

## THE ALL-OR-NOTHING MULTICOMMODITY FLOW PROBLEM\*

CHANDRA CHEKURI<sup>†</sup>, SANJEEV KHANNA<sup>‡</sup>, AND F. BRUCE SHEPHERD<sup>§</sup>

**Abstract.** We consider the *all-or-nothing multicommodity flow* problem in general graphs. We are given a capacitated undirected graph  $G = (V, E, u)$  and a set of  $k$  node pairs  $s_1t_1, s_2t_2, \dots, s_kt_k$ . Each pair has a unit demand. A subset  $S$  of  $\{1, 2, \dots, k\}$  is *routable* if there is a multicommodity flow in  $G$  that simultaneously sends one unit of flow between  $s_i$  and  $t_i$  for each  $i$  in  $S$ . Note that this differs from the *edge-disjoint path* problem (EDP) in that we do not insist on integral flows for the pairs. The objective is to find a *maximum* routable subset  $S$ . When  $G$  is a capacitated tree, the problem already generalizes  $b$ -matchings, and even in this case it is NP-hard and APX-hard to approximate. For trees, a 2-approximation is known for the cardinality case and a 4-approximation for the weighted case. In this paper we show that the natural linear programming relaxation for the all-or-nothing flow problem has a polylogarithmic integrality gap in general undirected graphs. This is in sharp contrast to EDP, where the gap is known to be  $\Theta(\sqrt{n})$ ; this ratio is also the best approximation ratio currently known for EDP. Our algorithm extends to the case where each pair  $s_it_i$  has a demand  $d_i$  associated with it and we need to completely route  $d_i$  to get credit for pair  $i$ ; we assume that the maximum demand of the pairs is at most the minimum capacity of the edges. We also consider the *online admission control* version where pairs arrive online and the algorithm has to decide immediately on its arrival whether to accept it and the accepted pairs have to be routed. We obtain a randomized algorithm which has a polylogarithmic competitive ratio for maximizing throughput of the accepted requests if it is allowed to violate edge capacities by a  $(2 + \epsilon)$  factor.

**Key words.** all-or-nothing, multicommodity flow, routing, undirected graphs, approximation

**AMS subject classifications.** 68Q25, 68W25, 90C27, 90C59

**DOI.** 10.1137/100796820

### 1. Introduction.

**1.1. Background.** A pervasive problem in communication networks is that of allocating bandwidth to satisfy a given collection of service requests. In situations where there is limited network capacity but an abundance of requests, one must optimize over the choice of which requests to satisfy. Such *maximization* problems arise, for instance, in the area of bandwidth trading, or when operators carve out subnets (so-called VPNs) within their network, for sale to interested enterprise customers. Constraints on how bandwidth may be allocated vary according to the type of requesting service, as well as the technology in the underlying network. In SONET networks, for instance, each request must reserve a path between its origin and destination, each link of the path supporting the traffic rate specified by the request. (In this paper, we ignore any restrictions imposed by requirements that traffic be protected against network failures.) In many cases, including in data networks, the scale

---

\*Received by the editors June 1, 2010; accepted for publication (in revised form) March 28, 2013; published electronically July 2, 2013. A preliminary version of this paper appeared in *Proceedings of ACM STOC*, 2004. The first and third authors did most of the work on this paper while at Lucent Bell Labs; they acknowledge support from basic research grant N00014-02-M-0125 from the Office of Naval Research to Bell Labs.

<http://www.siam.org/journals/sicomp/42-4/79682.html>

<sup>†</sup>Department of Computer Science, University of Illinois, Urbana, IL 61801 (chekuri@cs.illinois.edu). This author received recent support from NSF grant CCF-0728782.

<sup>‡</sup>Department of Computer and Information Science, University of Pennsylvania, Philadelphia, PA 19104 (sanjeev@cis.upenn.edu). This author was supported in part by an Alfred P. Sloan Research Fellowship, NSF Career Award CCR-0093117, and NSF grant CCF-1116961.

<sup>§</sup>Mathematics and Statistics, McGill University, 805 Sherbrooke West, Montreal, QC H3A 2K6, Canada (bruce.shepherd@mcgill.ca).

of demands is large relative to the fine granularity at which it can be managed and routed. In such settings, it makes sense to design the networks in terms of (fractional) flows, rather than (unsplittable) paths.

In this paper, we discuss and compare several fundamental models related to this class of optimization problems. In each model we are given an  $n$ -node capacitated graph  $G = (V, E, u)$  (undirected or directed); here  $u$  denotes a nonnegative integer edge (we use edge to also refer to a directed arc) capacity vector. In addition, we are supplied with a set of  $k$  (unit) demand node pairs  $s_1t_1, s_2t_2, \dots, s_kt_k$ , each possibly with its own weight  $w_i$ . We call the  $s_i$ 's and  $t_i$ 's *terminals*; note that they need not be distinct. A subset  $S$  of the demands  $\{1, 2, \dots, k\}$  is *routable* if the demands in  $S$  can be simultaneously "satisfied" while obeying the capacity of the graph. The objective is to find a largest (or maximum  $w$ -weight) routable subset of demands. Depending on what restrictions we place on how demands may be routed, we obtain several distinct models.

In this setting, a fundamental problem is the *edge-disjoint paths problem* (EDP). Here, a set  $S$  is routable if  $G$  contains edge-disjoint paths  $P_i$  for each  $i \in S$  such that  $P_i$  is a path from  $s_i$  to  $t_i$ . For this problem, the best approximation ratio known is  $O(\sqrt{n})$  for undirected graphs [13] and is  $O(\min((n \log n)^{2/3}, \sqrt{m}))$  for directed graphs [39]. The directed version of this problem is provably hard to approximate. In [22], it was shown that for any fixed  $\epsilon > 0$ , there is no  $O(n^{1/2-\epsilon})$ -approximation algorithm unless  $P = NP$ . The story for undirected EDP is incomplete, however. Even though the approximation ratio is polynomial in  $n$ , the best known inapproximability bound states that the problem is hard to approximate to within a factor of  $\Omega(\log^{1/2-\epsilon})$  unless  $NP \not\subseteq ZPTIME(n^{O(\text{polylog}(n))})$  [1]. Closing this gap is a fundamental open problem in approximation algorithms. Such divergence in our understanding of the approximability of undirected versus directed versions of a problem is not rare, and some other examples include generalized Steiner network and the multicut problem. Partly motivated by this gap in our understanding for undirected EDP, we focus on a class of maximization problems where demands only request a *fractional* unit flow in the network. This forms a relaxation of EDP and our positive bounds in this case improve on the best known for EDP. (We establish a polylogarithmic approximation.) Some of these ideas have also subsequently led to improved understanding for certain classes of undirected EDP [11, 12, 13, 14].

**1.2. All-or-nothing multicommodity flows.** For the remainder of the paper we study the *all-or-nothing multicommodity flow problem* which was introduced in [16]; we denote this problem by AN-MCF. For this version, a set  $S$  is *routable* if there is a multicommodity flow in  $G$  that satisfies every demand  $i \in S$ . In other words we want to find the largest weight subset of  $S$  such that the maximum concurrent multicommodity flow for the subset is at least 1. This differs from EDP problem in that we do not insist on integral flows for the pairs. We also observe that given  $S$ , the problem of deciding whether all of  $S$  can be routed is polynomial-time solvable via linear programming (LP). In contrast, for EDP this same decision problem is NP-complete; moreover, while for a fixed number of demands it is polynomial-time solvable by the work of Robertson and Seymour [37], it is NP-hard for undirected graphs even when the terminals are restricted to lie among a given set of four nodes [19].

In trees, the AN-MCF problem coincides with the maximum integer multicommodity flow problem, and this generalizes the general  $b$ -matching problem. This problem on trees is APX-hard [21] and a 2-approximation is known for the cardinality case (each  $w_i = 1$ ) [21] and a 4-approximation for the weighted case [16]. For

general graphs, where fractional routings come into play, there is no previous work on AN-MCF. On the other hand, recently it was shown that for general undirected graphs, AN-MCF is  $\Omega(\log^{1/2-\epsilon})$ -hard to approximate for any  $\epsilon > 0$  under the assumption that  $NP \not\subseteq ZPTIME(n^{O(\text{polylog}(n))})$  [1]. As noted earlier, this is also the strongest hardness result known for “EDP” [1].

A natural LP relaxation (which we denote by MCF-LP) for the AN-MCF is as follows. For each demand pair  $i$ , let  $\mathcal{P}_i$  denote the paths joining  $s_i, t_i$  in  $G$ . We then have a nonnegative variable  $x(P)$  for each  $P \in \mathcal{P}_i$  and each  $i = 1, 2, \dots, k$ . The objective is to maximize  $\sum_i w_i \sum_{P \in \mathcal{P}_i} x(P)$  subject to the constraints  $\sum_{P \in \mathcal{P}_i} x(P) \leq 1$  for each demand  $i$ , and  $\sum_{P: e \in P} x(P) \leq u(e)$  for each edge  $e$ . Note that this same LP is also valid for EDP, and in fact all known approximation algorithms for EDP obtain their ratios directly or indirectly via the lower bound provided by the LP. Throughout, we let OPT denote the optimal value of the LP relaxation for the instance under consideration. Our main result is the following.

**THEOREM 1.1.** *There is a polynomial-time algorithm which given an undirected  $n$ -node instance of AN-MCF problem returns a solution with weight  $\Omega(\frac{\text{OPT}}{\log^3 n \log \log n})$ .*

In contrast to the above, an  $\Omega(\sqrt{n})$  integrality gap is known for this LP when applied to EDP [21]; we mention that this gap is tight [13]. For the AN-MCF problem, in [1] an integrality gap of  $\log^{1/2-\epsilon} n$  is established for the natural flow based LP relaxation. Our proof of Theorem 1.1 uses the results of Räcke on hierarchically decomposing undirected graphs to construct *demand-oblivious* routings for multicommodity flow problems [33]. Specifically, for a given  $n$ -node capacitated graph  $G$ , Räcke constructs an oblivious routing by means of a *flow template* as follows. For each pair of nodes  $s, t$ , a template must specify a fractional unit flow between  $s$  and  $t$ . Räcke’s template has the following property. Given *any* demand matrix on the nodes, routing the flow according to the flow template (that is, if a pair  $st$  has a demand  $d$ , then  $d$  units of flow are routed along the fixed flow paths for the unit flow from  $s$  to  $t$  with each path getting  $d$  times its share of the unit flow), the congestion on any edge is  $O(\log^3 n)$  times the optimal congestion for routing the given demand matrix. In the present setting we combine these *oblivious* routings of Räcke with other routings to avoid bottlenecks in lightly capacitated regions of the network. Azar et al. [6] and Applegate and Cohen [2] show that the *optimal* (with respect to congestion) oblivious routing can be computed in polynomial time by solving an LP. Bansal et al. [7] give an online algorithm to compute a near optimal oblivious routing. However we need the tree-based hierarchical scheme of Räcke in our algorithm. We note that some other applications of oblivious routing schemes also need the tree based hierarchical decomposition. See Maggs et al. [32] for an application to speed up iterative solvers of linear systems.

We have thus far assumed that each demand  $i$  asks for a “unit” flow between its endpoints, but on occasion we may be supplied more generally with a demand  $d_i$  other than 1. We mention that our techniques equally apply to such *multicommodity demand flow problems* (AN-DMCF). Here, a subset  $S$  of the demands is *routable* if there exists a multicommodity flow in  $G$  that satisfies each demand  $i \in S$ . Note that we maintain the all-or-nothing aspect that we receive the credit  $w_i$  for demand  $i$  only if we fulfill the whole demand  $d_i$ . For an instance of the demand flow problem, we let  $d_{\max} = \max\{d_i : i \in S\}$  and  $u_{\min} = \min\{u(e) : e \in E\}$ . When  $d_{\max} \leq u_{\min}$  (*no bottleneck* case), our result for unit demand carries over with a loss of a constant factor in the approximation ratio by a result in [16], whose proof is based on the grouping and scaling technique of [28]. On the other hand, when  $d_{\max}$  and  $u_{\min}$  are allowed

to be arbitrary, it was noted in [22] that the demand flow problem cannot in general be approximated better than  $O(n^{1/2-\epsilon})$  if integer flows are required. However, the proof from [22] extends directly to show the same hardness even if fractional flows are permitted. In [9], the integrality gap of the LP for demand flow is shown to be  $\Omega(n)$  even when  $G$  is a path. Under the no-bottleneck assumption we obtain the following extension of Theorem 1.1. Here OPT is the value of an optimal fractional solution to flow based LP relaxation for AN-DMCF (described formally in section 3.5).

**THEOREM 1.2.** *There is a polynomial-time algorithm which given an undirected  $n$ -node instance of the multicommodity demand flow problem (AN-DMCF) returns a solution with weight at least  $\Omega(\frac{\text{OPT}}{\log^3 n \log \log n})$  such that the flow on each edge  $e$  is at most  $u(e)$  if  $d_{\max} \leq u_{\min}$ .*

**Admission control in the online setting.** We have defined the AN-MCF and AN-DMCF problems as offline optimization problems. We also consider the *online* versions of these problems. We are given a capacitated graph  $G$  up front; however, the pairs that need to be routed arrive online. When a pair  $s_i t_i$  is presented to the online algorithm, it has to immediately decide if it should accept or reject the pair, and if it accepts, it has to route one unit of flow from  $s_i$  to  $t_i$  ( $d_i$  units in the demand case). For the unit demands case, the competitive ratio of the algorithm is the worst-case ratio of the number of demands routed by the online algorithm to the optimal offline algorithm. For nonunit demands, we compare the sum total of demands routed (throughput) by the online algorithm to the flow routed by the optimal offline algorithm. We assume that once a pair is accepted it stays forever. In the routing literature this is referred to as the permanent connection model [35]. We have the following theorem.<sup>1</sup>

**THEOREM 1.3.** *For the AN-MCF and AN-DMCF problems there are randomized online algorithms that route  $\Omega(\frac{\text{OPT}}{\log^5 n \log \log n})$  flow such that the flow on each edge  $e$  is at most  $(2 + \epsilon)u(e)$  if  $d_{\max} \leq u_{\min}$ .*

The approximation ratios in this paper improve by a  $\log n$  factor for graphs that exclude minors of fixed size, in particular for planar graphs.

**1.3. Related work.** The all-or-nothing flow problem in general graphs was introduced in [16]. They obtain results in the context where the supply network is a tree. Of course, in the tree case, there is no advantage to allowing fractional routing: each commodity must route all its demand on a single path. General multicommodity flow problems that require each demand to be routed on a single path are called *unsplittable flow problems* (UFP). EDP is a special case of UFP where all demands are 1. Multicommodity flow and disjoint path problems, along with several variants and special cases, have been extensively studied both for their fundamental importance to combinatorial optimization and their applications to a variety of areas such as network routing, VLSI layout, parallel computing, and many others. We refer the reader to [36, 30, 20, 26, 35, 38, 22, 29, 33, 10] for some pointers. Here we confine ourselves to mentioning the directly relevant literature on EDP and UFP and their online variants.

We first consider the offline case. For both EDP and UFP the best known approximation ratio is  $O(\sqrt{n})$  for undirected graphs [13] and  $O(\min((n \log n)^{2/3}, \sqrt{m}))$  for directed graphs [39, 38]. EDP and UFP are  $\log^{1/2-\epsilon} n$ -hard in undirected graphs [1] and are  $\Omega(n^{1/2-\epsilon})$ -hard in directed graphs [22]. If all pairs share a source, then EDP reduces to the single commodity maximum flow problem and can be solved in

---

<sup>1</sup>The conference version of this paper incorrectly claimed slightly stronger results for the online problem.

polynomial time. Even for UFP, the single source case is tractable in that most variants have constant factor approximation algorithms [25, 28, 18]. For EDP and UFP, large capacities help. Randomized rounding [36] yields constant factor approximation algorithms if  $d_{\max} \leq u_{\min}/(\Omega(\log n))$ . More generally if  $d_{\max} \leq u_{\min}/B$  for integer  $B$ , then an approximation ratio of  $O(n^{1/B})$  is achievable [38, 5]. We mention that all the offline bounds for EDP and UFP also apply to the integrality gap of MCF-LP. It is known, however, that MCF-LP has an integrality gap of  $\Omega(\sqrt{n})$  [21] and this is tight [13].

Now we consider the online versions of EDP and UFP. These problems have applications in ATM networks and are usually referred to as admission control for virtual circuit routing. A variety of models exist based on whether the circuits (pairs) are permanent or temporary and in the temporary case whether their durations are known or unknown. We refer the reader to the survey [35] for more details. Here we confine ourselves to the case of permanent connections. In this case online EDP asks us to maximize the number of pairs accepted compared to the offline optimal. For UFP we are interested in maximizing the *throughput*, that is, the sum of demands accepted. For both these problems the best competitive ratios are comparable to the offline approximation ratios. However, if capacities are large relative to the demands (a realistic assumption in practice) the algorithm of Awerbuch, Azar, and Plotkin [4] is  $O(\log n)$ -competitive provided  $d_{\max} \leq u_{\min}/\Omega(\log n)$ .

In the context of EDP and UFP we can also consider the problem of routing a given set of demands  $S$  so as to minimize the *congestion* of the routing. Congestion of a routing is defined as the maximum over all edges of the ratio of the flow on the edge to its capacity. If flows can be split, then the minimum congestion routing can be computed by solving an LP. However, for integral flow paths or for unsplittable flow, randomized rounding is the only effective algorithm known and yields an  $O(\log n / \log \log n)$  approximation [36]. In directed graphs this gap is known to be tight [31], that is, there are instances where  $S$  can be fractionally routed with congestion 1 while any integral routing has  $\Omega(\log n / \log \log n)$  congestion. Recently Chuzhoy et al. [17] have shown that in directed graphs, the minimum congestion to route a given set of demands is hard to approximate within a factor of  $\Omega(\log n / \log \log n)$ , effectively closing the approximability of the problem. Despite this hardness, even in the case where demands to be routed arrive online, an  $O(\log n)$ -competitive ratio is possible [3]. More recently, Räcke [33] obtained a randomized  $O(\log^3 n)$ -competitive algorithm for minimizing congestion that is *oblivious*; this result of Räcke [33] is the starting point for our work. The ratio has been improved to  $O(\log^2 n \log \log n)$  by Harrelson, Hildrum, and Rao [24], and recently Räcke [34] has obtained an optimal  $O(\log n)$  bound. These polylogarithmic bounds on oblivious routing hold only for undirected graphs. Polynomial-factor lower bounds are known on the competitive ratio for oblivious routing in node-capacitated undirected graphs and consequently also in directed graphs [23].

**2. Preliminaries.** We consider multicommodity flow problems in a given capacitated graph  $G = (V, E, u)$  and we assume that each capacity  $u(e)$  is an integer. We let  $n = |V|$  and  $m = |E|$ . For any graph  $G$  and a node subset  $S \subseteq V$ , we denote by  $\delta_G(S)$ , or simply  $\delta(S)$  if  $G$  is clear from the context, the set of edges of  $G$  with exactly one endpoint in  $S$ .

An instance of a multicommodity flow problem consists of a graph  $G = (V, E, u)$  and a collection of undirected demand pairs  $s_i t_i$  for  $i = 1, 2, \dots, k$ , where each  $s_i, t_i \in V$ . Without loss of generality we assume that demand pairs are distinct. Otherwise, we can add dummy nodes to ensure this. In addition, we may have a weight  $w_i$  associated

with each demand. Let  $\mathcal{P}_i$  denote the set of paths joining  $s_i, t_i$  in  $G$ . For each path  $P$  in  $G$ , let  $x(P)$  denote a nonnegative variable. An assignment  $x$  is called a multicommodity flow and is said to satisfy demand  $i$  if  $\sum_{P \in \mathcal{P}_i} x(P) = 1$ . (If in addition, a demand has an associated value  $d_i$ , then the right-hand side changes from 1 to  $d_i$ .) The *load* of  $x$  on the edge  $e$  is  $l(e) = \sum_i \sum_{P \in \mathcal{P}_i: e \in P} x(P)$  and the congestion on  $e$  is  $l(e)/u(e)$ . The flow  $x$  is feasible if the congestion on each edge is at most 1. A set  $S$  of demands is routable if there is a feasible flow that satisfies each demand  $i \in S$ . We say that a demand  $i$  is routed on path  $P$  if  $x(P) > 0$ . We also refer to a flow as being obtained by routing a demand on certain paths.

The objective of the all-or-nothing multicommodity flow problem is to find a maximum weight routable set. The natural relaxation for this problem is to find a feasible flow that maximizes  $\sum_i \sum_{P \in \mathcal{P}_i} w_i x(P)$ . To facilitate further discussion, we introduce a new variable  $x_i$  for each pair  $i$ , where  $x_i = \sum_{P \in \mathcal{P}_i} x(P)$ , and write the LP relaxation below:

$$(2.1) \quad \begin{aligned} \max \quad & \sum_{i=1}^k w_i x_i \quad \text{s.t.} \\ & x_i - \sum_{P \in \mathcal{P}_i} x(P) = 0, \quad 1 \leq i \leq k, \\ & \sum_{i=1}^k \sum_{P \in \mathcal{P}_i: e \in P} x(P) \leq u(e), \quad e \in E, \\ & x_i, x(P) \in [0, 1], \quad 1 \leq i \leq k, P \in \cup_i \mathcal{P}_i. \end{aligned}$$

We prefer this path formulation for flows to the compact formulation for ease of presentation. Either formulation admits a polynomial-time algorithm to solve this relaxation. We let OPT denote the value of an optimum solution to the above relaxation. We assume that  $u(e) \leq m$  for each  $i$  because of the following reason. We can find a basic solution to (2.1) in polynomial time. One can easily argue that the number of demands for which  $x_i \in (0, 1)$  is at most  $m$ . If the demands for which  $x_i = 1$  form  $\Omega(1)$  of OPT, then we immediately have a good approximate solution. Thus the difficult case is when most of the LP's profit is accrued through the (at most  $m$ ) fractionally routed demands. In this case, we may restrict to these demands and note that the total load on any edge is at most  $m$ . We may thus assume from now on that the edge capacities have been appropriately reduced and are hence polynomially bounded in  $|V|, |E|$ . We also assume that each node  $v$  is the endpoint of at most one demand pair in the input instance; it is easy to add dummy nodes to ensure this.

Consider two nonnegative vectors  $\pi, \pi'$  defined on the nodes of  $G$  such that  $L = \sum_v \pi(v) = \sum_v \pi'(v)$ . By *distributing* (or *routing*)  $\pi(v)$  units of flow from  $v$  to  $\pi'$ , we refer to a (single-source) flow  $f$  with the property that for each  $v'$ , there is a flow of value  $\pi(v)\pi'(v')$  from  $v$  to  $v'$ . By distributing flow from  $\pi$  to  $\pi'$ , we implicitly refer to a multicommodity flow such that for each  $v$  we are distributing  $\pi(v)$  units of flow from  $v$  to  $\pi'$ . We sometimes say that we are sending flow from the distribution  $\pi$  to the distribution  $\pi'$ .

**2.1. Oblivious routing and Racke trees.** For a graph  $G$ , we call a capacitated tree  $T = (V_T, E_T, u_T)$  with a specified *root* node  $r_t$  a *Racke tree* for  $G$  if the leaves of  $T$  are precisely the nodes of  $G$ . In particular,  $r_t$  is not one of the leaves. Note that each edge  $e \in T$  determines, in a natural way, a partition  $(S, V - S)$  of  $G$ 's nodes, corresponding to the nodes of  $G$  that appear in the leaves of the two resulting

subtrees in  $T - e$ . The capacity of an edge  $e$  in  $T$ , denoted by  $u_T(e)$ , is defined as the total capacity on edges in the cut  $\delta_G(S)$ . We also let  $h(T)$  denote the height of  $T$ . Each Racke tree induces a canonical routing strategy for  $G$ , independent of the set of demands. These (demand) oblivious routings are discussed in more detail later; we first describe some of their basic properties.

An *oblivious routing* strategy consists of defining a unit flow  $F_{st}$  for each pair of nodes  $s, t$ ; this collection of flows is called a routing *template*. Given a collection of demands of, say,  $d_i$  between  $s_i$  and  $t_i$ , a template induces a flow as follows. For each  $i$  and each  $P \in \mathcal{P}_i$ , route  $F_{s_i t_i}(P)d_i$  amount of flow along  $P$ . The multicommodity flow obtained is said to be *routed according* to the oblivious routing. For each oblivious routing scheme one can consider the worst-case congestion of the scheme defined as follows. Consider a set  $X$  of ‘‘demands’’ in  $G$ , i.e.,  $X$  is a multiset of undirected edges  $uv$  with  $u, v \in V$ .  $X$  is routable in  $G$  if there is a feasible multicommodity flow for the demands of  $X$  in the capacities of the edges of  $G$ . The congestion of routing  $X$  according to a given oblivious routing scheme is the maximum, over all edges  $e$  of  $G$ , the total flow on  $e$  divided by the capacity of  $e$ . The worst-case congestion  $\alpha \geq 1$  of an oblivious routing scheme is the maximum, over all routable demand sets  $X$ , of the congestion caused by routing  $X$  according to the oblivious scheme.

Each Racke tree  $T$  induces a specific oblivious routing scheme (discussed further in section 2.2). We define  $\alpha(G, T)$ , the congestion caused by this scheme in the following specific way. Consider a set  $X$  of demands in  $G$ .  $X$  is routable in  $T$  if there is a feasible multicommodity flow for  $X$  in  $T$ : each demand of  $X$  is routed along the unique path connecting its endpoints in  $T$  and the total flow on each edge  $e$  of  $T$  is at most  $u_T(e)$ . The quantity  $\alpha(G, T)$  is then the maximum, over all demand sets  $X$  that are routable in  $T$ , of the congestion caused by routing  $X$  in  $G$  according to the oblivious routing scheme given by  $T$ .

**THEOREM 2.1** (Racke [33]). *Let  $G$  be an undirected graph on  $n$  nodes. There exists a Racke tree  $T$  with  $\alpha(G, T) = O(\log^3 n)$  and  $h(T) = O(\log n)$ .*

Racke’s original construction did not yield a polynomial-time algorithm to find such a  $T$ . Subsequently, two independent papers [8] and [24] obtained polynomial-time algorithms. In [24] it is shown how to construct  $T$  with  $\alpha(G, T) = O(\log^2 n \log \log n)$  and  $h(T) = O(\log n)$ . The bound on  $\alpha(G, T)$  improves by a logarithmic factor when  $G$  is a planar graph. In the following, we let  $\alpha(G)$  denote the minimum of  $\alpha(G, T)$  over some computable class of Racke trees  $T$  for  $G$ .

**2.2. The Racke routing scheme.** In this section we describe the oblivious routing strategy of Racke [33] based on a hierarchical decomposition of the given graph. We need a refinement of Racke’s result, and hence we delve into the details of how the oblivious routings are constructed. In doing so, we follow the approach and methods of [8] which simplified that of [33], in addition to giving a polynomial-time construction of the decomposition in [33].

Let  $T = (V_T, E_T, u_T)$  be a Racke tree for  $G = (V, E, u)$ . We denote by  $T_v$  the subtree of  $T$  rooted at  $v$  and let  $L_v$  denote the leaves of the subtree  $T_v$  and  $G_v$ , the subgraph of  $G$  induced by  $L_v$ . Recall then that  $u_T(e)$  equals the capacity of the cut  $\delta(L_v)$  in  $G$ . A key feature of Racke’s routing is that any demand  $st$  with  $s, t \in L_v$  is entirely routed within the subgraph  $G_v$ . We now describe the oblivious routing of Racke in detail.

First we set up some notation. Consider a nonleaf node  $v \in V_T$  with children  $v_1, v_2, \dots, v_k$ . We deal with two important sets of edges induced by such a node  $v$ . The first is the set of edges in the cut induced by the leaves  $L_v$ , that is,  $\delta_G(L_v)$ . Second, we

say an edge is separated by  $v$  if either it lies in  $\delta(L_v)$  or it has its endpoints in distinct subtrees  $T_{v_i}, T_{v_j}$ . Let  $S(v)$  denote the set of such edges. A measure of these edges' capacity is critical in the following; we denote by  $U(v)$  the sum  $\sum_i u(\delta(L_{v_i}))$ . Note that the capacity of each edge in  $S(v)$  is accounted for twice except for those in  $\delta(L_v)$ , which are accounted for once. For a node  $a \in V(G)$  such that  $a \in L_v$ , we denote by  $u_v^{\text{cut}}(a)$  the total capacity on the edges incident to  $a$  that lie in the cut  $\delta_G(L_v)$ , in other words, the quantity  $u(\delta_G(a) \cap \delta_G(L_v))$ . We denote by  $u_v^{\text{sep}}(a)$  the quantity  $u(\delta_G(a) \cap S(v))$ . We denote by  $\pi_v^{\text{cut}}(a)$  and  $\pi_v^{\text{sep}}(a)$  the quantities  $u_v^{\text{cut}}(a)/u(\delta_G(L_v))$  and  $u_v^{\text{sep}}(a)/U(v)$ , respectively. Note that  $\sum_{a \in L_v} \pi_v^{\text{cut}}(a) = \sum_{a \in L_v} \pi_v^{\text{sep}}(a) = 1$ . We sometimes refer to  $\pi_v^{\text{sep}}()$  as the distribution at  $v$ .

We describe a strategy for routing traffic from a node  $s \in V$  to a node  $t \in V$ . Let  $P = (s = v_0, v_1, \dots, v_p, v_{p+1}, \dots, v_l = t)$  be the unique path in  $T$  joining  $s$  and  $t$ , where  $v_p$  is the ‘‘high point’’ of the path (the least common ancestor of  $s$  and  $t$  in  $T$ ). We construct a unit flow from  $s$  to  $t$  by concatenating a collection of (multicommodity) flow vectors that act as building blocks. We describe these now.

Each edge and each internal node of  $P$  sponsor a transformation of flow as we now describe. An edge  $(y, v)$  (where  $y$  is a child of  $v$ ) transforms the node supply vector  $\pi_y^{\text{cut}}$  on nodes of  $L_y$  into a node demand vector  $\pi_v^{\text{sep}}$  on nodes of  $L_v$ . This is called a spreading transformation. In other words, it is distributing flow from  $\pi_y^{\text{cut}}$  to  $\pi_v^{\text{sep}}$ . This transformation is accomplished by a multicommodity flow denoted by  $f_{(y,v)}$  (called the spreading flow) in  $G_v$ . For each  $a \in L_y$  and  $b \in L_v$ , the flow  $f_{(y,v)}$  routes  $d_{(y,v)}(a, b) = \pi_y^{\text{cut}}(a) \cdot \pi_v^{\text{sep}}(b)$  flow between  $a$  and  $b$ . We also denote by  $f_{(v,y)}$  the flow obtained by reversing the orientation on the flow paths in  $f_{(y,v)}$ .

If  $v$  is an internal node of  $P$ , then  $v$  transforms the node supply vector  $\pi_v^{\text{sep}}$  on the nodes of  $L_v$  into a demand  $\pi_v^{\text{cut}}$  on the nodes of  $L_v$ . This is called a concentrating transformation and is accomplished by a multicommodity flow  $f_v$  (called the concentrating flow) on  $G_v$  that routes for  $a, b \in L_v$ , an amount of flow  $d_v(a, b) = \pi_v^{\text{sep}}(a) \cdot \pi_v^{\text{cut}}(b)$  between  $a$  and  $b$ .

We now return to constructing a unit flow from  $s$  to  $t$ . This flow, denoted by  $F_{s,t}$ , is obtained from the path  $P$  by merging the flows  $f_{(v_0,v_1)}, f_{v_1}, f_{(v_1,v_2)}, f_{v_2}, \dots, f_{v_{p-1}}, f_{(v_{p-1},v_p)}, f_{(v_p,v_{p+1})}, f_{v_{p+1}}, f_{(v_{p+1},v_{p+2})}, f_{v_{p+2}}, \dots, f_{(v_{l-1},v_l)}$ .

At this point, we have not mentioned anything about the actual flow paths for these flow vectors. A key aspect of Räcke’s strategy is that these flows are used as building blocks for many different demand pairs. He shows that the resulting template is efficient (with regard to congestion) as long as for each nonleaf node  $v$  in  $T$  there is an efficient routing for a canonical multicommodity problem in  $G_v$  called the exchange flow problem for  $v$ . For each non leaf node  $v \in T$ , we consider a multicommodity flow  $g_v$  in the graph  $G_v$  for the following exchange flow problem. For each pair of nodes  $a, b$  in  $L_v$ , the exchange flow must route  $D_v(a, b) = u_v^{\text{sep}}(a) \cdot u_v^{\text{sep}}(b)/U(v) = \pi_v^{\text{sep}}(a) \cdot \pi_v^{\text{sep}}(b) \cdot U(v)$  (we write  $D(a, b)$  if  $v$  is clear from the context) amount of flow between  $a$  and  $b$ . We denote by  $q_v \leq 1$  a throughput guarantee for this problem, indicating that we may route  $q_v D(a, b)$  times each of these demands simultaneously in the graph  $G_v$  (in other words, the maximum concurrent flow for the demand matrix  $D_v$ ).

The heart of Räcke’s technique is to show that we may obtain our flow vectors  $f_v, f_{(y,v)}$  from a solution for these exchange flow problems.<sup>2</sup> One may show that if some set of demands can be routed in  $T$  with congestion 1, then routing along

---

<sup>2</sup>The final congestion guarantee  $\alpha(G, T)$  is  $\geq 1/q_v$ ; however, the precise dependence varies based on the specifics of the construction of the tree  $T$  and is not relevant for our application.



the flow vectors  $F_{s,t}$  obtained through the exchange flows has congestion  $\alpha(G, T)$  in  $G$ .

LEMMA 2.2. *Let  $T$  be a Racke tree for  $G$  and let  $X$  be a set of demands that is routable in  $T$  and such that for each  $st \in X$ , the path joining  $s, t$  in  $T$  contains  $r$ . For  $v \in V$ , let  $X(v)$  be the number of demands in  $X$  incident to  $v$ . Then, in  $G$  we can distribute  $X(v)$  units of flow from each  $v \in V$  to the Racke distribution at  $r$  (that is, we can route  $X(v)\pi_r^{sep}(u)$  flow from  $v$  to each  $u \in L_r$ ) with congestion  $\alpha(G, T)$ .*

*Proof.* Consider the Racke routing of  $X$  in  $G$ . For any pair  $st \in X$ ,  $r$  is the least common ancestor. It follows that the Racke routing for  $st$  consists of distributing one unit of flow from  $s$  to its Racke distribution  $\pi_r^{sep}$  at  $r$  and similarly for  $t$ . Note that  $\pi_r^{sep}$  is agnostic to the terminal from which flow originates. Since  $X$  is routable in  $T$  with congestion 1, it follows that the congestion in  $G$  for Racke routing  $X$  is at most  $\alpha(G, T)$ .  $\square$

**3. The offline algorithm.** In this section, we present an  $O(\alpha(G) \log n)$  approximation algorithm for the AN-MCF problem. We first develop a scheme that allows us to route a large fraction of demands with low congestion. Specifically, we show that for any  $\epsilon(n) > 0$ , we can obtain a solution of weight  $\Omega(\frac{\epsilon(n)OPT}{\log^3 n \log \log n})$  such that the flow on any edge  $e$  is at most  $\epsilon(n)u(e) + 1$ . We then design a more sophisticated routing scheme that allows us to obtain a solution of weight  $\Omega(\frac{OPT}{\log^3 n \log \log n})$  without exceeding any edge capacities. We make use of the following known results for AN-MCF on trees with integer capacities. For the cardinality case there is a polynomial-time algorithm that routes at least  $OPT/2$  pairs [21], and for the weighted case there is an algorithm that routes pairs of weight at least  $OPT/4$ ; here  $OPT$  is the optimum solution value to the LP relaxation.

**3.1. Starting point.** We start by describing a simple algorithm that is the starting point of our approach:

- (a) Construct a Racke tree  $T$  for  $G$ .
- (b) Scale down all tree capacities by setting  $\bar{u}_t(e) = \lfloor u_T(e)/\alpha(G) \rfloor$ . Let  $\bar{T}$  denote this new tree.
- (c) Solve the all-or-nothing LP induced on  $\bar{T}$ , and let  $OPT(\bar{T})$  denote the value of an optimal LP solution.
- (d) Find a set  $X$  of weight at least  $\Omega(OPT(\bar{T}))$  that can be routed integrally [21, 16].
- (e) For each routed pair  $st \in X$ , send 1 unit of flow from  $s$  to  $t$  in  $G$  using the flow function  $F_{s,t}$  specified by the Racke routing.

LEMMA 3.1. *If all edges in  $T$  are of capacity at least  $\alpha(G)$ , there is a simple  $O(\alpha(G))$ -approximation algorithm for any instance of the AN-MCF problem in  $G$ .*

*Proof.* Let  $OPT(G)$  denote the optimal LP solution value for our instance of AN-MCF. Clearly  $OPT(T) \geq OPT(G)$ . Further, by our assumption that  $u_T(e) \geq \alpha(G)$  for all  $e \in T$ ,  $\lfloor u_T(e)/\alpha(G) \rfloor \geq u_T(e)/2\alpha(G)$  and hence  $OPT(\bar{T}) \geq \frac{OPT(T)}{2\alpha(G)}$ . Using [16], we can choose a subset of pairs  $X$  that can be feasibly routed in  $\bar{T}$  so that  $w(X)$  is at least  $OPT(\bar{T})/4$ ; in the cardinality case by [21] we can choose  $X$  so that  $|X| \geq OPT(\bar{T})/2$ . Note that if  $X$  is feasible on  $\bar{T}$ , then in  $T$  we can route  $\alpha$  units of demand for each  $i \in X$ . Thus in  $G$  we can route  $\alpha$  demand for each pair in  $X$  with congestion at most  $\alpha$ . Hence we can scale down and fractionally route a unit amount for each  $i$  in  $X$  with congestion 1.  $\square$

We subsequently often refer to an instance as consisting of a triple  $(G, T, X)$ , where  $X$  is a set of demands that can be feasibly routed on an associated Racke tree

$T$  for  $G$ . Recall also that each node in  $V$  is the endpoint of at most one demand pair in  $X$ .

**3.2. Preprocessing the instance.** The simple approach outlined above does not apply when there are edges in  $T$  of capacity less than  $\alpha(G)$ . In the following we assume (using [21, 16]) that we are starting with a set of demands  $X$  of weight  $\Omega(\text{OPT}(T))$  which is feasible for the tree  $T$ , i.e., can be integrally routed on the capacitated Racke tree  $T$ .

For a node  $v$  of  $T$ , we denote by  $\ell_v$ , its level or the distance from the root. For a demand  $st \in X$ , we define its level to be the level of the least common ancestor of  $s$  and  $t$  in  $T$ , denoted by  $\text{lca}(s, t)$ . An instance  $(G', T', X')$  of AN-MCF is called NICE if  $X'$  can be routed on  $T'$  with congestion 1 and the level of each demand is 0, that is, for each  $(s, t) \in X$ ,  $\text{lca}(s, t)$  is the root of  $T'$ .

LEMMA 3.2. *Given a feasible instance  $(G, T, X)$  of AN-MCF we can obtain in polynomial-time  $p \leq n$  NICE AN-MCF instances  $(G_i, T_i, X_i)$ ,  $1 \leq i \leq p$  such that (a) each  $G_i$  is a subgraph of  $G$  and each  $X_i \subseteq X$ , (b)  $G_i$ 's are pairwise node-disjoint, and (c)  $\sum_i w(X_i) = \Omega(\frac{w(X)}{h(T)})$ .*

*Proof.* By the pigeonhole principle, there exists a subset  $X' \subseteq X$  with  $w(X') \geq w(X)/h(T)$  such that all demands in  $X'$  have the same level. We focus on the set  $X'$  from here on and let  $\ell$  be the common level for pairs in  $X'$ . Let  $r_1, r_2, \dots, r_p$  be the nodes of  $T$  at level  $\ell$ . For each  $i$ , let  $T_i$  denote the tree  $T_{r_i}$  and  $X'_i$  denote the demands of  $X'$  with both ends in  $T_{r_i}$ . We also let  $G_i$  be the subgraph of  $G$  induced by the leaves of  $T_i$ . Note that the graphs  $G_i$  are node-disjoint and that  $T_i$  is a valid Racke tree for  $G_i$ . Thus the instances  $(G_i, T_i, X_i)$ ,  $1 \leq i \leq p$ , are disjoint and feasible, have the property that the level of each demand is 0, and collectively satisfy  $\sum_i w(X_i) = \Omega(\frac{w(X)}{h(T)})$ .  $\square$

**3.3. Low congestion routings for NICE instances.** Our next goal is to show that given a NICE instance  $(G, T, X)$  and any given  $\epsilon(n) > 0$ , we can find a large subset  $Z \subseteq X$  that can be routed in  $G$  with *low congestion*—a flow of  $\epsilon(n)u(e) + 1$  on each edge  $e$ . We start by establishing a simple lemma about grouping subsets of nodes in an edge-disjoint manner.

LEMMA 3.3. *Let  $G$  be a connected graph with a weight function  $\rho : V \rightarrow [0, W]$  such that  $\sum_{v \in V} \rho(v) \geq W$ . We can find  $p > (\sum_{v \in V} \rho(v)/2W - 1/2)$  edge-disjoint connected subgraphs  $H_1 = (V_1, E_1), H_2 = (V_2, E_2), \dots, H_p = (V_p, E_p)$  and node-disjoint subsets  $S_1, S_2, \dots, S_p$  such that for each  $i$ , (i)  $S_i \subseteq V_i$  and (ii)  $W \leq \sum_{v \in S_i} \rho(v) \leq 2W$ .*

*Proof.* The *weight* of a node subset  $X$  is  $\sum_{v \in X} \rho(v)$  and is denoted by  $\rho(X)$ .  $X$  is called *heavy* if  $\rho(X) \geq W$ . Let  $T'$  be a spanning tree of  $G$  rooted at an arbitrary node  $r$ . For any node  $x$  let  $T'_x$  be the subtree rooted at  $x$ . Choose some  $x$  such that  $V(T'_x)$  is heavy, but none of  $x$ 's children has this property ( $x$  exists since  $\rho(V) \geq W$ ). Let  $x_1, x_2, \dots, x_\ell$  be the children of  $x$  in  $T'$ . Let  $j$  be the smallest index such that  $\gamma = \rho(x) + \sum_{i=1}^j \rho(V(T'_{x_i})) \geq W$ . Since  $\rho(V(T'_{x_i})) < W$  for  $1 \leq i \leq \ell$  and  $\rho(x) \leq W$ ,  $W \leq \gamma < 2W$ . We obtain  $H_1$  by taking the tree obtained from the union of  $T'_{x_1}, T'_{x_2}, \dots, T'_{x_j}$  with  $x$ . We set  $S_1$  to be the nodes in  $H_1$  with nonzero weight. We remove  $T'_{x_1}, T'_{x_2}, \dots, T'_{x_j}$  from  $T'$  and set  $\rho(x) = 0$  to obtain a new tree  $T''$ . We iterate this process on  $T''$  until the total weight of the remaining nodes falls below  $W$ . Note that a node can participate in many subgraphs that we create; however, it can have positive weight only in the first subgraph. At the end of the process we have the desired edge-disjoint connected subgraphs  $H_1, H_2, \dots, H_p$  and disjoint sets of nodes  $S_1, S_2, \dots, S_p$  with  $S_i \subseteq V(H_i)$ . From the construction, for  $1 \leq i \leq p$ ,

$W \leq \rho(S_i) < 2W$ . The final tree has weight at most  $W$ . Therefore  $2pW + W > \rho(V)$ , and hence,  $p > \rho(V)/(2W) - 1/2$ .  $\square$

Let  $S$  be the set of terminals for  $X$ . Note that a terminal is the endpoint of a single demand in  $X$ .

LEMMA 3.4. *Let  $(G, T, X)$  be a nice instance. We can distribute  $\frac{\epsilon}{\alpha(G)}$  flow from each  $v \in S$  according to the Rucke distribution  $\pi_r^{sep}$  on  $V(G)$  such that the flow on any edge  $e \in G$  is at most  $\epsilon u(e)$ .*

*Proof.* Following the proof of Lemma 2.2, we know that  $X$  can be routed in  $G$ , using Rucke routing, with a congestion of  $\alpha(G)$ . Scaling down the flow by  $\alpha(G)/\epsilon$  and noting that each demand is distributing its flow to the distribution  $\pi_r^{sep}$  gives us the desired result.  $\square$

For each  $v \in V(G)$ , we define  $\beta(v) \in [0, 1)$  as follows: if  $v$  is a terminal  $\beta(v) = \epsilon/\alpha(G)$ ; else  $\beta(v) = 0$ . Let  $\beta = \sum_v \beta(v)$ . We will assume  $\beta \geq 1$  for otherwise we can simply route an arbitrary pair on a path connecting them and we obtain the desired approximation easily. We use Lemma 3.3 with  $W = 1$  to group the terminals  $S$  into disjoint clusters  $S_1, S_2, \dots, S_p$ , where  $p = \max(1, \Omega(\beta))$  such that  $\sum_{v \in S_i} \beta(v) \geq 1$ . Note that the lemma guarantees that each cluster  $S_i$  has at least  $\alpha(G)/\epsilon$  and at most  $2\alpha(G)/\epsilon$  terminals.

We now identify a subset of  $Z \subseteq X$  that can be routed in  $G$  with low congestion. We order the pairs  $st \in X$  in nonincreasing order of their weight, consider them one by one, and build a set of demands  $Z$  which we ultimately route. Initially  $Z = \emptyset$ . We also maintain a set of active clusters: initially all clusters are active. Let  $st$  be the current demand pair. We say that the pair  $st$  is feasible at  $s$  if  $s$  is in an active cluster. We add  $st$  to  $Z$  if both  $s$  and  $t$  are feasible. Otherwise we reject the pair  $st$ . If  $st$  is added to  $Z$  we mark the cluster containing  $s$  as inactive. We do similarly for  $t$ . Note that both  $s, t$  could be in the same cluster.

LEMMA 3.5. *The procedure produces  $Z \subseteq X$  such that  $w(Z) = \Omega(\frac{\epsilon}{\alpha(G)}w(X))$ . Each cluster  $S_i$  contains 0, 1, or 2 terminals for demands in  $Z$ . If it contains 2, then they are the endpoints of the same demand.*

*Proof.* We charge the rejected pairs in  $X \setminus Z$  to those in  $Z$ . Suppose  $st$  is rejected because the cluster containing  $s$  was already marked. We charge  $st$  to the pair  $s't' \in Z$  that marked the cluster of  $s$ . Note that the weight of pair  $s't'$  is no smaller than that of  $st$  since it was considered earlier in the ordering. Since each cluster has at most  $2\alpha(G)/\epsilon$  terminals, a pair in  $Z$  is charged by at most  $4\alpha(G)/\epsilon$  other pairs. This proves that  $w(Z) = \Omega(\frac{\epsilon}{\alpha(G)}w(X))$ . The second part of the lemma is easy to see from the marking procedure.  $\square$

LEMMA 3.6. *The subset  $Z$  can be routed in  $G$  such that on each edge  $e$ , the flow is at most  $1 + \epsilon u(e)$ .*

*Proof.* Consider a pair  $st$  in  $Z$ . Suppose first that  $s, t$  lie in a common cluster. In this case we simply connect  $s$  and  $t$  by a path in this cluster. Otherwise, each of  $s, t$  sends a unit of flow to the root's distribution  $\pi_r^{sep}$ . Let  $S_i, S_j$  be the active clusters that contained  $s$  and  $t$ , respectively, when  $st$  was added to  $Z$ . We distribute one unit of flow from  $s$  to the nodes in  $S_i$  such that each node  $v \in S_i$  gets a flow  $\gamma(v) \leq \beta(v)$ . This is feasible since  $\sum_{v \in S_i} \beta_v \geq 1$  and  $S_i$  is connected. Each node  $v \in S_i$  then sends  $\gamma(v)$  flow to the root's distribution  $\pi_r^{sep}$ . The terminal  $t$  similarly sends its unit of flow to the root's distribution  $\pi_r^{sep}$  using nodes in its cluster  $S_j$ .

We bound the congestion of an edge  $e$  as follows. The edge  $e$  can belong to at most one cluster. Within a cluster it can be used to route at most one unit of flow

from a terminal. Now consider flow on  $e$  from the routings from terminals to the root distribution. For  $v$ , let  $f(v)$  be the total flow that is routed from  $v$  to  $\pi_r^{\text{sep}}$ . From our routing above it is easy to see that  $f(v) \leq \beta(v) \leq \frac{\epsilon}{\alpha(G)}$ . Now we can apply Lemma 3.4 to claim that this flow on  $e$  is at most  $\epsilon u(e)$ . Hence the total flow on  $e$  is at most  $1 + \epsilon u(e)$ .  $\square$

**THEOREM 3.7.** *Given an instance of AN-MCF  $(G, X)$  and a Racke tree  $T$ , there is a polynomial-time algorithm that routes  $\Omega(\frac{\epsilon^{\text{OPT}}}{h(T)\alpha(G,T)})$  demands from  $X$  where OPT is the optimum LP solution for  $X$  on  $G$  such that the flow on any edge  $e$  of  $G$  is at most  $1 + \epsilon u(e)$ .*

Using [24] there is a Racke tree  $T$  for  $G$  such that  $\alpha(G, T) = O(\log^2 n \log \log n)$  and  $h(T) = O(\log n)$ ; hence we obtain an approximation ratio of  $O(\log^3 n \log \log n / \epsilon)$ . For planar graphs or graphs that exclude a constant size minor, [24] gives a Racke tree  $T$  such that  $\alpha(G, T) = O(\log n \log \log n)$  and  $h(T) = O(\log n)$  and hence we obtain an approximation ratio of  $O(\log^2 n \log \log n)$  for these graphs.

**3.4. Congestion 1 routings for NICE instances.** We now show how to find a routing that does not violate the capacities. To simplify the description we focus on the cardinality version of the problem; we mention, as needed, the simple modifications that are needed for the weighted version. Again we work with NICE instances. Let  $(G, T, X)$  be a NICE instance. We say that a capacitated graph  $G$  is 2-edge connected if the uncapacitated multigraph obtained by making  $u(e)$  copies of each edge  $e$  is 2-edge connected.

**THEOREM 3.8.** *Given a NICE instance  $(G, T, X)$  of AN-MCF, there is a polynomial-time algorithm that routes  $\max\{1, \Omega(|X|/\alpha(G, T))\}$  pairs from  $X$  in  $G$  without violating capacities.*

The rest of this section is devoted to the proof of the above theorem. We assume without loss of generality that  $u(e) = 1$  for all  $e$ . The basic idea for the proof of Theorem 3.8 is similar to that in section 3.3. Given a nice instance  $(G, T, X)$  each of the pairs in  $X$  can route  $\epsilon/\alpha(G)$  flow to the root distribution such that each edge  $e$  has a load of at most  $\epsilon u(e) = \epsilon$ . In section 3.3 we used tree-based clustering. Each terminal that we chose to route, then sent one unit of flow to  $\Omega(\alpha(G)/\epsilon)$  other terminals in its cluster. This additional routing required one unit of capacity on the edges in each cluster. In this section we present a more complicated clustering in which a terminal can send one unit of flow to  $\Omega(\alpha(G))$  other terminals using at most  $1/2$  unit of capacity on the edges of its cluster subgraph (which may no longer be a tree). Choosing  $\epsilon = 1/2$  ensures that no edge capacity is violated.

Recall that each node  $s \in V(G)$  is incident to at most one demand pair in  $X$ . From the given NICE instance  $(G, T, X)$  we define a node weight function  $\rho$  as follows:  $\rho(v) = 1/(2\alpha(G))$  if  $v$  is a terminal in  $X$  and  $\rho(v) = 0$  otherwise. Hence a node  $v$  is a terminal if and only if  $\rho(v) > 0$ . For a set  $S$ ,  $\rho(S) = \sum_{v \in S} \rho(v)$ . Therefore  $\rho(V) = |X|/\alpha(G)$ . For a graph  $H$  and weight vector  $\rho' : V(H) \rightarrow [0, 1]$ , we say that a node  $v$  is a center in  $H$  if the following property is true: there is a single-source flow of one unit that originates at  $v$  and such that the flow terminating at any node  $y \neq v$  is at most  $\rho'(y)$  and the flow on any edge is at most  $1/2$ . A node  $v$  is a center in  $H$  with respect to a subset  $S \subset V(H)$  if the flow is constrained to be terminated only at nodes in  $S$ .

**LEMMA 3.9.** *Let  $H = (V, E)$  be a connected graph with a node weight function  $\rho : V \rightarrow \mathcal{R}^+$  such that  $\rho(V) \geq 1$ . Then there is a center  $v$  in  $H$  with respect to  $\rho$ .*

*Proof.* We prove this for the case when  $H$  is a tree  $T$  and  $\rho(V) = 1$ . The proof is by induction on  $n = |V|$ . For  $n \leq 2$  the claim can be easily verified. Assume  $n > 2$ .

Then there exists a node  $v$  of degree at least 2. Consider the trees  $T_1, T_2, \dots, T_k$  produced by the removal of  $v$  from  $H$ . If  $\max_i \rho(V(T_i)) \leq 1/2$ , then  $v$  is the desired center. Otherwise let  $T_1$  be such that  $\rho(V(T_1)) > 1/2$ . Let  $T'$  be a tree obtained by removing  $T_2, \dots, T_k$  from  $H$  and setting  $\rho(v) = \rho(v) + \sum_{i=2}^k \rho(V(T_i))$ . By induction, the tree  $T'$  has a center  $v'$ . Note that  $v' \neq v$  since  $v$  is a leaf in  $T'$  and  $\rho(v) < 1/2$ , and hence  $v' \in V(T_1)$ . It can be easily checked that  $v'$  is a center for  $T$  as well.  $\square$

We need a slight relaxation of the definition of a center. We call a node  $v$  a *pseudocenter* in a subgraph  $H$  of  $G$  if either  $v$  is a center or if the following is true: there is a single source flow of one unit that originates at  $v$  and such that the total flow terminating at a node  $y \neq v$  is at most  $\rho(y)$  and the flow on any edge other than a cut edge of  $G$  is at most  $1/2$ . Note that the flow on a cut edge can be up to 1.

LEMMA 3.10. *Let  $G$  be a unit capacity graph and  $X$  a set of node pairs. Let  $C$  be the set of cut edges in  $G$ . Let  $f$  be a multicommodity flow in  $G$  that sends one unit of flow for each pair in  $X$  with the condition that the flow on any edge is strictly less than 2 if it is in  $C$  and at most 1 otherwise. Then there is another multicommodity flow  $f'$  that sends one unit of flow for each pair in  $X$  and such that the load on any  $e \in E$  is at most 1.*

*Proof.* Given  $f$  we decompose it into flow paths for each pair. We ensure that all the flow paths are simple paths. Let  $f'$  be the resulting flow. We claim that  $f'$  satisfies the desired properties. We have that the flow on any edge under  $f'$  is at most that under  $f$ , so all we need to show is that the flow under  $f'$  is at most 1 for each  $e \in C$ . Let  $G_1$  and  $G_2$  be the components of  $G - e$ . If a pair  $st \in X$  is such that  $s, t \in G_1$  or  $s, t \in G_2$ , then it is clear that the flow between  $s$  and  $t$  does not use  $e$  since the flow paths are simple. Hence the only flow that uses  $e$  is for a pair  $s't'$  such that  $s' \in V(G_1)$  and  $t' \in V(G_2)$ . For such a pair the whole unit of their flow crosses  $e$ . Since the load on any cut edge is strictly less than 2, there is at most one such pair, and hence the flow on  $e$  is at most 1.  $\square$

Our goal is to prove the following clustering lemma.

LEMMA 3.11. *Given a NICE instance  $(G, T, X)$  we can find a subset  $Z \subset X$  with  $|Z| = \Omega(\rho(V)) = \Omega(|X|/\alpha(G))$  with the following properties. There is a collection of connected edge-disjoint subgraphs  $H_1, H_2, \dots, H_p$  such that (i) each  $H_i$  contains terminals for at most one pair in  $Z$ , (ii) each  $H_i$  contains a set of nodes  $S_i$  such that the sets  $S_1, \dots, S_p$  are disjoint, and (iii) for each pair  $st \in Z$ , if  $s \in H_i$ , then either  $s$  is a pseudocenter of  $H_i$  with respect to  $\rho$  and  $S_i$  or  $t$  is also in  $H_i$  (and similarly for  $t$ ).*

**Outline of congestion 1 routing.** We prove Theorem 3.8 from the above lemma. A pair  $st \in Z$  is *mated* if there is an  $H_i$  with  $s, t \in H_i$ . From the properties of the above lemma, we can route all the mated pairs in an edge-disjoint fashion. We now show that the unmated pairs in  $Z$  can also be routed. It thus follows that we can route at least half the pairs in  $Z$  with congestion 1 in  $G$ .

Let  $st$  be an unmated pair. We have that  $s \in H_i$  and  $t \in H_j$  with  $i \neq j$  and further  $H_i$  and  $H_j$  contain no other terminals from  $Z$ . Since  $s$  is a pseudocenter in  $H_i$  with respect to  $S_i$ , it can send one unit of flow to  $S_i$  such that each node  $v \in S_i$  receives flow of at most  $\gamma(v) \leq \rho(v)$  and the total flow on each edge  $e \in H_i$  is at most  $1/2$ , unless  $e$  is a cut edge of  $G$  in which case the flow can be up to 1. Since the  $S_i$ 's are node disjoint, the total flow received by a node  $v$  is at most  $\gamma(v)$ . Now each node  $v$  sends  $\gamma(v)$  flow to the root's distribution  $\pi_r^{\text{sep}}$ . Lemma 3.4 implies that we can route this flow sending at most  $1/2$  unit of flow on each edge. Thus, all the terminals in unmated pairs can simultaneously route one unit of flow to the root's distribution

$\pi_r^{\text{sep}}$  such that the flow on any edge  $e$  is at most  $3/2$  if  $e$  is a cut edge of  $G$  and 1 otherwise. From Lemma 3.10, we conclude that all the unmated pairs can be routed without violating edge capacities.

In the weighted case, Lemma 3.11 can be generalized to ensure that  $w(Z) = \Omega(w(X)/\alpha(G))$ . The above outline can be easily extended to the weighted case.

The rest of the subsection is devoted to the proof of Lemma 3.11.

*Initial clustering.* We first describe an algorithm to find the clusters  $H_i$ . We start with a basic clustering as given by Lemma 3.3 with  $W = 1$  and  $\rho(v) = 1/(2\alpha(G))$  for terminals in  $X$ . Let  $H_1, \dots, H_p$  and  $S_1, \dots, S_p$  be as guaranteed by Lemma 3.3. We have that  $p = \Omega(|X|/\alpha(G))$ . Instead of using a tree for  $H_i$  as before, we define our clusters  $H_i$  to be the subgraph induced by  $V(H_i)$ . These larger clusters remain edge-disjoint since the proof of Lemma 3.3 actually guarantees that for  $i \neq j$ ,  $H_i$  and  $H_j$  have at most one node in common. In the following we assume for simplicity that the  $H_i$  are actually node disjoint and therefore that  $S_i = V(H_i)$ . This does not affect the proof and we will indicate later how this assumption can be dispensed with.

Lemma 3.9 implies that in each cluster  $H_i$ , we can choose a center node, denoted by  $c_i = c(H_i)$ . Note that there may be many choices for a cluster center. For a terminal  $s$ , we denote by  $H(s)$  the unique cluster  $H_i$  such that  $s \in S_i$ .

*Phase 1 demands.* The terminal pairs  $Z$  guaranteed by Lemma 3.11 are obtained in two phases, starting with the clusters  $H_1, H_2, \dots, H_p$ . Let  $X' = X$ . We create an initial set  $Z$  greedily as follows. We consider pairs one by one in order of nonincreasing weight. If  $st$  is the current pair, we add  $st$  to  $Z$  and remove from  $X'$  all pairs  $s't'$  if either  $s'$  or  $t'$  is in  $H(s) \cup H(t)$ . We repeat this until  $X'$  is empty. This procedure is very similar to the one in section 3.3. Using a counting argument similar to that of Lemma 3.5 we claim that  $|Z| = \Omega(p)$ . By construction, every cluster  $H$  contains either 0, 1, or 2 terminals from  $Z$ , and if it contains 2, then these are the endpoints of a pair from  $Z$ .

At this point we could proceed as outlined above if each of the terminals in  $Z$  is a center of its cluster. However, this need not be the case and for this we need to refine  $Z$  and also allow the merging of multiple clusters. The following characterization is useful for working with nodes that cannot be cluster centers.

LEMMA 3.12. *Let  $H$  be a connected graph with  $\rho(V(H)) \geq 1$ . For any  $v \in V(H)$ , if  $v$  is not a center, then there exists a cut edge  $e$  in  $H$  such that every center of  $H$  lies in a distinct component from  $v$  in  $H - e$ .*

*Proof.* Let  $H'$  be a graph obtained by adding a new node  $t$  to  $H$  and an edge of capacity  $\rho(w)$  from  $t$  to each  $w \in V(H)$  and setting the capacity of each edge in  $H$  to  $1/2$ . If  $v$  is not a center, then there exists a  $v - t$  cut  $\delta_{H'}(S)$  of capacity less than 1 in  $H'$ , where  $v$  is chosen to lie in  $S$ . Since  $\rho(H) \geq 1$ , if  $S = V(H)$ , then this cut has capacity at least 1. Hence  $V(H) - S$  is nonempty. But then  $|\delta_H(S)| = 1$  for otherwise the cut has capacity at least 1 in  $H'$ . By definition of a center node, any center cannot lie in  $S$ , hence the lemma.  $\square$

A node  $s$  may be a terminal in  $Z$  but not a center for  $H(s)$ . In this case the preceding lemma guarantees that there is a cut edge  $e$  such that in  $H(s) - e$ ,  $s$  is in a component  $F$  that does not contain  $c(H(s))$ . Let  $e_1, e_2, \dots, e_\ell$  be cut edges in  $H(s)$  that separate  $s$  from  $c(H(s))$  such that if  $F_i$  is the component containing  $s$  in  $H(s) - e_i$ , then  $F_i \subset F_{i+1}$  for  $1 \leq i < \ell$ . Call an edge  $e' = xy$  external if  $x \in H(s)$  and  $y \notin H(s)$ . Suppose  $F_\ell$  has an external edge incident to one of its nodes. Then let  $j$  be the smallest index such that there is an external edge incident to  $F_j$  and let  $e(s) = xy$  be this edge and let  $H'$  be the cluster that contains  $y$ . We refer to  $H'$  as

being tagged by  $s$  (or by  $H(s)$ ) and denote it by  $R(s)$ . We observe that  $e_1, \dots, e_{j-1}$  are actually cut edges in  $G$ , for otherwise it would contradict the minimality of  $j$ .

We need a method for turning terminals into centers in larger, merged clusters; this is captured in the following lemma.

**LEMMA 3.13** (tagging lemma). *Let  $s$  be a terminal that is not a center in  $H(s)$ . Suppose  $e(s) = xy$  exists and  $P$  is a path from  $x$  starting with  $e(s)$  to a cluster  $H(s')$  such that  $P$  is edge-disjoint from  $H(s)$ . Then,  $s$  is a pseudocenter in  $H(s) \cup P \cup H(s')$ . If  $e(s)$  does not exist, then  $s$  is a pseudocenter in  $H(s)$ .*

*Proof.* We first consider the case that  $e(s)$  exists. Let  $c$  and  $c'$  be centers of  $H(s)$  and  $H(s')$ , respectively. Let  $G' = H(s) \cup P \cup H(s')$ . To prove that  $s$  is a pseudocenter in  $G'$ , we add a new node  $t$  to  $G'$  and connect it to each node  $a$  in  $G'$  with an edge of capacity  $\rho(a)$ . We make the capacity of each edge in  $G'$  equal to  $1/2$  except the cut edges  $e_1, e_2, \dots, e_{j-1}$  to which we assign a capacity 1. Consider a minimum  $s$ - $t$  cut  $S$  in  $G'$  that contains  $s$ . By Lemma 3.12, neither  $c$  nor  $c'$  is in  $S$  since they are centers in  $G'$ . If one of  $e_1, \dots, e_{j-1}$  is in  $\delta_{G'}(S)$ , then this cut has capacity at least 1, so suppose this is not the case. Note that there are two paths in  $G'$  from  $s$  to  $c$  and  $c'$  such that each edge, other than  $e_1, e_2, \dots, e_{j-1}$ , is used in at most one of the paths. Since  $\delta_{G'}(S)$  contains no cut edges, it must contain at least two edges, and hence the cut has capacity at least 1.

In the case that  $e(s)$  does not exist, we have that  $e_1, \dots, e_\ell$  are cut edges in  $G$ . An argument similar to the above works by considering the auxiliary graph obtained by adding  $t$  to nodes in  $H(s)$ .  $\square$

We could apply a greedy procedure if each cluster is tagged only by a few terminals. However, a cluster can be tagged by many terminals in  $Z$ . We thus capture the tagging structure by creating a digraph  $D$  whose nodes are the clusters. We put an arc in  $D$  from  $H$  to  $H'$  if  $H$  tags  $H'$  and  $H$  does not contain both terminals of a pair in  $Z$ . If  $H$  does not tag any cluster and contains a terminal from  $Z$ , we simply put a loop arc from  $H$  to itself. Note that a cluster node has out-degree 1 if and only if it contains a single terminal from  $Z$  and 0 otherwise. Let  $D'$  be the digraph obtained from  $D$  by identifying to a single node the two clusters associated with each pair  $st \in Z$  and then eliminating any loops. One sees that this results in  $|Z|$  new nodes and each of these has out-degree at most 2. We call the nodes that represent the pairs in  $Z$  special. The remaining nodes, one each for cluster that does not have terminals in  $Z$ , have out-degree 0.

We create a stable set  $I$  of size  $|Z|/5$  among special nodes as follows. The total out-degree of nodes in  $D'$  is at most  $2|Z|$ . Thus there must be a special node  $x$  with in-degree at most 2. We add  $x$  to  $I$  and remove  $x$  and nodes that have arcs to and from  $x$  from  $D'$ . Repeating this process we get the desired set  $I$ . In the weighted case one needs more care. A similar argument as above can be used to show that  $D'$  can be colored with five colors where each color class is a stable set. We simply pick  $I$  to be the color class with the largest weight where the weight of a node in  $D'$  is equal to the weight of the corresponding pair in  $Z$ .

We trim down  $Z$  to the demand pairs corresponding to the special nodes in  $I$ . We now consider the clusters in  $D$  corresponding to the terminals in  $Z$ . Consider a pair  $st \in Z$ . There are three types of pairs in  $Z$ : (i) both  $s$  and  $t$  belong to the same cluster, (ii)  $s$  and  $t$  belong to different clusters and  $H(s)$  tags  $H(t)$  or vice versa, and (iii)  $s$  and  $t$  belong to different clusters and neither of the clusters tags the other. Let  $Z_1, Z_2$ , and  $Z_3$  be the pairs in  $Z$  according to the above classification. For a pair  $st$  in  $Z_1$ , we do not need to do anything since  $H(s)$  contains both  $s$  and  $t$ . For a pair

$st \in Z_2$  we amalgamate  $H(s)$  and  $H(t)$  into a single connected cluster since they have an edge between them. If  $Z_1 \cup Z_2$  contained at least half the weight of the pairs from  $Z$ , then we could simply connect these pairs by disjoint paths in their corresponding clusters. Hence we may focus on  $Z_3$  from now on.

We call a cluster a transit cluster if it is tagged by some terminal in  $Z_3$ . From the construction of  $I$ , there is no transit cluster among terminals in  $Z$ . The basic approach is captured by the tagging lemma. If clusters  $H, H'$  both tag a cluster  $C$ , then by connecting  $H$  and  $H'$  by a path  $P$  through  $C$ , we may have any node in  $H$  or  $H'$  act as a pseudocenter in  $H \cup H'$ . We show that we may pair up all clusters tagging  $C$  by edge-disjoint paths through  $C$  and hence by merging clusters we get a large number of clusters with pseudocenters. The following simple lemma shows how to match up the tagging clusters.

**LEMMA 3.14.** *Let  $T$  be a tree and  $A$  be some even multiset of nodes in  $V(T)$ . Then there exists  $|A|/2$  edge-disjoint paths in  $T$  such that each node  $v \in A$  is the endpoint of exactly  $n_v$  of these paths, where  $v$  occurs  $n_v$  times in  $A$ .*

*Proof.* If  $n_v > 1$ , then we may take as one of our paths the singleton path  $v$  and remove two copies of  $v$  from  $A$ . So we assume that  $n_v = 0$  or  $n_v = 1$  for each node of  $T$ . Let  $S$  be the set of nodes with  $n_v = 1$ . Consider rooting  $T$  at an arbitrary node  $r$ . For each node  $v$ , let  $T_v$  be the subtree of  $T$  rooted at  $v$ . Let  $T_x$  be minimal such that  $S \cap V(T_x) \geq 2$ . Let  $x_1, \dots, x_h$  be the children of  $x$ . By minimality of  $T_x$ , each  $T_{x_i}$  has at most one node in  $S$ . We can assume without loss of generality that each  $T_{x_i}$  has a node  $v_i$  in  $S$ ; otherwise we can remove  $T_{x_i}$  from  $T$ . If  $x \in S$ , then we connect  $v_1$  to  $x$  by a path and remove  $T_{x_1}$  from  $T$ . Otherwise we match up  $v_1$  to  $v_2$  by a path and remove  $T_{x_1}$  and  $T_{x_2}$  from  $T$ . We then repeat this process. It is easy to check that this inductively matches up the pairs in  $S$  by edge-disjoint paths.  $\square$

For a transit cluster  $C$  let  $d(C)$  denote the number of clusters that tag  $C$ . Without loss of generality let  $H_1, \dots, H_\ell$  be the clusters that tag  $C$ . Let  $e_i = u_i v_i$  denote the edge responsible for  $H_i$  tagging  $C$  with  $u_i \in H_i$  and  $v_i \in C$ . Note that the  $v_i$  need not be distinct. For  $v \in C$  let  $n_v$  be the number of  $e_i$  that are incident to  $v$ . We make use of Lemma 3.14, by considering a spanning tree of  $C$ , to connect the clusters that tag  $C$ . If  $d(C) = 1$  we do not do anything. Otherwise we connect  $2\lfloor d(C)/2 \rfloor$  clusters that tag  $C$  by edge-disjoint paths. Note that a path  $P$  that connects clusters  $H(s)$  and  $H(s')$  has  $e(s)$  and  $e(s')$  as its first and last edges, respectively. We call this the pairing procedure on  $C$ . See Figure 3.1. Based on this procedure we define a graph  $B$  on the terminals in  $Z_3$  as follows. We put a self-loop on a node  $s$  either if  $s$  is a center in  $H(s)$  or if the cluster  $C$  that  $H(s)$  tags has  $d(C) = 1$ . We add an edge between  $s$  and  $s'$  in  $B$  if  $R(s) = R(s') = C$  and are connected by a path from the pairing procedure on  $C$ .

If a node  $s$  in  $B$  does not have an edge incident to it, then  $d(R(s))$  is an odd number greater than 1 and  $H(s)$  was left unmatched in the pairing procedure on  $R(s)$ . Let  $Y$  be the pairs  $st$  in  $Z_3$  such that neither  $s$  nor  $t$  is isolated in the graph  $B$ ; these are the pairs that have a chance to have their clusters expanded appropriately. We claim that  $|Y| \geq |Z_3|/3$ . Let  $a$  be the number of isolated terminals. Since each isolated terminal is the unique unmatched terminal in a transit cluster of degree at least 3, the total number of terminals is at least  $3a$ , and the number of pairs with an isolated terminal is at most  $a$ . Hence  $|Y| \geq |Z_3| - a \geq |Z_3|/3$  as claimed. In the weighted case we need to ensure that the left-behind pairs are not of heavy weight. We can ensure this as follows. Lemma 3.14 guarantees that any even multiset of nodes in  $V(T)$  can be matched up. Thus in the pairing procedure for  $C$  we ensure that the



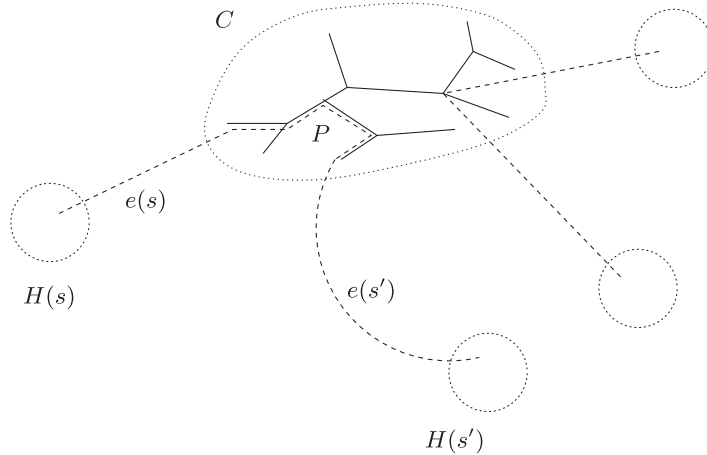


FIG. 3.1. Pairing using  $C$ .

unpaired terminal is the one with the lowest weight among all terminals tagging  $C$ . The rest of the details are simple yet tedious and we omit them.

From  $Y$  we obtain a subset  $Y'$  such that  $|Y'| \geq |Y|/3$  as follows. We consider pairs in  $Y$  one by one in nonincreasing weight order and add the current pair  $st$  to  $Y'$  if neither  $s$  nor  $t$  is connected by an edge in  $B$  to a terminal of some previously added pair in  $Y'$ . Since the degree of each node in  $B$  is at most 1, adding a pair to  $Y'$  can result in the removal at most two pairs from  $Y$  and hence  $|Y'| \geq |Y|/3$ . Moreover, it is easy to see that for a terminal  $s$  from a pair  $uv \in Y'$  ( $s$  could be  $u$  or  $v$ ) exactly one of the following properties holds: (i)  $s$  has a self-loop in  $B$ , (ii)  $ss'$  is an edge in  $B$  and  $s'$  is not an endpoint of a pair in  $Y'$ , or (iii)  $st$  is an edge in  $B$  and  $st$  is a pair from  $Y$ . We process each terminal  $s$  from the pairs in  $Y'$  as follows. If  $s$  is a pseudocenter in  $H(s)$  we do nothing. If  $R(s)$  has degree 1, then we amalgamate  $H(s)$  and  $R(s)$  into a single cluster via the edge  $e$  that  $H(s)$  uses to tag  $R(s)$  (note that in this case  $s$  has a self-loop in  $B$ ). If  $H(s)$  is connected to  $H(s')$  via a path  $P$  through the transit cluster  $R(s)$ , then we amalgamate  $H(s)$  and  $H(s')$  using the path  $P$ . If  $ss' \in Y'$ , then the amalgamated cluster has both endpoints of a pair; otherwise by Lemma 3.13,  $s$  is a pseudocenter in the amalgamated cluster.

After the processing above, the pairs that remain in  $Y'$  are in clusters that satisfy Lemma 3.11. Thus either  $|Y'| = \Omega(|Z|)$  or  $|Z_1 \cup Z_2| = \Omega(|Z|)$ , and in either case we obtain the conditions asserted in Lemma 3.11.

We now address the assumption that we made about the  $H_i$  being node disjoint. In the above proof the only place where we use this assumption is when we find an external edge to tag another cluster. The notion of tagging and the subsequent claims are easily verified if a cluster tags another cluster by a shared node instead of an edge between the clusters.

**3.5. Multicommodity demand flows.** We now consider AN-DMCF, which generalizes AN-MCF in the following way. Each pair  $s_i t_i$  has a nonnegative demand  $d_i$ . A subset  $S$  of the given pairs is routed if there is a feasible multicommodity flow for the pairs in  $S$  where the demand routed for pair  $s_i t_i$  in  $S$  is  $d_i$ . We can translate our result for AN-MCF to AN-DMCF at the expense of an additional constant factor in the approximation ratio. This works under the no-bottleneck assumption that  $d_{\max} \leq u_{\min}$  which

is commonly made in this context. Alternatively, one needs to allow an additive  $d_{\max}$  congestion. The translation is obtained via a grouping and scaling technique of [28] that works for a large class of column-restricted packing integer programs (CPIPs). We consider the following LP relaxation for AN-DMCF that is very similar to the one for AN-MCF in section 2 with the main difference being that we eliminated the variables  $x_i$ :

$$(3.1) \quad \max \sum_{i=1}^k w_i \sum_{P \in \mathcal{P}_i} x(P) \quad \text{s.t.}$$

$$\sum_{i=1}^k d_i \sum_{P \in \mathcal{P}_i: e \in P} x(P) \leq u(e), \quad e \in E,$$

$$\sum_{P \in \mathcal{P}_i} x(P) \leq 1, \quad 1 \leq i \leq k,$$

$$x(P) \in [0, 1], \quad 1 \leq i \leq k, P \in \cup_i \mathcal{P}_i.$$

We observe that the first set of constraints are packing constraints of the form  $Ax \leq b$ , where  $A \geq 0$ ; moreover the constraint matrix  $A$  is column-restricted in that each column has the value in each nonzero entry. If all demands are 1 this is a relaxation for AN-MCF and we have established a polylogarithmic bounds on the integrality gap. The techniques of [28] show that the integrality gap of a CPIP is within a constant factor of the integrality gap of its corresponding 0-1 PIP; in [16] the relevant theorems are stated in a transparent and directly applicable form. In particular we refer to Theorem 3.1. We note that (3.1) is technically not a CPIP due to the second set of constraints; however, the integrality gap translations can be done in the same way (for more details see section 4.3 in [16]). In our particular setting the 0-1 PIP corresponds to AN-MCF and the CPIP corresponds to AN-DMCF. Thus, the result for AN-MCF, with an additional constant factor, gives us Theorem 1.2.

**4. The online algorithm.** We now describe a randomized online algorithm for the AN-MCF problem. The reader would be served well if the earlier offline clustering arguments are fresh in mind. We focus here on the unit demand case and sketch the extension to routing arbitrary demands at the end of the section. Recall that we seek to maximize the throughput, which in the unit-demand case is equivalent to the number of routed pairs. For any  $\epsilon > 0$ , our algorithm achieves a competitive ratio of  $\Omega(h(T)^3 \alpha(G) / \epsilon)$  in expectation, provided we allow a congestion of  $(2 + \epsilon)$  on the edges. (Recall that  $h(T), \alpha(G)$  denote the height and congestion parameters associated with the Racke tree for  $G$ .) For the unit demand case, we in fact achieve a slightly stronger guarantee on the congestion, namely, the total flow on each edge is bounded by  $2 + \epsilon u(e)$ .

As before, our starting point is a Racke tree  $T$  for the input graph  $G$ , which we can assume is a connected graph. We compare the performance of the online algorithm against the maximum possible number of demands that can be routed in the Racke tree  $T$ . The latter is clearly an upper bound on the optimal solution value. We borrow the essential ideas from the offline algorithm; however, the clustering scheme requires nontrivial modifications since we do not have an a priori LP solution to work with.

We assume that all edges in the Racke tree are directed toward the root, and routing a pair  $(s, t)$  is thus equivalent to routing both  $s$  and  $t$  to the root in accordance with the Racke routing. We also assume, for technical reasons that will be clear soon,

that each demand pair  $(s, t)$  is an *ordered pair*, where  $s$  is the *source* node and  $t$  is the *destination* node; given an unordered pair  $st$  we can produce a consistent ordered pair by using an ordering of the nodes of  $G$  and using the lower numbered node in  $s, t$  as the source. For any node  $x$  in  $T$ , let  $T(x)$  denote the subtree of  $T$  rooted at  $x$ , and let  $G(x)$  denote the connected subgraph of  $G$  that corresponds to  $T(x)$ . We let  $L(x)$  denote the leaves in  $T(x)$ . For a leaf node  $z$  and an ancestor  $x$  in  $T$ , we let  $P[z \rightsquigarrow x]$  denote the path from  $z$  to  $x$  in  $T$ .

**4.1. The preprocessing phase.** We start with a preprocessing phase that allows the online algorithm to consider only demand pairs  $(s, t)$  for which  $\text{lca}(s, t)$  is the root of the tree. We lose a factor of  $h(T)$  in the (expected) competitive ratio in order to ensure this property. We also modify the capacities in the Racke tree to satisfy some useful properties; this alteration does not change the value of an optimal solution.

As a first preprocessing step, the algorithm guesses uniformly at random a level  $\ell^* \in \{1, \dots, h(T)\}$ . The algorithm considers routing a request pair  $(s, t)$  only if  $\text{lca}(s, t)$  is a node at level  $\ell^*$  in the tree. Thus the Racke tree is implicitly partitioned into disjoint trees  $T_1, T_2, \dots, T_q$  with roots  $r_1, r_2, \dots, r_q$ , and we only consider a pair  $(s, t)$  if  $\text{lca}(s, t) \in \{r_1, r_2, \dots, r_q\}$ . This process loses at most a factor of  $h(T)$  in expectation. Without loss of generality, from here on we focus our attention on a single Racke tree  $T = T_i$  with root node  $r = r_i$  and an underlying connected graph  $G = G_i$ . We abuse notation and assume that  $G = G_i$  and  $T = T_i$ .

The second preprocessing step reduces capacity of some tree edges so as to ensure the following two properties:

- (P1) Along any path from a leaf node to the root, the edge capacities are non-decreasing.
- (P2) For each edge  $e = (x, y)$  of the oriented  $T$ , the capacity  $u_T(e)$  is equal to the maximum amount of flow in  $T$  that can be routed to  $x$  from the leaves in the subtree of  $T(x)$ .

In doing this capacity reduction, we preserve the “routing capacity” of the tree (recall that at this point, we are only considering demands routed via the root) by proceeding as follows. We say that an edge  $e = (x, y)$  is a violating edge if the parent edge  $e' = (y, z)$  has  $u_T(e') < u_T(e)$  (property (P1) is violated) or the capacity of  $u_T(e)$  is greater than the maximum flow that can be routed from leaves of  $T(x)$  to  $x$  (property (P2) is violated). As long as there is a violation of either (P1) or (P2), we find the deepest violating edge and reduce its capacity by the least amount needed to fix the violation.

The following proposition is easy to see.

**PROPOSITION 4.1.** *The tree  $T'$  obtained by the capacity reduction steps satisfies properties (P1) and (P2). Moreover, if  $S$  is a set of request pairs such that the root is the least common ancestor of each pair in the set, then  $S$  is routable in  $T$  if and only if it is routable in  $T'$ .*

**4.2. The algorithm.** Recall that all edges in  $T$  are directed toward the root. For a leaf node  $z$ , we say that a node  $x$  is the *last bottleneck* node in  $T$  if and only if the capacity of every edge on the path  $P[z \rightsquigarrow x]$  from  $z$  to node  $x$  is smaller than  $\alpha(G)/\epsilon$ , and the capacity of every edge on the path  $P[x \rightsquigarrow r]$  is at least  $\alpha(G)/\epsilon$ . By property (P1), such a node  $x$  must exist. Note that the last bottleneck node  $x$  for a leaf node  $z$  may be the node  $z$  itself or the root  $r$ . If  $x = z$ , then every edge in  $P[z \rightsquigarrow r]$  has capacity at least  $\alpha(G)/\epsilon$ , and if  $x = r$ , then every edge in  $P[z \rightsquigarrow r]$  has capacity less than  $\alpha(G)/\epsilon$ .

*Overview of the algorithm.* The online algorithm chooses a pair of integers  $\ell_1, \ell_2 \in \{1, \dots, h(T)\}$  uniformly at random and routes only demands  $(s, t)$  such that the last bottleneck node for  $s$  is a node at level  $\ell_1$  and the last bottleneck node for  $t$  is a node at level  $\ell_2$ . This worsens the expected competitive ratio achieved by the online algorithm by a factor of  $h(T)^2$ . Prior to the arrival of any demands, the algorithm identifies once and for all a collection of edge-disjoint connected subgraphs of  $G$ , called clusters, such that each leaf node  $z$  has a unique cluster  $C_z$  assigned to it. In the following, we use some clustering arguments similar to those used in section 3.3. The cluster of  $z$  always contains  $z$  as a node. Note that multiple leaf nodes may map to the same cluster. If the algorithm accepts a demand pair  $(s, t)$  for routing, it first distributes a unit of flow from  $s$  to a cluster  $C_s$  in  $G(x)$ , where  $x$  is the last bottleneck node for  $s$ . We ensure that  $C_s$  has routing capacity of at least  $\alpha(G)/\epsilon$  to the root. Similarly, it distributes a unit of flow from  $t$  to a cluster  $C_t$  in  $G(y)$  whose routing capacity is at least  $\alpha(G)/\epsilon$ ; here  $y$  is the last bottleneck node for  $t$ . It then routes a unit of flow from  $C_s$  to  $r$ , and from  $C_t$  to  $r$ , using the Racke distribution. The cluster-to-cluster routing induces a congestion of  $\epsilon$  on the edges of  $G$ . The algorithm ensures that (i) no cluster gets used more than once, (ii) no edge appears in more than two clusters, and (iii) the cluster-to-cluster routing induces an overall congestion of  $\epsilon$ . Thus the overall congestion of the resulting routing solution is bounded by  $(2 + \epsilon)$ . It now remains to describe the clustering scheme.

*The clustering scheme.* Let  $\Gamma_1$  denote the set of all nodes at level  $\ell_1$ , and let  $\Gamma_2$  denote the set of all nodes at level  $\ell_2$  in the Racke tree  $T$ . Note that for any pair of distinct nodes  $x, x' \in \Gamma_1$ , the graphs  $G(x)$  and  $G(x')$  are node-disjoint connected subgraphs of  $G$ . Similarly for any pair of distinct nodes  $y, y' \in \Gamma_2$ , the graphs  $G(y)$  and  $G(y')$  are node-disjoint connected subgraphs of  $G$ . For each node  $x \in \Gamma_1$ , we do an edge-disjoint clustering using only edges in  $G(x)$ , and the resulting clusters are only used for routing demands originating at a leaf node in  $L(x)$ ; moreover, the leaf should have  $x$  as its last bottleneck node. Similarly, we perform a clustering in  $G(y)$  for each  $y \in \Gamma_2$ ; hence overall, each edge is in at most two clusters.

Fix a node  $x$  in  $T$ . We describe the clustering scheme for the leaves  $L(x)$  in the graph  $G(x)$ . If  $x$  is a leaf node, then  $x$  is in its own cluster and there is nothing further to do. We refer to such nodes  $x$  as self-clustered nodes. Otherwise, consider the capacitated tree  $T(x)$  and recall that it satisfies properties (P1) and (P2). We define a new capacity function  $u'_T$  for edges in  $T(x)$  as follows:  $u'_T(e) = \min\{u_T(e), \alpha(G)/\epsilon\}$ . Let  $f$  be a flow of value  $\sum_{i=1}^p u'_T(x_i, x)$  from the leaves of  $T(x)$  to the node  $x$  in accordance with the capacity function  $u'_T$ . Note that such a flow is guaranteed to exist by properties (P1) and (P2). We will slightly abuse the notation and for each leaf node  $z$  in  $T(x)$ , we let  $f(z)$  denote the amount of flow that  $z$  sends to  $x$  in this flow solution; note that  $f(z) \leq \alpha(G)/\epsilon$  for all  $z$  and  $\sum_{z \in L(x)} f(z) = \sum_{i=1}^p u'_T(x_i, x)$ . We will assume from here on that total flow arriving into  $x$ , namely,  $\sum_{i=1}^p u'_T(x_i, x)$ , is at least  $\alpha(G)/\epsilon$  for the following reason. If  $x$  is the root and  $\sum_{i=1}^p u'_T(x_i, x) < \alpha(G)/\epsilon$ , then an optimal algorithm can route at most  $\alpha(G)/\epsilon$  pairs via the root. Thus any algorithm that accepts a single pair to route is  $\alpha(G)/\epsilon$ -competitive. We ignore this trivial case. Otherwise, suppose  $x$  is an internal node of  $T$ . Then if  $\sum_{i=1}^p u'_T(x_i, x) < \alpha(G)/\epsilon$ , it follows that  $\sum_{i=1}^p u'_T(x_i, x) = \sum_{i=1}^p u_T(x_i, x) = u_T(x, x')$ , where  $x'$  is the parent of  $x$ . Since  $u_T(x, x') < \alpha(G)/\epsilon$ , no leaf in  $L(x)$  has  $x$  as its last bottleneck node and hence we do not need to cluster  $G(x)$ .

At a high level, our goal now is to cluster leaf nodes in  $G(x)$  into clusters of  $f$ -weight  $\Theta(\alpha(G)/\epsilon)$  such that each cluster satisfies an additional property, namely, at most  $\Theta(\alpha(G)/\epsilon)$  demands can originate from any cluster in any feasible solution. In

order to do this clustering, we consider a graph  $H$  on  $p$  nodes defined as follows. The graph  $H$  has a node  $v(x_i)$  for each  $G(x_i)$ , and there is an edge between two nodes  $v(x_i)$  and  $v(x_j)$  if and only if there is an edge between a node in  $G(x_i)$  and a node in  $G(x_j)$ . The node  $v(x_i)$  inherits the total  $f$ -weight of leaves in  $G(x_i)$ , that is, the  $f$ -weight of  $v(x_i)$  is  $\sum_{z \in L(x_i)} f(z)$ . We refer to the nodes  $v(x_1), \dots, v(x_p)$  as the hub nodes. Our clustering scheme, described next, will proceed in two steps. In the first step, it will first create node-disjoint trees whose  $f$ -weight is at least  $\alpha(G)/\epsilon$ . Some of the trees created in this step may have  $f$ -weight that far exceeds  $\alpha(G)/\epsilon$ ; the second step of the clustering scheme recovers edge-disjoint clusters of  $f$ -weight  $\Theta(\alpha(G)/\epsilon)$  from these trees.

Let  $\mathcal{T}$  be an arbitrary rooted spanning tree of  $H$ . While the  $f$ -weight of  $\mathcal{T}$  is  $\alpha(G)/\epsilon$  or more, we find a deepest node  $\eta$  with the property that the subtree rooted at  $\eta$  has  $f$ -weight at least  $\alpha(G)/\epsilon$ , remove the sub-tree rooted at  $\eta$  from  $\mathcal{T}$ , and label  $\eta$  as a heavy hub node. We repeat this process on the remaining tree; let  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_k$  denote the subtrees removed in this manner. The process terminates when the  $f$ -weight of the remaining tree is less than  $\alpha(G)/\epsilon$ ; we merge the remaining tree with  $\mathcal{T}_k$  and the root of the new tree is the root of  $\mathcal{T}_k$  before merging. Any node in  $H$  not labeled as heavy in this process is called a light hub node. Observe that only the roots of the subtrees produced are heavy. Note that the  $f$ -weight of a tree  $\mathcal{T}_i$  may be much greater than  $\alpha(G)/\epsilon$ . We remark that this process is slightly different from the clustering scheme in the proof of Lemma 3.3 which produces edge-disjoint clusters but allows a node to be in multiple clusters; we want node-disjoint clustering here since  $H$  is a virtual graph and each node represents a subgraph of  $G$ . We next describe a process that takes the trees  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_k$  and uses them to recover edge-disjoint clusters of  $f$ -weight  $\Theta(\alpha(G)/\epsilon)$  from each  $\mathcal{T}_i$ .

Assume without loss of generality that the nodes have been renumbered so that  $v(x_i)$  is the root of the tree  $\mathcal{T}_i$  for  $1 \leq i \leq k$ . We describe the creation of two sets of edge-disjoint clusters  $\mathcal{C}_i$  and  $\mathcal{C}'_i$  from each  $\mathcal{T}_i$ . From  $\mathcal{T}_i$  we create an expanded tree  $\mathcal{T}'_i$  for  $\mathcal{T}_i$  as follows. We replace the root node  $v(x_i)$  by an arbitrary rooted spanning tree of the graph  $G(x_i)$ , and any edge  $(v(x_\ell), v(x_i))$  in  $\mathcal{T}_i$  is replaced by an edge  $(v(x_\ell), v)$  where  $v$  is some node in  $G(x_i)$  that is connected to some node in  $G(x_\ell)$ . We emphasize that only the root node of  $\mathcal{T}_i$  is expanded into a spanning tree; the nodes of  $\mathcal{T}'_i$  are the nodes of  $G(x_i)$  and the nonroot nodes of  $\mathcal{T}_i$ . Thus  $\mathcal{T}'_i$  is a rooted tree that has the property that for each hub node  $v(x_\ell) \in \mathcal{T}'_i$ , the total weight of the subtree of  $\mathcal{T}'_i$  rooted under  $v(x_\ell)$  is at most  $\alpha(G)/\epsilon$ . The nodes of  $G(x_i)$  inherit their  $f$ -weight in the expanded tree  $\mathcal{T}'_i$  while the  $f$ -weight of each hub node of  $\mathcal{T}'_i$  is the same as its weight in  $\mathcal{T}_i$ . We now use the clustering scheme of Lemma 3.3 on  $\mathcal{T}'_i$  to create edge-disjoint clusters using the  $f$ -values as weights such that each cluster has an  $f$ -weight between  $\alpha(G)/\epsilon$  and  $2\alpha(G)/\epsilon$ . Let  $\mathcal{C}'_i$  denote the set of clusters that we create for the tree  $\mathcal{T}'_i$ . Since any hub node in  $\mathcal{T}'_i$  has the property that its subtree has  $f$ -weight at most  $\alpha(G)/\epsilon$ , the clustering scheme ensures that each light hub node appears in exactly one cluster. Given any cluster in  $C \in \mathcal{C}'_i$ , we can expand  $C$  into a cluster in  $G(x)$  by replacing any light hub node  $v(x_\ell)$  in the cluster  $C$  by a spanning tree of  $G(x_\ell)$ ; the expanded cluster  $C$  would be used for routing demands originating at the leaf nodes in  $L(x_\ell)$ .

We observe that the preceding process that created the clusters in  $\mathcal{C}'_i$  used up the graph  $G(x_i)$  (via its spanning tree expansion) to connect the light hub nodes in the tree  $\mathcal{T}_i$ . We also need to ensure that the leaf nodes in  $L(x_i)$  corresponding to the root node of  $\mathcal{T}_i$  are also assigned to a cluster. We create an alternate cluster for this purpose. We define  $C_i$  to be the single cluster defined by an arbitrarily chosen

spanning tree of the connected subgraph of  $G(x)$  induced by  $\mathcal{T}_i$  and let  $\mathcal{C}_i = \{C_i\}$ . Note that the  $f$ -weight of  $\mathcal{T}_i$  is at least  $\alpha(G)/\epsilon$ . We only assign the cluster to the leaf nodes  $L(x_i)$ , although the other nodes contribute to the  $f$ -weight of  $C_i$ .

For each  $i \in [1..k]$ , the clusters in  $\mathcal{C}'_i$  are edge-disjoint by construction, and there is only a single cluster in  $\mathcal{C}_i$ . However, the cluster  $C_i \in \mathcal{C}_i$  can share edges with clusters in  $\mathcal{C}'_i$ . The leaves  $L(x_i)$  of the root of  $\mathcal{T}_i$  participate in the clusters of  $\mathcal{C}'_i$  but are not assigned to them. Similarly, the leaves  $L(x_\ell)$  of a nonroot hub node  $v(x_\ell)$  participate in the single cluster of  $\mathcal{C}_i$  but are not assigned to it. Let  $\mathcal{C}(x) = \bigcup_{j=1}^k \mathcal{C}_i$  and  $\mathcal{C}'(x) = \bigcup_{j=1}^k \mathcal{C}'_i$ . Since the trees  $\mathcal{T}_1, \dots, \mathcal{T}_k$  are node-disjoint, the clusters in  $\mathcal{C}(x)$  are edge-disjoint, and the clusters in  $\mathcal{C}'(x)$  are edge-disjoint. In our online routing scheme, we will choose uniformly at random one of  $\mathcal{C}(x)$  or  $\mathcal{C}'(x)$  for routing, ensuring edge-disjointness among the clusters utilized in the routing.

*Remark.* The clustering scheme can be somewhat simplified if we settle for a congestion of  $4 + \epsilon$  instead of  $2 + \epsilon$ . We could allow an edge to be in up to two clusters and have a single set of clusters for each  $x$  instead of keeping track of  $\mathcal{C}'(x)$  and  $\mathcal{C}(x)$  and accounting for them separately.

The lemma below summarizes a key property of the clustering scheme defined above.

LEMMA 4.2. *Let  $C$  be any cluster in  $\mathcal{C}(x) \cup \mathcal{C}'(x)$ . Let  $S = \{(s'_1, t'_1), \dots, (s'_j, t'_j)\}$  be a set of requests routable in  $T$  such that for  $1 \leq i \leq j$ , (i)  $\text{lca}(s'_i, t'_i) = r$  and (ii) the last bottleneck node for  $s'_i$  is  $x$ . Then the number of requests in  $S$  with a source assigned to  $C$  is at most  $2\alpha(G)/\epsilon$ .*

*Proof.* Let  $L'(x)$  be the set of leaves in  $T(x)$  such that  $x$  is the last bottleneck node for them; note that  $L'(x) \subseteq L(x)$ . Any  $z \in L'(x)$  has an edge of capacity less than  $\alpha(G)/\epsilon$  on the path  $P[z \rightsquigarrow x]$ . By assumption,  $\{s'_1, \dots, s'_j\} \subseteq L'(x)$ . Since  $S$  is routable in  $T$ , for each pair  $(s'_i, t'_i)$  the path from  $s'_i$  to  $r$  goes via  $x$ . In computing the flow  $f$ , we truncated the edge-capacities to  $\alpha(G)/\epsilon$ ; this does not affect the routability of the pairs in  $S$ .

Now first consider the case that  $C \in \mathcal{C}(x)$ . Then  $C$  corresponds to some heavy hub node, say,  $v(x_i)$ , and is only assigned the leaf nodes in  $L(x_i)$ . Since  $u_{T'}(x_i, x) \leq \alpha(G)/\epsilon$  and  $S$  is routable in  $T$ , it follows that at most  $\alpha(G)/\epsilon$  requests in  $S$  can originate in  $C$ . Now consider a cluster  $C \in \mathcal{C}'(x)$ . Suppose  $C$  was created from the tree  $\mathcal{T}'_i$  for some  $1 \leq i \leq k$ . Then for any light hub node  $x_j$ , the cluster  $C$  either contains all leaves in  $L(x_j)$  or no leaves from  $L(x_j)$ , along with a subset of leaf nodes in the set  $L(x_i)$ . (Recall that  $v(x_i)$  is the heavy hub node identified with the tree  $\mathcal{T}'_i$ .) Let  $v(x_{\ell_1}), \dots, v(x_{\ell_q})$  be the light hub nodes whose leaves appear in  $C$ . Then by our clustering process,  $\sum_{j=1}^q u_{T'}(x_{j_\ell}, x) \leq 2(\alpha(G)/\epsilon)$ , and thus the number of demands in  $S$  that originate from  $C$  cannot exceed  $2\alpha(G)/\epsilon$ , as claimed in the lemma.  $\square$

*Overall algorithm.* Recall that in the preprocessing step we have already restricted to demands with a least common ancestor at some fixed level  $\ell^*$  in the Racke tree. We summarize below the online algorithm for routing demands whose ancestor is a particular node  $r$ . At the outset we fix two integers  $\ell_1, \ell_2 \in \{1, 2, \dots, h(T)\}$  which are chosen independently at random. Recall that we let  $\Gamma_1$  and  $\Gamma_2$  denote the set of all nodes at level  $\ell_1$  and level  $\ell_2$ , respectively, in the tree  $T$ . For each node in  $x \in \Gamma_1$ , we compute the set of clusters  $\mathcal{C}(x)$  and  $\mathcal{C}'(x)$ . We choose uniformly at random between the clustering  $\mathcal{C}(x)$  and  $\mathcal{C}'(x)$ . If  $\mathcal{C}(x)$  is chosen, each leaf node  $v$  such that  $v$  appears in the subtree of a heavy hub node is assigned the unique cluster in  $\mathcal{C}(x)$  in which it appears, and all other leaf nodes do not get a cluster assignment. If  $\mathcal{C}'(x)$  is chosen, then each leaf node  $v$  such that  $v$  appears in the subtree of a light

hub node is assigned the unique cluster in  $\mathcal{C}'(x)$  in which it appears, and all other leaf nodes do not get a cluster assignment. Thus each leaf node in  $G(x)$  gets assigned a cluster with probability  $1/2$ . We do the same random cluster assignment for each node  $y \in \Gamma_2$ .

- (a) Consider an arriving request  $(s, t)$ . Reject it outright if either of the following conditions holds:
  - $\text{lca}(s, t)$  is not  $r$  or
  - the levels of last bottleneck nodes of  $s$  and  $t$  are not  $\ell_1, \ell_2$ .
- (b) Let  $x, y$  be the last bottleneck nodes of  $s, t$ , respectively;  $x$  is at level  $\ell_1$  and  $y$  is at level  $\ell_2$ . We accept the request  $(s, t)$  if all three conditions below are satisfied:
  - (b1)  $s$  is either a self-clustered node (i.e., a leaf node) or it was assigned a  $\Gamma_1$ -cluster at the beginning of the algorithm, and  $t$  is either a self-clustered node (i.e., a leaf node) or it was assigned a  $\Gamma_2$ -cluster at the beginning of the algorithm.
  - (b2) If  $s$  is not a self-clustered node, then its assigned  $\Gamma_1$ -cluster is unused thus far, and if  $t$  is not a self-clustered node, then its  $\Gamma_2$ -cluster is unused thus far, and
  - (b3) There is a path in  $T$  of residual capacity at least  $\alpha(G)/\epsilon$  from each of  $x$  and  $y$  to  $r$ .
- (c) If the pair  $(s, t)$  is accepted, we proceed as follows. If node  $s$  is not self-clustered, it distributes one unit of flow to nodes in its cluster in accordance with the weight function  $f$  such that a node  $z$  with weight  $f(z)$  received a flow of at most  $f(z) \cdot \epsilon / \alpha(G)$ . Each of the nodes in the cluster in turn sends its flow to the root according to the Racke distribution at the root. If we route using a cluster, we mark it as used; this does not apply to a self-clustered node. We next remove a capacity of  $\alpha(G)/\epsilon$  from the edges along the path from  $x$  to  $r$ . We proceed in a similar manner for the node  $t$  and remove a capacity of  $\alpha(G)/\epsilon$  from the edges along the path from  $y$  to  $r$ .

*Analysis.* We first observe that the algorithm correctly routes the accepted pairs. For each routed pair  $(s, t)$ , the algorithm creates a unit flow from  $s$  to the root  $r$  in  $T$ . This corresponds to  $s$  distributing one unit of flow to the nodes in  $G$  according to the Racke distribution given by  $r$ . Similarly  $t$  distributes one unit of flow to the nodes in  $G$  according to the Racke distribution given by  $r$ . This ensures that each accepted pair has a unit flow sent from  $s$  to  $t$ .

We now bound the congestion on the edges induced by the routing of the accepted pairs.

**LEMMA 4.3.** *The algorithm routes the accepted pairs such that the flow on an edge  $e$  with capacity  $u(e)$  is at most  $2 + 2\epsilon u(e)$ .*

*Proof.* For each demand pair  $(s, t)$  routed by the online algorithm, the routing uses a cluster for the source  $s$ , a cluster for the destination  $t$ , and a routing to distribute flow from the nodes in the cluster of  $s$  to the nodes in the cluster of  $t$ . The routing inside the clusters for source nodes creates at most a congestion of 1 on each edge in  $G$  since for any pair of nodes  $x', x'' \in \Gamma_1$ , the graphs  $G(x')$  and  $G(x'')$  are node-disjoint, the clusters for a node  $x$  are edge-disjoint, and each cluster is used at most once. Similarly, the routing inside the clusters for destination nodes creates at most a congestion of 1 on each edge in  $G$ . Finally, we claim that for any edge  $e$  in  $G$ , the total capacity utilized in the Racke routings over all pairs is at most  $2\epsilon u(e)$ . Consider a leaf node  $z$  and let  $f'(z)$  be the total flow it receives from all end nodes of accepted pairs

that it then sends to the root using the Racke distribution. It suffices to prove that  $T$  supports a simultaneous flow of  $\alpha(G)/(2\epsilon) \cdot f'(z)$  from each leaf node  $z$  to the root  $r$ ; this implies that in  $G$  these routings induce a congestion of at most  $2\epsilon u(e)$ . To see this, focus on a node  $x \in \Gamma_1$ . If  $x$  is a leaf, then it belongs to its own cluster each time a pair with endpoint  $x$  is routed, we remove  $\alpha(G)/\epsilon$  capacity on the path  $P[x \rightsquigarrow r]$ . Now suppose  $x$  is an internal node. For each routed pair  $(s, t)$  with  $s \in L(x)$ , we remove a capacity of  $\alpha(G)/\epsilon$  on the path  $P[x \rightsquigarrow r]$ . Therefore, what remains is to ensure each leaf  $z \in L(x)$  can route  $\alpha(G)/(2\epsilon) \cdot f'(z)$  to  $x$ . Note that  $z$  is in at most one cluster in  $G(x)$  and since a cluster is used at most once,  $f'(z) \leq 2\epsilon/\alpha(G) \cdot f(z)$  by the algorithm’s routing. (The factor of 2 is to account for the fact that  $x$  may also be in  $\Gamma_2$  for the destination nodes.) Recall that, by definition, the leaves in  $L(x)$  can simultaneously send a flow of  $f(z)$  each to  $x$  in  $T$ .

Thus the overall flow on an edge  $e$  at most  $2 + 2\epsilon u(e)$ . □

We now show that the online algorithm routes in expectation a polylogarithmic fraction of demands routed by an optimal solution. Fix an optimal solution  $\text{OPT}$  where we only consider demands  $(s, t)$  with  $\text{lca}(s, t) = r$ . Let  $\mathcal{S}_{\text{OPT}}$  be the set of demand pairs  $st$  routed in  $\text{OPT}$  such that the levels of last bottleneck nodes of  $s$  and  $t$  are  $\ell_1$  and  $\ell_2$ , respectively. Note that  $\mathcal{S}_{\text{OPT}}$  is a random set and that  $\mathbf{E}[\mathcal{S}_{\text{OPT}}] \geq \text{OPT}/h(T)^2$ ; we use the fact that the pairs or ordered for this claim. Let  $\mathcal{S}_{\text{ON}}$  denote the set of demands routed by the online algorithm.

LEMMA 4.4. *The expected size of the set  $\mathcal{S}_{\text{OPT}} \setminus \mathcal{S}_{\text{ON}}$  is bounded by  $(3|\mathcal{S}_{\text{OPT}} \setminus \mathcal{S}_{\text{ON}}|)/4 + O(\alpha(G)/\epsilon) \cdot |\mathcal{S}_{\text{ON}}|$ . Thus in expectation, we have  $|\mathcal{S}_{\text{ON}}| = \Omega(\epsilon/\alpha(G)) \cdot |\mathcal{S}_{\text{OPT}}|$ .*

*Proof.* Consider any demand  $(s, t) \in \mathcal{S}_{\text{OPT}} \setminus \mathcal{S}_{\text{ON}}$ . The probability that  $(s, t)$  is rejected in step (b1) above is bounded by  $3/4$ , and thus the expected number of demands rejected in this manner is at most  $(3|\mathcal{S}_{\text{OPT}} \setminus \mathcal{S}_{\text{ON}}|)/4$ . We now account for the remaining demand pairs in  $\mathcal{S}_{\text{OPT}} \setminus \mathcal{S}_{\text{ON}}$  and show that their number can be bounded by  $O(\alpha(G)/\epsilon)|\mathcal{S}_{\text{ON}}|$ .

If  $(s, t)$  is rejected in step (b2) above, then  $\mathcal{S}_{\text{ON}}$  contains a demand  $(s', t')$  such that either  $s, s'$  share a cluster of node  $x \in \Gamma_1$  or  $t, t'$  share a cluster of a node  $y \in \Gamma_2$ . If  $s, s'$  share a cluster we charge  $(s, t)$  to the demand  $(s', t')$  routed in  $\mathcal{S}_{\text{ON}}$ . From Lemma 4.2, the total number of requests in  $\mathcal{S}_{\text{OPT}}$  that can originate in the clusters of  $s$  and  $t$  is bounded by  $2\alpha(G)/\epsilon$  each. Therefore we charge each demand in  $\mathcal{S}_{\text{ON}}$  at most  $4\alpha(G)/\epsilon$  times in this manner.

We now account for demands that are rejected in step (b3) above. We observe that if a demand routed in  $\mathcal{S}_{\text{OPT}}$  shares an edge in  $T$  with a demand routed in  $\mathcal{S}_{\text{ON}}$  along the path from  $x$  to  $r$ , then it must share all the edges on the path from  $x$  to  $r$  (since we only consider demands with least common ancestor  $r$ ). Similarly, if a demand routed in  $\mathcal{S}_{\text{OPT}}$  shares an edge with a demand routed in  $\mathcal{S}_{\text{ON}}$  along the path from  $y$  to  $r$ , then it must share all the edges on the path from  $y$  to  $r$ . Since we remove  $\alpha(G)/\epsilon$  capacity along the paths from  $x$  and  $y$  to  $r$  for each demand routed in  $\mathcal{S}_{\text{ON}}$ , we can account for all demands in  $\mathcal{S}_{\text{OPT}}$  rejected in step (b2) by charging at most  $2\alpha(G)/\epsilon$  demands to each demand in  $\mathcal{S}_{\text{ON}}$ .

Thus combining accounting of demands rejected in steps (b1), (b2), and (b3), we get  $|\mathcal{S}_{\text{OPT}} \setminus \mathcal{S}_{\text{ON}}| \leq (3|\mathcal{S}_{\text{OPT}} \setminus \mathcal{S}_{\text{ON}}|)/4 + 6(\alpha(G)/\epsilon)|\mathcal{S}_{\text{ON}}|$ . Hence  $|\mathcal{S}_{\text{ON}}| = \Omega(\epsilon/\alpha(G)) \cdot |\mathcal{S}_{\text{OPT}}|$ . □

Since  $\mathbf{E}[\mathcal{S}_{\text{OPT}}] \geq \text{OPT}/h(T)^2$ , we have that  $\mathbf{E}[\mathcal{S}_{\text{ON}}] = \Omega(\text{OPT} \cdot \epsilon/(\alpha(G)h(T)^2))$ . Finally, since we lose an additional factor of  $h(T)$  in the preprocessing phase, we obtain the following theorem.



**THEOREM 4.5.** *For any  $\epsilon > 0$ , the online algorithm routes in expectation an  $\Omega\left(\frac{\epsilon}{(h(T))^3 \alpha(G)}\right)$ -fraction of demands routed in an optimal solution. Moreover, the flow on each edge  $e$  in the graph  $G$  is bounded by  $2 + \epsilon u(e)$  in the solution generated by the online algorithm.*

**Extension to nonunit demands.** Now we consider the case when request pairs may have different demands. The ideas here are standard and have been used several times in the literature and therefore we only give a sketch of the arguments. Let  $d(st)$  denote the demand for pair  $st$ . The objective function is the amount of total demand routed; in other words routing a demand  $st$  gives a value  $d(st)$ . We assume that  $d_{\max} \leq u_{\min}$ . We assume by scaling that  $u_{\min} = 1$  and hence  $d(st) \leq 1$  for each pair. Call a demand pair  $st$  large if  $d(st) > 1/2$ ; otherwise it is small. The algorithm randomly chooses upfront whether to route only demands that are large or only demands that are small. This affects the expected profit only by a factor of 2. If all demands are large, then we can simply use the unit-demand online algorithm by pretending that each of the large demands has a unit demand. If all demands are less than  $1/2$  we can apply the same algorithm as the one described above for the unit demand problem with the following minor changes. When routing a demand on the tree we keep track of the capacity used by the demand. When using a cluster, we let multiple demands use a cluster as long as the total sum of the demands using the cluster does not exceed 1. The analysis is similar to the unit-demand problem. The reason to split into large and small demands follows standard ideas and avoids the following situation. Consider a graph consisting of a single edge of capacity 1. If a small demand of size  $\epsilon$  arrives first and the algorithm accepts, then the adversary gives a demand of size 1 which cannot be routed. If the algorithm rejects the first demand, then the adversary gives no further demands. Clearly no deterministic algorithm can overcome this situation. On the other hand if all demands are guaranteed to be at most  $1/2$ , then it can be seen that accepting the first demand is not a problem in the throughput measure.

**5. Conclusions.** We obtained a polylogarithmic approximation for the AN-MCF problem, a natural relaxation of the edge-disjoint path problem. A main technical tool is the hierarchical graph decomposition that Räcke developed for oblivious routing. In subsequent work [12] we introduce a general framework for approximating routing problems that relies on well-linked decompositions. The Räcke tree based algorithm in this paper can be viewed in that framework; the tree implicitly gives a well-linked decomposition (in fact the general framework is inspired by this observation). In this paper the approximation ratio we obtained for AN-MCF is  $O(\alpha(G, T)h(T))$ , where  $T$  is a Räcke tree for  $G$ ,  $\alpha(G, T)$  is the congestion bound, and  $h(T)$  is the height of the tree. If one is interested in the cardinality version of the problem, then an  $O(\alpha(G, T))$  bound can be obtained; thus the ratios that we give in this paper can be improved by a logarithmic factor for the cardinality version. In particular it can be shown that the tree  $T$  yields an  $O(\alpha(G, T))$ -well-linked decomposition. Well-linked decompositions can also be obtained for node-capacitated routing problems [12], while it is known that there are no oblivious routings with polylogarithmic congestion in the node-capacitated setting [23]. The well-linked framework and related ideas have led to several new algorithms [11, 12, 13, 14, 27, 15] for disjoint paths and AN-MCF. In some recent work Räcke [34] obtained an asymptotically optimal bound of  $O(\log n)$  for oblivious routing. This new result also relies on hierarchical graph decompositions but differs in an important technical way from the previous approach [33, 24]. The desired oblivious routing is constructed from a convex combination of hierarchical decompositions (trees); each tree in the

decomposition is a Racke tree. Several algorithmic applications that rely on a single hierarchical decomposition can be modified to rely on the convex decomposition (see [34]). However, it is not apparent yet whether the convex decomposition would yield a well-linked decomposition with an  $O(\log n)$  ratio; if it did, the ratio would match a lower bound (shown in the conference version of [15] but to appear in the journal version of [11]).

It is an interesting open problem to obtain improved online algorithms for the AN-MCF problem. We believe that a polylogarithmic competitive ratio can be obtained with congestion  $(1 + \epsilon)$ . Finally, is there a deterministic algorithm that achieves a polylogarithmic competitive ratio with constant congestion?

**Acknowledgments.** The authors thank two anonymous reviewers for their detailed comments, which helped improve the paper.

#### REFERENCES

- [1] M. ANDREWS, J. CHUZHUY, V. GURUSWAMI, S. KHANNA, K. TALWAR, AND L. ZHANG, *In-approximability of edge-disjoint paths and low congestion routing on undirected graphs*, *Combinatorica*, 30 (2010), pp. 485–520.
- [2] D. APPLEGATE AND E. COHEN, *Making intra-domain routing robust to changing and uncertain traffic demands: Understanding fundamental tradeoffs*, in *Proceedings of ACM SIGCOMM*, 2003, pp. 313–324.
- [3] J. ASPNES, Y. AZAR, A. FIAT, S. PLOTKIN, AND O. WAARTS, *On-line routing of virtual circuits with applications to load balancing and machine scheduling*, *J. ACM*, 44 (1997), pp. 486–504.
- [4] B. AWERBUCH, Y. AZAR, AND S. PLOTKIN, *Throughput-competitive online routing*, in *Proceedings of IEEE FOCS*, 1993, pp. 32–40.
- [5] Y. AZAR AND O. REGEV, *Combinatorial algorithms for the unsplittable flow problem*, *Algorithmica*, 44 (2006), pp. 49–66.
- [6] Y. AZAR, E. COHEN, A. FIAT, H. KAPLAN, AND H. RACKE, *Optimal oblivious routing in polynomial time*, *J. Comput. System Sci.*, 69 (2004), pp. 383–394.
- [7] N. BANSAL, A. BLUM, S. CHAWLA, AND A. MEYERSON, *Online oblivious routing*, in *Proceedings of SPAA*, 2003, pp. 44–49.
- [8] M. BIENKOWSKI, M. KORZENIOWSKI, AND H. RACKE, *A practical algorithm for constructing oblivious routing schemes*, in *Proceedings of SPAA*, 2003, pp. 24–33.
- [9] A. CHAKRABARTI, C. CHEKURI, A. GUPTA, AND A. KUMAR, *Approximation algorithms for the unsplittable flow problem*, *Algorithmica*, 47 (2007), pp. 53–78.
- [10] C. CHEKURI AND S. KHANNA, *Edge disjoint paths revisited*, *ACM Trans. Algorithms*, 3 (2007).
- [11] C. CHEKURI, S. KHANNA, AND F. B. SHEPHERD, *Edge-disjoint paths in planar graphs*, in *Proceedings of IEEE FOCS*, 2004, pp. 71–80.
- [12] C. CHEKURI, S. KHANNA, AND F. B. SHEPHERD, *Multicommodity flow, well-linked terminals, and routing problems*, in *Proceedings of ACM STOC*, 2005, pp. 183–192.
- [13] C. CHEKURI, S. KHANNA, AND F. B. SHEPHERD, *An  $O(\sqrt{n})$  approximation and integrality gap for disjoint paths and unsplittable flow*, *Theory Comput.*, 2 (2006), pp. 137–146.
- [14] C. CHEKURI, S. KHANNA, AND F. B. SHEPHERD, *A note on multiflows and treewidth*, *Algorithmica*, 54 (2009), pp. 400–412.
- [15] C. CHEKURI, S. KHANNA, AND F. B. SHEPHERD, *Edge-disjoint paths in planar graphs with constant congestion*, *SIAM J. Comput.*, 39 (2009), pp. 281–301.
- [16] C. CHEKURI, M. MYDLARZ, AND F. B. SHEPHERD, *Multicommodity demand flow in a tree and packing integer programs*, *ACM Trans. Algorithms*, 3 (2007).
- [17] J. CHUZHUY, V. GURUSWAMI, S. KHANNA, AND K. TALWAR, *Hardness of directed routing with congestion*, in *Proceedings of ACM STOC*, 2007, pp. 165–178.
- [18] Y. DINITZ, N. GARG, AND M. X. GOEMANS, *On the single source unsplittable flow problem*, *Combinatorica*, 19 (1999), pp. 17–41.
- [19] S. EVEN, A. ITAI, AND A. SHAMIR, *On the complexity of timetable and multicommodity flow problems*, *SIAM J. Comput.*, 5 (1976), pp. 691–703.
- [20] A. FRANK, *Packing paths, cuts, and circuits—a survey*, in *Paths, Flows and VLSI-Layout*, B. Korte, L. Lovasz, H. J. Promel, and A. Schrijver, eds., Springer-Verlag, Berlin, 1990, pp. 49–100.

- [21] N. GARG, V. VAZIRANI, AND M. YANNAKAKIS, *Primal-dual approximation algorithms for integral flow and multicut in trees*, *Algorithmica*, 18 (1997), pp. 3–20.
- [22] V. GURUSWAMI, S. KHANNA, R. RAJARAMAN, F. B. SHEPHERD, AND M. YANNAKAKIS, *Near-optimal hardness results and approximation algorithms for edge-disjoint paths and related problems*, *J. Comput. System Sci.*, 67 (2003), pp. 473–496.
- [23] M. T. HAJIAGHAYI, R. D. KLEINBERG, T. LEIGHTON, AND H. RÄCKE, *Oblivious routing on node-capacitated and directed graphs*, *ACM Trans. Algorithms*, 3 (2007).
- [24] C. HARRELSON, K. HILDRUM, AND S. RAO, *A polynomial-time tree decomposition to minimize congestion*, in *Proceedings of SPAA*, 2003, pp. 34–43.
- [25] J. M. KLEINBERG, *Single-source unsplittable flow*, in *Proceedings of IEEE FOCS*, 1996, pp. 68–77.
- [26] J. M. KLEINBERG, *Approximation Algorithms for Disjoint Paths Problems*, Ph.D. thesis, MIT, Cambridge, MA, 1996.
- [27] J. M. KLEINBERG, *An approximation algorithm for the disjoint paths problem in even-degree planar graphs*, in *Proceedings of IEEE FOCS*, 2005, pp. 627–636.
- [28] S. G. KOLLIPOULOS AND C. STEIN, *Improved approximation algorithms for unsplittable flow problems*, *SIAM J. Comput.*, 31 (2001), pp. 919–946.
- [29] P. KOLMAN AND S. SCHEIDLER, *Improved bounds for the unsplittable flow problem*, *J. Algorithms*, 61 (2006), pp. 20–44.
- [30] T. LEIGHTON AND S. RAO, *Multicommodity max-flow min-cut theorems and their use in designing approximation algorithms*, *J. ACM*, 46 (1999), pp. 787–832.
- [31] T. LEIGHTON, S. RAO, AND A. SRINIVASAN, *Multicommodity flow and circuit switching*, in *Proceedings of the Hawaii International Conference on System Sciences (HICSS)*, 1998, pp. 459–465.
- [32] B. MAGGS, G. MILLER, O. PAREKH, R. RAVI, AND S. L. M. WOO, *Finding effective support-tree preconditioners*, in *Proceedings of SPAA*, 2005, pp. 176–185.
- [33] H. RÄCKE, *Minimizing congestion in general networks*, in *Proceedings of IEEE FOCS*, 2002, pp. 43–52.
- [34] H. RÄCKE, *Optimal hierarchical decompositions for congestion minimization in networks*, in *Proceedings of ACM STOC*, 2008, pp. 255–264.
- [35] S. PLOTKIN, *Competitive routing of virtual circuits in ATM networks*, *IEEE J. Selected Areas in Communications*, 13 (1995), pp. 1128–1136.
- [36] P. RAGHAVAN AND C. D. THOMPSON, *Randomized rounding: A technique for provably good algorithms and algorithmic proofs*, *Combinatorica*, 7 (1987), pp. 365–374.
- [37] N. ROBERTSON AND P. D. SEYMOUR, *Graph Minors XIII. The Disjoint Paths Problem*, *J. Comb. Theory, Ser. B*, 63(1), 1995, pp. 65–110, <http://dx.doi.org/10.1006/jctb.1995.1006>.
- [38] A. SRINIVASAN, *Improved approximations for edge-disjoint paths, unsplittable flow, and related routing problems*, in *Proceedings of IEEE FOCS*, 1997, pp. 416–425.
- [39] K. VARADARAJAN AND G. VENKATARAMAN, *Graph decomposition and a greedy algorithm for edge-disjoint paths*, in *Proceedings of ACM-SIAM SODA*, 2004, pp. 379–380.