

# Efficient Algorithms for Geometric Partial Matching and Unbalanced Transportation problem

Pankaj K. Agarwal

Duke University, USA

[pankaj@cs.duke.edu](mailto:pankaj@cs.duke.edu)

Hsien-Chih Chang

Duke University, USA

[hsienchih.chang@duke.edu](mailto:hsienchih.chang@duke.edu)

Allen Xiao

Duke University, USA

[axiao@cs.duke.edu](mailto:axiao@cs.duke.edu)

---

## 1 Abstract

2 Let  $A$  and  $B$  be two point sets in the plane with uneven sizes  $r$  and  $n$  respectively (assuming  $r$  is at most  
3  $n$ ), and let  $k$  be a parameter. The geometric partial matching problem asks to find the minimum-cost  
4 size- $k$  matching between  $A$  and  $B$  under powers of  $L_p$  distances. Applying combinatorial algorithms  
5 for partial matching in general graphs to our setting naively requires quadratic time due to existence  
6 of many edges between point sets  $A$  and  $B$ . Most previous work for geometric matching has focused  
7 on the setting when  $k$ ,  $r$ , and  $n$  are equal. The best algorithm in this setting, due to Sharathkumar  
8 and Agarwal [STOC 2012], runs in time  $O(n \text{polylog } n \text{poly } \epsilon^{-1})$ , but is limited to matching objectives  
9 that are sum-of-distances.

10 We present the first set of efficient geometric algorithms which work for any powers of  $L_p$ -norm matching  
11 objective: An exact algorithm which runs in  $O((n + k^2) \text{polylog } n)$  time, and a  $(1 + \epsilon)$ -approximation  
12 which runs in  $O((n + k\sqrt{k}) \text{polylog } n \log \epsilon^{-1})$  time. Both algorithms are based on primal-dual flow  
13 augmentation scheme; the main improvements are obtained by using dynamic data structures to  
14 achieve efficient flow augmentations. Using similar techniques, we give an exact algorithm for the  
15 transportation problem in the plane, which runs in  $O(rn(r + \sqrt{n}) \text{polylog } n)$  time. This is the first  
16 sub-quadratic time exact algorithm when  $r = o(\sqrt{n})$ , which improves over the state-of-art quadratic  
17 time algorithm by Agarwal et al. [SOCG 2016].

2012 ACM Subject Classification Theory of computation → Design and analysis of algorithms

Keywords and phrases partial matchings, transportation, minimum-cost flow,

Lines 1017

VMS? geometric?

---

## 18 1 Introduction

### 19 «REWRITE AFTER TECHNICAL SECTIONS»

20 Consider the problem of finding a minimum-cost bichromatic matching between a set of red  
21 points  $A$  and a set of blue points  $B$  lying in the plane, where the cost of a matching edge  $(a, b)$  is  
22 the Euclidean distance  $\|a - b\|$ ; in other words, the minimum-cost bipartite matching problem  
23 on the Euclidean complete graph  $G = (A \cup B, A \times B)$ . Let  $r := |A|$  and  $n := |B|$ . Without loss of  
24 generality, assume that  $r \leq n$ . We consider the problem of *partial matching* (also called *imperfect*  
25 *matching*), where the task is to find a minimum-cost matching of size  $k \leq r$ . When  $k = r = n$ ,  
26 we say the matching instance is *balanced*. When  $k = r < n$  ( $A$  and  $B$  have different sizes, but  
27 the matching is maximal), we say the matching instance is *unbalanced*. We call the geometric



28 problem of finding a size  $k$  matching of point sets  $A$  and  $B$  the *geometric partial matching*  
 29 *problem*. **(talk about the near-linear time alg)**

### 30 1.1 Contributions

31 In this paper, we present two algorithms for geometric partial matching that are based on fitting  
 32 nearest-neighbor (NN) and geometric closest pair (BCP) oracles into primal-dual algorithms  
 33 for non-geometric bipartite matching and minimum-cost flow. This pattern is not new, see for  
 34 example (...)(**TODO cite**). Unlike these previous works, we focus on obtaining running time  
 35 dependencies on  $k$  or  $r$  instead of  $n$ , that is, faster for inputs with small  $r$  or  $k$ . We begin in  
 36 Section ?? by introducing notation for matching and minimum-cost flow.

37 First in Section 2, we show that the Hungarian algorithm [7] combined with a BCP oracle  
 38 solves geometric partial matching exactly in time  $O((n + k^2) \text{polylog } n)$ . Mainly, we show that we  
 39 can separate the  $O(n \text{polylog } n)$  preprocessing time for building the BCP data structure from the  
 40 augmenting paths' search time and update duals in a lazy fashion such that the number of dual  
 41 updates per augmenting path is  $O(k)$ .

42 ▶ **Theorem 1.1.** Let  $A$  and  $B$  be two point sets in the plane with  $|A| = r$  and  $|B| = n$  satisfying  
 43  $r \leq n$ , and let  $k$  be a parameter. A minimum-cost geometric partial matching of size  $k$  can be  
 44 computed between  $A$  and  $B$  in  $O((n + k^2) \text{polylog } n)$  time.

45 **(State the settings separately so no need to repeat in the theorem statement.)**

46 Next in Section 3, we apply a similar technique to the unit-capacity min-cost circulation  
 47 algorithm of Goldberg, Hed, Kaplan, and Tarjan [4]. The resulting algorithm finds a  $(1 + \varepsilon)$ -  
 48 approximation to the optimal geometric partial matching in  $O((n + k\sqrt{k}) \text{polylog } n \log(n/\varepsilon))$   
 49 time.

50 ▶ **Theorem 1.2.** Let  $A$  and  $B$  be two point sets in the plane with  $|A| = r$  and  $|B| = n$  satisfying  
 51  $r \leq n$ , and let  $k$  be a parameter. A  $(1 + \varepsilon)$  geometric partial matching of size  $k$  can be computed  
 52 between  $A$  and  $B$  in  $O((n + k\sqrt{k}) \text{polylog } n \log(1/\varepsilon))$  time.

53 Our third algorithm solves the transportation problem in the unbalanced setting. The trans-  
 54 portation problem is a weighted generalization of the matching problem. Each point of  $A$  is  
 55 weighted with an integer ~~or supply~~ and each point of  $B$  is weighted with integer ~~or demand~~ such that  
 56 the sum of supply and demand ~~are equal~~. The goal of the transportation problem is to find a  
 57 minimum-cost mapping of all supplies to demands, where the cost of moving a unit of supply at  
 58  $a \in A$  to satisfy a unit of demand at  $b \in B$  is  $\|a - b\|$ . For this, we use the strongly polynomial  
 59 uncapacitated min-cost flow algorithm by Orlin [8]. The result is an  $O(n^{3/2}r \text{polylog } n)$  time  
 60 algorithm for unbalanced transportation. This improves over the  $O(n^2 \text{polylog } n)$  time algorithm  
 61 of Agarwal et al. [1] when  $r = o(\sqrt{n})$ .

62 ▶ **Theorem 1.3.** Let  $A$  and  $B$  be two point sets in the plane with  $|A| = r$  and  $|B| = n$  satisfying  $r \leq n$ ,  
 63 with supplies and demands given by the function  $\lambda : (A \cup B) \rightarrow \mathbb{Z}$  such that  $\sum_{a \in A} \lambda(a) = \sum_{b \in B} \lambda(b)$ .  
 64 An optimal transportation map can be computed in  $O(rn(r/\sqrt{n} + \sqrt{n}) \text{polylog } n)$  time.

65 By nature of the BCP/NN oracles we use, these results generalize to when  $\|a - b\|$  is any  $L_p$   
 66 distance, and if we use  $p$ -th power costs  $c(a, b) = \|a - b\|^p$  for any  $1 \leq p < \infty$ .

## 67 2 Min-cost Partial Matching using Hungarian Algorithm

68 The Hungarian ~~A~~ is a primal-dual algorithm for min-cost bipartite matching in general  
 69 graphs and solves ~~partial~~ exactly if stopped after  $k$  iterations (see e.g. [?]). In this

*that can be adapted to terminate*

section, we prove Theorem 1.1 by implementing  $k$  iterations of the Hungarian Algorithm in  $O((n + k^2) \text{polylog } n)$  time for geometric partial matching.

## 2.1 Matching definitions

Let  $G$  be a bipartite graph between vertex sets  $V$  and  $W$  and edge set  $E$ , with costs  $c(v, w)$  for each edge  $(v, w) \in E$ . Let  $C := \max_{(v,w) \in E} c(v, w)$ . A **matching**  $M \subseteq E$  is a set of edges where no two edges share an endpoint. The **size** of a matching is the number of edges in the set, and the **cost** of a matching is the sum of costs of its edges. For a parameter  $k$ , the **minimum-cost partial matching problem (MPM)** asks to find a size- $k$  matching  $M^*$  of minimum cost. In geometric partial matching, we have  $V = A$ ,  $W = B$ ,  $E = A \times B$ , and  $c(a, b) = \|a - b\|_p^q$  for all  $(a, b) \in E$ .

The linear programming dual to the MPM linear program uses dual variables for each vertex, called **potentials**  $\pi : (V \cup W) \rightarrow \mathbb{R}$ . With potentials  $\pi$ , we define the **reduced cost**  $c_\pi(v, w) := c(v, w) - \pi(v) + \pi(w)$ . Potentials  $\pi$  are **feasible** if the reduced cost is nonnegative for all  $(v, w) \in E$ . We say  $(v, w) \in E$  is **admissible** under potentials  $\pi$  if  $c_\pi(v, w) = 0$ .

Consider a matching  $M$  of size less than  $r$ . An **alternating augmenting path** (or simply **augmenting path**)  $\Pi = (a_1, b_1, \dots, a_\ell, b_\ell)$  is an odd-length path with unmatched endpoints  $(a_1)$  and  $(b_\ell)$  and all other points matched. The edges of  $\Pi$  alternate between unmatched  $((a_i, b_i))$  and matched  $((b_i, a_{i+1}))$ . The symmetric difference  $M' \leftarrow M \oplus \Pi$  creates a new matching with size  $|M'| = |M| + 1$ . We say that  $M'$  is the result of **augmenting**  $M$  by  $\Pi$ .

**MOT**

## 2.2 The Hungarian Algorithm

The Hungarian Algorithm is initialized with  $M = \emptyset$  and  $\pi = 0$ , and maintains the following invariants: (i)  $\pi$  is feasible, (ii) all edges in  $M$  are admissible, (iii) unmatched vertices of  $A$  all have the same potential  $\alpha$  with  $\alpha \geq \pi(a)$  for any matched  $a \in A$ , and (iv) unmatched vertices of  $B$  all have the same potential  $\beta$  with  $\beta \leq \pi(b)$  for any matched  $b \in B$ . Ramshaw and Tarjan [?] show that these conditions are sufficient to **prove** that  $M$  is a minimum-cost size  $|M|$  matching. In each iteration, the Hungarian Algorithm augments  $M$  with an admissible augmenting path  $\Pi$ , discovered using a procedure called the **Hungarian search**. Once  $|M| = k$ ,  $M$  is an optimal size- $k$  matching.

The Hungarian search raises  $\pi$  on subsets of vertices to grow a set  $S$  of “vertices reachable from unmatched  $v \in V$  by alternating paths of admissible edges.” Once  $S$  contains an unmatched  $w \in W$ , an admissible augmenting path exists. Initially,  $S$  is the set of unmatched  $v \in V$ . Let the **frontier** of  $S$  be the edges of  $(V \cap S) \times (W \setminus S)$ . In each iteration, the Hungarian search first **relaxes** the minimum-reduced-cost edge  $(v, w)$  in the frontier, raising  $\pi(v) \leftarrow \pi(v) + c_\pi(v, w)$  for all  $v \in S$  to make  $(v, w)$  admissible, and adding  $w$  into  $S$ . It is easy to verify that this potential change preserves feasibility. As  $w \in W$  is added into  $S$ , we can store a backpointer to  $v$ , which can be used later to recover the admissible augmenting path through  $w$ . If  $w$  is matched, say to  $v'$ , then we also relax  $(v', w)$  by adding  $v'$  into  $S$  (no potential change needed, by invariant) with backpointer to  $w$ . If  $w$  is unmatched, the search finishes and we can recover an admissible augmenting path to  $w$  by following backpointers to an unmatched  $v \in V$ .

Each **alternating** augmenting path has length  $O(k)$ , as every other edge is a matching edge and  $|M| \leq k$ . Additionally, there are  $k$  augmentations throughout the Hungarian algorithm, so the total time spent on updating the matching (during augmentations) is  $O(k^2)$ . The remainder of this section describes an implementation of Hungarian search that **can run** in  $O(k \text{polylog } n)$  time after a one-time  $O(n \text{polylog } n)$  time preprocessing. With this, our implementation of the Hungarian Algorithm runs in  $O((n + k^2) \text{polylog } n)$  time.

## 114 2.3 Geometric implementation of Hungarian search

115 Observe that the Hungarian search makes  $O(k)$  relaxations, as each relaxation either leads to  
 116 an unmatched vertex (ending the search) or adds both vertices of a matching edge. We will  
 117 implement each relaxation step in  $O(\text{polylog } n)$  time, after preprocessing.

118 In general graphs, the expensive step is finding the minimum-reduced cost frontier edge — the  
 119 search must “look at every edge” e.g. by pushing them into a priority queue even if they are not  
 120 relaxed. Since the costs are geometric, we can replace this with a query to a dynamic **bichromatic**  
 121 **closest pair** (BCP) data structure. Given two point sets  $P$  and  $Q$  in the plane, the BCP is the pair  
 122 of points  $p \in P$  and  $q \in Q$  minimizing the additively weighted distance  $\|p - q\| - \omega(p) + \omega(q)$ , for  
 123 some real-valued vertex weights  $\omega(p)$ . Thus, the minimum-reduced cost frontier edge is precisely  
 124 the BCP of  $P = A \cap S$  and  $Q = B \setminus S$ , with  $\omega(p) = \pi(p)$ . **(Short history on BCP?)**

125 The state of the art dynamic BCP data structure from Kaplan, Mulzer, Roditty, Seiferth, and  
 126 Sharir [6] supports point insertions and deletions in  $O(\text{polylog } n)$  time, and answers queries in  
 127  $O(\log^2 n)$  time. During each relaxation, we perform at most one query and add a vertex to  $S$   
 128 incurring one BCP insertion or deletion. So ignoring the time to build  $P, Q$  at the beginning of the  
 129 search and the time needed for updating  $\pi$ , the running time for the search is  $O(k \text{ polylog } n)$ .

130 **Initial BCP sets by rewinding.** Recall that  $S$  is initialized to the set of unmatched vertices  
 131 of  $V = A$ , and therefore  $Q = B \setminus S$  has size  $n$ . We cannot afford to take  $O(n \text{ polylog } n)$  time to  
 132 construct the BCP sets at the beginning of every Hungarian search beyond the first. However, the  
 133 set of unmatched  $A$  vertices has changed by exactly one vertex since the last Hungarian search  
 134 — the augmentation newly matched one vertex  $a^* \in A$ . Thus, given the initial BCP sets  $P', Q'$   
 135 from the beginning of the last Hungarian search, we can construct the initial  $P$  and  $Q$  by taking  
 136  $Q \leftarrow Q' \cup \{a^*\}$  and  $P \leftarrow P' \setminus \{a^*\}$ . In other words, we only need a single BCP deletion ( $O(\text{polylog } n)$ )  
 137 time, given  $P'$  and  $Q'$ . **Take a list**  $P \leftarrow P' \setminus \{a^*\}$  **for the current**  $Q \leftarrow Q' \cup \{a^*\}$

138 To acquire  $P'$  and  $Q'$ , we keep a log of the points that are added to  $S$  over the course of the  
 139 Hungarian search, and **rewind** this log at the end of a Hungarian search by tracing over it in  
 140 reverse order. If a point added to  $S$  incurred a BCP insertion during the search, we now delete it  
 141 (resp. incurred deletion, insert it). The number of events in the log is  $O(k)$ , since the number of  
 142 relaxations per Hungarian search is  $O(k)$ . Thus, in  $O(k \text{ polylog } n)$  time, we can reconstruct  $P'$   
 143 and  $Q'$  for each Hungarian search beyond the first.

144 We refer to this as the **rewinding mechanism**. In summary, we can construct each initial BCP  
 145 set in  $O(k \text{ polylog } n)$  time, after  $O(n \text{ polylog } n)$  preprocessing time for constructing the very first  
 146 BCP sets.

147 **Potential updates by Vaidya's trick.** We modify a trick from Vaidya [10] for batching  
 148 potential updates. Potentials have a **stored value**, i.e. the currently recorded value of  $\pi(v)$ ,  
 149 and a **true value**, which may have changed from  $\pi(v)$ . The resulting algorithm queries the  
 150 minimum-reduced cost under the true values of  $\pi$  and updates the stored value rarely. **occasionally**

151 Throughout the entire Hungarian Algorithm, we maintain a scalar  $\delta \geq 0$  (initially 0) which  
 152 aggregates potential changes. Vertices  $a \in A$  that are added to  $S$  are inserted into BCP with weight  
 153  $\omega(a) \leftarrow \pi(a) - \delta$ , for whatever value  $\delta$  is at the time of insertion. Similarly, vertices  $b \in B$  that  
 154 are added to  $S$  have  $\omega(b) \leftarrow \pi(b) - \delta$  recorded ( $B \cap S$  points aren't added into a BCP set). When  
 155 the Hungarian search wants to raise the potentials of points in  $S$ ,  $\delta$  is increased by that amount  
 156 instead. Thus, true value for any potential of a point in  $S$  is  $\omega(p) + \delta$ . For points of  $(A \cup B) \setminus S$ , the  
 157 true potential is equal to the stored potential. Since all the points of  $S \cap A$  have weights uniformly  
 158 offsetted from their true potentials, the minimum edge returned by the BCP does not change.

+ as why?

159 Once a point is removed from  $S$  (i.e. by an augmentation or the rewinding mechanism), we  
 160 update its stored potential  $\pi(p) \leftarrow \omega(p) + \delta$ , again for the current value of  $\delta$ . Importantly,  $\delta$  is not  
 161 reset at the end of a Hungarian search, and persists through the entire algorithm. Thus, the initial  
 162 BCP sets constructed by the rewinding mechanism have true potentials accurately represented by  
 163  $\delta$  and  $\omega(p)$ .

164 We update  $\delta$  once per edge relaxations, so  $O(k)$  times per Hungarian search. There are  $O(k)$   
 165 stored values updated per Hungarian search, during the rewind. The time spent on potential  
 166 updates per Hungarian search is therefore  $O(k)$ . Combined with the previous discussion, we  
 167 can implement each Hungarian search in  $O(k \text{ polylog } n)$  time after a one-time  $O(n \text{ polylog } n)$   
 168 preprocessing.

In  
Summary

### 3 Approximating Min-Cost Partial Matching through Cost-Scaling

170 The goal of section is to prove Theorem 1.2; that is, to compute a size- $k$  geometric partial matching  
 171 between two point sets  $A$  and  $B$  in the plane, with cost at most  $(1 + \varepsilon)$  times the optimal matching,  
 172 in time  $O((n + k\sqrt{k}) \text{ polylog } n \log(1/\varepsilon))$ .

173 After introducing the necessary terminologies in Section 3.1, we reduce the partial matching  
 174 problem to computing an approximate minimum-cost flow on a unit-capacity reduction network  
 175 in Section 3.2. In Section 3.3 we outline the high-level overview of the cost-scaling algorithm.  
 176 We postpone the fast implementation using dynamic data structures to Section 4.

#### 3.1 Preliminaries on Network Flows

177 **Network.** Let  $G = (V, E)$  be a directed graph, augmented by edge costs  $c$  and capacities  $u$ , and  
 178 a supply-demand function  $\phi$  defined on the vertices. One can turn the graph  $G$  into a **network**  
 179  $N = (V, \vec{E})$ : For each directed edge  $(v, w)$  in  $E$ , insert two **arcs**  $v \rightarrow w$  and  $w \rightarrow v$  into the arc set  
 180  $\vec{E}$ ; the **forward arc**  $v \rightarrow w$  inherits the capacity and cost from the directed graph  $G$ , while the  
 181 **backward arc**  $w \rightarrow v$  satisfies  $u(w \rightarrow v) = 0$  and  $c(w \rightarrow v) = -c(v \rightarrow w)$ . This we ensure that the graph  
 182  $(V, \vec{E})$  is **symmetric** and the cost function  $c$  is **antisymmetric** on  $N$ . The positive values of  $\phi(v)$  are  
 183 referred to as **supply**, and the negative values of  $\phi(v)$  as **demand**. We assume that all capacities  
 184 are nonnegative, all supplies and demands are integers, and the sum of supplies and demands is  
 185 equal to zero. A **unit-capacity** network has all its edge capacities equal to 1. In this section we  
 186 assume all networks are of unit-capacity.

187 **Pseudoflows.** Given a network  $N := (V, \vec{E}, c, u, \phi)$ , a **pseudoflow** (or **flow** to be short)  $f : \vec{E} \rightarrow \mathbb{Z}$ <sup>1</sup>  
 188 on  $N$  is an antisymmetric function on the arcs of  $N$  satisfying  $f(v \rightarrow w) \leq u(v \rightarrow w)$  for every arc  
 189  $v \rightarrow w$ . We sometimes abuse the terminology by allowing pseudoflow to be defined on a directed  
 190 graph, in which case we are actually referring to the pseudoflow on the corresponding network by  
 191 extending the flow values antisymmetrically to the arcs. We say that  $f$  **saturates** an arc  $v \rightarrow w$   
 192 if  $f(v \rightarrow w) = u(v \rightarrow w)$ ; an arc  $v \rightarrow w$  is **residual** if  $f(v \rightarrow w) < u(v \rightarrow w)$ . The **support** of  $f$  in  $N$ ,  
 193 denoted as  $\text{supp}(f)$ , is the set of arcs with positive flows:

$$194 \text{supp}(f) := \{v \rightarrow w \in \vec{E} \mid f(v \rightarrow w) > 0\}.$$

188 <sup>1</sup> In general the pseudoflows are allowed to take real-values. Here under the unit-capacity assumption any optimal  
 189 flows are integer-valued. **«(cite integrality theorem?)»**

198 Given a pseudoflow  $f$ , we define the ***imbalance*** of a vertex (with respect to  $f$ ) to be

$$199 \quad \phi_f(v) := \phi(v) + \sum_{w \rightarrow v \in \vec{E}} f(w \rightarrow v) - \sum_{v \rightarrow w \in \vec{E}} f(v \rightarrow w).$$

200 We call positive imbalance ***excess*** and negative imbalance ***deficit***; and vertices with positive and  
201 negative imbalance ***excess vertices*** and ***deficit vertices***, respectively. A vertex is ***balanced*** if it  
202 has zero imbalance. If all vertices are balanced, the pseudoflow is a ***circulation***. The ***cost*** of a  
203 pseudoflow is defined to be

$$204 \quad \text{cost}(f) := \sum_{v \rightarrow w \in \text{supp}(f)} c(v \rightarrow w) \cdot f(v \rightarrow w).$$

205 The ***minimum-cost flow problem (MCF)*** asks to find a circulation of minimum cost inside a  
206 given directed graph.

207 **Residual graph.** Given a pseudoflow  $f$ , one can define the ***residual network*** as follows. Recall  
208 that the set of ***residual arcs***  $\vec{E}_f$  under  $f$  are those arcs  $v \rightarrow w$  satisfying  $f(v \rightarrow w) < u(v \rightarrow w)$ . In other  
209 words, an arc that is not saturated by  $f$  is a residual arc; similarly, given an arc  $v \rightarrow w$  with positive  
210 flow value, the backward arc  $w \rightarrow v$  is a residual arc.

211 Let  $N = (V, \vec{E}, c, u, \phi)$  be a network constructed from graph  $G$ , with a pseudoflow  $f$  on  $N$ .  
212 The ***residual graph***  $G_f$  of  $f$  has  $V$  as its vertex set and  $\vec{E}_f$  as its arc set. The ***residual capacity***  
213  $u_f$  with respect to pseudoflow  $f$  is defined to be  $u_f(v \rightarrow w) := u(v \rightarrow w) - f(v \rightarrow w)$ . Observe that  
214 the residual capacity is always nonnegative. We can define residual arcs differently using residual  
215 capacities:

$$216 \quad \vec{E}_f = \{v \rightarrow w \mid u_f(v \rightarrow w) > 0\}.$$

217 In other words, the set of residual arcs are precisely those arcs in the residual graph, each of  
218 which has nonzero residual capacity.

219 **LP-duality and admissibility.** To solve the minimum-cost flow problem, we focus on the  
220 primal-dual algorithms using linear programming. Let  $G = (V, E)$  be a given directed graph with  
221 the corresponding network  $N = (V, \vec{E}, c, u, \phi)$ . Formally, the ***potentials***  $\pi(v)$  are the variables of  
222 the linear program dual to the standard linear program for the minimum-cost flow problem with  
223 variables  $f(v, w)$  for each directed edge in  $E$ . Assignments to the primal variables satisfying the  
224 capacity constraints extend naturally into a pseudoflow on the network  $N$ . Let  $G_f = (V, \vec{E}_f)$  be the  
225 residual graph under pseudoflow  $f$ . The ***reduced cost*** of an arc  $v \rightarrow w$  in  $\vec{E}_f$  with respect to  $\pi$  is  
226 defined as

$$227 \quad c_\pi(v \rightarrow w) := c(v \rightarrow w) - \pi(v) + \pi(w).$$

228 Notice that the cost function  $c_\pi$  is also antisymmetric.

229 The ***dual feasibility constraint*** says that  $c_\pi(v \rightarrow w) \geq 0$  holds for every directed edge  $(v, w)$   
230 in  $E$ ; potentials  $\pi$  which satisfy this constraint are said to be ***feasible***. Suppose we relax the dual  
231 feasibility constraint to allow some small violation in the value of  $c_\pi(v \rightarrow w)$ . We say that a pair  
232 of pseudoflow  $f$  and potential  $\pi$  is  ***$\epsilon$ -optimal*** [?, ?] if  $c_\pi(v \rightarrow w) \geq -\epsilon$  for every residual arc  $v \rightarrow w$   
233 in  $\vec{E}_f$ . Pseudoflow  $f$  is  ***$\epsilon$ -optimal*** if it is  $\epsilon$ -optimal with respect to some potentials  $\pi$ ; potential  
234  $\pi$  is  ***$\epsilon$ -optimal*** if it is  $\epsilon$ -optimal with respect to some pseudoflow  $f$ . Given a pseudoflow  $f$  and  
235 potentials  $\pi$ , a residual arc  $v \rightarrow w$  in  $\vec{E}_f$  is ***admissible*** if  $c_\pi(v \rightarrow w) \leq 0$ . We say that a pseudoflow  
236  $g$  in  $G_f$  is ***admissible*** if all support arcs of  $g$  on  $G_f$  are admissible; in other words,  $g(v \rightarrow w) > 0$   
237 holds only on admissible arcs  $v \rightarrow w$ .

238 ▶ **Lemma 3.1.** *Let  $f$  be an  $\epsilon$ -optimal pseudoflow in  $G$  and let  $f'$  be an admissible flow in  $G_f$ . Then  
239  $f + f'$  is also  $\epsilon$ -optimal. **«Lemma 5.3 in [5]? Also, where is this used?»***

### 240 3.2 Reduction to Unit-Capacity Min-Cost Flow Problem

241 The goal of the subsection is to reduce the minimum-cost partial matching problem to the unit-  
 242 capacity minimum-cost flow problem with a polynomial bound on diameter. To this end we  
 243 first provide an upper bound on the size of support of an integral pseudoflow on the standard  
 244 reduction network between the two problems. This upper bound in turn provides an additive  
 245 approximation on the cost of an  $\varepsilon$ -optimal circulation. Next we employ a technique by Sharathku-  
 246 mar and Agarwal [9] to transform an additive  $\varepsilon$ -approximate solution into a multiplicative  
 247  $(1 + \varepsilon)$ -approximation for the geometric partial matching problem. The reduction does not work  
 248 out of the box, as Sharathkumar and Agarwal were tackling a similar but different problem on  
 249 geometric transportations.

250 ▶ **Lemma 3.2.** *Computing a  $(1 + \varepsilon)$ -approximate geometric partial matching can be reduced to  
 251 the following problem in  $O(n \text{polylog } n)$  time: Given a reduction network  $N$  over a point set with  
 252 diameter at most  $K \cdot kn^3$  for some constant  $K$ , compute a  $(K \cdot \varepsilon / 6k)$ -optimal circulation on  $N$ .*

253 **Additive approximation.** Given a bipartite graph  $G = (A, B, E_0)$  for the geometric partial  
 254 matching problem with cost function  $c$ , we construct the **reduction network**  $N_H$  as follows:  
 255 Direct the edges in  $E_0$  from  $A$  to  $B$ , and assign each directed edge with capacity 1. Now add a  
 256 dummy vertex  $s$  with directed edges to all vertices in  $A$ , and add a dummy vertex  $t$  with directed  
 257 edges from all vertices in  $B$ ; each edge added this way has cost 0 and capacity 1. Denote the new  
 258 graph with vertex set  $V = A \cup B \cup \{s, t\}$  and edge set  $E$  as the **reduction graph**  $H$ . Assign vertex  $s$   
 259 with supply  $k$  and vertex  $t$  with demand  $k$ ; the rest of the vertices in  $H$  have zero supply-demand  
 260 We call the network naturally corresponds to  $H$  as the **reduction network**, denoted by  $N_H$ .

261 It is straightforward to show that any integer circulation  $f$  on  $N_H$  uses exactly  $k$  of the  $A$ -to- $B$   
 262 arcs, which correspond to the edges of a size- $k$  matching  $M_f$ . Notice that the cost of the circulation  
 263  $f$  is equal to the cost of the corresponding matching  $M_f$ . In other words, a  $(1 + \varepsilon)$ -approximation  
 264 to the MCF problem on the reduction network  $N_H$  translates to a  $(1 + \varepsilon)$ -approximation to the  
 265 geometric matching problem on the input graph  $G$ .

266 First we show that the number of arcs used by any integer pseudoflow in  $N_H$  is asymptotically  
 267 bounded by the excess of the pseudoflow.

268 ▶ **Lemma 3.3.** *Let  $f$  be an integer circulation in the reduction network  $N_H$ . Then, the size of the  
 269 support of  $f$  is at most  $3k$ . As a corollary, the number of residual backward arcs is at most  $3k$ .*

270 Using the bound on the support size, we show that an  $\varepsilon$ -optimal integral circulation gives an  
 271 additive  $O(k\varepsilon)$ -approximation to the MCF problem.

272 ▶ **Lemma 3.4.** *Let  $f$  be an  $\varepsilon$ -optimal integer circulation in  $N_H$ , and  $f^*$  be an optimal integer  
 273 circulation for  $N_H$ . Then,  $\text{cost}(f) \leq \text{cost}(f^*) + 6k\varepsilon$ .*

274 **Multiplicative approximation.** Now we employ a technique from Sharathkumar and Agar-  
 275 wal [9] to convert the additive approximation into a multiplicative one. Here we sketch a proof  
 276 to Lemma 3.2; a complete proof can be found in the Appendix.

277 **Sketch Sharathkumar and Agarwal [9] with our modification.}}**

### 278 3.3 High-Level Description of Cost-Scaling Algorithm

282 Our main algorithm for the unit-capacity minimum-cost flow problem is based on the **cost-scaling**  
 283 technique, originally due to Goldberg and Tarjan [5]; Goldberg, Hed, Kaplan, and Tarjan [4]  
 284 applied the technique on unit-capacity networks. The algorithm finds  $\varepsilon$ -optimal circulations

285 for geometrically shrinking values of  $\varepsilon$ . Each fixed value of  $\varepsilon$  is called a *cost scale*. Once  $\varepsilon$  is  
 286 sufficiently small, the  $\varepsilon$ -optimal flow is a suitable approximation according to Lemma 3.2.<sup>2</sup>

287 The cost-scaling algorithm initializes the flow  $f$  and the potential  $\pi$  to be zero. Note that the  
 288 zero flow is trivially a  $kC$ -optimal flow. At the beginning of each scale starting at  $\varepsilon = kC$ ,

289 ■■■ SCALE-INIT takes the previous circulation (now  $2\varepsilon$ -optimal) and transforms it into an  $\varepsilon$ -optimal  
 290 pseudoflow with  $O(k)$  excess.

291 ■■■ REFINE then reduces the excess in the newly constructed pseudoflow to zero, making it an  
 292  $\varepsilon$ -optimal circulation.

293 Thus, the algorithm produces an  $\varepsilon^*$ -optimal circulation after  $O(\log(kC/\varepsilon^*))$  scales. Using the  
 294 reduction in Lemma 3.2, we have the diameter of the point set, thus maximum cost  $C$ , bounded  
 295 by  $O(K \cdot kn^3)$  for some value  $K$ . By setting  $\varepsilon^*$  to be  $K \cdot \varepsilon/6k$ , the number of cost scales is bounded  
 296 above by  $O(\log(n/\varepsilon))$ .

297 **Scale initialization.** Recall that  $H$  is the *reduction graph* and  $N_H$  is the *reduction network*, both  
 298 constructed in Section 3.2. The vertex set of  $H$  consists of two point sets  $A$  and  $B$ , as well as two  
 299 dummy vertices  $s$  and  $t$ . The directed edges in  $H$  are pointed from  $s$  to  $A$ , from  $A$  to  $B$ , and from  
 300  $B$  to  $t$ . We call those arcs in  $N_H$  whose direction is consistent with their corresponding directed  
 301 edges as *forward arcs*, and those arcs that points in the opposite direction as *backward arcs*.

302 The procedure SCALE-INIT transforms a  $2\varepsilon$ -optimal circulation from the previous cost scale  
 303 into an  $\varepsilon$ -optimal flow with  $O(k)$  excess, by raising the potentials  $\pi$  of all vertices in  $A$  by  $\varepsilon$ , those  
 304 in  $B$  by  $2\varepsilon$ , and the potential of  $t$  by  $3\varepsilon$ . The potential of  $s$  remains unchanged. Now the reduced  
 305 cost of every forward arc is dropped by  $\varepsilon$ , and thus all the forward arcs have reduced cost at  
 306 least  $-\varepsilon$ .

307 As for backward arcs, the procedure SCALE-INIT continues by setting the flow on  $v \rightarrow w$  to  
 308 zero for each backward arc  $w \rightarrow v$  violating the  $\varepsilon$ -optimality constraint. In other words, we set  
 309  $f(v \rightarrow w) = 0$  whenever  $c_\pi(w \rightarrow v) < -\varepsilon$ . This ensures that all such backward arcs are no longer  
 310 residual, and therefore the flow (now with excess) is  $\varepsilon$ -optimal.

311 Because the arcs are of unit-capacity in  $N_H$ , each arc desaturation creates one unit of excess.  
 312 By Lemma 3.3 the number of backward arcs is at most  $3k$ . Thus the total amount of excess  
 313 created is also  $O(k)$ .

314 In total, potential updates and backward arc desaturations, thus the whole procedure SCALE-  
 315 INIT, take  $O(n)$  time.

316 **Refinement.** The procedure REFINE is implemented using a primal-dual augmentation al-  
 317 gorithm, which sends flows on admissible arcs to reduce the total excess, like the Hungarian  
 318 Algorithm. Unlike the Hungarian Algorithm, it uses *blocking flows* instead of augmenting paths.  
 319 An *augmenting path* is a path in the residual network from an excess vertex to a deficit vertex.  
 320 We call a pseudoflow  $f$  on residual network  $N_g$  a *blocking flow* if  $f$  saturates at least one residual  
 321 arc in every augmenting path in  $N_g$ . In other words, there is no admissible augmenting path in  
 322  $N_{f+g}$  from an excess vertex to a deficit vertex.

323 Each iteration of REFINE finds an admissible blocking flow that is then added to the current  
 324 pseudoflow in two stages:

---

279 <sup>2</sup> When the costs are integers, an  $\varepsilon$ -optimal circulation for a sufficiently small  $\varepsilon$  (say less than  $1/n$ ) is itself an  
 280 optimal solution [4,5]. We present this algorithm without the integral-cost assumption because in the geometric  
 281 partial matching setting (with respect to  $L_p$  norms) the costs are generally not integers.

- 325 1. A *Hungarian search*, which increases the dual variables  $\pi$  of vertices that are reachable from  
 326 an excess vertex by at least  $\varepsilon$ , in a Dijkstra-like manner, until there is an excess-deficit path of  
 327 admissible edges.
- 328 2. A *depth-first search* through the set of admissible edges to construct an admissible blocking  
 329 flow. It suffices to repeatedly extract admissible augmenting paths until no more admissible  
 330 excess-deficit paths remain. By definition, the union of such paths is a blocking flow. **(Move**  
 331 **to where the blocking flow is introduced?)**

332 The algorithm continues until the total excess becomes zero and the  $\varepsilon$ -optimal flow is now a  
 333 circulation.

334 First we analyze the number of iterations executed by REFINE. The proof follows the strategy  
 335 in Goldberg *et al.* [4, Section 3.2]. **(and maybe §5 of Goldberg-Tarjan?)** To this end we need  
 336 a bound on the size of the support of  $f$  right before and throughout the execution of REFINE.

337 ▶ **Lemma 3.5.** *Let  $f$  be an integer pseudoflow in  $N_H$  with  $O(k)$  excess. Then, the size of the support  
 338 of  $f$  is at most  $O(k)$ .*

339 ▶ **Corollary 3.6.** *The size of  $\text{supp}(f)$  is at most  $O(k)$  for pseudoflow  $f$  right before or during the  
 340 execution of REFINE.*

341 ▶ **Lemma 3.7.** *Let  $f$  be a pseudoflow in  $N_H$  with  $O(k)$  excess. The procedure REFINE runs for  
 342  $O(\sqrt{k})$  iterations before the excess of  $f$  becomes zero.*

343 The goal of the next section is to show that after  $O(n \text{polylog } n)$  time preprocessing, each  
 344 Hungarian search and depth-first search can be implemented in  $O(k \text{polylog } n)$  time. Combined  
 345 with the  $O(\sqrt{k})$  bound on the number of iterations we just proved, the procedure REFINE can be  
 346 implemented in  $O((n + k\sqrt{k}) \text{polylog } n)$  time. Together with our analysis on scale initialiation and  
 347 the bound on number of cost scales, this concludes the proof to Theorem 1.2.

## 348 4 Fast Implementation

349 Both Hungarian search and depth-first search are implemented in a Dijkstra-like fashion, traversing  
 350 through the residual graph using admissible arcs starting from the excess vertices. Each step of  
 351 the search procedures *relaxes* a minimum-reduced-cost arc from the set of visited vertices to an  
 352 unvisited vertex, until a deficit vertex is reached. At a high level, our analysis strategy is to charge  
 353 the relaxation events to the support arcs of  $f$ , which has size at most  $O(k)$  by Corollary 3.6.

### 354 4.1 Null vertices and shortcut graph

355 **(A figure might be helpful for this section.)**

356 As it turns out, there are some vertices visited by a relaxation event which we cannot charge  
 357 to  $\text{supp}(f)$ . Unfortunately the number of such vertices can be as large as  $\Omega(n)$  (consider the  
 358 residual graph under the zero flow). To overcome this issue, we replace the residual graph  
 359 with an equivalent graph that excludes all the null vertices, and run the Hungarian search and  
 360 depth-first search on the resulting graph instead.

361 **Null vertices.** We say a vertex  $v$  in the residual graph  $N_f$  is a *null vertex* if  $\phi_f(v) = 0$  and no  
 362 arcs of  $\text{supp}(f)$  is incident to  $v$ . We use  $A_\emptyset$  and  $B_\emptyset$  to denote the null vertices  $A$  and  $B$  respectively.  
 363 Vertices that are not null are called *normal vertices*. A *null 2-path* is a length-2 subpath in  $N_f$   
 364 from a normal vertex to another normal vertex, passing through a null vertex. As every vertex in  
 365  $A$  has in-degree 1 and every vertex in  $B$  has out-degree 1 in the residual graph, the null 2-paths

366 must be of the form either  $(s, v, b)$  for some vertex  $b$  in  $B \setminus B_\emptyset$  or  $(a, v, t)$  for some vertex  $a$  in  
 367  $A \setminus A_\emptyset$ . In either case, we say that the null 2-path **passes through** null vertex  $v$ . Similarly, we  
 368 define the length-3 path from  $s$  to  $t$  that passes through two null vertices to be a **null 3-path**.  
 369 Because reduced costs telescope for residual paths, the reduced cost of any null 2-path or null  
 370 3-path does not depend on the null vertices it passes through.

371 **Shortcut graph.** We construct the **shortcut graph**  $\tilde{H}_f$  from the reduction network  $H$  by  
 372 removing all null vertices and their incident edges, followed by inserting an arc from the head of  
 373 each each null path  $\Pi$  to its tail, with cost equals to the sum of costs on the arcs. We call this  
 374 arc the **shortcut** of null path  $\Pi$ , denoted as **short**( $\Pi$ ). The resulting multigraph  $\tilde{H}_f$  contains only  
 375 normal vertices of  $H_f$ , and the reduced cost of any path between normal vertices are preserved.  
 376 We argue now that  $\tilde{H}_f$  is fine as a surrogate for  $H_f$ . Let  $\tilde{\pi}$  be an  $\varepsilon$ -optimal potential on  $\tilde{H}_f$ .  
 377 Construct potentials  $\pi$  on  $H_f$  which extends  $\tilde{\pi}$  to null vertices, by setting  $\pi(a) := \tilde{\pi}(s)$  for  $a \in A_\emptyset$   
 378 and  $\pi(b) := \tilde{\pi}(t)$  for  $b \in B_\emptyset$ .

379 ▶ **Lemma 4.1.** Consider  $\tilde{\pi}$  an  $\varepsilon$ -optimal potential on  $\tilde{H}_f$  and  $\pi$  the corresponding potential  
 380 constructed on  $H_f$ . Then,

- 381 1. potential  $\pi$  is  $\varepsilon$ -optimal on  $H_f$ , and
- 382 2. if arc  $\text{short}(\Pi)$  is admissible under  $\tilde{\pi}$ , then every arc in  $\Pi$  is admissible under  $\pi$ .

## 383 4.2 Dynamic data structures for search procedures

384 **Hungarian search.** *«Shortly describe the Hungarian search and depth-first search  
 385 implementations.»*

386 Conceptually, we are executing the Hungarian search on the shortcut graph  $\tilde{H}_f$ . We describe  
 387 how we can query the minimum-reduced cost arc leaving  $\tilde{S}$  in  $O(\text{polylog } n)$  time for the shortcut  
 388 graph, without constructing  $\tilde{H}_f$  explicitly. For this purpose, let  $S$  be a set of “reached” vertices  
 389 maintained, identical to  $\tilde{S}$  except whenever a shortcut is relaxed, we add the null vertices passed  
 390 by the corresponding null path to  $S$  in addition to its (normal) endpoints. Observe that the arcs  
 391 of  $\tilde{H}_f$  leaving  $\tilde{S}$  fall into  $O(1)$  categories.

392 By construction, the distance returned by each of the BCP data structure is equal to the  
 393 reduced cost of the shortcut, which is equal to the reduced cost of the corresponding null path.  
 394 Each of the above data structures requires one query per relaxation, and an update operation  
 395 whenever a new vertex moves into  $\tilde{S}$ . The data structures above can perform both queries and  
 396 updates in  $O(\text{polylog } n)$  time each, so the running time of the Hungarian search other than the  
 397 potential updates can be charged to the number of relaxation steps.

398 **Depth-first search.** The depth-first search is similar to Hungarian search in that it uses the  
 399 relaxation of minimum-reduced-cost arcs/null paths, this time to identify admissible arcs/null  
 400 paths in a depth-first manner. This requires some adjustments to the data structures for finding  
 401 the minimum-reduced-cost arc leaving  $v' \in \tilde{S}$ .

402 Each data structure performs a constant number of queries and updates per relaxation, each  
 403 of which can be implemented in  $O(\text{polylog } n)$  time []; so the running time is again bounded by  
 404  $O(\text{polylog } n)$  times the number of relaxations.

## 405 4.3 Time analysis

406 First we bound the number of relaxations performed by both the Hungarian search and the  
 407 depth-first search.

408 ► **Lemma 4.2.** Both Hungarian search and depth-first search performs  $O(k)$  relaxations before a  
409 deficit vertex is reached.

410 Now we complete the time analysis by showing that each Hungarian search and depth-  
411 first search can be implemented in  $O(k \text{polylog } n)$  time after a one-time  $O(n \text{polylog } n)$ -time  
412 preprocessing.

413 ► **Lemma 4.3.** After  $O(n \text{polylog } n)$ -time preprocessing, each Hungarian search and depth-first  
414 search can be implemented in  $O(k \text{polylog } n)$  time.

#### 415 4.4 Number of potential updates on null vertices

416 In our implementation of REFINE, we do not explicitly construct  $\tilde{H}_f$ ; instead we query its edges  
417 using BCP/NN oracles and min/max heaps on elements of  $H_f$ . Potentials on the null vertices are  
418 only required right before an augmentation sends a flow through a null path, making the null  
419 vertices it passes normal. We use the construction from Lemma 4.1 to obtain potential  $\pi$  on  $H_f$   
420 such that the flow  $f$  is both  $\epsilon$ -optimal and admissible with respect to  $\pi$ .

421 **Size of blocking flows.** Now we bound the total number arcs whose flow is updated by a  
422 blocking flow during the course of REFINE. This bounds both the time spent updating the flow on  
423 these arcs and also the time spent on null vertex potential updates (Lemma 4.5).

424 ► **Lemma 4.4.** The support of each blocking flow found in REFINE is of size  $O(k)$ . 3.5?

425 ► **Lemma 4.5.** The number of end-of-REFINE null vertex potential updates is  $O(n)$ . The number of  
426 augmentation-induced null vertex potential updates in each invocation of REFINE is  $O(k \log k)$ .

427 Now combining Lemma B.3, Lemma B.4, and Lemma 4.5 completes the proof of Theorem 1.2.

## 428 5 Unbalanced transportation Lemma 4.3

429 Finally, In this section, we give an exact algorithm which solves the transportation problem in  $O(rn(r +$   
430  $\sqrt{n}) \text{polylog } n)$  time, proving Theorem 1.3. This algorithm is a geometric implementation of  
431 the uncapacitated min-cost flow algorithm due to Orlin [8], combined with some of the tools  
432 developed in Sections 2 and 3. Mainly, we batch potential updates and use the rewinding  
433 mechanism to initialize each Hungarian search in time proportional to the previous Hungarian  
434 search.

435 Let  $A$  and  $B$  be points in the plane with  $r = |A|$  and  $n = |B|$ . Let  $\lambda : A \cup B \rightarrow \mathbb{Z}$  be a  
436 supply-demand function with positive value on points of  $A$ , negative value on points of  $B$ , and  
437  $\sum_{a \in A} \lambda(a) = -\sum_{b \in B} \lambda(b)$ . We use  $U := \max_{p \in A \cup B} |\lambda(p)|$ . A transportation map is a function  
438  $\tau : A \times B \rightarrow \mathbb{R}_{\geq 0}$ . A transportation map  $\tau$  is feasible if  $\sum_{b \in B} \tau(a, b) = \lambda(a)$  for all  $a \in A$ , and  
439  $\sum_{a \in A} \tau(a, b) = -\lambda(b)$  for all  $b \in B$ . In other words, the value  $\tau(a, b)$  describes how much supply  
440 at  $a$  should be sent to meet demands at  $b$ , and we require that all supplies are sent and all  
441 demands are met. We define the cost of  $\tau$  to be

$$442 \text{cost}(\tau) := \sum_{(a,b) \in A \times B} \|a - b\|_p^q \cdot \tau(a, b). \quad \text{the size of } A \leq B.$$

443 Given  $A$ ,  $B$ , and  $\lambda$ , the transportation problem asks to find a feasible transportation map of  
444 minimum cost. We focus on analyzing the unbalanced setting where  $r \leq n$ .

445 There is a simple reduction from the transportation problem to uncapacitated min-cost flow.  
446 Consider the complete bipartite graph  $G$  between  $A$  and  $B$  (with all edges directed  $A$  to  $B$ ), set

the costs  $c(a, b) \leq \|a - b\|_p^q$ , all capacities to infinity, and use  $\phi = \lambda$ . Any circulation  $f$  in the network  $N = (G, c, u, \phi)$  can be converted into a feasible transportation map  $\tau_f$  by taking  $\tau_f(a, b) := f(a \rightarrow b)$ . Furthermore,  $\text{cost}(f) = \text{cost}(\tau_f)$ .

## 5.1 Uncapacitated MCF by excess scaling

We give an outline of the strongly polynomial algorithm for uncapacitated MCF from Orlin [8]. Orlin's algorithm follows an *excess-scaling* paradigm originally due to Edmonds and Karp [2]. Consider the basic primal-dual skeleton used in the previous sections: The algorithm begins with  $f = 0, \pi = 0$ , then repeatedly runs a *Hungarian search* that raises potentials (while maintaining dual feasibility) to create an admissible augmenting excess-deficit path, on which it then augments flow. If supplies/demands are integral and each augmentation is at least one unit of flow, then such an algorithm terminates. In terms of cost,  $f$  is maintained to be 0-optimal with respect to  $\pi$  and augmentations over admissible edges preserve this by Lemma 3.1. Thus, the final circulation must be optimal. The excess-scaling paradigm tunes this skeleton by specifying (i) between which excess and deficit vertices we augment, and (ii) how much flow is sent by the augmentation.

The excess-scaling algorithm maintains a *scale parameter*  $\Delta$ , initially  $\Delta = U$ . Let vertices with  $|\phi_f(v)| \geq \Delta$  be *active*. Each augmenting path is chosen between an active excess vertex and an active deficit vertex. Once there are either no more active excess or no more active deficit vertices,  $\Delta$  is halved. Each sequence of augmentations where  $\Delta$  holds a constant value is called an *excess scale*. There are  $O(\log U)$  excess scales before  $\Delta < 1$  and, by integrality of supplies/demands,  $f$  is a circulation.

With some modifications to the excess-scaling algorithm, Orlin [8] obtains an algorithm with a strongly polynomial bound on the number of augmentations and excess scales. First, an *active* vertex is redefined to be one where  $|\phi_f(v)| \geq \alpha\Delta$ , for a parameter  $\alpha \in (1/2, 1)$ . Second, arcs which have  $f(v \rightarrow w) \geq 3n\Delta$  at the beginning of a scale are *contracted*, creating a new vertex  $\hat{v}$  which inherits all the arcs of  $v$  and  $w$ , and has  $\phi(\hat{v}) = \phi(v) + \phi(w)$ . We use  $\hat{G} = (\hat{V}, \hat{E})$  to denote the resulting *contracted graph*, where each  $\hat{v} \in \hat{V}$  is a contracted component of vertices from  $V$ . Intuitively, the flow is so high on contracted arcs that no set of future augmentations can remove the arc from  $\text{supp}(f)$ . Third,  $\Delta$  is lowered to  $\max_{v \in V} \phi_f(v)$  if there are no active excess vertices, and  $f(v, w) = 0$  on all non-contracted arcs  $(v, w) \in \hat{E}$ . Finally, flow values are not tracked within contracted components, but once an optimal circulation is found on  $\hat{G}$ , optimal potentials  $\pi^*$  can be *recovered* for  $G$  in linear time by sequentially undoing the contractions. The algorithm performs a post-processing step which finds the optimal circulation  $f^*$  on  $G$  by solving a max-flow problem on the set of admissible arcs under  $\pi^*$ .

► **Theorem 5.1** Orlin [8], Theorems 2 and 3. Orlin's algorithm finds optimal potentials after  $O(n \log n)$  scaling phases and  $O(n \log n)$  total augmentations.

The remainder of the section focuses on showing that each augmentation can be implemented in  $O(r(r/\sqrt{n} + \sqrt{n} \log n))$  time (after preprocessing). Additionally, we show that  $f^*$  can be recovered from  $\pi^*$  very quickly for our specific graph.

**Implementing contractions.** Following Agarwal et al. [1], our geometric data structures must deal with real points  $(A, B)$ , rather than the contracted components  $(\hat{V})$ . We will track the contracted components described in  $\hat{G}$  (e.g. with a disjoint-set data structure) and mark the arcs of  $\text{supp}(f)$  that get contracted. We maintain potentials on the points  $(A \cup B)$  directly, instead of the contracted components  $(\hat{V})$ .

When conducting the Hungarian search, we initialize  $S$  with all vertices from *active excess contracted components* who (in sum) meet the imbalance criteria. Upon relaxing any  $v \in \hat{V}$ ,

492 we immediately relax all the contracted  $\text{supp}(f)$  arcs which span  $\hat{v}$ . Since the input network  
 493 is uncapacitated, each contracted component is strongly connected in the residual network by  
 494 the admissible forward/reverse arcs of each contracted arc. For relaxing arcs of  $\hat{E}$ , we relax  
 495 support arcs before attempting to relax any non-support arcs. Relaxations of support arcs can be  
 496 performed without further potential changes, since they are admissible by invariant.

497 During ~~the~~ augmentations, contracted residual arcs are considered to have infinite capacity, and  
 498 we do not update the value of flow on these arcs. We allow augmenting paths to begin from  
 499 any ~~some~~  $a \in \hat{v} \cap A$  of an active excess component  $\hat{v}$ , and end in any  $b \in \hat{w} \cap B$  of an active deficit  
 500 component  $\hat{w}$ . ~~for some~~

501 **Recovering optimal flow.** ~~Summarize the results.~~ Naively recovering an optimal flow  
 502 requires  $O(rn^2)$  time. Use a strategy from Agarwal et al. [1], we can recover the optimal flow in  
 503 time  $O(n \text{polylog } n)$ . If furthermore the cost function is the  $p$ -norm (without the  $q$ -power), an  
 504 even stronger result stands: In this case, the set of admissible arcs under an optimal potential  
 505 forms a planar graph, and thus we can apply the planar maximum-flow algorithm [?, 3] which  
 506 runs in  $O(n \log n)$  time.

## 5.2 Dead vertices and support stars

507 Given Theorem 5.1, our goal is to implement each augmentation in  $O(r(r/\sqrt{n} + n^{1/2}) \text{polylog } n)$   
 508 time. To find an augmenting path, we again use a Hungarian search with geometric data structures  
 509 to perform relaxations quickly. Like in Section 3, there are vertices which cannot be charged to  
 510 the flow support. Even worse, the flow support may have size  $\Omega(n)$  (consider an instance with  
 511  $r=1$ , and demand uniformly distributed among the vertices of  $B$ ). Our strategy is summarized  
 512 as follows:

- 514   ■ Discard vertices which lead to dead ends in the search (are not on a path to a deficit vertex).
- 515   ■ Cluster parts of the flow support, such that the number of support arcs outside clusters is  
        $O(r)$ . The number of relaxations we perform is proportional to the number of support arcs  
       outside of clusters. *Lemma 2*
- 518   ■ Querying/updating clusters degrades our amortized time per relaxation from  $O(\text{polylog } n)$  to  
        $O((r/\sqrt{n} + \sqrt{n}) \text{polylog } n)$ .

520 Let  $E(\text{supp}(f)) := \{(v, w) \mid v \rightarrow w \in \text{supp}(f)\}$  be the set of undirected edges corresponding to  
 521 the arcs in  $\text{supp}(f)$ . Clearly,  $|\text{supp}(f)| = |E(\text{supp}(f))|$ . Let the **support degree** of a vertex be its  
 522 degree in  $E(\text{supp}(f))$ .

523 **Dead vertices.** We call a vertex  $b \in B$  **dead** if it has support degree 0 and is not an active  
 524 excess or deficit vertex; call it **live** otherwise. Dead vertices are essentially equivalent to the **null**  
 525 vertices of Section 3. Since the reduction in this section does not use a super-source/super-sink,  
 526 we can simply remove these from consideration during a Hungarian search — they will not  
 527 terminate the search, and have no outgoing residual arcs. Like null vertices, we ignore dead  
 528 vertex potentials and infer feasible potentials when they become live again. We use  $A_\ell$  and  $B_\ell$   
 529 to denote the living vertices of points in  $A$  and  $B$ , respectively. Note that being dead/alive is a  
 530 notion strictly for vertices, and not contracted components.

531 We say a dead vertex is **revived** when it stops meeting either condition of the definition. Dead  
 532 vertices are only revived after  $\Delta$  decreases (say, in a subsequent excess scale) as no augmenting  
 533 path will cross a dead vertex and they cannot meet the criteria for contractions. When a dead  
 534 vertex is revived, we must add it back into each of our data structures and give it a feasible  
 535 potential. For revived  $b \in B$ , a feasible choice of potential is  $\pi(b) \leftarrow \max_{a \in A} \pi(a) - c(a, b)$  which

536 we can query by maintaining a weighted nearest neighbor data structure on the points of  $A$ . The  
 537 total number of revivals is bounded above by the number of augmentations: since the final flow  
 538 is a circulation on  $\hat{G}$  and a newly revived vertex  $v$  has no adjacent supp( $v$ ) arcs and cannot be  
 539 contracted, there is at least one subsequent augmentation which uses  $v$  as its beginning or end.  
 540 Thus, the total number of revivals is  $O(n \log n)$ .

541 **Support stars.** The vertices of  $B$  with support degree 1 are partitioned into subsets  $\Sigma_a \subset B$  by  
 542 the  $a \in A$  lying on the other end of their single support arc. We call  $\Sigma_a$  the *support star* centered  
 543 around  $a \in A$ .

544 Roughly speaking, we would like to handle each support star as a single unit. When the  
 545 Hungarian search reaches  $a$  or any  $b \in \Sigma_a$ , then the entirety of  $\Sigma_a$  (as well as  $a$ ) is also admissible-  
 546 reachable and can be included into  $S$  without further potential updates. Additionally, the only  
 547 outgoing residual arcs of every  $b \in \Sigma_a$  lead to  $a$ , the only way to leave  $\Sigma_a \cup \{a\}$  is through an  
 548 arc leaving  $a$ . Once a relaxation step reaches some  $b \in \Sigma_a$  or  $a$  itself, we would like to quickly  
 549 update the state such that the rest of  $b \in \Sigma_a$  is also reached without performing relaxation steps  
 550 to each individual  $b \in \Sigma_a$ .

### 551 5.3 Implementation details

552 Before describing our workaround for support stars, we analyze the number of relaxation steps  
 553 for arcs outside of support stars.

554 ► **Lemma 5.2.** *Suppose we have stripped the graph of dead vertices. The number of relaxation  
 555 steps in a Hungarian search outside of support stars is  $O(r)$ .*

556 The running time of a Hungarian search will be  $O(r)$  times the time it takes us to implement  
 557 each relaxation.

558 **Relaxations outside support stars.** For relaxations that don't involve support star vertices,  
 559 we can once again maintain a BCP to query the minimum  $A_\ell$ -to- $B_\ell$  arc. To elaborate, this is the  
 560 BCP between  $P = A_\ell \cap S$  and  $Q = (B_\ell \setminus (\bigcup_{a \in A_\ell} \Sigma_a)) \setminus S$ , weighted by potentials. This can be  
 561 queried in  $O(\log n)$  time and updated in  $O(\text{polylog } n)$  time per point. Since it doesn't deal with  
 562 support stars, there is at most one insertion/deletion per relaxation step.

563 For  $B_\ell$ -to- $A_\ell$ , backward (support) arcs are kept admissible by invariant, so we relax them  
 564 immediately when they arrive on the frontier.

565 **Relaxing a support star.** We classify support stars into two categories: *big stars* are those  
 566 with  $|\Sigma_a| > \sqrt{n}$ , and *small stars* are those with  $|\Sigma_a| \leq \sqrt{n}$ . Let  $A_{\text{big}} \subseteq A$  denote the centers of big  
 567 stars and  $A_{\text{small}} \subseteq A$  denote the centers of small stars. We keep the following data structures  
 568 to manage support stars.

- 569 1. For each big star  $\Sigma_a$ , we use a data structure  $\mathcal{D}_{\text{big}}(a)$  to maintain BCP between  $P = A_\ell \cap S$  and  
 570  $Q = \Sigma_a$ , weighted by potentials. We query this until  $a \in S$  or any vertex of  $\Sigma_a$  is added to  $S$ .
- 571 2. All small stars are added to a single BCP data structure  $\mathcal{D}_{\text{small}}$  between  $P = A_\ell \cap S$  and  
 572  $Q = (\bigcup_{a \in A_{\text{small}}} \Sigma_a) \setminus S$ , weighted by potentials. When an  $a \in A_{\text{small}}$  or any vertex of its support  
 573 star is added to  $S$ , we remove the points of  $\Sigma_a$  from  $\mathcal{D}_{\text{small}}$  using  $|\Sigma_a|$  deletion operations.
- 574 We will update these data structures as each support star center is added into  $S$ . If a relaxation  
 575 step adds some  $b \in B_\ell$  and  $b$  is in a support star  $\Sigma_a$ , then we immediately relax  $b \rightarrow a$ , as all  
 576 support arcs are admissible. Relaxations of non-support star  $b \in B_\ell$  will not affect the support  
 577 star data structures.

578 Suppose a relaxation step adds some  $a \in A_\ell$  to  $S$ . For the support star data structures, we  
 579 must (i) remove  $a$  from every  $\mathcal{D}_{\text{big}}$ , (ii) remove  $a$  from  $\mathcal{D}_{\text{small}}$ . If  $a \in A_{\text{big}}$ , we also (iii) deactivate  
 580  $\mathcal{D}_{\text{big}}(a)$ . If  $a \in A_{\text{small}}$ , we also (iv) remove the points of  $\Sigma_a$  from  $\mathcal{D}_{\text{small}}$ . The operations (i–iii) can  
 581 be performed in  $O(\text{polylog } n)$  time each, but (iv) may take up to  $O(\sqrt{n} \text{ polylog } n)$  time.

582 On the other hand, there are now  $O(\sqrt{n})$  data structures to query during each relaxation step,  
 583 as there are  $O(n/\sqrt{n})$  data structures  $\mathcal{D}_{\text{big}}(\cdot)$ . Thus, the query time within each relaxation step is  
 584  $O(\sqrt{n} \log n)$ . We can now bound the time spent within the Hungarian search.

585 ▶ **Lemma 5.3.** *Hungarian search takes  $O(r\sqrt{n} \text{ polylog } n)$  time.*

586 **Updating support stars.** As the flow support changes, the membership of support stars may  
 587 shift and a big star may eventually become small (or vice versa). To efficiently support this,  
 588 introduce some fuzziness to when a star should be big or small. Standard charging argument  
 589 shows that the amortized update time is  $O(r\sqrt{n}(r + \sqrt{n}) \text{ polylog } n)$ .

590 Membership of support stars can only be changed by augmentations, so the number of  
 591 star membership changes by a single augmenting path is bounded above by twice its length  
 592 (each vertex may be removed from one star, and/or added to another star). Thus, individual  
 593 membership changes can be performed in  $O(\text{polylog } n)$  time each, and there are  $O(rn \log n)$  total.

594 **Preprocessing time.** To build the very first set of data structures, we take  $O(rn \text{ polylog } n)$   
 595 time. There are  $r|\Sigma_a|$  points in each  $\mathcal{D}_{\text{big}}(a)$ , but the  $\Sigma_a$  are disjoint, so the total points to insert  
 596 is  $O(rn)$ .  $\mathcal{D}_{\text{small}}$  also has at most  $O(rn)$  points. Each BCP data structure can be constructed in  
 597  $O(\text{polylog } n)$  times its size, so the total preprocessing time is  $O(rn \text{ polylog } n)$ .

598 **Between searches.** After an augmentation, we reset the above data structures to their initial  
 599 state plus the change from the augmentation using the rewinding mechanism. By reversing  
 600 the sequence of insertions/deletions to each data structure over the course of the Hungarian  
 601 search, we can recover the versions data structures as they were when the Hungarian search  
 602 began. This takes time proportional to the time of the Hungarian search,  $O(r\sqrt{n} \text{ polylog } n)$  by  
 603 Lemma 5.3. The most recent augmentation may have deactivated at most one active excess and  
 604 at most one active deficit, which we can update in the data structures in  $O(\sqrt{n} \text{ polylog } n)$  time.  
 605 Additionally, the augmentation may have changed the membership of some support stars, but  
 606 we analyzed the time for membership changes earlier. Finally, we note that an augmenting path  
 607 cannot reduce the support degree of a vertex to zero, and therefore no new dead vertices are  
 608 created by augmentation.

609 **Between excess scales.** When the excess scale changes, vertices that were previously inactive  
 610 may become active, and vertices that were dead may be revived (however, no active vertices  
 611 deactivate, and no live vertices die as the result of  $\Delta$  decreasing). If we have the data structures  
 612 built on the active excesses at the end of the previous scale, then we can add in each newly  
 613 active  $a \in A$  and charge this insertion to the (future) augmenting path or contraction which  
 614 eventually makes the vertex inactive, or absorbs it into another component. By Theorem 5.1,  
 615 there are  $O(n \log n)$  such newly active vertices. The time to perform data structure updates for  
 616 each of them is  $O(\sqrt{n} \text{ polylog } n)$ , so the total time spent bookkeeping newly active vertices is  
 617  $O(n^{3/2} \text{ polylog } n)$ .

618 **Putting it together.** After  $O(rn \text{ polylog } n)$  preprocessing, we spend  $O(r\sqrt{n} \text{ polylog } n)$  time  
 619 each Hungarian search by Lemma 5.3. After each augmentation, we spend the same amount

of time (plus  $O(\text{polylog } n)$  extra) to initialize data structures for the next Hungarian search. We spend up to  $O((rn + r^2\sqrt{n}) \text{polylog } n)$  total time making big-small star switching updates. We spend  $O(n^{3/2} \text{polylog } n)$  time activating and reviving vertices. Thus, the algorithm takes  $O(rn(r/\sqrt{n} + \sqrt{n}) \text{polylog } n)$  time to produce optimal potentials  $\pi^*$ , from which we can recover  $f^*$  in  $O(r\sqrt{n} \text{polylog } n)$  additional time. This completes the proof of Theorem 1.3.

**Acknowledgment.**

---

**References**

---

- 1 Pankaj K. Agarwal, Kyle Fox, Debmalya Panigrahi, Kasturi R. Varadarajan, and Allen Xiao. Faster algorithms for the geometric transportation problem. *33rd International Symposium on Computational Geometry, SoCG 2017, July 4-7, 2017, Brisbane, Australia*, 7:1–7:16, 2017. [⟨https://doi.org/10.4230/LIPIcs.SoCG.2017.7⟩](https://doi.org/10.4230/LIPIcs.SoCG.2017.7).
- 2 Jack Edmonds and Richard M. Karp. Theoretical improvements in algorithmic efficiency for network flow problems. *J. ACM* 19(2):248–264, 1972. [⟨https://doi.org/10.1145/321694.321699⟩](https://doi.org/10.1145/321694.321699).
- 3 Jeff Erickson. Maximum flows and parametric shortest paths in planar graphs. *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2010, Austin, Texas, USA, January 17-19, 2010*, 794–804, 2010. [⟨https://doi.org/10.1137/1.9781611973075.65⟩](https://doi.org/10.1137/1.9781611973075.65).
- 4 Andrew V. Goldberg, Sagi Hed, Haim Kaplan, and Robert E. Tarjan. Minimum-cost flows in unit-capacity networks. *Theory Comput. Syst.* 61(4):987–1010, 2017. [⟨https://doi.org/10.1007/s00224-017-9776-7⟩](https://doi.org/10.1007/s00224-017-9776-7).
- 5 Andrew V. Goldberg and Robert E. Tarjan. Finding minimum-cost circulations by successive approximation. *Math. Oper. Res.* 15(3):430–466, 1990. [⟨https://doi.org/10.1287/moor.15.3.430⟩](https://doi.org/10.1287/moor.15.3.430).
- 6 Haim Kaplan, Wolfgang Mulzer, Liam Roditty, Paul Seiferth, and Micha Sharir. Dynamic planar voronoi diagrams for general distance functions and their algorithmic applications. *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2017, Barcelona, Spain, Hotel Porta Fira, January 16-19, 2495–2504, 2017*. [⟨https://doi.org/10.1137/1.9781611974782.165⟩](https://doi.org/10.1137/1.9781611974782.165).
- 7 Harold W Kuhn. The Hungarian method for the assignment problem. *Naval Research Logistics (NRL)* 2(1-2):83–97. Wiley Online Library, 1955.
- 8 James B. Orlin. A faster strongly polynomial minimum cost flow algorithm. *Operations Research* 41(2):338–350, 1993. [⟨https://doi.org/10.1287/opre.41.2.338⟩](https://doi.org/10.1287/opre.41.2.338).
- 9 Lyle Ramshaw and Robert Endre Tarjan. A weight-scaling algorithm for min-cost imperfect matchings in bipartite graphs. *53rd Annual IEEE Symposium on Foundations of Computer Science, FOCS 2012, New Brunswick, NJ, USA, October 20-23, 2012*, 581–590, 2012. [⟨https://doi.org/10.1109/FOCS.2012.9⟩](https://doi.org/10.1109/FOCS.2012.9).
- 10 R. Sharathkumar and Pankaj K. Agarwal. Algorithms for the transportation problem in geometric settings. *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, 306–317, 2012. [⟨http://portal.acm.org/citation.cfm?id=2095145&CFID=63838676&CFTOKEN=79617016⟩](http://portal.acm.org/citation.cfm?id=2095145&CFID=63838676&CFTOKEN=79617016).
- 11 Pravin M. Vaidya. Geometry helps in matching. *SIAM J. Comput.* 18(6):1201–1225, 1989. [⟨https://doi.org/10.1137/0218080⟩](https://doi.org/10.1137/0218080).

663 **A Proofs from Section 3**

664 ► **Lemma 3.1.** Let  $f$  be an  $\varepsilon$ -optimal pseudoflow in  $G$  and let  $f'$  be an admissible flow in  $G_f$ . Then  
665  $f + f'$  is also  $\varepsilon$ -optimal. **(Lemma 5.3 in [GT90]? Also, where is this used?)**

666 **Proof.** Augmentation by  $f'$  will not change the potentials, so any previously  $\varepsilon$ -optimal arcs  
667 remain  $\varepsilon$ -optimal. However, it may introduce new arcs  $v \rightarrow w$  with  $u_{f+f'}(v \rightarrow w) > 0$ , that previously  
668 had  $u_f(v \rightarrow w) = 0$ . We will verify that these arcs satisfy the  $\varepsilon$ -optimality condition.

669 If an arc  $v \rightarrow w$  is newly introduced this way, then by definition of residual capacities  $f(v \rightarrow w) = u(v \rightarrow w)$ . At the same time,  $u_{f+f'}(v \rightarrow w) > 0$  implies that  $(f + f')(v \rightarrow w) < u(v \rightarrow w)$ . This means  
670 that  $f'$  augmented flow in the reverse direction of  $v \rightarrow w$  ( $f'(w \rightarrow v) > 0$ ). By assumption, the arcs  
671 of  $\text{supp}(f')$  are admissible, so  $w \rightarrow v$  was an admissible arc ( $c_\pi(w \rightarrow v) \leq 0$ ). By antisymmetry of  
672 reduced costs, this implies  $c_\pi(v \rightarrow w) \geq 0 \geq -\varepsilon$ . Therefore, all arcs with  $u_{f+f'}(v, w) > 0$  respect  
673 the  $\varepsilon$ -optimality condition, and thus  $f + f'$  is  $\varepsilon$ -optimal. ◀

675 ► **Lemma 3.3.** Let  $f$  be an integer circulation in the reduction network  $N_H$ . Then, the size of the  
676 support of  $f$  is at most  $3k$ . As a corollary, the number of residual backward arcs is at most  $3k$ .

677 **Proof.** Because  $f$  is a circulation,  $\text{supp}(f)$  can be decomposed into  $k$  paths from  $s$  to  $t$ . Each  
678  $s$ -to- $t$  path in  $N_H$  is of length three, so the size of  $\text{supp}(f)$  is at most  $3k$ . As every backward arc in  
679 the residual network must be induced by positive flow in the opposite direction, the total number  
680 of residual backward arcs is at most  $3k$ . ◀

681 ► **Lemma 3.4.** Let  $f$  be an  $\varepsilon$ -optimal integer circulation in  $N_H$ , and  $f^*$  be an optimal integer  
682 circulation for  $N_H$ . Then,  $\text{cost}(f) \leq \text{cost}(f^*) + 6k\varepsilon$ .

683 **Proof.** By Lemma 3.3, the total number of backward arcs in the residual network  $N_f$  is at most  
684  $3k$ . Consider the residual flow in  $N_f$  defined by the difference between  $f^*$  and  $f$ . Since both  $f$   
685 and  $f^*$  are both circulations and  $N_H$  has unit-capacity, the flow  $f - f^*$  is comprised of unit flows  
686 on a collection of edge-disjoint residual cycles  $\Gamma_1, \dots, \Gamma_\ell$ . Observe that each residual cycle  $\Gamma_i$  must  
687 have exactly half of its arcs being backward arcs, and thus we have  $\sum_i |\Gamma_i| \leq 6k$ .

688 Let  $\pi$  be some potential certifying that  $f$  is  $\varepsilon$ -optimal. Because  $\Gamma_i$  is a residual cycle, we have  
689  $c_\pi(\Gamma_i) = c(\Gamma_i)$  since the potential terms telescope. We then see that

$$690 \text{cost}(f) - \text{cost}(f^*) = \sum_i c(\Gamma_i) = \sum_i c_\pi(\Gamma_i) \geq \sum_i (-\varepsilon) \cdot |\Gamma_i| \geq -6k\varepsilon,$$

691 where the second-to-last inequality follows from the  $\varepsilon$ -optimality of  $f$  with respect to  $\pi$ . Rearranging  
692 the terms we have that  $\text{cost}(f) \leq \text{cost}(f^*) + 6k\varepsilon$ . ◀

693 Let  $T$  be the minimum spanning tree on input graph  $G$  and order its edges by increasing  
694 length as  $e_1, \dots, e_{r+n-1}$ . Let  $T_i$  denote the subgraph of  $T$  obtained by removing the heaviest  
695 edges in  $T$ . Let  $i$  be the largest index so that the optimal solution to the MPM problem has  
696 edges between components of  $T_i$ . Choose  $j$  to be the smallest index larger than  $i$  satisfying  
697  $c(e_j) \geq kn \cdot c(e_i)$ . For each component  $K$  of  $T_j$ , let  $G_K$  be the subgraph of  $G$  induced on vertices  
698 of  $K$ ; let  $A_K := K \cap A$  and  $B_K := K \cap B$ , respectively. We partition  $A$  and  $B$  into the collection of  
699 sets  $A_K$  and  $B_K$  according to the components  $K$  of  $T_j$ . Since  $j < i$ , the optimal partial matching in  
700  $G$  can be partitioned into edges between  $A_K$  and  $B_K$  within  $G_K$ ; no optimal matching edges lie  
701 between components.

702 ► **Lemma A.1 (Sharathkumar and Agarwal [SA12, §3.5]).** Let  $G = (A, B, E_0)$  be the input  
703 to MPM problem, and consider the partitions  $A_K$  and  $B_K$  defined as above. Let  $M^*$  be the optimal  
704 partial matching in  $G$ . Then,

- 705 (i)  $c(e_i) \leq \text{cost}(M^*) \leq kn \cdot c(e_i)$ , and  
 706 (ii) the diameter of  $G_K$  is at most  $kn^2 \cdot c(e_i)$  for every  $K \in T_j$ ,

707 To prove Lemma 3.2, we need to further modify the point set so that the cost of the optimal  
 708 solution does not change, while the diameter of the *whole* point set is bounded. Move the points  
 709 within each component in *translation* so that the minimum distances between points across  
 710 components are at least  $kn \cdot c(e_i)$  but at most  $O(n \cdot kn^2 \cdot c(e_i))$ . This will guarantee that the  
 711 optimal solution still uses edges within the components by Lemma A.1. The simplest way of  
 712 achieving this is by aligning the components one by one into a “straight line”, so that the distance  
 713 between the two farthest components is at most  $O(n)$  times the maximum diameter of the cluster.

714 Now one can prove Lemma 3.2 by computing an  $(\varepsilon c(e_i)/6k)$ -optimal circulation  $f$  on the  
 715 point set after translations using additive approximation from Lemma 3.4, together with the  
 716 bound  $c(e_i) \leq \text{cost}(M^*)$  from Lemma A.1.

717 One small problem remains: We need to show that such reduction can be performed in  
 718  $O(n \text{polylog } n)$  time. Sharathkumar and Agarwal [SA12] have shown that the partition of  $A$  and  
 719  $B$  into  $A_K$ s and  $B_K$ s can be computed in  $O(n \text{polylog } n)$  time, assuming that the indices  $i$  and  $j$  can  
 720 be determined in such time as well. However in our application the choice of index  $i$  depends on  
 721 the optimal solution of MPM problem which we do not know.

722 To solve this issue we perform a binary search on the edges  $e_1, \dots, e_{r+n-1}$ . **«Hmm, we have  
 723 no way to check Lemma 4.5(i); but in fact a polynomial bound is good enough.»**  
 724 **«UNRESOLVED ISSUE»**

725 ► **Lemma 3.5.** Let  $f$  be an integer pseudoflow in  $N_H$  with  $O(k)$  excess. Then, the size of the support  
 726 of  $f$  is at most  $O(k)$ .

727 **Proof.** Observe that the reduction graph  $H$  is a directed acyclic graph, and thus the support of  $f$   
 728 does not contain a cycle. Now  $\text{supp}(f)$  can be decomposed into a set of inclusion-maximal paths,  
 729 each of which contributes a single unit of excess to the flow if the path does not terminate at  $t$   
 730 or if more than  $k$  paths terminate at  $t$ . By assumption, there are  $O(k)$  units of excess to which  
 731 we can associate to the paths, and at most  $k$  paths (those that terminate at  $t$ ) that we cannot  
 732 associate with a unit of excess. The length of any such paths is at most three by construction of  
 733 the reduction graph  $H$ . Therefore we can conclude that the number of arcs in the support of  $f$  is  
 734  $O(k)$ . ◀

735 ► **Lemma 3.7.** Let  $f$  be a pseudoflow in  $N_H$  with  $O(k)$  excess. The procedure REFINE runs for  
 736  $O(\sqrt{k})$  iterations before the excess of  $f$  becomes zero.

737 **Proof.** Let  $f_0$  and  $\pi_0$  be the flow and potential at the start of the procedure REFINE. Let  $f$  and  $\pi$   
 738 be the current flow and the potential. Let  $d(v)$  defined to be the amount of potential increase at  
 739  $v$ , measured in units of  $\varepsilon$ ; in other words,  $d(v) := (\pi(v) - \pi_0(v))/\varepsilon$ .

740 Now divide the iterations executed by the procedure REFINE into two phases: The transition  
 741 from the first phase to the second happens when every excess vertex  $v$  has  $d(v) \geq \sqrt{k}$ . At most  
 742  $\sqrt{k}$  iterations belong to the first phase as each Hungarian search increases the potential  $\pi$  by at  
 743 least  $\varepsilon$  for each excess vertex (and thus increases  $d(v)$  by at least one).

744 The number of iterations belonging to the second phase is upper bounded by the amount  
 745 of total excess at the end of the first phase, because each subsequent push of a blocking flow  
 746 reduces the total excess by at least one. We now show that the amount of such excess is at most  
 747  $O(\sqrt{k})$ . Consider the set of arcs  $E^+ := \{v \rightarrow w \mid f(v \rightarrow w) < f_0(v \rightarrow w)\}$ . The total amount of excess is  
 748 upper bounded by the number of arcs in  $E^+$  that crosses an arbitrarily given cut  $X$  that separates  
 749 the excess vertices from the deficit vertices, when the network has unit-capacity [GHKT17,  
 750 Lemma 3.6]. Consider the set of cuts  $X_i := \{v \mid d(v) > i\}$  for  $0 \leq i < \sqrt{k}$ ; every such cut separates

751 the excess vertices from the deficit vertices at the end of first phase. Each arc in  $E^+$  crosses at  
 752 most 3 cuts of type  $X_i$  [GHKT17, Lemma 3.1]. So there is one  $X_i$  crossed by at most  $3|E^+|/\sqrt{k}$   
 753 arcs in  $E^+$ . The size of  $E^+$  is bounded by the sum of support sizes of  $f$  and  $f_0$ ; by Corollary 3.6  
 754 the size of  $E^+$  is  $O(k)$ . This implies an  $O(\sqrt{k})$  bound on the total excess after the first phase,  
 755 which in turn bounds the number of iterations in the second phase.  $\blacktriangleleft$

756 **B Proofs from Section 4**

757 ► **Lemma 4.1.** Consider  $\tilde{\pi}$  an  $\varepsilon$ -optimal potential on  $\tilde{H}_f$  and  $\pi$  the corresponding potential  
 758 constructed on  $H_f$ . Then,

- 759 1. potential  $\pi$  is  $\varepsilon$ -optimal on  $H_f$ , and  
 760 2. if arc  $\text{short}(\Pi)$  is admissible under  $\tilde{\pi}$ , then every arc in  $\Pi$  is admissible under  $\pi$ .

761 **Proof.** Reduced costs for any arc from a normal vertex another is unchanged under either  $\tilde{\pi}$  or  $\pi$ .  
 762 Recall that a null path is comprised of one  $A$ -to- $B$  arc, and one or two zero-cost arcs (connecting  
 763 the null vertex/vertices to  $s$  and/or  $t$ ). With our choice of null vertex potentials, we observe that  
 764 the zero-cost arcs still have zero reduced cost. It remains to prove that an arbitrary **((residual?))**  
 765 arc  $(a, b)$  **((arc or directed edge?))** satisfies the  $\varepsilon$ -optimality condition and admissibility when  
 766 either  $a$  or  $b$  is a null vertex.

767 By construction of the shortcut graph, there is always a null path  $\Pi$  that contains  $(a, b)$ .  
 768 Observe that  $c_\pi(a, b) = c_\pi(\Pi)$ , independent to the type of null path. Again by construction,  
 769  $c_\pi(\Pi) = c_{\tilde{\pi}}(\text{short}(\Pi))$ , so we have  $c_\pi(a, b) = c_{\tilde{\pi}}(\text{short}(\Pi)) \geq -\varepsilon$ . Additionally, if  $\text{short}(\Pi)$  is  
 770 admissible under  $\tilde{\pi}$ , then so is  $(a, b)$  under  $\pi$ . This proves the lemma.  $\blacktriangleleft$

- 
- 771 1. Non-shortcut backward arcs  $(v, w)$  with  $(w, v) \in \text{supp}(f)$ . For these, we can maintain a  
 772 min-heap on  $\text{supp}(f)$  arcs as each  $v$  arrives in  $\tilde{S}$ .  
 773 2. Non-shortcut  $A$ -to- $B$  forward arcs. For these, we can use a BCP data structure between  
 774  $(A \setminus A_\emptyset) \cap \tilde{S}$  and  $(B \setminus B_\emptyset) \setminus \tilde{S}$ , weighted by potential.  
 775 3. Non-shortcut forward arcs from  $s$ -to- $A$  and from  $B$ -to- $t$ . For  $s$ , we can maintain a min-heap  
 776 on the potentials of  $B \setminus \tilde{S}$ , queried while  $s \in \tilde{S}$ . For  $t$ , we can maintain a max-heap on the  
 777 potentials of  $A \cap \tilde{S}$ , queried while  $t \notin \tilde{S}$ .  
 778 4. Shortcut arcs  $(s, b)$  corresponding to null 2-paths from  $s$  to  $b \in (B \setminus B_\emptyset) \setminus S$ . For these, we  
 779 maintain a BCP data structure with  $P = A_\emptyset$ ,  $Q = (B \setminus B_\emptyset) \setminus S$  with weights  $\omega(p) = \pi(s)$  for all  
 780  $p \in P$ , and  $\omega(q) = \pi(q)$  for all  $q \in Q$ . A response  $(a, b)$  corresponds to the null 2-path  $(s, a, b)$ .  
 781 This is only queried while  $s \in S$ .  
 782 5. Shortcut arcs  $(a, t)$  corresponding to null 2-paths from  $a \in (A \setminus A_\emptyset) \cap S$  to  $t$ . For these, we  
 783 maintain a BCP data structure with  $P = (A \setminus A_\emptyset) \cap S$ ,  $Q = B_\emptyset \setminus S$  with weights  $\omega(p) = \pi(p)$  for all  
 784  $p \in P$ , and  $\omega(q) = \pi(t)$  for all  $q \in Q$ . A response  $(a, b)$  corresponds to the null 2-path  $(a, b, t)$ . This is only queried  
 785 while  $t \notin S$ .  
 786 6. Shortcut arcs  $(s, t)$  corresponding to null 3-paths. For these, we maintain in a BCP data  
 787 structure with  $P = A_\emptyset \setminus S$ ,  $Q = B_\emptyset \setminus S$  with weights  $\omega(p) = \pi(s)$  for all  $p \in P$ , and  $\omega(q) = \pi(t)$   
 788 for all  $q \in Q$ . A response  $(a, b)$  corresponds to the null 3-path  $(s, a, b, t)$ . This is only queried  
 789 while  $s \in S$  and  $t \notin S$ .

790 Given  $v' \in \tilde{S}$ , we would like to query:

- 791 1. Non-shortcut backward arcs  $(v', w')$  with  $(w', v') \in \text{supp}(f)$ . For these, we can maintain a  
 792 min-heap on  $(w', v') \in \text{supp}(f)$  arcs for each normal  $v' \in V$ .  
 793 2. Non-shortcut  $A$ -to- $B$  forward arcs. For these, we maintain a NN data structure over  $P =$   
 794  $(B \setminus B_\emptyset) \setminus \tilde{S}$ , with weights  $\omega(p) = \pi(p)$  for each  $p \in P$ . We subtract  $\pi(v')$  from the NN distance  
 795 to recover the reduced cost of the arc from  $v'$ .

- 796 3. Non-shortcut forward arcs from  $s$ -to- $A$  and from  $B$ -to- $t$ . For  $s$ , we can maintain a min-heap  
 797 on the potentials of  $B \setminus \tilde{S}$ , queried only if  $v' = s$ . For  $B$ -to- $t$  arcs, there is only one arc to check  
 798 if  $v' \in B$ , which we can examine manually.
- 799 4. Shortcut arcs  $(s, b)$  corresponding to null 2-paths from  $s$  to  $b \in (B \setminus B_\emptyset) \setminus S$ . For these, we  
 800 maintain a BCP data structure with  $P = A_\emptyset$ ,  $Q = (B \setminus B_\emptyset) \setminus S$  with weights  $\omega(p) = \pi(s)$  for all  
 801  $p \in P$ , and  $\omega(q) = \pi(q)$  for all  $q \in Q$ . A response  $(a, b)$  corresponds to the null 2-path  $(s, a, b)$ .  
 802 This is only queried if  $v' = s$ .
- 803 5. Shortcut arcs  $(a, t)$  corresponding to null 2-paths from  $a \in (A \setminus A_\emptyset) \cap S$  to  $t$ . For these, we  
 804 maintain a NN data structure over  $P = B_\emptyset \setminus S$  with weights  $\omega(p) = \pi(t)$  for each  $p \in P$ .  
 805 A response  $(v', b)$  corresponds to the null 2-path  $(v', b, t)$ . We subtract  $\pi(v')$  from the NN  
 806 distance to recover the reduced cost of the arc from  $v'$ . This is not queried if  $t \in \tilde{S}$ .
- 807 6. Shortcut arcs  $(s, t)$  corresponding to null 3-paths. For these, we maintain in a BCP data  
 808 structure with  $P = A_\emptyset \setminus S$ ,  $Q = B_\emptyset \setminus S$  with weights  $\omega(p) = \pi(s)$  for all  $p \in P$ , and  $\omega(q) = \pi(t)$   
 809 for all  $q \in Q$ . A response  $(a, b)$  corresponds to the null 3-path  $(s, a, b, t)$ . This is only queried  
 810 while  $v' = s$  and  $t \notin S$ .

811 ▶ **Lemma B.1.** *Hungarian search performs  $O(k)$  relaxations before a deficit vertex is reached.*

812 **Proof.** **«TO BE REWRITTEN.»** First we prove that there are  $O(k)$  non-shortcut relaxations.  
 813 Each edge relaxation adds a new vertex to  $S$ , and non-shortcut relaxations only add normal  
 814 vertices. The vertices of  $V \setminus S$  fall into several categories: (i)  $s$  or  $t$ , (ii) vertices of  $A$  or  $B$  with  
 815 0 imbalance, and (iii) deficit vertices of  $A$  or  $B$  ( $S$  contains all excess vertices). The number of  
 816 vertices in (i) and (iii) is  $O(k)$ , leaving us to bound the number of (ii) vertices.

817 An  $A$  or  $B$  vertex with 0 imbalance must have an even number of  $\text{supp}(f)$  edges. There is  
 818 either only one positive-capacity incoming arc (for  $A$ ) or outgoing arc (for  $B$ ), so this quantity  
 819 is either 0 or 2. Since the vertex is normal, this must be 2. We charge 0.5 to each of the two  
 820  $\text{supp}(f)$  arcs; the arcs of  $\text{supp}(f)$  have no more than 1 charge each. Thus, the number of type  
 821 (ii) vertex relaxations is  $O(|\text{supp}(f)|)$ . By Corollary 3.6,  $O(|\text{supp}(f)|) = O(k)$ .

822 Next we prove that there are  $O(k)$  shortcut relaxations. Recall the categories of shortcuts from  
 823 the list of datastructures above. We have shortcuts corresponding to (i) null 2-paths surrounding  
 824  $a \in A_\emptyset$ , (ii) null 2-paths surrounding  $b \in B_\emptyset$ , and (iii) null 3-paths, which go from  $s$  to  $t$ . There is  
 825 only one relaxation of type (iii), since  $t$  can only be added to  $S$  once. The same argument holds  
 826 for type (ii).

827 Each type (i) relaxation adds some normal  $b \in B \setminus B_\emptyset$  into  $S$ . Since  $b$  is normal, it must either  
 828 have deficit or an adjacent arc of  $\text{supp}(f)$ . We charge this relaxation to  $b$  if it is deficit, or the  
 829 adjacent arc of  $\text{supp}(f)$  otherwise. No vertex is charged more than once, and no  $\text{supp}(f)$  edge is  
 830 charged more than twice, therefore the total number of type (i) relaxations is  $O(|\text{supp}(f)|)$ . By  
 831 Corollary 3.6,  $O(|\text{supp}(f)|) = O(k)$ . ◀

832 Similarly we can prove that there are  $O(k)$  relaxations during the DFS.

833 ▶ **Corollary B.2.** *Depth-first search performs  $O(k)$  relaxations before a deficit vertex is reached.*

834 ▶ **Lemma B.3.** *After  $O(n \text{polylog } n)$ -time preprocessing, each Hungarian search can be implemented  
 835 in  $O(k \text{polylog } n)$  time.*

836 **Proof.** Each of the constant number of data structures used by the Hungarian search can be  
 837 constructed in  $O(n \text{polylog } n)$  time. For each data structure queried during a relaxation, the new  
 838 vertex moved into  $S$  causes a constant number of updates, each of which can be implemented  
 839 in  $O(\text{polylog } n)$  time. We first prove that the number of BCP operations during the Hungarian  
 840 search over is bounded by  $O(k)$ .

- 841 1. Let  $S^t$  denote the initial set  $S$  at the beginning of the  $t$ -th Hungarian search, Assume for now  
 842 that, at the beginning of the  $(t+1)$ -th Hungarian search, we have the set  $S^t$  from the previous  
 843 iteration. To construct  $S^{t+1}$ , we remove the vertices that had excess decreased to zero by the  
 844  $t$ -th blocking flow. Thus, we are able to initialize  $S$  at the cost of one BCP deletion per excess  
 845 vertex, which sums to  $O(k)$  over the entire course of REFINE. **Too strong as a bound? Is  
 846 it enough to look at one Hungarian search?**
- 847 2. During each Hungarian search, a vertex entering  $S$  may incur one BCP insertion/deletion.  
 848 We can charge the updates to the number of relaxations over the course of Hungarian search.  
 849 The number of relaxations in a Hungarian search is  $O(k)$  by Lemma B.1.
- 850 3. To obtain  $S^t$ , we keep track of the points added to  $S^t$  since the last Hungarian search. After  
 851 the augmentation, we remove those points added to  $S^t$ . By (2) there are  $O(k)$  such points to  
 852 be deleted, so reconstructing  $S^t$  takes  $O(k)$  BCP operations.

853 For potential updates, we use the same trick by Vaidya [Vai89] to lazily update potentials after  
 854 vertices leave  $S$  (similar to Lemma ??), but this time only for normal vertices. Normal vertices  
 855 are stored in each data structure with weight  $\omega(v) = \pi(v) - \delta$ , and  $\delta$  is increased in lieu of  
 856 increasing the potential of vertices in  $S$ . When a vertex leaves  $S$  (through the rewind mechanism  
 857 above), we restore its potential as  $\pi(v) \leftarrow \omega(v) + \delta$ . With lazy updates, the number of potential  
 858 updates on normal vertices is bounded by the number of relaxations in the Hungarian search,  
 859 which is  $O(k)$  by Lemma B.1. Note that null vertex potentials are not handled in the Hungarian  
 860 search. **then where? Lemma 4.5** ◀

861 There are no potentials to update within DFS, so the running time of DFS boils down to the  
 862 time spent to querying and updating the data structures.

863 ▶ **Lemma B.4.** *After  $O(n \text{polylog } n)$ -time preprocessing, each depth-first search can be implemented  
 864 in  $O(k \text{polylog } n)$  time.*

865 **Proof.** At the beginning of REFINE, we can initialize the  $O(1)$  data structures used in DFS in  
 866  $O(n \text{polylog } n)$  time. We use the same rewinding mechanism as in Hungarian search (Lemma B.3)  
 867 to avoid reconstructing the data structures across iterations of REFINE, so the total time spent  
 868 is bounded by the  $O(\text{polylog } n)$  times the number of relaxations. By Corollary B.2, the running  
 869 time for depth-first search is  $O(k \text{polylog } n)$ . ◀

870 ▶ **Lemma 4.4.** *The support of each blocking flow found in REFINE is of size  $O(k)$ .*

871 **Proof.** Let  $i$  be fixed and consider the invocation of DFS which produces the  $i$ -th blocking flow  
 872  $f_i$ . DFS constructs  $f_i$  as a sequence of admissible excess-deficit paths, which appear as path  $P$   
 873 in Algorithm ?? Every arc in  $P$  is an arc relaxed by DFS, so  $N_i$  is bounded by the number of  
 874 relaxations performed in DFS. Using Corollary B.2, we have  $N_i = O(k)$ . ◀

875 ▶ **Lemma 4.5.** *The number of end-of-REFINE null vertex potential updates is  $O(n)$ . The number of  
 876 augmentation-induced null vertex potential updates in each invocation of REFINE is  $O(k \log k)$ .*

877 **Proof.** The number of end-of-REFINE potential updates is  $O(n)$ . Each update due to flow aug-  
 878mentation involves a blocking flow sending positive flow through a null path, causing a potential  
 879 update on the passed null vertex. We charge this potential update to the edges of that null path,  
 880 which are in turn arcs with positive flow in the blocking flow. For each blocking flow, no positive  
 881 arc is charged more than twice. It follows that the number of augmentation-induced updates  
 882 is at most the size of support of the blocking flow, which is  $O(k)$  by Lemma 4.4. According to  
 883 Lemma 3.7 there are  $O(\sqrt{k})$  iterations of REFINE before it terminates. Summing up we have an  
 884  $O(k\sqrt{k})$  bound over the course of REFINE. ◀

## 885 C Proofs from Section 5

886 **Recovering the optimal flow.** *«use this one if we want to use exponent > 1. »*887 *«Move everything to appendix and left a pointer to the socg 2016 paper.»*

888 We use a strategy from Agarwal *et al.* [AFP<sup>+</sup>17]. Instead of finding a max flow in the entire  
 889 admissible network under  $\pi^*$ , we claim that it is sufficient to find a max flow in a *spanning tree* of  
 890 admissible arcs, e.g. a shortest path tree on reduced costs. There are some details to explain —  
 891 like where the tree should be rooted, how to ensure the underlying network is strongly connected  
 892 by admissible arcs — but we give the intuition first: Such a spanning tree is a maximal set of  
 893 linearly independent dual LP constraints for the optimal dual ( $\pi^*$ ), so there exists an optimal  
 894 primal solution ( $f^*$ ) with support only on these arcs. To see this, we can use a perturbation  
 895 argument: raising the cost of each non-tree edge by  $\varepsilon > 0$  does not change  $\text{cost}(\pi^*)$  or the  
 896 feasibility of  $\pi^*$ , but does raise the cost of any circulation  $f$  using non-tree edges. Strong duality  
 897 suggests that  $\text{cost}(f^*) = \text{cost}(\pi^*)$  is unchanged, therefore  $f^*$  must have support only on the tree  
 898 edges.

899 *«TODO: the SPT construction requires describing Dijkstra and promising strong  
 900 connectivity»*

901 **Recovering the optimal flow.** *«use this one if we only use exponent 1. »*

902 Instead of running a generic max-flow algorithm after finding the optimal potentials, we use  
 903 the following observation.

904 Up until now, we have not placed restrictions on coincidence between  $A$  and  $B$ , but for the  
 905 next proof it is useful to do so. We can assume that all points within  $A \cup B$  are distinct, otherwise  
 906 we can replace all points coincident at  $x \in \mathbb{R}^2$  with a single point whose supply/demand is  
 907  $\sum_{v \in A \cup B : v=x} \lambda(v)$ . This is roughly equivalent to transporting as much as we can between coincident  
 908 supply and demand, and is optimal by triangle inequality. So without loss of generality, we assume  
 909 all points of  $A \cup B$  are distinct.

910 Without loss of generality, assume  $\pi^*$  is nonnegative (raising  $\pi^*$  uniformly on all points does  
 911 not change the objective or feasibility). Recall that feasibility of  $\pi^*$  states that, for all  $a \in A$  and  
 912  $b \in B$

$$c_{\pi^*}(a, b) = \|a - b\|_p - \pi^*(a) + \pi^*(b) \geq 0.$$

914 An arc  $a \rightarrow b$  is admissible when

$$c_{\pi^*}(a, b) = \|a - b\|_p - \pi^*(a) + \pi^*(b) = 0.$$

916 We note that these definitions have a nice visual: Place disks  $D_q$  of radius  $\pi(q)$  at each  $q \in A \cup B$ .  
 917 Feasibility states that for all  $a \in A$  and  $b \in B$ ,  $D_a$  cannot contain  $D_b$  with a gap between their  
 918 boundaries. The arc  $a \rightarrow b$  is admissible when  $D_a$  contains  $D_b$  and their boundaries are tangent.

919 ▶ **Lemma C.1.** *Let  $\pi^*$  be a set of optimal potentials for the point sets  $A$  and  $B$ , under costs  
 920  $c(a, b) = \|a - b\|_p$ . Then, the set of admissible arcs under  $\pi^*$  form a planar graph.*

921 **Proof.** We assume the points of  $A \cup B$  are in general position (e.g. by symbolic perturbation)  
 922 such that no three points are collinear. Let  $a_1 \rightarrow b_1$  and  $a_2 \rightarrow b_2$  be any pair of admissible arcs under  
 923  $\pi^*$ . We will isolate them from the rest of the points, considering  $\pi^*$  restricted to the four points  
 924  $\{a_1, a_2, b_1, b_2\}$ . Clearly, this does not change whether the two arcs cross. Observe that we can  
 925 raise  $\pi^*(a_2)$  and  $\pi^*(b_2)$  uniformly, until  $c_{\pi^*}(a_2, b_1) = 0$ , without breaking feasibility or changing  
 926 admissibility of  $a_1 \rightarrow b_1$  and  $a_2 \rightarrow b_2$ . Henceforth, we assume that we have modified  $\pi^*$  in this way  
 927 to make  $a_2 \rightarrow b_1$  admissible. Given positions of  $a_1, a_2$ , and  $b_1$ , we now try to place  $b_2$  such that

928  $a_1 \rightarrow b_1$  and  $a_2 \rightarrow b_2$  cross. Specifically,  $b_2$  must be placed within a region  $\mathcal{F}$  that lies between the  
 929 rays  $\overrightarrow{a_2 a_1}$  and  $\overrightarrow{a_2 b_1}$ , and within the halfplane bounded by  $\overleftrightarrow{a_1 b_1}$  that does not contain  $a_2$ .

930 Let  $g_a(q) := \|a - q\| - \pi^*(a)$  for  $a \in A$  and  $q \in \mathbb{R}^2$ . Let the **bisector** between  $a_1$  and  $a_2$  be  
 931  $\beta := \{q \in \mathbb{R}^2 \mid g_{a_1}(q) = g_{a_2}(q)\}$ .  $\beta$  is a curve subdividing the plane into two open faces, one where  
 932  $g_{a_1}$  is minimized and the other where  $g_{a_2}$  is. From these definitions, admissibility of  $a_1 \rightarrow b_1$  and  
 933  $a_2 \rightarrow b_1$  imply that  $b_1$  is a point of the bisector.

934 We show that  $\mathcal{F}$  lies entirely on the  $g_{a_1}$  side of the bisector. First, we prove that the closed  
 935 segment  $\overline{a_1 b_1}$  lies entirely on the  $g_{a_1}$  side, except  $b_1$  which lies on  $\beta$ . Any  $q \in \overline{a_1 b_1}$  can be written  
 936 parametrically as  $q(t) = (1-t)b_1 + t a_1$  for  $t \in [0, 1]$ . Consider the single-variable functions  
 937  $g_{a_1}(q(t))$  and  $g_{a_2}(q(t))$ .

$$\begin{aligned} g_{a_1}(q(t)) &= (1-t)\|a_1 - b_1\| - \pi(a_1) \\ g_{a_2}(q(t)) &= \|(a_2 - b_1) - t(a_1 - b_1)\| - \pi(a_2) \end{aligned}$$

939 At  $t = 0$ , these expressions are equal. Observe that the derivative with respect to  $t$  of  $g_{a_1}(q(t))$   
 940 is less than  $g_{a_2}(q(t))$ . Indeed, the value of  $\frac{d}{dt}\|(a_2 - b_1) - t(a_1 - b_1)\|$  is at least  $-\|a_1 - b_1\| =$   
 941  $\frac{d}{dt}g_{a_1}(q(t))$ , which is realized iff  $\frac{(a_2 - b_1)}{\|a_2 - b_1\|} = \frac{(a_1 - b_1)}{\|a_1 - b_1\|}$ . This corresponds to  $\overrightarrow{a_2 b_1}$  and  $\overrightarrow{a_1 b_1}$  being  
 942 parallel, but this is disallowed since  $a_1, a_2, b_1$  are in general position. Thus,  $g_{a_1}(q(t)) \leq g_{a_2}(q(t))$   
 943 with equality only at  $b_1$ .

944 Now, we parameterize each point of  $\mathcal{F}$  in terms of points on  $\overline{a_1 b_1}$ . Every  $q \in \mathcal{F}$  can be written  
 945 as  $q(t') = q' + t'(\overrightarrow{q' - a_2})$  for some  $q' \in \overline{a_1 b_1}$  and  $t' \geq 0$ , i.e.  $q' = \overrightarrow{a_1 b_1} \cap \overrightarrow{a_2 q}$ . We call  $q'$  the  
 946 **projection** of  $q$  onto  $\overline{a_1 b_1}$ . We can write  $g_{a_1}$  and  $g_{a_2}$  in terms of  $t'$  and observe that  $\frac{d}{dt'}g_{a_1}(q(t')) \leq$   
 947  $\frac{d}{dt'}g_{a_2}(q(t'))$ , as the derivative of  $g_a(q(t'))$  is maximized if  $(q(t') - a)$  is parallel to  $(q(t') - a_2)$   
 948 and lower otherwise. Notably,  $q(t')$  with projection  $b_1$  have  $\frac{d}{dt'}g_{a_1}(q(t')) < \frac{d}{dt'}g_{a_2}(q(t'))$ , since  
 949  $a_1, a_2, b_1$  are in general position. Any  $q(t')$  with a different projection do not have strict inequality,  
 950 but the projection itself has  $g_{a_1}(q') < g_{a_2}(q')$  for  $q' \neq b_1$  since it lies on  $\overline{a_1 b_1}$ . Therefore, for all  
 951  $q \in \mathcal{F} \setminus \{b_1\}$ ,  $g_{a_1}(q') < g_{a_2}(q')$ , and  $\mathcal{F}$  lies on the  $g_{a_1}$  side of the bisector except for  $b_1$  which lies  
 952 on  $\beta$ . We can eliminate  $b_1$  as a candidate position for  $b_2$ , since points of  $B$  cannot coincide.

953 Observe that  $g_{a_1}(b) < g_{a_2}(b)$  for  $b \in B$  implies that  $c_\pi(a_1, b) < c_\pi(a_2, b)$ , and  $c_\pi(a_1, b) =$   
 954  $c_\pi(a_2, b)$  if and only if  $b$  lies on  $\beta$ . This holds for all  $b \in \mathcal{F}$  including our prospective  $b_2$ , but then  
 955  $c_\pi(a_1, b_2) < c_\pi(a_2, b_2) = 0$  since  $a_2 \rightarrow b_2$  is admissible. This violates feasibility of  $a_1 \rightarrow b_2$ , so there  
 956 is no feasible placement of  $b_2$  which also crosses  $a_1 \rightarrow b_1$  with  $a_2 \rightarrow b_2$ .  $\blacktriangleleft$

957 We can construct the entire set of admissible arcs by repeatedly querying the minimum-  
 958 reduced cost outgoing arc for each  $a \in A$  until the result is not admissible. By Lemma C.1 the  
 959 resulting arc set forms a planar graph, so by Euler's formula the number of arcs to query is  
 960  $O(n)$ . We can then find the maximum flow in time  $O(n \log n)$  time, using for example the planar  
 961 maximum-flow algorithm by Erickson [Eri10]. [?]

962 By prioritizing the relaxation of support arcs, we also have the following lemma.

963 **Lemma C.2 (Agarwal et al. [AFP<sup>+</sup>17]).** *If arcs of  $\text{supp}(f)$  are relaxed first as they arrive on  
 964 the frontier, then  $E(\text{supp}(f))$  is acyclic.*

965 **Proof.** Let  $f_i$  be the pseudoflow after the  $i$ -th augmentation, and let  $T_i$  be the forest of relaxed  
 966 arcs generated by the Hungarian search for the  $i$ -th augmentation. Namely, the  $i$ -th augmenting  
 967 path is an excess-deficit path in  $T_i$ , and all arcs of  $T_i$  are admissible by the time the augmentation  
 968 is performed. Let  $E(T_i)$  be the undirected edges corresponding to arcs of  $T_i$ . Notice that,  
 969  $E(\text{supp}(f_{i+1})) \subseteq E(\text{supp}(f_i)) \cup E(T_i)$ . We prove that  $E(\text{supp}(f_i)) \cup E(T_i)$  is acyclic by induction  
 970 on  $i$ ; as  $E(\text{supp}(f_{i+1}))$  is a subset of these edges, it must also be acyclic. At the beginning with  
 971  $f_0 = 0$ ,  $E(\text{supp}(f_0))$  is vacuously acyclic.

Let  $E(\text{supp}(f_i))$  be acyclic by induction hypothesis. Since  $T_i$  is a forest (thus, acyclic), any hypothetical cycle  $\Gamma$  that forms in  $E(\text{supp}(f_i)) \cup E(T_i)$  must contain edges from both  $E(\text{supp}(f_i))$  and  $E(T_i)$ . To give a visual analogy, we will color  $e \in \Gamma$  **purple** if  $e \in E(\text{supp}(f_i)) \cap E(T_i)$ , **red** if  $e \in E(\text{supp}(f_i))$  but  $e \notin E(T_i)$ , and **blue** if  $e \in E(T_i)$  but  $e \notin E(\text{supp}(f_i))$ . Then,  $\Gamma$  is neither entirely red nor entirely blue. We say that red and purple edges are **red-tinted**, and similarly blue and purple edges are **blue-tinted**. Roughly speaking, our implementation of the Hungarian search prioritizes relaxing red-tinted admissible arcs over pure blue arcs.

We can sort the blue-tinted edges of  $\Gamma$  by the order they were relaxed into  $S$  during the Hungarian search forming  $T_i$ . Let  $(v, w) \in \Gamma$  be the last pure blue edge relaxed, of all the blue-tinted edges in  $\Gamma$  — after  $(v, w)$  is relaxed, the remaining unrelaxed, blue-tinted edges of  $\Gamma$  are purple.

Let us pause the Hungarian search the moment before  $(v, w)$  is relaxed. At this point,  $v \in S$  and  $w \notin S$ , and the Hungarian search must have finished relaxing all frontier support arcs. By our choice of  $(v, w)$ ,  $\Gamma \setminus (v, w)$  is a path of relaxed blue edges and red-tinted edges which connect  $v$  and  $w$ . Walking around  $\Gamma \setminus (v, w)$  from  $v$  to  $w$ , we see that every vertex of the cycle must be in  $S$  already:  $v \in S$ , relaxed blue edges have both endpoints in  $S$ , and any unrelaxed red-tinted edge must have both endpoints in  $S$ , since the Hungarian search would have prioritized relaxing the red-tinted edges to grow  $S$  before relaxing  $(v, w)$  (a blue edge). It follows that  $w \in S$  already, a contradiction.

No such cycle  $\Gamma$  can exist, thus  $E(\text{supp}(f_i)) \cup E(T_i)$  is acyclic and  $E(\text{supp}(f_{i+1})) \subseteq E(\text{supp}(f_i)) \cup E(T_i)$  is acyclic. By induction,  $E(\text{supp}(f_i))$  is acyclic for all  $i$ .  $\blacktriangleleft$

Let  $E(\Sigma_a)$  **«only used once»** be the underlying edges of the support star centered at  $a$  and  $F := E(\text{supp}(f)) \setminus \bigcup_{a \in A} E(\Sigma_a)$ . Using Lemma C.2, we can show that the number of support arcs outside support stars ( $|F|$ ) is small.

► **Lemma C.3.**  $|B_\ell \setminus \bigcup_{a \in A} \Sigma_a| \leq r$ .

**Proof.**  $F$  is constructed from  $E(\text{supp}(f))$  by eliminating edges in support stars, therefore all edges in  $F$  must adjoin vertices in  $B$  of support degree at least 2. By Lemma C.2,  $E(\text{supp}(f))$  is acyclic and therefore forms a spanning forest over  $A \cup B_\ell$ , so  $F$  is also a bipartite forest. All leaves of  $F$  are therefore vertices of  $A$ .

Pick an arbitrary root for each connected component of  $F$  to establish parent-child relationships for each edge. As no vertex in  $B$  is a leaf, each vertex in  $B$  has at least one child. Charge each vertex in  $B$  to one of its children in  $F$ , which must belong to  $A$ . Each vertex in  $A$  is charged at most once. Thus, the number of  $B_\ell$  vertices outside of support stars is no more than  $r$ .  $\blacktriangleleft$

► **Lemma 5.2.** Suppose we have stripped the graph of dead vertices. The number of relaxation steps in a Hungarian search outside of support stars is  $O(r)$ .

**Proof.** If there are no dead vertices, then each non-support star relaxation step adds either (i) an active deficit vertex, (ii) a non-deficit vertex  $a \in A_\ell$ , or (iii) a non-deficit vertex  $b \in B_\ell$  of support degree at least 2. There is a single relaxation of type (i), as it terminates the search. The number of vertices of type (ii) is  $r$ , and the number of vertices of type (iii) is at most  $r$  by Lemma C.3. The lemma follows.  $\blacktriangleleft$

► **Lemma 5.3.** Hungarian search takes  $O(r\sqrt{n} \text{ polylog } n)$  time.

**Proof.** The number of relaxation steps outside of support stars is  $O(r)$  by Lemma 5.2. The time per relaxation outside of support stars is  $O(\sqrt{n} \text{ polylog } n)$ . The time spent processing relaxations within a support star is  $O(\sqrt{n} \text{ polylog } n)$ , and at most  $r$  are relaxed during the search. The total time is therefore  $O(r\sqrt{n} \text{ polylog } n)$ .  $\blacktriangleleft$

1017 Initially, we label stars big or small according to the  $\sqrt{n}$  threshold. A star that is currently big  
 1018 is turned into a small star once  $|\Sigma_a| \leq \sqrt{n}/2$ . A star that is currently small is turned into a big star  
 1019 once  $|\Sigma_a| \geq 2\sqrt{n}$ . This way, the time spent rebuilding/updating the respective data structures  
 1020 can be amortized to the insertions/deletions that preceded the switch, plus some  $O(r)$  extra  
 1021 work if the update is small-to-big.

1022 A star  $\Sigma_a$  that is switching from big-to-small has size  $|\Sigma_a| \leq \sqrt{n}/2$ . When switching, we delete  
 1023  $\mathcal{D}_{\text{big}}(a)$  and insert  $\Sigma_a$  into  $\mathcal{D}_{\text{small}}$ . Thus, the time spent for big-to-small update is  $O(\sqrt{n} \text{ polylog } n)$ ,  
 1024 and there were at least  $\sqrt{n}/2$  points removed from  $\Sigma_a$  since it was last big.

1025 A star  $\Sigma_a$  that is switching from small-to-big has size  $|\Sigma_a| = \sqrt{n} + x \geq 2\sqrt{n}$ , for some integer  
 1026  $x \geq \sqrt{n}$ . Rearranging, we have  $|\Sigma_a| \leq 2x$ . When switching, we delete all  $|\Sigma_a|$  points from  
 1027  $\mathcal{D}_{\text{small}}$  and construct a new  $\mathcal{D}_{\text{big}}(a)$ . Constructing  $\mathcal{D}_{\text{big}}(a)$  requires inserting  $O(r)$  points of  $A$  (into  
 1028  $P$ ) and the  $|\Sigma_a|$  points of the star (into  $Q$ ). Thus, the time spent for a small-to-big update is  
 1029  $O((r + x) \text{ polylog } n)$ , and there were at least  $x \geq \sqrt{n}$  points added to  $\Sigma_a$  since it was last small.

---

### 1030 References for the Appendix

---

- 1031 **1AFP<sup>+</sup>17** Pankaj K. Agarwal, Kyle Fox, Debmalya Panigrahi, Kasturi R. Varadarajan, and Allen Xiao. Faster  
 1032 algorithms for the geometric transportation problem. In *33rd International Symposium on Com-*  
 1033 *putational Geometry, SoCG 2017, July 4-7, 2017, Brisbane, Australia*, pages 7:1–7:16, 2017.
- 1034 **1Eri10** Jeff Erickson. Maximum flows and parametric shortest paths in planar graphs. In *Proceedings*  
 1035 *of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2010, Austin,*  
 1036 *Texas, USA, January 17-19, 2010*, pages 794–804, 2010.
- 1037 **1GHKT17** Andrew V. Goldberg, Sagi Hed, Haim Kaplan, and Robert E. Tarjan. Minimum-cost flows in  
 1038 unit-capacity networks. *Theory Comput. Syst.*, 61(4):987–1010, 2017.
- 1039 **1GT90** Andrew V. Goldberg and Robert E. Tarjan. Finding minimum-cost circulations by successive  
 1040 approximation. *Math. Oper. Res.*, 15(3):430–466, 1990.
- 1041 **1SA12** R. Sharathkumar and Pankaj K. Agarwal. Algorithms for the transportation problem in geo-  
 1042 metric settings. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete*  
 1043 *Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, pages 306–317, 2012.
- 1044 **1Vai89** Pravin M. Vaidya. Geometry helps in matching. *SIAM J. Comput.*, 18(6):1201–1225, 1989.