

# SCALABLE SINKHORN BACKPROPAGATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Optimal transport has recently gained increasing attention in the context of deep learning. A major contributing factor is the line of work on smooth relaxations that make the classical optimal transport problem differentiable. The most prominent example is entropy regularized optimal transport which can be optimized efficiently via an alternating scheme of Sinkhorn projections. We thus experienced a surge of deep learning techniques that use the Sinkhorn operator to learn matchings, permutations, sorting and ranking, or to construct a geometrically motivated loss function for generative models. The prevalent approach to training such a neural network is first-order optimization by algorithmic unrolling of the forward pass. Hence, the runtime and memory complexity of the backward pass increase linearly with the number of Sinkhorn iterations. This often makes it impractical when computational resources like GPU memory are scarce. A more efficient alternative is computing the derivative of a Sinkhorn layer via implicit differentiation. Our main contribution is deriving a simple and efficient algorithm that performs this backward pass in closed form. It is based on the Sinkhorn operator in its most general form – with learnable cost matrices and target capacities. We further provide a theoretical analysis with error bounds for approximate inputs. Finally, we demonstrate that, for a number of applications, replacing automatic differentiation with our module often improves the stability and accuracy of the obtained gradients while drastically reducing the computation cost.

## 1 INTRODUCTION

Computing similarities between probability distributions is a fundamental problem in machine learning. In the most general case, tools like the Kullback–Leibler (KL) divergence (Kullback & Leibler, 1951) establish a notion of distance between such measures  $\mu$  and  $\nu$ . In practice,  $\mu$  and  $\nu$  are often themselves embedded in a metric space  $(\Omega, d_\Omega)$ . Common examples include spaces of images, voxel grids, point clouds, 3D surface meshes, or generic Euclidean features. For this class of problems, optimal transport (OT) proves to be a powerful formalism (Villani, 2003; Peyré et al., 2019). The Kantorovich formulation (Kantorovich, 1942) of OT yields a well-defined distance function over probability measures on  $\Omega$  that respects the geometric structure induced by  $d_\Omega$ .

The pioneering work of Cuturi (2013) further proposes to add an entropy regularization term to the standard Kantorovich objective. This relaxation dramatically reduces the cost of computing the OT distance between discrete measures  $\mu, \nu$  from a super-cubic to a quadratic runtime complexity. A crucial contributing factor that popularized entropy regularized OT in the context of deep learning is that the resulting Sinkhorn divergence is differentiable. As a result, a number of works show how the Sinkhorn divergence can be used as a loss function to train a machine learning model (Luisse et al., 2018; Genevay et al., 2018; Klatt et al., 2020; Ablin et al., 2020). In this work, we consider an extension of this idea via the more general *Sinkhorn operator*. This operator is defined as the mapping  $(C, \mathbf{a}, \mathbf{b}) \mapsto P$  of a cost matrix  $C \in \mathbb{R}^{m \times n}$  and the source and target capacities  $\mathbf{a} \in \mathbb{R}^m, \mathbf{b} \in \mathbb{R}^n$  to the transportation plan  $P \in \mathbb{R}_+^{m \times n}$ . In comparison, the Sinkhorn divergence  $(\mathbf{a}, \mathbf{b}) \mapsto \langle C, P \rangle_F$  is a special case of this operator (Cuturi, 2013, Eq. (2)).

It is nowadays common practice to design a deep neural network as the composition of various mathematical functions or “layers”. Some of those layers may themselves be defined as the solution of an (inner) optimization problem, as long as the mapping induced by the inner optimization is well defined and differentiable. In order to train such a neural network, we need to solve a bilevel optimization problem. The most common approaches to that end are algorithmic unrolling of the inner

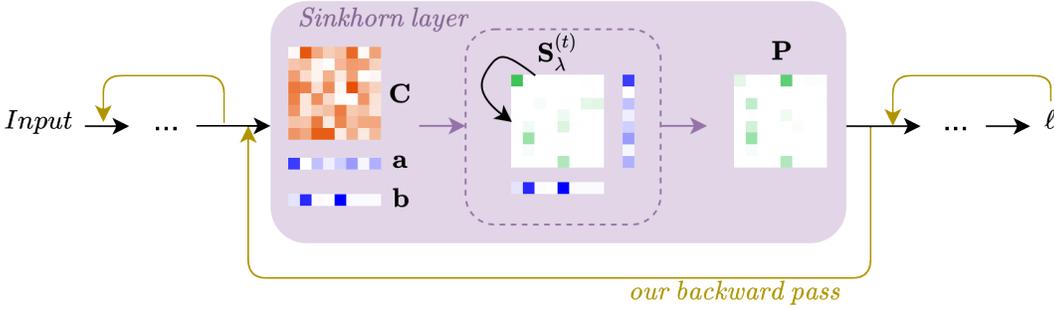


Figure 1: Overview of a typical workflow with an embedded Sinkhorn layer. We consider a neural network whose input are e.g. images, 3D point clouds, voxel grids, surface meshes, etc. The Sinkhorn layer maps the cost matrix  $C$  and marginals  $\mathbf{a}$ ,  $\mathbf{b}$  to the transportation plan  $P$  via iterative matrix scaling. In general, the whole network potentially contains learnable weights before and after the Sinkhorn layer. To compute the gradients ( $\nabla_C \ell$ ,  $\nabla_{\mathbf{a}} \ell$ ,  $\nabla_{\mathbf{b}} \ell$ ) during training, instead of using naive automatic differentiation, we propose a novel customised backward pass which computes gradients efficiently and in closed form, see Algorithm 1.

optimization loop and implicit differentiation (Amos & Kolter, 2017; Gould et al., 2019; Bai et al., 2019; Blondel et al., 2021). There is a vast literature on applications of combining the Sinkhorn operator with deep learning, see our related work discussion in Section 4. The majority of approaches unroll the forward optimization steps and apply automatic differentiation to compute gradients of the resulting iterative scheme. Since the computation graph for all iterations needs to be maintained, the resulting computation cost is often prohibitively high for GPU processing. To alleviate this cost, there are a number of approaches related to our work that apply implicit differentiation to specific subclasses of problems, like learning with a fixed cost matrix or immutable input capacities (Luise et al., 2018; Flamary et al., 2018; Campbell et al., 2020; Cuturi et al., 2020; Klatt et al., 2020). We can recover most such methods as special cases of our framework. In our work, we derive, from first principles, an efficient algorithm to compute gradients of a generic Sinkhorn layer see Fig. 1. Our contribution can be summarized as follows:

1. We leverage the implicit function theorem to compute gradients of the *Sinkhorn operator* in its most general form. We further use the Schur complement of the corresponding vector-Jacobian product (VJP) to construct an efficient backpropagation algorithm (Algorithm 1).
2. We provide theoretical guarantees for the accuracy of the resulting gradients, in dependency of the approximation error in the forward pass (Theorem 5).
3. Finally, we show that our algorithm can be used out of the box to enhance existing approaches that use a differentiable Sinkhorn layer. Plugging in our module often improves the quantitative results while using significantly less GPU memory.

## 2 BACKGROUND

**Optimal transport** Optimal transport enables us to compute the distance of two probability measures on the same domain  $\Omega \subset \mathbb{R}^d$ . In this work, we focus on the specific case of discrete probability measures  $\mu := \sum_{i=1}^m a_i \delta_{\mathbf{x}_i}$  and  $\nu := \sum_{j=1}^n b_j \delta_{\mathbf{y}_j}$ , defined over the sets of points  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  and  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ , where  $\delta_{\mathbf{x}_i}$  is the Dirac measure at  $\mathbf{x}_i$ . Such measures are fully characterized by the probability mass vectors  $\mathbf{a} \in \Delta_m$  and  $\mathbf{b} \in \Delta_n$  that lie on the probability simplex

$$\Delta_m = \{\mathbf{a} \in \mathbb{R}^m \mid a_i \geq 0, \mathbf{a}^\top \mathbf{1}_m = 1\}, \quad (1)$$

where  $\mathbf{1}_m \in \mathbb{R}^m$  is the vector of all ones. We can then define the distance between  $\mu$  and  $\nu$  as

$$d(\mu, \nu) := \min_{P \in \Pi(\mathbf{a}, \mathbf{b})} \langle P, C \rangle_F. \quad (2)$$

The transportation plan  $P \in \Pi(\mathbf{a}, \mathbf{b})$  determines a discrete probability measure on the product space  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\} \times \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ , whose marginal distributions coincide with  $\mu$  and  $\nu$ . Consequently,

$\mathbf{P}$  is contained in the transportation polytope  $\Pi(\mathbf{a}, \mathbf{b})$  defined as

$$\Pi(\mathbf{a}, \mathbf{b}) := \{\mathbf{P} \in \mathbb{R}_+^{m \times n} \mid \mathbf{P}\mathbf{1}_n = \mathbf{a}, \mathbf{P}^\top \mathbf{1}_m = \mathbf{b}\}. \quad (3)$$

The cost matrix  $\mathbf{C} \in \mathbb{R}^{m \times n}$  specifies the transportation cost from individual points  $\mathbf{x}_i$  to  $\mathbf{y}_j$ . Choosing  $\mathbf{C}_{i,j} := \|\mathbf{x}_i - \mathbf{y}_j\|_2^p$  for  $p \geq 1$ , e.g. yields the so-called Wasserstein distance  $d(\cdot, \cdot) = W_p^p(\cdot, \cdot)$ , see Villani (2003).

**Entropy regularization** Evaluating the distance  $d(\mu, \nu)$  in practice requires solving the linear assignment problem (LAP) from Equation 2. This can be done via specialized algorithms like the Hungarian algorithm (Kuhn, 1955) or the Auction algorithm (Bertsekas, 1979), as well as recent solvers (Rubner et al., 1997; Pele & Werman, 2009). However, most approaches are computationally heavy and slow in practice (Cuturi, 2013). A popular alternative is augmenting the LAP objective in Equation 2 with an additional entropy regularizer, giving rise to the *Sinkhorn operator*

$$S_\lambda(\mathbf{C}, \mathbf{a}, \mathbf{b}) := \arg \min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle_F - \lambda h(\mathbf{P}), \quad (4)$$

where  $\lambda > 0$ . The seminal work of Cuturi (2013) showed that the additional entropy regularization term  $h(\mathbf{P}) = -\sum_{i,j} P_{i,j}(\log P_{i,j} - 1)$  allows for an efficient minimization of Equation 4. Specifically, this can be done via a scheme of alternating Sinkhorn projections

$$\begin{aligned} \mathbf{S}_\lambda^{(0)} &:= \exp\left(-\frac{1}{\lambda}\mathbf{C}\right), \quad \text{and} \\ \mathbf{S}_\lambda^{(t+1)} &:= \mathcal{T}_c(\mathcal{T}_r(\mathbf{S}_\lambda^{(t)})). \end{aligned} \quad (5)$$

The operators  $\mathcal{T}_c(\mathbf{S}) := \mathbf{S} \odot (\mathbf{1}_m \mathbf{1}_m^\top \mathbf{S}) \odot (\mathbf{1}_m \mathbf{b}^\top)$  and  $\mathcal{T}_r(\mathbf{S}) := \mathbf{S} \odot (\mathbf{S} \mathbf{1}_n \mathbf{1}_n^\top) \odot (\mathbf{a} \mathbf{1}_n^\top)$  correspond to renormalizations of the columns and rows of  $\mathbf{S}_\lambda^{(t)}$ , where  $\odot$  denotes the Hadamard product and  $\odot$  denotes element-wise division. As shown by Cuturi (2013), in the limit this scheme converges to a minimizer  $\mathbf{S}_\lambda^{(t)} \xrightarrow{t \rightarrow \infty} \mathbf{S}_\lambda$  of Equation 4. In practice, we can use a finite number of iterations  $\tau \in \mathbb{N}$  to achieve a sufficiently small residual.

## 3 METHOD

### 3.1 PROBLEM FORMULATION

The integration of the Sinkhorn operator from Equation 4 into deep neural networks has become popular for solving a wide range of practical tasks, see our discussion in Section 4. A major contributing factor is that the entropy regularization makes the mapping  $S_\lambda : \mathbb{R}^{m \times n} \times \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$  differentiable. To allow for first-order-optimization, we need to compute

$$\begin{aligned} (\mathbf{C}, \mathbf{a}, \mathbf{b}) &\mapsto \mathbf{P}^* := S_\lambda(\mathbf{C}, \mathbf{a}, \mathbf{b}) \quad \text{and} & (6) \\ \nabla_{\mathbf{P}} \ell &\mapsto (\nabla_{\mathbf{C}} \ell, \nabla_{\mathbf{a}} \ell, \nabla_{\mathbf{b}} \ell), & (7) \end{aligned}$$

which denote the forward pass and the backpropagation of gradients, respectively. Those expressions arise in the context of a typical workflow within a deep neural network with a scalar loss  $\ell$  and learnable parameters before and/or after the Sinkhorn operator  $S_\lambda$ , see Fig. 1 for an overview.

A common strategy is to replace the exact forward pass  $S_\lambda(\mathbf{C}, \mathbf{a}, \mathbf{b})$  in Equation 6 by the approximate solution  $\mathbf{S}_\lambda^{(\tau)}$  from Equation 5. Like the original solution in Equation 4,  $\mathbf{S}_\lambda^{(\tau)}$  is differentiable w.r.t.  $(\mathbf{C}, \mathbf{a}, \mathbf{b})$ . Moreover, the mapping  $(\mathbf{C}, \mathbf{a}, \mathbf{b}) \mapsto \mathbf{S}_\lambda^{(\tau)}$  consists of a small number of matrix scaling operations that can be implemented in a few lines of code, see Equation 5.

For the backward pass in Equation 7, an overwhelming number of applications nowadays rely on automatic differentiation in environments like PyTorch (Paszke et al., 2019) or TensorFlow (Abadi et al., 2016). While this solution is straightforward and convenient, it has fundamental shortcomings: First of all, the computation time for the backward pass increases with the number of Sinkhorn iterations  $\tau$ . Moreover, the computation graph needs to be maintained for all  $\tau$  matrix scaling iterations. For many applications, this makes it infeasible or limits the resolution since the GPU memory is often the computational bottleneck in practice. Instead, to address these shortcomings, in the next section we utilize implicit differentiation in order to introduce a closed-form expression for computing the gradients of  $\ell$  w.r.t both the marginals  $\mathbf{a}, \mathbf{b}$  and the cost matrix  $\mathbf{C}$ .

### 3.2 BACKWARD PASS VIA IMPLICIT DIFFERENTIATION

The goal of this section is to derive the main result stated in Theorem 3, which is the key motivation of our algorithm in Sec. 3.3. To this end, we start by reframing the optimization problem in Equation 4 in terms of its Karush–Kuhn–Tucker (KKT) conditions:

**Lemma 1.** *The transportation plan  $\mathbf{P}^*$  is a global minimum of Equation 4 iff*

$$\mathcal{K}(\mathbf{c}, \mathbf{a}, \mathbf{b}, \mathbf{p}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) := \begin{bmatrix} \mathbf{c} + \lambda \log(\mathbf{p}^*) + \mathbb{1}_n \otimes \boldsymbol{\alpha}^* + \boldsymbol{\beta}^* \otimes \mathbb{1}_m \\ (\mathbb{1}_n^\top \otimes \mathbf{I}_m) \mathbf{p}^* - \mathbf{a} \\ (\mathbf{I}_n \otimes \mathbb{1}_m^\top) \mathbf{p}^* - \mathbf{b} \end{bmatrix} = \mathbf{0}_l, \quad (8)$$

where  $l := mn + m + n$ . Here,  $\boldsymbol{\alpha}^* \in \mathbb{R}^m$  and  $\boldsymbol{\beta}^* \in \mathbb{R}^n$  are the dual variables corresponding to the two equality constraints in Equation 3. We further define  $\mathbf{c}, \mathbf{p}^* \in \mathbb{R}^{mn}$  as the vectorized versions of  $\mathbf{C}, \mathbf{P}^* \in \mathbb{R}^{m \times n}$ , respectively, and assume  $\log(p) := -\infty, p \leq 0$ .

*Proof.* The function  $\mathcal{K}$  contains the KKT conditions corresponding to the optimization problem in Equation 4. The proposed identity therefore follows directly from the (strict) convexity of Equation 4, see Erlander & Stewart (1990). Apart from the two equality constraints, Equation 3 contains additional inequality constraints  $P_{i,j} \geq 0$ . Those are however inactive and can be dropped, because the entropy term in Equation 4 invariably yields transportation plans in the interior of the positive orthant  $P_{i,j} > 0$ , see Peyré et al. (2019, p. 68).  $\square$

Establishing this identity is an important first step towards computing a closed-form gradient for the backward pass in Equation 7. It reframes the optimization problem in Equation 4 as a root-finding problem  $\mathcal{K}(\cdot) = \mathbf{0}$ . In the next step, this then allows us to explicitly construct the derivative of the Sinkhorn operator  $S_\lambda(\cdot)$  via implicit differentiation:

**Lemma 2.** *The KKT conditions in Equation 8 implicitly define a continuously differentiable function  $(\mathbf{c}, \mathbf{a}, \tilde{\mathbf{b}}) \mapsto (\mathbf{p}, \boldsymbol{\alpha}, \tilde{\boldsymbol{\beta}})$  with the Jacobian matrix<sup>1</sup>*

$$\mathbf{J} := \frac{\partial [\mathbf{p}; \boldsymbol{\alpha}; \tilde{\boldsymbol{\beta}}]}{\partial [\mathbf{c}; -\mathbf{a}; -\tilde{\mathbf{b}}]} = - \underbrace{\begin{bmatrix} \lambda \text{diag}(\mathbf{p})^{-1} & \tilde{\mathbf{E}} \\ \mathbf{E}^\top & \mathbf{0} \end{bmatrix}^{-1}}_{:=\mathbf{K}} \in \mathbb{R}^{(l-1) \times (l-1)}. \quad (9)$$

Note that the last entry of  $\tilde{\mathbf{b}} := \mathbf{b}_{-n}$  and  $\tilde{\boldsymbol{\beta}} := \boldsymbol{\beta}_{-n}$  is removed. This is due to a surplus degree of freedom in the equality conditions from Equation 3, see part (b) of the proof. Likewise, for

$$\mathbf{E} = [\mathbb{1}_n \otimes \mathbf{I}_m \quad \mathbf{I}_n \otimes \mathbb{1}_m] \in \mathbb{R}^{mn \times (m+n)}, \quad (10)$$

the corresponding last column is removed  $\tilde{\mathbf{E}} := \mathbf{E}_{:, -(m+n)}$ .

*Sketch of the proof.* The basis for this proof is leveraging the implicit function theorem (IFT), which we show in two parts, (a) and (b). See Appendix C.1 for the full proof.

- (a) We first establish that the matrix  $\mathbf{K}$  is invertible. We can then prove by direct computation that  $\mathbf{K}$  corresponds to the partial derivative of  $\mathcal{K}$  from Equation 8 wrt.  $[\mathbf{p}; \boldsymbol{\alpha}; \tilde{\boldsymbol{\beta}}]$ . The identity from Equation 9 then follows directly from the IFT, since the partial derivative of  $\mathcal{K}$  w.r.t.  $[\mathbf{c}; -\mathbf{a}; -\tilde{\mathbf{b}}]$  is simply  $\frac{\partial \mathcal{K}}{\partial [\mathbf{c}; -\mathbf{a}; -\tilde{\mathbf{b}}]} = \mathbf{I}_l$ .
- (b) To show why the last condition from  $\mathbf{E}$  is removed, we first assert that the kernel  $\ker(\mathbf{E}^\top)$  is  $(m-1)(n-1)$  dimensional. The rank-nullity theorem then yields that the subspace spanned by the columns of  $\mathbf{E}$  is  $mn - (m-1)(n-1) = m+n-1$  dimensional. Removing the last condition  $\tilde{\mathbf{E}} := \mathbf{E}_{:, -(m+n)}$  ensures that the  $m+n-1$  columns of  $\tilde{\mathbf{E}}$  are linearly independent and  $\mathbf{K}$  is invertible.  $\square$

<sup>1</sup>For brevity we use the short hand notation  $[\mathbf{v}; \mathbf{u}] := [\mathbf{v}^\top, \mathbf{u}^\top]^\top$  for stacking vectors  $\mathbf{v}, \mathbf{u}$  vertically.

In principle, we can use Lemma 2 directly to compute the backward pass from Equation 7. However, the computational cost of inverting the matrix  $\mathbf{K}$  in Equation 9 is prohibitive. In fact, even storing the Jacobian  $\mathbf{J}$  in the working memory of a typical machine is problematic, since it is a dense matrix with  $\mathcal{O}(mn)$  rows and columns, where  $m, n > 1000$  in practice. Instead, we observe that computing Equation 7 merely requires us to compute vector-Jacobian products (VJP) of the form  $\mathbf{v}^\top \mathbf{J}$ . The main results from this section can therefore be summarized as follows:

**Theorem 3** (Backward pass). *For  $\mathbf{P} = \mathbf{P}^*$ , the backward pass in Equation 7 can be computed in closed form by solving the following linear system:*

$$\begin{bmatrix} \lambda \text{diag}(\mathbf{p})^{-1} & \tilde{\mathbf{E}} \\ \mathbf{E}^\top & \mathbf{0} \end{bmatrix} \begin{bmatrix} \nabla_{\mathbf{c}} \ell \\ -\nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell \end{bmatrix} = \begin{bmatrix} -\nabla_{\mathbf{p}} \ell \\ \mathbf{0} \end{bmatrix}. \quad (11)$$

*Proof.* This identity follows trivially from Lemma 2 by applying the chain rule

$$\nabla_{[\mathbf{c}; -\mathbf{a}; -\tilde{\mathbf{b}}]} \ell = \left( \frac{\partial [\mathbf{p}; \boldsymbol{\alpha}; \tilde{\boldsymbol{\beta}}]}{\partial [\mathbf{c}; -\mathbf{a}; -\tilde{\mathbf{b}}]} \right)^\top \nabla_{[\mathbf{p}; \boldsymbol{\alpha}; \tilde{\boldsymbol{\beta}}]} \ell = -\mathbf{K}^{-1} \nabla_{[\mathbf{p}; \boldsymbol{\alpha}; \tilde{\boldsymbol{\beta}}]} \ell = \mathbf{K}^{-1} \begin{bmatrix} -\nabla_{\mathbf{p}} \ell \\ \mathbf{0} \end{bmatrix}. \quad (12)$$

The last equality holds, since the loss  $\ell$  does not depend on the dual variables  $\boldsymbol{\alpha}$  and  $\tilde{\boldsymbol{\beta}}$ , see Fig. 1.  $\square$

### 3.3 ALGORITHM

In the previous section, we derived a closed-form expression of the Sinkhorn backward pass in Theorem 3. This requires solving the sparse linear system in Equation 11, which has  $\mathcal{O}(mn)$  rows and columns, and thus amounts to a worst-case complexity of  $\mathcal{O}(m^3n^3)$ . We can further reduce the computation cost by exploiting the specific block structure of  $\mathbf{K}$ , which leads to our algorithm:

---

#### Algorithm 1: Sinkhorn operator backward

---

**Input** :  $\nabla_{\mathbf{p}} \ell, \mathbf{P}, \mathbf{a}, \mathbf{b}$   
**Output**:  $\nabla_{\mathbf{c}} \ell, \nabla_{\mathbf{a}} \ell, \nabla_{\mathbf{b}} \ell$

- 1  $\mathbf{T} \leftarrow \mathbf{P} \odot \nabla_{\mathbf{p}} \ell$ .
- 2  $\tilde{\mathbf{T}} \leftarrow \mathbf{T}_{:, -n}, \tilde{\mathbf{P}} \leftarrow \mathbf{P}_{:, -n} \in \mathbb{R}^{m \times n-1}$ .
- 3  $\mathbf{t}^{(a)} \leftarrow \mathbf{T} \mathbb{1}_n, \tilde{\mathbf{t}}^{(b)} \leftarrow \tilde{\mathbf{T}}^\top \mathbb{1}_m$ .
- 4  $\begin{bmatrix} \nabla_{\mathbf{a}} \ell \\ \nabla_{\tilde{\mathbf{b}}} \ell \end{bmatrix} \leftarrow \begin{bmatrix} \text{diag}(\mathbf{a}) & \tilde{\mathbf{P}} \\ \tilde{\mathbf{P}}^\top & \text{diag}(\tilde{\mathbf{b}}) \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{t}^{(a)} \\ \tilde{\mathbf{t}}^{(b)} \end{bmatrix}$ .
- 5  $\nabla_{\mathbf{b}} \ell \leftarrow [\nabla_{\tilde{\mathbf{b}}} \ell; \mathbf{0}]$ .
- 6  $\mathbf{U} \leftarrow \nabla_{\mathbf{a}} \ell \mathbb{1}_n^\top + \mathbb{1}_m \nabla_{\mathbf{b}} \ell^\top$ .
- 7  $\nabla_{\mathbf{c}} \ell \leftarrow -\lambda^{-1}(\mathbf{T} - \mathbf{P} \odot \mathbf{U})$ .

---

See Appendix A for a PyTorch implementation of this algorithm. We now show that the resulting gradients  $\nabla_{\mathbf{c}} \ell, \nabla_{\mathbf{a}} \ell, \nabla_{\mathbf{b}} \ell$  from Algorithm 1 are indeed solutions of the linear system in Theorem 3.

**Theorem 4.** *Let  $\mathbf{a}, \mathbf{b}$  be two input marginals and  $\mathbf{P} = \mathbf{P}^*$  the transportation plan resulting from the forward pass in Equation 6, then Algorithm 1 solves the backward pass defined in Equation 7.*

*Sketch of the proof.* The main idea of this proof is showing that Algorithm 1 yields a solution  $\nabla_{[\mathbf{c}; \mathbf{a}; \tilde{\mathbf{b}}]} \ell$  of the linear system from Equation 11. To that end, we leverage the Schur complement trick which yields the following two expressions:

$$\nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell = (\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p}) \tilde{\mathbf{E}})^{-1} \tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p}) \nabla_{\mathbf{p}} \ell. \quad (13a)$$

$$\nabla_{\mathbf{c}} \ell = -\lambda^{-1}(\text{diag}(\mathbf{p}) \nabla_{\mathbf{p}} \ell - \text{diag}(\mathbf{p}) \tilde{\mathbf{E}} \nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell). \quad (13b)$$

In Appendix C.2 we further show that these two identities in their vectorized form are equivalent to Algorithm 1 in matrix notation.  $\square$

### 3.4 PRACTICAL CONSIDERATIONS

**Error bounds** Theorem 4 proves that Algorithm 1 computes the exact gradients  $\nabla_{\mathbf{C}}\ell$ ,  $\nabla_{\mathbf{a}}\ell$ ,  $\nabla_{\mathbf{b}}\ell$ , given that  $\mathbf{P} = \mathbf{P}^*$  is the exact solution of Equation 4. In practice, the operator  $S_\lambda$  in Equation 6 is replaced by the Sinkhorn approximation  $S_\lambda^{(\tau)}$  from Equation 5 for a fixed, finite  $\tau \in \mathbb{N}$ . This small discrepancy in the approximation  $\mathbf{P} = S_\lambda^{(\tau)} \approx \mathbf{P}^*$  propagates to the backward pass as follows:

**Theorem 5** (Error bounds). *Let  $\mathbf{P}^* := S_\lambda(\mathbf{C}, \mathbf{a}, \mathbf{b})$  be the exact solution of Equation 4 and let  $\mathbf{P}^{(\tau)} := S_\lambda^{(\tau)}$  be the Sinkhorn estimate from Equation 5. Further, let  $\sigma_+, \sigma_-, C_1, C_2, \epsilon > 0$ , s.t.  $\|\mathbf{P}^* - \mathbf{P}^{(\tau)}\|_F < \epsilon$  and that for all  $\mathbf{P}$  for which  $\|\mathbf{P} - \mathbf{P}^*\|_F < \epsilon$  we have  $\min_{i,j} P_{i,j} \geq \sigma_-$ ,  $\max_{i,j} P_{i,j} \leq \sigma_+$  and the loss  $\ell$  has bounded derivatives  $\|\nabla_{\mathbf{p}}\ell\|_2 \leq C_1$  and  $\|\nabla_{\mathbf{p}}^2\ell\|_F \leq C_2$ . For  $\kappa = \|\tilde{\mathbf{E}}^\dagger\|_2$ , where  $\tilde{\mathbf{E}}^\dagger$  indicates the Moore-Penrose inverse of  $\tilde{\mathbf{E}}$ , the difference between the gradients  $\nabla_{\mathbf{C}}\ell^*$ ,  $\nabla_{\mathbf{a}}\ell^*$ ,  $\nabla_{\mathbf{b}}\ell^*$  of the exact  $\mathbf{P}^*$  and the gradients  $\nabla_{\mathbf{C}}\ell^{(\tau)}$ ,  $\nabla_{\mathbf{a}}\ell^{(\tau)}$ ,  $\nabla_{\mathbf{b}}\ell^{(\tau)}$  of the approximate  $\mathbf{P}^{(\tau)}$ , obtained via Algorithm 1, satisfy*

$$\|\nabla_{[\mathbf{a};\mathbf{b}]} \ell^* - \nabla_{[\mathbf{a};\mathbf{b}]} \ell^{(\tau)}\|_F \leq \kappa \sqrt{\frac{\sigma_+}{\sigma_-}} \left( \frac{1}{\sigma_-} C_1 + C_2 \right) \|\mathbf{P}^* - \mathbf{P}^{(\tau)}\|_F, \quad \text{and} \quad (14a)$$

$$\|\nabla_{\mathbf{C}} \ell^* - \nabla_{\mathbf{C}} \ell^{(\tau)}\|_F \leq \lambda^{-1} \sigma_+ \left( \frac{1}{\sigma_-} C_1 + C_2 \right) \|\mathbf{P}^* - \mathbf{P}^{(\tau)}\|_F. \quad (14b)$$

We provide a proof in Appendix C.3, as well as an empirical evaluation in Appendix B.1.

**Computation cost** In comparison to automatic differentiation (AD), the computation cost of Algorithm 1 is independent of the number of Sinkhorn iterations  $\tau$ . For square matrices,  $m = n$ , the runtime and memory complexities of AD are  $\mathcal{O}(\tau n^2)$ . On the other hand, our approach has a runtime and memory complexity of  $\mathcal{O}(n^3)$  and  $\mathcal{O}(n^2)$  respectively. We show empirical comparisons between the two approaches in Sec. 5.1. Another compelling feature of our approach is that none of the operations in Algorithm 1 explicitly convert the matrices  $\mathbf{P}, \nabla_{\mathbf{p}}\ell, \nabla_{\mathbf{C}}\ell, \dots \in \mathbb{R}^{m \times n}$  into their vectorized form  $\mathbf{p}, \nabla_{\mathbf{p}}\ell, \nabla_{\mathbf{C}}\ell, \dots \in \mathbb{R}^{mn}$ . This makes it computationally more efficient since GPU processing favors small, dense matrix operations over the large, sparse linear system in Equation 11.

**Marginal probability invariance** As discussed in Lemma 2, the last element of  $\tilde{\mathbf{b}}$  needs to be removed to make  $\mathbf{K}$  invertible. However, setting the last entry of the gradient  $\nabla_{b_n} \ell = 0$  to zero still yields exact gradients: By definition, the full marginal  $\mathbf{b}$  is constrained to the probability simplex  $\Delta_n$ , see Equation 1. In practice, we apply an a priori softmax to  $\mathbf{b}$  (and analogously  $\mathbf{a}$ ). For some applications,  $\mathbf{b}$  can be assumed to be immutable, if we only want to learn the cost matrix  $\mathbf{C}$  and not the marginals  $\mathbf{a}$  and  $\mathbf{b}$ . Overall, this means that the gradient of  $\mathbf{b}$  is effectively indifferent to constant offsets of all entries, and setting  $\nabla_{b_n} \ell = 0$  does not contradict the statement of Theorem 3.

## 4 RELATED WORK

There is a vast literature on computational optimal transport (OT) (Villani, 2003; Peyré et al., 2019). In the following, we provide an overview of related machine learning applications. Our approach is based on entropy regularized optimal transport pioneered by Cuturi (2013). The resulting differentiable Sinkhorn divergence can be used as a loss function for training machine learning models (Frogner et al., 2015; Feydy et al., 2019; Chizat et al., 2020). To allow for first-order optimization, two common approaches for computing gradients are implicit differentiation (Luise et al., 2018; Cuturi et al., 2020; Klatt et al., 2020) and automatic differentiation (Genevay et al., 2018; Ablin et al., 2020). Relevant applications of the Sinkhorn divergence include computing Wasserstein barycenters (Cuturi & Doucet, 2014; Solomon et al., 2015; Luise et al., 2019), dictionary learning (Schmitz et al., 2018), as well as using a geometrically meaningful loss function for autoencoders (Patrini et al., 2020) or generative adversarial networks (GAN) (Genevay et al., 2018; Salimans et al., 2018).

More recently, several approaches emerged that use the Sinkhorn operator as a differentiable transportation layer in a neural network. Potential applications include permutation learning (Santa Cruz et al., 2017; Mena et al., 2018), ranking (Adams & Zemel, 2011; Cuturi et al., 2019), sorting via reinforcement learning (Emami & Ranka, 2018), discriminant analysis (Flamary et al., 2018) and

computing matchings between images (Sarlin et al., 2020), point clouds (Yew & Lee, 2020; Yang et al., 2020; Liu et al., 2020) or triangle meshes (Eisenberger et al., 2020; Pai et al., 2021). Most of these approaches rely on automatic differentiation of the Sinkhorn algorithm to address the resulting bilevel optimization problem. In our work, we follow the recent trend of using implicit differentiation for the inner optimization layer (Amos & Kolter, 2017; Gould et al., 2019; Blondel et al., 2021). Similar ideas are pursued in (Luise et al., 2018; Klatt et al., 2020; Campbell et al., 2020) for the special subcases of learning with loss functions that are independent of the cost matrix or the marginals. Other approaches compute the input cost matrix via Bayesian inverse modeling (Stuart & Wolfram, 2020) or smooth the OT linear assignment problem (LAP) directly (Pogan et al., 2019).

## 5 EXPERIMENTS

We now provide practical evidence for the merits of our Algorithm 1. In Sec. 5.1, we compare its computation cost to automatic differentiation (AD). In Sec. 5.2 and Sec. 5.3, we show results on two common classes of applications where we want to learn the marginals  $\mathbf{a}$  and the cost matrix  $\mathbf{C}$  respectively. For all experiments, we assume a fixed GPU memory (VRAM) budget of 24GB – any setting that exceeds this limit is deemed out of memory (OOM).

### 5.1 COMPUTATION COST

We empirically compare the computation cost of our algorithm with the standard automatic differentiation approach, see Fig. 2. All results were computed on a single NVIDIA Quadro RTX 8000 graphics card. In practice, the computation cost of both approaches primarily depends on the parameters  $m, n, \tau$ . It is for the most part indifferent to other hyperparameters and the actual values of  $\mathbf{C}, \mathbf{a}, \mathbf{b}$ . We therefore use random (log normal distributed) cost matrices  $\ln \mathbf{C}_{i,j} \sim \mathcal{N}(0, 1)$  and uniform marginals  $\mathbf{a} = \mathbf{b} = \frac{1}{n} \mathbb{1}_n$  with  $m = n \in \{10, 100, 1000\}$ . For each setting, we report the cost of the forward and backward pass averaged over 1k iterations. Depending on  $m, n$ , our approach is faster for  $\tau \gtrsim 40, 50, 90$  iterations. Note that our backward pass is independent of the number of forward iterations  $\tau$ . Finally, the memory requirements are dramatically higher for AD, since it needs to maintain the computation graph of all  $\tau$  forward iterations. In practice, this often limits the admissible batch size or input resolution, see Sec. 5.2 and Sec. 5.3.

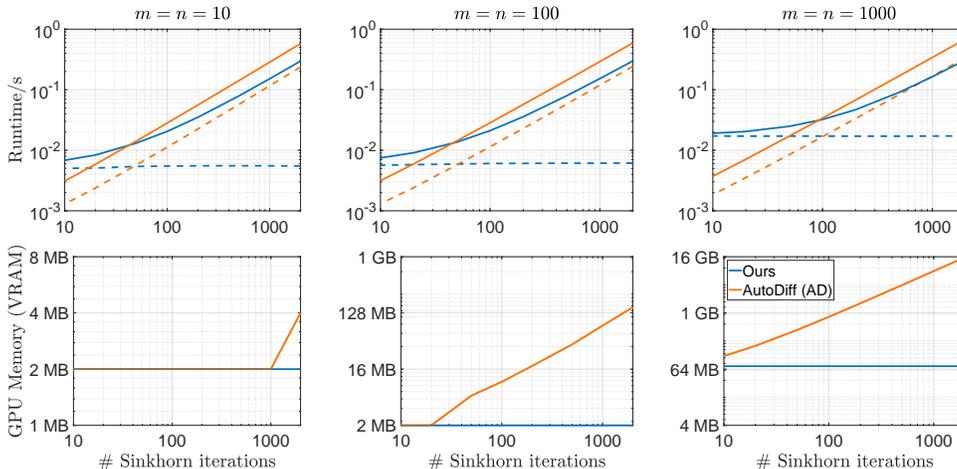


Figure 2: **Computational complexity.** We compare the runtime per iteration (top row) and GPU memory requirements (bottom row) of our approach (blue) and automatic differentiation (orange). We consider a broad range of settings with quadratic cost matrices of size  $m = n \in \{10, 100, 1000\}$  and  $\tau \in [10, 2000]$  Sinkhorn iterations. For the runtime, we show both the total time (solid lines) and the time of only the backward pass (dashed lines). Both ours and AD were implemented in the PyTorch (Paszke et al., 2019) framework, where memory is allocated in discrete units, which leads to a large overlap for the minimum allocation size of 2MB (bottom row, left plot).

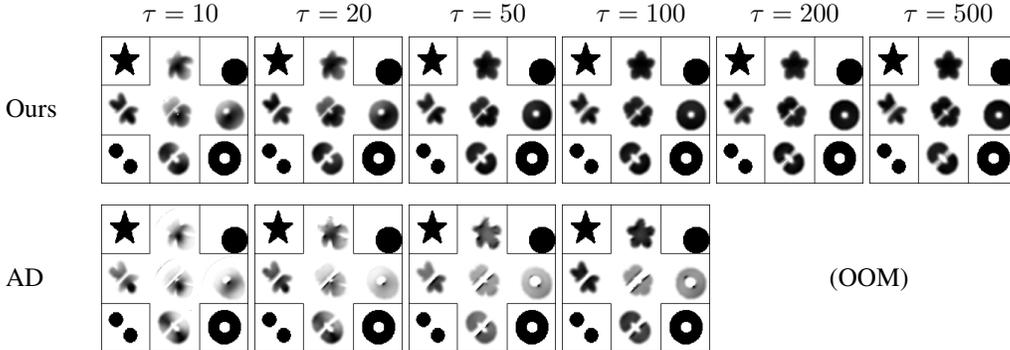


Figure 3: **Wasserstein barycenter.** A comparison between our method (top row) and automatic differentiation (bottom row) on the application of image barycenter computation. In each cell, we show 5 centroids of 4 input images (corners) with bilinear interpolation weights. The predictions of our approach are more stable, even for very few Sinkhorn iterations  $\tau$ . Moreover, AD is out of memory for  $\tau \geq 200$ . Here, the input images have a resolution of  $n = 64^2$  and we set  $\lambda = 0.002$ .

## 5.2 WASSERSTEIN BARYCENTERS

Barycenter computation is a standard example of an inverse problem based on the Sinkhorn operator in Equation 4. The main idea is to interpolate between a collection of objects  $\{b_1, \dots, b_k\} \subset \mathbb{R}^n$  as a convex combination with weights that lie on the probability simplex  $w \in \Delta_n$ , see Equation 1. Specifically, the goal is to optimize for the marginal  $a^*$  via the expression

$$a^* := \arg \min_{a \in \Delta_n} \sum_{i=1}^k w_i d(a, b_i) \quad \text{with} \quad d(a, b) := \min_{P \in \Pi(a, b)} \langle P, D \rangle_F - \lambda h(P), \quad (15)$$

where  $D \in \mathbb{R}^{n \times n}$  denotes the squared pairwise distance matrix between the domains of  $a$  and  $b$ . We use the Adam optimizer (Kingma & Ba, 2014) for the outer optimization in Equation 15. The inner optimization is a special case of Equation 4. Overall, Equation 15 allows us to compute geometrically meaningful interpolations in arbitrary metric spaces. We consider the explicit tasks of interpolating between images in Fig. 3 and functions on manifolds in Appendix B.2. Note that there are a number of specialized algorithms that minimize Equation 15 in a highly efficient manner (Cuturi & Doucet, 2014; Solomon et al., 2015; Luise et al., 2019). We mainly consider the barycenter problem a useful toy example to assess the stability of our algorithm in comparison to AD. On the other hand, the interpretation as a generic optimization problem is overall more flexible and allows us to trivially extend it to related tasks like image clustering, see Appendix B.2.

## 5.3 PERMUTATION LEARNING AND MATCHING

**Number sorting** The Sinkhorn operator is nowadays a standard tool to parameterize approximate, learnable permutations within a neural network, see the second paragraph of our related work discussion in Sec. 4. One work that clearly demonstrates the effectiveness of this approach is the Gumbel-Sinkhorn (GS) method (Mena et al., 2018). The main idea is to learn the natural ordering of sets of input elements  $\{x_1, \dots, x_n\}$ , see Appendix B.3 for more details. Here, we consider the concrete example of learning to sort real numbers from the unit interval  $x_i \in [0, 1]$  for  $n \in \{200, 500, 1000\}$  numbers. We insert our Sinkhorn module in the GS network and compare the performance with the vanilla GS method in Fig. 4. Without further modifications, our method significantly decreases the error at test time, defined as the proportion of incorrectly sorted elements.

**Point cloud registration** A number of recent methods use the Sinkhorn operator as a differentiable, bijective matching layer for deep learning (Sarlin et al., 2020; Yew & Lee, 2020; Yang et al., 2020; Liu et al., 2020; Eisenberger et al., 2020). Here, we consider the concrete application of rigid point cloud registration (Yew & Lee, 2020) and show that we can improve the performance with our backward algorithm, see Table 1. While our results on the clean test data are comparable but

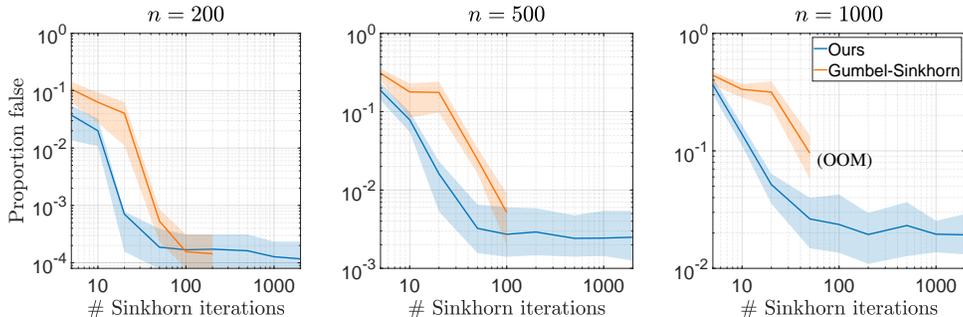


Figure 4: **Number sorting.** We show that we can improve the Gumbel-Sinkhorn method (Mena et al., 2018) directly with our Sinkhorn module. Specifically, we consider the task of permutation learning to sort random number sequences of length  $n \in \{200, 500, 1000\}$ , see Mena et al. (2018, Sec 5.1) for more details. We replace AD in the GS network with our module (blue curves) and compare the obtained results to the vanilla GS architecture (orange curves). Our approach yields more accurate permutations while using much less computational resources – GS is out of memory for  $\tau > 200, 100, 50$  forward iterations, respectively. For all settings, we show the mean proportion of correct test set predictions (solid lines), as well as the 10 and 90 percentiles (filled areas).

slightly worse than the vanilla RPM-Net (Yew & Lee, 2020), our module generalizes more robustly to partial and noisy observations. This indicates that, since computing gradients with our method is less noisy than AD, it helps to learn a robust matching policy with less overfitting. We provide more details on the RPM-Net baseline and qualitative comparisons in Appendix B.3.

## 6 DISCUSSION

Our experiments clearly demonstrate that, for a broad range of applications, combining our Sinkhorn module with existing approaches often improves the performance. On the other hand, AD has an empirically faster runtime for very few Sinkhorn iterations  $\tau \approx 10$  and high  $n \geq 1000$ . In our experiments, however, we see a clear trend that, to a certain point, increasing  $\tau$  leads to more accurate results. Choosing  $\tau$  is generally subject to a trade-off between the computation cost and accuracy. The main advantage of implicit differentiation is that it proves to be much more scalable than AD, since the backward pass is independent of  $\tau$ . For the same reason, it also has a dramatically lower GPU memory demand, see Fig. 2. We therefore believe that the theoretical and empirical insights we provide have the potential to open up new avenues for future research, for which naive backpropagation schemes are computationally intractable or not sufficiently accurate.

		clean data	partial			noisy		
			90%	80%	70%	$\sigma = 0.001$	$\sigma = 0.01$	$\sigma = 0.1$
Rot. MSE	RPM	0.2329	63.7670	69.6706	74.4227	37.2575	50.3618	66.1996
	Ours	0.3231	10.3078	24.5810	40.7793	1.7912	2.2110	4.0208
Trans. MSE	RPM	0.0014	0.2659	0.3079	0.3462	0.1603	0.2100	0.2740
	Ours	0.0033	0.0797	0.1582	0.2420	0.0149	0.0181	0.0303
Chamf. dist.	RPM	0.0005	4.3413	4.6829	4.9581	2.2077	3.0492	4.6935
	Ours	0.0054	0.5498	1.4291	2.2080	0.0783	0.1237	0.4562

Table 1: **Point cloud registration.** We compare the quantitative performance of RPM-Net (Yew & Lee, 2020) and our approach on ModelNet40 (Wu et al., 2015). The two architectures are identical except for the altered Sinkhorn module. For all results, we follow the training protocol described in Yew & Lee (2020, Sec. 6). Moreover, we assess the ability of the obtained networks to generalize to partial and noisy inputs at test time. For the former, we follow Yew & Lee (2020, Sec. 6.6) and remove up to 70% of the input point clouds from a random half-space. For the noisy test set, we add Gaussian white noise  $\mathcal{N}(0, \sigma)$  with different variances  $\sigma \in \{0.001, 0.01, 0.1\}$ . For all settings, we report the rotation and translation errors, as well as the Chamfer distance to the reference surface.

## REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*, pp. 265–283, 2016.
- Pierre Ablin, Gabriel Peyré, and Thomas Moreau. Super-efficiency of automatic differentiation for functions defined as a minimum. In *International Conference on Machine Learning*, pp. 32–41. PMLR, 2020.
- Ryan Prescott Adams and Richard S Zemel. Ranking via sinkhorn propagation. *arXiv preprint arXiv:1106.1925*, 2011.
- Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. In *International Conference on Machine Learning*, pp. 136–145. PMLR, 2017.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. Deep equilibrium models. *Advances in Neural Information Processing Systems*, 32:690–701, 2019.
- Dimitri P Bertsekas. A distributed algorithm for the assignment problem. *Lab. for Information and Decision Systems Working Paper, MIT*, 1979.
- Mathieu Blondel, Quentin Berthet, Marco Cuturi, Roy Frostig, Stephan Hoyer, Felipe Llinares-López, Fabian Pedregosa, and Jean-Philippe Vert. Efficient and modular implicit differentiation. *arXiv preprint arXiv:2105.15183*, 2021.
- Ethan D Bolker. Transportation polytopes. *Journal of Combinatorial Theory, Series B*, 13(3):251–262, 1972.
- Dylan Campbell, Liu Liu, and Stephen Gould. Solving the blind perspective-n-point problem end-to-end with robust differentiable geometric optimization. In *European Conference on Computer Vision*, pp. 244–261. Springer, 2020.
- Lenaïc Chizat, Pierre Roussillon, Flavien Léger, François-Xavier Vialard, and Gabriel Peyré. Faster wasserstein distance estimation with the sinkhorn divergence. *Advances in Neural Information Processing Systems*, 33, 2020.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pp. 2292–2300, 2013.
- Marco Cuturi and Arnaud Doucet. Fast computation of wasserstein barycenters. In *International conference on machine learning*, pp. 685–693. PMLR, 2014.
- Marco Cuturi, Olivier Teboul, and Jean-Philippe Vert. Differentiable ranking and sorting using optimal transport. *Advances in Neural Information Processing Systems*, 32, 2019.
- Marco Cuturi, Olivier Teboul, Jonathan Niles-Weed, and Jean-Philippe Vert. Supervised quantile normalization for low-rank matrix approximation. *arXiv preprint arXiv:2002.03229*, 2020.
- Marvin Eisenberger, Aysim Toker, Laura Leal-Taixé, and Daniel Cremers. Deep shells: Unsupervised shape correspondence with optimal transport. In *Advances in Neural Information Processing Systems*, volume 33, pp. 10491–10502. Curran Associates, Inc., 2020.
- Patrick Emami and Sanjay Ranka. Learning permutations with sinkhorn policy gradient. *arXiv preprint arXiv:1805.07010*, 2018.
- Sven Erlander and Neil F Stewart. *The gravity model in transportation analysis: theory and extensions*, volume 3. Vsp, 1990.
- Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trounev, and Gabriel Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2681–2690. PMLR, 2019.

- Rémi Flamary, Marco Cuturi, Nicolas Courty, and Alain Rakotomamonjy. Wasserstein discriminant analysis. *Machine Learning*, 107(12):1923–1945, 2018.
- Charlie Frogner, Chiyuan Zhang, Hossein Mobahi, Mauricio Araya-Polo, and Tomaso Poggio. Learning with a wasserstein loss. In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, pp. 2053–2061, 2015.
- Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pp. 1608–1617. PMLR, 2018.
- Stephen Gould, Richard Hartley, and Dylan Campbell. Deep declarative networks: A new hope. *arXiv preprint arXiv:1909.04866*, 2019.
- L Kantorovich. On the transfer of masses: Doklady akademii nauk ussr. 1942.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Marcel Klatt, Carla Tameling, and Axel Munk. Empirical regularized optimal transport: Statistical theory and applications. *SIAM Journal on Mathematics of Data Science*, 2(2):419–443, 2020.
- Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Liu Liu, Dylan Campbell, Hongdong Li, Dingfu Zhou, Xibin Song, and Ruigang Yang. Learning 2d-3d correspondences to solve the blind perspective-n-point problem. *arXiv preprint arXiv:2003.06752*, 2020.
- Giulia Luise, Alessandro Rudi, Massimiliano Pontil, and Carlo Ciliberto. Differential properties of sinkhorn approximation for learning with wasserstein distance. *Advances in Neural Information Processing Systems*, 31:5859–5870, 2018.
- Giulia Luise, Saverio Salzo, Massimiliano Pontil, and Carlo Ciliberto. Sinkhorn barycenters with free support via frank-wolfe algorithm. *Advances in Neural Information Processing Systems*, 32: 9322–9333, 2019.
- G Mena, J Snoek, S Linderman, and D Belanger. Learning latent permutations with gumbel-sinkhorn networks. In *ICLR 2018 Conference Track*, volume 2018. OpenReview, 2018.
- Gautam Pai, Jing Ren, Simone Melzi, Peter Wonka, and Maks Ovsjanikov. Fast sinkhorn filters: Using matrix scaling for non-rigid shape correspondence with functional maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 384–393, 2021.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32: 8026–8037, 2019.
- Giorgio Patrini, Rianne van den Berg, Patrick Forre, Marcello Carioni, Samarth Bhargav, Max Welling, Tim Genewein, and Frank Nielsen. Sinkhorn autoencoders. In *Uncertainty in Artificial Intelligence*, pp. 733–743. PMLR, 2020.
- Ofir Pele and Michael Werman. Fast and robust earth mover’s distances. In *2009 IEEE 12th international conference on computer vision*, pp. 460–467. IEEE, 2009.
- Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.

- Marin Vlastelica Pogan, Anselm Paulus, Vit Musil, Georg Martius, and Michal Rolinek. Differentiation of blackbox combinatorial solvers. In *International Conference on Learning Representations*, 2019.
- Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 652–660, 2017.
- Yossi Rubner, Leonidas J Guibas, and Carlo Tomasi. The earth mover’s distance, multi-dimensional scaling, and color-based image retrieval. In *Proceedings of the ARPA image understanding workshop*, volume 661, pp. 668, 1997.
- Tim Salimans, Dimitris Metaxas, Han Zhang, and Alec Radford. Improving gans using optimal transport. In *6th International Conference on Learning Representations, ICLR 2018*, 2018.
- Rodrigo Santa Cruz, Basura Fernando, Anoop Cherian, and Stephen Gould. Deeppermnet: Visual permutation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3949–3957, 2017.
- Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4938–4947, 2020.
- Morgan A Schmitz, Matthieu Heitz, Nicolas Bonneel, Fred Ngole, David Coeurjolly, Marco Cuturi, Gabriel Peyré, and Jean-Luc Starck. Wasserstein dictionary learning: Optimal transport-based unsupervised nonlinear dictionary learning. *SIAM Journal on Imaging Sciences*, 11(1):643–678, 2018.
- Justin Solomon, Fernando De Goes, Gabriel Peyré, Marco Cuturi, Adrian Butscher, Andy Nguyen, Tao Du, and Leonidas Guibas. Convolutional wasserstein distances: Efficient optimal transportation on geometric domains. *ACM Transactions on Graphics (TOG)*, 34(4):1–11, 2015.
- Andrew M Stuart and Marie-Therese Wolfram. Inverse optimal transport. *SIAM Journal on Applied Mathematics*, 80(1):599–619, 2020.
- Cédric Villani. *Topics in optimal transportation*. Number 58. American Mathematical Soc., 2003.
- Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1912–1920, 2015.
- Lei Yang, Wenxi Liu, Zhiming Cui, Nenglun Chen, and Wenping Wang. Mapping in a cycle: Sinkhorn regularized unsupervised learning for point cloud shapes. In *European Conference on Computer Vision*, pp. 455–472. Springer, 2020.
- Zi Jian Yew and Gim Hee Lee. Rpm-net: Robust point matching using learned features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11824–11833, 2020.

## A PYTORCH IMPLEMENTATION

Our proposed algorithm can be easily integrated into existing deep learning architectures. In the following, we provide our PyTorch (Paszke et al., 2019) implementation of this module. The forward pass is the standard Sinkhorn matrix scaling algorithm (Cuturi, 2013) with a robust log-space implementation. The backward function contains an implementation of our Algorithm 1.

---

```
import torch

class Sinkhorn(torch.autograd.Function):
    @staticmethod
    def forward(ctx, c, a, b, num_sink, lambd_sink):
        log_p = -c / lambd_sink
        log_a = torch.log(a).unsqueeze(dim=1)
        log_b = torch.log(b).unsqueeze(dim=0)
        for _ in range(num_sink):
            log_p -= (torch.logsumexp(log_p, dim=0, keepdim=True) - log_b)
            log_p -= (torch.logsumexp(log_p, dim=1, keepdim=True) - log_a)
        p = torch.exp(log_p)

        ctx.save_for_backward(p, torch.sum(p, dim=1), torch.sum(p, dim=0))
        ctx.lambd_sink = lambd_sink
        return p

    @staticmethod
    def backward(ctx, grad_p):
        p, a, b = ctx.saved_tensors
        m, n = p.shape

        grad_p *= -1 / ctx.lambd_sink * p
        K = torch.cat((
            torch.cat((torch.diag(a), p), dim=1),
            torch.cat((p.T, torch.diag(b)), dim=1)),
            dim=0)[: -1, : -1]
        )
        t = torch.cat((
            grad_p.sum(dim=1),
            grad_p[:, : -1].sum(dim=0)),
            dim=0).unsqueeze(1)
        grad_ab, _ = torch.solve(t, K)
        grad_a = grad_ab[:m, :]
        grad_b = torch.cat((grad_ab[m:, :], torch.zeros([1, 1],
            device=device, dtype=torch.float32)), dim=0)
        U = grad_a + grad_b.T
        grad_p -= p * U
        grad_a = -ctx.lambd_sink * grad_a.squeeze(dim=1)
        grad_b = -ctx.lambd_sink * grad_b.squeeze(dim=1)

        return grad_p, grad_a, grad_b, None, None
```

---

## B ADDITIONAL EXPERIMENTS

### B.1 GRADIENT ACCURACY

Theorem 5 states that the error of our backward pass in Algorithm 1 can be bounded by the error of the forward pass in Equation 6. We now assess the magnitude of this error in practice by revisiting the experiments on image barycenter computation and number sorting from Sec. 5. Since the ground truth gradients are unknown in general, we define the (approximate) ground truth gradients as  $\nabla_{\mathcal{C}} \ell^* := \nabla_{\mathcal{C}} \ell^{(\tau_{\max})}$ ,  $\nabla_{\mathbf{a}} \ell^* := \nabla_{\mathbf{a}} \ell^{(\tau_{\max})}$  for  $\tau_{\max} := 10,000$ . In Fig. 5, we compare the error of our approach to AD, averaged over all settings and iterations from the experiments in Fig. 3 and Fig. 4 respectively. These results show clearly that, on average, our approach produces

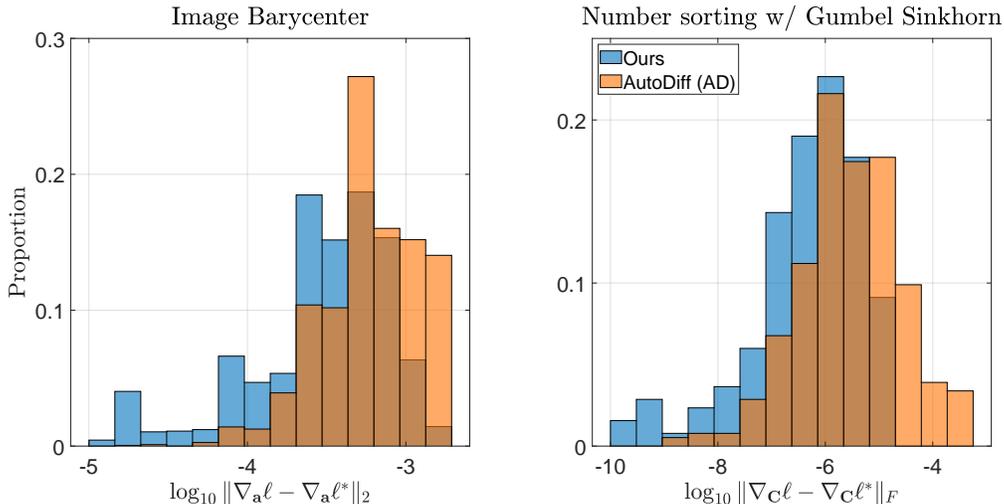


Figure 5: **Gradient accuracy.** We empirically assess the accuracy of the error bound discussed in Theorem 5. Specifically, we show the accuracy of the gradients  $\nabla_{\mathbf{a}}\ell$  for the image barycenter experiment in Sec. 5.2 and  $\nabla_{\mathbf{C}}\ell$  for the number sorting experiment in Sec. 5.3. While both distributions have a large overlap, the gradients from our approach are noticeably more accurate on average. Note that both comparisons show histograms on a log scale.

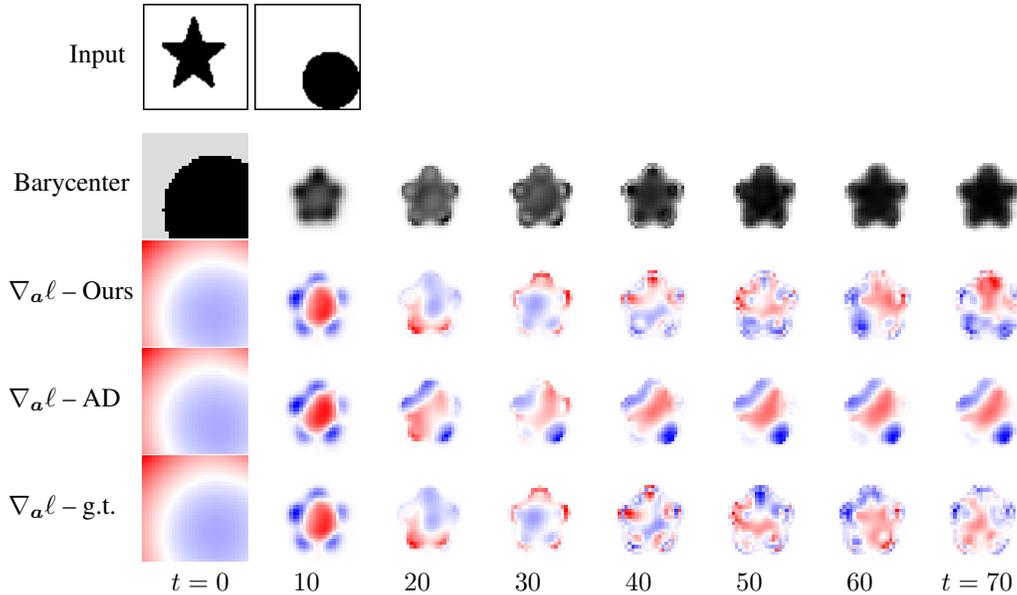


Figure 6: **Image barycenter gradients.** A qualitative comparison of our gradients (3rd row), the AD gradients (4th row), and the ground truth gradients (last row) for the image barycenter experiment from Sec. 5.2. Specifically, we consider the task of interpolating between two input images (1st row) with uniform interpolation weights  $w_1 = w_2 = 0.5$ . We show intermediate snapshots of the obtained barycenter image (2nd row) for different numbers of gradient descent iterations  $t \in \{0, \dots, 70\}$  that result from minimizing the energy in Equation 15.

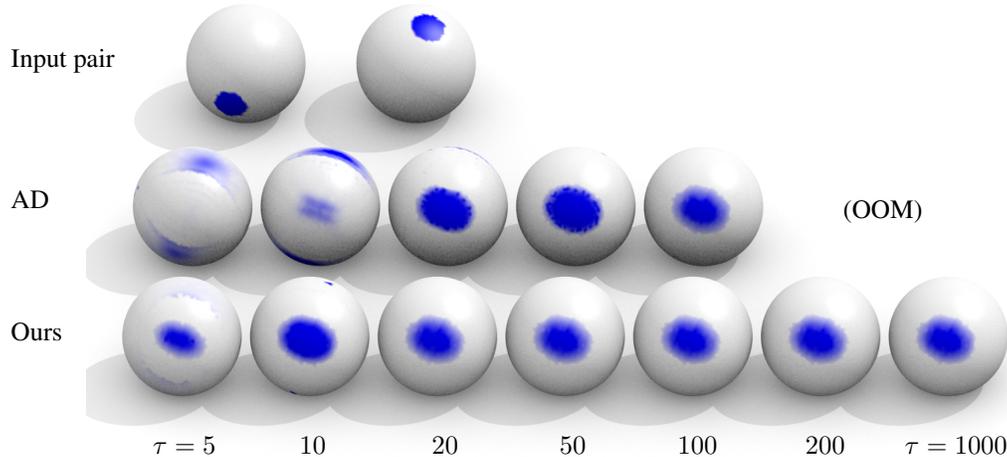


Figure 7: **Manifold barycenter.** We compute barycenters of two circular input distributions on the surface of a sphere (first row). Specifically, we compare the results of minimizing Equation 15 with AD (second row) and our gradients (third row). The sphere is discretized as a triangular mesh with 5000 vertices. On this resolution, AD is out of memory for  $\tau \geq 200$  Sinkhorn iterations whereas ours is still feasible for  $\tau = 1000$ . The obtained interpolations produce the slightly elongated shape of an ellipse since the surface of the sphere has a constant positive Gaussian curvature.

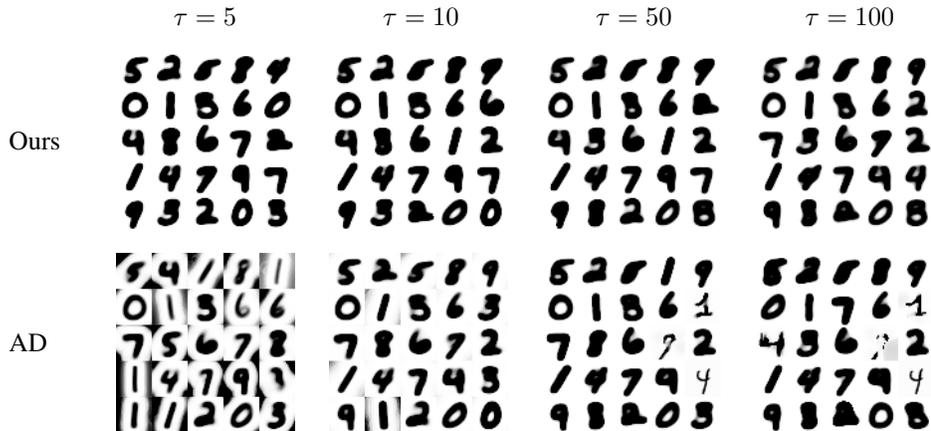


Figure 8: **MNIST k-means clustering.** A comparison of our approach and automatic differentiation on the task of k-means image clustering. For both approaches, we show the predicted  $k = 25$  clusters for  $\tau \in \{5, 10, 50, 100\}$  Sinkhorn iterations. We choose more than 10 clusters to capture several different appearances and styles for each digit. To make individual results more comparable, we use an identical random initialization of the cluster centroids for all settings. For AD, the maximum permissible batch sizes for the 4 settings are 512, 256, 64, 32, whereas ours consistently allows for a batch size of 1024.

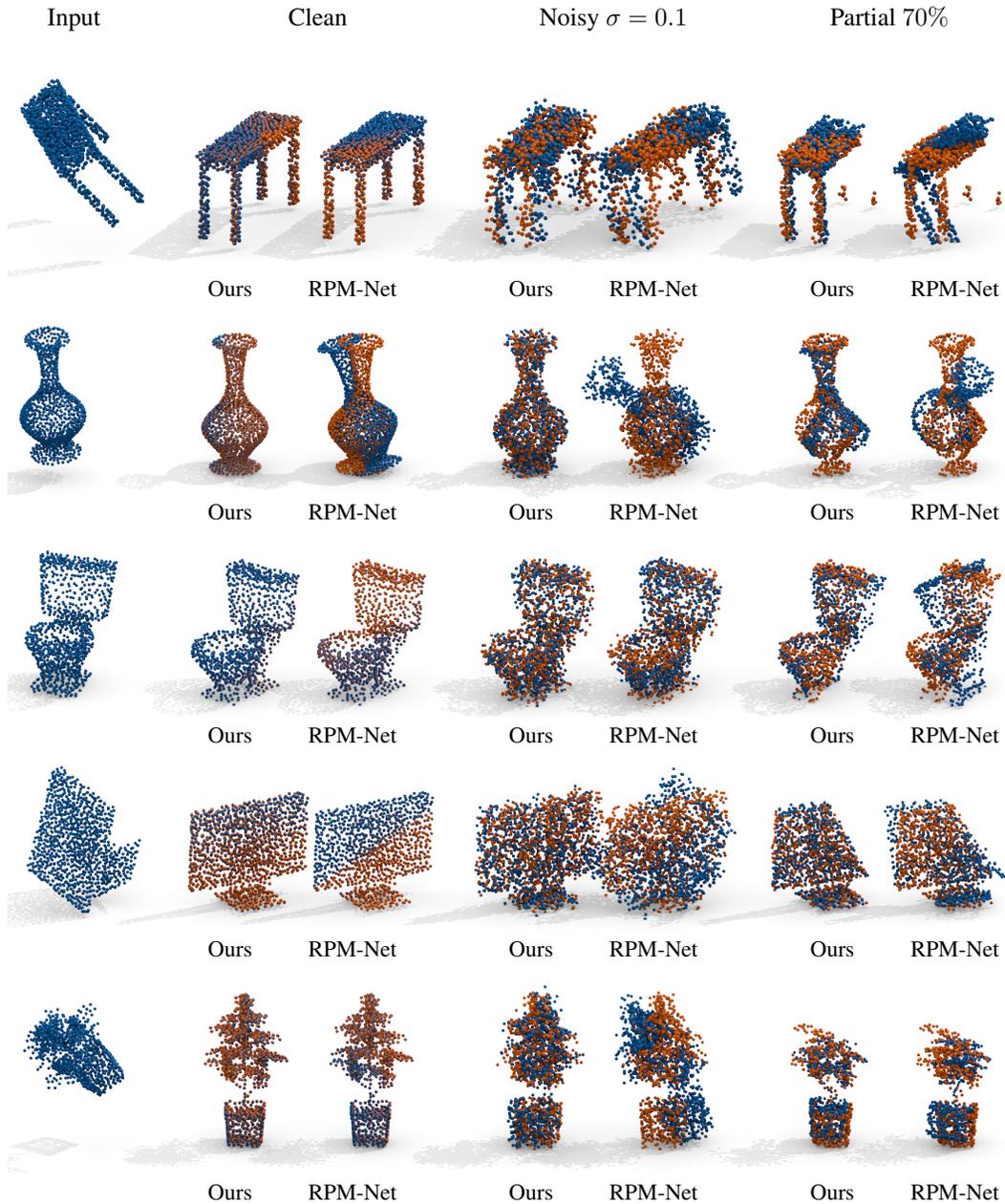


Figure 9: **Point cloud registration.** A qualitative comparison of RPM-Net (Yew & Lee, 2020) and the augmented version with our custom backward pass, see Alg. 1. The increased stability of the gradients predicted by our algorithm directly translates into more robust generalization results: Both methods are trained and tested on separate subsets of the 40 object categories in ModelNet40 (Wu et al., 2015), see Yew & Lee (2020, Section 6) for more details. Both methods yield accurate predictions for the clean test data, as indicated by the corresponding quantitative results in Table 1. On the other hand, our approach shows significant improvements when generalizing to noisy test data and partial inputs. In each row, we show a different test pair with the input pose  $X$  (1st column, blue), as well as the overlap of the reference pose  $Y$  (orange) and the predicted pose (blue) for the clean, noisy, and partial settings.

more accurate gradients than AD. Intuitively, ours yields optimal gradients for a suboptimal forward pass  $\mathbf{P}^* \approx \mathcal{S}_\lambda^{(\tau)}$ , see Theorem 5, whereas AD yields approximate gradients for an approximate forward pass. To further illustrate this point, we show a qualitative comparison of the gradients of both approaches on the task of image barycenter computation in Fig. 6. These results demonstrate the usefulness of our custom backward pass. Our gradients and the ground truth start disagreeing in certain regions for  $t \geq 50$ , but only after the barycenter optimization converges. Compared to the vanilla AD approach, the gradients are much closer to the ground truth and therefore more useful.

## B.2 MANIFOLD BARYCENTERS AND IMAGE CLUSTERING

In Section 5.2, we show that our approach can be used to compute Wasserstein barycenters of images. We can leverage the same approach to interpolate between distributions on more general manifold domains. Specifically, we can minimize Equation 15 for a set of input marginals  $\{\mathbf{b}_1, \dots, \mathbf{b}_k\}$  and use squared pairwise geodesic distances  $\mathbf{D} \in \mathbb{R}^{m \times n}$  as the input cost matrix. In Fig. 7, we show the predicted barycenters for our method and AD with two circular input distributions on a sphere. Moreover, we can use a similar approach and leverage Equation 15 to predict a clustering of a set of input images. To that end, we apply the k-means algorithm which alternates between computing cluster centroids (by minimizing Equation 15) and assigning all images  $\mathbf{b}_i$  to their respective corresponding centroid (defined as the centroid with minimum distance  $d$ , right side of Equation 15). We show results on the 60k images from the MNIST dataset (LeCun, 1998) in Fig. 8. While our algorithm is robust to varying Sinkhorn iterations  $\tau$ , AD’s results vary significantly. For small  $\tau \leq 10$ , the Sinkhorn approximation is not sufficiently exact. For higher  $\tau$ , memory constraints force AD to use a smaller batch size which again leads to instabilities.

## B.3 DETAILS ON PERMUTATION BASELINES

**Gumbel-Sinkhorn networks** As outlined in Sec. 5.3, the goal of the Gumbel-Sinkhorn method (Mena et al., 2018) is to learn how to sort a set of input elements  $\{x_1, \dots, x_n\}$ . To this end, the cost matrix  $\mathbf{C}$  is parameterized via a permutation-invariant neural network architecture (set encoder), conditioned on the input elements  $\{x_1, \dots, x_n\}$ . The matrix  $\mathbf{C}$ , together with marginals  $\mathbf{a} = \mathbf{b} = \mathbb{1}_n$  are then passed through a differentiable Sinkhorn layer.<sup>2</sup> The final output  $\mathbf{P}$  is a bistochastic matrix which encodes an approximate permutation. The training objective is a geometric loss, minimizing the distance of the sorted elements  $x_i$  to their natural ground-truth ordering. At test time, the Hungarian algorithm (Kuhn, 1955) is applied to the learned cost matrix  $\mathbf{C}$  to obtain an exact, hard permutation. In Sec. 5.3, we show the concrete application of sorting  $n$  real numbers, sampled randomly from the uniform distribution  $x_i \sim \mathcal{U}(0, 1)$ . More specifically, we consider separate training and test sets of 4096 and 128 random sequences each and report the error, defined as the proportion of incorrectly placed numbers in the test set, see Fig. 4. To provide a more complete picture, we report quantitative results on the task of generalizing to different test sets here. Specifically, we train the vanilla GS method and our approach for  $\tau = 200$  (the maximum number for which GS is below the GPU memory limit) for  $n = 200$  numbers sampled from  $\mathcal{U}(0, 1)$ . We then investigate the error on test sets sampled from different distributions  $x_i \sim \mathcal{U}(s, t)$  with  $s < t$  in Table 2. Even though the variance is quite high, these results indicate that our method evidently helps to reduce overfitting. Note, that the improved generalization is observed for the setting  $\tau = 200$ ,  $n = 200$  where the performance of both methods on the standard test set is almost identical, see Fig. 4.

**RPM-Net** A number of recent works use the Sinkhorn operator as a differentiable matching module for deep learning (Sarlin et al., 2020; Yew & Lee, 2020; Yang et al., 2020; Liu et al., 2020; Eisenberger et al., 2020). The standard strategy of such methods is to use a learnable feature extractor to obtain descriptors  $\mathbf{F}^X \in \mathbb{R}^{m \times p}$ ,  $\mathbf{F}^Y \in \mathbb{R}^{n \times p}$  on pairs of input objects  $X$  and  $Y$  in a siamese manner. We can then define the cost matrix  $\mathbf{C} := \mathbf{D}$  as the squared pairwise distance matrix  $D_{i,j} = \|\mathbf{F}_{:,i}^X - \mathbf{F}_{:,j}^Y\|_2^2$ . For fixed, uniform marginals  $\mathbf{a}$  and  $\mathbf{b}$ , the Sinkhorn operator then yields a soft matching  $\mathbf{P} = \mathcal{S}_\lambda(\mathbf{C}, \mathbf{a}, \mathbf{b}) \in \mathbb{R}^{m \times n}$  between the two input objects. The baseline method RPM-Net we consider in Sec. 5.3 uses this methodology to obtain a matching between pairs of input

<sup>2</sup>Strictly speaking, the choice of marginals  $\mathbf{a} = \mathbf{b} = \mathbb{1}_n$  does not fit in our framework, since we require  $\mathbf{a}, \mathbf{b} \in \Delta_n$ , see Equation 1. However, we can simply use  $\mathbf{a} = \mathbf{b} = \frac{1}{n} \mathbb{1}_n$  and scale the resulting transportation plan  $\mathbf{P}$  by a constant factor of  $n$ .

	$\mathcal{U}(1, 2)$	$\mathcal{U}(10, 11)$	$\mathcal{U}(100, 101)$	$\mathcal{U}(-1, 0)$
Ours	$0.2964 \pm 0.0744$	$0.3340 \pm 0.3059$	$0.3531 \pm 0.1380$	$0.3552 \pm 0.2116$
Gumbel-Sinkhorn	$0.3163 \pm 0.0976$	$0.3620 \pm 0.3179$	$0.3955 \pm 0.1478$	$0.4678 \pm 0.2526$

Table 2: **Number sorting generalization.** We assess the capability of our approach (first row) and the vanilla Gumbel-Sinkhorn method (second row) (Mena et al., 2018) to generalize to various test sets  $\mathcal{U}(s, t)$  with  $s < t$ . We train both methods to sort sets of  $n = 200$  numbers  $x_i$  from the uniform distribution on the unit interval  $\mathcal{U}(0, 1)$  with  $\tau = 200$  Sinkhorn iterations. For each test set, we show the mean proportion of false predictions, as well as the empirical standard deviation, obtained from 50 test runs per setting.

point clouds. As a feature extractor, it uses PointNet (Qi et al., 2017) with a custom input task point cloud, see Yew & Lee (2020, Sec. 5.1) for more details. Using the obtained soft correspondences, a simple SVD transformation then yields the optimal rigid transformation between the two input point clouds. In order to train their model, RPM-Net uses automatic differentiation of the Sinkhorn layer. We can now demonstrate that replacing AD with our backward algorithm improves the performance. To that end, we revisit the experiments from Yew & Lee (2020, Sec. 6): We use 20 separate object identities of the ModelNet40 dataset (Wu et al., 2015) as our train and test set respectively. In Fig 9, we now show a qualitative comparison corresponding to the results in Table 1 in the main paper. On the standard test set, both approaches produce comparable, high-quality results. On the other hand, our method significantly improves the robustness when generalizing to noisy data or partial views.

## C PROOFS

In the following, we provide proofs of Lemma 2, Theorem 4 and Theorem 5 from the main paper.

### C.1 PROOF OF LEMMA 2

*Proof.* As mentioned above, the key for proving this statement is applying the implicit function theorem. We start by (a) showing that this indeed yields the identity from Equation 9, and then (b) justify why removing the last  $(m + n)$ -th equality condition from  $\mathbf{E}$  is necessary.

- (a) First of all, we can verify by direct computation that the matrix  $\mathbf{K}$  from Equation 9 corresponds to the partial derivatives of the KKT conditions from Equation 8

$$\mathbf{K} = \left( \frac{\partial \mathcal{K}(\mathbf{c}, \mathbf{a}, \mathbf{b}, \mathbf{p}, \boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial [\mathbf{p}; \boldsymbol{\alpha}; \boldsymbol{\beta}]} \right)_{-l, -l}, \quad (16)$$

where the notation  $(\cdot)_{-l, -l}$  means that the last row and column is removed and where  $l = mn + m + n$ . Furthermore,  $\mathbf{K}$  is invertible (since the solution lies in the interior  $\mathbf{P}_{i,j} > 0$ , see Peyré et al. (2019, p. 68)), and the  $m + n - 1$  columns of  $\tilde{\mathbf{E}}$  are linearly dependent, see (b). Consequently, the implicit function theorem states that  $\mathcal{K}$  implicitly defines a mapping  $(\mathbf{c}, \mathbf{a}, \mathbf{b}) \mapsto (\mathbf{p}, \boldsymbol{\alpha}, \boldsymbol{\beta})$  whose Jacobian is

$$\mathbf{J} = \frac{\partial [\mathbf{p}; \boldsymbol{\alpha}; \tilde{\boldsymbol{\beta}}]}{\partial [\mathbf{c}; -\mathbf{a}; -\tilde{\mathbf{b}}]} = - \left( \frac{\partial \mathcal{K}}{\partial [\mathbf{p}; \boldsymbol{\alpha}; \boldsymbol{\beta}]} \right)_{-l, -l}^{-1} \underbrace{\left( \frac{\partial \mathcal{K}}{\partial [\mathbf{c}; -\mathbf{a}; -\tilde{\mathbf{b}}]} \right)_{-l, -l}}_{=\mathbf{I}_{l-1}} = -\mathbf{K}^{-1}. \quad (17)$$

- (b) As part of the proof in (a), we use the fact that the columns of  $\tilde{\mathbf{E}}$  are linearly independent. Verifying this statement also provides insight as to why removing the last row of  $\tilde{\mathbf{b}}, \tilde{\boldsymbol{\beta}}$  and the last column of  $\tilde{\mathbf{E}}$  is necessary. Intuitively, the columns of  $\mathbf{E}$  contain one redundant condition: The identity

$$\sum_{i=1}^m \mathbf{P}_{i,n} = 0 \iff \mathbf{E}_{:,m+n}^\top \mathbf{p} = 0, \quad (18)$$

follows directly from the other  $m+n-1$  conditions. More formally, we can take the kernel

$$\ker(\mathbf{E}^\top) = \{\mathbf{P} \in \mathbb{R}^{m \times n} \mid \mathbf{E}^\top \mathbf{p} = 0\}, \quad (19)$$

of  $\mathbf{E}^\top \in \mathbb{R}^{m+n \times mn}$  and observe that  $\dim(\ker(\mathbf{E}^\top)) = (m-1)(n-1)$ , see Bolker (1972, Sec. 1). The rank-nullity theorem then implies that the dimension of the subspace spanned by the columns of  $\mathbf{E}$  is of dimension  $mn - (m-1)(n-1) = m+n-1$ . Consequently, removing the redundant condition in Equation 19 from  $\mathbf{E}$  yields the reduced  $\tilde{\mathbf{E}} \in \mathbb{R}^{mn \times m+n-1}$  with  $m+n-1$  linearly independent columns. □

## C.2 PROOF OF THEOREM 4

*Proof.* We want to show that the gradients  $\nabla_{\mathbf{C}}\ell, \nabla_{\mathbf{a}}\ell, \nabla_{\mathbf{b}}\ell$  obtained with Algorithm 1 are equivalent to the solution of Equation 11. To that end, we start by applying the Schur complement trick to the block matrix  $\mathbf{K}$ . This yields the expression

$$(\mathbf{K}^{-1})_{:,1:mn} = \begin{bmatrix} \lambda^{-1} \text{diag}(\mathbf{p})(\mathbf{I}_{mn} + \tilde{\mathbf{E}}(\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p})\tilde{\mathbf{E}})^{-1}\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p})) \\ (\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p})\tilde{\mathbf{E}})^{-1}\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p}) \end{bmatrix}, \quad (20)$$

for the first  $mn$  columns of its inverse. In the next step, we can insert this expression in Equation 11 and invert the linear system of equations

$$\begin{bmatrix} \nabla_{\mathbf{C}}\ell \\ -\nabla_{[\mathbf{a};\mathbf{b}]} \ell \end{bmatrix} = \mathbf{K}^{-1} \begin{bmatrix} -\nabla_{\mathbf{p}}\ell \\ \mathbf{0} \end{bmatrix} = -(\mathbf{K}^{-1})_{:,1:mn} \nabla_{\mathbf{p}}\ell. \quad (21)$$

Further simplification yields the following identities for the gradients of  $\mathbf{C}$ ,  $\mathbf{a}$  and  $\mathbf{b}$

$$\begin{aligned} \begin{bmatrix} \nabla_{\mathbf{C}}\ell \\ \nabla_{[\mathbf{a};\mathbf{b}]} \ell \end{bmatrix} &= \begin{bmatrix} -\lambda^{-1} \text{diag}(\mathbf{p})(\mathbf{I}_{mn} - \tilde{\mathbf{E}}(\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p})\tilde{\mathbf{E}})^{-1}\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p}))\nabla_{\mathbf{p}}\ell \\ (\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p})\tilde{\mathbf{E}})^{-1}\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p})\nabla_{\mathbf{p}}\ell \end{bmatrix} \\ &= \begin{bmatrix} -\lambda^{-1}(\text{diag}(\mathbf{p})\nabla_{\mathbf{p}}\ell - \text{diag}(\mathbf{p})\tilde{\mathbf{E}}\nabla_{[\mathbf{a};\mathbf{b}]} \ell) \\ (\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p})\tilde{\mathbf{E}})^{-1}\tilde{\mathbf{E}}^\top \text{diag}(\mathbf{p})\nabla_{\mathbf{p}}\ell \end{bmatrix}, \end{aligned} \quad (22)$$

where the latter equality results from substituting the obtained expression for  $\nabla_{[\mathbf{a};\mathbf{b}]} \ell$  in the first block row. In the remainder of this proof, we can show line by line that these expressions yield Algorithm 1. The main idea is to first compute the second block row identity in Equation 22 and then use the result to eventually obtain  $\nabla_{\mathbf{C}}\ell$  from the first block row:

1. 1 The first line defines the matrix  $\mathbf{T} := \mathbf{P} \odot \nabla_{\mathbf{p}}\ell$  via the Hadamard product  $\odot$ . In vectorized form it corresponds to the expression

$$\mathbf{t} = \text{diag}(\mathbf{p})\nabla_{\mathbf{p}}\ell. \quad (23)$$

1. 2 As detailed in Lemma 2, we remove the last equality condition from  $\mathbf{E}$  to obtain  $\tilde{\mathbf{E}}$ . Equivalent considerations require us to introduce the truncated versions  $\tilde{\mathbf{T}}, \tilde{\mathbf{P}}$  of  $\mathbf{T}, \mathbf{P}$ .

1. 3 The operator  $\mathbf{E}^\top$  then maps  $\mathbf{t}$  to the vector

$$\mathbf{E}^\top \mathbf{t} = [\mathbb{1}_n \otimes \mathbf{I}_m \quad \mathbf{I}_n \otimes \mathbb{1}_m]^\top \mathbf{t} = \begin{bmatrix} (\mathbb{1}_n^\top \otimes \mathbf{I}_m)\mathbf{t} \\ (\mathbf{I}_n \otimes \mathbb{1}_m^\top)\mathbf{t} \end{bmatrix} = \begin{bmatrix} \mathbf{T}\mathbb{1}_n \\ \tilde{\mathbf{T}}^\top \mathbb{1}_m \end{bmatrix}, \quad (24)$$

that contains its row and column sums. In terms of the truncated  $\tilde{\mathbf{E}}$ , the last row of Equation 24 gets removed, thus

$$\tilde{\mathbf{E}}^\top \mathbf{t} = \begin{bmatrix} \mathbf{T}\mathbb{1}_n \\ \tilde{\mathbf{T}}^\top \mathbb{1}_m \end{bmatrix} = \begin{bmatrix} \mathbf{t}^{(a)} \\ \tilde{\mathbf{t}}^{(b)} \end{bmatrix}. \quad (25)$$

1.4 A direction computation reveals that

$$\begin{aligned} \mathbf{E}^\top \text{diag}(\mathbf{p}) \mathbf{E} &= \begin{bmatrix} (\mathbb{1}_n^\top \otimes \mathbf{I}_m) \text{diag}(\mathbf{p})(\mathbb{1}_n \otimes \mathbf{I}_m) & (\mathbb{1}_n^\top \otimes \mathbf{I}_m) \text{diag}(\mathbf{p})(\mathbf{I}_n \otimes \mathbb{1}_m) \\ (\mathbf{I}_n \otimes \mathbb{1}_m^\top) \text{diag}(\mathbf{p})(\mathbb{1}_n \otimes \mathbf{I}_m) & (\mathbf{I}_n \otimes \mathbb{1}_m^\top) \text{diag}(\mathbf{p})(\mathbf{I}_n \otimes \mathbb{1}_m) \end{bmatrix} \\ &= \begin{bmatrix} \text{diag}(\mathbf{P} \mathbb{1}_n) & \mathbf{P} \\ \mathbf{P}^\top & \text{diag}(\mathbf{P}^\top \mathbb{1}_m) \end{bmatrix} = \begin{bmatrix} \text{diag}(\mathbf{a}) & \mathbf{P} \\ \mathbf{P}^\top & \text{diag}(\mathbf{b}) \end{bmatrix}. \end{aligned} \quad (26)$$

The linear system in 1.4 of Algorithm 1 therefore yields the gradients  $\nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell$  by inserting Equation 23, Equation 25 and (the reduced version of) Equation 26 into the second block row identity from Equation 22.

- 1.5 We can expand  $\nabla_{\tilde{\mathbf{b}}} \ell$  to  $\nabla_{\mathbf{b}} \ell$  by appending zero  $\nabla_{b_n} \ell = 0$  as the last entry. Since  $\mathbf{b}$  is constrained to the probability simplex  $\Delta_n$  this gradient is exact for all entries, see the discussion in Section 3.4.
- 1.6 Having computed the gradients  $\nabla_{[\mathbf{a}; \mathbf{b}]} \ell$ , we can now insert them in the first row of Equation 22. Here, the reduced and the original expressions are equivalent

$$\tilde{\mathbf{E}} \nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell = \mathbf{E} \nabla_{[\mathbf{a}; \mathbf{b}]} \ell \quad (27)$$

because 1.5 specifies  $\nabla_{b_n} \ell = 0$ . Thus,

$$\tilde{\mathbf{E}} \nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell = [\mathbb{1}_n \otimes \mathbf{I}_m \quad \mathbf{I}_n \otimes \mathbb{1}_m] \nabla_{[\mathbf{a}; \mathbf{b}]} \ell = \mathbb{1}_n \otimes \nabla_{\mathbf{a}} \ell + \nabla_{\mathbf{b}} \ell \otimes \mathbb{1}_m =: \mathbf{u}, \quad (28)$$

defines the vectorized version of  $\mathbf{U}$  from 1.6.

- 1.7 Putting everything together, we insert the identities from Equation 23 and Equation 27 into the first block row of Equation 22

$$\nabla_{\mathbf{c}} \ell = -\lambda^{-1} (\text{diag}(\mathbf{p}) \nabla_{\mathbf{p}} \ell - \text{diag}(\mathbf{p}) \tilde{\mathbf{E}} \nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell) = -\lambda^{-1} (\mathbf{t} - \text{diag}(\mathbf{p}) \mathbf{u}), \quad (29)$$

which is equivalent to the matrix-valued expression in 1.7.

□

### C.3 PROOF OF THEOREM 5

*Proof.* The key for constructing the error bounds in Equation 14a and Equation 14b is finding a bound for the first-order derivatives  $\frac{\partial \nabla_{\mathbf{c}} \ell}{\partial \mathbf{p}}$  and  $\frac{\partial \nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell}{\partial \mathbf{p}}$ . For brevity, we introduce the short-hand notation  $\tilde{\mathbf{P}} := \text{diag}(\mathbf{p})$ . Furthermore, we define the projection of  $\mathbf{x}$  onto the column space of  $\tilde{\mathbf{E}}$  as

$$\Pi_{\tilde{\mathbf{E}}} \mathbf{x} := \arg \min_{\mathbf{y} \in \text{span}(\tilde{\mathbf{E}})} \|\mathbf{x} - \mathbf{y}\|_{\tilde{\mathbf{P}}}^2 = \tilde{\mathbf{E}} \arg \min_{\mathbf{z} \in \mathbb{R}^{m+n-1}} \|\mathbf{x} - \tilde{\mathbf{E}} \mathbf{z}\|_{\tilde{\mathbf{P}}}^2, \quad (30)$$

where  $\langle \cdot, \cdot \rangle_{\tilde{\mathbf{P}}} := \langle \tilde{\mathbf{P}}^{\frac{1}{2}} \cdot, \tilde{\mathbf{P}}^{\frac{1}{2}} \cdot \rangle_2$ . In matrix notation, Equation 30 reads

$$\Pi_{\tilde{\mathbf{E}}} = \tilde{\mathbf{E}} (\tilde{\mathbf{E}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{E}})^{-1} \tilde{\mathbf{E}}^\top \tilde{\mathbf{P}}. \quad (31)$$

Using Equation 13a and Equation 13b from the proof of Theorem 4, we can then rewrite the backward pass compactly as

$$\nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell = \tilde{\mathbf{E}}^\dagger \Pi_{\tilde{\mathbf{E}}} \nabla_{\mathbf{p}} \ell, \quad \text{and} \quad (32a)$$

$$\nabla_{\mathbf{c}} \ell = -\lambda^{-1} (\tilde{\mathbf{P}} (\mathbf{I} - \Pi_{\tilde{\mathbf{E}}}) \nabla_{\mathbf{p}} \ell), \quad (32b)$$

The first identity follows from  $\tilde{\mathbf{E}}^\dagger \tilde{\mathbf{E}} = \mathbf{I}$ , since the columns of  $\tilde{\mathbf{E}}$  are linearly independent (see part (b) of the proof of Lemma 2 in Appendix C.1). Direct substitution of Equation 13a into Equation 13b immediately yields Equation 32b. To differentiate  $\nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell$  and  $\nabla_{\mathbf{c}} \ell$ , we apply the chain rule which in turn requires a closed-form expression for the derivative of the projection operator  $\Pi_{\tilde{\mathbf{E}}}$ . Since it is defined as the solution of an optimization problem, we apply the implicit function theorem to the gradient of the objective in Equation 30, i.e.

$$\nabla_{\mathbf{z}} \left( \frac{1}{2} \|\mathbf{x} - \tilde{\mathbf{E}} \mathbf{z}\|_{\tilde{\mathbf{P}}}^2 \right) = \tilde{\mathbf{E}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{E}} \mathbf{z} - \tilde{\mathbf{E}}^\top \tilde{\mathbf{P}} \mathbf{x} = \mathbf{0}. \quad (33)$$

The Jacobian of the mapping  $\mathbf{p} \mapsto \mathbf{\Pi}_E \mathbf{x}$  can therefore be written in terms of the IFT as

$$\frac{\partial \mathbf{\Pi}_E \mathbf{x}}{\partial \mathbf{p}} = \tilde{\mathbf{E}} (\tilde{\mathbf{E}}^\top \bar{\mathbf{P}} \tilde{\mathbf{E}})^{-1} \tilde{\mathbf{E}}^\top \text{diag}(\mathbf{x} - \mathbf{\Pi}_E \mathbf{x}), \quad (34)$$

This auxiliary result implies that the Jacobians of the mappings  $\mathbf{p} \mapsto \nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell$  and  $\mathbf{p} \mapsto \nabla_{\mathbf{c}} \ell$  defined in Equation 32a and Equation 32b are

$$\begin{aligned} \frac{\partial \nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell}{\partial \mathbf{p}} &= \tilde{\mathbf{E}}^\dagger \frac{\partial}{\partial \mathbf{p}} \mathbf{\Pi}_E \nabla_{\mathbf{p}} \ell = (\tilde{\mathbf{E}}^\top \bar{\mathbf{P}} \tilde{\mathbf{E}})^{-1} \tilde{\mathbf{E}}^\top \text{diag}((\mathbf{I} - \mathbf{\Pi}_E) \nabla_{\mathbf{p}} \ell) + \tilde{\mathbf{E}}^\dagger \mathbf{\Pi}_E \nabla_{\mathbf{p}}^2 \ell, \quad \text{and} \quad (35a) \\ \frac{\partial \nabla_{\mathbf{c}} \ell}{\partial \mathbf{p}} &= -\lambda^{-1} \left( \text{diag}((\mathbf{I} - \mathbf{\Pi}_E) \nabla_{\mathbf{p}} \ell) - \mathbf{\Pi}_E^\top \text{diag}((\mathbf{I} - \mathbf{\Pi}_E) \nabla_{\mathbf{p}} \ell) + \bar{\mathbf{P}} (\mathbf{I} - \mathbf{\Pi}_E) \nabla_{\mathbf{p}}^2 \ell \right) \\ &= -\lambda^{-1} \left( (\mathbf{I} - \mathbf{\Pi}_E^\top) \text{diag}((\mathbf{I} - \mathbf{\Pi}_E) \nabla_{\mathbf{p}} \ell) + \bar{\mathbf{P}} (\mathbf{I} - \mathbf{\Pi}_E) \nabla_{\mathbf{p}}^2 \ell \right). \end{aligned} \quad (35b)$$

In order to bound the errors of these two gradients, we first derive an upper bound for the norm of the operator  $\mathbf{\Pi}_E$ . An important insight is that we can precondition  $\mathbf{\Pi}_E$  via  $\bar{\mathbf{P}}^{\frac{1}{2}}$

$$\bar{\mathbf{P}}^{\frac{1}{2}} \mathbf{\Pi}_E \bar{\mathbf{P}}^{-\frac{1}{2}} \mathbf{x} = \arg \min_{\mathbf{y} \in \text{span}(\bar{\mathbf{P}}^{\frac{1}{2}} \tilde{\mathbf{E}})} \|\mathbf{x} - \mathbf{y}\|_2, \quad (36)$$

which results in an orthogonal projection  $\bar{\mathbf{P}}^{\frac{1}{2}} \mathbf{\Pi}_E \bar{\mathbf{P}}^{-\frac{1}{2}}$ . Since such projections have a spectral radius of at most 1, we can bound the norm of  $\mathbf{\Pi}_E$  as

$$\|\mathbf{\Pi}_E\|_2 \leq \|\bar{\mathbf{P}}^{-\frac{1}{2}}\|_2 \|\bar{\mathbf{P}}^{\frac{1}{2}} \mathbf{\Pi}_E \bar{\mathbf{P}}^{-\frac{1}{2}}\|_2 \|\bar{\mathbf{P}}^{\frac{1}{2}}\|_2 \leq \|\bar{\mathbf{P}}^{-\frac{1}{2}}\|_2 \|\bar{\mathbf{P}}^{\frac{1}{2}}\|_2, \quad (37)$$

and equivalently show for the complementary projector  $\mathbf{I} - \mathbf{\Pi}_E$  that

$$\|\mathbf{I} - \mathbf{\Pi}_E\|_2 \leq \|\bar{\mathbf{P}}^{-\frac{1}{2}}\|_2 \|\bar{\mathbf{P}}^{\frac{1}{2}} (\mathbf{I} - \mathbf{\Pi}_E) \bar{\mathbf{P}}^{-\frac{1}{2}}\|_2 \|\bar{\mathbf{P}}^{\frac{1}{2}}\|_2 \leq \|\bar{\mathbf{P}}^{-\frac{1}{2}}\|_2 \|\bar{\mathbf{P}}^{\frac{1}{2}}\|_2. \quad (38)$$

The Jacobians of the backward pass can then be bounded as

$$\begin{aligned} \left\| \frac{\partial \nabla_{[\mathbf{a}; \tilde{\mathbf{b}}]} \ell}{\partial \mathbf{p}} \right\|_F &\leq \|(\tilde{\mathbf{E}}^\top \bar{\mathbf{P}} \tilde{\mathbf{E}})^{-1} \tilde{\mathbf{E}}^\top\|_2 \|(\mathbf{I} - \mathbf{\Pi}_E) \nabla_{\mathbf{p}} \ell\|_2 + \|\tilde{\mathbf{E}}^\dagger \mathbf{\Pi}_E \nabla_{\mathbf{p}}^2 \ell\|_F \\ &= \|\tilde{\mathbf{E}}^\dagger \mathbf{\Pi}_E \bar{\mathbf{P}}^{-1}\|_2 \|(\mathbf{I} - \mathbf{\Pi}_E) \nabla_{\mathbf{p}} \ell\|_2 + \|\tilde{\mathbf{E}}^\dagger \mathbf{\Pi}_E \nabla_{\mathbf{p}}^2 \ell\|_F \\ &\leq \|\tilde{\mathbf{E}}^\dagger\|_2 \|\bar{\mathbf{P}}^{-\frac{1}{2}}\|_2^2 \|\mathbf{I} - \mathbf{\Pi}_E\|_2 \|\nabla_{\mathbf{p}} \ell\|_2 + \|\tilde{\mathbf{E}}^\dagger\|_2 \|\mathbf{\Pi}_E\|_2 \|\nabla_{\mathbf{p}}^2 \ell\|_F \\ &\leq \|\tilde{\mathbf{E}}^\dagger\|_2 \|\bar{\mathbf{P}}^{-\frac{1}{2}}\|_2 \|\bar{\mathbf{P}}^{\frac{1}{2}}\|_2 \left( \|\bar{\mathbf{P}}^{-\frac{1}{2}}\|_2^2 \|\nabla_{\mathbf{p}} \ell\|_2 + \|\nabla_{\mathbf{p}}^2 \ell\|_F \right) \end{aligned} \quad (39a)$$

$$\begin{aligned} \left\| \frac{\partial \nabla_{\mathbf{c}} \ell}{\partial \mathbf{p}} \right\|_F &\leq \lambda^{-1} \|\mathbf{I} - \mathbf{\Pi}_E^\top\|_2 \|\mathbf{I} - \mathbf{\Pi}_E\|_2 \|\nabla_{\mathbf{p}} \ell\|_2 + \lambda^{-1} \|\bar{\mathbf{P}} (\mathbf{I} - \mathbf{\Pi}_E)\|_2 \|\nabla_{\mathbf{p}}^2 \ell\|_F \\ &\leq \lambda^{-1} \|\bar{\mathbf{P}}^{-\frac{1}{2}}\|_2^2 \|\bar{\mathbf{P}}^{\frac{1}{2}}\|_2^2 \|\nabla_{\mathbf{p}} \ell\|_2 + \lambda^{-1} \|\bar{\mathbf{P}}^{\frac{1}{2}}\|_2^2 \|\nabla_{\mathbf{p}}^2 \ell\|_F \\ &\leq \lambda^{-1} \sigma_+ \left( \frac{1}{\sigma_-} C_1 + C_2 \right), \end{aligned} \quad (39b)$$

where the constants  $\sigma_-, \sigma_+, C_1, C_2 > 0$  are as defined in Theorem 5, and where we use the identity

$$\|\mathbf{A} \text{diag}(\mathbf{b})\|_F \leq \|\mathbf{A}\|_2 \|\text{diag}(\mathbf{b})\|_F = \|\mathbf{A}\|_2 \|\mathbf{b}\|_2. \quad (40)$$

As a direct consequence, we obtain the bounds from Equation 14a and Equation 14b, since the bounded derivatives imply the Lipschitz continuity of the differentiable map  $\mathbf{p} \mapsto \nabla_{[\mathbf{c}; \mathbf{a}; \tilde{\mathbf{b}}]} \ell$ .

□