we take an arbitrary (loop-free) path p carrying a flow of value ≥ 1 from s^+ to s^- : the first hop represents the chosen chunk type, the second hop the chosen replica, and the last but one hop represents the server: we will assign the replica to an arbitrary unused node on this server. Having found this pair, we reduce the flow along the path p by one unit. We continue the pairing process until every chunk type is assigned.

Analysis. The correctness of our approach follows from our construction of T^* , using integer capacities (in our case $\lfloor \text{cap}(e)/b_1 \rfloor$), and the fact that cost optimal integral solutions always exist [1]. The runtime of our algorithm consists of four parts: construction of T^* , computation of the min-cost-max-flow, flow rounding, and decomposition. The dominant term in the asymptotic runtime is the flow computation. Using the state-of-the-art min-cost-max-flow algorithms [19, 34] we get a runtime of $\mathcal{O}(n_S^2 \cdot \log \log \min\{U, \tau\})$ where U is the maximal link capacity; note that in networks with high capacity and uncapacitated networks, we can simply set $U = \tau$.

3.2. Matching Algorithms (RS + MA + NI and MA + NI + BW)

This section presents faster algorithms to solve the two problem variants RS + MA + NI and MA + NI + BW which can also be solved with the flow approach introduced above. In general, we refer to the algorithms presented in this section as matching approaches.

3.2.1. RS + MA + NI

Let us first consider the RS + MA + NI variant. Recall that in this problem, we are given a set of redundant chunks (RS) and a set of nodes at fixed locations. The number of chunk types is larger than the number of nodes (MA), and each node needs to be connected to its chunks as well as to other nodes (NI). Our goal is to minimize the resource footprint F , consisting of the bandwidth reservations in the access network and the inter-connect.

Algorithm. Due to Observation 1, RS + MA + NI degenerates to RS + MA. In order to solve the RS+MA problem variant, we construct a bipartite graph between the set V of nodes and the set of chunks. Concretely, we clone

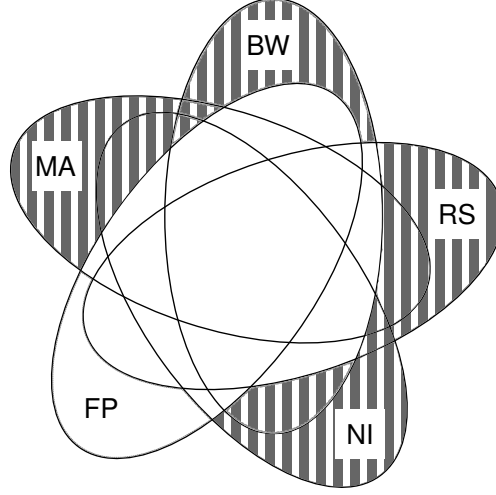


Figure 4: Variants solved by matching approaches.

each node m times, as each node needs to process m chunk types, and we collect all copies of a given chunk type in a single “super-node”. We connect each node to all chunk types using the *lowest hop count* to one of the copies as the cost metric (the link weight). On the resulting bipartite graph, we can now compute a *Minimum Weight Perfect Matching* [17]: the resulting matching describes the optimal assignment of chunks to nodes.

Example. Before analyzing our algorithm, let us consider a small example. Figure 5 illustrates an instance where two nodes are cloned into $m = 2$ nodes each, resulting in a total of four nodes in the matching problem representation. The two replicas of each chunk type are aggregated into a single chunk type vertex c_j in the matching problem; this gives a total of four chunk type vertices in the matching graph. The costs on the links between all clones of a specific vertex and a chunk type are set to the minimum distance. We can observe this for instance at the edges connecting the two clones of v_1 to c_2 : both weights are 0.

Analysis. The correctness of our algorithm follows from the construction and the optimal solution of the minimum matching. The runtime consists of two parts: the construction of the matching graph and the actual matching computation. The constructed graph consists of $m \cdot n_V \cdot \tau$ many edges, and for each edge we need to compute its cost, i.e., the shortest distance which in a tree we can compute in time n_S ; thus, the overall construction time is $\mathcal{O}(n_S \cdot \tau^2)$. The state of the art algorithm to compute matchings are based

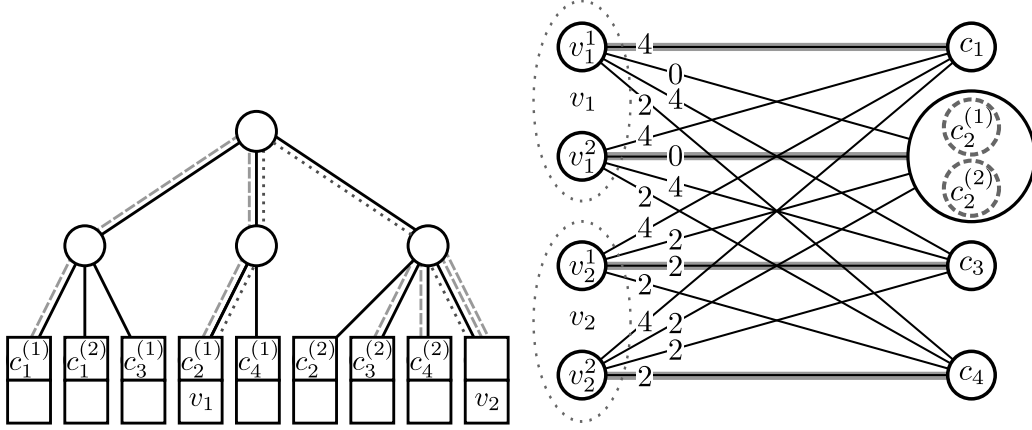


Figure 5: The RS + MA problem on the *left* is converted into a matching problem on the *right*. Since each node has to process two chunks, the nodes are replicated in the matching representation. The two replicas of each chunk type are represented by a single node, and all edges connecting to this node have a weight according to the shorter distance to one of the replicas. This is visualized for c_2 .

on scaling techniques [13]. The runtime translates to $\mathcal{O}(\tau^{5/2} \cdot \log(\tau \cdot n_S))$; recall that $\tau = m \cdot n_V$.

3.2.2. Faster MA + NI and MA + NI + BW

We now show that we can solve MA + NI even faster, by exploiting locality. Moreover, we will show that we can even solve MA + NI + BW problem variants by simply verifying feasibility. In the following, due to Observation 1, we can focus on the MA resp. MA + BW problem.

We first introduce the following definition.

Definition 1 (Local Assignment (LA)). We define an assignment μ to be local in a specific subtree T' , iff μ assigns the maximum number of chunks in the subtree to nodes in the same subtree. We define μ to be local when it is local with respect to all possible subtrees of the substrate network.

Example. Figure 6 illustrates the concept of local assignment: The closest chunk to v_2 is c_1 , and the closest node to c_1 is v_2 . However, a subtree T' exists such that $v_1 \in T'$ and $c_1 \in T'$, but $v_2 \notin T'$. Therefore, a local assignment cannot assign c_1 to v_2 .

We will see later that optimal solutions to MA have a local assignment. We exploit this in our algorithms described in the following.

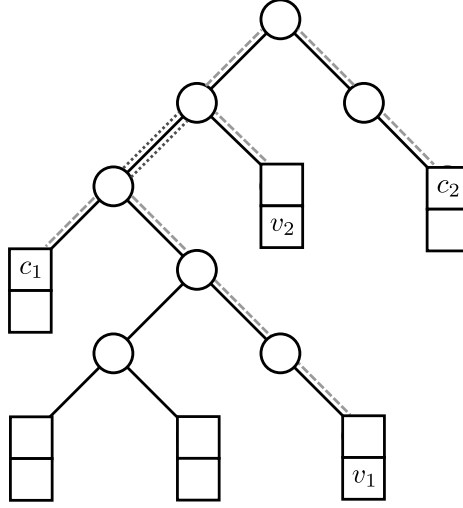


Figure 6: Illustration of local assignment: The dashed lines indicate bandwidth allocations, which occur independently of the chosen assignment. The dotted lines indicate bandwidth allocation which occur only if c_2 is assigned to v_1 .

Algorithm. Our proposed algorithm for MA proceeds in a bottom-up fashion, traversing the substrate network T from the leaves toward the root. For each subtree T' , we maintain two sets S_1, S_2 in order to match unmatched chunks S_1 in the subtree T' to unmatched nodes S_2 in T' . Both sets are initially empty.

We first process all the leaves, in an arbitrary order; subsequently, we process arbitrary inner vertices of T , whenever all their children have been processed. We process any leaf ℓ by adding any nodes or chunks which are located on ℓ to the corresponding sets S_1 and S_2 . A non-leaf vertex u is processed in the following way: we take the union of the sets of u 's children, i.e., the sets contain the unmatched chunks and nodes in this subtree. For both leaves and inner nodes, whenever both sets are non-empty, we greedily match an arbitrary chunk in S_1 with an arbitrary node in S_2 , and remove them from the sets.

Analysis. On a given vertex u , emptying one of the sets, results in a *local assignment* (cf Definition 1) in the subtree rooted at u . The bottom-up strategy ensures that this works for every subtree in the substrate, rendering the resulting assignment local. The complexity of this construction is low: For each vertex in the substrate graph, we build the union of the children's sets, and since each vertex can only be the child of one vertex, the amortized

runtime per vertex is constant; and hence the overall runtime $\mathcal{O}(n_S)$. The sum of all remove operations, is equal to the number of chunk types $\mathcal{O}(\tau)$. Hence the overall complexity of this construction amounts to $\mathcal{O}(n_S + \tau)$.

It remains to prove optimality of such local assignments. We first characterize the bandwidth allocation on uplinks of subtrees.

Lemma 1. *Given an MA problem and a subtree T' containing x chunks and y nodes, the minimal bandwidth allocation of any assignment μ on the uplink of T' is $|x - y \cdot m| \cdot b_1$.*

Proof. In case the number of chunk types equals the processing capacities of the nodes in the given subtree, the bandwidth allocation inflicted by the chunk access network on the uplink can be zero, since we can assign all chunks to nodes in the same subtree. Otherwise, we distinguish between two cases: Recall, that in instances without RS, all chunks have to be processed. In case there are more chunks in the subtree, at least all of the excess chunks have to be transferred to a different subtree, which will inflict costs b_1 per excess chunk on the uplink connecting T' with the remaining parts of T , which will inflict costs b_1 per excess chunk on the uplink of root of T' . Similarly, if the processing capabilities exceed the amount of available chunks, excess chunks from other subtrees will have to be transferred to nodes in the subtree T' , inflicting bandwidth costs of b_1 each. Hence, the minimum bandwidth allocation for the chunk access on the uplink is the difference between the number of chunks and the processing capabilities of the subtree $|x - y \cdot m|$ times the amount of bandwidth needed, for a single transfer b_1 . \square

Theorem 2. *Given an MA + NI problem instance, a feasible assignment μ is optimal iff it is local.*

Proof. Local assignments generate exactly the minimal allocations on all links, as the assignments which generate the minimal bandwidth allocations described in the proof of Lemma 1 are local in the given subtree. Hence each local assignment has to be optimal. A non-local assignment, has at least one subtree, in which it is not local. This subtree will have a higher allocation on the uplink. Since the local assignment has minimal allocations on all other links, the non local assignment has a larger footprint. \square

Combined with a simple postprocessing step, this approach can also solve MA + BW. The central idea of this extension, is that *local* assignments allocate the minimal bandwidth on each individual edge. In consequence, each bandwidth constraint which is lower than the allocation of a

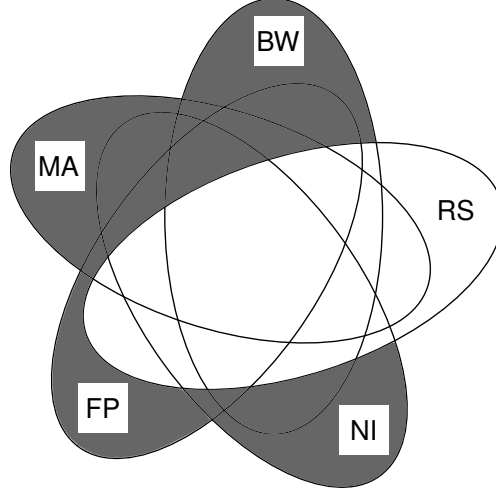


Figure 7: Variants solved by dynamic programming approach.

local assignment on one link, renders the problem infeasible. Hence, it is sufficient to temporarily omit the bandwidth limitations, compute an optimal assignment for an MA instance, and verify that the resulting allocations do not violate any capacities. The postprocessing step scales linearly with the number of edges in the substrate graph.

3.3. Dynamic Programming (MA + FP + NI + BW)

We now show how to solve the MA + FP + NI + BW problem variant in polynomial time. Note that this problem variant requires to find a tradeoff between the desire to place nodes as close as possible to each other (in order to minimize communication costs), and the desire to place nodes as close as possible to the chunk locations.

Example. Figure 8 shows an example: one extreme solution is to minimize the distance between chunks and nodes, see mapping π_1 in Figure 8 (*left*): the four nodes are all collocated with chunks, resulting in a zero-cost chunk access network. As a result, the paths between the individual nodes are longer than in alternative node placements: each node has a distance of two hops to one other node, and four hops to two other nodes. Hence the resulting allocations for the node interconnect sum up to $20 \cdot b_2$.

Figure 8 (*right*) shows a different node mapping π_2 , which seeks to minimize the communication costs between the nodes, and places all nodes in one subtree. The distance between all nodes is two, which results in a total

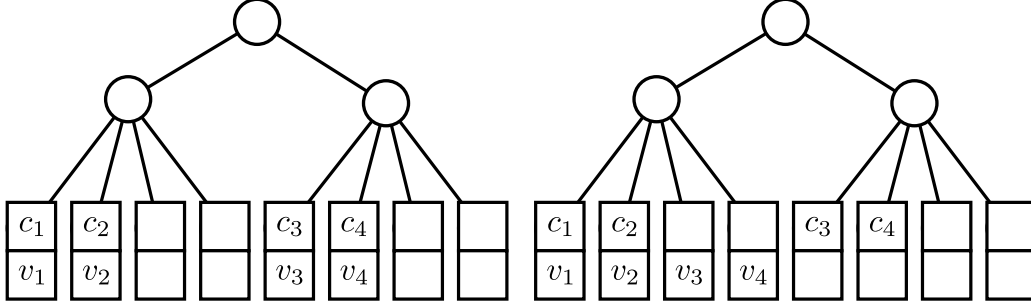


Figure 8: Two different node placements for the same substrate graph and chunk locations. For $b_1 = b_2$, both solutions have an identical footprint. In other cases, one solution outperforms the other.

bandwidth allocation of $12 \cdot b_2$ for the interconnect. However, this reduced price comes at additional costs in the access network: c_3 and c_4 have to be communicated to v_3 and v_4 , which requires a total bandwidth allocation of $8 \cdot b_1$.

Basic ideas. Our proposed approach is based on dynamic programming, and leverages the *optimal substructure property* of MA + FP + NI + BW: as we will see, optimal solutions for subproblems (namely subtrees) can efficiently be combined into optimal solutions for larger problems. Indeed, the MA + FP + NI + BW problem exhibits such a structure, and we show how to exploit it to compute efficient embeddings, even in scenarios where multiple chunks need to be assigned to flexibly placeable nodes.

For ease of presentation we will transform the substrate network T into a binary tree, using binarization: we clone every higher-degree node, iteratively attaching additional clones as right children and original children as left descendants.

As usual in dynamic programs, we define, over the structure of the tree, a recursive formula f for the minimal cost solution *given* any possible number of nodes embedded in a given subtree. The actual set does not matter, due to symmetry arguments. Our approach is to evaluate this function in a bottom-up manner. To finally compute the actual optimal embedding, we traverse the computed minimal-cost path backwards (according to the optimal values found for f during the bottom-up computation).

Concretely, the first argument to function f is a subtree T' , containing a given number of chunks $y(T')$, and the second argument is the number of nodes to be embedded in the subtree. Function f is evaluated in a bottom up manner. We initialize the function at each leaf ℓ , by $f(T_\ell, x) = \infty$ for

all numbers of nodes x which are larger than the server capacity $cap(\ell)$; to calculate $f(T_\ell, x)$, for $x \leq cap(\ell)$, we compute the bandwidth allocation on the uplink of T_ℓ , referred to by the function $bw(T_\ell, x)$: $bw(T_\ell, x) = b_1 \cdot |x - y(T_\ell)| + b_2 \cdot (n_V - x) \cdot x$, which accounts for the bandwidth allocation on the uplink of T_ℓ . The first term represents the required bandwidth for the communication between the x nodes on ℓ , and the $n_V - x$ nodes in the remaining parts of the substrate network. The second term represents the bandwidth, which is necessary to transport the chunks from their location to the node which should process the data (see Lemma 1 for more details).

After initialization, we proceed to compute f for non-leaf nodes in a bottom-up manner: We split the x nodes into two positive integer values, and we put r on the right and $x - r$ on the left subtree. That is, we take the optimal cost (given recursively) of placing r nodes in the right subtree $RI(T')$ of T' and $x - r$ nodes in left subtree $LE(T')$ of T' . Given the cheapest combination, we add the bandwidth requirements on the uplink of T' to generate the overall costs for placing x nodes in T' . Therefore, $f(T', x) = \min_{0 \leq r \leq x} \{f(LE(T'), x - r) + f(RI(T'), r)\} + bw(T', x)$. Again, we set $f(T', x)$ to infinity if the required bandwidth bw exceeds the capacity cap of the uplink of T' .

Analysis. The correctness and optimality of our dynamic program is due to the decoupling of the costs induced by the tree structure of T and the substructure optimality property. The substructure optimality follows from the observation that costs can be accounted on the uplink, and the fact that we check each possible node distribution. For each substrate vertex (n_S many) we have to check the cost of all possible splits, resulting in an overall complexity of $\mathcal{O}(n_S \cdot n_V^2)$. The runtime to binarize T is asymptotically negligible.

3.4. Simple Problems

For the sake of completeness, we also observe that there are several problems which allow for a trivial solution. Concretely, problems with FP plus any combination of RS and BW (but without MA and NI) can easily be solved by mapping nodes to chunk locations. Figure 9 shows a Venn diagram of the trivial property combinations.

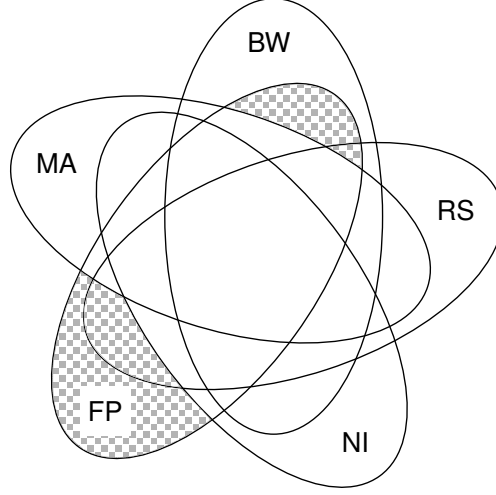


Figure 9: Trivially solvable problem variants.

4. NP-Hardness Results

We have seen that even problems with multiple dimensions of flexibility can be solved optimally in polynomial time. This section now points out fundamental limitations in terms of computational tractability. In particular, we will show that problems become NP-hard if flexibly placeable nodes (FP) have to be assigned to one of multiple replicas (*RS*), either with multiple chunks per node (MA in Section 4.2) or with communication among nodes (NI in Section 4.3). Both results hold even in uncapacitated networks, and even in small-diameter substrate networks (namely two- or three-level trees [2]). The hardness of FP+RS+MA and FP+RS+NI imply the hardness of four additional, more general models, as summarized in Figure 10:

$$\begin{array}{ccc}
 \mathbf{RS+FP+NI} & \implies & \mathbf{RS+FP+NI+BW} \\
 \Downarrow & & \Downarrow \\
 \mathbf{RS+MA+FP+NI} & \implies & \mathbf{RS+MA+FP+NI+BW} \\
 \Uparrow & & \Uparrow \\
 \mathbf{RS+MA+FP} & \implies & \mathbf{RS+MA+FP+BW}
 \end{array}$$

Figure 10: The NP-hardness of 2 variants implies the hardness of 4 other variants.

4.1. Introduction to 3D Perfect Matching

Both the hardness of $\text{FP} + \text{RS} + \text{MA}$ and $\text{FP} + \text{RS} + \text{NI}$ are shown by a reduction from the NP-complete problem of *3D Perfect Matching* [9], which we can see as a generalization of bipartite matchings to 3-uniform hypergraphs. We will refer to this problem by 3-DM, and for completeness, review it quickly: 3-DM is defined as follows. We are given three finite and disjoint sets X , Y , and Z of cardinality k , as well as a subset of triples $T \subseteq X \times Y \times Z$. Set $M \subseteq T$ is a 3-dimensional matching if and only if, for any two distinct triples $t_1 = (x_1, y_1, z_1) \in M$ and $t_2 = (x_2, y_2, z_2) \in M$, it holds that $x_1 \neq x_2$, $y_1 \neq y_2$, and $z_1 \neq z_2$. Our goal is to decide if we can construct a $M \subseteq T$ which is *perfect*, that is, a subset which covers all elements of $X \times Y \times Z$ exactly once.

4.2. Multi-Assignments are hard ($\text{FP} + \text{RS} + \text{MA}$)

Our proof that $\text{FP} + \text{RS} + \text{MA}$ is NP-hard is based on the following main ideas. We encode a 3-DM instance as an $\text{FP} + \text{RS} + \text{MA}$ instance as follows:

- For every element in the universe $X \cup Y \cup Z$, we create a chunk type. Intuitively, in 3-DM, each element must be covered, which corresponds to the requirement of $\text{FP} + \text{RS} + \text{MA}$ that each chunk type is processed.
- We will encode each triple as gadget with three leaves in a substrate tree T . The three leaves are close to each other in T , and the placement of chunk replicas in $\text{FP} + \text{RS} + \text{MA}$ corresponds to the elements of the triples in these leaves.
- The node placement will correspond to the choice of triples, independently of which leaf the node is mapped to. A node will process its collocated chunk, as well as the chunks in other two leaves of the same gadget.
- In order to turn the optimization problem into a decision problem, we will use a cost threshold Th . The cost threshold will be met by all assignments which assign all three chunks of each triple to a node which is collocated with one of the chunks. Assignments which connect a chunk to a node in a different triple, will have a larger footprint, and are considered to be infeasible.

Construction. Given an instance I of 3-DM in which k triples have to be chosen, we construct an instance I' of $\text{FP} + \text{RS} + \text{MA}$ as follows:

- *Tree Construction:* We create a tree consisting of a root, and for each triple, we create a gadget which we directly attach as child of the root. The gadget is of height 2, and has the following form: The gadget of each triple consists of an inner node (a router) and three leaves.
- *Chunks and chunk replicas:* For each element in X , Y and Z , we create a chunk type ($3 \cdot k$ in total). Every gadget contains three chunk replicas, corresponding to the elements of the triple. Each leaf in a gadget, contains exactly one replica.
- *Other properties:* We set the number of to-be-embedded nodes to k , b_1 to 1, and the number of chunk slots in each node to the multi-assignment factor $m = 3$. We use a threshold $Th = 4 \cdot k$.

Example. Figure 11 shows an example of our construction: An instance I of 3-DM is given: The disjoint sets X , Y and Z have a cardinality $k = 2$. We will refer to the two elements in X as x_1 and x_2 , and use the same notation for the other two sets. T contains the three triples (x_1, y_1, z_1) , (x_2, y_1, z_2) , and (x_2, y_2, z_2) . The goal of 3-DM is to find a subset $M \subseteq T$, which contains each element in each of the three sets exactly once. This instance only has one solution: $M = \{(x_1, y_1, z_1), (x_2, y_2, z_2)\}$.

To construct the corresponding instance I' of $FP + RS + MA$, we create a gadget for each triple in T . For each variable which occurs in a triple, the corresponding gadget contains a chunk of the type of the variable. The triple (x_2, y_1, z_2) of the instance is represented by the middle gadget in Figure 11. The objective of I' is to spawn $k = 2$ nodes, with the smallest possible footprint. If the total footprint is $\leq 2 \cdot 2 \cdot k$, we can construct a solution to I from the solution to I' . The footprint consists of the costs which occur when a node is embedded in a gadget, and the three chunks of that gadget which are assigned to that node: one of the chunks is collocated with the node, the other two have to be transferred via two hops, inflicting unitary costs on each hop.

Correctness. Given these concepts, we can now show the computational hardness.

Theorem 3. $FP + RS + MA$ is NP-hard.

Proof. Let I be an instance of 3-DM and let I' be an instance of $FP + RS + MA$ constructed as described above. We prove that I' has a solution of cost $\leq Th$ if (\Rightarrow) and only if (\Leftarrow) I has a matching of size k .

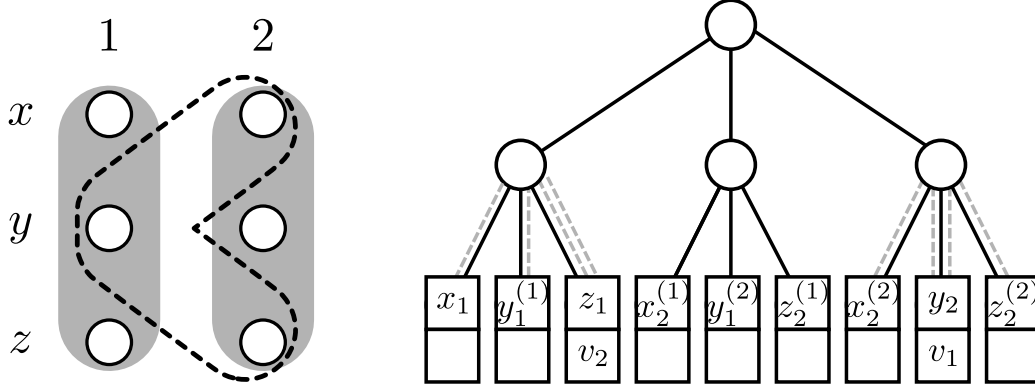


Figure 11: *Left:* A 3-DM instance with three triples: (x_1, y_1, z_1) , (x_2, y_1, z_2) , and (x_2, y_2, z_2) . The solution is indicated by the grey triples; the dashed triple is not used for the solution. *Right:* The corresponding problem and solution of FP + MA + RS.

(\Rightarrow) Let us take a solution to 3-DM. We place a node in every gadget that corresponds to the chosen triples. In each of the corresponding gadgets, we match every chunk to the node in this gadget. This solution has cost exactly Th . As every element of the universe is covered, every chunk type is processed.

(\Leftarrow) Let us take a solution to FP + RS + MA of cost $\leq Th$. We choose triples that correspond to gadgets where there are nodes. Since all chunks are processed, every element of X , Y and Z is matched. Each node must process chunks that correspond to the triple, otherwise the cost must be larger than Th (high costs for chunk transportation). \square

4.3. Inter-connects are hard (FP + RS + NI)

Next, we prove that the joint optimization of node placement and replica selection is NP-hard if an inter-connect has to be established between nodes. In our terminology, this is the FP + RS + NI problem.

The proof is similar in spirit to the proof of FP + RS + MA, however, we modify the construction to account for the absence of MA: we choose a high value for b_1 , such that nodes will be directly collocated with their assigned chunks. We leverage the fact that any solution which does not assign 0 or 3 chunks to each gadget, will have higher communication costs.

Construction. Let I be an instance of 3-DM. We will create an instance I' for FP + RS + NI as follows:

- We will construct the same tree as in previous reduction with chunk replicas placed in the same way.
- The communication cost in the inter-connect is set to $b_2 = 1$.
- The number of nodes (virtual machines) is $n_V = 3 \cdot k$, where k is the set cardinality.
- Only solutions which place a node in each leaf of k gadgets, can be converted into solutions for the 3-DM problem. We use the cost threshold $Th = 6 \cdot k + 18 \cdot (k - 1) \cdot k$, to verify whether a solution achieves this, transforming FP + RS + NI into a decision problem. A detailed explanation of this value can be found in the proof of Theorem 5.
- We set the access cost b_1 to a chunk replica to a high value W . This will force nodes to be colocated with the replica. One example of sufficient (and polynomial but not necessarily minimal) W is the value of the threshold $Th + 1$. Any solution not assigning chunks to colocated nodes, have cost $> Th$: communicating a chunk inflicts costs $W = Th + 1$ over every link.

We focus on instances with unit server capacities.

Proof of correctness. Intuitively, in order to minimize embedding costs, nodes should be placed on near-by replicas. We use the following helper lemma.

Lemma 4. *In every valid solution of I' of cost $\leq Th$, each gadget falls in one of two categories: k gadgets have exactly 3 nodes, and $n - k$ gadgets remain empty.*

Proof. Since W is large enough, the $3 \cdot k$ nodes have to be placed directly on different chunks, resulting in 0 costs for the access network. Consider any pair of nodes communicating over the inter-connect; due to our construction, the communication cost for each such pair is either 2 hops (if they belong to the same gadget) or 4 hops (if they belong to different gadgets). The lemma then follows from the observation that Th is chosen such that it is never possible to distribute nodes among more than k gadgets. \square

Theorem 5. FP + RS + NI is NP-hard.

Proof. Let I be an instance of 3-DM and let I' be an instance of FP+RS+NI constructed as described above. We prove that I' has solution of cost $\leq Th$ if (\Rightarrow) and only if (\Leftarrow) I has a solution.

(\Rightarrow) In order to compute a solution for I' given a solution for I , we proceed as follows. Given an exact covering set of triples $S = \{t_1, t_2, \dots, t_k\}$, we place three nodes in each gadget that corresponds to every triple of S . Chunks are matched to the nodes which are located on the same server.

The solution has the following cost: (1) the communication cost inside a gadget is $2 \cdot \binom{3}{2}$, as every pair contributes two hops; (2) the communication cost from each gadget to all other gadgets is $4 \cdot 3 \cdot 3 \cdot (k-1)/2$, where the factor 4 is for the communication over 4 hops, the factor 3 corresponds to the number of nodes per gadget, and $3 \cdot (k-1)$ is the number of nodes in remote gadgets; as we count each pair twice, we need to divide by two in the end. Summing up over all k gadgets, we get exactly Th .

(\Leftarrow) Given a solution for I' , we can exploit Lemma 4 to construct a solution for I . We know that in any solution of cost at most Th , k gadgets contain exactly 3 nodes. These gadgets correspond to a valid 3D Perfect Matching: exactly one replica of every chunk type is processed and hence every element is covered exactly once. \square

5. A Detailed Study of Replica Selection Hardness

We have seen that replica selection flexibilities can render embeddings computationally hard. We will now provide a more detailed look at this hardness result and explore the minimal requirements for rendering replica selection hard. In particular, we will show that already two replicas for each chunk type are sufficient to introduce intractability.

Across the next sections, we use the following notation:

1. $n := |X|$ (remember that $|X| = |Y| = |Z|$)
2. e – an element of the universe ($e \in X \cup Y \cup Z$)
3. T – set of triples
4. t – a triple ($t \in T$)
5. T_e – set of triples that contain element e
6. $\deg(e) := |T_e|$
7. $\{e_X(t), e_Y(t), e_Z(t)\}$ – elements of triple t
8. V – set of nodes to be spawned in embedding instance

5.1. Two Replicas without Bandwidth Constraints

We now show that the 2-replica selection problem is even NP-hard without capacity constraints. In particular, we consider the problem variant $\text{FP} + \text{RS}(2) + \text{MA}(4)$ with at most two replicas of each chunk type and assignment factor four. There are no capacity constraints on links.

Construction.

Let's take any instance $I_{3\text{DPM}}$ of 3-DM and create a $\text{FP} + \text{RS}(2) + \text{MA}(4)$ instance I_{VCEMB} as follows.

Chunks. We construct three types of chunks. The first type corresponds to covers of elements by triples, with two replicas each. The other two types are chunk types with one replica only, therefore called *unique*. We construct two types of unique chunks, distinguished by a different role in the construction. For unique chunks we simply annotate the chunk type with chunk replica. Thus, the constructed instance uses at most two replicas of each chunk type.

1. For each triple $t \in T$, we construct 3 chunk types, with two replicas each. We construct different chunk types for each triple t , which contain element e . We refer to those replicas by $ch_1(e, t)$ and $ch_2(e, t)$.
2. We construct n additional chunk types named u_1, \dots, u_n , with one replica each.
3. For each element $e \in X \cup Y \cup Z$, we construct additional $3 \cdot (\deg(e) - 1)$ chunks, with one replica each. We call this set \mathcal{U}_e .

Tree. We construct the following tree.

1. The physical network consists of two subtrees connected to the root: A *Matching Subtree* and a *Cover Subtree*. The *Matching Subtree* consists of $|T|$ *Triple Gadgets*, one per each triple $t \in T$ and n *Unique Gadgets*. The *Cover Subtree* consist of n *Element Gadgets*, one for each element $e \in X \cup Y \cup Z$.
2. *Triple Gadget* consists of four vertices: three leaves and the root of the gadget.
3. *Unique Gadget* consists of two vertices: the leaf and the root of the gadget. Note that we construct nodes not only to keep the tree balanced, but also to keep leaves of *Unique Gadgets* far from leaves of other *Unique Gadgets*.
4. *Element Gadget* of element e has a structure that depends on the number of triples that cover e . The *Element Gadget* consists of the root, and $4 \cdot (\deg(e) - 1) + 1$ leaves.

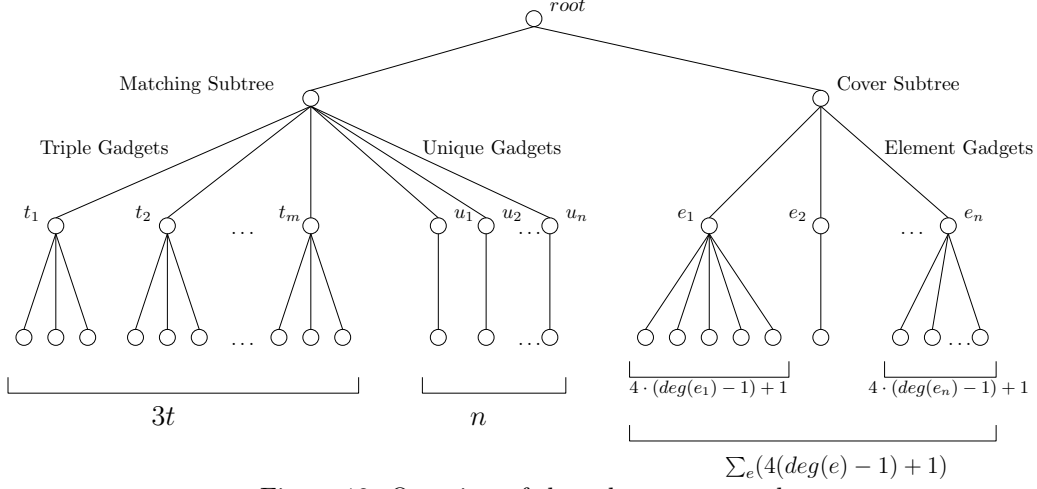


Figure 12: Overview of the substrate network

Chunk Placement. The chunks are placed as follows:

1. *Chunks in matching subtree:* In *Triple Gadget* of triple t we put three replicas: $ch_1(e_X(t), t), ch_1(e_Y(t), t), ch_1(e_Z(t), t)$, one per each leaf.
2. *Chunks in unique subtree:* We place replicas u_1, \dots, u_n at the leaves of *Unique Gadgets*.
3. *Chunks in element gadget:* In leaves *Element Gadget* of element e we put replicas $ch_2(t, e)$ for each $t \in T_e$. Additionally, we put replicas \mathcal{U}_e .

Other properties of the instance.

1. *Multiple assignment:* We set the assignment factor (number of processed chunks by each node) to four.
2. *Number of nodes:* We allow to spawn $|V| := n + \sum_e (\deg(e) - 1)$ nodes.
3. *Threshold:* We set the following threshold for cost of a solution to be feasible. $Th := 8 \cdot n + 6 \cdot \sum_e (\deg(e) - 1)$

Reduction. We can make the following observations.

Observation 2. *The construction fulfills the following properties:*

1. *There is exactly one replica of each chunk type in the Matching Subtree.*
2. *There is exactly one replica of each chunk type in the Cover Subtree.*
3. *There are at most two replicas of each chunk type in the substrate network.*

For the reduction we proceed as follows:

1. Let's take any instance $I_{3\text{-DM}}$ of 3-DM and produce an instance I as described in the construction section.
2. Let's take any feasible solution $S_{3\text{-DM}}$ to $I_{3\text{-DM}}$ and place n (remember that $n = |X| = |Y| = |Z|$) nodes at any leaf of the *Triple Gadget* chosen in $S_{3\text{-DM}}$ solution. We match each such node to three chunks that are put in its *Triple Gadget* (it is co-located with one of the chunks), and we match it to an arbitrary unmatched chunk in any *Unique Gadget*.
3. We put $\deg(e) - 1$ nodes in each *Element Subtree*. We match all remaining chunks from this *Element Subtree* to those nodes in an arbitrary way.
4. This solution has cost $\leq Th$ (easy to see by summing up the costs of transporting chunks to each node).
5. The feasibility of the constructed solution follows from the fact that each chunk type was processed; each node processes four chunks.

Proof of correctness.

Let $\widehat{\text{cost}}(u) := u \cdot 12 + (t - u) \cdot 4 + (4 \cdot |V| - |V| - 4 \cdot u - (t - u)) \cdot 2$, and $\widehat{\text{cost}}'(u) := 4 \cdot t + (4 \cdot |V| - |V| - t) \cdot 2$.

Lemma 6. *For each $u \in \{1, \dots, t - 1\}$ we have that $\widehat{\text{cost}}(u) > Th$ and $\widehat{\text{cost}}'(u) > Th$.*

Proof. We have that $\widehat{\text{cost}}(u + 1) \geq \widehat{\text{cost}}(u)$ and $\widehat{\text{cost}}'(u + 1) \geq \widehat{\text{cost}}'(u)$, for $1 \leq u \leq t - 1$. It is easy to verify that $\widehat{\text{cost}}(1) > Th$ and $\widehat{\text{cost}}'(1) > Th$. \square

Lemma 7. *Take any feasible solution S to the instance I of $\text{FP} + \text{RS}(2) + \text{MA}(4)$ (as constructed above). If the cost of S is at most Th , then no node is spawned in the *Unique Subtree*.*

Proof. For the sake of contradiction, let us assume a feasible solution S with at least one node spawned in *Unique Subtree*. We show that in this case, the cost of solution S is greater than Th . Let ℓ be the number of nodes spawned in the *Unique Subtree*. We know that $1 \leq \ell \leq |T|$. In S we have exactly $4 \cdot |V|$ chunk transportations, incurring cost 0, 2, 4 or 6 (the tree has an edge-height of 3). At most $|V|$ transportations are of cost 0. Note that the leaves of the *Unique Subtree* are separated from other leaves of the tree by at least 4 edges. The cost of chunk transportation to nodes

$v \in \sigma$ is at least 12. The chunks in *Unique Subtree* are unique, therefore the solution transports $|T| - \ell$ chunks to nodes outside the *Unique Subtree*, incurring cost 4 for each chunk. Therefore $\text{cost}(S) \geq \widehat{\text{cost}}(\ell)$, and by Lemma 6 we conclude that $\text{cost}(S) > Th$. \square

Lemma 8. *Given any feasible solution S to a given instance I of $\text{FP} + \text{RS}(2) + \text{MA}(4)$ (as constructed above). If the cost of S is at most Th , then exactly n nodes are spawned in the Matching Subtree.*

Proof. Let δ be the set of nodes spawned in *Matching Subtree*. Let us assume the contrary, namely that $|\delta| \neq n$. First, we use Theorem 7 to restrict the placement of nodes in *Unique Subtree*. Then we consider two cases:

1. **Case $\ell \leq n$:** There are at least $u := 4 \cdot (n - \ell)$ chunks in the *Matching Subtree* that are not processed in the *Matching Subtree*, each incurring transportation cost of 6. From the structure of the substrate network and placement of chunks we know that $\text{cost}(S) \geq \widehat{\text{cost}}'(\ell)$, and by Lemma 6 we conclude that $\text{cost}(S) > Th$.
2. **Case $\ell > n$:** We use the fact that there are not enough nodes in *Cover Subtree*, and we need to transport over 6 hops at least 3 unique chunks for each node missing from *Cover Subtree*.

Theorem 9. $\text{RS}(2) + \text{MA} + \text{FP}$ is NP-hard.

Proof. Let's take an instance I of 3DPM and construct an instance I' of $\text{RS}(2) + \text{MA} + \text{FP}$ in the way described in the construction section. We show that I' has a solution of cost $\leq Th$ if and only if $I \in \text{3DPM}$ (there exists a perfect 3D matching).

(\Leftarrow) Let's take any feasible solution S to I . We construct a solution S' to I' in the following way:

1. We place n nodes in n *Triple Gadgets* (one per gadget) that correspond to triples in S . We match each such node to chunks in the gadget it is placed, as well as one arbitrary chunk in *Unique Subtree*.
2. In each *Element Gadget* that corresponds to element e , we place $\deg(e) - 1$ nodes and match them to arbitrary chunks in this gadget, which are not yet matched in any *Triple Gadget*.

We can observe that every chunk type was processed. By simple calculations we see that S' indeed has cost $\leq Th$.

(\Rightarrow) Let's take any feasible solution S' to I' . We construct the solution S to I in the following way: We call the *Triple Gadget* *active*, if they contain a node at any leaf. We call active node in *Triple Gadgets* *active nodes*. We construct a 3DPM solution from triples that correspond to active *Triple Gadgets*. We observe the following properties of S :

1. From Lemmas 7 and 8, we know that exactly n *Triple Gadgets* are *active*.
2. In S , only one node is spawned in an active *Triple Gadget*.
3. Each *active node* v processes the 3 chunks that are placed in v 's *Triple Gadget*, as well as one chunk in an *Unique Gadget*.
4. Every chunk type is covered.
5. In each *Element Gadget* for element e , one chunk instance of set $t(e)$ is not processed. Let's call this chunk instance $\gamma(e)$, and let's call $\gamma = \cup_e \gamma(e)$. Note that $|\gamma| = n$. The set γ is covered by *active nodes*

From above observations we conclude that M_S is indeed feasible. \square

5.2. Two replicas without multiple assignment

We now show that RS(2) + FP + NI + BW is even NP-hard without multiple assignment.

Construction.

Chunk Types. We construct the following chunk types: For each element $e \in X \cup Y \cup Z$, we construct $\deg(e)$ chunk types with two replicas. Additionally, we construct $\max\{3 \cdot |T| + 3 \cdot n + 1, \sum_e (2 \cdot \deg(e) - 1)\}$ chunk types called *unique chunks*. We refer to the set of unique chunks by U .

The substrate network.

1. The physical network consists of three subtrees connected to the root: A *Matching Subtree*, a *Cover Subtree*, and a *Unique Subtree*. In the *Matching Subtree* we put $|T|$ *Triple Gadgets* (remember that $|T|$ is the number of triples in 3DPM instance). *Cover Subtree* consist of n element gadgets.
2. The *Unique Subtree* consist of $|U|$ leaves, and two middle nodes: a lower and an upper middle node. Note that this is different from RS(2) + FP + MA(4) NP-completeness proof, where *Unique Subtree* was placed in the *Matching Subtree*.

3. *Triple Gadget*: For each triple, we create a subtree consisting of four vertices: three leaves and one triple root. We attach the root of the triple to the root of the matching subtree.
4. *Element Gadget*: For each element $e \in X \cup Y \cup Z$, we construct a subtree consisting of the root of the element (attached to the root of the cover subtree), and $4 \cdot (\deg(e) - 1) + 1$ leaves.

Chunk placement. The chunks are placed as follows:

1. *Chunks in matching subtree*: For each triple t we put three chunks at the leaves of the corresponding *Triple Gadget*, $e_1(t), e_2(t), e_3(t)$.
2. *Chunks in unique subtree*: We place unique chunks U at the leaves of *Unique Subtree*.
3. *Chunks in element gadget*: For each element $e \in X \cup Y \cup Z$, we place the chunks $t(e)$ at the leaves of each element subtree.

Bandwidth constraints. We use bandwidth constraints of the form $\text{BW}(k) := k \cdot (|V| - k)$. Namely, we set the bandwidth constraints of an uplink of an *Element Gadget* for each element e to $\text{BW}(\deg(e) - 1)$, the bandwidth of an uplink of a *Matching Subtree* to $\text{BW}(n)$, and an uplink of a *Cover Subtree* to $\text{BW}(\sum_e (\deg(e) - 1))$.

The threshold value and other properties of the instance. We set the cost threshold for any solution to the following value:

$$\begin{aligned}
Th = & 2 \cdot (3 \cdot |V| + \sum_e (\deg(e) - 1)) && \text{(over 2 hops)} \\
& + 4 \cdot (n \cdot (3 \cdot (3 \cdot |V| - 3))/2) && \text{(over 4 hops in Matching Subtree)} \\
& + 4 \cdot (\sum_e ((\deg(e) - 1) \cdot (\sum_{f \neq e} \deg(f) - 1)/2)) && \text{(over 4 hops in Cover Subtree)} \\
& + 6 \cdot (3 \cdot |V| \cdot \sum_e (\deg(e) - 1)) && \text{(between Matching Subtree and Cover Subtree)} \\
& + |U| \cdot (|U| - 1)/2 && \text{(inside Unique Subtree)} \\
& + |U| \cdot (3 \cdot n + \sum_e (\deg(e) - 1)) && \text{(Unique Subtree to other nodes)}
\end{aligned}$$

We set b_1 , the cost of chunk transportation to $Th + 1$ (so that no chunk transportation happens in any feasible solution), $b_2 = 1$, and we host only one node per machine. We set the number of machines to spawn to: $|V| := 3 \cdot n + \sum_e (\deg(e) - 1) + |U|$.

Properties of the substrate network.

Lemma 10. Assume we have a RS(2) + FP + NI + BW instance I with a subtree T' with l leaves and the bandwidth capacity on uplink of T' is $\text{BW}(k)$. Assume that no chunk transportation is allowed ($b_1 = \infty$, so every node must be collocated with the chunk it processes in every feasible solution), and $b_2 = 1$. Then in any feasible solution the number of nodes spawned in T (we name it s) satisfies $s \leq k \vee n - s \leq k$, and $s \leq l$.

Proof. It holds that $s \leq l$ as we cannot spawn more nodes than leaves. The bandwidth allocation on the uplink of T' is $\text{uplink}(s, T) := s \cdot (|V| - s)$, as no chunk transportation is allowed ($b_1 = \infty$), and every node in T has to communicate over T' 's uplink with nodes spawned outside of T' . Therefore, in every feasible solution we have: $\text{uplink}(s, T') \leq \text{BW}(k)$. Let's define the remaining bandwidth on the uplink of T' $\text{remainBw}(s) := \text{BW}(k) - \text{uplink}(s, T') = s^2 - s \cdot |V| - k^2 + k \cdot |V|$. Every feasible solution fulfills $\text{remainBw}(s) \geq 0$, which is true for $s \leq k \vee |V| - s \leq k$ (follows from the properties of the quadratic function). \square

Next, we show how to precisely control the number of nodes in the constructed subtree.

Observation 3. In every feasible solution we have exactly $|U|$ nodes spawned in a Unique Subtree (no chunk transportation is allowed, and every chunk type must be processed).

Lemma 11. Consider an instance I of 3DPM. We construct the RS(2) + FP + NI + BW instance I' as described above. Then we have that in I' :

1. The number of nodes spawned in a Matching Subtree is $3 \cdot n$.
2. The number of nodes spawned in a Cover Subtree is $\sum_e (\deg(e) - 1)$

Proof. From Observation 3 we know that we have $|U|$ nodes in the Unique Subtree. Let's refer to the number of nodes spawned in a Matching Subtree by M , and to the number of nodes spawned in Cover Subtree by C . By applying Lemma 10 to Matching Subtree, we know that: $M \leq 3 \cdot n \vee M \geq |V| - 3 \cdot n$. We observe that $|V| - 3 \cdot n$ is greater than the number of leaves in a Matching Subtree. By applying Lemma 10 to the Cover Subtree we know that: $C \leq \sum_e (\deg(e) - 1) \vee C \geq |V| - \sum_e (\deg(e) - 1)$. We observe that $|V| - \sum_e (\deg(e) - 1)$ is greater than the number of leaves in the Cover Subtree. We also know that $|V| = |U| + C + M$. Therefore, by the pigeon-hole principle $C = \sum_e (\deg(e) - 1)$ and $M = 3 \cdot n$. \square

Lemma 12. *Assume an instance I of 3DPM. We construct the $\text{RS}(2) + \text{FP} + \text{NI} + \text{BW}$ instance I' as described above. Then we have that in I' the number of nodes spawned in Triple Gadget of element e is $\deg(e) - 1$.*

Proof. Let's call the number of nodes spawned in the *Triple Gadget* of element e x_e . From Lemma 10, we know that $x_e \leq \deg(e) - 1 \vee x_e \geq |V| - \deg(e) + 1$. We observe that $|V| - \deg(e) + 1$ is greater than the number of leaves of the gadget, which is $\deg(e)$. From Lemma 11, we know that the number of nodes spawned in the entire *Cover Subtree* is $\sum_e (\deg(e) - 1)$. Therefore, by the pigeon-hole principle, we have that $x_e = \deg(e) - 1$. \square

From the above lemmas we know the precise number of nodes spawned in certain parts of the tree. Feasible solutions only differ in the choice of the $\deg(e) - 1$ out of $\deg(e)$ chunks in each *Triple Gadget*, and the placement of nodes in the *Matching Subtree*.

Similar in spirit to the NP-completeness proof of $\text{RS}(2) + \text{MA}(4) + \text{FP}$, we call the *Triple Gadget* active if it contains exactly three nodes. Similarly, we call the *Triple Gadget* inactive if it does not contain spawned nodes, and *partially active* if it has one or two spawned nodes.

Lemma 13. *Consider an $\text{RS}(2) + \text{FP} + \text{NI} + \text{BW}$ instance I . Assume that chunk transportation is not allowed, and $b_2 = 1$. In every feasible solution to I , we have exactly n active Triple Gadgets.*

Proof. Since I is feasible, we know that it has a solution S of cost $\leq Th$. By Lemma 11, we know that there are exactly $3 \cdot n$ spawned nodes in the *Matching Subtree*. Therefore, by the pigeon-hole principle, we know that we have at most n active *Triple Gadgets*. It remains to show that there are no partially active *Triple Gadgets* in the solution of cost $\leq Th$. Using Lemma 12, we conclude that the communication cost of nodes in the *Cover Subtree* is the same for every feasible solution (let's name that cost P). We also know that the communication cost between nodes in *Cover Subtree* and *Matching Subtree* is the same for every feasible solution (let's name it Q). Let's call the would-be cost of communication in the *Matching Subtree*, if there were exactly n active gadgets, R . The threshold value was chosen so that $Th = P + Q + R$. If we have at least one partially active gadget, then the cost of communication in *Matching Subtree* is greater than R , because we increase the number of 4-hop communications by at least one per each partially active gadget in comparison to a solution where we have exactly n active gadgets. \square

The reduction.

Theorem 14. $\text{RS}(2) + \text{FP} + \text{NI} + \text{BW}$ is NP-hard.

Proof. Let's take an instance I of 3DPM and construct an instance I' of $\text{RS}(2) + \text{FP} + \text{NI} + \text{BW}$ in the way described above. We show that I' has solution of cost $\leq Th$ if and only if $I \in \text{3DPM}$ (there exists a perfect 3D matching).

(\Leftarrow) Let's take any feasible solution S to I and produce a solution S' to I' . We show that the cost of S' is indeed $\leq Th$. For each triple t_1, \dots, t_n in S , we put 3 nodes at leaves of triple gadgets corresponding to those triples. In each element gadget (that corresponds to element e), we put $\deg(e) - 1$ nodes. In each element gadget there is only one leaf without the node placed in it: this node contains the chunk replica that is processed in the *Matching Subtree*. It is easy to see that S' has cost exactly Th and no bandwidth constraint is violated. Each chunk type is processed.

(\Rightarrow) Let's take any feasible solution S' to I' and produce a solution S to I by taking triples that correspond to active triple gadgets. Using Lemma 13, we conclude that there are exactly n active triple gadgets. By feasibility of S' , we know that each chunk type is processed. From Lemma 12, we know that out of $\deg(e)$ chunk types that correspond to $x \in A \cup B \cup C$, exactly one is processed in the *Matching Subtree*, and therefore each element of $A \cup B \cup C$ is covered. \square

6. Related Work

There has recently been much interest in programming models and distributed system architectures for the processing and analysis of big data (e.g. [3, 10, 36]). The model studied in this paper is motivated by MapReduce [10] like batch-processing applications, also known from the popular open-source implementation *Apache Hadoop*. These applications generate large amounts of network traffic [7, 26, 38], and over the last years, several systems have been proposed which provide a provable network performance, also in shared cloud environments, by supporting relative [27, 28, 32] or, as in the case of our paper, *absolute* [5, 22, 29, 30, 35] bandwidth reservations between the virtual machines.

The most popular virtual network abstraction for batch-processing applications today is the *virtual cluster*, introduced in the Oktopus paper [5],

and later studied by many others [26, 16, 31, 35]. In particular, Proteus [35] improves upon the Oktopus [5] embedding algorithm of fat-trees and makes the case for a time-adaptive embedding. The Kraken system [16] is based on an optimal embedding algorithm of fat-trees and allows to elastically scale both link as well as node resources. In [31], Rost et al. show that the virtual cluster abstraction can even be embedded on general graphs in polynomial time, and initiate the algorithmic study of a Hose interpretation of the virtual cluster abstraction.

Several heuristics have been developed to compute “good” embeddings of virtual clusters: embeddings with small footprints (minimal bandwidth reservation costs). The virtual network embedding problem has also been studied for more general graph abstractions (e.g., motivated by wide-area networks). [8, 14]

From a theoretical perspective, the virtual network embedding problem can be seen as a generalization of classic VPN graph embedding problems [20, 23], in the sense that in virtual network embedding problems, also the embedding endpoints are flexible. In this respect, the virtual network embedding problem can also be seen as a generalization of the classic NP-hard Minimum Linear Arrangement problem which asks for the embedding of guest graphs on a simple *line topology* (rather than tree-like topologies as studied in this paper) [11, 12].

However, to the best of our knowledge, we are the first to provide an algorithmic study of the virtual cluster embedding problem which takes into account data locality as well as the possibility to select replicas—aspects which so far have only been studied from a best-effort perspective and using coarse-grained metrics (e.g., same rack or same server), thus limiting the flexibility of the system [4, 24, 37].

Bibliographic Note. A preliminary version of this paper appeared at the 23rd IEEE International Conference on Network Protocols (ICNP), 2015 [15].

7. Summary and Conclusion

At the heart of locality and replica aware virtual cluster embeddings lie fundamental algorithmic problems. This paper has shown that despite the multiple dimensions of flexibility in terms of chunk assignment and node placement, and despite the large scale of modern datacenters, many problems can be solved efficiently. However, we have also shown that several embedding

NP-hard	5 combinations	RS + MA + FP + NI + BW
	4 combinations	RS + MA + FP + NI; RS + MA + FP + BW; RS + FP + NI + BW
	3 combinations	RS + MA + FP; RS + FP + NI
Flow	4 combinations	RS + MA + NI + BW
	3 combinations	RS + NI + BW; RS + MA + BW
	2 combinations	RS + BW
DP	4 combinations	MA + FP + NI + BW
	3 combinations	MA + FP + NI; MA + FP + BW; FP + NI + BW
	2 combinations	MA + FP; FP + NI;
Matching	3 combinations	RS + MA + NI; MA + NI + BW
	2 combinations	RS + MA; RS + NI; MA + NI; MA + BW; NI + BW
	1 combinations	RS; MA; NI; BW
0 Cost	3 combinations	RS + FP + BW
	2 combinations	RS + FP; FP + BW
	1 combinations	FP

Table 1: Fastest algorithms for different respective problem variants.

problems are NP-hard already in two- and three-level trees—a practically relevant result given today’s datacenter topologies [2]).

Our results are summarized in Table 1. One interesting takeaway from this figure regards the question which properties render the problem NP-hard. For instance, we see that, BW does not influence the hardness of any problem variant, while RS is crucial for NP-hardness. MA only affects hardness if combined with RS. NI is trivial without FP, and FP requires more sophisticated algorithms when combined with NI or MA; in combination with RS and MA or NI, FP renders the problem NP-hard.

Acknowledgments. We would like to thank Paolo Costa for many dis-

cussions. This research is in part supported by Polish National Science Centre grant DEC-2013/09/B/ST6/01538, the EU project Bigfoot FP7-ICT-317858, as well as by the German BMBF Software Campus grant 01IS12056.

- [1] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin. *Network Flows: Theory, Algorithms, and Applications*. Prentice-Hall, Inc., 1993.
- [2] M. Al-Fares, A. Loukissas, and A. Vahdat. A scalable, commodity data center network architecture. In *Proc. ACM SIGCOMM*, pages 63–74, 2008.
- [3] I. Alagiannis, R. Borovica, M. Branco, S. Idreos, and A. Ailamaki. Nodb: Efficient query execution on raw data files. In *Proc. ACM SIGMOD*, pages 241–252, 2012.
- [4] G. Ananthanarayanan, S. Agarwal, S. Kandula, A. Greenberg, I. Stoica, D. Harlan, and E. Harris. Scarlett: Coping with skewed content popularity in mapreduce clusters. In *Proc. EuroSys*, pages 287–300, 2011.
- [5] H. Ballani, P. Costa, T. Karagiannis, and A. Rowstron. Towards predictable data-center networks. In *Proc. ACM SIGCOMM*, 2011.
- [6] R. Chaiken, B. Jenkins, P.-A. Larson, B. Ramsey, D. Shakib, S. Weaver, and J. Zhou. Scope: Easy and efficient parallel processing of massive data sets. *Proc. VLDB Endow.*, 1(2), 2008.
- [7] M. Chowdhury, M. Zaharia, J. Ma, M. I. Jordan, and I. Stoica. Managing Data Transfers in Computer Clusters with Orchestra. In *Proc. ACM SIGCOMM*, 2011.
- [8] N. M. K. Chowdhury and R. Boutaba. A survey of network virtualization. *Comput. Netw.*, 54(5):862–876, 2010.
- [9] P. Crescenzi, V. Kann, M. Halldorsson, M. Karpinski, and G. Woeginger. Maximum 3-dimensional matching. *A Compendium of NP Optimization Problems*, 2000.
- [10] J. Dean and S. Ghemawat. Mapreduce: Simplified data processing on large clusters. In *Proc. USENIX OSDI*, pages 137–150, 2004.
- [11] N. R. Devanur, S. A. Khot, R. Saket, and N. K. Vishnoi. Integrality gaps for sparsest cut and minimum linear arrangement problems. In *Proc. ACM STOC*, 2006.
- [12] J. Díaz, J. Petit, and M. Serna. A survey of graph layout problems. *ACM Comput. Surv.*, 34(3):313–356, 2002.
- [13] R. Duan and H.-H. Su. A scaling algorithm for maximum weight matching in bipartite graphs. In *Proc. 23rd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1413–1424, 2012.
- [14] A. Fischer, J. Botero, M. Beck, H. DeMeer, and X. Hesselbach. Virtual network embedding: A survey. 2013.

- [15] C. Fuerst, M. Pacut, P. Costa, and S. Schmid. How hard can it be? understanding the complexity of replica aware virtual cluster embeddings. In *Proc. 23rd IEEE International Conference on Network Protocols (ICNP)*, 2015.
- [16] C. Fuerst, S. Schmid, L. Suresh, and P. Costa. Kraken: Online and elastic resource reservations for multi-tenant datacenters. In *Proc. 35th IEEE Conference on Computer Communications (INFOCOM)*, 2016.
- [17] H. Gabow. A scaling algorithm for weighted matching on general graphs. In *Proc. IEEE FOCS*, 1985.
- [18] S. Ghemawat, H. Gobioff, and S.-T. Leung. The google file system. In *Proc. ACM SOSP*, pages 29–43, 2003.
- [19] A. V. Goldberg and R. E. Tarjan. Finding minimum-cost circulations by canceling negative cycles. *J. ACM*, 36(4):873–886, 1989.
- [20] N. Goyal, N. Olver, and F. B. Shepherd. The VPN conjecture is true. *Proc. ACM STOC*, 2008.
- [21] A. Greenberg, J. R. Hamilton, N. Jain, S. Kandula, C. Kim, P. Lahiri, D. A. Maltz, P. Patel, and S. Sengupta. VL2: a scalable and flexible data center network. In *Proc. ACM SIGCOMM*, 2009.
- [22] C. Guo, G. Lu, H. J. Wang, S. Yang, C. Kong, P. Sun, W. Wu, and Y. Zhang. Second-Net: A data center network virtualization architecture with bandwidth guarantees. In *Proc. 6th CoNEXT*, 2010.
- [23] A. Gupta, J. Kleinberg, A. Kumar, R. Rastogi, and B. Yener. Provisioning a virtual private network. In *Proc. ACM STOC*, New York, New York, USA, 2001.
- [24] M. Isard, V. Prabhakaran, J. Currey, U. Wieder, K. Talwar, and A. Goldberg. Quincy: Fair scheduling for distributed computing clusters. In *Proc. ACM SOSP*, pages 261–276, 2009.
- [25] Z. Kiraly and P. Kovacs. Efficient implementations of minimum-cost flow algorithms. In *ArXiv Technical Report 1207.6381*, 2012.
- [26] J. C. Mogul and L. Popa. What we talk about when we talk about cloud network performance. *ACM SIGCOMM CCR*, 2012.
- [27] L. Popa, A. Krishnamurthy, S. Ratnasamy, and I. Stoica. Faircloud: Sharing the network in cloud computing. In *Proc. HotNets-X*, 2011.
- [28] L. Popa, P. Yalagandula, S. Banerjee, J. C. Mogul, Y. Turner, and J. R. Santos. Elasticswitch: Practical work-conserving bandwidth guarantees for cloud computing. In *Proc. ACM SIGCOMM*, pages 351–362, 2013.

- [29] B. Raghavan, K. Vishwanath, S. Ramabhadran, K. Yocum, and A. C. Snoeren. Cloud control with distributed rate limiting. In *Proc. SIGCOMM*, 2007.
- [30] H. Rodrigues, J. R. Santos, Y. Turner, P. Soares, and D. Guedes. Gatekeeper: Supporting bandwidth guarantees for multi-tenant datacenter networks. In *Proc. 3rd Conference on I/O Virtualization (WIOV)*, 2011.
- [31] M. Rost, C. Fuerst, and S. Schmid. Beyond the stars: Revisiting virtual cluster embeddings. In *Proc. ACM SIGCOMM Computer Communication Review (CCR)*, 2015.
- [32] A. Shieh, S. Kandula, A. Greenberg, and C. Kim. Seawall: Performance isolation for cloud datacenter networks. In *Proc. USENIX HotCloud*, 2010.
- [33] K. Shvachko, H. Kuang, S. Radia, and R. Chansler. The hadoop distributed file system. In *Proc. IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST)*, pages 1–10, 2010.
- [34] E. Tardos. A strongly polynomial minimum cost circulation algorithm. *Combinatorica*, 5(3):247–255, July 1985.
- [35] D. Xie, N. Ding, Y. C. Hu, and R. Kompella. The only constant is change: incorporating time-varying network reservations in data centers. *ACM SIGCOMM Computer Communication Review (CCR)*, 2012.
- [36] R. S. Xin, J. Rosen, M. Zaharia, M. J. Franklin, S. Shenker, and I. Stoica. Shark: Sql and rich analytics at scale. In *Proce. ACM SIGMOD*, pages 13–24, 2013.
- [37] M. Zaharia, D. Borthakur, J. Sen Sarma, K. Elmeleegy, S. Shenker, and I. Stoica. Delay scheduling: A simple technique for achieving locality and fairness in cluster scheduling. In *Proc. EuroSys*, pages 265–278, 2010.
- [38] Measuring EC2 system performance. <http://goo.gl/V5zhEd>.