

# Geometric Partial Matchings 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

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

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

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

**Keywords and phrases** partial matching, transportation, minimum-cost flow, rms-distance, bichromatic closest pair, cost scaling, excess scaling, primal-dual

**Lines** 1070

## 18 1 Introduction

### 19 «REWRITE AFTER TECHNICAL SECTIONS»

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



© Pankaj K. Agarwal, Hsien-Chih Chang, Allen Xiao;  
licensed under Creative Commons License CC-BY

The 35th International Symposium on Computational Geometry (SOCG 2019).



Leibniz International Proceedings in Informatics

LIPIC Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



28 *unbalanced*. We call the geometric problem of finding a size  $k$  matching of point sets  $A$  and  
 29  $B$  the *geometric partial matching problem*. **«talk about the near-linear time algorithm»**

## 30 1.1 Contributions

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

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

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

46 **«State the settings separately so no need to repeat in the theorem statement.»**

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

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

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

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

$O(rn^{3/2} \text{polylog } n)$

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

## 2 Minimum-Cost Partial Matchings using Hungarian Algorithm

The Hungarian algorithm [7] is a primal-dual algorithm for min-cost bipartite matching in general graphs that can be adapted to solve the partial matching problem exactly if one terminates the algorithm after  $k$  iterations (see e.g. [9]). In this section, we prove Theorem 1.1 by implementing the Hungarian algorithm in  $O((n + k^2) \text{polylog } n)$  time.

### 2.1 Matching Terminologies

Let  $G$  be a bipartite graph between vertex sets  $A$  and  $B$  and edge set  $E$ , with costs  $c(v, w)$  for each edge  $(v, w)$  in  $G$ . Let  $C := \max_{(v, w) \in E} c(v, w)$ . **«When is this used?»** A *matching*  $M \subseteq E$  is a set of edges where no two edges share an endpoint. A vertex  $v$  is *matched* by  $M$  if  $v$  is the endpoint of some matching edge in  $M$ ; otherwise  $v$  is *unmatched*. 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 of minimum cost. In the geometric partial matching setting, we have  $E = A \times B$  and  $c(a, b) = (\|a - b\|_2)^q$  for every edge  $(a, b)$  in  $G$ . **What is  $q$ ?**

This should be part of problem

The linear program dual to the standard linear program for MPM has dual variables for each vertex, called *potentials*  $\pi$ . Given potentials  $\pi$ , we can define the *reduced cost* on the edges to be  $c_\pi(v, w) := c(v, w) - \pi(v) + \pi(w)$ . Potentials  $\pi$  are *feasible* if the reduced costs are nonnegative for all edges in  $G$ . We say that an edge  $(v, w)$  is *admissible* under potentials  $\pi$  if  $c_\pi(v, w) = 0$ .

Consider a matching  $M$  of size less than  $r$ . An *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 edges outside and inside of matching  $M$ . The symmetric difference  $M \oplus \Pi$  creates a new matching of size  $|M| + 1$ . We say that  $M \oplus \Pi$  is the result of *augmenting*  $M$  by  $\Pi$ .

### 2.2 The Hungarian Algorithm

The Hungarian algorithm is initialized with  $M = \emptyset$  and  $\pi = 0$ . It 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  $\alpha$  satisfying  $\alpha \geq \pi(a)$  for any matched vertex  $a \in A$ , and
- (iv) unmatched vertices of  $B$  all have the same potential  $\beta$  satisfying  $\beta \leq \pi(b)$  for any matched vertex  $b \in B$ .

Ramshaw and Tarjan [9] show that these conditions are sufficient to guarantee that  $M$  is a *minimum-cost matching*. Each iteration the Hungarian algorithm augments  $M$  with an admissible augmenting path  $\Pi$ , discovered using a procedure called the *Hungarian search*. The algorithm terminates once  $M$  has size  $k$ , in which case  $M$  is an optimal matching.

The Hungarian search tries grow a set of *reachable vertices*  $S$  by augmenting paths consisting of admissible edges. Initially,  $S$  is the set of unmatched vertices in  $A$ . Let the *frontier* of  $S$  be the edges in  $(A \cap S) \times (B \setminus S)$ . In each iteration, the Hungarian search first *relaxes* the minimum-reduced-cost edge  $(a, b)$  in the frontier, raising  $\pi(a)$  by  $c_\pi(a, b)$  for all  $a \in S$  to make  $(a, b)$  admissible, and adding  $b$  into  $S$ . It is easy to verify that this potential change preserves feasibility. As  $b \in B$  is added into  $S$ , we can store a backpointer to  $a$ , which can be used later to recover the admissible augmenting path through  $b$ . If  $b$  is matched, say

Indetailed.

to vertex  $a'$ , then we also relax  $(a', b)$  by adding  $a'$  into  $S$  (no potential change needed, by invariant) with backpointer to  $b$ . If  $b$  is unmatched, the search finishes and we can recover an admissible augmenting path to  $b$  by following backpointers to an unmatched vertex  $a \in A$ .

Each 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 runs in  $O(k \text{ polylog } n)$  time after an  $O(n \text{ polylog } n)$  time preprocessing. With this, our implementation of the Hungarian algorithm runs in  $O((n + k^2) \text{ polylog } n)$  time, which proves Theorem 1.1.

## 2.3 Geometric implementation of Hungarian search

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

In general graphs, the most expensive step is to find the minimum-reduced-cost frontier edge — the search must “look at every edge”. However, in the geometric setting, we can achieve this with a query to a dynamic *bichromatic closest pair* (BCP) data structure. Given two point sets  $P$  and  $Q$  in the plane, the bichromatic closest pair are two points  $p \in P$  and  $q \in Q$  minimizing the additively weighted distance  $\|p - q\| - \omega(p) + \omega(q)$  for some real-valued vertex weights  $\omega$ . Thus, the minimum reduced-cost among the frontier edges is precisely the cost of the BCP of point sets  $P = A \cap S$  and  $Q = B \setminus S$ , with  $\omega(p) = \pi(p)$ .

The state of the art dynamic BCP data structure from Kaplan, Mulzer, Roditty, Seifertl, and Sharir [6] supports point insertions and deletions in  $O(\text{polylog } n)$  time, and answers queries in  $O(\log^2 n)$  time. During each relaxation, we perform at most one query and add a vertex to  $S$  incurring one BCP insertion or deletion. Thus the running time for the search is  $O(k \text{ polylog } n)$ .

**Initial BCP sets by rewinding.** Recall that  $S$  is initialized to the set of unmatched vertices of  $V = A$ , and therefore  $Q = B \setminus S$  has size  $n$ . We cannot afford to take  $O(n \text{ polylog } n)$  time to initialize the BCP data structure at the beginning of every Hungarian search beyond the first. However, the set of unmatched  $A$  vertices has changed by exactly one vertex since the last Hungarian search — the augmentation newly matched one vertex  $a^* \in A$ . Thus, given the initial BCP sets  $P', Q'$  from the beginning of the last Hungarian search, we can construct  $P$  and  $Q$  for the current iteration using a single BCP deletion in  $O(\text{polylog } n)$  time.

To acquire  $P'$  and  $Q'$ , we keep track of a list of the points added to  $S$  over the course of the Hungarian search. At the end of each Hungarian search we *rewind* the BCP data structure by tracing the list in reverse order. The number of points in the list is at most  $O(k)$  as it is bounded by the number of relaxations per Hungarian search. Thus, in  $O(k \text{ polylog } n)$  time, we can reconstruct  $P'$  and  $Q'$  for each Hungarian search beyond the first. We refer to this procedure as the *rewinding mechanism*.

**Potential updates by Vaidya's trick.** We modify a trick from Vaidya [11] for batching potential updates. Potentials have a *stored value*, i.e. the currently recorded value of  $\pi(v)$ , and a *true value*, which may have changed from  $\pi(v)$ . The resulting algorithm queries the minimum-reduced-cost under the true values of  $\pi$  and updates the stored value occasionally.

Throughout the entire Hungarian algorithm, we maintain a nonnegative scalar  $\delta$  (initially set to 0) which aggregates potential changes. Vertices  $a \in A$  that are added to  $S$  are inserted into BCP with weight  $\omega(a) \leftarrow \pi(a) - \delta$ , for whatever value  $\delta$  is at the time of insertion.

Vaidya resets the stored value of all potentials to their true value at the end of each Hungarian search, but we cannot afford to do it and so instead we proceed as follows.

It is needed to be relaxed

Need to say under what metric / distance function in  $\mathbb{R}^2$

also observed in [7]

in the beginning of each iteration,

Similarly, vertices  $b \in B$  that are added to  $S$  have  $\omega(b) \leftarrow \pi(b) - \delta$  recorded ( $B \cap S$  points aren't added into a BCP set). When the Hungarian search wants to raise the potentials of points in  $S$ ,  $\delta$  is increased by that amount instead. Thus, true value for any potential of a point in  $S$  is always  $\omega(p) + \delta$ . For points of  $(A \cup B) \setminus S$ , the true potential is equal to the stored potential. Since all the points of  $A \cap S$  have weights uniformly offset from their true potentials, the minimum edge returned by the BCP does not change. **«why?»**

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

We update  $\delta$  once per edge relaxations; thus  $O(k)$  times in total per Hungarian search. There are  $O(k)$  stored values updated per Hungarian search during the rewinding process. The time spent on potential updates per Hungarian search is therefore  $O(k)$ .

**Lemma:** In summary, we can implement each Hungarian search in  $O(k \text{ polylog } n)$  time after a one-time  $O(n \text{ polylog } n)$  preprocessing.

*State of Lemma:*

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

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

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

#### 3.1 Preliminaries on Network Flows

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

**Pseudoflows.** Given a network  $N := (V, \vec{E}, c, u, \phi)$ , a *pseudoflow* (or *flow* to be short)  $f: \vec{E} \rightarrow \mathbb{Z}^1$  on  $N$  is an antisymmetric function on the arcs of  $N$  satisfying  $f(v \rightarrow w) \leq u(v \rightarrow w)$  for every arc  $v \rightarrow w$ . We sometimes abuse the terminology by allowing pseudoflow to be defined on a directed graph, in which case we are actually referring to the pseudoflow on the

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

corresponding network by extending the flow values antisymmetrically to the arcs. We say that  $f$  *saturates* an arc  $v \rightarrow w$  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  $N$ , denoted as  $\text{supp}(f)$ , is the set of arcs with positive flows:

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

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

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

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

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

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

**Residual graph.** Given a pseudoflow  $f$ , one can define the *residual network* as follows. Recall 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 words, an arc that is not saturated by  $f$  is a residual arc; similarly, given an arc  $v \rightarrow w$  with positive flow value, the backward arc  $w \rightarrow v$  is a residual arc.

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

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

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

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

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

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

The *dual feasibility constraint* says that  $c_\pi(v \rightarrow w) \geq 0$  holds for every directed edge  $(v, w)$  in  $E$ ; potentials  $\pi$  which satisfy this constraint are said to be *feasible*. Suppose we relax the dual feasibility constraint to allow some small violation in the value of  $c_\pi(v \rightarrow w)$ . We say that



a pair of pseudoflow  $f$  and potential  $\pi$  is  $\varepsilon$ -optimal [?, ?] if  $c_\pi(v \rightarrow w) \geq -\varepsilon$  for every residual arc  $v \rightarrow w$  in  $\vec{E}_f$ . Pseudoflow  $f$  is  $\varepsilon$ -optimal if it is  $\varepsilon$ -optimal with respect to some potentials  $\pi$ ; potential  $\pi$  is  $\varepsilon$ -optimal if it is  $\varepsilon$ -optimal with respect to some pseudoflow  $f$ . Given a pseudoflow  $f$  and 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  $g$  in  $G_f$  is *admissible* if all support arcs of  $g$  on  $G_f$  are admissible; in other words,  $g(v \rightarrow w) > 0$  holds only on admissible arcs  $v \rightarrow w$ .

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

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

The goal of the subsection is to reduce the minimum-cost partial matching problem to the unit-capacity minimum-cost flow problem with a polynomial bound on diameter. To this end we first provide an upper bound on the size of support of an integral pseudoflow on the standard reduction network between the two problems. This upper bound in turn provides an additive approximation on the cost of an  $\varepsilon$ -optimal circulation. Next we employ a technique by Sharathkumar and Agarwal [10] to transform an additive  $\varepsilon$ -approximate solution into a multiplicative  $(1 + \varepsilon)$ -approximation for the geometric partial matching problem. The reduction does not work out of the box, as Sharathkumar and Agarwal were tackling a similar but different problem on geometric transportations.

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

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

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

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

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

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

What is diameter of  $N_H$ ?

Don't really need that?

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

283 **Multiplicative approximation.** Now we employ a technique from Sharathkumar and Agar-  
 284 wal [10] to convert the additive approximation into a multiplicative one. Here we sketch a  
 285 proof to Lemma 3.2; a complete proof can be found in the Appendix.

286 «Sketch Sharathkumar and Agarwal [10] with our modification.»

### 287 3.3 High-Level Description of Cost-Scaling Algorithm

291 Our main algorithm for the unit-capacity minimum-cost flow problem is based on the *cost-*  
 292 *scaling* technique, originally due to Goldberg and Tarjan [5]; Goldberg, ~~et al.~~ <sup>R</sup> ~~Kaplan, and~~  
 293 ~~Tarjan~~ [4] applied the technique on unit-capacity networks. The algorithm finds  $\varepsilon$ -optimal  
 294 circulations for geometrically shrinking values of  $\varepsilon$ . Each fixed value of  $\varepsilon$  is called a *cost*  
 295 *scale*. Once  $\varepsilon$  is sufficiently small, the  $\varepsilon$ -optimal flow is a suitable approximation according  
 296 to Lemma 3.2.<sup>2</sup>

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

- 300 ■ SCALE-INIT takes the previous circulation (now  $2\varepsilon$ -optimal) and transforms it into an  
 301  $\varepsilon$ -optimal pseudoflow with  $O(k)$  excess.
- 302 ■ REFINE then reduces the excess in the newly constructed pseudoflow to zero, making it  
 303 an  $\varepsilon$ -optimal circulation.

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

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

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

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

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



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

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

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

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

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

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

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

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

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

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

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

## 4 Fast Implementation of Refine?

Both Hungarian search and depth-first search are implemented in a Dijkstra-like fashion, traversing through the residual graph using admissible arcs starting from the excess vertices. Each step of the search procedures *relaxes* a minimum-reduced-cost arc from the set of visited

1:10

## Geometric Partial Matchings and Unbalanced Transportation Problem

cannot cross  
I assume this is not the  
analysis but also an efficient  
implementation. State  
what's the main idea  
behind implementation

vertices to an unvisited vertex, until a deficit vertex is reached. At a high level, our analysis strategy is to charge the relaxation events to the support arcs of  $f$ , which has size at most  $O(k)$  by Corollary 3.6.

### 4.1 Null vertices and shortcut graph

«A figure might be helpful for this section.»

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

**Null vertices.** We say a vertex  $v$  in the residual graph  $N_f$  is a *null vertex* if  $\phi_f(v) = 0$  and no 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. Vertices that are not null are called *normal vertices*. A *null 2-path* is a length-2 subpath in  $N_f$  from a normal vertex to another normal vertex, passing through a null vertex. As every vertex in  $A$  has in-degree 1 and every vertex in  $B$  has out-degree 1 in the residual graph, the null 2-paths 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  $A \setminus A_\emptyset$ . In either case, we say that the null 2-path *passes through* null vertex  $v$ . Similarly, we define the length-3 path from  $s$  to  $t$  that passes through two null vertices to be a *null 3-path*. Because reduced costs telescope for residual paths, the reduced cost of any null 2-path or null 3-path does not depend on the null vertices it passes through.

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

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

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

### 4.2 Dynamic data structures for search procedures

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

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

By construction, the distance returned by each of the BCP data structure is equal to the reduced cost of the shortcut, which is equal to the reduced cost of the corresponding null

What's the size of  
the short-cut graph.  
I assume this is one of  
new ideas.

408 path. Each of the above data structures requires one query per relaxation, and an update  
 409 operation whenever a new vertex moves into  $\tilde{S}$ . The data structures above can perform both  
 410 queries and updates in  $O(\text{polylog } n)$  time each, so the running time of the Hungarian search  
 411 other than the potential updates can be charged to the number of relaxation steps.

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

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

### 419 4.3 Time analysis

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

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

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

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

### 429 4.4 Number of potential updates on null vertices

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

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

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

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

442 Now combining Lemma B.3, Lemma B.4, and Lemma 4.5 completes the proof of The-  
 443 orem 1.2.

Do we need this  
 Lemma?  
 May be combine  
 Section 4.3 and 4.4  
 these two sections  
 have a sequence of  
 lemmas but wider.

## 5 Unbalanced Transportation

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

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

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

Given  $A$ ,  $B$ , and  $\lambda$ , the *transportation problem* asks to find a feasible transportation map of minimum cost. We focus on the case when the point sets are *unbalanced*; that is, the given point sets have different sizes  $r$  and  $n$ ; without loss of generality assuming  $r \leq n$ .

There is a simple reduction from the transportation problem to the uncapacitated min-cost flow problem. Consider the complete bipartite graph  $G$  between  $A$  and  $B$  (with all edges directed from  $A$  to  $B$ ). Set the costs  $c(a, b)$  to be  $\|a - b\|_p^q$ , all capacities  $u(a, b)$  to infinity, and the supply-demand function  $\phi = \lambda$ . Any circulation  $f$  in the network  $N = (G, c, u, \phi)$  can be converted into a feasible transportation map  $\tau_f$  by taking  $\tau_f(a, b) := f(a \rightarrow b)$  for every edges  $(a, b)$ . One simply has  $\text{cost}(f) = \text{cost}(\tau_f)$ .

### 5.1 Uncapacitated MCF by excess scaling

We give an outline of the strongly polynomial algorithm for uncapacitated min-cost flow problem from Orlin [8]. Orlin's algorithm follows an *excess-scaling* paradigm originally due to Edmonds and Karp [2]. Consider the basic primal-dual framework used in the previous sections: The algorithm begins with both flow  $f$  and potentials  $\pi$  set to zero. Repeatedly runs a *Hungarian search* that raises potentials (while maintaining dual feasibility) to create an admissible augmenting excess-deficit path, on which we perform flow augmentations. In terms of cost,  $f$  is maintained to be 0-optimal with respect to  $\pi$  and each augmentation over admissible edges preserve such property by Lemma 3.1. Thus, the final circulation must be optimal. The excess-scaling paradigm builds on top of this skeleton by specifying (i) between which excess and deficit vertices we send flows, and (ii) how much flow is sent by the augmentation.

The excess-scaling algorithm maintains a *scale parameter*  $\Delta$ , initially set to  $U$ . A vertex  $v$  with  $|\phi_f(v)| \geq \Delta$  is called *active*. Each augmenting path is chosen between an active excess vertex and an active deficit vertex. Once there are no more active excess or 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 a strongly polynomial bound on the number of augmentations and excess scales. First, an *active* vertex

This should be part of intro.

is redefined to be one satisfying  $|\phi_f(v)| \geq \alpha\Delta$ , for a fixed parameter  $\alpha \in (0.5, 1)$ . Second, arcs with flow value at least  $3n\Delta$  at the beginning of a scale are *contracted* to create a new vertex, whose supply-demand is the sum of those on the two endpoints of the contracted arc. 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, in addition to halving,  $\Delta$  is aggressively lowered to  $\max_{v \in V} \phi_f(v)$  if there are no active excess vertices and  $f(v \rightarrow w) = 0$  holds for every arc  $v \rightarrow w \in \hat{E}$ . 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  $G$  by sequentially undoing the contractions. The algorithm then 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 a set of 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  $O(r\sqrt{n} \text{polylog } n)$  time (after preprocessing). Additionally, we show that  $f^*$  can be recovered from  $\pi^*$  very quickly in our setting.

**Implementing contractions.** Following Agarwal *et al.* [1], our geometric data structures must deal with real points in the plane instead of the the contracted components. 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 are contracted. We maintain potentials on the points  $A$  and  $B$  directly, instead of the contracted components.

When conducting the Hungarian search, we initialize  $S$  to be the set of vertices from *active excess contracted components* who (in sum) meet the imbalance criteria. **«unclear»** Upon relaxing any  $v \in \hat{v}$ , we immediately relax all the contracted support arcs which span  $\hat{v}$ . Since the input network is uncapacitated, each contracted component is strongly connected in the residual network by the admissible forward/backward arcs of each contracted arc. **«unparsable»** To relax arcs in  $\hat{E}$ , we relax the support arcs before attempting to relax any non-support arcs. **«mention the reason to make support acyclic»** Relaxations of support arcs can be performed without further potential changes, since they are admissible by invariant.

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

**Recovering optimal flow.** Rewinding contracted components to recover an optimal flow naïvely takes  $O(rn^2)$  time. Use a strategy from Agarwal *et al.* [1], we can recover the optimal flow in time  $O(n \text{polylog } n)$ . If furthermore the cost function is just the  $p$ -norm (without the  $q$ th-power), an even stronger result stands: In this case, the set of admissible arcs under an optimal potential forms a planar graph, and thus we can apply the planar maximum-flow algorithm [?, 3] which runs in  $O(n \log n)$  time. For details see the appendix.

## 5.2 Dead vertices and support stars

Our goal is to implement each augmentation in  $O(r\sqrt{n} \text{polylog } n)$  time. To find an augmenting path, we again use a Hungarian search with geometric data structures to perform relaxations

quickly. Like in Section 3, there are vertices which cannot be charged to the flow support. Even worse, the flow support for the transportation problem may have size  $\Omega(n)$  (consider when  $A$  has one point, and demands are uniformly distributed among the vertices of  $B$ ). Our strategy is summarized as follows:

- Discard vertices which lead to dead ends in the search (not on a path to a deficit vertex).
- 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.

Querying/updating clusters degrades our amortized time per relaxation from  $O(\text{polylog } n)$  to  $O(\sqrt{n} \text{ polylog } n)$ . Thus overall each augmentation takes  $O(r\sqrt{n} \text{ polylog } n)$  time.

**Dead vertices.** Let the *support degree* of a vertex be its degree in the graph induced by the underlying edges of  $\text{supp}(f)$ . We call a vertex  $b \in B$  *dead* if  $b$  has support degree 0 and is not an active excess or deficit vertex; call it *live* otherwise. Dead vertices are essentially equivalent to the *null vertices* of Section 3. However, since the reduction in this section does not use a super-source/super-sink, we can simply remove these from consideration during a Hungarian search — they will not terminate the search, and have no outgoing residual arcs. Like the null vertices, we ignore dead vertex potentials and infer feasible potentials when they become live again. We use  $A_\ell$  and  $B_\ell$  to denote the living vertices of points in  $A$  and  $B$ , respectively. Note that being dead/alive is a notion strictly defined only for vertices, and not for contracted components.

We say a dead vertex is *revived* when it stops meeting either condition of the definition. Dead vertices are only revived after  $\Delta$  decreases (at the start of a subsequent excess scale) as no augmenting path will cross a dead vertex and they cannot meet the criteria for contractions. When a dead vertex is revived, we must add it back into each of our data structures and give it a feasible 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 we can query by maintaining a weighted nearest neighbor data structure on the points of  $A$ . The total number of revivals is bounded above by the number of augmentations: since the final flow is a circulation on  $\hat{G}$  and a newly revived vertex  $v$  has no incident arcs in  $\text{supp}(f)$  and cannot be contracted, there is at least one subsequent augmentation which uses  $v$  as its beginning or end. Thus, the total number of revivals is  $O(n \log n)$ .

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

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

### 5.3 Implementation details

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



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

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

580 **Relaxations outside support stars.** For relaxations that don't involve support star vertices,  
 581 we can once again maintain a BCP to query the minimum  $A_\ell$ -to- $B_\ell$  arc. To elaborate, this is  
 582 the 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  
 583 can be queried in  $O(\log n)$  time and updated in  $O(\text{polylog } n)$  time per point. Since it doesn't  
 584 deal with support stars, there is at most one insertion/deletion per relaxation step.

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

587 **Relaxing a support star.** We classify support stars into two categories: *big stars* are those  
 588 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  
 589 of big stars and  $A_{\text{small}} \subseteq A$  denote the centers of small stars. We keep the following  
 590 data structures to manage support stars.

- 591 1. For each big star  $\Sigma_a$ , we use a data structure  $\mathcal{D}_{\text{big}}(a)$  to maintain BCP between  $P = A_\ell \cap S$   
 592 and  $Q = \Sigma_a$ , weighted by potentials. We query this until  $a \in S$  or any vertex of  $\Sigma_a$  is  
 593 added to  $S$ .
- 594 2. All small stars are added to a single BCP data structure  $\mathcal{D}_{\text{small}}$  between  $P = A_\ell \cap S$  and  
 595  $Q = (\bigcup_{a \in A_{\text{small}}} \Sigma_a) \setminus S$ , weighted by potentials. When an  $a \in A_{\text{small}}$  or any vertex of its  
 596 support star is added to  $S$ , we remove the points of  $\Sigma_a$  from  $\mathcal{D}_{\text{small}}$  using  $|\Sigma_a|$  deletion  
 597 operations.

598 We will update these data structures as each support star center is added into  $S$ . If a  
 599 relaxation step adds some  $b \in B_\ell$  and  $b$  is in a support star  $\Sigma_a$ , then we immediately relax  
 600  $b \rightarrow a$ , as all support arcs are admissible. Relaxations of non-support star  $b \in B_\ell$  will not  
 601 affect the support star data structures.

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

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

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

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

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

**Preprocessing time.** To build the very first set of data structures, we take  $O(rn \text{ polylog } n)$  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 is  $O(rn)$ .  $\mathcal{D}_{\text{small}}$  also has at most  $O(rn)$  points. Each BCP data structure can be constructed in  $O(\text{polylog } n)$  times its size, so the total preprocessing time is  $O(rn \text{ polylog } n)$ .

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

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

**Putting it together.** After  $O(rn \text{ polylog } n)$  preprocessing, we spend  $O(r\sqrt{n} \text{ polylog } n)$  time 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.**

Thank Haim Kaplan for suggesting to use [4] for useful discussion.

## 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>).
- 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>).
- 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*,

- 663     Austin, Texas, USA, January 17-19, 2010, 794–804, 2010. <https://doi.org/10.1137/1.9781611973075.65>.
- 664
- 665     **4**     Andrew V. Goldberg, Sagi Hed, Haim Kaplan, and Robert E. Tarjan. Minimum-cost  
666     flows in unit-capacity networks. *Theory Comput. Syst.* 61(4):987–1010, 2017. <https://doi.org/10.1007/s00224-017-9776-7>.
- 667
- 668     **5**     Andrew V. Goldberg and Robert E. Tarjan. Finding minimum-cost circulations by success-  
669     ive approximation. *Math. Oper. Res.* 15(3):430–466, 1990. [https://doi.org/10.1287/](https://doi.org/10.1287/moor.15.3.430)  
670     [moor.15.3.430](https://doi.org/10.1287/moor.15.3.430).
- 671     **6**     Haim Kaplan, Wolfgang Mulzer, Liam Roditty, Paul Seiferth, and Micha Sharir. Dynamic  
672     planar voronoi diagrams for general distance functions and their algorithmic applications.  
673     *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms,*  
674     *SODA 2017, Barcelona, Spain, Hotel Porta Fira, January 16-19*, 2495–2504, 2017. <https://doi.org/10.1137/1.9781611974782.165>.
- 675
- 676     **7**     Harold W Kuhn. The Hungarian method for the assignment problem. *Naval Research*  
677     *Logistics (NRL)* 2(1-2):83–97. Wiley Online Library, 1955.
- 678     **8**     James B. Orlin. A faster strongly polynomial minimum cost flow algorithm. *Operations*  
679     *Research* 41(2):338–350, 1993. <https://doi.org/10.1287/opre.41.2.338>.
- 680     **9**     Lyle Ramshaw and Robert Endre Tarjan. A weight-scaling algorithm for min-cost im-  
681     perfect matchings in bipartite graphs. *53rd Annual IEEE Symposium on Foundations of*  
682     *Computer Science, FOCS 2012, New Brunswick, NJ, USA, October 20-23, 2012*, 581–590,  
683     2012. <https://doi.org/10.1109/FOCS.2012.9>.
- 684     **10**     R. Sharathkumar and Pankaj K. Agarwal. Algorithms for the transportation problem in geo-  
685     metric settings. *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Dis-*  
686     *crete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, 306–317, 2012. <http://portal.acm.org/citation.cfm?id=2095145&CFID=63838676&CFTOKEN=79617016>.
- 687
- 688     **11**     Pravin M. Vaidya. Geometry helps in matching. *SIAM J. Comput.* 18(6):1201–1225, 1989.  
689     <https://doi.org/10.1137/0218080>.

## A Proofs from Section 3

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

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

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 that  $f'$  augmented flow in the reverse direction of  $v \rightarrow w$  ( $f'(w \rightarrow v) > 0$ ). By assumption, the arcs of  $\text{supp}(f')$  are admissible, so  $w \rightarrow v$  was an admissible arc ( $c_\pi(w \rightarrow v) \leq 0$ ). By antisymmetry of reduced costs, this implies  $c_\pi(v \rightarrow w) \geq 0 \geq -\varepsilon$ . Therefore, all arcs with  $u_{f+f'}(v, w) > 0$  respect the  $\varepsilon$ -optimality condition, and thus  $f + f'$  is  $\varepsilon$ -optimal. ◀

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

**Proof.** Because  $f$  is a circulation,  $\text{supp}(f)$  can be decomposed into  $k$  paths from  $s$  to  $t$ . Each  $s$ -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 the residual network must be induced by positive flow in the opposite direction, the total number of residual backward arcs is at most  $3k$ . ◀

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

775 The number of iterations belonging to the second phase is upper bounded by the amount  
 776 of total excess at the end of the first phase, because each subsequent push of a blocking

flow reduces the total excess by at least one. We now show that the amount of such excess is at most  $O(\sqrt{k})$ . 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 upper bounded by the number of arcs in  $E^+$  that crosses an arbitrarily given cut  $X$  that separates the excess vertices from the deficit vertices, when the network has unit-capacity [GHKT17, 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 the excess vertices from the deficit vertices at the end of first phase. Each arc in  $E^+$  crosses at most 3 cuts of type  $X_i$  [GHKT17, Lemma 3.1]. So there is one  $X_i$  crossed by at most  $3|E^+|/\sqrt{k}$  arcs in  $E^+$ . The size of  $E^+$  is bounded by the sum of support sizes of  $f$  and  $f_0$ ; by Corollary 3.6 the size of  $E^+$  is  $O(k)$ . This implies an  $O(\sqrt{k})$  bound on the total excess after the first phase, which in turn bounds the number of iterations in the second phase. ◀

## B Proofs from Section 4

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

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

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

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

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



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

- 823 1. Non-shortcut backward arcs  $(v', w')$  with  $(w', v') \in \text{supp}(f)$ . For these, we can maintain  
824 a min-heap on  $(w', v') \in \text{supp}(f)$  arcs for each normal  $v' \in V$ .
- 825 2. Non-shortcut  $A$ -to- $B$  forward arcs. For these, we maintain a NN data structure over  
826  $P = (B \setminus B_\emptyset) \setminus \tilde{S}$ , with weights  $\omega(p) = \pi(p)$  for each  $p \in P$ . We subtract  $\pi(v')$  from the  
827 NN distance to recover the reduced cost of the arc from  $v'$ .
- 828 3. Non-shortcut forward arcs from  $s$ -to- $A$  and from  $B$ -to- $t$ . For  $s$ , we can maintain a  
829 min-heap on the potentials of  $B \setminus \tilde{S}$ , queried only if  $v' = s$ . For  $B$ -to- $t$  arcs, there is only  
830 one arc to check if  $v' \in B$ , which we can examine manually.
- 831 4. Shortcut arcs  $(s, b)$  corresponding to null 2-paths from  $s$  to  $b \in (B \setminus B_\emptyset) \setminus S$ . For these, we  
832 maintain a BCP data structure with  $P = A_\emptyset$ ,  $Q = (B \setminus B_\emptyset) \setminus S$  with weights  $\omega(p) = \pi(s)$   
833 for all  $p \in P$ , and  $\omega(q) = \pi(q)$  for all  $q \in Q$ . A response  $(a, b)$  corresponds to th null  
834 2-path  $(s, a, b)$ . This is only queried if  $v' = s$ .
- 835 5. Shortcut arcs  $(a, t)$  corresponding to null 2-paths from  $a \in (A \setminus A_\emptyset) \cap S$  to  $t$ . For these,  
836 we maintain a NN data structure over  $P = B_\emptyset \setminus S$  with weights  $\omega(p) = \pi(t)$  for each  
837  $p \in P$ . A response  $(v', b)$  corresponds to th null 2-path  $(v', b, t)$ . We subtract  $\pi(v')$  from  
838 the NN distance to recover the reduced cost of the arc from  $v'$ . This is not queried if  
839  $t \in \tilde{S}$ .
- 840 6. Shortcut arcs  $(s, t)$  corresponding to null 3-paths. For these, we maintain in a BCP data  
841 structure with  $P = A_\emptyset \setminus S$ ,  $Q = B_\emptyset \setminus S$  with weights  $\omega(p) = \pi(s)$  for all  $p \in P$ , and  
842  $\omega(q) = \pi(t)$  for all  $q \in Q$ . A response  $(a, b)$  corresponds to th null 3-path  $(s, a, b, t)$ . This  
843 is only queried while  $v' = s$  and  $t \notin S$ .

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

## C Proofs from Section 5

**Recovering the optimal flow.** *«use this one if we want to use exponent  $> 1$ . »*

*«Move everything to appendix and left a pointer to the socg 2016 paper.»*

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

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

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

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

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

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

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

An arc  $a \rightarrow b$  is admissible when

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

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

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

**Proof.** We assume the points of  $A \cup B$  are in general position (e.g. by symbolic perturbation) 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  $\pi^*$ . We will isolate them from the rest of the points, considering  $\pi^*$  restricted to the four points  $\{a_1, a_2, b_1, b_2\}$ . Clearly, this does not change whether the two arcs cross. Observe that we can raise  $\pi^*(a_2)$  and  $\pi^*(b_2)$  uniformly, until  $c_\pi(a_2, b_1) = 0$ , without breaking feasibility or changing admissibility of  $a_1 \rightarrow b_1$  and  $a_2 \rightarrow b_2$ . Henceforth, we assume that we have modified  $\pi^*$  in this way 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  $a_1 \rightarrow b_1$  and  $a_2 \rightarrow b_2$  cross. Specifically,  $b_2$  must be placed within a region  $\mathcal{F}$  that lies between the rays  $\overrightarrow{a_2 a_1}$  and  $\overrightarrow{a_2 b_1}$ , and within the halfplane bounded by  $\overrightarrow{a_1 b_1}$  that does not contain  $a_2$ .

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  $\beta := \{q \in \mathbb{R}^2 \mid g_{a_1}(q) = g_{a_2}(q)\}$ .  $\beta$  is a curve subdividing the plane into two open faces, one where  $g_{a_1}$  is minimized and the other where  $g_{a_2}$  is. From these definitions, admissibility of  $a_1 \rightarrow b_1$  and  $a_2 \rightarrow b_1$  imply that  $b_1$  is a point of the bisector.

We show that  $\mathcal{F}$  lies entirely on the  $g_{a_1}$  side of the bisector. First, we prove that the closed 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 parametrically as  $q(t) = (1 - t)b_1 + ta_1$  for  $t \in [0, 1]$ . Consider the single-variable functions  $g_{a_1}(q(t))$  and  $g_{a_2}(q(t))$ .

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

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

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 as  $q(t') = q' + t'(q' - a_2)$  for some  $q' \in \overline{a_1 b_1}$  and  $t \geq 0$ , i.e.  $q' = \overline{a_1 b_1} \cap \overline{a_2 q}$ . We call  $q'$  the *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  $\frac{d}{dt'}g_{a_1}(q(t')) \leq \frac{d}{dt'}g_{a_2}(q(t'))$ , as the derivative of  $g_a(q(t'))$  is maximized if  $(q(t') - a)$  is parallel to  $(q(t') - a_2)$  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  $a_1, a_2, b_1$  are in general position. Any  $q(t')$  with a different projection do not have strict inequality, 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  $q \in \mathcal{F} \setminus \{b_1\}$ ,  $g_{a_1}(q') < g_{a_2}(q')$ , and  $\mathcal{F}$  lies on the  $g_{a_1}$  side of the bisector except for  $b_1$  which lies on  $\beta$ . We can eliminate  $b_1$  as a candidate position for  $b_2$ , since points of  $B$  cannot coincide.

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

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

1000 By prioritizing the relaxation of support arcs, we also have the following lemma.

1001 ► **Lemma C.2 (Agarwal et al. [AFP<sup>+</sup>17]).** *If arcs of  $\text{supp}(f)$  are relaxed first as they arrive*  
 1002 *on the frontier, then  $E(\text{supp}(f))$  is acyclic.*

1003 **Proof.** Let  $f_i$  be the pseudoflow after the  $i$ -th augmentation, and let  $T_i$  be the forest of  
 1004 relaxed arcs generated by the Hungarian search for the  $i$ -th augmentation. Namely, the  $i$ -th  
 1005 augmenting path is an excess-deficit path in  $T_i$ , and all arcs of  $T_i$  are admissible by the time  
 1006 the augmentation is performed. Let  $E(T_i)$  be the undirected edges corresponding to arcs of  
 1007  $T_i$ . Notice that,  $E(\text{supp}(f_{i+1})) \subseteq E(\text{supp}(f_i)) \cup E(T_i)$ . We prove that  $E(\text{supp}(f_i)) \cup E(T_i)$   
 1008 is acyclic by induction on  $i$ ; as  $E(\text{supp}(f_{i+1}))$  is a subset of these edges, it must also be  
 1009 acyclic. At the beginning with  $f_0 = 0$ ,  $E(\text{supp}(f_0))$  is vacuously acyclic.

1010 Let  $E(\text{supp}(f_i))$  be acyclic by induction hypothesis. Since  $T_i$  is a forest (thus, acyclic),  
 1011 any hypothetical cycle  $\Gamma$  that forms in  $E(\text{supp}(f_i)) \cup E(T_i)$  must contain edges from  
 1012 both  $E(\text{supp}(f_i))$  and  $E(T_i)$ . To give a visual analogy, we will color  $e \in \Gamma$  *purple* if  
 1013  $e \in E(\text{supp}(f_i)) \cap E(T_i)$ , *red* if  $e \in E(\text{supp}(f_i))$  but  $e \notin E(T_i)$ , and *blue* if  $e \in E(T_i)$  but  
 1014  $e \notin E(\text{supp}(f_i))$ . Then,  $\Gamma$  is neither entirely red nor entirely blue. We say that red and purple  
 1015 edges are *red-tinted*, and similarly blue and purple edges are *blue-tinted*. Roughly speaking,  
 1016 our implementation of the Hungarian search prioritizes relaxing red-tinted admissible arcs  
 1017 over pure blue arcs.

1018 We can sort the blue-tinted edges of  $\Gamma$  by the order they were relaxed into  $S$  during the  
 1019 Hungarian search forming  $T_i$ . Let  $(v, w) \in \Gamma$  be the last pure blue edge relaxed, of all the  
 1020 blue-tinted edges in  $\Gamma$  — after  $(v, w)$  is relaxed, the remaining unrelaxed, blue-tinted edges  
 1021 of  $\Gamma$  are purple.

1022 Let us pause the Hungarian search the moment before  $(v, w)$  is relaxed. At this point,  
 1023  $v \in S$  and  $w \notin S$ , and the Hungarian search must have finished relaxing all frontier support  
 1024 arcs. By our choice of  $(v, w)$ ,  $\Gamma \setminus (v, w)$  is a path of relaxed blue edges and red-tinted edges  
 1025 which connect  $v$  and  $w$ . Walking around  $\Gamma \setminus (v, w)$  from  $v$  to  $w$ , we see that every vertex of  
 1026 the cycle must be in  $S$  already:  $v \in S$ , relaxed blue edges have both endpoints in  $S$ , and any  
 1027 unrelaxed red-tinted edge must have both endpoints in  $S$ , since the Hungarian search would  
 1028 have prioritized relaxing the red-tinted edges to grow  $S$  before relaxing  $(v, w)$  (a blue edge).  
 1029 It follows that  $w \in S$  already, a contradiction.

1030 No such cycle  $\Gamma$  can exist, thus  $E(\text{supp}(f_i)) \cup E(T_i)$  is acyclic and  $E(\text{supp}(f_{i+1})) \subseteq$   
 1031  $E(\text{supp}(f_i)) \cup E(T_i)$  is acyclic. By induction,  $E(\text{supp}(f_i))$  is acyclic for all  $i$ . ◀

1032 Let  $E(\Sigma_a)$  **◀only used once▶** be the underlying edges of the support star centered at  $a$   
 1033 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  
 1034 support arcs outside support stars ( $|F|$ ) is small.

1035 ► **Lemma C.3.**  $|B_\ell \setminus \bigcup_{a \in A} \Sigma_a| \leq r$ .

1036 **Proof.**  $F$  is constructed from  $E(\text{supp}(f))$  by eliminating edges in support stars, therefore all  
 1037 edges in  $F$  must adjoin vertices in  $B$  of support degree at least 2. By Lemma C.2,  $E(\text{supp}(f))$   
 1038 is acyclic and therefore forms a spanning forest over  $A \cup B_\ell$ , so  $F$  is also a bipartite forest.  
 1039 All leaves of  $F$  are therefore vertices of  $A$ .

1040 Pick an arbitrary root for each connected component of  $F$  to establish parent-child  
 1041 relationships for each edge. As no vertex in  $B$  is a leaf, each vertex in  $B$  has at least one  
 1042 child. Charge each vertex in  $B$  to one of its children in  $F$ , which must belong to  $A$ . Each  
 1043 vertex in  $A$  is charged at most once. Thus, the number of  $B_\ell$  vertices outside of support  
 1044 stars is no more than  $r$ . ◀

1045 ► **Lemma 5.2.** *Suppose we have stripped the graph of dead vertices. The number of relaxation*  
 1046 *steps in a Hungarian search outside of support stars is  $O(r)$ .*

1047 **Proof.** If there are no dead vertices, then each non-support star relaxation step adds either  
 1048 (i) an active deficit vertex, (ii) a non-deficit vertex  $a \in A_\ell$ , or (iii) a non-deficit vertex  $b \in B_\ell$   
 1049 of support degree at least 2. There is a single relaxation of type (i), as it terminates the  
 1050 search. The number of vertices of type (ii) is  $r$ , and the number of vertices of type (iii) is at  
 1051 most  $r$  by Lemma C.3. The lemma follows. ◀

1052 ► **Lemma 5.3.** *Hungarian search takes  $O(r\sqrt{n} \text{ polylog } n)$  time.*

1053 **Proof.** The number of relaxation steps outside of support stars is  $O(r)$  by Lemma 5.2. The  
 1054 time per relaxation outside of support stars is  $O(\sqrt{n} \text{ polylog } n)$ . The time spent processing  
 1055 relaxations within a support star is  $O(\sqrt{n} \text{ polylog } n)$ , and at most  $r$  are relaxed during the  
 1056 search. The total time is therefore  $O(r\sqrt{n} \text{ polylog } n)$ . ◀

1057 Initially, we label stars big or small according to the  $\sqrt{n}$  threshold. A star that is currently  
 1058 big is turned into a small star once  $|\Sigma_a| \leq \sqrt{n}/2$ . A star that is currently small is turned into  
 1059 a big star once  $|\Sigma_a| \geq 2\sqrt{n}$ . This way, the time spent rebuilding/updating the respective  
 1060 data structures can be amortized to the insertions/deletions that preceded the switch, plus  
 1061 some  $O(r)$  extra work if the the update is small-to-big.

1062 A star  $\Sigma_a$  that is switching from big-to-small has size  $|\Sigma_a| \leq \sqrt{n}/2$ . When switching, we  
 1063 delete  $\mathcal{D}_{\text{big}}(a)$  and insert  $\Sigma_a$  into  $\mathcal{D}_{\text{small}}$ . Thus, the time spent for big-to-small update is  
 1064  $O(\sqrt{n} \text{ polylog } n)$ , and there were at least  $\sqrt{n}/2$  points removed from  $\Sigma_a$  since it was last big.

1065 A star  $\Sigma_a$  that is switching from small-to-big has size  $|\Sigma_a| = \sqrt{n} + x \geq 2\sqrt{n}$ , for some  
 1066 integer  $x \geq \sqrt{n}$ . Rearranging, we have  $|\Sigma_a| \leq 2x$ . When switching, we delete all  $|\Sigma_a|$   
 1067 points from  $\mathcal{D}_{\text{small}}$  and construct a new  $\mathcal{D}_{\text{big}}(a)$ . Constructing  $\mathcal{D}_{\text{big}}(a)$  requires inserting  
 1068  $O(r)$  points of  $A$  (into  $P$ ) and the  $|\Sigma_a|$  points of the star (into  $Q$ ). Thus, the time spent for  
 1069 a small-to-big update is  $O((r+x) \text{ polylog } n)$ , and there were at least  $x \geq \sqrt{n}$  points added  
 1070 to  $\Sigma_a$  since it was last small.

## 1071 — References for the Appendix —

- 1072 **AFP<sup>+</sup>17** Pankaj K. Agarwal, Kyle Fox, Debmalaya Panigrahi, Kasturi R. Varadarajan, and Allen  
 1073 Xiao. Faster algorithms for the geometric transportation problem. In *33rd International*  
 1074 *Symposium on Computational Geometry, SoCG 2017, July 4-7, 2017, Brisbane, Australia*,  
 1075 pages 7:1–7:16, 2017.
- 1076 **Eri10** Jeff Erickson. Maximum flows and parametric shortest paths in planar graphs. In *Proceed-*  
 1077 *ings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA*  
 1078 *2010, Austin, Texas, USA, January 17-19, 2010*, pages 794–804, 2010.
- 1079 **GHKT17** Andrew V. Goldberg, Sagi Hed, Haim Kaplan, and Robert E. Tarjan. Minimum-cost flows  
 1080 in unit-capacity networks. *Theory Comput. Syst.*, 61(4):987–1010, 2017.
- 1081 **GT90** Andrew V. Goldberg and Robert E. Tarjan. Finding minimum-cost circulations by success-  
 1082 ive approximation. *Math. Oper. Res.*, 15(3):430–466, 1990.
- 1083 **SA12** R. Sharathkumar and Pankaj K. Agarwal. Algorithms for the transportation problem in  
 1084 geometric settings. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium*  
 1085 *on Discrete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, pages 306–317,  
 1086 2012.
- 1087 **Vai89** Pravin M. Vaidya. Geometry helps in matching. *SIAM J. Comput.*, 18(6):1201–1225, 1989.