

# LP-rounding Algorithms for the Fault-Tolerant Facility Placement Problem

Li Yan and Marek Chrobak  
Department of Computer Science  
University of California at Riverside

## Abstract

The Fault-Tolerant Facility Placement problem (FTFP) is a generalization of the classic Uncapacitated Facility Location Problem (UFL). In FTFP we are given a set of facility sites and a set of clients. Opening a facility at site  $i$  costs  $f_i$  and connecting client  $j$  to a facility at site  $i$  costs  $d_{ij}$ . We assume that the connection costs (distances)  $d_{ij}$  satisfy the triangle inequality. Multiple facilities can be opened at any site. Each client  $j$  has a demand  $r_j$ , which means that it needs to be connected to  $r_j$  different facilities (some of which could be located on the same site). The goal is to minimize the sum of facility opening cost and connection cost.

The main result of this paper is a 1.575-approximation algorithm for FTFP, based on LP-rounding. The algorithm first reduces the demands to values polynomial in the number of sites. Then it uses a technique that we call adaptive partitioning, which partitions the instance by splitting clients into unit demands and creating a number of (not yet opened) facilities at each site. It also partitions the optimal fractional solution to produce a fractional solution for this new instance. The partitioned instance satisfies a number of properties that allow us to exploit existing LP-rounding methods for UFL to round our partitioned solution to an integral solution, preserving the approximation ratio. In particular, our 1.575-approximation algorithm is based on the ideas from the 1.575-approximation algorithm for UFL by Byrka *et al.*, with changes necessary to satisfy the fault-tolerance requirement.

# 1 Introduction

In the *Fault-Tolerant Facility Placement* problem (FTFP), we are given a set  $\mathbb{F}$  of *sites* at which facilities can be built, and a set  $\mathbb{C}$  of *clients* with some demands that need to be satisfied by different facilities. A client  $j \in \mathbb{C}$  has demand  $r_j$ . Building one facility at a site  $i \in \mathbb{F}$  incurs a cost  $f_i$ , and connecting one unit of demand from client  $j$  to a facility at site  $i$  costs  $d_{ij}$ . Throughout the paper we assume that the connection costs (distances)  $d_{ij}$  form a metric, that is, they are symmetric and satisfy the triangle inequality. In a feasible solution, some number of facilities, possibly zero, are opened at each site  $i$ , and demands from each client are connected to those open facilities, with the constraint that demands from the same client have to be connected to different facilities. Note that any two facilities at the same site are considered different.

It is easy to see that if all  $r_j = 1$  then FTFP reduces to the classic Uncapacitated Facility Location problem (UFL). If we add a constraint that each site can have at most one facility built on it, then the problem becomes equivalent to the Fault-Tolerant Facility Location problem (FTFL). One implication of the one-facility-per-site restriction in FTFL is that  $\max_{j \in \mathbb{C}} r_j \leq |\mathbb{F}|$ , while in FTFP the values of  $r_j$ 's can be much bigger than  $|\mathbb{F}|$ .

The UFL problem has a long history; in particular, great progress has been achieved in the past two decades in developing techniques for designing constant-ratio approximation algorithms for UFL. Shmoys, Tardos and Aardal [16] proposed an approach based on LP-rounding, that they used to achieve a ratio of 3.16. This was then improved by Chudak [5] to 1.736, and later by Sviridenko [17] to 1.582. The best known “pure” LP-rounding algorithm is due to Byrka *et al.* [3] with ratio 1.575. Byrka and Aardal [2] gave a hybrid algorithm that combines LP-rounding and dual-fitting (based on [10]), achieving a ratio of 1.5. Recently, Li [13] showed that, with a more refined analysis and randomizing the scaling parameter used in [2], the ratio can be improved to 1.488. This is the best known approximation result for UFL. Other techniques include the primal-dual algorithm with ratio 3 by Jain and Vazirani [11], the dual fitting method by Jain *et al.* [10] that gives ratio 1.61, and a local search heuristic by Arya *et al.* [1] with approximation ratio 3. On the hardness side, UFL is easily shown to be NP-hard, and it is known that it is not possible to approximate UFL in polynomial time with ratio less than 1.463, provided that  $\text{NP} \not\subseteq \text{DTIME}(n^{O(\log \log n)})$  [6]. An observation by Sviridenko strengthened the underlying assumption to  $\text{P} \neq \text{NP}$  (see [19]).

FTFL was first introduced by Jain and Vazirani [12] and they adapted their primal-dual algorithm for UFL to obtain a ratio of  $3 \ln(\max_{j \in \mathbb{C}} r_j)$ . All subsequently discovered constant-ratio approximation algorithms use variations of LP-rounding. The first such algorithm, by Guha *et al.* [7], adapted the approach for UFL from [16]. Swamy and Shmoys [18] improved the ratio to 2.076 using the idea of pipage rounding introduced in [17]. Most recently, Byrka *et al.* [4] improved the ratio to 1.7245 using dependent rounding and laminar clustering.

FTFP is a natural generalization of UFL. It was first studied by Xu and Shen [20], who extended the dual-fitting algorithm from [10] to give an approximation algorithm with a ratio claimed to be 1.861. However their algorithm runs in polynomial time only if  $\max_{j \in \mathbb{C}} r_j$  is polynomial in  $O(|\mathbb{F}| \cdot |\mathbb{C}|)$  and the analysis of the performance guarantee in [20] is flawed<sup>1</sup>. To date, the best approximation ratio for FTFP in the literature is 4, established by Yan and Chrobak [21], while the only known lower bound is the 1.463 lower bound for UFL from [6], as UFL is a special case of FTFP. If all demand values  $r_j$  are equal, the problem can be solved by simple scaling and applying LP-rounding

---

<sup>1</sup>Confirmed through private communication with the authors.

algorithms for UFL. This does not affect the approximation ratio, thus achieving ratio 1.575 for this special case (see also [14]).

The main result of this paper is an LP-rounding algorithm for FTFP with approximation ratio 1.575, matching the best ratio for UFL achieved via the LP-rounding method [3] and significantly improving our earlier bound in [21]. In Section 3 we prove that, for the purpose of LP-based approximations, the general FTFP problem can be reduced to the restricted version where all demand values are polynomial in the number of sites. This *demand reduction* trick itself gives us a ratio of 1.7245, since we can then treat an instance of FTFP as an instance of FTFL, by creating a sufficient (but polynomial) number of facilities at each site and using the algorithm from [4].

The reduction to polynomial demands suggests an approach where clients' demands are split into unit demands. These unit demands can be thought of as "unit-demand clients", and a natural approach would be to adapt LP-rounding methods from [8, 5, 3] to this new set of unit-demand clients. Roughly, these algorithms iteratively pick a client that minimizes a certain cost function (that varies for different algorithms) and open one facility in the neighborhood of this client. The remaining clients are then connected to these open facilities. In order for this to work, we also need to convert the optimal fractional solution  $(\mathbf{x}^*, \mathbf{y}^*)$  of the original instance into a solution  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  of the modified instance which then can be used in the LP-rounding process. This can be thought of as partitioning the fractional solution, as each connection value  $x_{ij}^*$  must be divided between the  $r_j$  unit demands of client  $j$  in some way. In Section 4 we formulate a set of properties required for this partitioning to work. For example, one property guarantees that we can connect demands to facilities so that two demands from the same client are connected to different facilities. Then we present our *adaptive partitioning* technique that computes a partitioning with all the desired properties. Using adaptive partitioning we were able to extend the algorithms for UFL from [8, 5, 3] to FTFP. We illustrate the fundamental ideas of our approach in Section 5, showing how they can be used to design an LP-rounding algorithm with ratio 3. In Section 6 we refine the algorithm to improve the approximation ratio to  $1 + 2/e \approx 1.736$ . Finally, in Section 7, we improve it even further to 1.575 – the main result of this paper.

Summarizing, our contributions are two-fold: One, we show that the existing LP-rounding algorithms for UFL can be extended to a much more general problem FTFP, retaining the approximation ratio. We believe that, should even better LP-rounding algorithms be developed for UFL in the future, using our demand reduction and adaptive partitioning methods, it should be possible to extend them to FTFP. In fact, some improvement of the ratio should be achieved by randomizing the scaling parameter  $\gamma$  used in our algorithm, as Li showed in [13] for UFL. (Since the ratio 1.488 for UFL in [13] uses also dual-fitting algorithms [15], we would not obtain the same ratio for FTFP yet using only LP-rounding.)

Two, our ratio of 1.575 is significantly better than the best currently known ratio of 1.7245 for the closely-related FTFL problem. This suggests that in the fault-tolerant scenario the capability of creating additional copies of facilities on the existing sites makes the problem easier from the point of view of approximation.

## 2 The LP Formulation

The FTFP problem has a natural Integer Programming (IP) formulation. Let  $y_i$  represent the number of facilities built at site  $i$  and let  $x_{ij}$  represent the number of connections from client  $j$  to facilities at site  $i$ . If we relax the integrality constraints, we obtain the following LP:

$$\begin{aligned}
& \text{minimize} && \text{cost}(\mathbf{x}, \mathbf{y}) = \sum_{i \in \mathbb{F}} f_i y_i + \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij} \\
& \text{subject to} && y_i - x_{ij} \geq 0 \quad \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \sum_{i \in \mathbb{F}} x_{ij} \geq r_j \quad \forall j \in \mathbb{C} \\
& && x_{ij} \geq 0, y_i \geq 0 \quad \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned} \tag{1}$$

The dual program is:

$$\begin{aligned}
& \text{maximize} && \sum_{j \in \mathbb{C}} r_j \alpha_j \\
& \text{subject to} && \sum_{j \in \mathbb{C}} \beta_{ij} \leq f_i \quad \forall i \in \mathbb{F} \\
& && \alpha_j - \beta_{ij} \leq d_{ij} \quad \forall i \in \mathbb{F}, j \in \mathbb{C} \\
& && \alpha_j \geq 0, \beta_{ij} \geq 0 \quad \forall i \in \mathbb{F}, j \in \mathbb{C}
\end{aligned} \tag{2}$$

In each of our algorithms we will fix some optimal solutions of the LPs (1) and (2) that we will denote by  $(\mathbf{x}^*, \mathbf{y}^*)$  and  $(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$ , respectively.

With  $(\mathbf{x}^*, \mathbf{y}^*)$  fixed, we can define the optimal facility cost as  $F^* = \sum_{i \in \mathbb{F}} f_i y_i^*$  and the optimal connection cost as  $C^* = \sum_{i \in \mathbb{F}, j \in \mathbb{C}} d_{ij} x_{ij}^*$ . Then  $\text{LP}^* = \text{cost}(\mathbf{x}^*, \mathbf{y}^*) = F^* + C^*$  is the joint optimal value of (1) and (2). We can also associate with each client  $j$  its fractional connection cost  $C_j^* = \sum_{i \in \mathbb{F}} d_{ij} x_{ij}^*$ . Clearly,  $C^* = \sum_{j \in \mathbb{C}} C_j^*$ . Throughout the paper we will use notation OPT for the optimal integral solution of (1). OPT is the value we wish to approximate, but, since  $\text{OPT} \geq \text{LP}^*$ , we can instead use  $\text{LP}^*$  to estimate the approximation ratio of our algorithms.

**Completeness and facility splitting.** Define  $(\mathbf{x}^*, \mathbf{y}^*)$  to be *complete* if  $x_{ij}^* > 0$  implies that  $x_{ij}^* = y_i^*$  for all  $i, j$ . In other words, each connection either uses a site fully or not at all. As shown by Chudak and Shmoys [5], we can modify the given instance by adding at most  $|\mathbb{C}|$  sites to obtain an equivalent instance that has a complete optimal solution, where “equivalent” means that the values of  $F^*$ ,  $C^*$  and  $\text{LP}^*$  are not affected. Roughly, the argument is this: We notice that, without loss of generality, for each client  $k$  there exists at most one site  $i$  such that  $0 < x_{ik}^* < y_i^*$ . We can then perform the following *facility splitting* operation on  $i$ : introduce a new site  $i'$ , let  $y_{i'}^* = y_i^* - x_{ik}^*$ , redefine  $y_i^*$  to be  $x_{ik}^*$ , and then for each client  $j$  redistribute  $x_{ij}^*$  so that  $i$  retains as much connection value as possible and  $i'$  receives the rest. Specifically, we set

$$\begin{aligned}
y_{i'}^* &\leftarrow y_i^* - x_{ik}^*, \quad y_i^* \leftarrow x_{ik}^*, \quad \text{and} \\
x_{i'j}^* &\leftarrow \max(x_{ij}^* - x_{ik}^*, 0), \quad x_{ij}^* \leftarrow \min(x_{ij}^*, x_{ik}^*) \quad \text{for all } j \neq k.
\end{aligned}$$

This operation eliminates the partial connection between  $k$  and  $i$  and does not create any new partial connections. Each client can split at most one site and hence we shall have at most  $|\mathbb{C}|$  more sites.

By the above paragraph, without loss of generality we can assume that the optimal fractional solution  $(\mathbf{x}^*, \mathbf{y}^*)$  is complete. This assumption will in fact greatly simplify some of the arguments in the paper. Additionally, we will frequently use the facility splitting operation described above in our algorithms to obtain fractional solutions with desirable properties.

### 3 Reduction to Polynomial Demands

This section presents a *demand reduction* trick that reduces the problem for arbitrary demands to a special case where demands are bounded by  $|\mathbb{F}|$ , the number of sites. (The formal statement is a little more technical – see Theorem 2.) Our algorithms in the sections that follow process individual demands of each client one by one, and thus they critically rely on the demands being bounded polynomially in terms of  $|\mathbb{F}|$  and  $|\mathbb{C}|$  to keep the overall running time polynomial.

The reduction is based on an optimal fractional solution  $(\mathbf{x}^*, \mathbf{y}^*)$  of LP (1). From the optimality of this solution, we can also assume that  $\sum_{i \in \mathbb{F}} x_{ij}^* = r_j$  for all  $j \in \mathbb{C}$ . As explained in Section 2, we can assume that  $(\mathbf{x}^*, \mathbf{y}^*)$  is complete, that is  $x_{ij}^* > 0$  implies  $x_{ij}^* = y_i^*$  for all  $i, j$ . We split this solution into two parts, namely  $(\mathbf{x}^*, \mathbf{y}^*) = (\hat{\mathbf{x}}, \hat{\mathbf{y}}) + (\dot{\mathbf{x}}, \dot{\mathbf{y}})$ , where

$$\begin{aligned} \hat{y}_i &\leftarrow \lfloor y_i^* \rfloor, & \hat{x}_{ij} &\leftarrow \lfloor x_{ij}^* \rfloor & \text{ and} \\ \dot{y}_i &\leftarrow y_i^* - \lfloor y_i^* \rfloor, & \dot{x}_{ij} &\leftarrow x_{ij}^* - \lfloor x_{ij}^* \rfloor \end{aligned}$$

for all  $i, j$ . Now we construct two FTFP instances  $\hat{\mathcal{I}}$  and  $\dot{\mathcal{I}}$  with the same parameters as the original instance, except that the demand of each client  $j$  is  $\hat{r}_j = \sum_{i \in \mathbb{F}} \hat{x}_{ij}$  in instance  $\hat{\mathcal{I}}$  and  $\dot{r}_j = \sum_{i \in \mathbb{F}} \dot{x}_{ij} = r_j - \hat{r}_j$  in instance  $\dot{\mathcal{I}}$ . It is obvious that if we have integral solutions to both  $\hat{\mathcal{I}}$  and  $\dot{\mathcal{I}}$  then, when added together, they form an integral solution to the original instance. Moreover, we have the following lemma.

**Lemma 1.** (i)  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  is a feasible integral solution to instance  $\hat{\mathcal{I}}$ .

(ii)  $(\dot{\mathbf{x}}, \dot{\mathbf{y}})$  is a feasible fractional solution to instance  $\dot{\mathcal{I}}$ .

(iii)  $\dot{r}_j \leq |\mathbb{F}|$  for every client  $j$ .

*Proof.* (i) For feasibility, we need to verify that the constraints of LP (1) are satisfied. Directly from the definition, we have  $\hat{r}_j = \sum_{i \in \mathbb{F}} \hat{x}_{ij}$ . For any  $i$  and  $j$ , by the feasibility of  $(\mathbf{x}^*, \mathbf{y}^*)$  we have  $\hat{x}_{ij} = \lfloor x_{ij}^* \rfloor \leq \lfloor y_i^* \rfloor = \hat{y}_i$ .

(ii) From the definition, we have  $\dot{r}_j = \sum_{i \in \mathbb{F}} \dot{x}_{ij}$ . It remains to show that  $\dot{y}_i \geq \dot{x}_{ij}$  for all  $i, j$ . If  $x_{ij}^* = 0$ , then  $\dot{x}_{ij} = 0$  and we are done. Otherwise, by completeness, we have  $x_{ij}^* = y_i^*$ . Then  $\dot{y}_i = y_i^* - \lfloor y_i^* \rfloor = x_{ij}^* - \lfloor x_{ij}^* \rfloor = \dot{x}_{ij}$ .

(iii) From the definition of  $\dot{x}_{ij}$  we have  $\dot{x}_{ij} < 1$ . Then the bound follows from the definition of  $\dot{r}_j$ .  $\square$

Notice that our construction relies on the completeness assumption; in fact, it is easy to give an example where  $(\dot{\mathbf{x}}, \dot{\mathbf{y}})$  would not be feasible if we used a non-complete optimal solution  $(\mathbf{x}^*, \mathbf{y}^*)$ . Note also that the solutions  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  and  $(\dot{\mathbf{x}}, \dot{\mathbf{y}})$  are in fact optimal for their corresponding instances, for if a better solution to  $\hat{\mathcal{I}}$  or  $\dot{\mathcal{I}}$  existed, it could give us a solution to  $\mathcal{I}$  with a smaller objective value.

**Theorem 2.** Suppose that there is a polynomial-time algorithm  $\mathcal{A}$  that, for any instance of FTFP with maximum demand bounded by  $|\mathbb{F}|$ , computes an integral solution that approximates the fractional optimum of this instance within factor  $\rho \geq 1$ . Then there is a  $\rho$ -approximation algorithm  $\mathcal{A}'$  for FTFP.

*Proof.* Given an FTFP instance with arbitrary demands, Algorithm  $\mathcal{A}'$  works as follows: it solves the LP (1) to obtain a fractional optimal solution  $(\mathbf{x}^*, \mathbf{y}^*)$ , then it constructs instances  $\hat{\mathcal{I}}$  and

$\hat{\mathcal{I}}$  described above, applies algorithm  $\mathcal{A}$  to  $\hat{\mathcal{I}}$ , and finally combines (by adding the values) the integral solution  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  of  $\hat{\mathcal{I}}$  and the integral solution of  $\hat{\mathcal{I}}$  produced by  $\mathcal{A}$ . This clearly produces a feasible integral solution for the original instance  $\mathcal{I}$ . The solution produced by  $\mathcal{A}$  has cost at most  $\rho \cdot \text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ , because  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  is feasible for  $\hat{\mathcal{I}}$ . Thus the cost of  $\mathcal{A}'$  is at most

$$\text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \rho \cdot \text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \leq \rho(\text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \text{cost}(\hat{\mathbf{x}}, \hat{\mathbf{y}})) = \rho \cdot \text{LP}^* \leq \rho \cdot \text{OPT},$$

where the first inequality follows from  $\rho \geq 1$ . This completes the proof.  $\square$

## 4 Adaptive Partitioning

In this section we develop our second technique, which we call *adaptive partitioning*. Given an FTFP instance and an optimal fractional solution  $(\mathbf{x}^*, \mathbf{y}^*)$  to LP (1), we split each client  $j$  into  $r_j$  individual *unit demand points* (or just *demands*), and we split each site  $i$  into no more than  $|\mathbb{F}| + 2R|\mathbb{C}|^2$  *facility points* (or *facilities*), where  $R = \max_{j \in \mathbb{C}} r_j$ . We denote the demand set by  $\overline{\mathbb{C}}$  and the facility set by  $\overline{\mathbb{F}}$ , respectively. We will also partition  $(\mathbf{x}^*, \mathbf{y}^*)$  into a fractional solution  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  for the split instance. We will typically use symbols  $\nu$  and  $\mu$  to index demands and facilities respectively, that is  $\bar{\mathbf{x}} = (\bar{x}_{\mu\nu})$  and  $\bar{\mathbf{y}} = (\bar{y}_\mu)$ . As before, the *neighborhood of a demand*  $\nu$  is  $\overline{N}(\nu) = \{\mu \in \overline{\mathbb{F}} : \bar{x}_{\mu\nu} > 0\}$ . We will use notation  $\nu \in j$  to mean that  $\nu$  is a demand of client  $j$ ; similarly,  $\mu \in i$  means that facility  $\mu$  is on site  $i$ . Different demands of the same client (that is,  $\nu, \nu' \in j$ ) are called *siblings*. Further, we use the convention that  $f_\mu = f_i$  for  $\mu \in i$ ,  $\alpha_\nu^* = \alpha_j^*$  for  $\nu \in j$  and  $d_{\mu\nu} = d_{\mu j} = d_{ij}$  for  $\mu \in i$  and  $\nu \in j$ . We define  $C_\nu^{\text{avg}} = \sum_{\mu \in \overline{N}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu} = \sum_{\mu \in \overline{\mathbb{F}}} d_{\mu\nu} \bar{x}_{\mu\nu}$ . One can think of  $C_\nu^{\text{avg}}$  as the average connection cost of demand  $\nu$ , if we chose a connection to facility  $\mu$  with probability  $\bar{x}_{\mu\nu}$ . In our partitioned fractional solution we guarantee for every  $\nu$  that  $\sum_{\mu \in \overline{\mathbb{F}}} \bar{x}_{\mu\nu} = 1$ .

Some demands in  $\overline{\mathbb{C}}$  will be designated as *primary demands* and the set of primary demands will be denoted by  $P$ . In addition, we will use the overlap structure between demand neighborhoods to define a mapping that assigns each demand  $\nu \in \overline{\mathbb{C}}$  to some primary demand  $\kappa \in P$ . As shown in the rounding algorithms in later sections, for each primary demand we guarantee exactly one open facility in its neighborhood, while for a non-primary demand, there is constant probability that none of its neighbors open. In this case we estimate its connection cost by the distance to the facility opened in its assigned primary demand's neighborhood. For this reason the connection cost of a primary demand must be “small” compared to the non-primary demands assigned to it. We also need sibling demands assigned to different primary demands to satisfy the fault-tolerance requirement. Specifically, this partitioning will be constructed to satisfy a number of properties that are detailed below.

(PS) *Partitioned solution.* Vector  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  is a partition of  $(\mathbf{x}^*, \mathbf{y}^*)$ , with unit-value demands, that is:

1.  $\sum_{\mu \in \overline{\mathbb{F}}} \bar{x}_{\mu\nu} = 1$  for each demand  $\nu \in \overline{\mathbb{C}}$ .
2.  $\sum_{\mu \in i, \nu \in j} \bar{x}_{\mu\nu} = x_{ij}^*$  for each site  $i \in \mathbb{F}$  and client  $j \in \mathbb{C}$ .
3.  $\sum_{\mu \in i} \bar{y}_\mu = y_i^*$  for each site  $i \in \mathbb{F}$ .

(CO) *Completeness.* Solution  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  is complete, that is  $\bar{x}_{\mu\nu} \neq 0$  implies  $\bar{x}_{\mu\nu} = \bar{y}_\mu$ , for all  $\mu \in \overline{\mathbb{F}}, \nu \in \overline{\mathbb{C}}$ .

(PD) *Primary demands.* Primary demands satisfy the following conditions:

1. For any two different primary demands  $\kappa, \kappa' \in P$  we have  $\overline{N}(\kappa) \cap \overline{N}(\kappa') = \emptyset$ .
2. For each site  $i \in \mathbb{F}$ ,  $\sum_{\mu \in i} \sum_{\kappa \in P} \bar{x}_{\mu\kappa} \leq y_i^*$ .
3. Each demand  $\nu \in \overline{\mathbb{C}}$  is assigned to one primary demand  $\kappa \in P$  such that
  - (a)  $\overline{N}(\nu) \cap \overline{N}(\kappa) \neq \emptyset$ , and
  - (b)  $C_\nu^{\text{avg}} + \alpha_\nu^* \geq C_\kappa^{\text{avg}} + \alpha_\kappa^*$ .

(SI) *Siblings.* For any pair  $\nu, \nu'$  of different siblings we have

1.  $\overline{N}(\nu) \cap \overline{N}(\nu') = \emptyset$ .
2. If  $\nu$  is assigned to a primary demand  $\kappa$  then  $\overline{N}(\nu') \cap \overline{N}(\kappa) = \emptyset$ . In particular, by Property (PD.3(a)), this implies that different sibling demands are assigned to different primary demands.

As we shall demonstrate in later sections, these properties allow us to extend known UFL rounding algorithms to obtain an integral solution to our FTFP problem with a matching approximation ratio. Our partitioning is “adaptive” in the sense that it is constructed one demand at a time, and the connection values for the demands of a client depend on the choice of earlier demands, of this or other clients, and their connection values. We would like to point out that the adaptive partitioning process for the 1.575-approximation algorithm (Section 7) is more subtle than that for the 3-approximation (Section 5) and the 1.736-approximation algorithms (Section 6), due to the introduction of close and far neighborhood.

**Implementation of Adaptive Partitioning.** We now describe an algorithm for partitioning the instance and the fractional solution so that the properties (PS), (CO), (PD), and (SI) are satisfied. Recall that  $\overline{\mathbb{F}}$  and  $\overline{\mathbb{C}}$ , respectively, denote the sets of facilities and demands that will be created in this stage, and  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  is the partitioned solution to be computed.

The adaptive partitioning algorithm consists of two phases: Phase 1 is called the partitioning phase and Phase 2 is called the augmenting phase. Phase 1 is done in iterations, where in each iteration we find the “best” client  $j$  and create a new demand  $\nu$  out of it. This demand either becomes a primary demand itself, or it is assigned to some existing primary demand. We call a client  $j$  *exhausted* when all its  $r_j$  demands have been created and assigned to some primary demands. Phase 1 completes when all clients are exhausted. In Phase 2 we ensure that every demand has a total connection values  $\bar{x}_{\mu\nu}$  equal to 1, that is condition (PS.1).

For each site  $i$  we will initially create one “big” facility  $\mu$  with initial value  $\bar{y}_\mu = y_i^*$ . While we partition the instance, creating new demands and connections, this facility may end up being split into more facilities to preserve completeness of the fractional solution. Also, we will gradually decrease the fractional connection vector for each client  $j$ , to account for the demands already created for  $j$  and their connection values. These decreased connection values will be stored in an auxiliary vector  $\tilde{\mathbf{x}}$ . The intuition is that  $\tilde{\mathbf{x}}$  represents the part of  $\mathbf{x}^*$  that still has not been allocated to existing demands and future demands can use  $\tilde{\mathbf{x}}$  for their connections. For technical reasons,  $\tilde{\mathbf{x}}$  will be indexed by facilities (rather than sites) and clients, that is  $\tilde{\mathbf{x}} = (\tilde{x}_{\mu j})$ . At the beginning, we set  $\tilde{x}_{\mu j} \leftarrow x_{ij}^*$  for each  $j \in \mathbb{C}$ , where  $\mu \in i$  is the single facility created initially at site  $i$ . At each step, whenever we create a new demand  $\nu$  for a client  $j$ , we will define its values  $\bar{x}_{\mu\nu}$  and appropriately

reduce the values  $\tilde{x}_{\mu j}$ , for all facilities  $\mu$ . We will deal with two types of neighborhoods, with respect to  $\tilde{\mathbf{x}}$  and  $\bar{\mathbf{x}}$ , that is  $\tilde{N}(j) = \{\mu \in \bar{\mathbb{F}} : \tilde{x}_{\mu j} > 0\}$  for  $j \in \mathbb{C}$  and  $\bar{N}(\nu) = \{\mu \in \bar{\mathbb{F}} : \bar{x}_{\mu\nu} > 0\}$  for  $\nu \in \bar{\mathbb{C}}$ . During this process we preserve the completeness (CO) of the fractional solutions  $\tilde{\mathbf{x}}$  and  $\bar{\mathbf{x}}$ . More precisely, the following properties will hold for every facility  $\mu$  after every iteration:

- (c1) For each demand  $\nu$  either  $\bar{x}_{\mu\nu} = 0$  or  $\bar{x}_{\mu\nu} = \bar{y}_\mu$ . This is the same condition as condition (CO), yet we repeat it here as (c1) needs to hold after every iteration, while condition (CO) only applies to the final partitioned fractional solution  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ .
- (c2) For each client  $j$ , either  $\tilde{x}_{\mu j} = 0$  or  $\tilde{x}_{\mu j} = \bar{y}_\mu$ .

A full description of the algorithm is given in Pseudocode 1. Initially, the set  $U$  of non-exhausted clients contains all clients, the set  $\bar{\mathbb{C}}$  of demands is empty, the set  $\bar{\mathbb{F}}$  of facilities consists of one facility  $\mu$  on each site  $i$  with  $\bar{y}_\mu = y_i^*$ , and the set  $P$  of primary demands is empty (Lines 1–4). In one iteration of the while loop (Lines 5–8), for each client  $j$  we compute a quantity called  $\text{tcc}(j)$  (tentative connection cost), that represents the average distance from  $j$  to the set  $\tilde{N}_1(j)$  of the nearest facilities  $\mu$  whose total connection value to  $j$  (the sum of  $\tilde{x}_{\mu j}$ 's) equals 1. This set is computed by Procedure NEARESTUNITCHUNK() (see Pseudocode 2, Lines 1–9), which adds facilities to  $\tilde{N}_1(j)$  in order of nondecreasing distance, until the total connection value is exactly 1. (The procedure actually uses the  $\bar{y}_\mu$  values, which are equal to the connection values, by the completeness condition (c2).) This may require splitting the last added facility and adjusting the connection values so that conditions (c1) and (c2) are preserved.

The next step is to pick a client  $p$  with minimum  $\text{tcc}(p) + \alpha_p^*$  and create a demand  $\nu$  for  $p$  (Lines 9–10). If  $\tilde{N}_1(p)$  overlaps the neighborhood of some existing primary demand  $\kappa$  (if there are multiple such  $\kappa$ 's, pick any of them), we assign  $\nu$  to  $\kappa$ , and  $\nu$  acquires all the connection values  $\tilde{x}_{\mu p}$  between client  $p$  and facility  $\mu$  in  $\tilde{N}(p) \cap \bar{N}(\kappa)$  (Lines 11–13). Note that although we check for overlap with  $\tilde{N}_1(p)$ , we then move all facilities in the intersection with  $\tilde{N}(p)$ , a bigger set, into  $\bar{N}(\nu)$ . The other case is when  $\tilde{N}_1(p)$  is disjoint from the neighborhoods of all existing primary demands. Then, in Lines 15–16,  $\nu$  becomes itself a primary demand and we assign  $\nu$  to itself. It also inherits the connection values to all facilities  $\mu \in \tilde{N}_1(p)$  from  $p$  (recall that  $\tilde{x}_{\mu p} = \bar{y}_\mu$ ), with all other  $\bar{x}_{\mu\nu}$  values set to 0.

At this point all primary demands satisfy Property (PS.1), but this may not be true for non-primary demands. For those demands we still may need to adjust the  $\bar{x}_{\mu\nu}$  values so that the total connection value for  $\nu$ , that is  $\text{conn}(\nu) \stackrel{\text{def}}{=} \sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu\nu}$ , is equal 1. This is accomplished by Procedure AUGMENTTOUNIT() (definition in Pseudocode 2, Lines 10–21) that allocates to  $\nu \in j$  some of the remaining connection values  $\tilde{x}_{\mu j}$  of client  $j$  (Lines 19–21). AUGMENTTOUNIT() will repeatedly pick any facility  $\eta$  with  $\tilde{x}_{\eta j} > 0$ . If  $\tilde{x}_{\eta j} \leq 1 - \text{conn}(\nu)$ , then the connection value  $\tilde{x}_{\eta j}$  is reassigned to  $\nu$ . Otherwise,  $\tilde{x}_{\eta j} > 1 - \text{conn}(\nu)$ , in which case we split  $\eta$  so that connecting  $\nu$  to one of the created copies of  $\eta$  will make  $\text{conn}(\nu)$  equal 1, and we'll be done.

Notice that we start with  $|\mathbb{F}|$  facilities and in each iteration of the while loop in Line 5 each client causes at most one split. We have a total of no more than  $R|\mathbb{C}|$  iterations as in each iteration we create one demand. (Recall that  $R = \max_j r_j$ .) In Phase 2 we do an augment step for each demand  $\nu$  and this creates no more than  $R|\mathbb{C}|$  new facilities. So the total number of facilities we created will be at most  $|\mathbb{F}| + R|\mathbb{C}|^2 + R|\mathbb{C}| \leq |\mathbb{F}| + 2R|\mathbb{C}|^2$ , which is polynomial in  $|\mathbb{F}| + |\mathbb{C}|$  due to our earlier bound on  $R$ .



---

**Pseudocode 1** Algorithm: Adaptive Partitioning

---

**Input:**  $\mathbb{F}, \mathbb{C}, (\mathbf{x}^*, \mathbf{y}^*)$ **Output:**  $\bar{\mathbb{F}}, \bar{\mathbb{C}}, (\bar{\mathbf{x}}, \bar{\mathbf{y}})$  $\triangleright$  Unspecified  $\bar{x}_{\mu\nu}$ 's and  $\tilde{x}_{\mu j}$ 's are assumed to be 0

```
1:  $\tilde{\mathbf{r}} \leftarrow \mathbf{r}, U \leftarrow \mathbb{C}, \bar{\mathbb{F}} \leftarrow \emptyset, \bar{\mathbb{C}} \leftarrow \emptyset, P \leftarrow \emptyset$   $\triangleright$  Phase 1
2: for each site  $i \in \mathbb{F}$  do
3:   create a facility  $\mu$  at  $i$  and add  $\mu$  to  $\bar{\mathbb{F}}$ 
4:    $\bar{y}_\mu \leftarrow y_i^*$  and  $\tilde{x}_{\mu j} \leftarrow x_{ij}^*$  for each  $j \in \mathbb{C}$ 
5: while  $U \neq \emptyset$  do
6:   for each  $j \in U$  do
7:      $\tilde{N}_1(j) \leftarrow \text{NEARESTUNITCHUNK}(j, \bar{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}})$   $\triangleright$  see Pseudocode 2
8:      $\text{tcc}(j) \leftarrow \sum_{\mu \in \tilde{N}_1(j)} d_{\mu j} \cdot \tilde{x}_{\mu j}$ 
9:      $p \leftarrow \arg \min_{j \in U} \{\text{tcc}(j) + \alpha_j^*\}$ 
10:    create a new demand  $\nu$  for client  $p$ 
11:    if  $\tilde{N}_1(p) \cap \bar{N}(\kappa) \neq \emptyset$  for some primary demand  $\kappa \in P$  then
12:      assign  $\nu$  to  $\kappa$ 
13:       $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$  and  $\tilde{x}_{\mu p} \leftarrow 0$  for each  $\mu \in \tilde{N}(p) \cap \bar{N}(\kappa)$ 
14:    else
15:      make  $\nu$  primary,  $P \leftarrow P \cup \{\nu\}$ , assign  $\nu$  to itself
16:      set  $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$  and  $\tilde{x}_{\mu p} \leftarrow 0$  for each  $\mu \in \tilde{N}_1(p)$ 
17:       $\bar{\mathbb{C}} \leftarrow \bar{\mathbb{C}} \cup \{\nu\}, \tilde{r}_p \leftarrow \tilde{r}_p - 1$ 
18:      if  $\tilde{r}_p = 0$  then  $U \leftarrow U \setminus \{p\}$ 
19: for each client  $j \in \mathbb{C}$  do  $\triangleright$  Phase 2
20:   for each demand  $\nu \in j$  do  $\triangleright$  each client  $j$  has  $r_j$  demands
21:    if  $\sum_{\mu \in \bar{N}(\nu)} \bar{x}_{\mu\nu} < 1$  then  $\text{AUGMENTTOUNIT}(\nu, j, \bar{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}})$   $\triangleright$  see Pseudocode 2
```

---

*Correctness.* We now show that all the required properties (PS), (CO), (PD) and (SI) are satisfied by the above construction.

Properties (PS) and (CO) follow directly from the algorithm. (CO) is implied by the completeness condition (c1) that the algorithm maintains after each iteration. Condition (PS.1) is a result of calling Procedure AUGMENTTOUNIT() in Line 21. To see that (PS.2) holds, note that at each step the algorithm maintains the invariant that, for every  $i \in \mathbb{F}$  and  $j \in \mathbb{C}$ , we have  $\sum_{\mu \in i} \sum_{\nu \in j} \bar{x}_{\mu\nu} + \sum_{\mu \in i} \tilde{x}_{\mu j} = x_{ij}^*$ . In the end, we will create  $r_j$  demands for each client  $j$ , with each demand  $\nu \in j$  satisfying (PS.1), and thus  $\sum_{\nu \in j} \sum_{\mu \in \mathbb{F}} \bar{x}_{\mu\nu} = r_j$ . This implies that  $\tilde{x}_{\mu j} = 0$  for every facility  $\mu \in \bar{\mathbb{F}}$ , and PS(2) follows. PS(3) holds because every time we split a facility  $\mu$  into  $\mu'$  and  $\mu''$ , the sum of  $\bar{y}_{\mu'}$  and  $\bar{y}_{\mu''}$  is equal to the old value of  $\bar{y}_\mu$ .

Now we deal with properties in group (PD). First, (PD.1) follows directly from the algorithm, Pseudocode 1 (Lines 14–16), since every primary demand has its neighborhood fixed when created, and that neighborhood is disjoint from those of the existing primary demands.

Property (PD.2) follows from (PD.1), (CO) and (PS.3). In more detail, it can be justified as follows. By (PD.1), for each  $\mu \in i$  there is at most one  $\kappa \in P$  with  $\bar{x}_{\mu\kappa} > 0$  and we have  $\bar{x}_{\mu\kappa} = \bar{y}_\mu$  due to (CO). Let  $K \subseteq i$  be the set of those  $\mu$ 's for which such  $\kappa \in P$  exists, and denote this  $\kappa$  by  $\kappa_\mu$ . Then, using conditions (CO) and (PS.3), we have  $\sum_{\mu \in i} \sum_{\kappa \in P} \bar{x}_{\mu\kappa} = \sum_{\mu \in K} \bar{x}_{\mu\kappa_\mu} = \sum_{\mu \in K} \bar{y}_\mu \leq \sum_{\mu \in i} \bar{y}_\mu = y_i^*$ .

---

**Pseudocode 2** Helper functions used in Pseudocode 1

---

```

1: function NEARESTUNITCHUNK( $j, \bar{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}}$ )                                 $\triangleright$  upon return,  $\sum_{\mu \in \tilde{N}_1(j)} \tilde{x}_{\mu j} = 1$ 
2:   Let  $\tilde{N}(j) = \{\mu_1, \dots, \mu_q\}$  where  $d_{\mu_1 j} \leq d_{\mu_2 j} \leq \dots \leq d_{\mu_q j}$ 
3:   Let  $l$  be such that  $\sum_{k=1}^l \bar{y}_{\mu_k} \geq 1$  and  $\sum_{k=1}^{l-1} \bar{y}_{\mu_k} < 1$ 
4:   Create a new facility  $\sigma$  at the same site as  $\mu_l$  and add it to  $\bar{\mathbb{F}}$                                  $\triangleright$  split  $\mu_l$ 
5:   Set  $\bar{y}_\sigma \leftarrow \sum_{k=1}^l \bar{y}_{\mu_k} - 1$  and  $\bar{y}_{\mu_l} \leftarrow \bar{y}_{\mu_l} - \bar{y}_\sigma$ 
6:   For each  $\nu \in \mathbb{C}$  with  $\bar{x}_{\mu_l \nu} > 0$  set  $\bar{x}_{\mu_l \nu} \leftarrow \bar{y}_{\mu_l}$  and  $\bar{x}_{\sigma \nu} \leftarrow \bar{y}_\sigma$ 
7:   For each  $j' \in \mathbb{C}$  with  $\tilde{x}_{\mu_l j'} > 0$  (including  $j$ ) set  $\tilde{x}_{\mu_l j'} \leftarrow \bar{y}_{\mu_l}$  and  $\tilde{x}_{\sigma j'} \leftarrow \bar{y}_\sigma$ 
8:   (All other new connection values are set to 0)
9:   return  $\tilde{N}_1(j) = \{\mu_1, \dots, \mu_{l-1}, \mu_l\}$ 

10: function AUGMENTTOUNIT( $\nu, j, \bar{\mathbb{F}}, \tilde{\mathbf{x}}, \bar{\mathbf{x}}, \bar{\mathbf{y}}$ )                                 $\triangleright \nu$  is a demand of client  $j$ 
11:   while  $\sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu \nu} < 1$  do                                 $\triangleright$  upon return,  $\sum_{\mu \in \bar{N}(\nu)} \bar{x}_{\mu \nu} = 1$ 
12:     Let  $\eta$  be any facility such that  $\tilde{x}_{\eta j} > 0$ 
13:     if  $1 - \sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu \nu} \geq \tilde{x}_{\eta j}$  then
14:        $\bar{x}_{\eta \nu} \leftarrow \tilde{x}_{\eta j}, \tilde{x}_{\eta j} \leftarrow 0$ 
15:     else
16:       Create a new facility  $\sigma$  at the same site as  $\eta$  and add it to  $\bar{\mathbb{F}}$                                  $\triangleright$  split  $\eta$ 
17:       Let  $\bar{y}_\sigma \leftarrow 1 - \sum_{\mu \in \bar{\mathbb{F}}} \bar{x}_{\mu \nu}, \bar{y}_\eta \leftarrow \bar{y}_\eta - \bar{y}_\sigma$ 
18:       Set  $\bar{x}_{\sigma \nu} \leftarrow \bar{y}_\sigma, \bar{x}_{\eta \nu} \leftarrow 0, \tilde{x}_{\eta j} \leftarrow \bar{y}_\eta, \tilde{x}_{\sigma j} \leftarrow 0$ 
19:       For each  $\nu' \neq \nu$  with  $\bar{x}_{\eta \nu'} > 0$  set  $\bar{x}_{\eta \nu'} \leftarrow \bar{y}_\eta, \bar{x}_{\sigma \nu'} \leftarrow \bar{y}_\sigma$ 
20:       For each  $j' \neq j$  with  $\tilde{x}_{\eta j'} > 0$  set  $\tilde{x}_{\eta j'} \leftarrow \bar{y}_\eta, \tilde{x}_{\sigma j'} \leftarrow \bar{y}_\sigma$ 
21:       (All other new connection values are set to 0)

```

---

Property (PD.3(a)) follows from the way the algorithm assigns primary demands. When demand  $\nu$  of client  $p$  is assigned to a primary demand  $\kappa$  in Lines 11–13 of Pseudocode 1, we move all facilities in  $\tilde{N}(p) \cap \bar{N}(\kappa)$  (the intersection is nonempty) into  $\bar{N}(\nu)$ , and we never remove a facility from  $\bar{N}(\nu)$ . We postpone the proof for (PD.3(b)) to Lemma 5.

Finally we argue that the properties in group (SI) hold. (SI.1) is easy, since for any client  $j$ , each facility  $\mu$  is added to the neighborhood of at most one demand  $\nu \in j$ , by setting  $\bar{x}_{\mu \nu}$  to  $\bar{y}_\mu$ , while other siblings  $\nu'$  of  $\nu$  have  $\bar{x}_{\mu \nu'} = 0$ . Note that right after a demand  $\nu \in p$  is created, its neighborhood is disjoint from the neighborhood of  $p$ , that is  $\bar{N}(\nu) \cap \tilde{N}(p) = \emptyset$ , by Lines 11–13 of the algorithm. Thus all demands of  $p$  created later will have neighborhoods disjoint from the set  $\bar{N}(\nu)$  before the augmenting phase 2. Furthermore, Procedure AUGMENTTOUNIT() preserves this property, because when it adds a facility to  $\bar{N}(\nu)$  then it removes it from  $\tilde{N}(p)$ , and in case of splitting, one resulting facility is added to  $\bar{N}(\nu)$  and the other to  $\tilde{N}(p)$ . Property (SI.2) is shown below in Lemma 3.

It remains to show Properties (PD.3(b)) and (SI.2). We show them in the lemmas below, thus completing the description of our adaptive partition process.

**Lemma 3.** *Property (SI.2) holds after the Adaptive Partitioning stage.*

*Proof.* Let  $\nu_1, \dots, \nu_{r_j}$  be the demands of a client  $j \in \mathbb{C}$ , listed in the order of creation, and, for each  $q = 1, 2, \dots, r_j$ , denote by  $\kappa_q$  the primary demand that  $\nu_q$  is assigned to. After the completion of Phase 1 of Pseudocode 1 (Lines 5–18), we have  $\bar{N}(\nu_s) \subseteq \bar{N}(\kappa_s)$  for  $s = 1, \dots, r_j$ . Since any two

primary demands have disjoint neighborhoods, we have  $\overline{N}(\nu_s) \cap \overline{N}(\kappa_q) = \emptyset$  for any  $s \neq q$ , that is Property (SI.2) holds right after Phase 1.

After Phase 1 all neighborhoods  $\overline{N}(\kappa_s)$ ,  $s = 1, \dots, r_j$  have already been fixed and they do not change in Phase 2. None of the facilities in  $\tilde{N}(j)$  appear in any of  $\overline{N}(\kappa_s)$  for  $s = 1, \dots, r_j$ , by the way we allocate facilities in Lines 13 and 16. Therefore during the augmentation process in Phase 2, when we add facilities from  $\tilde{N}(j)$  to  $\overline{N}(\nu)$ , for some  $\nu \in j$  (Line 19–21 of Pseudocode 1), all the required disjointness conditions will be preserved.  $\square$

We need one more lemma before proving our last property (PD.3(b)). For a client  $j$  and a demand  $\nu$ , we use notation  $\text{tcc}^\nu(j)$  for the value of  $\text{tcc}(j)$  at the time when  $\nu$  was created. (It is not necessary that  $\nu \in j$  but we assume that  $j$  is not exhausted at that time.)

**Lemma 4.** *Let  $\eta$  and  $\nu$  be two demands, with  $\eta$  created not later than  $\nu$ , and let  $j \in \mathbb{C}$  be a client that is not exhausted when  $\nu$  is created. Then we have*

- (a)  $\text{tcc}^\eta(j) \leq \text{tcc}^\nu(j)$ , and
- (b) if  $\nu \in j$  then  $\text{tcc}^\eta(j) \leq C_\nu^{\text{avg}}$ .

*Proof.* We focus first on the time when demand  $\eta$  is about to be created, right after the call to `NEARESTUNITCHUNK()` in Pseudocode 1, Line 7. Let  $\tilde{N}(j) = \{\mu_1, \dots, \mu_q\}$  with all facilities  $\mu_s$  ordered according to nondecreasing distance from  $j$ . Consider the following linear program:

$$\begin{aligned} & \text{minimize} && \sum_s d_{\mu_s j} z_s \\ & \text{subject to} && \sum_s z_s \geq 1 \\ & && 0 \leq z_s \leq \tilde{x}_{\mu_s j} \quad \text{for all } s \end{aligned}$$

This is a fractional minimum knapsack covering problem (with knapsack size equal 1) and its optimal fractional solution is the greedy solution, whose value is exactly  $\text{tcc}^\eta(j)$ .

On the other hand, we claim that  $\text{tcc}^\nu(j)$  can be thought of as the value of some feasible solution to this linear program, and that the same is true for  $C_\nu^{\text{avg}}$  if  $\nu \in j$ . Indeed, each of these quantities involves some later values  $\tilde{x}_{\mu j}$ , where  $\mu$  could be one of the facilities  $\mu_s$  or a new facility obtained from splitting. For each  $s$ , however, the sum of all values  $\tilde{x}_{\mu j}$ , over the facilities  $\mu$  that were split from  $\mu_s$ , cannot exceed the value  $\tilde{x}_{\mu_s j}$  at the time when  $\eta$  was created, because splitting facilities preserves this sum and creating new demands for  $j$  can only decrease it. Therefore both quantities  $\text{tcc}^\nu(j)$  and  $C_\nu^{\text{avg}}$  (for  $\nu \in j$ ) correspond to some choice of the  $z_s$  variables (adding up to 1), and the lemma follows.  $\square$

**Lemma 5.** *Property (PD.3(b)) holds after the Adaptive Partitioning stage.*

*Proof.* Suppose that demand  $\nu \in j$  is assigned to some primary demand  $\kappa \in p$ . Then

$$C_\kappa^{\text{avg}} + \alpha_\kappa^* = \text{tcc}^\kappa(p) + \alpha_p^* \leq \text{tcc}^\kappa(j) + \alpha_j^* \leq C_\nu^{\text{avg}} + \alpha_\nu^*.$$

We now justify this derivation. By definition we have  $\alpha_\kappa^* = \alpha_p^*$ . Further, by the algorithm, if  $\kappa$  is a primary demand of client  $p$ , then  $C_\kappa^{\text{avg}}$  is equal to  $\text{tcc}(p)$  computed when  $\kappa$  is created, which is exactly  $\text{tcc}^\kappa(p)$ . Thus the first equation is true. The first inequality follows from the choice of  $p$  in Line 9 in Pseudocode 1. The last inequality holds because  $\alpha_j^* = \alpha_\nu^*$  (due to  $\nu \in j$ ), and because  $\text{tcc}^\kappa(j) \leq C_\nu^{\text{avg}}$ , which follows from Lemma 4.  $\square$

We have thus proved that all properties (PS), (CO), (PD) and (SI) hold for our partitioned fractional solution  $(\bar{x}, \bar{y})$ . In the following sections we show how to use these properties to round the fractional solution to an approximate integral solution. For the 3-approximation algorithm (Section 5) and the 1.736-approximation algorithm (Section 6), the first phase of the algorithm is exactly the same partition process as described above. However, the 1.575-approximation algorithm (Section 7) demands a more sophisticated partitioning process as the interplay between close and far neighborhood of sibling demands result in more delicate properties that our partitioned fractional solution must satisfy.

## 5 Algorithm EGUP with Ratio 3

With the partitioned FTFP instance and its associated fractional solution in place, we now begin to introduce our rounding algorithms. The algorithm we describe in this section achieves ratio 3. Although this is still quite far from our best ratio 1.575 that we derive later, we include this algorithm in the paper to illustrate, in a relatively simple setting, how the properties of our partitioned fractional solution are used in rounding it to an integral solution with cost not too far away from an optimal solution. The rounding approach we use here is an extension of the corresponding method for UFL described in [8].

**Algorithm EGUP.** At a high level, we would open exactly one facility for each primary demand  $\kappa$ , and each non-primary demand is connected to the facility opened for the primary demand it was assigned to.

More precisely, we apply a rounding process, guided by the fractional values  $(\bar{y}_\mu)$  and  $(\bar{x}_{\mu\nu})$ , that produces an integral solution. This integral solution is obtained by choosing a subset of facilities in  $\bar{\mathbb{F}}$  to open, and for each demand in  $\bar{\mathbb{C}}$ , specifying an open facility that this demand will be connected to. For each primary demand  $\kappa \in P$ , we want to open one facility  $\phi(\kappa) \in \bar{N}(\kappa)$ . To this end, we use randomization: for each  $\mu \in \bar{N}(\kappa)$ , we choose  $\phi(\kappa) = \mu$  with probability  $\bar{x}_{\mu\kappa}$ , ensuring that exactly one  $\mu \in \bar{N}(\kappa)$  is chosen. Note that  $\sum_{\mu \in \bar{N}(\kappa)} \bar{x}_{\mu\kappa} = 1$ , so this distribution is well-defined. We open this facility  $\phi(\kappa)$  and connect to  $\phi(\kappa)$  all demands that are assigned to  $\kappa$ .

In our description above, the algorithm is presented as a randomized algorithm. It can be derandomized using the method of conditional expectations, which is commonly used in approximation algorithms for facility location problems and standard enough that presenting it here would be redundant. Readers less familiar with this field are recommended to consult [5], where the method of conditional expectations is applied in a context very similar to ours.

**Analysis.** We now bound the expected facility cost and connection cost by establishing the two lemmas below.

**Lemma 6.** *The expectation of facility cost  $F_{\text{EGUP}}$  of our solution is at most  $F^*$ .*

*Proof.* By Property (PD.1), the neighborhoods of primary demands are disjoint. Also, for any primary demand  $\kappa \in P$ , the probability that a facility  $\mu \in \bar{N}(\kappa)$  is chosen as the open facility  $\phi(\kappa)$

is  $\bar{x}_{\mu\kappa}$ . Hence the expected total facility cost is

$$\begin{aligned}\text{Exp}[F_{\text{EGUP}}] &= \sum_{\kappa \in P} \sum_{\mu \in \bar{N}(\kappa)} f_{\mu} \bar{x}_{\mu\kappa} \\ &= \sum_{\kappa \in P} \sum_{\mu \in \mathbb{F}} f_{\mu} \bar{x}_{\mu\kappa} \\ &= \sum_{i \in \mathbb{F}} f_i \sum_{\mu \in i} \sum_{\kappa \in P} \bar{x}_{\mu\kappa} \\ &\leq \sum_{i \in \mathbb{F}} f_i y_i^* = F^*,\end{aligned}$$

where the inequality follows from Property (PD.2).  $\square$

**Lemma 7.** *The expectation of connection cost  $C_{\text{EGUP}}$  of our solution is at most  $C^* + 2 \cdot \text{LP}^*$ .*

*Proof.* For a primary demand  $\kappa$ , its expected connection cost is  $C_{\kappa}^{\text{avg}}$  because we choose facility  $\mu$  with probability  $\bar{x}_{\mu\kappa}$ .

Consider a non-primary demand  $\nu$  assigned to a primary demand  $\kappa \in P$ . Let  $\mu$  be any facility in  $\bar{N}(\nu) \cap \bar{N}(\kappa)$ . Since  $\mu$  is in both  $\bar{N}(\nu)$  and  $\bar{N}(\kappa)$ , we have  $d_{\mu\nu} \leq \alpha_{\nu}^*$  and  $d_{\mu\kappa} \leq \alpha_{\kappa}^*$  (This follows from the complementary slackness conditions since  $\alpha_{\nu}^* = \beta_{\mu\nu}^* + d_{\mu\nu}$  for each  $\mu \in \bar{N}(\nu)$ ). Thus, applying the triangle inequality, for any fixed choice of facility  $\phi(\kappa)$  we have

$$d_{\phi(\kappa)\nu} \leq d_{\phi(\kappa)\kappa} + d_{\mu\kappa} + d_{\mu\nu} \leq d_{\phi(\kappa)\kappa} + \alpha_{\kappa}^* + \alpha_{\nu}^*.$$

Therefore the expected distance from  $\nu$  to its facility  $\phi(\kappa)$  is

$$\begin{aligned}\text{Exp}[d_{\phi(\kappa)\nu}] &\leq C_{\kappa}^{\text{avg}} + \alpha_{\kappa}^* + \alpha_{\nu}^* \\ &\leq C_{\nu}^{\text{avg}} + \alpha_{\nu}^* + \alpha_{\nu}^* = C_{\nu}^{\text{avg}} + 2\alpha_{\nu}^*,\end{aligned}$$

where the second inequality follows from Property (PD.3(b)). From the definition of  $C_{\nu}^{\text{avg}}$  and Property (PS.2), for any  $j \in \mathbb{C}$  we have

$$\begin{aligned}\sum_{\nu \in j} C_{\nu}^{\text{avg}} &= \sum_{\nu \in j} \sum_{\mu \in \mathbb{F}} d_{\mu\nu} \bar{x}_{\mu\nu} \\ &= \sum_{i \in \mathbb{F}} d_{ij} \sum_{\nu \in j} \sum_{\mu \in i} \bar{x}_{\mu\nu} \\ &= \sum_{i \in \mathbb{F}} d_{ij} x_{ij}^* = C_j^*.\end{aligned}$$

Thus, summing over all demands, the expected total connection cost is

$$\begin{aligned}\text{Exp}[C_{\text{EGUP}}] &\leq \sum_{j \in \mathbb{C}} \sum_{\nu \in j} (C_{\nu}^{\text{avg}} + 2\alpha_{\nu}^*) \\ &= \sum_{j \in \mathbb{C}} (C_j^* + 2r_j \alpha_j^*) = C^* + 2 \cdot \text{LP}^*,\end{aligned}$$

completing the proof of the lemma.  $\square$

**Theorem 8.** *Algorithm EGUP is a 3-approximation algorithm.*

*Proof.* By Property (SI.2), different demands from the same client are assigned to different primary demands, and by (PD.1) each primary demand opens a different facility. This ensures that our solution is feasible, namely each client  $j$  is connected to  $r_j$  different facilities (some possibly located on the same site). As for the total cost, Lemma 6 and Lemma 7 imply that the total cost is at most  $F^* + C^* + 2 \cdot \text{LP}^* = 3 \cdot \text{LP}^* \leq 3 \cdot \text{OPT}$ .  $\square$

## 6 Algorithm ECHS with Ratio 1.736

In this section we improve the approximation ratio to  $1 + 2/e \approx 1.736$ . The improvement comes from a slightly modified rounding process and refined analysis. Note that the facility opening cost of Algorithm EGUP does not exceed that of the fractional optimum solution, while the connection cost could be far from the optimum, since we connect a non-primary demand to a facility in the neighborhood of its assigned primary demand and then estimate the distance using the triangle inequality. The basic idea to improve the estimate of the connection cost, following the approach of Chudak and Shmoys [5], is to connect each non-primary demand to its nearest neighbor when one is available, and to only use the facility opened by its assigned primary demand when none of its neighbors is open.

**Algorithm ECHS.** As before, the algorithm starts by solving the linear program and applying the adaptive partitioning algorithm described in Section 4 to obtain a partitioned solution  $(\bar{x}, \bar{y})$ . Then we apply the rounding process to compute an integral solution (see Pseudocode 3).

We start, as before, by opening exactly one facility  $\phi(\kappa)$  in the neighborhood of each primary demand  $\kappa$  (Line 2). For any non-primary demand  $\nu$  assigned to  $\kappa$ , we refer to  $\phi(\kappa)$  as the *target* facility of  $\nu$ . In Algorithm EGUP,  $\nu$  was connected to  $\phi(\kappa)$ , but in Algorithm ECHS we may be able to find an open facility in  $\nu$ 's neighborhood and connect  $\nu$  to this facility. Specifically, the two changes in the algorithm are as follows:

- (1) Each facility  $\mu$  that is not in the neighborhood of some primary demand is opened, independently, with probability  $\bar{y}_\mu$  (Lines 4–5). Notice that if  $\bar{y}_\mu > 0$  then, due to completeness of the partitioned fractional solution, we have  $\bar{y}_\mu = \bar{x}_{\mu\nu}$  for some demand  $\nu$ . This implies that  $\bar{y}_\mu \leq 1$ , because  $\bar{x}_{\mu\nu} \leq 1$ , by (PS.1).
- (2) When connecting demands to facilities, a primary demand  $\kappa$  is connected to the only facility  $\phi(\kappa)$  opened in its neighborhood, as before (Line 3). For a non-primary demand  $\nu$ , if its neighborhood  $\bar{N}(\nu)$  has an open facility, we connect  $\nu$  to the closest open facility in  $\bar{N}(\nu)$  (Line 8). Otherwise, we connect  $\nu$  to its target facility (Line 10).

---

### Pseudocode 3 Algorithm ECHS: Constructing Integral Solution

---

```

1: for each  $\kappa \in P$  do
2:   choose one  $\phi(\kappa) \in \bar{N}(\kappa)$ , with each  $\mu \in \bar{N}(\kappa)$  chosen as  $\phi(\kappa)$  with probability  $\bar{y}_\mu$ 
3:   open  $\phi(\kappa)$  and connect  $\kappa$  to  $\phi(\kappa)$ 
4: for each  $\mu \in \bar{\mathbb{F}} - \bigcup_{\kappa \in P} \bar{N}(\kappa)$  do
5:   open  $\mu$  with probability  $\bar{y}_\mu$  (independently)
6: for each non-primary demand  $\nu \in \bar{\mathbb{C}}$  do
7:   if any facility in  $\bar{N}(\nu)$  is open then
8:     connect  $\nu$  to the nearest open facility in  $\bar{N}(\nu)$ 
9:   else
10:    connect  $\nu$  to  $\phi(\kappa)$  where  $\kappa$  is  $\nu$ 's assigned primary demand

```

---

**Analysis.** We shall first argue that the integral solution thus constructed is feasible, and then we bound the total cost of the solution. Regarding feasibility, the only constraint that is not explicitly enforced by the algorithm is the fault-tolerance requirement; namely that each client  $j$  is connected to  $r_j$  different facilities. Let  $\nu$  and  $\nu'$  be two different sibling demands of client  $j$  and let their assigned primary demands be  $\kappa$  and  $\kappa'$  respectively. Due to (SI.2) we know  $\kappa \neq \kappa'$ . From (SI.1) we have  $\overline{N}(\nu) \cap \overline{N}(\nu') = \emptyset$ . From (SI.2), we have  $\overline{N}(\nu) \cap \overline{N}(\kappa') = \emptyset$  and  $\overline{N}(\nu') \cap \overline{N}(\kappa) = \emptyset$ . From (PD.1) we have  $\overline{N}(\kappa) \cap \overline{N}(\kappa') = \emptyset$ . It follows that  $(\overline{N}(\nu) \cup \overline{N}(\kappa)) \cap (\overline{N}(\nu') \cup \overline{N}(\kappa')) = \emptyset$ . Since the algorithm connects  $\nu$  to some facility in  $\overline{N}(\nu) \cup \overline{N}(\kappa)$  and  $\nu'$  to some facility in  $\overline{N}(\nu') \cup \overline{N}(\kappa')$ ,  $\nu$  and  $\nu'$  will be connected to different facilities.

We now show that the expected cost of the computed solution is bounded by  $(1 + 2/e) \cdot \text{LP}^*$ . By (PD.1), every facility may appear in at most one primary demand's neighborhood, and the facilities open in Line 4–5 of Pseudocode 3 do not appear in any primary demand's neighborhood. Therefore, by linearity of expectation, the expected facility cost of Algorithm ECHS is

$$\text{Exp}[F_{\text{ECHS}}] = \sum_{\mu \in \mathbb{F}} f_\mu \bar{y}_\mu = \sum_{i \in \mathbb{F}} f_i \sum_{\mu \in i} \bar{y}_\mu = \sum_{i \in \mathbb{F}} f_i y_i^* = F^*,$$

where the third equality follows from (PS.3).

To bound the connection cost, we adapt an argument of Chudak and Shmoys [5]. Consider a demand  $\nu$ . This demand can either get connected directly to some facility in  $\overline{N}(\nu)$  or indirectly to its target facility  $\phi(\kappa) \in \overline{N}(\kappa)$ , where  $\kappa$  is the primary demand to which  $\nu$  is assigned.

We first estimate the expected cost  $d_{\phi(\kappa)\nu}$  of the indirect connection. Let  $\Lambda_\nu$  denote the event that none of the facilities in  $\overline{N}(\nu)$  is opened. Then

$$\text{Exp}[d_{\phi(\kappa)\nu} | \Lambda_\nu] = D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \nu), \quad (3)$$

where  $D(A, \sigma) \stackrel{\text{def}}{=} \sum_{\mu \in A} d_{\mu\sigma} \bar{y}_\mu / \sum_{\mu \in A} \bar{y}_\mu$ , for any set  $A$  of facilities and a demand  $\sigma$ . Note that  $C_\nu^{\text{avg}} = D(\overline{N}(\nu), \nu)$ , and that  $\Lambda_\nu$  implies that  $\overline{N}(\kappa) \setminus \overline{N}(\nu) \neq \emptyset$ , since  $\overline{N}(\kappa)$  contains exactly one open facility, namely  $\phi(\kappa)$ .

**Lemma 9.** *Let  $\nu$  be a demand assigned to a primary demand  $\kappa$ , and assume that  $\overline{N}(\kappa) \setminus \overline{N}(\nu) \neq \emptyset$ . Then  $\text{Exp}[d_{\phi(\kappa)\nu} | \Lambda_\nu] \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*$ .*

*Proof.* By (3), we need to show that  $D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \nu) \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*$ . There are two cases to consider.

Case 1: There exists some  $\mu' \in \overline{N}(\kappa) \cap \overline{N}(\nu)$  such that  $d_{\mu'\kappa} \leq C_\kappa^{\text{avg}}$ . In this case, for every  $\mu \in \overline{N}(\kappa) \setminus \overline{N}(\nu)$ , we have

$$d_{\mu\nu} \leq d_{\mu\kappa} + d_{\mu'\kappa} + d_{\mu'\nu} \leq \alpha_\kappa^* + C_\kappa^{\text{avg}} + \alpha_\nu^* \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*,$$

using the triangle inequality, complementary slackness, and (PD.3(b)). By summing over all  $\mu \in \overline{N}(\kappa) \setminus \overline{N}(\nu)$ , it follows that  $D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \nu) \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*$ .

Case 2: Every  $\mu' \in \overline{N}(\kappa) \cap \overline{N}(\nu)$  has  $d_{\mu'\kappa} > C_\kappa^{\text{avg}}$ . Since  $C_\kappa^{\text{avg}} = D(\overline{N}(\kappa), \kappa)$ , this implies that  $D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \kappa) \leq C_\kappa^{\text{avg}}$ . Therefore, choosing an arbitrary  $\mu' \in \overline{N}(\kappa) \cap \overline{N}(\nu)$ , we obtain

$$D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \nu) \leq D(\overline{N}(\kappa) \setminus \overline{N}(\nu), \kappa) + d_{\mu'\kappa} + d_{\mu'\nu} \leq C_\kappa^{\text{avg}} + \alpha_\kappa^* + \alpha_\nu^* \leq C_\nu^{\text{avg}} + 2\alpha_\nu^*,$$

where we again use the triangle inequality, complementary slackness, and (PD.3(b)).

Since the lemma holds in both cases, the proof is now complete.  $\square$

We now continue our estimation of the connection cost. Denote by  $C_\nu$  the random variable representing the connection cost for  $\nu$ . The next step of our analysis is to show that

$$\text{Exp}[C_\nu] \leq C_\nu^{\text{avg}} + \frac{2}{e}\alpha_\nu^*. \quad (4)$$

The argument is divided into three cases. The first, easy case is when  $\nu$  is a primary demand  $\kappa$ . According to the algorithm (see Pseudocode 3, Line 2), we have  $C_\kappa = d_{\mu\kappa}$  with probability  $\bar{y}_\mu$ , for  $\mu \in \bar{N}(\kappa)$ . Therefore  $\text{Exp}[C_\kappa] = C_\kappa^{\text{avg}}$ , so (4) holds.

Next, we consider a non-primary demand  $\nu$ . Let  $\kappa$  be the primary demand that  $\nu$  is assigned to. We first deal with the sub-case when  $\bar{N}(\kappa) \setminus \bar{N}(\nu) = \emptyset$ , which is the same as  $\bar{N}(\kappa) \subseteq \bar{N}(\nu)$ . Property (CO) implies that  $\bar{x}_{\mu\nu} = \bar{y}_\mu = \bar{x}_{\mu\kappa}$  for every  $\mu \in \bar{N}(\kappa)$ , so we have  $\sum_{\mu \in \bar{N}(\kappa)} \bar{x}_{\mu\nu} = \sum_{\mu \in \bar{N}(\kappa)} \bar{x}_{\mu\kappa} = 1$ , due to (PS.1). On the other hand, we have  $\sum_{\mu \in \bar{N}(\nu)} \bar{x}_{\mu\nu} = 1$ , and  $\bar{x}_{\mu\nu} > 0$  for all  $\mu \in \bar{N}(\nu)$ . Therefore  $\bar{N}(\kappa) = \bar{N}(\nu)$  and  $C_\nu$  has exactly the same distribution as  $C_\kappa$ . So this case reduces to the first case, namely we have  $\text{Exp}[C_\nu] = C_\nu^{\text{avg}}$ , and (4) holds.

The last, and only non-trivial case is when  $\bar{N}(\kappa) \setminus \bar{N}(\nu) \neq \emptyset$ . Let  $\bar{N}(\nu) = \{\mu_1, \dots, \mu_l\}$  and let  $d_s = d_{\mu_s\nu}$  and  $y_s = \bar{y}_{\mu_s}$  for  $s = 1, \dots, l$ . By reordering, we can assume that  $d_1 \leq d_2 \leq \dots \leq d_l$ .

**Lemma 10.** *Assume that  $\bar{N}(\kappa) \setminus \bar{N}(\nu) \neq \emptyset$ . Then, using the notation from the paragraph above, we have*

$$\text{Exp}[C_\nu] \leq \sum_{u=1}^l d_u y_u \prod_{s=1}^{u-1} (1 - y_s) + \text{Exp}[d_{\phi(\kappa)\nu} \mid \Lambda_\nu] \cdot \prod_{s=1}^l (1 - y_s).$$

The formal proof of Lemma 10 appears in [5]. The main idea of the proof is to consider a provably worse random process, which considers facilities in  $\bar{N}(\nu)$  by grouping them according to the neighborhood of primary demands. That is, two facilities  $\mu_1$  and  $\mu_2$  belong to the same group if both are in  $\bar{N}(\kappa)$  for some primary demand  $\kappa$ . The groups are then ordered by nondecreasing average distance to  $\nu$ , where the average is  $\sum_{\mu \in G} d_{\mu\nu} \bar{y}_\mu / \sum_{\mu \in G} \bar{y}_\mu$ . The process first consider each group in this order and select a group  $G$  independently with probability  $\sum_{\mu \in G} \bar{y}_\mu$ . Stop as soon as one group is selected, or we end up with none selected. If one group  $G$  is selected, we open a facility in  $G$  by choosing each  $\mu$  in  $G$  with probability  $\bar{y}_\mu / \sum_{\mu \in G} \bar{y}_\mu$ . We then connect demand  $\nu$  to this facility  $\mu$ . Otherwise no group is selected and we connect  $\nu$  to  $\phi(\kappa)$ .

Using Lemma 10, we estimate the (unconditional) expected connection cost of  $\nu$  as follows:

$$\begin{aligned} \text{Exp}[C_\nu] &\leq \sum_{u=1}^l d_u y_u \prod_{s=1}^{u-1} (1 - y_s) + \text{Exp}[d_{\phi(\kappa)\nu} \mid \Lambda_\nu] \cdot \prod_{s=1}^l (1 - y_s) \\ &\leq \left(1 - \prod_{s=1}^l (1 - y_s)\right) \cdot \sum_{u=1}^l d_u y_u + \text{Exp}[d_{\phi(\kappa)\nu} \mid \Lambda_\nu] \cdot \prod_{s=1}^l (1 - y_s) \end{aligned} \quad (5)$$

$$\leq \left(1 - \prod_{s=1}^l (1 - y_s)\right) \cdot C_\nu^{\text{avg}} + (C_\nu^{\text{avg}} + 2\alpha_\nu^*) \cdot \prod_{s=1}^l (1 - y_s) \quad (6)$$

$$\begin{aligned} &\leq \left(1 - \frac{1}{e}\right) C_\nu^{\text{avg}} + \frac{1}{e} (C_\nu^{\text{avg}} + 2\alpha_\nu^*) \\ &= C_\nu^{\text{avg}} + \frac{2}{e} \alpha_\nu^*, \end{aligned} \quad (7)$$

where inequality (5) is shown in the appendix (see also the proof in [5]), inequality (6) follows from Lemma 9, and inequality (7) follows from  $\prod_{s=1}^l (1 - y_s) \leq 1/e$ . This completes the proof of the bound (4) on the connection cost of  $\nu$ .

Summing over all demands of a client  $j$ , we bound the expected connection cost of client  $j$ :

$$\text{Exp}[C_j] = \sum_{\nu \in j} \text{Exp}[C_\nu] \leq \sum_{\nu \in j} (C_\nu^{\text{avg}} + \frac{2}{e} \alpha_\nu^*) = C_j^* + \frac{2}{e} r_j \alpha_j^*.$$



Finally, summing over all clients  $j$ , we obtain our bound on the expected connection cost,

$$\text{Exp}[C_{\text{ECHS}}] \leq C^* + \frac{2}{e} \text{LP}^*.$$

Therefore, we have established that our algorithm constructs a feasible integral solution with an overall expected cost

$$\text{Exp}[F_{\text{ECHS}} + C_{\text{ECHS}}] \leq F^* + C^* + \frac{2}{e} \cdot \text{LP}^* = (1 + 2/e) \cdot \text{LP}^* \leq (1 + 2/e) \cdot \text{OPT}.$$

Summarizing, we obtain the main result of this section.

**Theorem 11.** *Algorithm ECHS is a  $(1 + 2/e)$ -approximation algorithm for FTFP.*

## 7 Algorithm EBGs with Ratio 1.575

In this section we give our main result, a 1.575-approximation algorithm for FTFP, where 1.575 is the value of  $\min_{\gamma \geq 1} \max\{\gamma, 1 + 2/e^\gamma, \frac{1/e + 1/e^\gamma}{1 - 1/\gamma}\}$ , rounded to three decimal digits. This matches the ratio of the best known LP-rounding algorithm for UFL by Byrka *et al.* [3].

Recall that in Section 6 we showed how to compute an integral solution with facility cost bounded by  $F^*$  and connection cost bounded by  $C^* + 2/e \cdot \text{LP}^*$ . A natural idea is to balance these two costs, by reducing the connection cost at the expense of a slight increase in the facility cost. One way to do this would be to increase the probability of opening facilities, from  $\bar{y}_\mu$  (used in Algorithm ECHS) to, say,  $\gamma \bar{y}_\mu$ , for some  $\gamma > 1$ . This increases the expected facility cost by a factor of  $\gamma$  but, as it turns out, it also reduces the probability that an indirect connection occurs for a non-primary demand to  $1/e^\gamma$  (from the previous value  $1/e$ ). Moreover, for each primary demand  $\kappa$ , the new algorithm will select a facility from the nearest facilities  $\mu$  in  $\bar{N}(\kappa)$  such that the connection values  $\bar{x}_{\mu\nu}$  sum up to  $1/\gamma$ , instead of 1. It is easily seen that this will improve the estimate on connection cost for primary demands. These two changes, along with a more refined analysis, are the essence of the approach in [3], expressed in our terminology.

Our approach can be thought of as a combination of the above ideas with the techniques of demand reduction and adaptive partitioning that we introduced earlier. However, our adaptive partitioning technique needs to be carefully modified, because now we will be using a more intricate neighborhood structure, with the neighborhood of each demand divided into two disjoint parts, and with some restrictions on how parts from different demands can overlap.

We begin by describing properties that our partitioned fractional solution  $(\bar{x}, \bar{y})$  needs to satisfy. As mentioned earlier, the neighborhood  $\bar{N}(\nu)$  of each demand  $\nu$  will be divided into two disjoint parts. The first part, called the *close neighborhood* and denoted  $\bar{N}_{\text{cls}}(\nu)$ , contains the facilities in  $\bar{N}(\nu)$  nearest to  $\nu$  with the total connection value equal  $1/\gamma$ , that is  $\sum_{\mu \in \bar{N}_{\text{cls}}(\nu)} \bar{x}_{\mu\nu} = 1/\gamma$ . The second part, called the *far neighborhood* and denoted  $\bar{N}_{\text{far}}(\nu)$ , contains the remaining facilities in  $\bar{N}(\nu)$  (so  $\sum_{\mu \in \bar{N}_{\text{far}}(\nu)} \bar{x}_{\mu\nu} = 1 - 1/\gamma$ ). We restate these definitions formally below in Property (NB). Recall that for any set  $A$  of facilities and a demand  $\nu$ , by  $D(A, \nu)$  we denote the average distance between  $\nu$  and the facilities in  $A$ , that is  $D(A, \nu) = \sum_{\mu \in A} d_{\mu\nu} \bar{y}_\mu / \sum_{\mu \in A} \bar{y}_\mu$ . We will use notations  $C_{\text{cls}}^{\text{avg}}(\nu) = D(\bar{N}_{\text{cls}}(\nu), \nu)$  and  $C_{\text{far}}^{\text{avg}}(\nu) = D(\bar{N}_{\text{far}}(\nu), \nu)$  for the average distances from  $\nu$  to its close and far neighborhoods, respectively. By the definition of these sets and the completeness property (CO), these distances can be expressed as

$$C_{\text{cls}}^{\text{avg}}(\nu) = \gamma \sum_{\mu \in \bar{N}_{\text{cls}}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu} \quad \text{and} \quad C_{\text{far}}^{\text{avg}}(\nu) = \frac{\gamma}{\gamma - 1} \sum_{\mu \in \bar{N}_{\text{far}}(\nu)} d_{\mu\nu} \bar{x}_{\mu\nu}.$$

We will also use notation  $C_{\text{cls}}^{\max}(\nu) = \max_{\mu \in \overline{N}_{\text{cls}}(\nu)} d_{\mu\nu}$  for the maximum distance from  $\nu$  to its close neighborhood.

Our partitioned solution  $(\bar{x}, \bar{y})$  must satisfy the same partitioning and completeness properties as before, namely properties (PS) and (CO) in Section 4. In addition, it must satisfy a new neighborhood property (NB) and modified properties (PD') and (SI'), listed below.

(NB) *Neighborhoods.* For each demand  $\nu \in \overline{\mathbb{C}}$ , its neighborhood is divided into *close* and *far* neighborhood, that is  $\overline{N}(\nu) = \overline{N}_{\text{cls}}(\nu) \cup \overline{N}_{\text{far}}(\nu)$ , where

- $\overline{N}_{\text{cls}}(\nu) \cap \overline{N}_{\text{far}}(\nu) = \emptyset$ ,
- $\sum_{\mu \in \overline{N}_{\text{cls}}(\nu)} \bar{x}_{\mu\nu} = 1/\gamma$ , and
- if  $\mu \in \overline{N}_{\text{cls}}(\nu)$  and  $\mu' \in \overline{N}_{\text{far}}(\nu)$  then  $d_{\mu\nu} \leq d_{\mu'\nu}$ .

Note that the first two conditions, together with (PS.1), implies that  $\sum_{\mu \in \overline{N}_{\text{far}}(\nu)} \bar{x}_{\mu\nu} = 1 - 1/\gamma$ .

(PD') *Primary demands.* Primary demands satisfy the following conditions:

1. For any two different primary demands  $\kappa, \kappa' \in P$  we have  $\overline{N}_{\text{cls}}(\kappa) \cap \overline{N}_{\text{cls}}(\kappa') = \emptyset$ .
2. For each site  $i \in \mathbb{F}$ ,  $\sum_{\kappa \in P} \sum_{\mu \in i \cap \overline{N}_{\text{cls}}(\kappa)} \bar{x}_{\mu\kappa} \leq y_i^*$ . In the summation, as before, we overload notation  $i$  to stand for the set of facilities created on site  $i$ .
3. Each demand  $\nu \in \overline{\mathbb{C}}$  is assigned to one primary demand  $\kappa \in P$  such that
  - (a)  $\overline{N}_{\text{cls}}(\nu) \cap \overline{N}_{\text{cls}}(\kappa) \neq \emptyset$ , and
  - (b)  $C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\max}(\nu) \geq C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\max}(\kappa)$ .

(SI') *Siblings.* For any pair  $\nu, \nu' \in \overline{\mathbb{C}}$  of different siblings we have

1.  $\overline{N}(\nu) \cap \overline{N}(\nu') = \emptyset$ .
2. If  $\nu$  is assigned to a primary demand  $\kappa$  then  $\overline{N}(\nu') \cap \overline{N}_{\text{cls}}(\kappa) = \emptyset$ . In particular, by Property (PD'.3(a)), this implies that different sibling demands are assigned to different primary demands, since  $\overline{N}_{\text{cls}}(\nu')$  is a subset of  $\overline{N}(\nu')$ .

**Modified adaptive partitioning.** To obtain a fractional solution with the above properties, we employ a modified adaptive partitioning algorithm. As in Section 4, we have two phases. In Phase 1 we split clients into demands and create facilities on sites, while in Phase 2 we augment each demand's connection values  $\bar{x}_{\mu\nu}$  by setting  $\bar{x}_{\mu\nu}$  from 0 to  $\bar{y}_\mu$ , so that the total value is 1. As the partitioning algorithm proceeds, for any demand  $\nu$ ,  $\overline{N}(\nu)$  denotes the set of facilities with  $\bar{x}_{\mu\nu} > 0$ , hence  $\overline{N}(\nu)$  are dynamic sets which get fixed once the partitioning algorithm concludes both Phase 1 and Phase 2. On the other hand,  $\overline{N}_{\text{cls}}(\nu)$  and  $\overline{N}_{\text{far}}(\nu)$  refer to the close and far neighborhoods once  $\overline{N}(\nu)$  is fixed.

Similar to the algorithm in Section 4, Phase 1 runs in iterations. Fix some iteration and consider any client  $j$ . As before,  $\tilde{N}(j)$  is the neighborhood of  $j$  with respect to the yet unpartitioned solution, namely the set of facilities  $\mu$  such that  $\tilde{x}_{\mu j} > 0$ . Order the facilities in this set as  $\tilde{N}(j) = \{\mu_1, \dots, \mu_q\}$  with non-decreasing distance from  $j$ , that is  $d_{\mu_1 j} \leq d_{\mu_2 j} \leq \dots \leq d_{\mu_q j}$ , where  $q = |\tilde{N}(j)|$ . Without loss of generality, there is an index  $l$  for which  $\sum_{s=1}^l \tilde{x}_{\mu_s j} = 1/\gamma$ , since we can always split one

facility to achieve this. Then we define  $\tilde{N}_{\text{cls}}(j) = \{\mu_1, \dots, \mu_l\}$ . In case of ties, we use a tie-breaking rule that is explained in the proof of Lemma 12 (the only place where the rule is needed). We also use notation

$$\text{tcc}_{\text{cls}}(j) = D(\tilde{N}_{\text{cls}}(j), j) = \gamma \sum_{\mu \in \tilde{N}_{\text{cls}}(j)} d_{\mu j} \tilde{x}_{\mu j} \quad \text{and} \quad \text{dmax}_{\text{cls}}(j) = \max_{\mu \in \tilde{N}_{\text{cls}}(j)} d_{\mu j}.$$

When the iteration starts, we first find a not-yet-exhausted client  $p$  that minimizes the value of  $\text{tcc}_{\text{cls}}(p) + \text{dmax}_{\text{cls}}(p)$  and create a new demand  $\nu$  for  $p$ . Now we have two cases:

Case 1:  $\tilde{N}_{\text{cls}}(p) \cap \overline{N}(\kappa) \neq \emptyset$  for some existing primary demand  $\kappa \in P$ . In this case we assign  $\nu$  to  $\kappa$ . As before, if there are multiple such  $\kappa$ , we pick any of them. We also fix  $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$  and  $\tilde{x}_{\mu p} \leftarrow 0$  for each  $\mu \in \tilde{N}(p) \cap \overline{N}(\kappa)$ . Note that, as before, although we check for overlap between  $\tilde{N}_{\text{cls}}(p)$  and  $\overline{N}(\kappa)$ , the facilities we actually move into  $\overline{N}(\nu)$  include all facilities in the intersection of  $\tilde{N}(p)$ , a bigger set, with  $\overline{N}(\kappa)$ .

At this time, the total connection value  $\bar{x}_{\mu\nu}$  between  $\nu$  and  $\mu \in \overline{N}(\nu)$  is at most  $1/\gamma$  since the sum of  $\bar{y}_\mu$  for all  $\mu \in \overline{N}(\kappa)$  is  $1/\gamma$  and we have  $\overline{N}(\nu) \subseteq \overline{N}(\kappa)$  at this point. Later in Phase 2 we will add additional facilities from  $\tilde{N}(p)$  to  $\overline{N}(\nu)$  to make  $\nu$ 's total connection value equal to 1.

Case 2:  $\tilde{N}_{\text{cls}}(p) \cap \overline{N}(\kappa) = \emptyset$  for all existing primary demands  $\kappa \in P$ . In this case we make  $\nu$  a primary demand (that is, add it to  $P$ ) and assign it to itself. We then move the facilities from  $\tilde{N}_{\text{cls}}(p)$  to  $\overline{N}(\nu)$ , that is for  $\mu \in \tilde{N}_{\text{cls}}(p)$  we set  $\bar{x}_{\mu\nu} \leftarrow \tilde{x}_{\mu p}$  and  $\tilde{x}_{\mu p} \leftarrow 0$ .

It is easy to see that the total connection value of  $\nu$  to  $\overline{N}(\nu)$  is now exactly  $1/\gamma$ . Moreover, facilities remaining in  $\tilde{N}(p)$  are all farther away from  $\nu$  than those in  $\overline{N}(\nu)$ . As we add only facilities from  $\tilde{N}(p)$  to  $\overline{N}(\nu)$  in Phase 2, the final  $\overline{N}_{\text{cls}}(\nu)$  contains the same set of facilities as  $\overline{N}(\nu)$  at this point, taking into account of possible facility split.

Once all clients are exhausted, that is, each client  $j$  has  $r_j$  demands created, Phase 1 concludes. We then run Phase 2, the augmenting phase, following the same steps as in Section 4. For each client  $j$  and each demand  $\nu \in j$  with total connection value  $\bar{x}_{\mu\nu}$  for  $\mu \in \overline{N}(\nu)$  less than 1, we use our AUGMENTTOUNIT() procedure to add additional facilities (possibly split, if necessary) from  $\tilde{N}(j)$  to  $\overline{N}(\nu)$  to make the total connection value of  $\overline{N}(\nu)$  equal 1.

This completes the description of the partitioning algorithm. Summarizing, for each client  $j \in \mathbb{C}$  we have created  $r_j$  demands on the same point as  $j$  and we created a number of facilities at each site  $i \in \mathbb{F}$ . Thus computed sets of demands and facilities are denoted  $\overline{\mathbb{C}}$  and  $\overline{\mathbb{F}}$ , respectively. For each facility  $\mu \in i$  we defined its fractional opening value  $\bar{y}_\mu$ ,  $0 \leq \bar{y}_\mu \leq 1$ , and for each demand  $\nu \in j$  we defined its fractional connection value  $\bar{x}_{\mu\nu} \in \{0, \bar{y}_\mu\}$ . The connections with  $\bar{x}_{\mu\nu} > 0$  define the neighborhood  $\overline{N}(\nu)$ . The subset of  $\overline{N}(\nu)$  that are closest to  $\nu$  with total connection value  $1/\gamma$  form  $\overline{N}_{\text{cls}}(\nu)$  and the rest form  $\overline{N}_{\text{far}}(\nu)$ . We might split one facility to make the nearest facilities having total connection value exactly  $1/\gamma$ . It remains to show that this partitioning satisfies all the desired properties.

**Correctness of partitioning.** We now argue that our partitioned fractional solution  $(\bar{x}, \bar{y})$  satisfies all the stated properties. Properties (PS), (CO) and (NB) are directly enforced by the algorithm.

(PD'.1) holds because for each primary demand  $\kappa \in p$ ,  $\bar{N}_{\text{cls}}(\kappa)$  is the same set as  $\tilde{N}_{\text{cls}}(p)$  at the time when  $\kappa$  was created, and  $\tilde{N}_{\text{cls}}(p)$  is removed from  $\tilde{N}(p)$  right after this step. Further, the partitioning algorithm makes  $\kappa$  a primary demand only if  $\tilde{N}_{\text{cls}}(p)$  is disjoint from the set  $\bar{N}(\kappa')$  of all existing primary demands  $\kappa'$  at that iteration, but these neighborhoods are the same as the final close neighborhoods  $\bar{N}_{\text{cls}}(\kappa')$ .

The justification of (PD'.2) is similar to that for (PD.2) from Section 4. All close neighborhoods of primary demands are disjoint, due to (PD'.1), so each facility  $\mu \in i$  can appear in at most one  $\bar{N}_{\text{cls}}(\kappa)$ , for some  $\kappa \in P$ . Condition (CO) implies that  $\bar{y}_\mu = \bar{x}_{\mu\kappa}$  as we are only interested in  $\bar{x}_{\mu\kappa} > 0$ . As a result, the summation on the left-hand side is not larger than  $\sum_{\mu \in i} \bar{y}_\mu = y_i^*$ .

Regarding (PD'.3(a)), at first glance this property seems to follow directly from the algorithm, as we only assign a demand  $\nu$  to a primary demand  $\kappa$  when  $\bar{N}(\nu)$  at that iteration overlaps with  $\bar{N}(\kappa)$  (which is equal to the final value of  $\bar{N}_{\text{cls}}(\kappa)$ ). However, it is a little more subtle, as the final  $\bar{N}_{\text{cls}}(\nu)$  may contain facilities added to  $\bar{N}(\nu)$  in Phase 2. Those facilities may turn out to be closer to  $\nu$  than some facilities in  $\bar{N}(\kappa) \cap \tilde{N}(j)$  (not  $\tilde{N}_{\text{cls}}(j)$ ) that we added to  $\bar{N}(\nu)$  in Phase 1. If the final  $\bar{N}_{\text{cls}}(\nu)$  consists only facilities added in Phase 2, we no longer have the desired overlap of  $\bar{N}_{\text{cls}}(\kappa)$  and  $\bar{N}_{\text{cls}}(\nu)$ . Luckily this bad scenario never occurs. We postpone the proof of this property to Lemma 12. The proof of (PD'.3(b)) is similar to that of Lemma 5, and we defer it to Lemma 13.

(SI'.1) follows directly from the algorithm because for each demand  $\nu \in j$ , all facilities added to  $\bar{N}(\nu)$  are immediately removed from  $\tilde{N}(j)$ .

The proof of (SI'.2) is similar to that of Lemma 3. If  $\kappa = \nu$  then (SI'.2) follows from (SI'.1), so we can assume that  $\kappa \neq \nu$ . Suppose that  $\nu' \in j$  is assigned to  $\kappa' \in P$  and consider the situation after Phase 1. By the way we reassign facilities in Case 1, at this time we have  $\bar{N}(\nu) \subseteq \bar{N}(\kappa) = \bar{N}_{\text{cls}}(\kappa)$  and  $\bar{N}(\nu') \subseteq \bar{N}(\kappa') = \bar{N}_{\text{cls}}(\kappa')$ , so  $\bar{N}(\nu') \cap \bar{N}_{\text{cls}}(\kappa) = \emptyset$ , by (PD'.1). Moreover, we have  $\tilde{N}(j) \cap \bar{N}_{\text{cls}}(\kappa) = \emptyset$  after this iteration, because any facility that were also in  $\bar{N}_{\text{cls}}(\kappa)$  were removed from  $\tilde{N}(j)$  when  $\nu$  was created. In Phase 2, augmentation does not change  $\bar{N}_{\text{cls}}(\kappa)$  and all facilities added to  $\bar{N}(\nu')$  are from the set  $\tilde{N}(j)$  at the end of Phase 1, which is a subset of the set  $\tilde{N}(j)$  after this iteration, since  $\tilde{N}(j)$  is shrinking. So the condition (SI'.2) will remain true.

**Lemma 12.** *Property (PD'.3(a)) holds.*

*Proof.* Let  $j$  be the client for which  $\nu \in j$ . We consider an iteration when we create  $\nu$  from  $j$  and assign it to  $\kappa$ , and throughout the proof, notation  $\tilde{N}_{\text{cls}}(j)$  and  $\tilde{N}(j)$  will refer to the value of the sets at this particular time. At this time,  $\bar{N}(\nu)$  is initialized to  $\tilde{N}(j) \cap \bar{N}(\kappa)$ . Recall that  $\bar{N}(\kappa)$  is now equal to the final  $\bar{N}_{\text{cls}}(\kappa)$  (taking into account facility split). We would like to show that the set  $\tilde{N}_{\text{cls}}(j) \cap \bar{N}_{\text{cls}}(\kappa)$  (which is not empty) will be included in  $\bar{N}_{\text{cls}}(\nu)$  at the end. Technically speaking, this will not be true due to facility splitting, so we need to rephrase this claim and the proof in terms of the set of facilities obtained after the algorithm completes.

We define the sets  $A$ ,  $B$ ,  $E^-$  and  $E^+$  as the subsets of  $\bar{\mathbb{F}}$  (the final set of facilities) that were obtained from splitting facilities in the sets  $\tilde{N}(j)$ ,  $\tilde{N}_{\text{cls}}(j) \cap \bar{N}_{\text{cls}}(\kappa)$ ,  $\tilde{N}_{\text{cls}}(j) - \bar{N}_{\text{cls}}(\kappa)$  and  $\tilde{N}(j) - \tilde{N}_{\text{cls}}(j)$ , respectively. (See Figure 1.) We claim that at the end  $B \subseteq \bar{N}_{\text{cls}}(\nu)$ , with the caveat that the ties in the definition of  $\bar{N}_{\text{cls}}(\nu)$  are broken in favor of the facilities in  $B$ . (This is the tie-breaking rule that we mentioned in the definition of  $\bar{N}_{\text{cls}}(\nu)$ .) This will be sufficient to prove the lemma because  $B \neq \emptyset$ , by the algorithm.

We now prove this claim. Note first that  $B \subseteq \bar{N}(\nu) \subseteq A$ , because we never remove facilities from  $\bar{N}(\nu)$  and we only add facilities from  $\tilde{N}(j)$ . Also,  $B \cup E^-$  represents the facilities obtained from  $\tilde{N}_{\text{cls}}(j)$ , so  $\sum_{\mu \in B \cup E^-} \bar{y}_\mu = 1/\gamma$ . This and  $B \subseteq \bar{N}(\nu)$  implies that the total connection value

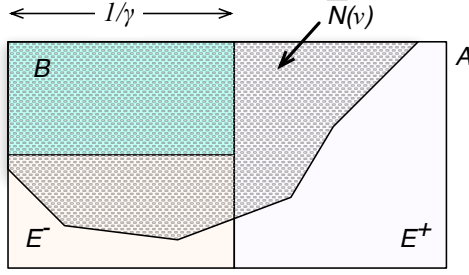


Figure 1: Illustration of the sets  $\bar{N}(\nu)$ ,  $A$ ,  $B$ ,  $E^-$  and  $E^+$  in the proof of Lemma 12. We have  $A = \tilde{N}(j)$ ,  $B = \tilde{N}_{\text{cls}}(j) \cap \bar{N}_{\text{cls}}(\kappa)$ ,  $E^- = \tilde{N}_{\text{cls}}(j) - \bar{N}_{\text{cls}}(\kappa)$ ,  $E^+ = \tilde{N}(j) - \tilde{N}_{\text{cls}}(j)$ .

of  $B \cup (\bar{N}(\nu) \cap E^-)$  to  $\nu$  is at most  $1/\gamma$ . But all facilities in  $B \cup (\bar{N}(\nu) \cap E^-)$  are closer to  $\nu$  than those in  $E^+ \cap \bar{N}(\nu)$ . It follows that  $B \subseteq \bar{N}_{\text{cls}}(\nu)$ , completing the proof.  $\square$

**Lemma 13.** *Property (PD'.3(b)) holds.*

*Proof.* This proof is similar to that for Lemma 5. For a client  $j$  and demand  $\eta$ , we will write  $\text{tcc}_{\text{cls}}^\eta(j)$  and  $\text{dmax}_{\text{cls}}^\eta(j)$  to denote the values of  $\text{tcc}_{\text{cls}}(j)$  and  $\text{dmax}_{\text{cls}}(j)$  at the time when  $\eta$  was created. (Here  $\eta$  may or may not be a demand of client  $j$ ).

Suppose  $\nu \in j$  is assigned to a primary demand  $\kappa \in p$ . By the way primary demands are constructed in the partitioning algorithm,  $\tilde{N}_{\text{cls}}(p)$  becomes  $\bar{N}_{\text{cls}}(\kappa)$ , so we have  $C_{\text{cls}}^{\text{avg}}(\kappa) = \text{tcc}_{\text{cls}}^\kappa(p)$  and  $C_{\text{cls}}^{\text{max}}(\kappa) = \text{tcc}_{\text{cls}}^\kappa(p)$ . Further, since we choose  $p$  to minimize  $\text{tcc}_{\text{cls}}(p) + \text{dmax}_{\text{cls}}(p)$ , we thus have that  $\text{tcc}_{\text{cls}}^\kappa(p) + \text{dmax}_{\text{cls}}^\kappa(p) \leq \text{tcc}_{\text{cls}}^\kappa(j) + \text{dmax}_{\text{cls}}^\kappa(j)$  at the time when demand  $\kappa$  was created.

Using an argument analogous to that in the proof of Lemma 4, our modified partitioning algorithm guarantees that  $\text{tcc}_{\text{cls}}^\kappa(j) \leq \text{tcc}_{\text{cls}}^\nu(j)$  and  $\text{dmax}_{\text{cls}}^\kappa(j) \leq \text{dmax}_{\text{cls}}^\nu(j)$  since  $\nu$  was created later. Therefore, we have

$$\begin{aligned} C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa) &= \text{tcc}_{\text{cls}}^\kappa(p) + \text{dmax}_{\text{cls}}^\kappa(p) \\ &\leq \text{tcc}_{\text{cls}}^\kappa(j) + \text{dmax}_{\text{cls}}^\kappa(j) \leq \text{tcc}_{\text{cls}}^\nu(j) + \text{dmax}_{\text{cls}}^\nu(j) = C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu), \end{aligned}$$

completing the proof.  $\square$

Now we have completed the proof that the computed partitioning satisfies all the required properties.

**Algorithm EBGs.** The complete algorithm starts with solving the linear program and computing the partitioning described earlier in this section. Given the partitioned fractional solution  $(\bar{x}, \bar{y})$  with the desired properties, we start the process of opening facilities and making connections to obtain an integral solution. To this end, for each primary demand  $\kappa \in P$ , we open exactly one facility  $\phi(\kappa)$  in  $\bar{N}_{\text{cls}}(\kappa)$ , where each  $\mu \in \bar{N}_{\text{cls}}(\kappa)$  is chosen as  $\phi(\kappa)$  with probability  $\gamma \bar{y}_\mu$ . For all facilities  $\mu \in \bar{\mathbb{F}} - \bigcup_{\kappa \in P} \bar{N}_{\text{cls}}(\kappa)$ , we open them independently, each with probability  $\gamma \bar{y}_\mu$ . Due to our choice of  $\gamma$  being  $1 < \gamma < 2$ , we have  $1 - 1/\gamma \leq 1/\gamma$  so property (CO) actually implies all  $\bar{y}_\mu \leq 1/\gamma$  since there exists some  $\nu$  such that  $\bar{x}_{\mu\nu} = \bar{y}_\mu$  and any  $\mu \in \bar{N}_{\text{cls}}(\nu)$  satisfies  $\bar{x}_{\mu\nu} \leq 1/\gamma$ , and any  $\mu \in \bar{N}_{\text{far}}(\nu)$  satisfies  $\bar{x}_{\mu\nu} \leq 1 - 1/\gamma$ . Then all  $\gamma \bar{y}_\mu$  are valid probability values.

Next, we connect demands to facilities. Each primary demand  $\kappa \in P$  will connect to the only open facility  $\phi(\kappa)$  in  $\overline{N}_{\text{cls}}(\kappa)$ . For each non-primary demand  $\nu \in \overline{\mathcal{C}} - P$ , if there is an open facility in  $\overline{N}(\nu)$  then we connect  $\nu$  to the nearest such facility. Otherwise, we connect  $\nu$  to its *target facility*  $\phi(\kappa)$ , where  $\kappa$  is the primary demand that  $\nu$  is assigned to.

**Analysis.** By the algorithm, for each client  $j$ , all its  $r_j$  demands are connected to open facilities. If two different siblings  $\nu, \nu' \in j$  are assigned, respectively, to primary demands  $\kappa, \kappa'$  then, by Properties (SI'.1), (SI'.2), and (PD'.1) we have

$$(\overline{N}(\nu) \cup \overline{N}_{\text{cls}}(\kappa)) \cap (\overline{N}(\nu') \cup \overline{N}_{\text{cls}}(\kappa')) = \emptyset.$$

This condition guarantees that  $\nu$  and  $\nu'$  are assigned to different facilities, regardless whether they are connected to a neighbor facility or to its target facility. Therefore the computed solution is feasible.

We now estimate the cost of the solution computed by Algorithm EBGs. The lemma below bounds the expected facility cost.

**Lemma 14.** *The expected facility cost  $F_{\text{EBGS}}$  of Algorithm EBGs is at most  $\gamma F^*$ .*

*Proof.* By the algorithm, each facility  $\mu \in \overline{\mathbb{F}}$  is opened with probability  $\gamma \bar{y}_\mu$ , independently of whether it belongs to the close neighborhood of a primary demand or not. Therefore, by linearity of expectation, we have that the expected facility cost is

$$F_{\text{EBGS}} = \sum_{\mu \in \overline{\mathbb{F}}} f_\mu \gamma \bar{y}_\mu = \gamma \sum_{i \in \mathbb{F}} f_i \sum_{\mu \in i} \bar{y}_\mu \stackrel{*}{=} \gamma \sum_{i \in \mathbb{F}} f_i y_i^* = \gamma F^*,$$

where the (\*) equality follows from (PS.3).  $\square$

To bound the connection cost, we start with a lemma that estimates the cost of indirect connections, between a non-primary demand and its target facility. The proof of this lemma relies on Properties (PD'.3(a)) and (PD'.3(b)) of modified partitioning and follows the reasoning in the proof of a similar lemma in [3, 2]. For the sake of completeness, we include a proof in the appendix.

**Lemma 15.** *In Algorithm EBGs, let  $\gamma$  be a constant such that  $1 < \gamma < 2$ , for a non-primary demand  $\nu$ , in the event that no facility in  $\overline{N}(\nu)$  opens, the expected connection cost is no more than  $C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu)$ .*

We are now ready to bound the overall connection cost of Algorithm EBGs.

**Lemma 16.** *Assuming that  $1 < \gamma < 2$ , the expected connection cost  $C_{\text{EBGS}}$  of Algorithm EBGs is at most  $C^* \cdot \max\{\frac{1/e+1/e^\gamma}{1-1/\gamma}, 1 + \frac{2}{e^\gamma}\}$ .*

*Proof.* The argument is similar to that in [3]. Fix some non-primary demand  $\nu$ . Let  $\overline{N}_{\text{cls}}(\nu) = \{\mu_1, \dots, \mu_l\}$  and  $\overline{N}_{\text{far}}(\nu) = \{\mu_{l+1}, \dots, \mu_{l'}\}$ , where facilities are ordered by non-decreasing distance to  $\nu$ , and  $d_1, \dots, d_l, d_{l+1}, \dots, d_{l'}$  be the corresponding distance to  $\nu$ . Let  $\Lambda_\nu$  be the event that none of  $\overline{N}(\nu)$  opens in algorithm EBGs's integral solution and  $\kappa$  be the primary demand that  $\nu$  was assigned to, we then have  $\text{Exp}[d_{\phi(\kappa)\nu} \mid \Lambda_\nu] \leq C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu)$  by Lemma 15.

Let  $C_\nu$  be the random variable of the distance from  $\nu$  to the facility that it connects to in the integral solution. A similar argument to estimate the distance by a provably worse random process,

as in the analysis of algorithm ECHS, will then show that the expected distance  $\text{Exp}[C_\nu]$  is no more than the estimate that uses the nearest facility in  $\bar{N}_{\text{cls}}(\nu)$  if one opened, then use the nearest in  $\bar{N}_{\text{far}}(\nu)$  if there is one, and only use  $\phi(\kappa)$  if none in  $\bar{N}(\nu)$  opened. A detailed derivation follows.

$$\begin{aligned}
\text{Exp}[C_\nu] &\leq \left(1 - \prod_{s=1}^l (1 - \gamma \bar{y}_s)\right) \sum_{s=1}^l d_s \gamma \bar{y}_s \\
&\quad + \left(\prod_{s=1}^l (1 - \gamma \bar{y}_s) - \prod_{s=1}^{l'} (1 - \gamma \bar{y}_s)\right) \left(\sum_{s=l+1}^{l'} d_s \bar{y}_s / \sum_{s=l+1}^{l'} \bar{y}_s\right) \\
&\quad + \left(\prod_{s=1}^{l'} (1 - \gamma \bar{y}_s)\right) \text{Exp}[d_{\phi(\kappa)\nu} \mid \Lambda_\nu] \\
&\leq (1 - 1/e) C_{\text{cls}}^{\text{avg}}(\nu) + (1/e - 1/e^\gamma) C_{\text{far}}^{\text{avg}}(\nu) + 1/e^\gamma (C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu)) \\
&\leq C_{\text{cls}}^{\text{avg}}(\nu) (1 - \frac{1}{e}) + C_{\text{far}}^{\text{avg}}(\nu) (\frac{1}{e} - \frac{1}{e^\gamma}) + (C_{\text{cls}}^{\text{avg}}(\nu) + 2C_{\text{far}}^{\text{avg}}(\nu)) \frac{1}{e^\gamma} \\
&\leq C^{\text{avg}}(\nu) ((1 - \rho_\nu) (\frac{1/e + 1/e^\gamma}{1 - 1/\gamma}) + \rho_\nu (1 + \frac{2}{e^\gamma})) \\
&\leq C^{\text{avg}}(\nu) \cdot \max \left\{ \frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma} \right\},
\end{aligned}$$

where  $\rho_\nu = C_{\text{cls}}^{\text{avg}}(\nu)/C^{\text{avg}}(\nu)$ . It is easy to see that  $\rho_\nu$  is between 0 and 1. In the second inequality we used the fact that  $\sum_{s=1}^l \bar{y}_s = 1/\gamma$  and  $\sum_{s=1}^{l'} \bar{y}_s = 1$ . Since  $\sum_{\nu \in j} C^{\text{avg}}(\nu) = \sum_{\nu \in j} \sum_{\mu \in \mathbb{F}} d_{\mu\nu} \bar{x}_{\mu\nu} = \sum_{i \in \mathbb{F}} d_{ij} x_{ij}^* = C_j^*$ , summing over all clients  $j$  we have total connection cost bounded by  $C^* \max\{\frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma}\}$ .  $\square$

Recall that the expected facility cost is bounded by  $\gamma F^*$ , as argued earlier. Hence the total cost is bounded by  $\max\{\gamma, \frac{1/e + 1/e^\gamma}{1 - 1/\gamma}, 1 + \frac{2}{e^\gamma}\} \cdot \text{LP}^*$ . Picking  $\gamma = 1.575$  we obtain the desired ratio.

**Theorem 17.** *Algorithm EBGs is a 1.575-approximation algorithm for FTFP.*

## 8 Final Comments

In this paper we show a sequence of LP-rounding approximation algorithms for FTFP, with the best algorithm achieving ratio 1.575. As we mentioned earlier, we believe that our techniques of demand reduction and adaptive partitioning are very flexible and should be useful in extending other LP-rounding methods for UFL to obtain matching bounds for FTFP.

One of the main open problems in this area is whether FTFL can be approximated with the same ratio as UFL, and our work was partly motivated by this question. The techniques we introduced are not directly applicable to FTFL, mainly because our partitioning approach involves facility splitting that could result in several sibling demands being served by facilities on the same site. Nonetheless, we hope that further refinements of our construction might get around this issue and lead to new algorithms for FTFL with improved ratios.

## References

- [1] Vijay Arya, Naveen Garg, Rohit Khandekar, Adam Meyerson, Kamesh Munagala, and Vinayaka Pandit. Local search heuristics for k-median and facility location problems. *SIAM J. Comput.*, 33(3):544–562, 2004.
- [2] Jaroslaw Byrka and Karen Aardal. An optimal bifactor approximation algorithm for the metric uncapacitated facility location problem. *SIAM J. Comput.*, 39(6):2212–2231, 2010.
- [3] Jaroslaw Byrka, MohammadReza Ghodsi, and Aravind Srinivasan. LP-rounding algorithms for facility-location problems. *CoRR*, abs/1007.3611, 2010.
- [4] Jaroslaw Byrka, Aravind Srinivasan, and Chaitanya Swamy. Fault-tolerant facility location, a randomized dependent LP-rounding algorithm. In *Proceedings of the 14th Integer Programming and Combinatorial Optimization, IPCO '10*, pages 244–257, 2010.
- [5] Fabián Chudak and David Shmoys. Improved approximation algorithms for the uncapacitated facility location problem. *SIAM J. Comput.*, 33(1):1–25, 2004.
- [6] Sudipto Guha and Samir Khuller. Greedy strikes back: improved facility location algorithms. In *Proceedings of the 9th ACM-SIAM Symposium on Discrete Algorithms, SODA '98*, pages 649–657, 1998.
- [7] Sudipto Guha, Adam Meyerson, and Kamesh Munagala. Improved algorithms for fault tolerant facility location. In *Proceedings of the 12th ACM-SIAM Symposium on Discrete Algorithms, SODA '01*, pages 636–641, 2001.
- [8] Anupam Gupta. Lecture notes: CMU 15-854b, spring 2008, 2008.
- [9] Godfrey Harold Hardy, John Edensor Littlewood, and George Pólya. *Inequalities*. Cambridge University Press, 1988.
- [10] Kamal Jain, Mohammad Mahdian, Evangelos Markakis, Amin Saberi, and Vijay Vazirani. Greedy facility location algorithms analyzed using dual fitting with factor-revealing LP. *J. ACM*, 50(6):795–824, 2003.
- [11] Kamal Jain and Vijay Vazirani. Approximation algorithms for metric facility location and k-median problems using the primal-dual schema and lagrangian relaxation. *J. ACM*, 48(2):274–296, 2001.
- [12] Kamal Jain and Vijay V. Vazirani. An approximation algorithm for the fault tolerant metric facility location problem. *Algorithmica*, 38(3):433–439, 2003.
- [13] Shi Li. A 1.488 approximation algorithm for the uncapacitated facility location problem. In *Proceedings of the 38th International Conference on Automata, Languages and Programming, ICALP '11*, volume 6756, pages 77–88, 2011.
- [14] Kewen Liao and Hong Shen. Unconstrained and constrained fault-tolerant resource allocation. In *Proceedings of the 17th Annual International Conference on Computing and Combinatorics, COCOON'11*, pages 555–566, 2011.



- [15] Mohammad Mahdian, Yinyu Ye, and Jiawei Zhang. Approximation algorithms for metric facility location problems. *SIAM J. Comput.*, 36(2):411–432, 2006.
- [16] David Shmoys, Éva Tardos, and Karen Aardal. Approximation algorithms for facility location problems (extended abstract). In *Proceedings of the 29th Annual ACM Symposium on Theory of Computing, STOC '97*, pages 265–274, 1997.
- [17] Maxim Sviridenko. An improved approximation algorithm for the metric uncapacitated facility location problem. In *Proceedings of the 9th International IPCO Conference on Integer Programming and Combinatorial Optimization, IPCO '02*, pages 240–257, 2002.
- [18] Chaitanya Swamy and David Shmoys. Fault-tolerant facility location. *ACM Trans. Algorithms*, 4(4):1–27, 2008.
- [19] Jens Vygen. *Approximation Algorithms for Facility Location Problems*. Forschungsinst. für Diskrete Mathematik, 2005.
- [20] Shihong Xu and Hong Shen. The fault-tolerant facility allocation problem. In *Proceedings of the 20th International Symposium on Algorithms and Computation, ISAAC '09*, pages 689–698, 2009.
- [21] Li Yan and Marek Chrobak. Approximation algorithms for the fault-tolerant facility placement problem. *Inf. Process. Lett.*, 111(11):545–549, 2011.

## A Proof of Lemma 15

Lemma 15 provides a bound on the expected connection cost of a demand  $\nu$  when none of  $\bar{N}(\nu)$  opens by the algorithm. The formal proof is very similar to that in [2], and we sketch the main steps here for completeness.

Let  $K_1 = \bar{N}_{\text{cls}}(\kappa) \setminus \bar{N}(\nu)$ ,  $K_2 = \bar{N}_{\text{cls}}(\kappa) \cap \bar{N}_{\text{cls}}(\nu)$  and  $K_3 = \bar{N}_{\text{cls}}(\kappa) \cap \bar{N}_{\text{far}}(\nu)$ , then we shall show that  $D(K_1, \nu) \leq C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu)$ , which implies  $D(K_1, \nu) \leq C_{\text{cls}}^{\text{avg}}(\nu) + C_{\text{cls}}^{\text{max}}(\nu) + C_{\text{far}}^{\text{avg}}(\nu)$  due to (PD'.3(b)). Recall that  $D(A, \nu) \stackrel{\text{def}}{=} \sum_{\mu \in A} d_{\mu\nu} \bar{y}_\mu$  is the average distance of facilities in set  $A$  to demand  $\nu$ .

We can dismiss a few simple cases: If  $D(K_1, \kappa) \leq C_{\text{cls}}^{\text{avg}}(\kappa)$ , then we can use any facility  $\mu \in K_2$  to show that  $D(K_1, \nu) \leq D(K_1, \kappa) + d_{\mu\kappa} + d_{\mu\nu} \leq C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu)$  because  $d_{\mu\kappa} \leq C_{\text{cls}}^{\text{max}}(\kappa)$  and  $d_{\mu\nu} \leq C_{\text{cls}}^{\text{max}}(\nu) \leq C_{\text{far}}^{\text{avg}}(\nu)$ , and we are done. So we can assume  $D(K_1, \kappa) > C_{\text{cls}}^{\text{avg}}(\kappa)$ . Our second assumption is that every  $\mu \in K_2$  has  $d_{\mu\kappa} > C_{\text{cls}}^{\text{avg}}(\kappa)$ , for otherwise we can use  $C_{\text{cls}}^{\text{max}}(\kappa)$  to bound  $D(K_1, \kappa)$  and there exists a facility  $\mu \in K_2$  such that  $d_{\mu\kappa} \leq C_{\text{cls}}^{\text{avg}}(\kappa)$ . Since  $\mu \in K_2$ , we infer that  $d_{\mu\nu} \leq C_{\text{cls}}^{\text{max}}(\nu) \leq C_{\text{far}}^{\text{avg}}(\nu)$ , and we are done as well.

Given the two assumptions we can now assume that  $K_3$  is not empty and there exists some finite number  $\delta > 0$  such that  $D(K_3, \kappa) = C_{\text{cls}}^{\text{avg}}(\kappa) - \delta$ . This in turn implies  $D(K_3, \nu) \geq C_{\text{far}}^{\text{avg}}(\nu) + \delta$ , or we will be done. Because  $K_3$  is a subset of  $\bar{N}_{\text{far}}(\nu)$ , this gives us a bound on  $D(\bar{N}_{\text{far}}(\nu) \setminus K_3, \nu)$ , which in turn, is an upper bound on  $C_{\text{cls}}^{\text{max}}(\nu)$ . Let  $\hat{y} = \gamma \sum_{\mu \in K_3} \bar{y}_\mu$ , using the fact that  $C_{\text{far}}^{\text{avg}}(\nu) = D(K_3, \nu) \hat{y} / (\gamma - 1) + D(\bar{N}_{\text{far}}(\nu) \setminus K_3) (\gamma - 1 - \hat{y}) / (\gamma - 1)$ , we are able to show that

$$C_{\text{cls}}^{\text{max}}(\nu) \leq D(\bar{N}_{\text{far}}(\nu) \setminus K_3, \nu) \leq C_{\text{far}}^{\text{avg}}(\nu) - \hat{y}\delta / (\gamma - 1 - \hat{y}) \leq C_{\text{far}}^{\text{avg}}(\nu) - \hat{y}\delta / (1 - \hat{y}),$$

since  $\gamma \leq 2$  (we pick  $\gamma = 1.575$  in the end). On the other hand, since  $\bar{N}_{\text{cls}}(\kappa) = K_1 \cup K_2 \cup K_3$ , we have  $C_{\text{cls}}^{\text{avg}}(\kappa) = (1 - \hat{y}) D(K_1 \cup K_2, \kappa) + \hat{y} D(K_3, \kappa)$ , which implies that

$$D(K_1 \cup K_2, \kappa) = C_{\text{cls}}^{\text{avg}}(\kappa) + \hat{y}\delta / (1 - \hat{y}).$$

Now we are essentially done. If there exists some  $\mu \in K_2$  such that  $d_{\mu\kappa} \leq C_{\text{cls}}^{\text{avg}}(\kappa) + \hat{y}\delta / (1 - \hat{y})$ , then we have

$$\begin{aligned} D(K_1, \nu) &\leq D(K_1, \kappa) + d_{\mu\kappa} + d_{\mu\nu} \\ &\leq C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{avg}}(\kappa) + \hat{y}\delta / (1 - \hat{y}) + C_{\text{cls}}^{\text{max}}(\nu) \\ &\leq C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu) \end{aligned}$$

Otherwise we must have  $D(K_1, \kappa) \leq C_{\text{cls}}^{\text{avg}}(\kappa) + \hat{y}\delta / (1 - \hat{y})$ . It follows that

$$\begin{aligned} D(K_1, \nu) &\leq D(K_1, \kappa) + d_{\mu\kappa} + d_{\mu\nu} \\ &\leq C_{\text{cls}}^{\text{avg}}(\kappa) + \hat{y}\delta / (1 - \hat{y}) + C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{max}}(\nu) \\ &\leq C_{\text{cls}}^{\text{max}}(\kappa) + C_{\text{cls}}^{\text{avg}}(\kappa) + C_{\text{far}}^{\text{avg}}(\nu). \end{aligned}$$

This concludes the proof of Lemma 15, and we have the desired bound on the expected distance when the algorithm connects a demand  $\nu$  indirectly to its target facility, that is, the facility opened by the primary demand that  $\nu$  was assigned.

## B Proof of an Inequality for Analysis of Algorithm ECHS

In the  $1 + 2/e = 1.736$ -approximation in Section 6 we need to show the following inequality

$$\sum_{u=1}^l d_u y_u \prod_{s=1}^{u-1} (1 - y_s) \leq \left(1 - \prod_{s=1}^l (1 - y_s)\right) \cdot \sum_{u=1}^l d_u y_u \quad (8)$$

for  $d_1 \leq d_2 \leq \dots \leq d_l$  and  $\sum_{s=1}^l y_s = 1, y_s \geq 0$ .

In this section we give a new proof of this inequality, much simpler than the existing proof in [5], and also simpler than the argument by Sviridenko [17]. We derive this inequality from the following generalized version of the Chebyshev Sum Inequality:

$$\sum_i p_i \sum_j p_j a_j b_j \leq \sum_i p_i a_i \sum_j p_j b_j, \quad (9)$$

where each summation below runs from 1 to  $l$  and the sequences  $(a_i)$ ,  $(b_i)$  and  $(p_i)$  satisfy the following conditions:  $p_i \geq 0, a_i \geq 0, b_i \geq 0$  for all  $i$ ,  $a_1 \leq a_2 \leq \dots \leq a_l$ , and  $b_1 \geq b_2 \geq \dots \geq b_l$ .

Given inequality (9), we can obtain our inequality (8) by simple substitution

$$p_i \leftarrow y_i, a_i \leftarrow d_i, b_i \leftarrow \prod_{s=1}^{i-1} (1 - y_s).$$

For the sake of completeness, we include the proof of inequality (9), due to Hardy, Littlewood and Polya [9]. The idea is to evaluate the following sum:

$$\begin{aligned} S &= \sum_i p_i \sum_j p_j a_j b_j - \sum_i p_i a_i \sum_j p_j b_j \\ &= \sum_i \sum_j p_i p_j a_j b_j - \sum_i \sum_j p_i a_i p_j b_j \\ &= \sum_j \sum_i p_j p_i a_i b_i - \sum_j \sum_i p_j a_j p_i b_i \\ &= \frac{1}{2} \cdot \sum_i \sum_j (p_i p_j a_j b_j - p_i a_i p_j b_j + p_j p_i a_i b_i - p_j a_j p_i b_i) \\ &= \frac{1}{2} \cdot \sum_i \sum_j p_i p_j (a_i - a_j)(b_i - b_j) \leq 0. \end{aligned}$$

The last inequality holds because  $(a_i - a_j)(b_i - b_j) \leq 0$ , since the sequences  $(a_i)$  and  $(b_i)$  are ordered oppositely.