

1 **FAST SINKHORN II: COLLINEAR TRIANGULAR MATRIX AND**
2 **LINEAR TIME ACCURATE COMPUTATION OF OPTIMAL**
3 **TRANSPORT***

4 QICHEN LIAO[†], ZIHAO WANG[‡], JING CHEN[§], BO BAI[¶], SHI JIN^{||}, AND HAO WU[#]

5 **Abstract.** In our previous work [arXiv:2202.10042], the complexity of Sinkhorn iteration is
6 reduced from $O(N^2)$ to the optimal $O(N)$ by leveraging the special structure of the kernel matrix.
7 In this paper, we explore the special structure of kernel matrices by defining and utilizing the proper-
8 ties of the **Lower-ColLinear Triangular Matrix** (L-CoLT matrix) and **Upper-ColLinear Triangular**
9 **Matrix** (U-CoLT matrix). We prove that (1) L/U-CoLT matrix-vector multiplications can be carried
10 out in $O(N)$ operations; (2) both families of matrices are closed under the Hadamard product and
11 matrix scaling. These properties help to alleviate two key difficulties for reducing the complexity
12 of the Inexact Proximal point method (IPOT), and allow us to significantly reduce the number of
13 iterations to $O(N)$. This yields the Fast Sinkhorn II (FS-2) algorithm for accurate computation of
14 optimal transport with low algorithm complexity and fast convergence. Numerical experiments are
15 presented to show the effectiveness and efficiency of our approach.

16 **Key words.** Optimal Transport, Wasserstein-1 metric, Sinkhorn algorithm, IPOT method,
17 FS-2 algorithm

18 **MSC codes.** 49M25; 65K10

19 **1. Introduction.** The Wasserstein metric, broadly used in optimal transport
20 theory with applications in many fields including machine learning, quantifies the
21 dissimilarity between two probabilistic distributions. Many methods have been pro-
22 posed to compute the Wasserstein metrics directly, such as the linear programming
23 methods [30, 22, 36], combinatorial methods [33], solving the Monge-Ampère equa-
24 tions [15, 14, 3], via Benamou-Brenier formulation [2, 21] and the proximal splitting
25 methods [8, 28]. In recent years, several approximation techniques in optimal trans-
26 port for high-dimensional distributions have also been proposed [26, 27].

27 The Sinkhorn algorithm [10, 34] is a popular $O(N^2)$ algorithm to approximate
28 the Wasserstein metric [31] by minimizing the entropy regularized optimal transport
29 (OT) problem. In [24], the FS-1 algorithm is proposed to solve entropy regularized
30 OT in $O(N)$ time by leveraging the special structure of the Sinkhorn kernel matrix
31 of the Wasserstein-1 metric. The solution of entropy regularized OT approximates
32 the accurate OT solution only if the regularization parameter is sufficiently small.
33 However, small regularization parameters lead to numerical instability and excessive

*Submitted to the editors DATE.

Funding: This work was supported by National Natural Science Foundation of People’s Republic of China Grant Nos.11871297 and 12031013, Tsinghua University Initiative Scientific Research Program, and Shanghai Municipal Science and Technology Major Project 2021SHZDZX0102.

[†]Department of Mathematical Sciences, Tsinghua University, Beijing 100084, People’s Republic of China; Theory Lab, Central Research Institute, 2012 Labs, Huawei Technologies Co. Ltd., Hong Kong SAR, People’s Republic of China (lqc20@mails.tsinghua.edu.cn).

[‡]Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong SAR, People’s Republic of China (zwanggc@cse.ust.hk).

[§]School of Physical and Mathematical Sciences, Nanyang Technological University, Singapore 639798 (jing.chen@ntu.edu.sg).

[¶]Theory Lab, Central Research Institute, 2012 Labs, Huawei Technologies Co. Ltd., Hong Kong SAR, People’s Republic of China (baibo8@huawei.com).

^{||}School of Mathematical Sciences, Institute of Natural Sciences, and MOE-LSC Shanghai Jiao Tong University, Shanghai 200240, People’s Republic of China (shijin-m@sjtu.edu.cn).

[#]Corresponding author. Department of Mathematical Sciences, Tsinghua University, Beijing 100084, People’s Republic of China (hwu@tsinghua.edu.cn).

34 iterations [13]. This causes the slow convergence of the Sinkhorn algorithm.

35 The Inexact Proximal point method [35] for the Optimal Transport problem
 36 (IPOT) has been proposed to address this challenge. It regularizes the original OT
 37 by introducing the proximal point term and solves a series of successive subproblems.
 38 Only fairly mild regularization parameters are required to ensure the method's fast
 39 convergence to the accurate OT solution in an $O(N^2)$ algorithm. The goal of this
 40 paper is to construct a new method to accurately compute OT solutions with good
 41 convergence behavior and $O(N)$ algorithm complexity by combing the IPOT method
 42 and the FS-1 algorithm. Note the two key steps in the IPOT method make it hard to
 43 reduce the complexity to $O(N)$: the matrix Hadamard product (Algorithm 2.1, line 4)
 44 and the matrix scaling (Algorithm 2.1, line 8). For general matrices, the complexity
 45 of the above operations is both $O(N^2)$. Moreover, these operations may destroy the
 46 special structure of the kernel matrix [24], making it impossible for us to implement
 47 matrix-vector multiplication with $O(N)$ cost.

48 We will explore the special structure of kernel matrices by defining and exploit-
 49 ing the properties of the **Lower-CoLinear Triangular Matrix** (L-CoLT matrix) and
 50 **Upper-CoLinear Triangular Matrix** (U-CoLT matrix). For these matrices, we can
 51 realize the matrix-vector multiplication with $O(N)$ cost by using the idea of dynamic
 52 programming similar to [24]. Next, we show that each L/U-CoLT matrix can be rep-
 53 resented by two vectors of dimension N . Furthermore, we prove the closure of families
 54 of L/U-CoLT matrices to matrix Hadamard product and matrix scaling. This means
 55 that the special structure of the kernel matrix is preserved by matrix Hadamard
 56 product and matrix scaling, so we can still implement matrix-vector multiplication
 57 (Algorithm 2.1, lines 6-7) with $O(N)$ cost. On the other hand, by updating two
 58 representation vectors of the L/U-CoLT matrix, we can also implement matrix Ha-
 59 damard product (Algorithm 2.1, line 4) and matrix scaling (Algorithm 2.1, line 8)
 60 with $O(N)$ cost. Consequently, the Fast Sinkhorn II (FS-2) algorithm is developed,
 61 which integrates the advantages of both IPOT and FS-1. Moreover, we also find that
 62 the FS-2 algorithm has the advantage in reducing the space complexity since all the
 63 matrices can be represented by vectors. Due to these benefits, one can expect that our
 64 FS-2 could be applied in various fields, e.g., machine learning [26, 27, 16, 25], image
 65 processing [32, 29], inverse problems [6, 12, 37, 18], density function theory[19, 5, 9].

66 The rest of the paper is organized as follows. In section 2, the basics of the
 67 Wasserstein-1 metric and the IPOT method are briefly reviewed. After presenting
 68 the definition, properties, and fast matrix-vector multiplications of the L/U-CoLT
 69 matrix in section 3, we apply them to accelerate the IPOT method, thus developing
 70 the FS-2 algorithm in section 4. In section 5, the FS-2 algorithm is extended to high
 71 dimensions. The numerical experiments are performed to verify our conclusions in
 72 section 6. We conclude the paper in section 7.

73 **2. The Wasserstein-1 metric and the IPOT method.** Given two unit dis-
 74 crete distributions \mathbf{u} and \mathbf{v} ,

$$75 \quad \mathbf{u} = (u_1, u_2, \dots, u_N)^\top \in \mathbb{R}^N, \quad \mathbf{v} = (v_1, v_2, \dots, v_N)^\top \in \mathbb{R}^N,$$

76 where $u_i \geq 0$, $v_j \geq 0$, and $\sum_i u_i = \sum_j v_j = 1$. The Wasserstein-1 distance between
 77 them is defined as [31]

$$78 \quad (2.1) \quad W_1(\mathbf{u}, \mathbf{v}) = \min_{\Gamma \mathbf{1} = \mathbf{u}, \Gamma^\top \mathbf{1} = \mathbf{v}} \langle C, \Gamma \rangle,$$

79 where $C = [c_{ij}] \in \mathbb{R}^{N \times N}$ is the cost matrix. The element $c_{ij} = \|\mathbf{x}_i - \mathbf{y}_j\|_1$ represents
 80 the cost of transporting the unit mass from position \mathbf{x}_i to position \mathbf{y}_j and the variable
 81 $\Gamma = [\gamma_{ij}] \in \mathbb{R}^{N \times N}$ to be optimized is the transport plan. Here, the Frobenius inner
 82 product $\langle A, B \rangle = \sum_{i,j} a_{ij} b_{ij}$, where $A = [a_{ij}]$, $B = [b_{ij}]$ are real-valued matrices.

83 The Sinkhorn algorithm [10, 34] solves an entropy regularized OT problem to
 84 obtain an approximate result of (2.1). However, the small regular parameter required
 85 by good approximation leads to a slow convergence rate and numerical instability. To
 86 avoid this problem, the proximal point iteration (2.2) is developed to solve (2.1) ac-
 87 curately [35]. It begins with a transport map $\Gamma^{(0)}$ and iteratively solves the following
 88 minimization problem

$$89 \quad (2.2) \quad \Gamma^{(t+1)} = \arg \min_{\Gamma \mathbf{1} = \mathbf{u}, \Gamma^T \mathbf{1} = \mathbf{v}} \langle C, \Gamma \rangle + \delta^{(t)} D_h \left(\Gamma, \Gamma^{(t)} \right),$$

90 where D_h is Bregman divergence, taken in the form of the KL divergence in [35],

$$91 \quad D_h(A, B) = \sum_{i,j} \left(a_{ij} \ln \frac{a_{ij}}{b_{ij}} - a_{ij} + b_{ij} \right)$$

92 and $\delta^{(t)}$ is the regular parameter. The Lagrangian of the above equation writes

$$93 \quad L(\Gamma, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \langle C, \Gamma \rangle + \delta^{(t)} D_h \left(\Gamma, \Gamma^{(t)} \right) + \boldsymbol{\alpha}^T (\Gamma \mathbf{1} - \mathbf{u}) + \boldsymbol{\beta}^T (\Gamma^T \mathbf{1} - \mathbf{v}).$$

94 Taking derivative of the Lagrangian with respect of γ_{ij} directly leads to

$$95 \quad \gamma_{ij} = e^{-\alpha_i / \delta^{(t)}} Q_{ij}^{(t)} e^{-\beta_j / \delta^{(t)}}, \quad \text{where } Q_{ij}^{(t)} = \gamma_{ij}^{(t)} e^{-c_{ij} / \delta^{(t)}} > 0.$$

96 Denoting \odot as the Hadamard product, $Q^{(t)} = K \odot \Gamma^{(t)}$ and $K = [e^{-c_{ij} / \delta^{(t)}}] \in \mathbb{R}^{N \times N}$
 97 is the kernel matrix. Letting $\phi_i = e^{-\alpha_i / \delta^{(t)}}$, $\psi_j = e^{-\beta_j / \delta^{(t)}}$, and vectors $\boldsymbol{\phi} = (\phi_i)$ and
 98 $\boldsymbol{\psi} = (\psi_j)$, one obtains

$$99 \quad (2.3) \quad \text{diag}(\boldsymbol{\phi}) Q^{(t)} \text{diag}(\boldsymbol{\psi}) \mathbf{1} = \mathbf{u}, \quad \text{diag}(\boldsymbol{\psi}) Q^{(t)\top} \text{diag}(\boldsymbol{\phi}) \mathbf{1} = \mathbf{v}.$$

100 By iteratively updating vectors $\boldsymbol{\phi}$ and $\boldsymbol{\psi}$

$$101 \quad (2.4) \quad \boldsymbol{\psi}^{(t, \ell+1)} = \mathbf{v} \oslash (Q^{(t)\top} \boldsymbol{\phi}^{(t, \ell)}), \quad \boldsymbol{\phi}^{(t, \ell+1)} = \mathbf{u} \oslash (Q^{(t)} \boldsymbol{\psi}^{(t, \ell+1)}),$$

102 one can obtain an accurate solution for the original OT problem (2.1). Here \oslash repre-
 103 sents pointwise division, t is the proximal iteration step (outer iteration) and ℓ it the
 104 Sinkhorn-type iteration step (inner iteration). The pseudo-code of IPOT is shown in
 105 Algorithm 2.1.

106 3. The Collinear Triangular Matrix.

107 3.1. Definition and Fast Matrix-Vector Multiplication.

108 DEFINITION 3.1 (Lower/Upper-Collinear Triangular Matrix). *A lower triangu-*
 109 *lar matrix is called a **Lower-Collinear Triangular Matrix** (L-CoLT matrix) if*
 110 *its corresponding entries on any two rows (columns) have the same-column (row)*
 111 *independent-ratio except those dividing by 0. Specifically, the N -dimensional L-CoLT*
 112 *matrix set is defined as follows:*

$$113 \quad (3.1) \quad \mathcal{C}_L^N = \left\{ M \in \mathbb{R}^{N \times N} \mid m_{i+1,j} / m_{i,j} = r_i, j \leq i; m_{i,j} = 0, i < j, \mathbf{r} \in (\mathbb{R} \setminus \{0\})^{N-1} \right\}.$$

Algorithm 2.1 IPOT**Input:** $\mathbf{u}, \mathbf{v} \in \mathbb{R}^N$; $K = e^{-C/\delta} \in \mathbb{R}^{N \times N}$; $L, \text{itr_max} \in \mathbb{N}^+$ **Output:** $W_1(\mathbf{u}, \mathbf{v})$

```

1:  $\phi, \psi \leftarrow \frac{1}{N} \mathbf{1}_N$ 
2:  $\Gamma = \mathbf{1}_N \mathbf{1}_N^T$ 
3: for  $t = 1 : \text{itr\_max}$  do
4:    $Q \leftarrow K \odot \Gamma$ 
5:   for  $\ell = 1 : L$  do
6:      $\psi \leftarrow \mathbf{v} \odot (Q^T \phi)$ 
7:      $\phi \leftarrow \mathbf{u} \odot (Q\psi)$ 
8:    $\Gamma \leftarrow \text{diag}(\phi) Q \text{diag}(\psi)$ 
return  $W_1(\mathbf{u}, \mathbf{v})$ 

```

114 Similarly, we define **Upper-Collinear Triangular Matrix** (*U-CoLT matrix*), which
 115 is a strictly upper triangular matrix:

(3.2)

$$116 \quad \mathcal{C}_U^N = \left\{ M \in \mathbb{R}^{N \times N} \mid m_{i-1,j}/m_{i,j} = r'_{i-1}, i < j; m_{i,j} = 0, i \geq j, \mathbf{r}' \in (\mathbb{R} \setminus \{0\})^{N-2} \right\}.$$

117 We call the vectors r and r' in (3.1)-(3.2) the **ratio vectors** of the collinear triangular
 118 matrix.

119 The matrices introduced in Definition 3.1 are termed as collinear triangular
 120 matrices (CoLT), due to the following collinearity between columns:

$$121 \quad m_{i,j}/m_{i,j+1} = m_{k,j}/m_{k,j+1} \iff m_{i,j}/m_{k,j} = m_{i,j+1}/m_{k,j+1}.$$

122 **THEOREM 3.2** (Vector Representation of Collinear Triangular Matrix). *Any L-*
 123 *CoLT matrix M_L can be represented by its diagonal elements γ and the ratio vector*
 124 *\mathbf{r} in Equation (3.1). Any U-CoLT matrix M_U can be represented by its superdiagonal*
 125 *elements γ' and the ratio vector \mathbf{r}' in (3.2).*

126 *Proof.* For any L-CoLT matrix $M_L \in \mathcal{C}_L^N$, if its corresponding γ and \mathbf{r} are given,
 127 then $m_{i,j} = \gamma_j \prod_{k=i}^{N-1} r_k$. The proof of U-CoLT is similar. \square

128 In the following, we use L-CoLT(γ, \mathbf{r}), $\gamma \in \mathbb{R}^N, \mathbf{r} \in \mathbb{R}^{N-1}$ and U-CoLT(γ', \mathbf{r}'),
 129 $\gamma' \in \mathbb{R}^{N-1}, \mathbf{r}' \in \mathbb{R}^{N-2}$ to represent a L-CoLT matrix and a U-CoLT matrix, respec-
 130 tively. A specific correspondence of L-CoLT and U-CoLT is shown as follow:

131 For the L-CoLT matrix M_L

$$132 \quad M_L = \text{L-CoLT}(\gamma, \mathbf{r}) = \begin{pmatrix} \gamma_1 & & & & & \\ \gamma_1 r_1 & \gamma_2 & & & & \\ \gamma_1 r_1 r_2 & \gamma_2 r_2 & \gamma_3 & & & \\ \vdots & \vdots & \vdots & \ddots & & \\ \gamma_1 \prod_{i=1}^{N-1} r_i & \gamma_2 \prod_{i=2}^{N-1} r_i & \gamma_3 \prod_{i=3}^{N-1} r_i & \cdots & \gamma_N & \end{pmatrix}$$

133 Similarly, for the U-CoLT matrix M_U

$$134 \quad M_U = \text{U-CoLT}(\boldsymbol{\gamma}', \mathbf{r}') = \begin{pmatrix} 0 & \gamma'_1 & \gamma'_2 r'_1 & \cdots & \gamma'_{N-1} \prod_{i=1}^{N-2} r'_i \\ & 0 & \gamma'_2 & \cdots & \gamma'_{N-1} \prod_{i=2}^{N-2} r'_i \\ & & \ddots & \ddots & \vdots \\ & & & 0 & \gamma'_{N-1} \\ & & & & 0 \end{pmatrix}$$

135 The special nature of the L-CoLT and U-CoLT matrices allows us to compute
136 matrix-vector multiplications in $O(N)$ operations.

137 For any $M_L = \text{L-CoLT}(\boldsymbol{\gamma}, \mathbf{r})$ and vector $\mathbf{y} \in \mathbb{R}^N$, the matrix-vector multiplication
138 $M_L \mathbf{y}$ is written as

$$139 \quad (3.3) \quad M_L \mathbf{y} = \begin{pmatrix} \gamma_1 y_1 & + & 0 & + & 0 & \cdots & + & 0 \\ \gamma_1 r_1 y_1 & + & \gamma_2 y_2 & + & 0 & \cdots & + & 0 \\ \gamma_1 r_1 r_2 y_1 & + & \gamma_2 r_2 y_2 & + & \gamma_3 y_3 & \cdots & + & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \gamma_1 \prod_{i=1}^{N-1} r_i y_1 & + & \gamma_2 \prod_{i=2}^{N-1} r_i y_2 & + & \gamma_3 \prod_{i=3}^{N-1} r_i y_3 & \cdots & + & \gamma_N y_N \end{pmatrix}.$$

140 Denote \mathbf{p}_k as the summation of the k -th row in (3.3), then one has

$$141 \quad p_1 = \gamma_1 y_1, \quad p_k = r_{k-1} p_{k-1} + \gamma_k y_k, \quad k = 2, \dots, N.$$

142 Based on this recursion formula, a fast implementation is proposed in Algorithm 3.1.

Algorithm 3.1 Fast L-CoLT Matrix-Vector multiplication

Input: input vector \mathbf{y} of size N , input matrix $M_L = \text{L-CoLT}(\boldsymbol{\gamma}, \mathbf{r})$

Output: $\mathbf{p} = M_L \mathbf{y}$

1: **procedure** LCMV($\mathbf{y}, \boldsymbol{\gamma}, \mathbf{r}$)

2: $p_1 = \gamma_1 y_1$

3: **for** $i = 1 : N - 1$ **do**

4: $p_{i+1} = r_i p_i + \gamma_{i+1} y_{i+1}$

return \mathbf{p}

143

144 Similarly, the fast matrix-vector multiplication for U-CoLT matrices is shown in
145 Algorithm 3.2.

146 Next, we denote the set \mathcal{C}^N as the direct sum of \mathcal{C}_L^N and \mathcal{C}_U^N , defined by

DEFINITION 3.3.

$$147 \quad (3.4) \quad \mathcal{C}^N = \mathcal{C}_L^N + \mathcal{C}_U^N = \{A + B \mid A \in \mathcal{C}_L^N, B \in \mathcal{C}_U^N\}.$$

148 Due to the linearity of matrix-vector multiplication, we can further develop the fast
149 matrix-vector multiplication algorithm for matrices in \mathcal{C}^N , which is given in Algo-
150 rithm 3.3.

151 The space and time complexities of these algorithms are $O(N)$, which is much
152 better than the original matrix-vector multiplication.

Algorithm 3.2 Fast U-CoLT Matrix-Vector multiplication**Input:** input vector \mathbf{y} of size N , input matrix $M = \text{U-CoLT}(\boldsymbol{\gamma}', \mathbf{r}')$ **Output:** $\mathbf{q} = M\mathbf{y}$

```

1: procedure UCMV( $\mathbf{y}$ ,  $\boldsymbol{\gamma}'$ ,  $\mathbf{r}'$ )
2:    $q_N = 0$ ,  $q_{N-1} = \gamma'_{N-1}y_N$ 
3:   for  $i = 2 : N - 1$  do
4:      $q_{N-i} = r'_{N-i}q_{N-i+1} + \gamma'_{N-i}y_{N-i+1}$ 
   return  $\mathbf{q}$ 

```

Algorithm 3.3 CoLT Matrix-Vector multiplication**Input:** input vector \mathbf{x} of size N , diagonal elements $\boldsymbol{\gamma}$, $\boldsymbol{\gamma}'$ and the ratio vector \mathbf{r} , \mathbf{r}' **Output:** $\mathbf{p} + \mathbf{q} = M\mathbf{y}$

```

1: procedure CMV( $\mathbf{y}$ ,  $\mathbf{c}^L$ ,  $\mathbf{c}^U$ ,  $\boldsymbol{\gamma}$ ,  $\boldsymbol{\gamma}'$ )
2:    $\mathbf{p} = \text{LCMV}(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{r})$ 
3:    $\mathbf{q} = \text{UCMV}(\mathbf{y}, \boldsymbol{\gamma}', \mathbf{r}')$ 
   return  $\mathbf{p} + \mathbf{q}$ 

```

153 **3.2. Some Basic Properties.** In this subsection, we justify some basic prop-
 154 erties of those matrices involved, which will be used in our algorithm.

155 **THEOREM 3.4.** (\mathcal{C}_L^N, \odot) and (\mathcal{C}_U^N, \odot) are Abelian groups, where \odot is the Hada-
 156 mard product.

157 *Proof.* We only prove the theorem for (\mathcal{C}_L^N, \odot) . It suffices to show that (\mathcal{C}_L^N, \odot)
 158 has the following properties:

159 **Closure:** For any two matrices $A = \text{L-CoLT}(\hat{\boldsymbol{\gamma}}, \hat{\mathbf{r}})$ and $B = \text{L-CoLT}(\tilde{\boldsymbol{\gamma}}, \tilde{\mathbf{r}})$, we
 160 set $D = A \odot B$. Since $d_{i,j} = a_{i,j}b_{i,j}$, one has

$$161 \quad (3.5) \quad d_{i,j}/d_{i+1,j} = (a_{i,j}b_{i,j}) / (a_{i+1,j}b_{i+1,j}) = \hat{r}_i \tilde{r}_i, \quad j = 1, 2, \dots, i,$$

162 and the strictly upper triangle part of D is obviously 0, which means $D = \text{L-CoLT}(\hat{\boldsymbol{\gamma}} \odot$
 163 $\tilde{\boldsymbol{\gamma}}, \hat{\mathbf{r}} \odot \tilde{\mathbf{r}}) \in \mathcal{C}_L^N$.

164 **Identity and Inverses:** Let

$$165 \quad (3.6) \quad E = \begin{pmatrix} 1 & & & & \\ 1 & 1 & & & \\ 1 & 1 & 1 & & \\ \vdots & \vdots & \vdots & \ddots & \\ 1 & 1 & 1 & \cdots & 1 \end{pmatrix} \in \mathcal{C}_L^N,$$

166 then for any $A = \text{L-CoLT}(\boldsymbol{\gamma}, \mathbf{r})$, $A \odot E = E \odot A = A$, which means E is the identity
 167 element. Let

$$168 \quad B = \begin{pmatrix} 1/a_{11} & & & & \\ 1/a_{21} & 1/a_{22} & & & \\ 1/a_{31} & 1/a_{32} & 1/a_{33} & & \\ \vdots & \vdots & \vdots & \ddots & \\ 1/a_{n1} & 1/a_{n2} & 1/a_{n3} & \cdots & 1/a_{nn} \end{pmatrix}.$$

169 Since

$$170 \quad b_{i,j}/b_{i+1,j} = a_{i+1,j}/a_{i,j} = 1/r_i, \quad j = 1, 2, \dots, i, \quad \square$$

171 then $B \in \mathcal{C}_L^N$. And obviously, $A \odot B = B \odot A = E$, which means B is the inverse of
 172 A .

173 Commutativity and Associativity: The commutativity and associativity can be
 174 derived from the commutative and associative law of real number multiplication.

175 Based on the above theorem, we can deduce directly

176 COROLLARY 3.5. (\mathcal{C}^N, \odot) is an abelian group with identity element $\mathbf{1}_N \mathbf{1}_N^T$.

177 THEOREM 3.6. For any vector $\mathbf{x} \in (\mathbb{R} \setminus \{0\})^N$, $f_{\mathbf{x}}(M) = (\text{diag}(\mathbf{x}))M$ and
 178 $g_{\mathbf{x}}(M) = M(\text{diag}(\mathbf{x}))$ are permutations in \mathcal{C}_L^N and \mathcal{C}_U^N .

179 *Proof.* We only prove the theorem for \mathcal{C}_L^N .

180 Closure: For any vector $\mathbf{x} \in (\mathbb{R} \setminus \{0\})^N$, and $M_L = \text{L-COLT}(\boldsymbol{\gamma}, \mathbf{r})$, let E be the
 181 one defined in Equation (3.6). Since

$$182 \quad E_1 = (\text{diag}(\mathbf{x}))E = \begin{pmatrix} x_1 & & & & \\ x_2 & x_2 & & & \\ x_3 & x_3 & x_3 & & \\ \vdots & \vdots & \vdots & \ddots & \\ x_n & x_n & x_n & \cdots & x_n \end{pmatrix} \in \mathcal{C}_L^N,$$

$$183 \quad E_2 = E(\text{diag}(\mathbf{x})) = \begin{pmatrix} x_1 & & & & \\ x_1 & x_2 & & & \\ x_1 & x_2 & x_3 & & \\ \vdots & \vdots & \vdots & \ddots & \\ x_1 & x_2 & x_3 & \cdots & x_n \end{pmatrix} \in \mathcal{C}_L^N,$$

185 we have

$$186 \quad (3.7) \quad \begin{aligned} f_{\mathbf{x}}(M) &= (\text{diag}(\mathbf{x}))M = (\text{diag}(\mathbf{x}))E \odot M = E_1 \odot M \in \mathcal{C}_L^N; \\ g_{\mathbf{x}}(M) &= M(\text{diag}(\mathbf{x})) = M \odot E(\text{diag}(\mathbf{x})) = M \odot E_2 \in \mathcal{C}_L^N. \end{aligned}$$

187 The last set membership can be derived by the closure of (\mathcal{C}_L^N, \odot) proved in Theo-
 188 rem 3.4. Hence, $f_{\mathbf{x}}$ and $g_{\mathbf{x}}$ are maps from \mathcal{C}_L^N to itself.

189 Injectiveness: For any two matrices $A = \text{L-COLT}(\hat{\boldsymbol{\gamma}}, \hat{\mathbf{r}})$ and $B = \text{L-COLT}(\tilde{\boldsymbol{\gamma}}, \tilde{\mathbf{r}})$,
 190 let D_1 be the inverse of E_1 and D_2 be the inverse of E_2 . If $f_{\mathbf{x}}(A) = f_{\mathbf{x}}(B)$, then

$$191 \quad A = D_1 \odot E_1 \odot A = D_1 \odot f_{\mathbf{x}}(A) = D_1 \odot f_{\mathbf{x}}(B) = D_1 \odot E_1 \odot B = B.$$

192 If $g_{\mathbf{x}}(A) = g_{\mathbf{x}}(B)$, then

$$193 \quad A = A \odot E_2 \odot D_2 = g_{\mathbf{x}}(A) \odot D_2 = g_{\mathbf{x}}(B) \odot D_2 = B \odot E_2 \odot D_2 = B,$$

194 which means $f_{\mathbf{x}}(\cdot)$ and $g_{\mathbf{x}}(\cdot)$ are injective functions.

195 Surjectiveness: For any $M \in \mathcal{C}_L^N$, let $Q_1 = D_1 \odot M$ and $Q_2 = M \odot D_2$, then

$$196 \quad \begin{aligned} f_{\mathbf{x}}(Q_1) &= E_1 \odot D_1 \odot M = E \odot M = M; \\ g_{\mathbf{x}}(Q_2) &= M \odot D_2 \odot E_2 = M \odot E = M, \end{aligned}$$

197 which means $f_{\mathbf{x}}(\cdot)$ and $g_{\mathbf{x}}(\cdot)$ are surjective functions. \square

198 COROLLARY 3.7. $f_{\mathbf{x}}(\cdot)$ and $g_{\mathbf{x}}(\cdot)$ are permutations in \mathcal{C}^N .

199 **4. The Fast Sinkhorn II.** In this section, we will discuss the implementation
 200 details to accelerate IPOT. In Algorithm 2.1, three parts lead to $O(N^2)$ algorithm
 201 complexity, i.e., the matrix Hadamard product (line 4), the matrix-vector multiplica-
 202 tion (lines 6-7), and the matrix scaling (line 8). They all rely on the representation
 203 and manipulation of L/U CoLT matrices.

204 For two discrete distributions on a 1D uniform mesh grid with a grid spacing of
 205 h , by introducing the notation $\lambda = e^{-h/\delta}$, the kernel matrix K is written as

$$206 \quad (4.1) \quad K = \begin{pmatrix} 1 & \lambda & \lambda^2 & \cdots & \lambda^{N-1} \\ \lambda & 1 & \lambda & \cdots & \lambda^{N-2} \\ \lambda^2 & \lambda & 1 & \cdots & \lambda^{N-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \lambda^{N-1} & \lambda^{N-2} & \lambda^{N-3} & \cdots & 1 \end{pmatrix} \in \mathcal{C}^N.$$

207 Below we discuss step by step of IPOT (Algorithm 2.1) to reduce the complexity:

- 208 • line 2: $\Gamma = \mathbf{1}_N \mathbf{1}_N^T \in \mathcal{C}^N$, we only need four vectors $(\boldsymbol{\gamma}, \boldsymbol{\gamma}', \mathbf{r}, \mathbf{r}')$ to represent
 209 Γ according to Theorem 3.2.
- 210 • line 4: the matrix Hadamard product $Q = K \odot \Gamma \in \mathcal{C}^N$ according to The-
 211 orem 3.4 and Corollary 3.5. By updating the four representation vectors
 212 $(\boldsymbol{\gamma}, \boldsymbol{\gamma}', \mathbf{r}, \mathbf{r}')$, we can obtain Q with $O(N)$ cost.
- 213 • lines 6-7: the matrix-vector multiplication $Q^T \boldsymbol{\phi}$ and $Q \boldsymbol{\psi}$ can be implemented
 214 with $O(N)$ cost according to Algorithm 3.3.
- 215 • line 8: the matrix scaling $\Gamma = \text{diag}(\boldsymbol{\phi}) Q \text{diag}(\boldsymbol{\psi}) \in \mathcal{C}^N$ according to The-
 216 orem 3.6 and Corollary 3.7. By updating the four representation vectors
 217 $(\boldsymbol{\gamma}, \boldsymbol{\gamma}', \mathbf{r}, \mathbf{r}')$, we can obtain Γ with $O(N)$ cost.

218 Based on the above discussions, we proposed the FS-2 algorithm with $O(N)$
 complexity. The pseudo-code is presented in Algorithm 4.1.

Algorithm 4.1 1D FS-2 Algorithm

Input: $\mathbf{u}, \mathbf{v} \in \mathbb{R}^N$; $L, \text{itr_max} \in \mathbb{N}^+$; $h, \delta \in \mathbb{R}$

Output: $W_1(\mathbf{u}, \mathbf{v})$

- 1: $\lambda \leftarrow e^{-h/\delta}$; $\boldsymbol{\phi}, \boldsymbol{\psi} \leftarrow \frac{1}{N} \mathbf{1}_N$; $\mathbf{r}, \mathbf{s} \leftarrow \mathbf{0}_N$
 - 2: $\boldsymbol{\alpha}^L, \boldsymbol{\beta}^L, \boldsymbol{\alpha}^U, \boldsymbol{\beta}^U \leftarrow \lambda \mathbf{1}_{N-1}$; $\boldsymbol{\gamma} \leftarrow \mathbf{1}_N$; $\boldsymbol{\gamma}' \leftarrow \lambda \mathbf{1}_{N-1}$
 - 3: **for** $t = 1 : \text{itr_max}$ **do**
 - 4: **for** $\ell = 1 : L$ **do**
 - 5: $\mathbf{r} \leftarrow \text{CMV}(\boldsymbol{\phi}, \boldsymbol{\beta}^L, \boldsymbol{\beta}^U, \boldsymbol{\gamma}, \boldsymbol{\gamma}')$
 - 6: $\boldsymbol{\psi} \leftarrow \mathbf{v} \odot \mathbf{r}$
 - 7: $\mathbf{s} \leftarrow \text{CMV}(\boldsymbol{\psi}, \boldsymbol{\alpha}^L, \boldsymbol{\alpha}^U, \boldsymbol{\gamma}, \boldsymbol{\gamma}')$
 - 8: $\boldsymbol{\phi} \leftarrow \mathbf{u} \odot \mathbf{s}$
 - 9: **for** $i = 1 : N - 1$ **do**
 - 10: $\alpha_i^L \leftarrow \lambda \alpha_i^L (\phi_{i+1}/\phi_i)$, $\beta_i^L \leftarrow \lambda \beta_i^L (\psi_{i+1}/\psi_i)$
 - 11: $\gamma_i' \leftarrow \lambda \gamma_i' \phi_i \psi_{i+1}$
 - 12: **for** $i = 1 : N - 2$ **do**
 - 13: $\alpha_i^U \leftarrow \lambda \alpha_i^U (\phi_i/\phi_{i+1})$, $\beta_i^U \leftarrow \lambda \beta_i^U (\psi_i/\psi_{i+1})$
 - 14: $\boldsymbol{\gamma} \leftarrow \boldsymbol{\phi} \odot \boldsymbol{\psi} \odot \boldsymbol{\gamma}$
 - return** $W_1(\mathbf{u}, \mathbf{v})$
-

219

220 There is a minor flaw in the above algorithm. The computational cost of $W_1(\mathbf{u}, \mathbf{v})$
 221 is still $O(N^2)$ in the last step. This was also ignored in our previous paper [24]. Now,

222 we would like to discuss this issue. The computation of $W_1(\mathbf{u}, \mathbf{v}) = \langle C, \Gamma \rangle$ can be
 223 regarded as the summation of all elements of the following matrix.

(4.2)

$$224 \quad C \odot \Gamma = \begin{pmatrix} 0 & h\gamma'_1 & 2h\gamma'_2 r'_1 & \cdots & (N-1)h\gamma'_{N-1} \prod_{i=1}^{N-2} r'_i \\ h\gamma_1 r_1 & 0 & h\gamma'_2 & \cdots & (N-2)h\gamma'_{N-1} \prod_{i=2}^{N-2} r'_i \\ 2h\gamma_1 r_1 r_2 & h\gamma_2 r_2 & 0 & \cdots & (N-3)h\gamma'_{N-1} \prod_{i=3}^{N-2} r'_i \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ (N-1)h\gamma_1 \prod_{i=1}^{N-1} r_i & (N-2)h\gamma_2 \prod_{i=2}^{N-1} r_i & (N-3)h\gamma_3 \prod_{i=3}^{N-1} r_i & \cdots & 0 \end{pmatrix}.$$

225 We separate the summation of the matrix to the lower and strictly upper triangular
 226 parts. Thus, the k -th line summation of two parts can be written as

$$227 \quad p_k = \sum_{i=1}^k \omega_{ki}, \quad q_k = \sum_{i=k+1}^N \omega_{ki}.$$

228 We can consider the following recursive computation

$$229 \quad (4.3) \quad \begin{aligned} p_1 &= 0, & p_2 &= h\gamma_1 r_1, & p'_2 &= h\gamma_1 r_1 + h\gamma_2, \\ p_i &= r_{i-1} (p_{i-1} + p'_{i-1}), & p'_i &= r_{i-1} p'_{i-1} + h\gamma_i, & i &= 3, 4, \dots, N. \\ q_N &= 0, & q_{N-1} &= h\gamma'_{N-1}, & q'_{N-1} &= h\gamma'_{N-1} + h, \\ q_j &= r'_j (q_{j+1} + q'_{j+1}), & q'_j &= r'_j q'_{j+1} + h, & j &= 1, 2, \dots, N-2. \end{aligned}$$

230 Thus, the Wasserstein-1 metric can be finally obtained with $O(N)$ cost

$$231 \quad W_1(\mathbf{u}, \mathbf{v}) = \langle C, \Gamma \rangle = \sum_{i=1}^N (p_i + q_i).$$

232 **5. Extension to high dimension.** In this section, we illustrate how the FS-
 233 2 algorithm generalizes to higher dimensions using the two-dimensional case as an
 234 example.

235 **5.1. Block Collinear Triangular Matrix.** Hereinafter, for $A \in \mathbb{R}^{MN \times MN}$,
 236 we break it into M^2 uniform blocks with size $N \times N$:

$$237 \quad A = \begin{pmatrix} A_{1,1} & A_{1,2} & A_{1,3} & \cdots & A_{1,M} \\ A_{2,1} & A_{2,2} & A_{2,3} & \cdots & A_{2,M} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ A_{M,1} & A_{M,2} & A_{M,3} & \cdots & A_{M,M} \end{pmatrix}.$$

238 And for vectors $\mathbf{x} \in \mathbb{R}^{kN}$, we break it into k uniform blocks as $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k)^T$, in
 239 which

$$240 \quad \mathbf{x}_i = (x_{1+(i-1)N}, x_{2+(i-1)N}, \dots, x_{iN})^T, \quad i = 1, 2, \dots, k.$$

241 To carry out the fast implementation of the matrix-vector multiplication of the block
242 matrix above, we generalize the definition of \mathcal{C}^N in (3.4) to the block case as

DEFINITION 5.1.

$$243 \quad \mathcal{C}^{N,M} = \{A \in \mathbb{R}^{MN \times MN} \mid A_{k,k} \in \mathcal{C}^N; \mathbf{r}^L, \mathbf{r}^U \in (\mathbb{R} \setminus \{0\})^{(M-1)N}; \\ A_{i+1,j} = (\text{diag}(\mathbf{r}_i^L)) A_{i,j}, j \leq i; A_{i-1,j} = (\text{diag}(\mathbf{r}_{i-1}^U)) A_{i,j}, i \leq j\}.$$

244 Since $\mathcal{C}^{N,M}$ is a generalization of \mathcal{C}^N , we can also use the strategy of Algorithm 3.3
245 in blocks to reduce the computational cost of matrix-vector multiplications. For a
246 vector $\mathbf{x} \in \mathbb{R}^{NM}$, the matrix-vector multiplication $A\mathbf{x}$ is written as

$$247 \quad (5.1) \quad A\mathbf{x} = \begin{pmatrix} A_{1,1}\mathbf{x}_1 + A_{1,2}\mathbf{x}_2 + A_{1,3}\mathbf{x}_3 + \cdots + A_{1,M}\mathbf{x}_M \\ A_{2,1}\mathbf{x}_1 + A_{2,2}\mathbf{x}_2 + A_{2,3}\mathbf{x}_3 + \cdots + A_{2,M}\mathbf{x}_M \\ A_{3,1}\mathbf{x}_1 + A_{3,2}\mathbf{x}_2 + A_{3,3}\mathbf{x}_3 + \cdots + A_{3,M}\mathbf{x}_M \\ \vdots \\ A_{M,1}\mathbf{x}_1 + A_{M,2}\mathbf{x}_2 + A_{M,3}\mathbf{x}_3 + \cdots + A_{M,M}\mathbf{x}_M \end{pmatrix}.$$

248 We separate the summation of row k to the lower triangular part \mathbf{p}_k and the strictly
249 upper triangular part \mathbf{q}_k . Then computing $A\mathbf{x}$ is formulated as

$$250 \quad A\mathbf{x} = \mathbf{p} + \mathbf{q}, \quad \mathbf{p}_k = \sum_{i=1}^k A_{k,i}\mathbf{x}_i, \quad \mathbf{q}_k = \sum_{i=k+1}^M A_{k,i}\mathbf{x}_i, \quad k = 1, \dots, M.$$

251 If A is in $\mathcal{C}^{N,M}$ with \mathbf{r}^L and \mathbf{r}^U , instead of directly calculating \mathbf{p}_k and \mathbf{q}_k , a successive
252 computation is used

$$253 \quad (5.2) \quad \mathbf{p}_1 = A_{1,1}\mathbf{x}_1, \quad \mathbf{p}_k = \mathbf{r}_{k-1}^L \odot \mathbf{p}_{k-1} + A_{k,k}\mathbf{x}_k, \quad k = 2, \dots, M, \\ \mathbf{q}_M = \mathbf{0}_N, \quad \mathbf{q}_{k-1} = \mathbf{r}_{k-1}^U \odot (\mathbf{q}_k + A_{k,k}\mathbf{x}_k), \quad k = M, M-1, \dots, 2.$$

254 Since the computation of $A_{k,k}\mathbf{x}_k$ can be carried out with $O(N)$ complexity by using
255 Algorithm 3.3, the whole computation is of $O(NM)$ complexity.

256 Similar to Theorem 3.4 and Theorem 3.6, $\mathcal{C}^{N,M}$ is closed under the Hadamard
257 product and matrix scaling.

258 THEOREM 5.2. $(\mathcal{C}^{N,M}, \odot)$ is an Abelian group; Matrix scaling operations are per-
259 mutations in $\mathcal{C}^{N,M}$.

260 *Proof.* For any $A \in \mathcal{C}^{N,M}$, since the diagonal blocks of A are in \mathcal{C}^N , by Corol-
261 lary 3.7, all blocks in A are in \mathcal{C}^N . Then the two properties can be proved in a similar
262 way as in subsection 3.2. \square

263 **5.2. The 2D FS-2 Algorithm.** Consider two discretized probabilistic distri-
264 butions

$$265 \quad \mathbf{u} = (u_{11}, u_{21}, \dots, u_{N1}, u_{12}, \dots, u_{i_1 j_1}, \dots, u_{NM}), \\ 266 \quad \mathbf{v} = (v_{11}, v_{21}, \dots, v_{N1}, v_{12}, \dots, v_{i_2 j_2}, \dots, v_{NM}),$$

268 on a uniform 2D mesh of size $N \times M$ with a vertical spacing of h_1 and a horizontal
 269 spacing of h_2 . The corresponding kernel matrix is written as

$$270 \quad K = \begin{pmatrix} K_0 & \lambda_2 K_0 & \lambda_2^2 K_0 & \cdots & \lambda_2^{M-1} K_0 \\ \lambda_2 K_0 & K_0 & \lambda_2 K_0 & \cdots & \lambda_2^{M-2} K_0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \lambda_2^{M-1} K_0 & \lambda_2^{M-2} K_0 & \lambda_2^{M-3} K_0 & \cdots & K_0 \end{pmatrix},$$

271 where the sub-matrix

$$272 \quad K_0 = \begin{pmatrix} 1 & \lambda_1 & \cdots & \lambda_1^{N-1} \\ \lambda_1 & 1 & \cdots & \lambda_1^{N-2} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{N-1} & \lambda_1^{N-2} & \cdots & 1 \end{pmatrix},$$

273 and

$$274 \quad \lambda_1 = e^{-h_1/\delta}, \quad \lambda_2 = e^{-h_2/\delta}.$$

275 Obviously, the original 2D kernel K contains blocks which are multiples of the
 276 1D kernel, hence belongs to $\mathcal{C}^{N,M}$. By an analysis similar to that in section 4 and
 277 using Theorem 5.2, the matrices Q and Γ are in $\mathcal{C}^{N,M}$ throughout the course of the
 278 iteration, which means that all the matrix-vector multiplications can be carried out
 279 by using recursion (5.2). Thus, the total cost of matrix-vector multiplication of our
 280 FS-2 algorithm for 2D Wasserstein-1 metric is reduced to $O(NM)$.

281 In the 2D FS-2 algorithm, we use $\hat{\cdot}$ to distinguish whether it is a coefficient of the
 282 block or the inner sub-matrix, and update them simultaneously after an inner loop.
 283 The pseudo-code is presented in Algorithm 5.1. Updating the coefficients of inner
 284 sub-matrices and computation of $W_1(\mathbf{u}, \mathbf{v})$ are omitted in the pseudo-code since they
 285 are similar to the 1D case, which we have described in detail.

286 **6. Numerical Experiments.** In this section, we carry out three numerical ex-
 287 periments to evaluate the FS-2 algorithm, including one 1D example and two 2D
 288 examples. The true Wasserstein metric W_{LP} is obtained by solving the original OT
 289 (2.1) using interior-point methods [11, 20]. In our experiments, for both IPOT and
 290 FS-2, the number of inner loops is set as $L = 20$ and the regularization param-
 291 eter $\delta^{(t)}$ is set to 1. The number of iterations here is the total number of loops:
 292 $\#iteration = itr_max \times L$. In order to deal with the difficulties caused by zeros, we
 293 utilize the rescaling method in [23]:

$$294 \quad (6.1) \quad D(f, g) = W_1 \left(\frac{\frac{f}{\|f\|} + \eta}{1 + N\eta}, \frac{\frac{g}{\|g\|} + \eta}{1 + N\eta} \right),$$

295 In the following, we refer to formula (6.1) for numerical stability with $\eta = 10^{-5}$.
 296 All the experiments are conducted on a platform with 128G RAM, and one Intel(R)
 297 Xeon(R) Gold 5117 CPU @2.00GHz with 14 cores.

Algorithm 5.1 2D FS-2 Algorithm

Input: $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{NM}$; $L, \text{itr_max} \in \mathbb{N}^+$; $h_1, h_2, \delta \in \mathbb{R}$

Output: $W_1(\mathbf{u}, \mathbf{v})$

- 1: $\lambda_1 \leftarrow e^{-h_1/\delta}; \lambda_2 \leftarrow e^{-h_2/\delta}; \boldsymbol{\phi}^{(0)}, \boldsymbol{\psi}^{(0)} \leftarrow \frac{1}{NM} \mathbf{1}_{NM}; \mathbf{p}^{(0)}, \mathbf{r}^{(0)}, \mathbf{q}^{(0)}, \mathbf{s}^{(0)} \leftarrow \mathbf{0}_{NM}$
- 2: $\boldsymbol{\gamma} \leftarrow \mathbf{1}_{NM}, \boldsymbol{\gamma}', \boldsymbol{\alpha}^L, \boldsymbol{\beta}^L, \lambda_1 \mathbf{1}_{(N-1)M}, \boldsymbol{\alpha}^U, \boldsymbol{\beta}^U \leftarrow \lambda_1 \mathbf{1}_{(N-2)M}$
- 3: $\widehat{\boldsymbol{\alpha}}^L, \widehat{\boldsymbol{\beta}}^L, \widehat{\boldsymbol{\alpha}}^U, \widehat{\boldsymbol{\beta}}^U \leftarrow \lambda \mathbf{1}_{N(M-1)}$
- 4: **while** $t = 1 : \text{itr_max}$ **do**
- 5: **for** $\ell = 1 : L$ **do**
- 6: $\mathbf{r}_1 \leftarrow \text{CMV}(\boldsymbol{\phi}_1, \boldsymbol{\beta}_1), \mathbf{s}_M \leftarrow \mathbf{0}_N$
- 7: **for** $i = 1 : M - 1$ **do**
- 8: $\mathbf{r}_{i+1} \leftarrow \widehat{\boldsymbol{\beta}}_i^L \odot \mathbf{r}_i + \text{CMV}(\boldsymbol{\phi}_{i+1}, \boldsymbol{\beta}_{i+1}^L, \boldsymbol{\beta}_{i+1}^U, \boldsymbol{\gamma}_{i+1}, \boldsymbol{\gamma}'_{i+1})$
- 9: $\mathbf{s}_{N-i} \leftarrow \widehat{\boldsymbol{\beta}}_{N-i}^U \odot (\mathbf{s}_{N-i+1} + \text{CMV}(\boldsymbol{\phi}_{N-i+1}, \boldsymbol{\beta}_{N-i+1}^L, \boldsymbol{\beta}_{N-i+1}^U, \boldsymbol{\gamma}_{N-i+1}, \boldsymbol{\gamma}'_{N-i+1}))$
- 10: $\boldsymbol{\psi} \leftarrow \mathbf{v} \odot (\mathbf{r} + \mathbf{s})$
- 11: $\mathbf{p}_1 \leftarrow \text{CMV}(\boldsymbol{\psi}_1, \boldsymbol{\alpha}_1), \mathbf{q}_M \leftarrow \mathbf{0}_N$
- 12: **for** $i = 1 : M - 1$ **do**
- 13: $\mathbf{p}_{i+1} \leftarrow \widehat{\boldsymbol{\alpha}}_i^L \odot \mathbf{p}_i + \text{CMV}(\boldsymbol{\psi}_{i+1}, \boldsymbol{\alpha}_{i+1}^L, \boldsymbol{\alpha}_{i+1}^U, \boldsymbol{\gamma}_{i+1}, \boldsymbol{\gamma}'_{i+1})$
- 14: $\mathbf{q}_{N-i} \leftarrow \widehat{\boldsymbol{\alpha}}_{N-i}^U \odot (\mathbf{q}_{N-i+1} + \text{CMV}(\boldsymbol{\psi}_{N-i+1}, \boldsymbol{\alpha}_{N-i+1}^L, \boldsymbol{\alpha}_{N-i+1}^U, \boldsymbol{\gamma}_{N-i+1}, \boldsymbol{\gamma}'_{N-i+1}))$
- 15: $\boldsymbol{\phi} \leftarrow \mathbf{u} \odot (\mathbf{p} + \mathbf{q})$
- 16: **for** $i = 1 : M - 1$ **do**
- 17: $\widehat{\boldsymbol{\alpha}}_i^L \leftarrow \lambda_2 (\boldsymbol{\phi}_{i+1}/\boldsymbol{\phi}_i) \odot \widehat{\boldsymbol{\alpha}}_i^L, \widehat{\boldsymbol{\beta}}_i^L \leftarrow \lambda_2 (\boldsymbol{\psi}_{i+1}/\boldsymbol{\psi}_i) \odot \widehat{\boldsymbol{\beta}}_i^L$
- 18: $\widehat{\boldsymbol{\alpha}}_i^U \leftarrow \lambda_2 (\boldsymbol{\phi}_i/\boldsymbol{\phi}_{i+1}) \odot \widehat{\boldsymbol{\alpha}}_i^U, \widehat{\boldsymbol{\beta}}_i^U \leftarrow \lambda_2 (\boldsymbol{\psi}_i/\boldsymbol{\psi}_{i+1}) \odot \widehat{\boldsymbol{\beta}}_i^U$
- 19: Update $\boldsymbol{\gamma}, \boldsymbol{\gamma}', \boldsymbol{\alpha}^L, \boldsymbol{\beta}^L, \boldsymbol{\alpha}^U, \boldsymbol{\beta}^U$

return $W_1(\mathbf{u}, \mathbf{v})$

298 **6.1. 1D Gaussian distributions.** We consider the Wasserstein-1 metric be-
 299 tween two mixtures of 1D Gaussian distributions: $0.4\mathcal{N}(60, 64) + 0.6\mathcal{N}(40, 36)$ and
 300 $0.5\mathcal{N}(35, 81) + 0.5\mathcal{N}(70, 81)$, which is the experiment setting in [35]. Input vectors
 301 \mathbf{u} and \mathbf{v} are generated by integration on the uniform discretization of interval $[0, 100]$
 302 with node size N .

303 We first compare the convergence of FS-1 and FS-2 for $N = 1000$. We tested 100
 304 experiments, and each experiment was performed for 10,000 iterations. In Figure 1,
 305 the differences of the Wasserstein-1 metric between the true solution W_{LP} and the
 306 numerical solutions generated by FS-1 and FS-2 are depicted. As expected, as ε
 307 decreases, the error of FS-1 decreases gradually after the iterations converge. We can
 308 observe this for $\varepsilon = 1/20$ and $\varepsilon = 1/80$. However, for $\varepsilon = 1/320$, we can not observe
 309 convergence. In fact, the error does not drop over 10,000 iterations. This is because
 310 ε is too small, making updates extremely slow. In fact, after 20,000 iterations, the
 311 result of $\varepsilon = 1/320$ will continue to drop, and the final error is smaller than that of
 312 $\varepsilon = 1/20$ and $\varepsilon = 1/80$. However, in any case, the results of FS-1 are far inferior to
 313 those of FS-2, both in terms of accuracy and convergence rate.

314 The averaged computational time of the IPOT method and the FS-2 algorithm is
 315 given in Table 1 and Figure 2 (left). Apparently, the FS-2 algorithm has achieved an
 316 overwhelming advantage over the IPOT method in terms of computational speed, and
 317 ensures that the transport plans of the two are almost the same. This replicates the
 318 advantages of the FS-1 algorithm over the Sinkhorn algorithm. According to the data

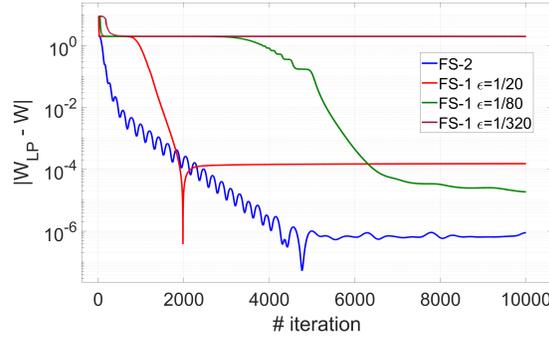


FIG. 1. The 1D Gaussian distribution problem. The errors between the numerical results generated by FS-1 or FS-2 and the true Wasserstein-1 metric w.r.t. number of iterations.

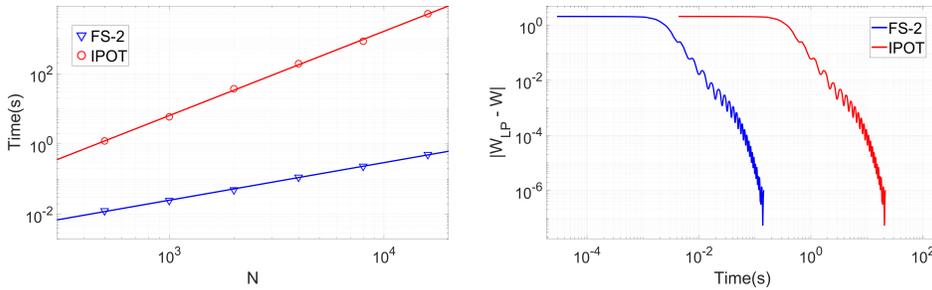


FIG. 2. The 1D Gaussian distribution problem. Left: The comparison of computational time between the FS-2 algorithm and the IPOT method with different numbers of grid points N . Right: The computational time required to reach the absolute error of the Wasserstein-1 metric.

319 fitting results, the empirical complexity of the FS-2 algorithm is $O(N^{1.07})$, which is
 320 much smaller than the $O(N^{2.40})$ complexity of the IPOT method. At last, we show the
 321 computational time required to reach the absolute error of the Wasserstein-1 metric
 322 for $N = 1,000$ in Figure 2 (right). Clearly, the FS-2 algorithm has an advantage of
 323 two orders of magnitude in computational time compared to the IPOT method.

324 **6.2. 2D Random distributions.** Next, we compute the Wasserstein-1 metric
 325 between two $N \times N$ dimensional random vectors whose elements obey the uniform
 326 distribution on $(0, 1)$. Without loss of generality, we set $h_1 = h_2 = 0.1$. We also
 327 tested 100 experiments, and each experiment was performed for 10,000 iterations.
 328 We hope to test the performance of the FS-2 algorithm in 2D through this example.
 329 The differences in the Wasserstein-1 metric between the true solution W_{LP} and the
 330 numerical solutions generated by FS-1 and FS-2 are shown in Figure 3. From this, we
 331 can observe that FS-1 converges quickly for $\varepsilon = 1/20$, but the error is large. When
 332 $\varepsilon = 1/80$, the iteration converges at about 5,000 steps. The error keep decreasing
 333 even after 10,000 steps for $\varepsilon = 1/320$. However, their errors and convergence speed
 334 are not as good as FS-2.

335 The averaged computational time of the IPOT method and the FS-2 algorithm
 336 is given in Table 2 and Figure 4 (left). According to the data fitting results, the
 337 empirical complexity of the FS-2 algorithm is $O(N^{2.05})$, which is much smaller than

TABLE 1

The 1D Gaussian distribution problem. The comparison between the IPOT method and the FS-2 algorithm with the different number of grid points N . Columns 2-4 are the averaged computational time of the two algorithms and the speed-up ratio of the FS-2 algorithm. Column 5 is the Frobenius norm of the difference between the transport plan computed by the two algorithms.

N	Computational time (s)		Speed-up ratio	$\ P_{FS} - P\ _F$
	FS-2	IPOT		
500	1.25×10^{-2}	1.22×10^0	9.76×10^1	2.09×10^{-15}
2000	4.95×10^{-2}	3.73×10^1	7.52×10^2	6.65×10^{-16}
8000	2.32×10^{-1}	8.41×10^2	3.63×10^3	8.57×10^{-16}

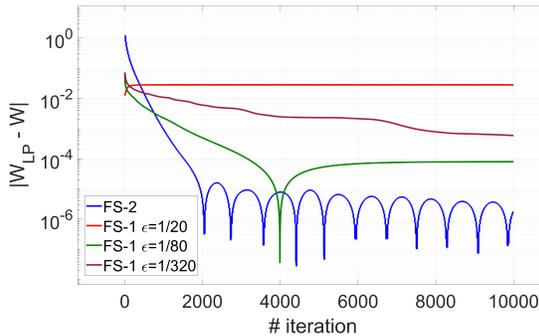


FIG. 3. The 2D random distribution problem. The errors between the numerical results generated by FS-1 or FS-2 and the true Wasserstein-1 metric w.r.t. number of iterations.

338 the $O(N^{4.98})$ complexity of the IPOT method. The computational time required to
 339 reach the absolute error of the Wasserstein-1 metric for $N \times N = 32 \times 32$ is also
 340 presented in Figure 4 (right). Similar to the previous subsection, we can also observe
 341 the huge computational efficiency of the FS-2 algorithm over the IPOT method.

342 **6.3. Image matching problem.** The final experiment tests the performance
 343 of our FS-2 algorithm for high-resolution image matching. This is a successful appli-
 344 cation of the Optimal Transport [4, 17, 7]. We select two images from the DIV2K
 345 dataset [1]. Through a process similar to Subsection 5.4 in the manuscript [24], we
 346 compute the Wasserstein-1 metric between the two images. The differences in the
 347 Wasserstein-1 metric between the true solution W_{LP} and the numerical solutions gen-
 348 erated by FS-1 and FS-2 are depicted in Figure 6. We also present the averaged
 349 computational time of the IPOT method and the FS-2 algorithm in Table 3. More-
 350 over, the computational time required to reach the absolute error of the Wasserstein-1
 351 metric for $N \times N = 32 \times 32$ is shown in Figure 7. From these results, we can get
 352 the same conclusion as before, that is, the FS-2 algorithm seems to be the numerical
 353 algorithm with the fastest convergence and the lowest complexity for computing the
 354 Wasserstein-1 metric.

355 **7. Conclusion.** As the follow-up of the FS-1 paper, we generalize the result of
 356 matrix-vector multiplication at $O(N)$ costs for the special matrix to the more general
 357 L/U-CoLT matrix. We illustrate that only two vectors are required to represent any
 358 L/U-CoLT matrix. Moreover, we also prove the closure of families of L/U-CoLT
 359 matrices to matrix Hadamard product and matrix scaling. Therefore, the above

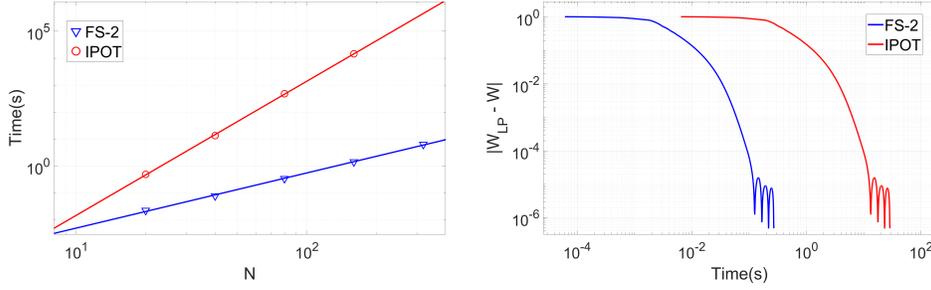


FIG. 4. The 2D random distribution problem. Left: The comparison of computational time between the FS-2 algorithm and the IPOT method with different numbers of grid points N . Right: The computational time required to reach the absolute error of the Wasserstein-1 metric.

TABLE 2

The 2D random distribution problem. The comparison between the IPOT method and the FS-2 algorithm with different total number of grid nodes $N \times N$. Columns 2-4 are the averaged computational time of the two algorithms and the speed-up ratio of the FS-2 algorithm. Column 5 is the Frobenius norm of the difference between the transport plan computed by the two algorithms.

$N \times N$	Computational time (s)		Speed-up ratio	$\ P_{FS} - P\ _F$
	FS-2	IPOT		
20×20	2.24×10^{-2}	4.88×10^{-1}	2.18×10^1	2.40×10^{-16}
40×40	7.74×10^{-2}	1.34×10^1	1.73×10^2	1.52×10^{-16}
80×80	3.38×10^{-1}	4.79×10^2	1.42×10^3	1.40×10^{-16}
160×160	1.42×10^0	1.46×10^4	1.03×10^4	8.51×10^{-17}
320×320	6.31×10^0	—	—	—



FIG. 5. The image matching problem. Illustration of images.

TABLE 3

The image matching problem. The comparison between the IPOT method and the FS-2 algorithm with the different total number of grid nodes $N \times N$. Columns 2-4 are the averaged computational time of the two algorithms and the speed-up ratio of the FS-2 algorithm. Column 5 is the Frobenius norm of the difference between the transport plan computed by the two algorithms.

$N \times N$	Computational time (s)		Speed-up ratio	$\ P_{FS} - P\ _F$
	FS-2	IPOT		
100×100	5.16×10^{-1}	1.44×10^3	2.79×10^3	6.44×10^{-17}
200×200	2.31×10^0	3.69×10^4	1.60×10^4	4.65×10^{-17}
400×400	9.69×10^0	—	—	—
800×800	4.18×10^1	—	—	—

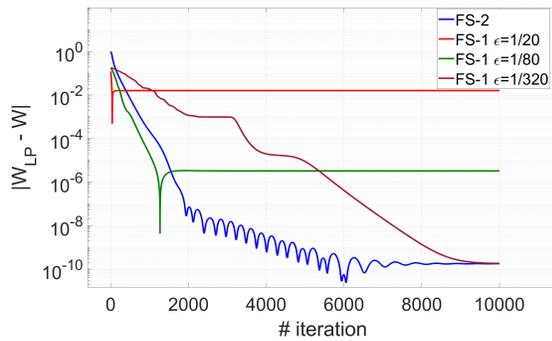


FIG. 6. The image matching problem. The errors between the numerical results generated by FS-1 or FS-2 and the true Wasserstein-1 metric w.r.t. number of iterations.

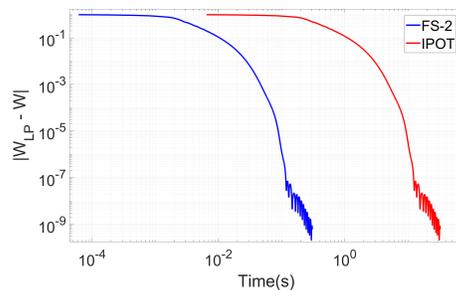


FIG. 7. The image matching problem. The computational time required to reach the absolute error of the Wasserstein-1 metric.

360 matrix operations are essentially updating the representation vectors, which reduce
 361 both time and space complexity to $O(N)$. These results can be directly applied to
 362 the Inexact Proximal point method for Optimal Transport problem and reduce the
 363 overall computational complexity to $O(N)$. From this, we develop the Fast Sinkhorn
 364 II algorithm. It does not seem to be an overstatement that for the computation of
 365 the Wasserstein-1 metric, we have probably obtained the most competitive method,
 366 both in terms of convergence speed and computational complexity.

367

REFERENCES

- 368 [1] E. AGUSTSSON AND R. TIMOFTE, *Ntire 2017 challenge on single image super-resolution:*
 369 *Dataset and study*, in The IEEE Conference on Computer Vision and Pattern Recognition
 370 (CVPR) Workshops, July 2017.
- 371 [2] J.-D. BENAMOU AND Y. BRENIER, *A computational fluid mechanics solution to the Monge-*
 372 *Kantorovich mass transfer problem*, Numer. Math., 84 (2000), pp. 375–393.
- 373 [3] J.-D. BENAMOU, B. D. FROESE, AND A. M. OBERMAN, *Numerical solution of the optimal*
 374 *transportation problem using the Monge–Ampère equation*, J. Comput. Phys., 260 (2014),
 375 pp. 107–126.
- 376 [4] M. BURGER, M. FRANEK, AND C.-B. SCHÖNLIEB, *Regularized regression and density estimation*
 377 *based on optimal transport*, Appl. Math. Res. Express, 2012 (2012), pp. 209–253.
- 378 [5] G. BUTTAZZO, L. DE PASCALE, AND P. GORI-GIORGI, *Optimal-transport formulation of elec-*
 379 *tronic density-functional theory*, Phys. Rev. A, 85 (2012), p. 062502.
- 380 [6] J. CHEN, Y. CHEN, H. WU, AND D. YANG, *The quadratic Wasserstein metric for earthquake*
 381 *location*, J. Comput. Phys., 373 (2018), pp. 188–209.
- 382 [7] P. CLARYSSE, B. DELHAY, M. PICQ, AND J. POUSIN, *Optimal extended optical flow subject to*

- 383 *a statistical constraint*, J. Comput. Appl. Math., 234 (2010), pp. 1291–1302.
- 384 [8] P. L. COMBETTES AND J.-C. PESQUET, *Proximal splitting methods in signal processing*, in
385 Fixed-Point Algorithms for Inverse Problems in Science and Engineering, Springer, 2011,
386 pp. 185–212.
- 387 [9] C. COTAR, G. FRIESECKE, AND C. KLÜPPELBERG, *Density functional theory and optimal trans-*
388 *portation with Coulomb cost*, Comm. Pure Appl. Math., 66 (2013), pp. 548–599.
- 389 [10] M. CUTURI, *Sinkhorn distances: Lightspeed computation of optimal transport*, in Advances in
390 Neural Information Processing Systems, vol. 26, 2013, pp. 2292–2300.
- 391 [11] I. DIKIN, *Iterative solution of problems of linear and quadratic programming*, Dokl. Akad. Nauk,
392 174 (1967), pp. 747–748.
- 393 [12] B. ENGQUIST, K. REN, AND Y. YANG, *The quadratic Wasserstein metric for inverse data*
394 *matching*, Inverse Problems, 36 (2020), p. 055001.
- 395 [13] J. FRANKLIN AND J. LORENZ, *On the scaling of multidimensional matrices*, Linear Algebra
396 Appl., 114 (1989), pp. 717–735.
- 397 [14] B. D. FROESE, *Numerical methods for the elliptic Monge-Ampère equation and optimal trans-*
398 *port*, PhD thesis, Simon Fraser University, Burnaby, BC, Canada, 2012.
- 399 [15] B. D. FROESE AND A. M. OBERMAN, *Convergent finite difference solvers for viscosity solutions*
400 *of the elliptic Monge-Ampère equation in dimensions two and higher*, SIAM J. Numer.
401 Anal., 49 (2011), pp. 1692–1714.
- 402 [16] I. GOODFELLOW, J. POUGET-ABADIE, M. MIRZA, B. XU, D. WARDE-FARLEY, S. OZAIR,
403 A. COURVILLE, AND Y. BENGIO, *Generative adversarial nets*, in Advances in Neural Infor-
404 mation Processing Systems, vol. 27, 2014.
- 405 [17] S. HAKER, L. ZHU, A. TANNENBAUM, AND S. ANGENENT, *Optimal mass transport for registra-*
406 *tion and warping*, Int. J. Comput. Vis., 60 (2004), pp. 225–240.
- 407 [18] H. HEATON, S. W. FUNG, A. T. LIN, S. OSHER, AND W. YIN, *Wasserstein-based projections*
408 *with applications to inverse problems*, arXiv preprint arXiv:2008.02200, (2020).
- 409 [19] Y. HU, H. CHEN, AND X. LIU, *A global optimization approach for multi-marginal optimal*
410 *transport problems with Coulomb cost*, arXiv preprint arXiv:2110.07352, (2021).
- 411 [20] N. KARMARKAR, *A new polynomial-time algorithm for linear programming*, in Proceedings of
412 the sixteenth annual ACM symposium on Theory of computing, 1984, pp. 302–311.
- 413 [21] W. LI, E. K. RYU, S. OSHER, W. YIN, AND W. GANGBO, *A parallel method for earth mover’s*
414 *distance*, J. Sci. Comput., 75 (2018), pp. 182–197.
- 415 [22] X. LI, D. SUN, AND K.-C. TOH, *An asymptotically superlinearly convergent semismooth New-*
416 *ton augmented Lagrangian method for linear programming*, SIAM J. Optim., 30 (2020),
417 pp. 2410–2440.
- 418 [23] Z. LI, Y. TANG, J. CHEN, AND H. WU, *The quadratic Wasserstein metric with squaring scaling*
419 *for seismic velocity inversion*, arXiv preprint arXiv:2201.11305, (2022).
- 420 [24] Q. LIAO, J. CHEN, Z. WANG, B. BAI, S. JIN, AND H. WU, *Fast Sinkhorn I: An $O(N)$ algorithm*
421 *for the Wasserstein-1 metric*, Commun. Math. Sci., accepted, (2022).
- 422 [25] A. T. LIN, W. LI, S. OSHER, AND G. MONTÚFAR, *Wasserstein proximal of GANs*, in Interna-
423 tional Conference on Geometric Science of Information, Springer, 2021, pp. 524–533.
- 424 [26] C. MENG, Y. KE, J. ZHANG, M. ZHANG, W. ZHONG, AND P. MA, *Large-scale optimal transport*
425 *map estimation using projection pursuit*, in Advances in Neural Information Processing
426 Systems, vol. 32, 2019, pp. 8118–8129.
- 427 [27] C. MENG, J. YU, J. ZHANG, P. MA, AND W. ZHONG, *Sufficient dimension reduction for clas-*
428 *sification using principal optimal transport direction*, in Advances in Neural Information
429 Processing Systems, vol. 33, 2020.
- 430 [28] L. MÉTIVIER, R. BROSSIER, Q. MERIGOT, É. OUDET, AND J. VIRIEUX, *An optimal trans-*
431 *port approach for seismic tomography: Application to 3D full waveform inversion*, Inverse
432 Problems, 32 (2016), p. 115008.
- 433 [29] O. MUSEYKO, M. STIGLMAYR, K. KLAMROTH, AND G. LEUGERING, *On the application of*
434 *the Monge-Kantorovich problem to image registration*, SIAM J. Imaging Sci., 2 (2009),
435 pp. 1068–1097.
- 436 [30] O. PELE AND M. WERMAN, *Fast and robust earth mover’s distances*, in 2009 IEEE 12th Inter-
437 national Conference on Computer Vision, IEEE, 2009, pp. 460–467.
- 438 [31] G. PEYRÉ, M. CUTURI, ET AL., *Computational optimal transport: With applications to data*
439 *science*, Found. Trends Mach. Learn., 11 (2019), pp. 355–607.
- 440 [32] Y. RUBNER, C. TOMASI, AND L. J. GUIBAS, *The earth mover’s distance as a metric for image*
441 *retrieval*, Int. J. Comput. Vis., 40 (2000), pp. 99–121.
- 442 [33] F. SANTAMBROGIO, *Optimal transport for applied mathematicians: Calculus of variations, pdes,*
443 *and modeling*, Progr. Nonlinear Differential Equations Appl., Birkhäuser, Basel, 2015.
- 444 [34] R. SINKHORN, *Diagonal equivalence to matrices with prescribed row and column sums*, Amer.

- 445 Math. Monthly, 74 (1967), pp. 402–405.
446 [35] Y. XIE, X. WANG, R. WANG, AND H. ZHA, *A fast proximal point method for computing exact*
447 *Wasserstein distance*, in *Uncertainty in artificial intelligence*, PMLR, 2020, pp. 433–453.
448 [36] L. YANG, J. LI, D. SUN, AND K.-C. TOH, *A fast globally linearly convergent algorithm for the*
449 *computation of Wasserstein barycenters.*, *J. Mach. Learn. Res.*, 22 (2021), pp. 1–37.
450 [37] Y. YANG, B. ENGQUIST, J. SUN, AND B. F. HAMFELDT, *Application of optimal transport and the*
451 *quadratic Wasserstein metric to full-waveform inversion*, *Geophysics*, 83 (2018), pp. R43–
452 R62.