

Efficient Algorithms for Geometric Partial Matching*

Pankaj K. Agarwal[†]

Hsien-Chih Chang[†]

Allen Xiao[†]

March 22, 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*.¹

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 τ 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.

*Work on this paper was supported by NSF under grants CCF-15-13816, CCF-15-46392, and IIS-14-08846, by an ARO grant W911NF-15-1-0408, and by BSF Grant 2012/229 from the U.S.-Israel Binational Science Foundation.

[†]Department of Computer Science, Duke University. Email: {pankaj, hc252, axiao}@cs.duke.edu.

¹Partial matching is also called *imperfect matching* or *imperfect assignment* [?, ?].

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 [?] and the $O(m \min\{\sqrt{m}, n^{2/3}\})$ time algorithm by Even and Tarjan [?], 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 [?], which uses an interior-point algorithm, runs in $O(m^{10/7} \text{polylog } n)$ time.

The Hungarian algorithm [?] 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 [?] and the improved $O(m\sqrt{n} \log C)$ time algorithm by Duan *et al.* [?] assuming the edge costs are integral; here C is the maximum cost of an edge. Ramshaw and Tarjan [?] 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.* [?] 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 $O(n^2 \text{polylog } n)$ [?] (see also [?, ?]). 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 [?, ?, ?, ?]. 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 [?] 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 [?]; 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 uncapacitated min-cost flow algorithm by Orlin [?] solves the transportation problem in $O((m + n \log n)n \log n)$ time. Lee and Sidford [?] 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.* [?, ?] 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 [?] 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.* [?] and Kaplan *et al.* [?]. 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.* [?].

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.* [?] for $r = o(\sqrt{n})$. Similar to their algorithm, we also use the strongly polynomial uncapacitated minimum-cost flow algorithm by Orlin [?], 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_ℓ) 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 Π . 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 π , we can define the *reduced cost* of the edges to be $c_\pi(v, w) := c(v, w) - \pi(v) + \pi(w)$. Potential π is *feasible* on edge (v, w) if $c_\pi(v, w)$ is nonnegative. Potential π is *feasible* if π are feasible on every edge in G . We say that an edge (v, w) is *admissible* under potential π 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 Π , discovered using a procedure called the *Hungarian search*. The algorithm

terminates after k augmentations, exactly when $|M| = k$; Ramshaw and Tarjan [?] 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_π on other edges within X , and keep π 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 [?, ?]. 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 π .

The dynamic BCP data structure by Kaplan *et al.* [?] 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 [?] 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 δ , which is initially 0 at the very beginning of the algorithm. Whenever we would raise the potential of all nodes in X , we raise δ 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 δ , and use that stored value as its BCP weight. Since the BCP weights are uniformly offset from $\pi(v)$ by δ , 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 δ . Unlike Vaidya [?], we do not reset δ across Hungarian searches for the sake of rewinding.

There are $O(k)$ relaxations and thus $O(k)$ updates to δ 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.* [?] 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 [?] 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 $O(k)$ excess, the size of $\text{supp}(f)$ is $O(k)$.*

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 π 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_π 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 π 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 π is *θ -optimal* [?, ?] if $c_\pi(v \rightarrow w) \geq -\theta$ for every residual arc $v \rightarrow w$ in \vec{E}_f . Pseudoflow f is *θ -optimal* if it is θ -optimal with respect to some potential π ; potential π is *θ -optimal* if it is θ -optimal with respect to some pseudoflow f . Given a pseudoflow f and potential π , 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.² We will use the following well-known property of θ -optimality.

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

Using Lemma 3.1, the following lemma can be proved about θ -optimality:

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

²The 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 $c_\pi(v \rightarrow w) = 0$.

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 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 $v \in A \cup B$.*

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.* [?]. The algorithm works in scales. In the beginning of each scale, we fix a *cost scaling parameter* θ and maintain potential π with the following property:

- (*) There exists a 2θ -optimal integral circulation in N with respect to π .

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 θ -optimal pseudoflow f . In the *refinement* stage, REFINES converts f into a θ -optimal (integral) circulation g . In both stages, π 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 begin with a 2θ -optimal circulation f in N . First, we raise the potential of all nodes in A by θ , all nodes in B by 2θ , and t by 3θ . 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 θ -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 θ -optimal pseudoflow with $O(k)$ excess into a θ -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 θ be the scaling parameter and π_0 the potential function at the beginning of a scale, such that there exists an integral 2θ -optimal circulation with respect to π_0 . Let f be a θ -optimal pseudoflow with excess $O(k)$. Then REFINE terminates within $O(\sqrt{k})$ iterations.*

Proof. We sketch the proof, which is adapted from Goldberg *et al.* [?]. Let f_0 be the assumed 2θ -optimal integral circulation with respect to π_0 , and let π 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 θ . 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 θ , 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 [?, Lemma 4]. The proof examines cuts $Y_i := \{v \mid d(v) > i\}$ for $0 \leq i \leq \sqrt{k}$. By θ -optimality of f and 2θ -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 3.8). 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. Due to lack of space, we only describe the former.

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 Π be the sum of c_π over its arcs. We say Π 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 $c(\Pi)$ and add every node of Π 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. Thus, we can query the minimum alive path using a partial assignment of π on only the alive nodes, leaving π 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 (additively weighted) bichromatic closest pair (BCP) data structure between $A^* \cap X$ and $B^* \setminus X$, with reduced cost as the weighted pair distance. Such BCP data structure can be implemented so that insertions/deletions can be performed in $O(\text{polylog } k)$ time [?].

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 3.7 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 3.7 again, there are $O(k)$ alive/dead updates caused by an augmentation, so the time for these updates per Hungarian search is $O(k \text{ polylog } n)$.

Updating potential. Potential updates for alive nodes can be handled in a batched fashion as in Section 2. The three data structures above have no dependency on the dead node potential; we leave them untracked as described before. The Hungarian search remains intact since alive nodes are visited in the same order as when using arc-by-arc relaxations. However, we need values of π 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 can reconstruct a “valid” potential in these situations. To recover potential for $v \in A \setminus A^*$ we set $\pi(v) \leftarrow \pi(s)$, and for $v \in B \setminus B^*$ we set $\pi(v) \leftarrow \pi(t)$. Straightforward calculation shows that such potential (1) preserves θ -optimality, and (2) makes Π (arc-wise) admissible for any weakly admissible alive path Π . Hence, a blocking flow composed of weakly admissible alive paths is admissible under the recovered potential.

The following lemma is crucial to the analysis of running time for the Hungarian search, bounding both the number of relaxations and potential update/recovery operations.

Lemma 3.7. *Both Hungarian search and augmentation stages explore $O(k)$ nodes, and the blocking flow found in augmentation stage is incident to $O(k)$ nodes.*

Augmentation can also be implemented in $O(k \text{ polylog } n)$ time, after $O(n \text{ polylog } n)$ time preprocessing, using similar data structures. We thus obtain the following:

Lemma 3.8. *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.* [?] 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 [?] as well as some of the ideas from Agarwal *et al.* [?] 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.

4.1 Overview of the algorithm

Orlin's algorithm follows an excess-scaling paradigm and the primal-dual framework. It maintains a *scale parameter* Δ , a flow function f , and potential π on the nodes. Initially Δ 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 Δ along this path.³ Repeat the process until either active excess or deficit nodes are gone; when this happens, Δ is halved. The sequence of augmentations with a fixed value of Δ 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 $\phi(z) = \phi(v) + \phi(w)$.⁴ Second, if there are no active excess nodes and $f(v \rightarrow w) = 0$ for every arc $v \rightarrow w$, then Δ is aggressively lowered to $\max_v \phi(v)$.

When the algorithm terminates, an optimal circulation in the contracted network is found. We use the algorithm described in Agarwal *et al.* [?] 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.

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 Δ 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 Section 2, along with lazy potential updates and the rewinding mechanism. The total time spent

³Note that this augmentation may convert an excess node into a deficit node.

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

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. While proof is omitted, 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.*

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^* . 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 [?], 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 discussion and suggestion to use Goldberg *et al.* [?] algorithm. We thank Debmalya Panigrahi for helpful discussions.

References