

# Efficient Algorithms for Geometric Partial Matching

Pankaj K. Agarwal

Hsien-Chih Chang

Allen Xiao

April 1, 2019

## Abstract

Let  $A$  and  $B$  be two point sets in the plane of sizes  $r$  and  $n$  respectively (assume  $r \leq n$ ), and let  $k$  be a parameter. A matching between  $A$  and  $B$  is a family of pairs in  $A \times B$  so that any point of  $A \cup B$  appears in at most one pair. Given two positive integers  $p$  and  $q$ , we define the cost of matching  $M$  to be  $c(M) = \sum_{(a,b) \in M} \|a - b\|_p^q$  where  $\|\cdot\|_p$  is the  $L_p$ -norm. The geometric partial matching problem asks to find the minimum-cost size- $k$  matching between  $A$  and  $B$ .

We present efficient algorithms for geometric partial matching problem that work for any powers of  $L_p$ -norm matching objective: An exact algorithm that runs in  $O((n + k^2) \text{polylog } n)$  time, and a  $(1 + \varepsilon)$ -approximation algorithm that runs in  $O((n + k\sqrt{k}) \text{polylog } n \cdot \log \varepsilon^{-1})$  time. Both algorithms are based on the primal-dual flow augmentation scheme; the main improvements involve using dynamic data structures to achieve efficient flow augmentations. With similar techniques, we give an exact algorithm for the planar transportation problem running in  $O(\min\{n^2, rn^{3/2}\} \text{polylog } n)$  time.

## 1 Introduction

Given two point sets  $A$  and  $B$  in the plane, we consider the problem of finding the minimum-cost partial matching between  $A$  and  $B$ . Formally, suppose  $A$  has size  $r$  and  $B$  has size  $n$  where  $r \leq n$ . Let  $G(A, B)$  be the undirected complete bipartite graph between  $A$  and  $B$ , and let the cost of edge  $(a, b)$  be  $c(a, b) = \|a - b\|_p^q$ , for some positive integers  $p$  and  $q$ . A *matching*  $M$  in  $G(A, B)$  is a set of edges sharing no endpoints. The *size* of  $M$  is the number of edges in  $M$ . The cost of matching  $M$ , denoted  $c(M)$ , is defined to be the sum of costs of edges in  $M$ . For a parameter  $k$ , the problem of finding the minimum-cost size- $k$  matching in  $G(A, B)$  is called the *geometric partial matching problem*. We call the corresponding problem in general bipartite graphs (with arbitrary edge costs) the *partial matching problem*.<sup>1</sup>

We also consider the following generalization of bipartite matching. Let  $\phi : A \cup B \rightarrow \mathbb{Z}$  be an integral *supply-demand function* with positive value on points of  $A$  and negative value on points of  $B$ , satisfying  $\sum_{a \in A} \phi(a) = -\sum_{b \in B} \phi(b)$ . Let  $U := \max_{p \in A \cup B} |\phi(p)|$ . A *transportation map* is a function  $\tau : A \times B \rightarrow \mathbb{R}_{\geq 0}$  such that  $\sum_{b \in B} \tau(a, b) = \phi(a)$  for all  $a \in A$  and  $\sum_{a \in A} \tau(a, b) = -\phi(b)$  for all  $b \in B$ . We define the cost of  $\tau$  to be

$$c(\tau) := \sum_{(a,b) \in A \times B} c(a, b) \cdot \tau(a, b).$$

The *transportation problem* asks to compute a transportation map of minimum cost.

---

<sup>1</sup>Partial matching is also called *imperfect matching* or *imperfect assignment* [10, 18].

**Related work.** Maximum-size bipartite matching is a classical problem in graph algorithms. Upper bounds include the  $O(m\sqrt{n})$  time algorithm by Hopcroft and Karp [11] and the  $O(m \min\{\sqrt{m}, n^{2/3}\})$  time algorithm by Even and Tarjan [8], where  $n$  is the number of nodes and  $m$  is the number of edges. The first improvement in over thirty years was made by Mądry [16], which uses an interior-point algorithm, runs in  $O(m^{10/7} \text{polylog } n)$  time.

The Hungarian algorithm [14] computes a minimum-cost maximum matching in a bipartite graph in roughly  $O(mn)$  time. Faster algorithms have been developed, such as the  $O(m\sqrt{n} \log(nC))$  time algorithms by Gabow and Tarjan [9] and the improved  $O(m\sqrt{n} \log C)$  time algorithm by Duan *et al.* [7] assuming the edge costs are integral; here  $C$  is the maximum cost of an edge. Ramshaw and Tarjan [18] showed that the Hungarian algorithm can be extended to compute a minimum-cost partial matching of size  $k$  in  $O(km + k^2 \log r)$  time, where  $r$  is the size of the smaller side of the bipartite graph. They also proposed a cost-scaling algorithm for partial matching that runs in time  $O(m\sqrt{k} \log(kC))$ , again assuming that costs are integral. By reduction to unit-capacity min-cost flow, Goldberg *et al.* [10] developed a cost-scaling algorithm for partial matching with an identical running time  $O(m\sqrt{k} \log(kC))$ , again only for integral edge costs.

In geometric settings, the Hungarian algorithm can be implemented to compute an optimal perfect matching between  $A$  and  $B$  (assuming equal size) in time  $O(n^2 \text{polylog } n)$  [12] (see also [1, 22]). This algorithm computes an optimal size- $k$  matching in time  $O(kn \text{polylog } n)$ . Faster approximation algorithms have been developed for computing perfect matchings in geometric settings [4, 19, 22, 23]. Recall that the cost of the edges are the  $q$ th power of their  $L_p$ -distances. When  $q = 1$ , the best algorithm to date by Sharathkumar and Agarwal [20] computes  $(1 + \varepsilon)$ -approximation to the value of optimal perfect matching in  $O(n \text{polylog } n \cdot \text{poly } \varepsilon^{-1})$  expected time with high probability. Their algorithm can also compute a  $(1 + \varepsilon)$ -approximate partial matching within the same time bound. For  $q > 1$ , the best known approximation algorithm to compute a perfect matching runs in  $O(n^{3/2} \text{polylog } n \log(1/\varepsilon))$  time [19]; it is not obvious how to extend this algorithm to the partial matching setting.

The transportation problem can also be formulated as an instance of the minimum-cost flow problem. The strongly polynomial uncappeditated min-cost flow algorithm by Orlin [17] solves the transportation problem in  $O((m + n \log n)n \log n)$  time. Lee and Sidford [15] give a weakly polynomial algorithm that runs in  $O(m\sqrt{n} \text{polylog}(n, U))$  time, where  $U$  is the maximum amount of node supply-demand. Agarwal *et al.* [2, 3] showed that Orlin's algorithm can be implemented to solve 2D transportation in time  $O(n^2 \text{polylog } n)$ . In case of  $O(1)$ -dimension Euclidean space, by adapting the Lee-Sidford algorithm, they developed a  $(1 + \varepsilon)$ -approximation algorithm that runs in  $O(n^{3/2} \text{poly } \varepsilon^{-1} \text{polylog}(n, U))$  time. They also gave a Monte-Carlo algorithm that computes an  $O(\log^2(1/\varepsilon))$ -approximate solution in  $O(n^{1+\varepsilon})$  time with high probability. Recently, Khesin, Niklov, and Paramonov [13] obtained a  $(1 + \varepsilon)$ -approximation in low-dimensional Euclidean space that runs in randomized  $O(n \text{poly } \varepsilon^{-1} \text{polylog}(n, U))$  time.

**Our results.** There are three main results in this paper. First in Section 2 we present an efficient algorithm for computing an optimal partial matching in the plane.

**Theorem 1.1.** *Given two point sets  $A$  and  $B$  in the plane each of size at most  $n$  and an integer  $k \leq n$ , a minimum-cost matching of size  $k$  between  $A$  and  $B$  can be computed in  $O((n + k^2) \text{polylog } n)$  time.*

We use *bichromatic closest pair (BCP)* data structures to implement the Hungarian algorithm efficiently, similar to Agarwal *et al.* [1] and Kaplan *et al.* [12]. But unlike their algorithms which take  $\Omega(n)$  time to find an augmenting path, we show that after  $O(n \text{polylog } n)$  time preprocessing,

an augmenting path can be found in  $O(k \text{ polylog } n)$  time. The key is to recycle (rather than rebuild) our data structures from one augmentation to the next. We refer to this idea as the *rewinding mechanism*.

Next in Sections 3, we obtain a  $(1+\varepsilon)$ -approximation algorithm for the geometric partial matching problem in the plane by providing an efficient implementation of the unit-capacity min-cost flow algorithm by Goldberg *et al.* [10].

**Theorem 1.2.** *Given two point sets  $A$  and  $B$  in  $\mathbb{R}^2$  each of size at most  $n$ , an integer  $k \leq n$ , and a parameter  $\varepsilon > 0$ , a  $(1 + \varepsilon)$ -approximate min-cost matching of size  $k$  between  $A$  and  $B$  can be computed in  $O((n + k\sqrt{k}) \text{ polylog } n \cdot \log \varepsilon^{-1})$  time.*

The main challenge here is how to deal with the *dead nodes*, which neither have excess/deficit nor have flow passing through them, but still contribute to the size of the graph. We show that the number of *alive nodes* is only  $O(k)$ , and then represent the dead nodes implicitly so that the Hungarian search and computation of a blocking flow can be implemented in  $O(k \text{ polylog } n)$  time.

Finally in Section 4 we present a faster algorithm for the transportation problem in  $\mathbb{R}^2$  when the two point sets are unbalanced.

**Theorem 1.3.** *Given two point sets  $A$  and  $B$  in  $\mathbb{R}^2$  of sizes  $r$  and  $n$  respectively with  $r \leq n$ , along with supply-demand function  $\phi : A \cup B \rightarrow \mathbb{Z}$ , an optimal transportation map between  $A$  and  $B$  can be computed in  $O(\min\{n^2, rn^{3/2}\} \text{ polylog } n)$  time.*

Our result improves over the  $O(n^2 \text{ polylog } n)$  time algorithm by Agarwal *et al.* [3] for  $r = o(\sqrt{n})$ . Similar to their algorithm, we also use the strongly polynomial uncapacitated minimum-cost flow algorithm by Orlin [17], but additional ideas are needed for efficient implementation. Unlike in the case of matchings, the support of the transportation problem may have size  $\Omega(n)$  even when  $r$  is a constant; so naively we can no longer spend time proportional to the size of support of the transportation map. However, with careful implementation we ensure that the support is acyclic, and one can find an augmenting path in  $O(r\sqrt{n} \text{ polylog } n)$  time with proper data structures, assuming  $r \leq \sqrt{n}$ .

## 2 Minimum-cost partial matchings using Hungarian algorithm

In this section, we solve the geometric partial matching problem and prove Theorem 1.1 by implementing the Hungarian algorithm for partial matching in  $O((n + k^2) \text{ polylog } n)$  time.

A node  $v$  is *matched* by matching  $M$  if  $v$  is the endpoint of some edge in  $M$ ; otherwise  $v$  is *unmatched*. Given a matching  $M$ , an *augmenting path*  $\Pi = (a_1, b_1, \dots, a_\ell, b_\ell)$  is an odd-length path with unmatched endpoints ( $a_1$  and  $b_\ell$ ) that alternates between edges outside and inside of  $M$ . The symmetric difference  $M \oplus \Pi$  creates a new matching of size  $|M| + 1$ , called the *augmentation* of  $M$  by  $\Pi$ . The dual to the standard linear program for partial matching has one dual variable for each node  $v$ , called the *potential*  $\pi(v)$  of  $v$ . Given potential  $\pi$ , we can define the *reduced cost* of the edges to be  $c_\pi(v, w) := c(v, w) - \pi(v) + \pi(w)$ . Potential  $\pi$  is *feasible* on edge  $(v, w)$  if  $c_\pi(v, w)$  is nonnegative. Potential  $\pi$  is *feasible* if  $\pi$  are feasible on every edge in  $G$ . We say that an edge  $(v, w)$  is *admissible* under potential  $\pi$  if  $c_\pi(v, w) = 0$ .

**Fast implementation of Hungarian search.** The Hungarian algorithm is initialized with  $M \leftarrow \emptyset$  and  $\pi \leftarrow 0$ . Each iteration of the Hungarian algorithm augments  $M$  by an admissible augmenting path  $\Pi$ , discovered using a procedure called the *Hungarian search*. The algorithm

terminates after  $k$  augmentations, exactly when  $|M| = k$ ; Ramshaw and Tarjan [18] showed that  $M$  is guaranteed to be an optimal partial matching.

The Hungarian search grows a set of *reachable nodes*  $X$  from all unmatched  $v \in A$  using augmenting paths of admissible edges. Initially,  $X$  is the set of unmatched nodes in  $A$ . Let the *frontier* of  $X$  be the edges in  $(A \cap X) \times (B \setminus X)$ .  $X$  is grown by *relaxing* an edge  $(a, b)$  in the frontier: Add  $b$  into  $X$ , modify potential to make  $(a, b)$  admissible, preserve  $c_\pi$  on other edges within  $X$ , and keep  $\pi$  feasible on edges outside of  $X$ . Specifically, the algorithm relaxes the min-reduced-cost frontier edge  $(a, b)$ , and then raises  $\pi(v)$  by  $c_\pi(a, b)$  for all  $v \in X$ . If  $b$  is already matched, then we also relax the matching edge  $(a', b)$  and add  $a'$  into  $X$ . The search finishes when  $b$  is unmatched, and an admissible augmenting path now can be recovered.

In the geometric setting, we find the min-reduced-cost frontier edge using a dynamic *bichromatic closest pair* (BCP) data structure, similar to [3, 22]. Given two point sets  $P$  and  $Q$  in the plane and a weight function  $\omega : P \cup Q \rightarrow \mathbb{R}$ , the BCP is two points  $a \in P$  and  $b \in Q$  minimizing the additively weighted distance  $c(a, b) - \omega(a) + \omega(b)$ . Thus, a minimum reduced-cost frontier edge is precisely the BCP of point sets  $P = A \cap X$  and  $Q = B \setminus X$ , with  $\omega = \pi$ . Note that the “state” of this BCP is parameterized by  $X$  and  $\pi$ .

The dynamic BCP data structure by Kaplan *et al.* [12] supports point insertions and deletions in  $O(\text{polylog } n)$  time and answers queries in  $O(\log^2 n)$  time for our setting. Each relaxation in the Hungarian search requires one query, one deletion, and at most one insertion (aside from the potential updates). As  $|M| \leq k$  throughout, there are at most  $2k$  relaxations in each Hungarian search, and the BCP can be used to implement each Hungarian search in  $O(k \text{ polylog } n)$  time.

**Rewinding mechanism.** We observe that exactly one node of  $A$  is newly matched after an augmentation. Thus (modulo potential changes), we can obtain the initial state of the BCP for the  $(i + 1)$ -th Hungarian search from the  $i$ -th one with a single BCP deletion.

If we remember the sequence of points added to  $X$  in the  $i$ -th Hungarian search, then at the start of the  $(i + 1)$ -th Hungarian search we can *rewind* this sequence by applying the opposite insert/delete operation to each BCP update in reverse order to obtain the initial state of the  $i$ -th BCP. With one additional BCP deletion, we have the initial state of the  $(i + 1)$ -th BCP. The number of insertions/deletions is bounded by the number of relaxations per Hungarian search which is  $O(k)$ . Therefore we can recover, in  $O(k \text{ polylog } n)$  time, the initial BCP data structure for each Hungarian search beyond the first. We refer to this procedure as the *rewinding mechanism*.

**Potential updates.** We modify a trick from Vaidya [22] to batch potential updates. Potential is tracked with a *stored value*  $\gamma(v)$ , while the *true value* of  $\pi(v)$  may have changed since  $\gamma(v)$  was last recorded. This is done by aggregating potential changes into a variable  $\delta$ , which is initially 0 at the very beginning of the algorithm. Whenever we would raise the potential of all nodes in  $X$ , we raise  $\delta$  by that amount instead. We maintain the following invariant:  $\pi(v) = \gamma(v)$  for  $v \notin X$ , and  $\pi(v) = \gamma(v) + \delta$  for  $v \in X$ .

At the beginning of the algorithm,  $X$  is empty and stored values are equal to true values. When  $a \in A$  is added to  $X$ , we update its stored value to  $\pi(a) - \delta$  for the current value of  $\delta$ , and use that stored value as its BCP weight. Since the BCP weights are uniformly offset from  $\pi(v)$  by  $\delta$ , the pair reported by the BCP is still minimum. When  $b \in B$  is added to  $X$ , we update its stored value to  $\pi(b) - \delta$  (although it won't be added to a BCP set). When a node is removed from  $X$  (e.g. by augmentation or rewinding), we update the stored potential  $\gamma(v) \leftarrow \pi(v) + \delta$ , again for the current value of  $\delta$ . Unlike Vaidya [22] and for the sake of rewinding, we do not reset  $\delta$  across Hungarian searches.

There are  $O(k)$  relaxations and thus  $O(k)$  updates to  $\delta$  per Hungarian search.  $O(k)$  stored values are updated per rewinding, so the time spent on potential updates per Hungarian search is  $O(k)$ . Putting everything together, our implementation of the Hungarian algorithm runs in  $O((n + k^2) \text{polylog } n)$  time. This proves Theorem 1.1.

### 3 Approximating min-cost partial matching through cost-scaling

In this section we describe an approximation algorithm for computing a min-cost partial matching. We reduce the problem to computing a min-cost circulation in a flow network (Section 3.1). We adapt the cost-scaling algorithm by Goldberg *et al.* [10] for computing min-cost flow of a unit-capacity network (Section 3.2). Finally, we show how their algorithm can be implemented in  $O((n + k^{3/2}) \text{polylog}(n) \log(1/\varepsilon))$  time in our setting (Section 3.3).

#### 3.1 From matching to circulation

Given a bipartite graph  $G$  with node sets  $A$  and  $B$ , we construct a flow network  $N = (V, \vec{E})$  in a standard way [18] so that a min-cost matching in  $G$  corresponds to a min-cost integral circulation in  $N$ .

**Flow network.** Each node in  $G$  becomes a node in  $N$  and each edge  $(a, b)$  in  $G$  becomes an arc  $a \rightarrow b$  in  $N$ ; we refer to these nodes and arcs as *bipartite nodes* and *bipartite arcs*. We also include a *source* node  $s$  and *sink* node  $t$  in  $N$ . For each  $a \in A$ , we add a *left dummy arc*  $s \rightarrow a$  and for each  $b \in B$  a *right dummy arc*  $b \rightarrow t$ . The cost  $c(v \rightarrow w)$  is equal to  $c(v, w)$  if  $v \rightarrow w$  is a bipartite arc and 0 if  $v \rightarrow w$  is a dummy arc. All arcs in  $N$  have unit capacity.

Let  $\phi : V \rightarrow \mathbb{Z}$  be an integral supply/demand function on nodes of  $N$ . The positive values of  $\phi(v)$  are referred to as *supply*, and the negative values of  $\phi(v)$  as *demand*. A *pseudoflow*  $f : \vec{E} \rightarrow [0, 1]$  is a function on arcs of  $N$ . The *support* of  $f$  in  $N$ , denoted as  $\text{supp}(f)$ , is the set of arcs with positive flow. Given a pseudoflow  $f$ , the *imbalance* of a node is

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

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

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

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

If we set  $\phi(s) = k$ ,  $\phi(t) = -k$ , and  $\phi(v) = 0$  for all  $v \in A \cup B$ , then an integral circulation  $f$  corresponds to a partial matching  $M$  of size  $k$  and vice versa. Moreover,  $c(M) = c(f)$ . Hence, the problem of computing a min-cost matching of size  $k$  in  $G$  transforms to computing an integral circulation in  $N$ . The following lemma will be useful for our algorithm.

**Lemma 3.1.** *Let  $N$  be the network constructed from the bipartite graph  $G$  above.*

- (i) *For any integral circulation  $g$  in  $N$ , the size of  $\text{supp}(g)$  is at most  $3k$ .*
- (ii) *For any integral pseudoflow  $f$  in  $N$  with  $K$  excess, the size of  $\text{supp}(f)$  is at most  $3(K + k)$ .*

*Proof.* (i) Because  $g$  is a circulation,  $\text{supp}(g)$  can be decomposed into  $k$  paths from  $s$  to  $t$ . Each  $s$ -to- $t$  path in  $N$  is of length three, so  $|\text{supp}(g)| \leq 3k$ .

(ii) The graph in  $N$  is a directed acyclic graph, so  $\text{supp}(f)$  does not form a directed cycle. Thus,  $\text{supp}(f)$  can be decomposed into a set of inclusion-maximal paths, each of which contributes a single unit of excess to the flow if the path does not terminate at  $t$  or if more than  $k$  paths terminate at  $t$ . By assumption, there are  $K$  units of excess to which we can associate to the paths, and at most  $k$  paths terminating at  $t$  that we cannot associate with a unit of excess. Each inclusion-maximal path has length at most 3, so  $|\text{supp}(f)| \leq 3(K + k)$ .  $\square$

### 3.2 A cost-scaling algorithm

Before describing the algorithm, we need to introduce a few more concepts.

**Residual network and admissibility.** If  $f$  is an integral pseudoflow on  $N$  (that is,  $f(v \rightarrow w) \in \{0, 1\}$  for every arc in  $\vec{E}$ ), then each arc  $v \rightarrow w$  in  $N$  is either *idle* with  $f(v \rightarrow w) = 0$  or *saturated* with  $f(v \rightarrow w) = 1$ .

Given a pseudoflow  $f$ , the *residual network*  $N_f = (V, \vec{E}_f)$  is defined as follows. For each idle arc  $v \rightarrow w$  in  $\vec{E}$ , we add a *forward* residual arc  $v \rightarrow w$  in  $N_f$ . For each saturated arc  $v \rightarrow w$  in  $\vec{E}$ , we add a *backward* residual arc  $w \rightarrow v$  in  $N_f$ . The set of residual arcs in  $N_f$  is therefore

$$\vec{E}_f := \{v \rightarrow w \mid f(v \rightarrow w) = 0\} \cup \{w \rightarrow v \mid f(v \rightarrow w) = 1\}.$$

The cost of a forward residual arc  $v \rightarrow w$  is  $c(v \rightarrow w)$ , while the cost of a backward residual arc  $w \rightarrow v$  is  $-c(v \rightarrow w)$ . Each arc in  $N_f$  also has unit capacity. By Lemma 3.1,  $N_f$  has  $O(k)$  backward arcs if  $f$  has  $O(k)$  excess.

A *residual pseudoflow*  $g$  in  $N_f$  can be used to change  $f$  into a different pseudoflow on  $N$  by *augmentation*. For simplicity, we only describe augmentation for the case where  $f$  and  $g$  are integral. Specifically, augmenting  $f$  by  $g$  produces a pseudoflow  $f'$  in  $N$  where

$$f'(v \rightarrow w) = \begin{cases} 0 & w \rightarrow v \in \vec{E}_f \text{ and } g(w \rightarrow v) = 1, \\ 1 & v \rightarrow w \in \vec{E}_f \text{ and } g(v \rightarrow w) = 1, \\ f(v \rightarrow w) & \text{otherwise.} \end{cases}$$

Using LP duality for min-cost flow, we assign *potential*  $\pi(v)$  to each node  $v$  in  $N$ . The *reduced cost* of an arc  $v \rightarrow w$  in  $N$  with respect to  $\pi$  is defined as

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

Similarly we define the reduced cost of arcs in  $N_f$ : the reduced cost of a forward residual arc  $v \rightarrow w$  in  $N_f$  is  $c_\pi(v \rightarrow w)$ , and the reduced cost of a backward residual arc  $w \rightarrow v$  in  $N_f$  is  $-c_\pi(v \rightarrow w)$ . Abusing the notation, we also use  $c_\pi$  to denote the reduced cost of arcs in  $N_f$ .

The *dual feasibility constraint* asks that  $c_\pi(v \rightarrow w) \geq 0$  holds for every arc  $v \rightarrow w$  in  $\vec{E}$ ; potential  $\pi$  that satisfy this constraint is 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  *$\theta$ -optimal* [5, 21] if  $c_\pi(v \rightarrow w) \geq -\theta$  for every residual arc  $v \rightarrow w$  in  $\vec{E}_f$ . Pseudoflow  $f$  is  *$\theta$ -optimal* if it is  $\theta$ -optimal with respect to some potential  $\pi$ ; potential  $\pi$  is  *$\theta$ -optimal* if it is  $\theta$ -optimal with respect to some pseudoflow  $f$ . Given a pseudoflow  $f$  and potential  $\pi$ , 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  $N_f$  is *admissible* if  $g(v \rightarrow w) > 0$  only on admissible arcs  $v \rightarrow w$ , and  $g(v \rightarrow w) = 0$  otherwise.<sup>2</sup> We will use the following well-known property of  $\theta$ -optimality.

**Lemma 3.2.** *Let  $f$  be an  $\theta$ -optimal pseudoflow in  $N$  and let  $g$  be an admissible pseudoflow in  $N_f$ . Then  $f$  augmented by  $g$  is also  $\theta$ -optimal in  $N$ .*

Using Lemma 3.1, the following lemma can be proved about  $\theta$ -optimality:

**Lemma 3.3.** *Let  $f$  be a  $\theta$ -optimal integer circulation in  $N$ , and  $f^*$  be an optimal integer circulation for  $N$ . Then,  $c(f) \leq c(f^*) + 6k\theta$ .*

*Proof.* Taking the symmetric difference between  $f$  and  $f^*$  gives a residual pseudoflow  $g$  of  $N_f$ , where augmenting  $f$  by  $g$  produces  $f^*$ . Then,  $c(f^*) = c(f) + c(g)$ . Since both  $f$  and  $f^*$  are circulations,  $g$  is comprised of cycle flows in  $N_f$ . Using Lemma 3.1, we have  $\text{supp}(g) \leq 6k$ , as exactly half the arcs in  $\text{supp}(g)$  are from  $\text{supp}(f)$  and half are from  $\text{supp}(f^*)$ . For any residual cycle, its cost is equal to its reduced cost as the potential telescopes. Thus, we can lower bound  $c(g)$  using the reduced cost of its cycles. By applying  $\theta$ -optimality of  $f$ ,  $c(g) \geq -6k\theta$ . Rearranging gives the lemma.  $\square$

**Estimating the value of  $c(f^*)$ .** We now describe a procedure for estimating  $c(f^*)$  within a polynomial factor, which is used to initialize the cost-scaling algorithm.

Let  $T$  be a minimum spanning tree of  $A \cup B$  under the cost function  $c$ . Let  $e_1, e_2, \dots, e_{n-1}$  be the edges of  $T$  sorted in nondecreasing order of length. Let  $T_i$  be the forest consisting of the nodes of  $A \cup B$  and edges  $e_1, \dots, e_i$ . We call a matching  $M$  *intra-cluster* if both endpoints of each edge in  $M$  lie in the same connected component of  $T_i$ . The following lemma will be used by our cost-scaling algorithm:

**Lemma 3.4.** *Let  $A$  and  $B$  be two point sets in the plane. Define  $i^*$  to be the smallest index  $i$  such that there is an intra-cluster matching of size  $k$  in  $T_{i^*}$ . Set  $\bar{\theta} := n^q \cdot c(e_{i^*})$ . Then*

- (i) *The value of  $c(e_{i^*})$  can be computed in  $O(n \log n)$  time.*
- (ii)  *$c(e_{i^*}) \leq c(f^*) \leq \bar{\theta}$ .*
- (iii) *There is a  $\bar{\theta}$ -optimal circulation in the network  $N$  with respect to the all-zero potential, assuming  $\phi(s) = k$ ,  $\phi(t) = -k$ , and  $\phi(v) = 0$  for all  $v \in A \cup B$ .*

*Proof.* (i) Observe that the MST under  $c(\cdot) = \|\cdot\|_p^q$  is identical (in terms of edges used) to the MST under  $\|\cdot\|_p$ , the latter of which can be computed in  $O(n \log n)$  time [6]. After computing the MST, we can sort its edges by cost in  $O(n \log n)$  time.

Suppose we have the sorted sequence  $e_1, \dots, e_{n-1}$ . For a cluster  $K \in T_i$ , let  $A_K = A \cap K$  and  $B_K = B \cap K$  respectively. For any  $i$ , largest intra-cluster partial matching of  $T_i$  has size

$$\#(T_i) := \sum_{K \in T_i} \min\{|A_K|, |B_K|\},$$

i.e., matches either  $A_K$  or  $B_K$  entirely for each cluster  $K$ . By definition,  $i^*$  is the smallest index for which  $\#(T_{i^*}) \geq k$ .

---

<sup>2</sup>The same admissibility/feasibility definitions will be used later in Section 4. However, the algorithm in Section 4 maintains a 0-optimal  $f$  and therefore admissible residual arcs always have  $c_\pi(v \rightarrow w) = 0$ .

To compute  $\#(T_i)$  efficiently, we maintain the  $|A_K|$  and  $|B_K|$  for each  $K \in T_i$ . Suppose that we have this per-cluster information for  $T_{i-1}$  and have computed  $\#(T_{i-1})$  previously.  $T_i$  is constructed from  $T_{i-1}$  by merging two clusters  $K_1, K_2 \in T_{i-1}$  to form a cluster  $K_3 \in T_i$ . Thus, we have  $|A_{K_3}| = |A_{K_1}| + |A_{K_2}|$  and similarly for  $|B_{K_3}|$ . Furthermore,

$$\#(T_i) = \#(T_{i-1}) - \min\{|A_{K_1}|, |B_{K_1}|\} - \min\{|A_{K_2}|, |B_{K_2}|\} + \min\{|A_{K_3}|, |B_{K_3}|\}.$$

This way, we can compute  $\#(T_i)$  from the  $T_{i-1}$  information in  $O(1)$  time. It takes  $O(n)$  time to compute the per-cluster information of  $T_0$ , the all-singletons clustering. The time to compute  $\#(T_i)$  for all  $i$  is therefore  $O(n)$ , and the total time to find  $c(e_{i^*})$  is  $O(n \log n)$ .

- (ii) Let  $M^*$  be the optimal matching. At least one edge in  $M^*$  must be between clusters of  $T_{i^*}$ , thus by definition of  $i^*$  at least one edge has cost at least  $c(e_{i^*})$ . The statement follows as  $c(M^*) = c(f^*)$ .
- (iii) The intra-cluster size- $k$  matching in  $T_{i^*}$  uses edges of cost at most  $\bar{\theta} = n^q \cdot c(e_{i^*})$ . This matching corresponds to a circulation  $f_{\text{intra}}$  supported on arcs of cost at most  $\bar{\theta}$ . In the resulting residual network, backward arcs have reduced cost  $c_\pi(w \rightarrow v) = -c(v \rightarrow w) - 0 + 0 \geq -\bar{\theta}$ . Reduced cost of residual forward arcs is positive, so overall  $f_{\text{intra}}$  is  $\bar{\theta}$ -optimal with respect to  $\pi = 0$ .  $\square$

As a consequence of Lemmas 3.4(ii) and 3.3, we have:

**Corollary 3.5.** *The cost of a  $\underline{\theta}$ -optimal integral circulation in  $N$  is at most  $(1 + \varepsilon)c(f^*)$ , where  $\underline{\theta} := \frac{\varepsilon}{6k} \cdot c(e_{i^*})$ .*

**Overview of the algorithm.** We are now ready to describe our algorithm, which closely follows Goldberg *et al.* [10]. The algorithm works in rounds called *scales*. Within each scale, we fix a *cost scaling parameter*  $\theta$  and maintain potential  $\pi$  with the following property:

- (\*) There exists a  $2\theta$ -optimal integral circulation in  $N$  with respect to  $\pi$ .

For the initial scale, we set  $\theta \leftarrow \bar{\theta}$  and  $\pi \leftarrow 0$ . By Lemma 3.4(iii), property (\*) is satisfied initially. Each scale of the algorithm consists of two stages. In the *scale initialization* stage, SCALE-INIT computes a  $\theta$ -optimal pseudoflow  $f$ . In the *refinement* stage, REFINES converts  $f$  into a  $\theta$ -optimal (integral) circulation  $g$ . In both stages,  $\pi$  is updated as necessary. If  $\theta \leq \underline{\theta}$ , we return  $g$ . Otherwise, we set  $\theta \leftarrow \theta/2$  and start the next scale. Note that property (\*) is satisfied in the beginning of each scale.

By Corollary 3.5, when the algorithm terminates, it returns an integral circulation  $\tilde{f}$  in  $N$  of cost at most  $(1 + \varepsilon)c(f^*)$ , which corresponds to a  $(1 + \varepsilon)$ -approximate min-cost matching of size  $k$  in  $G$ . The algorithm terminates in  $\log_2(\bar{\theta}/\underline{\theta}) = O(\log(n/\varepsilon))$  scales.

**Scale initialization.** In the first scale, we compute a  $\bar{\theta}$ -optimal pseudoflow by simply setting  $f(v \rightarrow w) \leftarrow 0$  for all arcs in  $\vec{E}$ . For subsequent scales, we initialize  $f$  to the  $2\theta$ -optimal circulation of the previous scale. First, we raise the potential of all nodes in  $A$  by  $\theta$ , all nodes in  $B$  by  $2\theta$ , and  $t$  by  $3\theta$ . The potential of  $s$  is unchanged. Such potential change increases the reduced cost of all forward arcs to at least  $-\theta$ .

Next, for each backward arc  $w \rightarrow v$  in  $N_f$  with  $c_\pi(w \rightarrow v) < -\theta$ , we set  $f(v \rightarrow w) \leftarrow 0$  (that is, make arc  $v \rightarrow w$  idle), which replaces the backward arc  $w \rightarrow v$  in  $N_f$  with forward arc  $v \rightarrow w$  of positive



reduced cost. After this step, the resulting pseudoflow must be  $\theta$ -optimal as all arcs of  $N_f$  have reduced cost at least  $-\theta$ .

The desaturation of each backward arc creates one unit of excess. Since there are at most  $3k$  backward arcs, the pseudoflow has at most  $3k$  excess after SCALE-INIT. There are  $O(n)$  potential updates and  $O(k)$  arcs to desaturate, so the time required for SCALE-INIT is  $O(n)$ .

**Refinement.** The procedure REFINE converts a  $\theta$ -optimal pseudoflow with  $O(k)$  excess into a  $\theta$ -optimal circulation, using a primal-dual augmentation algorithm. A path in  $N_f$  is an *augmenting path* if it begins at an excess node and ends at a deficit node. We call an admissible pseudoflow  $g$  in  $N_f$  an *admissible blocking flow* if  $g$  saturates at least one arc in every admissible augmenting path in  $N_g$ . In other words, there is no admissible excess-deficit path in the residual network after augmentation by  $g$ . Each iteration of REFINE finds an admissible blocking flow to be added to the current pseudoflow in two steps:

1. *Hungarian search*: a Dijkstra-like search that begins at the set of excess nodes and raises potential until there is an excess-deficit path of admissible arcs in  $N_f$ .
2. *Augmentation*: construct an admissible blocking flow by performing depth-first search on the set of admissible arcs of  $N_f$ . It suffices to repeatedly extract admissible augmenting paths until no more admissible excess-deficit paths remain.

The algorithm repeats these steps until the total excess becomes zero. The following lemma bounds the number of iterations in the REFINE procedure at each scale.

**Lemma 3.6.** *Let  $\theta$  be the scaling parameter and  $\pi_0$  the potential function at the beginning of a scale, such that there exists an integral  $2\theta$ -optimal circulation with respect to  $\pi_0$ . Let  $f$  be a  $\theta$ -optimal pseudoflow with excess  $O(k)$ . Then REFINE terminates within  $O(\sqrt{k})$  iterations.*

*Proof.* **«TODO flesh out»** We sketch the proof, which is adapted from Goldberg *et al.* [10]. Let  $f_0$  be the assumed  $2\theta$ -optimal integral circulation with respect to  $\pi_0$ , and let  $\pi$  be the potential maintained during REFINE. Let  $d(v) := (\pi(v) - \pi_0(v))/\theta$ , that is, the increase in potential at  $v$  in units of  $\theta$ . We divide the iterations of REFINE into two phases: before and after every (remaining) excess node has  $d(v) \geq \sqrt{k}$ . Each Hungarian search raises excess potential by at least  $\theta$ , since we use blocking flows. Thus, the first phase lasts at most  $\sqrt{k}$  iterations.

At the start of the second phase, consider the set of arcs  $E^+ := \{v \rightarrow w \in \vec{E} \mid f(v \rightarrow w) < f_0(v \rightarrow w)\}$ . One can argue that the remaining excess with respect to  $f$  is bounded above by the size of any cut separating the excess and deficit nodes [10, Lemma 4]. The proof examines cuts  $Y_i := \{v \mid d(v) > i\}$  for  $0 \leq i \leq \sqrt{k}$ . By  $\theta$ -optimality of  $f$  and  $2\theta$ -optimality of  $f_0$ , one can show that each arc in  $E^+$  crosses at most 3 cuts. Furthermore, the size of  $E^+$  is  $O(k)$ , bounded by the support size of  $f$  and  $f_0$ . Averaging, there is a cut among  $Y_i$ s of size at most  $3k/\sqrt{k}$ , so the total excess remaining is  $O(\sqrt{k})$ . Each iteration of REFINE eliminates at least one unit of excess, so the number of second phase iterations is also at most  $O(\sqrt{k})$ .  $\square$

In the next subsection we show that after  $O(n \text{ polylog } n)$  time preprocessing, an iteration of REFINE can be performed in  $O(k \text{ polylog } n)$  time (Lemma 3.9). By Lemma 3.6 and the fact the algorithm terminates in  $O(\log(n/\varepsilon))$  scales, the overall running time of the algorithm is  $O((n + k^{3/2}) \text{ polylog } n \log(1/\varepsilon))$ , as claimed in Theorem 1.2.

### 3.3 Fast implementation of refinement stage

We now describe a fast implementation of REFINE. The Hungarian search and augmentation steps are similar: each traversing through the residual network using admissible arcs starting from the excess nodes. We describe the Hungarian search first.

At a high level, let  $X$  be the subset of nodes visited by the Hungarian search so far. Initially  $X$  is the set of excess nodes. At each step, the algorithm finds a minimum-reduced-cost arc  $v \rightarrow w$  in  $N_f$  from  $X$  to  $V \setminus X$ . If  $v \rightarrow w$  is not admissible, the potential of all nodes in  $X$  is increased by  $\lceil c_\pi(v \rightarrow w)/\theta \rceil$  to make  $v \rightarrow w$  admissible. If  $w$  is a deficit node, the search terminates. Otherwise,  $w$  is added to  $X$  and the search continues.

Implementing the Hungarian search efficiently is more difficult than in Section 2 because (a) excess nodes may show up in  $A$  as well as in  $B$ , (b) a balanced node may become imbalanced later in the scales, and (c) the potential of excess nodes may be non-uniform. We therefore need a more complex data structure.

We call a node  $v$  of  $N$  *dead* if  $\phi_f(v) = 0$  and no arc of  $\text{supp}(f)$  is incident to  $v$ ; otherwise  $v$  is *alive*. Note that  $s$  and  $t$  are always alive. Let  $A^*$  denote the set of alive nodes in  $A$ ; define  $B^*$  similarly. There are only  $O(k)$  alive nodes, as each can be charged to its adjoining  $\text{supp}(f)$  arcs or its imbalance. We treat alive and dead nodes separately to implement the Hungarian search efficiently. By definition, dead nodes only adjoin forward arcs in  $N_f$ . Thus, the in-degree (resp. out-degree) of a node in  $A \setminus A^*$  (resp.  $B \setminus B^*$ ) is 1, and any path passing through a dead node has a subpath of the form  $s \rightarrow v \rightarrow b$  for some  $b \in B$  or  $a \rightarrow v \rightarrow t$  for some  $a \in A$ . Consequently, a path in  $N_f$  may have at most two consecutive dead nodes, and in the case of two consecutive dead nodes there is a subpath of the form  $s \rightarrow v \rightarrow w \rightarrow t$  where  $v \in A \setminus A^*$  and  $w \in B \setminus B^*$ . We call such paths, from an alive node to an alive node with only dead interior nodes, *alive paths*. Let the reduced cost  $c_\pi(\Pi)$  of an alive path  $\Pi$  be the sum of  $c_\pi$  over its arcs. We say  $\Pi$  is *weakly admissible* if  $c_\pi(\Pi) \leq 0$ .

We find the min-reduced-cost alive path of lengths 1, 2, and 3 leaving  $X$ , then relax the cheapest among them (raise potential of  $X$  by  $\lceil c_\pi(\Pi)/\theta \rceil$  and add every node of  $\Pi$  into  $X$ ). Essentially, relaxing alive paths “skips over” dead nodes. Since reduced costs telescope on paths, weak admissibility of an alive path depends only on the potential of its alive endpoints.

$$c_\pi(\Pi = v_1 \rightarrow \cdots \rightarrow v_\ell) = \sum_{i=1}^{\ell-1} c(v_i \rightarrow v_{i+1}) - \pi(v_i) + \pi(v_{i+1}) = c(\Pi) - \pi(v_1) + \pi(v_\ell)$$

Thus, we can query the minimum alive path using a partial assignment of  $\pi$  on only the alive nodes, leaving  $\pi$  over the dead nodes untracked. We now describe a data structure for each path length. Note that our “time budget” per Hungarian search is  $O(k \text{ polylog } n)$ .

**Finding length-1 paths.** This data structure finds a min-reduced-cost arc from an alive node of  $X$  to an alive node of  $V \setminus X$ . There are  $O(k)$  backward arcs, so the minimum among backward arcs can be maintained explicitly in a priority queue and retrieved in  $O(1)$  time.

There are three types of forward arcs:  $s \rightarrow a$  for some  $a \in A^*$ ,  $b \rightarrow t$  for some  $b \in B^*$ , and bipartite arc  $a \rightarrow b$  with two alive endpoints. Arcs of the first (resp. second) type can be found by maintaining  $A^* \setminus X$  (resp.  $B^* \cap X$ ) in a priority queue, but should only be queried if  $s \in X$  (resp.  $t \notin X$ ). The cheapest arc of the third type can be maintained using a dynamic BCP data structure between  $A^* \cap X$  and  $B^* \setminus X$ , with reduced cost as the weighted pair distance. Such a data structure can be implemented so that insertions/deletions can be performed in  $O(\text{polylog } k)$  time [12].

**Finding length-2 paths.** We describe how to find a cheapest path of the form  $s \rightarrow v \rightarrow b$  where  $v$  is dead and  $b \in B^*$ . A cheapest path  $a \rightarrow v \rightarrow t$  can be found similarly. Similar to length-1 paths, we

only query paths starting at  $s$  if  $s \in X$  and paths ending at  $t$  if  $t \notin X$ .

Note that  $c_\pi(s \rightarrow v \rightarrow b) = c(v, b) + \pi(b) - \pi(s)$ . Since  $\pi(s)$  is common in all such paths, it suffices to find a pair  $(v, w)$  between  $A \setminus A^*$  and  $B^* \setminus X$  minimizing  $c(v, w) + \pi(w)$ . This is done by maintaining a dynamic BCP data structure between  $A \setminus A^*$  and  $B^* \setminus X$  with the cost of a pair  $(v, w)$  being  $c(v, w) + \pi(w)$ . We may require an update operation for each alive node added to  $X$  during the Hungarian search, of which there are  $O(k)$ , so the time spent during a search is  $O(k \text{ polylog } n)$ .

Since the size of  $A \setminus A^*$  is at least  $r - k$ , we cannot construct this BCP from scratch at the beginning of each iteration. To resolve this, we use the idea of rewinding from Section 2, with a slight twist. There are now *two* ways that the initial BCP may change across consecutive Hungarian searches: (1) the initial set  $X$  may change as nodes lose excess through augmentation, and (2) the set of alive/dead nodes in  $A$  may change. The first is identical to the situation in Section 2; the number of excess depletions is  $O(k)$  over the course of REFINE. For the second, the alive/dead status of a node can change only if the blocking flow found passes through the node. By Lemma 3.8 below, there are  $O(k)$  such changes per Hungarian search, which can be done in  $O(k \text{ polylog } n)$  time.

**Finding length-3 paths.** We now describe how to find the cheapest path of the form  $s \rightarrow v \rightarrow w \rightarrow t$  where  $v \in A \setminus A^*$  and  $w \in B \setminus B^*$ . Note that  $c_\pi(s \rightarrow v \rightarrow w \rightarrow t) = c(v \rightarrow w) - \pi(s) + \pi(t)$ . A pair  $(v, w)$  between  $A \setminus A^*$  and  $B \setminus B^*$  minimizing  $c(v, w)$  can be found by maintaining a dynamic BCP data structure similar to the case of length-2 paths.

This BCP data structure has no dependency on  $X$ —the only update required comes from membership changes to  $A^*$  or  $B^*$  after an augmentation. Applying Lemma 3.8 again, there are  $O(k)$  alive/dead updates caused by an augmentation, so the time for these updates per Hungarian search is  $O(k \text{ polylog } n)$ .

**Correctness — the shortcut network.** «TODO figure, check the writing of this part» Before continuing to describe potential updates, we prove that relaxing alive paths keeps the Hungarian search intact. Specifically, our process of relaxing alive paths in  $N$  corresponds to arc-by-arc relaxations in an equivalent (analysis-only) network we call the *shortcut network*  $\tilde{N}$ , ultimately creating a suitable potential.

The shortcut network is constructed from  $N$  and  $f$  as follows: the nodes of  $\tilde{N}$  are the alive nodes of  $N$ , and for each length-2, -3, and forward length-1 alive path  $\Pi = v \rightarrow \dots \rightarrow w$  in  $N$ , we add an arc  $v \rightarrow w$  in  $\tilde{N}$  with  $c(v \rightarrow w) \leftarrow c(\Pi)$ . This arc is the *shortcut* of alive path  $\Pi$ , and denoted as  $\text{short}(\Pi)$ . All arcs of  $\text{supp}(f)$  (in  $N$ ) are themselves length-1 alive paths, so  $f$  defines an identical pseudoflow on the shortcuts in  $\tilde{N}$ . We abuse notation by using  $f$  to refer to both the pseudoflow in  $N$  as well as  $\tilde{N}$ . With this, any excess-deficit path in  $N_f$  can be mapped one-to-one to an excess-deficit path in  $\tilde{N}_f$  of the same cost, by mapping each alive path in  $N$  to its shortcut in  $\tilde{N}$ . In fact, if we define potentials on the alive nodes, both paths will have the same reduced cost (excess/deficit nodes are always alive).

As a result, a Hungarian search performing arc-by-arc relaxations on  $\tilde{N}$  is equivalent to our process of relaxing alive paths in  $N$ . For this to be helpful, the potential constructed on  $\tilde{N}$  must produce a useful potential for  $N$ . The next lemma gives one such construction.

**Lemma 3.7.** *For any  $\theta$ -optimal potential  $\tilde{\pi}$  defined on  $\tilde{N}$ , let  $\pi$  be potential for  $N$  constructed by extending  $\tilde{\pi}$  to dead nodes as follows:*

$$\pi(v) \leftarrow \begin{cases} \tilde{\pi}(v) & v \in A^* \cup B^* \\ \tilde{\pi}(s) & v \in A \setminus A^* \\ \tilde{\pi}(t) & v \in B \setminus B^* \end{cases}$$

Then,  $\pi$  is  $\theta$ -optimal for  $N$  and any admissible shortcut  $\text{short}(\Pi)$  corresponds to an admissible alive path  $\Pi$ .

*Proof.* We emphasize that the final statement is about arc-wise admissibility of  $\Pi$ , and not its weak admissibility. Indeed, even without extending  $\tilde{\pi}$  to dead nodes, an admissible  $\text{short}(\Pi)$  implies weak admissibility of  $\Pi$  under  $\tilde{\pi}$ .

To prove the  $\theta$ -optimality of  $\pi$ , we need to check reduced costs of length-2 and -3 alive paths, since it is immediate for length-1. The analysis is nearly identical for both types of length-2 paths and length-3, so we demonstrate it for alive paths of the form  $s \rightarrow v \rightarrow b$  only. For a length-2 alive path of the form  $\Pi = s \rightarrow v \rightarrow b$ , for  $b \in B^*$ , we have  $c_\pi(s \rightarrow v) = 0$  so we turn our eye to the bipartite arc  $v \rightarrow b$ . Observe that the reduced cost of  $v \rightarrow b$ , under  $\pi$ , is equal to the reduced cost of  $\text{short}(\Pi)$  under  $\tilde{\pi}$ .

$$\begin{aligned} c_\pi(v \rightarrow b) &= c(v \rightarrow b) - \pi(v) + \pi(b) \\ &= c(\Pi) - \pi(s) + \pi(b) \\ &= c(\text{short}(\Pi)) - \tilde{\pi}(s) + \tilde{\pi}(b) \\ &= c_{\tilde{\pi}}(\text{short}(\Pi)) \end{aligned} \tag{1}$$

It follows that  $v \rightarrow b$  is  $\theta$ -optimal under  $\pi$ , by  $\theta$ -optimality of  $\text{short}(\Pi)$  under  $\tilde{\pi}$ .

$$c_\pi(v \rightarrow b) = c_{\tilde{\pi}}(\text{short}(\Pi)) \geq -\theta$$

For the second statement, suppose a shortcut  $\text{short}(\Pi)$  is admissible. Once again, the proof for each length is nearly identical, so we demonstrate the proof only for  $\Pi = s \rightarrow v \rightarrow b$ . The non-bipartite arc of  $\Pi$  ( $s \rightarrow v$  with dead  $v$ ) has reduced cost 0 under  $\pi$ ; it is admissible. By (1), admissibility of  $v \rightarrow b$  under  $\pi$  follows from admissibility of  $\text{short}(\Pi)$  under  $\tilde{\pi}$ .

$$c_\pi(v \rightarrow b) = c_{\tilde{\pi}}(\text{short}(\Pi)) \leq 0$$

Thus, all arcs of  $\Pi$  are admissible, and  $\Pi$  is admissible.  $\square$

**Updating potential.** The query data structures for min-reduced-cost alive path have no dependency on dead node potential; we leave them untracked as described before. By Lemma 3.7, a blocking flow supported on weakly-admissible alive paths is one that is arc-wise admissible — with the right extension of potential to dead nodes. Indeed, we require values of  $\pi$  on all nodes at the end of a scale (for the next SCALE-INIT) and for individual dead nodes whenever they become alive (after augmentation). We use the construction in Lemma 3.7 to reconstruct a potential for dead nodes in these situations. Potential updates for alive nodes can be handled in a batched fashion as in Section 2.

**Implementing the augmentation stage.** «TODO check the writing of this part» Augmentation can also be implemented in  $O(k \text{ polylog } n)$  time, after  $O(n \text{ polylog } n)$  time preprocessing, using similar data structures. The augmentation step finds an admissible blocking flow as a sequence of disjoint admissible augmenting paths, using a depth-first search (DFS) through the admissible residual arcs. The DFS is initialized at each excess node and continues until it reaches a deficit node with remaining deficit (an augmenting path is found), or returns to the excess node (no path was found). We attempt DFS from a new excess node once no path is found, or we have found enough augmenting paths to balance the current excess. In each step of the DFS, we discover admissible residual arcs by querying the min-reduced-cost arc leaving the current node  $x \in X$  into  $V \setminus X$ , where  $X$  is the set of excess nodes plus nodes already visited by the DFS. If this minimum arc

is admissible, we can search across it. If it is not admissible, then there are no more unexplored admissible arcs leaving  $x$ . Like the Hungarian search, the DFS is reduced to a problem of querying the min-reduced-cost arc.

Applying Lemma 3.7, we can DFS through weakly-admissible alive paths instead of individual admissible arcs. In implementation, we use data structures for dynamic *nearest neighbor* (NN) in addition to BCP. The nearest neighbor problem is as follows: preprocess a point set  $P$  to answer, given a query point  $a$ , the point  $b \in P$  minimizing  $c(a, b) - \omega(a) + \omega(b)$ . Kaplan *et al.* [12] solve the 2D dynamic NN problem with  $O(\text{polylog } n)$  update and  $O(\log^2 n)$  query time. Similar to the Hungarian search implementation, we search across the length-1, -2, and -3 alive paths that start at  $x$ . Alive paths of the form  $s \rightarrow v \rightarrow b$  and  $s \rightarrow v \rightarrow w \rightarrow t$  are still queried using BCPs, while those of the form  $a \rightarrow b$  and  $a \rightarrow v \rightarrow t$  use an NN query with query point  $x = a$  instead. The total number of min-heaps and NN data structures is  $O(1)$ , and therefore rewinding for the augmentation data structures can be done in  $O(k \text{ polylog } n)$  time given Lemma 3.8.

**Bounding relaxations and blocking flow support size.** The following lemma is crucial to the analysis of running time for the Hungarian search and augmentation, bounding both the number of relaxations and potential update/recovery operations.

**Lemma 3.8.** *Both Hungarian search and augmentation stages explore  $O(k)$  nodes, and the blocking flow found in augmentation stage is incident to  $O(k)$  nodes.*

*Proof.* Both processes explore by adding  $O(1)$  nodes to  $X$  in each step, including at least one alive node. Since the number of alive nodes is  $O(k)$ , the number of explored nodes is  $O(k)$ .

For the second statement, recall that our blocking flow is a union of path flows along admissible augmenting paths. Each path flow is composed of alive paths explored by the depth-first search in the augmentation stage, of which there are  $O(k)$ . Since each such alive path has  $O(1)$  length, the number of incident nodes is  $O(k)$ .  $\square$

We thus obtain the following:

**Lemma 3.9.** *After  $O(n \text{ polylog } n)$  time preprocessing, each iteration of REFINE can be implemented in  $O(k \text{ polylog } n)$  time.*

## 4 Transportation algorithm

Given two point sets  $A$  and  $B$  in  $\mathbb{R}^2$  of sizes  $r$  and  $n$  respectively and a supply-demand function  $\phi : A \cup B \rightarrow \mathbb{Z}$  as defined in the introduction, we present an  $O(rn^{3/2} \text{ polylog } n)$  time algorithm for computing an optimal transport map between  $A$  and  $B$ . By applying this algorithm in the case of  $r \leq \sqrt{n}$  and the one by Agarwal *et al.* [3] when  $r > \sqrt{n}$ , we prove Theorem 1.3. We use a standard reduction to the uncapacitated min-cost flow problem and use Orlin’s algorithm [17] as well as some of the ideas from Agarwal *et al.* [3] for efficient implementation under the geometric settings. We first present an overview of the algorithm and then describe its fast implementation that achieves the desired running time.

**From transportation to circulation** Given a transportation instance between point sets  $A$  and  $B$ , we generate a min-cost flow network  $N$  as follows: add an arc  $a \rightarrow b$  for each  $a \in A$  and  $b \in B$ , and copy the transportation supply/demand function to use as the MCF supply/demand function  $\phi$ . A circulation  $f$  in  $N$  directly corresponds to a transportation map  $\tau$  of the same cost, by setting  $\tau(v, w) = f(v \rightarrow w)$  for all arcs.

## 4.1 Overview of the algorithm

Orlin's algorithm follows an excess-scaling paradigm and the primal-dual framework. It maintains a *scale parameter*  $\Delta$ , a flow function  $f$ , and potential  $\pi$  on the nodes. Initially  $\Delta$  is equal to the total supply,  $f = 0$ , and  $\pi = 0$ . We fix a constant parameter  $\alpha \in (0.5, 1)$ . A node  $v$  is called *active* if the magnitude of imbalance of  $v$  is at least  $\alpha\Delta$ . At each step, using the Hungarian search, the algorithm finds an admissible excess-to-deficit path between active nodes in the residual network and pushes a flow of amount  $\Delta$  along this path.<sup>3</sup> Repeat the process until either active excess or deficit nodes are gone; when this happens,  $\Delta$  is halved. The sequence of augmentations with a fixed value of  $\Delta$  is called an *excess scale*.

The algorithm also performs two preprocessing steps at the beginning of each excess scale. First, if  $f(v \rightarrow w) \geq 3n\Delta$ ,  $v \rightarrow w$  is contracted to a single node  $z$  with  $\phi(z) = \phi(v) + \phi(w)$ .<sup>4</sup> Second, if there are no active excess nodes and  $f(v \rightarrow w) = 0$  for every arc  $v \rightarrow w$ , then  $\Delta$  is aggressively lowered to  $\max_v \phi(v)$ .

When the algorithm terminates, an optimal circulation in the contracted network is found. **«TODO include a description of the contraction details»** We use the algorithm described in Agarwal *et al.* [3] to recover an optimal circulation for the original network in  $O(n \text{ polylog } n)$  time. Orlin showed that the algorithm terminates within  $O(n \log n)$  scales and performs a total of  $O(n \log n)$  augmentations. In the next subsection, we describe an algorithm that, after  $O(n \text{ polylog } n)$  time preprocessing, finds an admissible excess-to-deficit path in  $O(r\sqrt{n} \text{ polylog } n)$  amortized time. Summing this cost over all augmentations, we obtain the desired running time.

## 4.2 An efficient implementation

Recall in the previous sections that we could bound the running time of the Hungarian search by the size of  $\text{supp}(f)$ . Here, the number of active imbalanced nodes at any scale is  $O(r)$ , and the length of an augmenting path is also  $O(r)$ . Therefore one might hope to find an augmenting path in  $O(r \text{ polylog } n)$  time, by adapting the algorithms described in Sections 2 and 3. The challenge is that  $\text{supp}(f)$  may have  $\Omega(n)$  size, therefore an algorithm which runs in time proportional to the support size is no longer sufficient. Still, we manage to implement Hungarian search in time  $O(r\sqrt{n} \text{ polylog } n)$ , by exploiting a few properties of  $\text{supp}(f)$  as described below.

We note that each arc of  $\text{supp}(f)$  is admissible with reduced cost 0, so we prioritize the relaxation of support arcs as soon as they arrive in  $X \times (V \setminus X)$ , over the non-support arcs. This strategy ensures the following crucial property.

**«TODO reference the lemma/proof from arxiv instead»**

**Lemma 4.1.** *If the support arcs  $\text{supp}(f)$  are relaxed as soon as possible,  $\text{supp}(f)$  is acyclic.*

Next, similar to Section 3, we call node  $u$  *alive* if (a)  $u$  is an active imbalanced node or (b) if  $u$  is incident to an arc of  $\text{supp}(f)$ ;  $u$  is *dead* otherwise. Unlike in Section 3, once a node becomes alive it cannot be dead again. Furthermore, a dead node may become alive only at the beginning of a scale (after the value of  $\Delta$  is reduced). Also, an augmenting path cannot pass through a dead node. Therefore, we can ignore all dead nodes during Hungarian search, and update the set of alive/dead nodes at the beginning of a scale.

Let  $B^* \subseteq B^\circ$  be the set of nodes that are either (a) active imbalanced nodes or (b) incident to *exactly one* arc of  $\text{supp}(f)$ . Lemma 4.1 implies that  $B^\circ \setminus B^*$  has size  $O(r)$ . We can therefore find the min-reduced-cost arc between  $X \cap A^\circ$  and  $B^\circ \setminus (B^* \cup X)$  using a BCP data structure as in

<sup>3</sup>Note that this augmentation may convert an excess node into a deficit node.

<sup>4</sup>Intuitively,  $f(v \rightarrow w)$  is so high that future scales cannot deplete the flow on  $v \rightarrow w$ .



Section 2, along with lazy potential updates and the rewinding mechanism. The total time spent by Hungarian search on the nodes of  $B^* \setminus B^*$  will be  $O(r \text{ polylog } n)$ . We subsequently focus on handling  $B^*$ .

**Handling  $B^*$ .** We now describe how we query a min-reduced-cost arc between  $X \cap A^*$  and  $B^* \setminus X$ . Each node  $b \in B^*$  is incident to exactly one arc in  $\text{supp}(f)$ . We partition these nodes into clusters depending on their unique neighbor in  $N_f$ . That is, for a node  $a \in A^*$ , let  $B_a^* := \{b \in B^* \mid a \rightarrow b \in \text{supp}(f)\}$ . We refer to  $B_a^*$  as the *star* of  $a$ .

The crucial observation is that  $a$  is the only node in  $N_f$  reachable from each  $b \in B_a^*$ , so once the Hungarian search reaches a node of  $B_a^*$  and thus  $a$  (recall we prioritize relaxing support arcs), the Hungarian search need not visit any other nodes of  $B_a^*$ , as they will only lead to  $a$ . Hence, as soon as one node of  $B_a^*$  is reached, all other nodes of  $B_a^*$  can be discarded from further consideration. Using this observation, we handle  $B^*$  as follows.

We classify each  $a \in A^*$  as *light* or *heavy*: heavy if  $|B_a^*| \geq \sqrt{n}$ , and light if  $|B_a^*| \leq 2\sqrt{n}$ . Note that if  $\sqrt{n} \leq |B_a^*| \leq 2\sqrt{n}$  then  $a$  may be classified as light or heavy. We allow this flexibility to implement reclassification in a lazy manner. Namely, a light node is reclassified as heavy once  $|B_a^*| > 2\sqrt{n}$ , and a heavy node is reclassified as light once  $|B_a^*| < \sqrt{n}$ . This scheme ensures that the star of  $a$  has gone through at least  $\sqrt{n}$  updates between two successive reclassifications, and these updates will pay for the time spent in updating the data structure when  $a$  is re-classified.

For each heavy node  $a \in A^* \setminus X$ , we maintain a BCP data structure between  $B_a^*$  and  $X \cap A^*$ . Next, for all light nodes in  $A^* \setminus X$ , we collect their stars into a single set  $B_{<}^* := \bigcup_{a \text{ light}} B_a^*$ . We maintain one single BCP data structure between  $B_{<}^*$  and  $A^* \cap X$ . Thus, at most  $r$  different BCP data structures are maintained for stars.

Using these data structures, we can compute and relax a min-reduced-cost arc  $v \rightarrow w$  between  $A^* \cap X$  and  $B^* \setminus X$ . If  $w$  lies in some star  $B_a^*$ , then we also add  $a$  into  $X$ . If  $a$  is light, then we delete  $B_a^*$  from  $B_{<}^*$  and update the BCP data structure of  $B_{<}^*$ . If  $a$  is heavy, then we stop querying the BCP data structure of  $B_a^*$  for the remainder of the search. Finally, since  $a$  becomes part of  $X$ ,  $a$  is added to all  $O(r)$  BCP data structures. Recall that  $r \leq \sqrt{n}$  by assumption. Adding arc  $v \rightarrow w$  thus involves performing  $O(\sqrt{n})$  insertion/deletion operations in various BCP data structures, thereby taking  $O(\sqrt{n} \text{ polylog } n)$  time.

**Putting it together.** The following lemma bounds the running time of the Hungarian search.

**Lemma 4.2.** *Assuming all BCP data structures are initialized correctly, the Hungarian search terminates within  $O(r)$  steps, and takes  $O(r\sqrt{n} \text{ polylog } n)$  time.*

*Proof.* **«TODO»** □

Once an augmenting path is found and the augmentation is performed, the set of imbalanced nodes and the support arcs change. We thus need to update the sets  $B^*$ ,  $B_a^*$ s, and  $B_{<}^*$ . This can be accomplished in  $O(r \text{ polylog } n)$  amortized time. When we begin a new Hungarian search, we use the rewinding mechanism to set various BCP data structures in the right initial state. Finally, when we move from one scale to another, we also update the sets  $A^*$  and  $B^*$ . **«TODO details?»** Omitting all the details, we conclude the following.

**Lemma 4.3.** *Each Hungarian search can be performed in  $O(r\sqrt{n} \text{ polylog } n)$  time.*

Since there are  $O(n \log n)$  augmentations and the flow in the original network can be recovered from that in the contracted network in  $O(n \text{ polylog } n)$  time [3], the total running time of the algorithm is  $O(rn^{3/2} \text{ polylog } n)$ , as claimed in Theorem 1.3.

**Acknowledgment.** We thank Haim Kaplan for useful discussion and suggesting to use Goldberg *et al.* [10] for our approximation algorithm.

## References

- [1] Pankaj K. Agarwal, Alon Efrat, and Micha Sharir. Vertical decomposition of shallow levels in 3-dimensional arrangements and its applications. *SIAM J. Comput.* 29(3):912–953, 1999. <https://doi.org/10.1137/S0097539795295936>.
- [2] Pankaj K. Agarwal, Kyle Fox, Debmalya Panigrahi, Kasturi R. Varadarajan, and Allen Xiao. Faster algorithms for the geometric transportation problem. *Proc. 33rd Int. Sympos. Comput. Geom. (SoCG)*, 7:1–7:16, 2017. <https://doi.org/10.4230/LIPIcs.SoCG.2017.7>.
- [3] Pankaj K. Agarwal, Kyle Fox, Debmalya Panigrahi, Kasturi R. Varadarajan, and Allen Xiao. Faster algorithms for the geometric transportation problem. Preprint, 2019. <http://arxiv.org/abs/1903.08263>.
- [4] Pankaj K. Agarwal and Kasturi R. Varadarajan. A near-linear constant-factor approximation for Euclidean bipartite matching? *Proc. 20th Annu. Sympos. Comput. Geom. (SoCG)*, 247–252, 2004. <https://doi.org/10.1145/997817.997856>.
- [5] D. Bertsekas and D. El Baz. Distributed asynchronous relaxation methods for convex network flow problems. *SIAM J. Control and Opt.* 25(1):74–85, 1987. <https://doi.org/10.1137/0325006>.
- [6] L. Paul Chew and Robert L. (Scot) Drysdale III. Voronoi diagrams based on convex distance functions. *Proc. 1st Annu. Sympos. Comput. Geom. (SoCG)*, 235–244, 1985. <https://doi.org/10.1145/323233.323264>.
- [7] Ran Duan, Seth Pettie, and Hsin-Hao Su. Scaling algorithms for weighted matching in general graphs. *ACM Trans. Algorithms* 14(1):8:1–8:35, 2018. <https://doi.org/10.1145/3155301>.
- [8] Shimon Even and Robert E. Tarjan. Network flow and testing graph connectivity. *SIAM J. Comput.* 4(4):507–518, 1975. <https://doi.org/10.1137/0204043>.
- [9] Harold N. Gabow and Robert E. Tarjan. Faster scaling algorithms for network problems. *SIAM J. Comput.* 18(5):1013–1036, 1989. <https://doi.org/10.1137/0218069>.
- [10] Andrew V. Goldberg, Sagi Hed, Haim Kaplan, and Robert E. Tarjan. Minimum-cost flows in unit-capacity networks. *Theoret. Comput. Sci.* 61(4):987–1010, 2017. <https://doi.org/10.1007/s00224-017-9776-7>.
- [11] John E. Hopcroft and Richard M. Karp. An  $n^{5/2}$  algorithm for maximum matchings in bipartite graphs. *SIAM J. Comput.* 2(4):225–231, 1973. <https://doi.org/10.1137/0202019>.
- [12] Haim Kaplan, Wolfgang Mulzer, Liam Roditty, Paul Seiferth, and Micha Sharir. Dynamic planar Voronoi diagrams for general distance functions and their algorithmic applications. *Proc. 28th Annu. ACM-SIAM Sympos. Discrete Algorithms (SODA)*, 2495–2504, 2017. <https://doi.org/10.1137/1.9781611974782.165>.
- [13] Andrey Boris Khesin, Aleksandar Nikolov, and Dmitry Paramonov. Preconditioning for the geometric transportation problem. Preprint, 2019. <http://arxiv.org/abs/1902.08384>.

- [14] Harold W. Kuhn. The Hungarian method for the assignment problem. *Naval Research Logistics (NRL)* 2(1-2):83–97. John Wiley & Sons, 1955.
- [15] Yin Tat Lee and Aaron Sidford. Path finding methods for linear programming: Solving linear programs in  $\tilde{O}(\sqrt{\text{rank}})$  iterations and faster algorithms for maximum flow. *55th Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, 424–433, 2014. [⟨https://doi.org/10.1109/FOCS.2014.52⟩](https://doi.org/10.1109/FOCS.2014.52).
- [16] Aleksander Mądry. Navigating central path with electrical flows: From flows to matchings, and back. *54th Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, 253–262, 2013. [⟨https://doi.org/10.1109/FOCS.2013.35⟩](https://doi.org/10.1109/FOCS.2013.35).
- [17] James B. Orlin. A faster strongly polynomial minimum cost flow algorithm. *Operations Research* 41(2):338–350, 1993. [⟨https://doi.org/10.1287/opre.41.2.338⟩](https://doi.org/10.1287/opre.41.2.338).
- [18] Lyle Ramshaw and Robert E. Tarjan. A weight-scaling algorithm for min-cost imperfect matchings in bipartite graphs. *Proc. 53rd Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, 581–590, 2012. [⟨https://doi.org/10.1109/FOCS.2012.9⟩](https://doi.org/10.1109/FOCS.2012.9).
- [19] R. Sharathkumar and Pankaj K. Agarwal. Algorithms for the transportation problem in geometric settings. *Proc. 23rd Annu. ACM-SIAM Sympos. Discrete Algorithms (SODA)*, 306–317, 2012. [⟨https://dl.acm.org/citation.cfm?id=2095116.2095145⟩](https://dl.acm.org/citation.cfm?id=2095116.2095145).
- [20] R. Sharathkumar and Pankaj K. Agarwal. A near-linear time  $\epsilon$ -approximation algorithm for geometric bipartite matching. *Proc. 44th Annu. ACM Sympos. Theory Comput. (STOC)*, 385–394, 2012. [⟨https://doi.org/10.1145/2213977.2214014⟩](https://doi.org/10.1145/2213977.2214014).
- [21] Éva Tardos. A strongly polynomial minimum cost circulation algorithm. *Combinatorica* 5(3):247–256, 1985. [⟨https://doi.org/10.1007/BF02579369⟩](https://doi.org/10.1007/BF02579369).
- [22] Pravin M. Vaidya. Geometry helps in matching. *SIAM J. Comput.* 18(6):1201–1225, 1989. [⟨https://doi.org/10.1137/0218080⟩](https://doi.org/10.1137/0218080).
- [23] Kasturi R. Varadarajan. A divide-and-conquer algorithm for min-cost perfect matching in the plane. *Proc. 39th Annu. IEEE Sympos. Found. Comput. Sci. (FOCS)*, 320–331, 1998. [⟨https://doi.org/10.1109/SFCS.1998.743466⟩](https://doi.org/10.1109/SFCS.1998.743466).