

The Self-Optimal-Transport Feature Transform

Daniel Shalam and Simon Korman

University of Haifa, Israel

Abstract. The Self-Optimal-Transport (SOT) feature transform is designed to upgrade the set of features of a data instance to facilitate downstream matching or grouping related tasks. The transformed set encodes a rich representation of high order relations between the instance features. Distances between transformed features capture their *direct* original similarity and their *third party* ‘agreement’ regarding similarity to other features in the set. A particular min-cost-max-flow fractional matching problem, whose entropy regularized version can be approximated by an optimal transport (OT) optimization, results in our transductive transform which is efficient, differentiable, equivariant, parameterless and probabilistically interpretable. Empirically, the transform is highly effective and flexible in its use, consistently improving networks it is inserted into, in a variety of tasks and training schemes. We demonstrate its merits through the problem of unsupervised clustering and its efficiency and wide applicability for few-shot-classification, with state-of-the-art results, and large-scale person re-identification.

1 Introduction

In this work, we reassess the design and functionality of features for *instance-specific* problems. In such problems, typically, features computed at test time are mainly compared relative to one another, and less so to the features seen at training time. For such problems the standard practice of learning a generic feature extractor during training and applying it at test time might be sub-optimal.

We aim at finding training and inference schemes that take into account these considerations, being able to exploit large corpuses of training data to learn features that can easily adapt, or be relevant, to the test time task. Our approach to doing so will be in the form of a *feature transform* that jointly re-embeds the set of features of an instance in a way that resembles how recently popular self-attention mechanisms and Transformers [29,22,26,16] re-embed sets of features.

Being at the low-to-mid-level of most relevant architectures, advances in such feature re-embeddings have a direct impact and wide applicability in instance-specific problems such as few-shot classification [30], clustering [37], patch matching [19] and person re-identification [43], to name but a few.

The general idea of the Self-Optimal-Transport (SOT) feature transform that we propose is depicted and explained in Fig. 1, as part of the general design of networks that work on sets which we illustrate in Fig. 2.

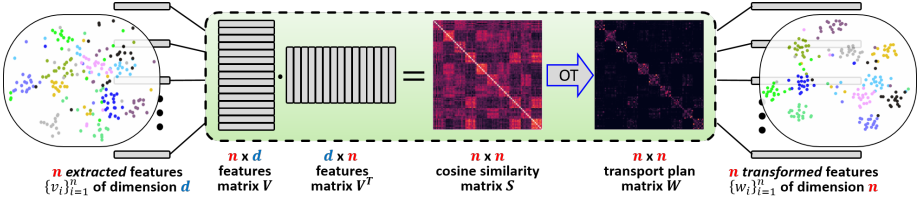


Fig. 1: **The SOT transform:** Its input is a set of n d -dimensional features (each shown as a horizontal gray rectangle, and as a colored point in the input embedding space where color depicts class label or equivalent). Processing is as follows: The unit length (normalized) features are arranged in an $n \times d$ matrix for computing a pairwise $n \times n$ cosine similarity matrix S . Then, the transport-plan matrix W (given a specific OT instance that depends on S) is computed using several Sinkhorn [7] iterations. Finally, the transformed output features are basically the rows of the matrix W . As we claim and observe in real results, the features are re-embedded in a way that is consistently superior for downstream grouping and matching tasks (observed the better formation of the embedded points, e.g. towards applying a linear classifier or an off-the-shelf clustering procedure).

1.1 Overview

We are given an instance of some inference problem, in the form of a set of n items $\{x_i\}_{i=1}^n$, represented as vectors in \mathbb{R}^D , for a fixed dimension D . A generic neural-network (Fig. 2 Left) typically uses a feature embedding (extractor) $F: \mathbb{R}^D \rightarrow \mathbb{R}^d$ (with $d \ll D$), which is applied independently on each input item, to obtain a set of features $V = \{v_i\}_{i=1}^n = \{F(x_i)\}_{i=1}^n$. The features V might be of high quality (concise, unique, descriptive), but are limited in representation since they are extracted based on knowledge acquired for similar examples at train time, with no context of the test time instance that they are part of.

We adapt a rather simple framework (Fig. 2 Right) in which some *transform* acts on the entire set of instance features. The idea is to jointly process the set of features to output an updated set (one for each input feature), that re-embeds each feature in light of the joint statistics of the entire instance. The proposed features transform can be seen as a special case of an attention mechanism [29] specialized to features of instance-specific tasks, with required adaptations. Techniques developed here borrow from and might lend to those used in set-to-set [44,42,25], self-attention [29,26] and transformer [22,16] architectures.

1.2 Contributions

We propose a parameter-less transform T , which can be used as a drop-in addition that can convert a conventional network to an instance-aware one (e.g. from Fig. 2 Left to Right). We propose an optimal-transport based feature transform which is shown to have the following attractive set of qualities. (i) *efficiency*: having real-time inference; (ii) *differentiability*: allowing end-to-end training of the entire ‘embedding-transform-inference’ pipeline of Fig. 2 Right; (iii) *equivariance*: ensuring that the embedding works coherently under any order of the

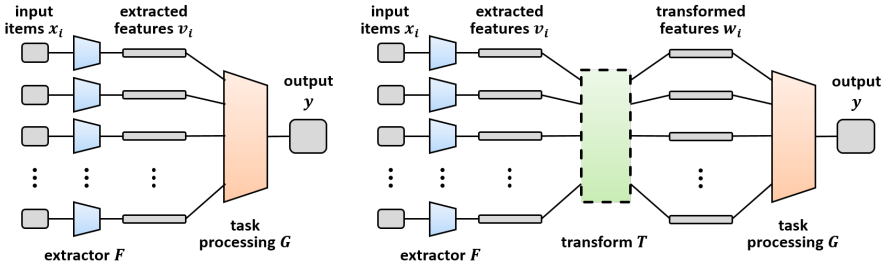


Fig. 2: **Generic designs of networks that act on sets of inputs.** These cover relevant architectures, e.g. for few-shot-classification and clustering. **Left:** A generic network for processing a set of input items typically follows the depicted structure: (i) Each item *separately* goes through a common feature extractor F . (ii) The set of extracted features is the input to a downstream task processing module G . ; **Right:** A more general structure in which the extracted features undergo a *joint* processing by a transform T . Our SOT transform (as well as other attention mechanisms) is of this type and its high-level design (within the ‘green’ module) is detailed in Fig. 1.

input items; (iv) *capturing relative similarity*: The comparison of embedded vectors will include both direct and indirect (third-party) similarity information between the input features; (v) *probabilistic interpretation*: each embedded feature will encode its distribution of similarities to all other features, by conforming to a doubly-stochastic constraint; (vi) *instance-aware dimensionality*: embedding dimension (capacity) is adaptive to input size (complexity).

We provide a detailed analysis of our method and show its flexibility and ease of application to a wide variety of tasks, by incorporating it in leading methods of each kind. A controlled experiment on unsupervised clustering is used to verify its performance, with a detailed analysis. For few-shot-classification we perform an extensive comparison to existing work on several benchmarks, showing that SOT achieves new state-of-art results. Finally, we show that SOT is easily applicable to large-scale benchmarks by using the person re-identification task, for which it consistently improves state-of-art networks that it is incorporated into.

2 Related Work

2.1 Related Techniques

Set-to-set or set-to-feature functions Our method can clearly be categorized along with recent techniques that act jointly on a set of items (typically features) to output an updated set (or a single feature), which are typically used for downstream inference tasks on the items individually, or as a set. The pioneering Deep-Sets [44] formalized fundamental requirements from architectures that process sets. Point-Net [27] presented an influential design that learns local and global features on 3D point-clouds, while Maron *et.al.* [25] study layer

designs that approximate equivariant and invariant functions. Unlike the proposed SOT transform, the joint processing in these methods is very limited, amounting to (Siamese) weight-sharing between separate processes and simple joint aggregations like average pooling.

Self-Attention The introduction of Relational Networks [32] and transformers [38] and their initial applications in vision models [29] have lead to a surge of following successful works [16], many of which are dedicated to few-shot-learning, such as ReNet [15], DeepEMD [45] and FEAT[42]. Different from these methods, SOT is parameterless, and hence can work at test-time on any pre-trained network. In addition, SOT is the only method that provides an explicit probabilistic global interpretation of the instance data.

Optimal Transport Optimal transport (OT) problems are tightly related to measuring and calculating distances between distributions or sets of features. In [7] Cuturi popularized the Sinkhorn algorithm which is a simple, differentiable and fast approximation of entropy-regularized OT problems. The Set transformer [22] uses an OT-based clustering algorithm, SuperGlue [33] uses OT in an end-to-end manner for feature-point matching, and many state-of-the-art methods in few-shot learning, which we review next, have adopted the Sinkhorn algorithm to model relations between features and class representations. The differentiability and efficiency of regularized OT solvers has recently been shown useful in related domains, to derive a differentiable ‘top-k’ operator [41] or for style transfer applications, by viewing styles as a distributions between which distances are approximated [18]. In this work we focus on *self* applications of OT, which enables concise modelings of the relative similarities within a set of items.

2.2 Few-Shot-Classification (FSC)

Few-Shot-Classification [39] is a branch of few-shot-learning in which a classifier needs to learn to recognize classes unseen given a limited number of labeled examples. A FSC task is a self-contained instance that includes both support (labeled) and query (unlabeled) items, hence is a clear instance-specific setup which SOT can handle.

Some leading FSC approaches follow the *meta-learning* (or “learn-to-learn”) principle in which the training data is split into tasks (or episodes) mimicking the test time tasks to which the learner is required to generalize. The celebrated MAML [10] “learns to fine-tune” by learning a network initialization from which it can adapt to a novel set of classes with very few gradient update steps on the labeled examples. In ProtoNet [34], a learner is meta-trained to predict query feature classes, based on distances from support (labeled) class-prototypes in the embedding space. The trainable version of SOT is a meta-learning algorithm, but unlike the above, it is transductive (see ahead) and exploits the task items as a set, while directly assessing the relative similarity relations between its items.

Subsequent works [5,9] have questioned the benefits of meta-learning, advocating the standard transfer learning procedure of fine-tuning pre-trained networks. In particular, they demonstrate the advantages of using larger and more

powerful feature-encoding architectures, as well as the employment of *transductive* inference, which fully exploits the data of the inference task, including unlabeled images. As mentioned, SOT is a purely transductive method, but it is significantly more flexible in its assumptions, since the transform is based on a general probabilistic grouping action. It does not make any assumptions on (nor does it need to know) the number of classes and the number of items per class in an instance.

More recently, *attention mechanisms* were shown to be effective for FSC. We have reviewed some relevant works of this line in the previous section.

Finally, a large number of works have adopted the Sinkhorn Algorithm [7] as a parameterless unsupervised classifier that computes fractional matchings between query embeddings and class centers. Many leading FSC works use this approach, including Laplacian-Shot [50], CentroidNet [13] and PT-MAP [12]. The current state-of-the-art is set by the recent Sill-Net [46], which augments training samples with illumination features that are separated from the images in feature space and by PT-MAP-sf [6], who propose a DCT-based feature embedding network, encoding detailed frequency-domain information that complements the standard spatial domain features. Both methods are based on PT-MAP [12]. SOT uses Sinkhorn to solve an entirely different OT problem - that of matching the set of features to itself, rather than against class representations. Nevertheless, SOT can be incorporated into these methods, immediately after their feature extraction stage.

2.3 Unsupervised Clustering and Person Re-Identification (Re-ID)

These domains are not at the focus of this work therefore we only briefly give some useful pointers for the sake of brevity.

Unsupervised image clustering is an active area of research, with standardised evaluation protocols (from Cifar-10 [20] to different subsets of ImageNet [8]). Prominent works in this area include Deep Adaptive Clustering (DAC) [4], Invariant Information Clustering (IIC) [14] and SCAN [37]. Clustering has recently gained popularity as a means for self-supervision in feature learning, showing excellent results on unsupervised image classification. See for example Deep-Cluster [2] and SWAV [3]. Clustering is a clear case instance specific problem, since most information is relative and unrelated directly to other training data. Our transform can hence be used to upgrade the feature representation quality.

We chose the Re-ID application as another instance-specific problem, which from our point of view differs from the others considered in two main aspects which we find attractive: (i) The tasks are of larger scale - querying thousands of identities against a target set of (tens of) thousands. (ii) The data is much more real-world compared to the carefully curated classification and clustering tasks. See [43] for an excellent recent and comprehensive survey on the topic.

3 Method

Assume we are given a task which consists of an inference problem over a set of n items $\{x_i\}_{i=1}^n$, where each of the items belongs to a space of input items $\Omega \subseteq \mathbb{R}^D$. The inference task can be modeled as $f_\theta(\{x_i\}_{i=1}^n)$, using a learned function f_θ , which acts on the set of input items and is parameterized by a set of parameters θ .

Typically, such functions combine an initial feature extraction stage that is applied independently to each input item, with a subsequent stage of (separate or joint) processing of the feature vectors (see Fig. 2 Left or Right, respectively).

That is, the function f_θ takes the form $f_\theta(\{x_i\}_{i=1}^n) = G_\psi(\{F_\phi(x_i)\}_{i=1}^n)$, where F_ϕ is the feature extractor (or embedding network) and G_ψ is the task inference function, parameterized by ϕ and ψ respectively, where $\theta = \phi \cup \psi$.

The feature embedding $F : \mathbb{R}^D \rightarrow \mathbb{R}^d$, usually in the form of a neural-network (with $d \ll D$), could be either pre-trained, or trained in the context of the task function f , along with the inference function G .

For an input $\{x_i\}_{i=1}^n$, let us define the set of features $\{v_i\}_{i=1}^n = \{F(x_i)\}_{i=1}^n$. In the following, we consider these sets of input vectors and features as real-valued row-stacked matrices $\mathcal{X} \in \mathbb{R}^{n \times D}$ and $\mathcal{V} \in \mathbb{R}^{n \times d}$.

We suggest a novel re-embedding of the feature set \mathcal{V} , using a transform that we denote by T , in order to obtain a new set of features $\mathcal{W} = T(\mathcal{V})$, where $\mathcal{W} \in \mathbb{R}^{n \times n}$. The new feature set \mathcal{W} has an explicit probabilistic interpretation, which is specifically suited for tasks related to classification, matching or grouping of items in the input set \mathcal{X} . In particular, \mathcal{W} will be a symmetric, doubly-stochastic matrix, where the entry w_{ij} (for $i \neq j$) gives the probability that items x_i and x_j belong to the same class or cluster.

The proposed transform $T : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times n}$ (see Fig. 1) acts on the original feature set \mathcal{V} as follows. It begins by computing the squared Euclidean pairwise distances matrix \mathcal{D} , namely, $d_{ij} = \|v_i - v_j\|^2$, which can be computed efficiently as $d_{ij} = 2(1 - \cos(v_i, v_j)) = 2(1 - v_i \cdot v_j^T)$, assuming that the rows of \mathcal{V} are unit normalized. Or in a compact form, $\mathcal{D} = 2(\mathbf{1} - \mathcal{S})$, where $\mathbf{1}$ is the all ones $n \times n$ matrix and $\mathcal{S} = \mathcal{V} \cdot \mathcal{V}^T$ is the cosine similarity matrix of \mathcal{V} .

\mathcal{W} will be computed as the optimal transport (OT) plan matrix between the n -dimensional all-ones vector $\mathbf{1}_n$ and itself, under the cost matrix \mathcal{D}_∞ , which is the distance matrix \mathcal{D} with a very (infinitely) large scalar replacing each of the entries on its diagonal (which were all zero). Explicitly, let $\mathcal{D}_\infty = \mathcal{D} + \alpha I$, where α is a very (infinitely) large constant and I is an $n \times n$ identity matrix.

\mathcal{W} is defined to be the doubly-stochastic matrix¹ that is the minimizer of the functional

$$\mathcal{W} = \arg \min_{\mathcal{W} \in B_n} \langle \mathcal{D}_\infty, \mathcal{W} \rangle \quad (1)$$

where B_n is the set (known as the Birkhoff polytope) of $n \times n$ doubly-stochastic matrices and $\langle \cdot, \cdot \rangle$ stands for the Frobenius (standard) dot-product.

¹ a square ($n \times n$) matrix of non-negative real values, each of whose rows and columns sums to 1

This objective can be minimized using simplex or interior point methods with complexity $\Theta(n^3 \log n)$. In practice, we use the highly efficient Sinkhorn-Knopp method [7], which is an iterative scheme that optimizes an entropy-regularized version of the problem, where each iteration takes $\Theta(n^2)$. Namely:

$$\mathcal{W} = \arg \min_{\mathcal{W} \in B_n} \langle \mathcal{D}_\infty, \mathcal{W} \rangle - \frac{1}{\lambda} h(\mathcal{W}) \quad (2)$$

where $h(\mathcal{W}) = -\sum_{i,j} w_{ij} \log(w_{ij})$ is the Shannon entropy of \mathcal{W} and λ is the entropy regularization parameter.

The *transport-plan* matrix \mathcal{W} that is the minimizer of Eq. (2) is the result of our transform, i.e. $\mathcal{W} = T(\mathcal{V})$ and each of its rows is the re-embedding of each of the corresponding features (rows) in \mathcal{V} . Recall that \mathcal{W} is doubly-stochastic and note that it is symmetric². We next explain its probabilistic interpretation.

The optimization problem in Eq. (1) can be written more explicitly as follows:

$$\min_{\mathcal{W}} \langle \mathcal{D}_\infty, \mathcal{W} \rangle \quad \text{s.t.} \quad \mathcal{W} \cdot \mathbf{1}_n = \mathcal{W}^T \cdot \mathbf{1}_n = \mathbf{1}_n \quad (3)$$

which can be seen to be the same as:

$$\begin{aligned} \min_{\mathcal{W}} \langle \mathcal{D}, \mathcal{W} \rangle \quad \text{s.t.} \quad & \mathcal{W} \cdot \mathbf{1}_n = \mathcal{W}^T \cdot \mathbf{1}_n = \mathbf{1}_n \\ & w_{ii} = 0 \quad \text{for } i = 1, \dots, n \end{aligned} \quad (4)$$

since the use of the infinite weights on the diagonal of \mathcal{D}_∞ is equivalent to using the original \mathcal{D} with a constraint of zeros along the diagonal of \mathcal{W} .

The optimization problem in Eq. (4) is in fact a fractional matching instance between the set of n original features and itself. It can be posed as a bipartite-graph min-cost max-flow instance. The graph has n nodes on each side, representing the original features $\{v_i\}_{i=1}^n$ (the rows of \mathcal{V}). Across the two sides, the cost of the edge (v_i, v_j) is the distance d_{ij} and the edges of the type (v_i, v_i) have a cost of infinity (or can simply be removed). Each ‘left’ node is connected to a ‘source’ node by an edge of cost 0 and similarly each ‘right’ node is connected to a ‘target’ (sink) node by an edge of cost 0. All edges in the graph have a capacity of 1 and the goal is to find an optimal fractional self matching, by finding a min-cost max-flow from source to sink. Note that the maximum flow can easily be seen to be n , but a min-cost flow is sought among the max-flows.

In this set-to-itself matching view, each vector v_i is fractionally matched to the set of all other vectors $\mathcal{V} - \{v_i\}$ based on the pairwise distances, but importantly taking into account the fractional matches of the rest of the vectors in order to satisfy the double-stochasticity constraint³. Therefore, the i th transformed (re-embedded) feature w_i (i th row of \mathcal{W}) is a *distribution* (non-negative entries, summing to 1), where $w_{ii} = 0$ and w_{ij} is the relative belief that features i and j belong to the same ‘class’.

² The symmetry of \mathcal{W} is as a result of the symmetry of \mathcal{D} and the double-stochasticity of \mathcal{W} .

³ The construction constrains the maximum flow to exactly have a total outgoing flow of 1 from each ‘left’ node and a total incoming flow of 1 from each ‘right’ node.

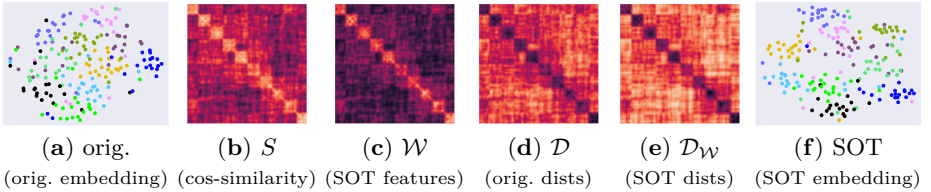


Fig. 3: A close look at the SOT transform as it operates on a 10-way 20-shot supervised clustering task: The input is a set of 200 33-dimensional unit-length feature vectors that are visualized on the plane in (a) using a t-SNE dimension reduction [36], where colors refer to the 10 classes. In (b) is the pairwise cosine similarity matrix \mathcal{S} , which is linearly related to the Euclidean pairwise distances \mathcal{D} shown in (d). Next, in (c) we show the SOT matrix \mathcal{W} whose rows (or columns, symmetrically) consist of our new embedding of the features. These 200-dimensional features are shown again on the plane in (f). Notice the visually apparent improvement in point gathering by class, from (a) to (f), which can be explained by comparing the matrices \mathcal{D} and $\mathcal{D}_{\mathcal{W}}$, which are the self-pairwise distances of the original and SOT embedding respectively. Notice the greater contrast in $\mathcal{D}_{\mathcal{W}}$ between inter- and intra- cluster points. Note, that like in the visualizations of Fig. 1, we show the matrices with row/col order based on the true classes, purely for ease of visualization.

Our final set of features \mathcal{W} is obtained by replacing the diagonal entries from 0s to 1s, namely $\mathcal{W} = \mathcal{W} + I$, where I is the $n \times n$ identity matrix. Please refer to Fig. 3 for a close look at the application of SOT to a toy clustering problem, where we demonstrate visually the improved embedding obtained through examining the pairwise distances before and after the transform. We can now point out some important properties of this new embedding \mathcal{W} :

Direct and Indirect similarity encoding: Each embedded feature encodes its distribution of similarities to all other features. An important property of our embedding is that the comparison of the embedded vectors w_i and w_j includes both *direct* and *indirect* information about the similarity between the features. Please refer to Fig. 4 for a detailed explanation of this property. If we look at the different coordinates k of the absolute difference vector $a = |w_i - w_j|$, SOT captures (i) *direct similarity*: For k which is either i or j , it holds that $a_k = 1 - w_{ij} = 1 - w_{ji}$ ⁴. This amount measures how high (*i.e.* close to 1) is the mutual belief of features i and j about one another. (ii) *indirect (3rd-party) similarity*: For $k \notin \{i, j\}$, we have $a_k = |w_{ik} - w_{jk}|$, which is a comparison of the beliefs of features i and j regarding the (third-party) feature k .

Parameterless-ness: Our proposed transform is parameterless, giving it the flexibility to be used in other pipelines, directly over different kinds of embeddings, without the harsh requirement of retraining the entire pipeline⁵.

⁴ Note: (i) $w_{ii} = w_{jj} = 1$; (ii) $w_{ij} = w_{ji}$ from the symmetry of \mathcal{W} ; (iii) all elements of \mathcal{W} are ≤ 1 and hence the $|\cdot|$ can be dropped;

⁵ Retraining is certainly possible, and beneficial in many situations, but not mandatory, as our experiments work quite well without it.

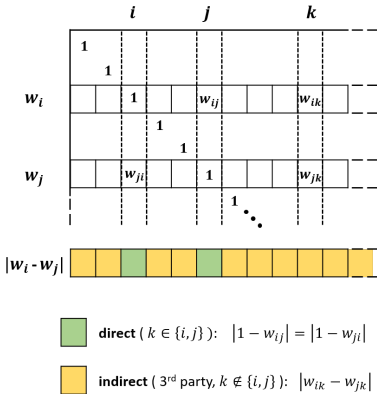


Fig. 4: **The (symmetric) embedding matrix W and the absolute difference between its i th and j th rows:** We examine the vector $|w_i - w_j|$: (i) Its i th and j th coordinates equal $|1 - w_{ij}| = |1 - w_{ji}|$, giving the *direct* similarity between the original features, since this amount (in green) is high (closer to 1). ; (ii) Its k th coordinate (for any $k \notin \{i, j\}$) gives $|w_{ik} - w_{jk}|$ which is an *indirect* (third-party) comparison between the original features through the k th feature. Similarity (in yellow) is stronger when features i and j have a similar belief regarding feature k , i.e. w_{ik} and w_{jk} are close.

Differentiability: Due to the differentiability of Cuturi’s [7] version of Sinkhorn, back-propagating through the SOT can be done naturally, hence it is possible to (re-)train the hosting network to adapt to the SOT, if desired.

Equivariance: The embedding works coherently with respect to any change of order of the input items (features). This can be shown by construction, since min-cost max-flow solvers as well as the Sinkhorn OT solver are equivariant with respect to permutations of their inputs.

Explainability: The non-parametric nature gives SOT an advantage over other set-to-set methods such as transformers in that its output is interpretable (e.g. by visually inspecting the transport-plan matrix W), with a clear probabilistic characterization of the relations it had found.

Task-Aware Dimensionality: SOT has the unique property that the dimension of the embedded feature depends on (equals) the number of features. On the one hand, this is a desired property, since it is only natural that the feature dimensionality (capacity) depends on the complexity of the task, which typically grows with the number of features (think of the inter-relations which are more complex to model). On the other hand, it might impose a problem in situations in which the downstream calculation that follows the feature embedding expects a fixed input size, for example a pre-trained non-convolutional layer. Nevertheless, in many situations the downstream computation has the flexibility to work with varying input dimensions. Also, in most benchmarks the instance set sizes are fixed, allowing for a single setting of sizes to work throughout.

4 Implementation details

Datasets: We consider three different applications to evaluate the performance of our method. For *unsupervised clustering* we designed a specialized synthetic data set with the goal of enabling controlled experiments over a wide range of difficulties, which are determined by data dimensionality and in-cluster spread.

For *few-shot classification* we use the standard benchmarks in the literature. The *MiniImagenet* [39] dataset is a subset of *Imagenet* [31] that contains 100 classes and 600 images of size 84×84 per class. We follow the standard setup of using 64 classes for training and 16 and 20 novel classes for validation and testing. The *CIFAR-FS* [1] dataset includes 100 classes with 600 images of size 32×32 per-class. We used the same splits as in *MiniImagenet* for this dataset. The *CUB* [40] dataset includes 200 classes of bird species and has 11,788 images of size 84×84 pixels in total. We followed the split suggested in [11] into 100 base classes, 50 validation classes and 50 novel classes.

For *person re-identification* (ReID) we use two common large-scale datasets. The *Market-1501* [47] and *CUHK03* [23] dataset consists of 1,501 and 1,467 identities and a total of 32,668 and 14,097 images taken from 6 cameras. We use the validation and test sets according to the splits in [49].

Pre-training: We pre-trained ProtoNet [34] with a 4-layer Convolution network adapting the procedures of [34] for training both with and without SOT, training on a 5-way (5/1)-shot 15-query task, using ADAM [17] with learning rate 0.01 and step size of 20 over 100 episodes (tasks) per epoch.

Fine-tuning: We perform fine-tuning on two types of backbone residual networks - a resnet-12 as used in [42] and a WRN-28-10 as used in [24]. For ProtoNet [34] and ProtoNet-SOT, we fine-tune the base network with parameters taken from [42]. For PTMAP-SOT, we use meta-training with batches of a single 10-way 5-shot 15-query task per batch. We use ADAM with learning rate $5e-5$ that decreases with step size 10 for 25 epochs. We train the WRN-28-10 and the resnet-12 backbones for 800 and 100 episodes respectively per epoch.

Hyper-parameters: SOT has two hyper-parameters which were chosen through cross-validation and were kept fixed for each of the applications over all datasets.

(i) The number of Sinkhorn iterations for computing the optimal transport plan was fixed to 10. (ii) The entropy regularization parameter λ (Eq. (3)) was set to 0.1 for clustering and few-shot-learning experiments and to 1.0 for the ReID experiments. We further ablate these in the supplementaries.

5 Results

5.1 Clustering on the Sphere

We first demonstrate the effectiveness of SOT using a controlled synthetically generated clustering experiment, with $k = 10$ cluster centers that are distributed uniformly at random on a d -dimensional unit-sphere, and 20 points per cluster (200 in total) that are perturbed around the cluster centers by Gaussian noise of increasing standard deviation, of up to 0.75, followed by a re-projection back to the sphere by dividing each vector by its L_2 magnitude. We also apply dimensionality reduction with PCA to $d = 50$, for dimensions above 50.

We performed the experiment over a logarithmic 2D grid of combinations of data dimensionalities d in the range [10, 1234] and Gaussian in-cluster noise STD in the range [0.1, 0.75]. Refer to Fig. 9 (i) for a visualization of the data generation process.

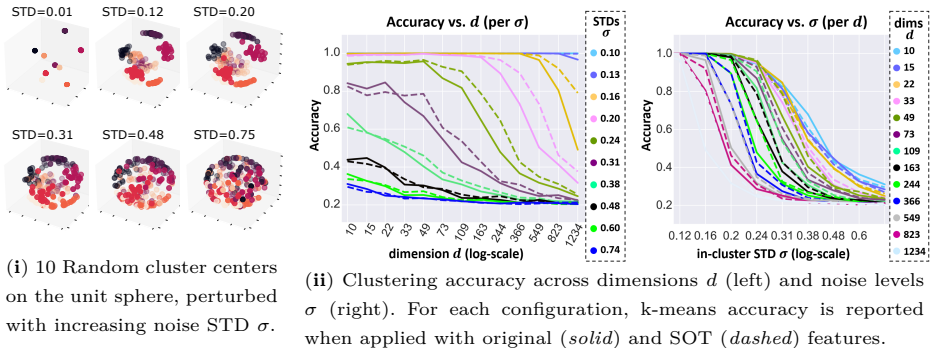


Fig. 5: **Clustering on the d -dimensional sphere.** Left (i): the data generation process (illustrated for the 3D case). Right (ii): detailed k-means accuracy results. The SOT (dashed) features give superior results throughout a majority of the space of settings.

Each point is represented by its d -dimensional euclidean coordinates vector, where the baseline clustering is obtained by running k-means on these location features. In addition, we run k-means on the set of features that has undergone SOT. Hence, the benefits of the transform (embedding) are measured indirectly through the accuracy⁶ achieved by running k-means on the embedded vs. original vectors. Evaluation results are reported in Fig. 9 (ii) as averages over 10 runs, by plotting accuracy vs. dimensionality (for different noise STDs) and accuracy vs. noise STDs (for different dimensionalities). The results show (i) general accuracy gains and robustness to wide ranges of data dimensionality (ii) the ability of SOT to find meaningful representations that enable clustering quality to degrade gracefully with the increase in cluster noise level. Note that the levels of noise are rather high, as they are relative to a unit radius sphere (a 3-dimensional example is shown at the top of the figure). We provide further details on this experiment in the supplementaries.

5.2 Few-Shot Classification (FSC)

Our main experiment is a comprehensive evaluation on the standard few-shot classification benchmarks *MiniImagenet* [39], *CIFAR-FS* [1], and *CUB* [40], with detailed results in Tables 1 and 2. For *MiniImagenet* (Table 1) we report on both versions “SOT_p” and “SOT_t” over a range of backbone architectures, while for the smaller datasets *CIFAR-FS* and *CUB* (Table 2) we focus on the ‘drop-in’ version “SOT_p” and only the strongest wrn-28-10 architecture.

One goal here is to show that we can achieve new state-of-the-art FSC results, when we build on current state-of-the-art. But more importantly, we demonstrate the flexibility and simplicity of applying SOT in this setup, with improvements in the entire range of testing, including: (i) when building on different ‘hosting’ methods; (ii) when working above different feature embeddings of different

⁶ Accuracy is measured by comparison with the optimal permutation of the predicted labels, found by the Hungarian Algorithm [21].

<i>method</i>	<i>transductive</i>	<i>backbone</i>	5way-1shot	5way-5shot
MAML(*) [10]	X	conv-4	46.47	62.71
RelationNet(*) [35]	X	conv-4	49.31	66.60
ProtoNet(#) [34]	X	conv-4	49.10	66.79
FEAT(\$) [42]	X	conv-4	55.15	71.61
ProtoNet-SOT _p	✓	conv-4	<u>54.01</u> (+10.2%)	69.39 (+3.9%)
ProtoNet-SOT _t	✓	conv-4	53.70 (+9.3%)	<u>70.40</u> (+5.4%)
ProtoNet(#) [34]	X	resnet-12	62.39	80.33
DeepEMD(\$) [45]	X	resnet-12	65.91	82.41
FEAT(\$) [42]	X	resnet-12	66.78	82.05
RENet(\$) [15]	X	resnet-12	67.60	82.58
PTMAP(#) [12]	✓	resnet-12	76.90	85.20
ProtoNet-SOT _p	✓	resnet-12	67.34 (+7.9%)	81.84 (+1.6%)
ProtoNet-SOT _t	✓	resnet-12	67.90 (+8.8%)	83.09 (+3.2%)
PTMAP-SOT _p	✓	resnet-12	78.35 (+1.9%)	86.01 (+1.0%)
PTMAP-SOT _t	✓	resnet-12	<u>77.30</u> (+0.5%)	<u>85.49</u> (+0.3%)
ProtoNet(&) [34]	X	wrn-28-10	62.60	79.97
PTMAP(\$) [12]	✓	wrn-28-10	82.92	88.80
Sill-Net(\$) [46]	✓	wrn-28-10	82.99	89.14
PTMAP-SF(\$) [6]	✓	wrn-28-10	<u>84.81</u>	<u>90.62</u>
PTMAP-COSINE	✓	wrn-28-10	74.60 (-10.0%)	84.68 (-4.6%)
PTMAP-SOFTMAX	✓	wrn-28-10	80.08 (-3.4%)	83.83 (-5.6%)
PTMAP-SOT _p	✓	wrn-28-10	83.19 (+0.3%)	89.56 (+0.9%)
PTMAP-SOT _t	✓	wrn-28-10	84.18 (+1.5%)	90.51 (+1.9%)
Sill-Net-SOT _p	✓	wrn-28-10	83.35 (+0.4%)	89.65 (+0.6%)
PTMAP-SF-SOT _p	✓	wrn-28-10	85.59 (+0.9%)	91.34 (+0.8%)

Table 1: **Few-Shot Classification (FSC)** accuracy on **MiniImagenet** [39]. The improvements introduced by the variants of SOT (percentages in brackets) are in comparison with each respective baseline hosting method. **Bold** and underline notations highlight best and second best results per backbone. (*) = from [5] ; (&) = from [50] ; (\$) = from the method’s paper itself ; (#) = our implementation ;

complexity backbones; and (iii) whether retraining the hosting network or just dropping-in SOT and performing standard inference.

To evaluate the performance of the proposed SOT, we applied it to previous FSC methods including the very recent state-of-the-art (PT-MAP [12], Sill-NET [46] and PT-MAP-SF [6]) as well as a to more conventional methods like the popular ProtoNet [34]. The detailed results are presented in Tables 1 and 2) for the different datasets. Note that SOT is by nature a transductive method⁷, hence we marked its results as so, regardless of whether the hosting network is transductive or not. In the following, we discuss the two modes in which our transform can be used in existing FSC methods.

⁷ SOT is transductive in the sense that it needs to jointly process the data, but importantly, unlike other methods it does not gain its benefit in being so from making limiting assumptions about the structure of the instance, like knowing the number of classes, or the number of items per class.

<i>FSC benchmark</i>	<i>CIFAR-FS</i> [1]		<i>CUB</i> [40]	
<i>method</i>	5way-1shot	5way-5shot	5way-1shot	5way-5shot
PTMAP(\$) ^[12]	87.69	90.68	91.55	93.99
Sill-Net(\$) ^[46]	87.73	91.09	94.73	96.28
PTMAP-SF(\$) ^[6]	89.39	92.08	95.45	96.70
PTMAP-SOT _p	87.37 (-0.4%)	91.12 (+0.5%)	91.90 (+0.4%)	94.63 (+0.7%)
Sill-Net-SOT _p	87.30 (-0.5%)	91.40 (+0.3%)	94.86 (+0.1%)	96.61 (+0.3%)
PTMAP-SF-SOT _p	89.94 (+0.6%)	92.83 (+0.8%)	95.80 (+0.4%)	97.12 (+0.4%)

Table 2: **Few-Shot Classification (FSC)** accuracy on *CIFAR-FS* [1] and *CUB* [40].

SOT insertion *without* network retraining (notated by SOT_p in Tables 1 and 2). Recall that the proposed transform is *non-parametric*. As such, we can simply apply it to a trained network at inference, without the need to re-train. This basic ‘drop-in’ use of SOT consistently, and in many cases also significantly, improved the performance of the tested methods, including state-of-the-art, across all benchmarks and backbones. SOT_p gave improvements of around 3.5% and 1.5% on 1 and 5 shot *MiniImagenet* tasks. This improvement without re-training the embedding backbone network shows SOT’s effectiveness in capturing meaningful relationships between features in a very general sense.

SOT insertion *with* network retraining (notated by SOT_t in Table 1). Due to its *differentiability* property, the proposed method can be applied while training and hence we expect an adaptation of the hosting network’s parameters to the presence of the transform with a potential for improvement. To evaluate this mode, we focused on the *MiniImagenet* benchmark [39], specifically on the same configurations that we used without re-training, to enable a direct comparison. The results in Table 1 show additional improvements in almost every method. SOT_t gave improvements of around 5% and 3% on 1 and 5 shot *MiniImagenet* tasks, further improving on the pre-trained counterpart. This result indicates the effectiveness of training with SOT in an end-to-end fashion.

Ablations Within the context of few-shot learning on *MiniImagenet*, we performed several ablation studies. In Table 1, the networks ‘PTMAP-COSINE’ and ‘PTMAP-SOFTMAX’ stand for the obvious baseline attempts (found to be unsuccessful) that work in the line of our approach, without the specialized OT-based transform. In the former, we take the output features to be the rows of the (un-normalized) matrix \mathcal{S} (rather than those of \mathcal{W}) and in the latter we also normalize its rows using soft-max. In the supplementaries we ablate on SOT’s two parameters - the number of Sinkhorn iterations and the entropy term λ .

5.3 Person re-Identification (Re-ID)

In this section, we explore the possibility of using SOT on large-scale datasets by considering the Person re-Identification task. Given a set of *query* images and a large set of *gallery* images, the task is to rank the similarities of each single query against the gallery. This is done by computing specialized image features among which similarities are based on Euclidean distances. SOT is applied to such pre-

<i>ReID benchmark</i>	<i>CUHK03-detected</i> [23]		<i>Market-1501</i> [47]	
<i>network</i>	mAP	Rank-1	mAP	Rank-1
TopDBNet [28]	72.9	75.7	85.7	94.3
TopDBNet-rerank [28]	87.1	87.1	94.0	95.3
TopDBNet-SOT _p	77.9 (+6.9%)	80.4 (+6.2%)	88.1 (+2.8%)	94.4 (+0.1%)
TopDBNet-rerank-SOT _p	87.9 (+0.9%)	88.0 (+1.0%)	94.0 (0.0%)	95.0 (-0.3%)

Table 3: **Re-ID** results on *CUHK03* [23] and *Market-1501* [47]

computed image features, refining them with the strong relative information that it is able to capture by applying it on the union of all query and gallery features. We adapted a pre-trained standard resnet-50 architecture [49] and the popular TopDBNet [28], which we tested on the large-scale ReID benchmarks *CUHK03* [23] (on the ‘detected’ version and similar results on the ‘labeled’ version in the supplementaries) and *Market-1501* [47], with and without the re-ranking [48] procedure. For evaluation, we followed their conventions and compare results using the mAP (mean Average Precision) and Rank-1 metrics.

The results in Table 3 show a consistent benefit in using SOT within the different networks. For *CUHK03*, the results improved by a large margin of +6.8% in mAP for the best configuration. These results demonstrate that the proposed SOT scales well to large-scale problems (with number of features in the thousands) and is attractive for a variety of applications. ReID is not the main focus of this work, hence, we did not re-train the hosting networks with SOT included. Further research is required to measure the possible effects of doing so.

6 Conclusions, Limitations and Future Work

In this paper, we explored the idea of utilizing global information of features, for instance-specific problems such as clustering, few-shot learning, and person re-identification. We proposed a novel module: the Self-Optimal-Transport (SOT) - a features transform that is non-parametric, differentiable and which can capture high-level relationships between data points in problems of this nature. The proposed method outperforms state-of-the-art networks on popular few-shot classification benchmarks and shows consistent improvements on tested ReID benchmarks. Based on these promising results, we believe that exploring its full potential can lead to improvements in a variety of fields and open new possibilities.

In future work, we plan to address some current limitations. (i) Regarding the output dimensionality of the embedding, which is dictated by the input set size. We will aim at being able to obtain an arbitrary dimension, for increased usage flexibility; (ii) We plan to investigate the usage of SOT in unsupervised settings, which would be possible by utilizing its informative representation for self-supervision; (iii) It would likely be beneficial to have a variant of SOT in which the transform is enriched with learnable parameters, similar to transformers, to extend its modeling capacity even further; (iv) SOT is purely transductive. We plan to explore non-transductive variants, possibly by comparing each sample separately to the support or gallery sets.

References

1. Luca Bertinetto, Joao F. Henriques, Philip Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. In *International Conference on Learning Representations (ICLR)*, 2019. 10, 11, 13
2. Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 5
3. Mathilde Caron, Ishan Misra, J. Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *ArXiv*, abs/2006.09882, 2020. 5
4. Jianlong Chang, Lingfeng Wang, Gaofeng Meng, Shiming Xiang, and Chunhong Pan. Deep adaptive image clustering. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017. 5
5. Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *International Conference on Learning Representations (ICLR)*, 2018. 4, 12
6. Xiangyu Chen and Guanghui Wang. Few-shot learning by integrating spatial and frequency representation. *arXiv preprint arXiv:2105.05348*, 2021. 5, 12, 13
7. Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2013. 2, 4, 5, 7, 9, 18
8. Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2009. 5
9. Guneet S Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. In *International Conference on Learning Representations (ICLR)*, 2020. 4
10. Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, 2017. 4, 12
11. Yuqing Hu, Vincent Gripon, and Stéphane Pateux. Exploiting unsupervised inputs for accurate few-shot classification. *ArXiv*, abs/2001.09849, 2020. 10
12. Yuqing Hu, Vincent Gripon, and Stéphane Pateux. Leveraging the feature distribution in transfer-based few-shot learning. In *arXiv preprint arXiv:2006.03806*, 2020. 5, 12, 13
13. Gabriel Huang, Hugo Larochelle, and Simon Lacoste-Julien. Are few-shot learning benchmarks too simple? solving them without task supervision at test-time. *arXiv preprint arXiv:1902.08605*, 2019. 5
14. Xu Ji, Joao F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 5
15. Dahyun Kang, Heeseung Kwon, Juhong Min, and Minsu Cho. Relational embedding for few-shot classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 4, 12
16. Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *arXiv preprint arXiv:2101.01169*, 2021. 1, 2, 4
17. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 10

18. Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. Style transfer by relaxed optimal transport and self-similarity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 4
19. Simon Korman and Shai Avidan. Coherency sensitive hashing. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2015. 1
20. Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009. 5
21. Harold W Kuhn. The hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 2, 1955. 11
22. Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International Conference on Machine Learning (ICML)*, 2019. 1, 2, 4
23. Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang. Deepreid: Deep filter pairing neural network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 10, 14
24. Puneet Mangla, Nupur Kumari, Abhishek Sinha, Mayank Singh, Balaji Krishnamurthy, and Vineeth N Balasubramanian. Charting the right manifold: Manifold mixup for few-shot learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2020. 10
25. Haggai Maron, Or Litany, Gal Chechik, and Ethan Fetaya. On learning sets of symmetric elements. In *International Conference on Machine Learning (ICML)*, 2020. 2, 3
26. Grégoire Mialon, Dexiong Chen, Alexandre d’Aspremont, and Julien Mairal. A trainable optimal transport embedding for feature aggregation and its relationship to attention. In *International Conference on Learning Representations (ICLR)*, 2021. 1, 2
27. 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 (CVPR)*, 2017. 3
28. Rodolfo Quispe and Helio Pedrini. Top-db-net: Top dropblock for activation enhancement in person re-identification. *25th International Conference on Pattern Recognition (ICPR)*, 2020. 14
29. Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 1, 2, 4
30. Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2017. 1
31. Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 2015. 10
32. Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 4
33. 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 (CVPR)*, 2020. 4

34. Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. [4](#), [10](#), [12](#), [18](#), [20](#)
35. Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [12](#)
36. Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research (JMLR)*, 9(11), 2008. [8](#)
37. Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. Scan: Learning to classify images without labels. In *European Conference on Computer Vision (ECCV)*. Springer, 2020. [1](#), [5](#)
38. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. [4](#)
39. Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (NeurIPS)*, 2016. [4](#), [10](#), [11](#), [12](#), [13](#), [18](#)
40. Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge J. Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. [10](#), [11](#), [13](#)
41. Yujia Xie, Hanjun Dai, Minshuo Chen, Bo Dai, Tuo Zhao, Hongyuan Zha, Wei Wei, and Tomas Pfister. Differentiable top-k with optimal transport. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. [4](#)
42. Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation with set-to-set functions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#), [4](#), [10](#), [12](#)
43. Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi. Deep learning for person re-identification: A survey and outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2021. [1](#), [5](#)
44. Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. [2](#), [3](#)
45. Chi Zhang, Yujun Cai, Guosheng Lin, and Chunhua Shen. Deepemd: Few-shot image classification with differentiable earth mover’s distance and structured classifiers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [4](#), [12](#)
46. Haipeng Zhang, Zhong Cao, Ziang Yan, and Changshui Zhang. Sill-net: Feature augmentation with separated illumination representation. *arXiv preprint arXiv:2102.03539*, 2021. [5](#), [12](#), [13](#)
47. Liang Zheng, Liye Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015. [10](#), [14](#), [18](#)
48. Zhun Zhong, Liang Zheng, Donglin Cao, and Shaozi Li. Re-ranking person re-identification with k-reciprocal encoding. 2017. [14](#)
49. Kaiyang Zhou and Tao Xiang. Torchreid: A library for deep learning person re-identification in pytorch. *arXiv preprint arXiv:1910.10093*, 2019. [10](#), [14](#)
50. Imtiaz Masud Ziko, Jose Dolz, Eric Granger, and Ismail Ben Ayed. Laplacian regularized few-shot learning. In *International Conference on Machine Learning (ICML)*, 2020. [5](#), [12](#)

Appendix

A ablation studies

A.1 Sinkhorn iterations

In Table 4 we ablate the number of normalization iterations in the Sinkhorn-Knopp (SK) [7] algorithm at test-time. We measured accuracy on the validation set of *MiniImagenet* [39], using ProtoNet-SOT_p (which is the non-fine-tuned drop-in version of SOT within ProtoNet [34]). As was reported in prior works following [7], we empirically observe that a very small number of iterations (around 5) provide rapid convergence. We observed similar behavior for other hosting methods, and therefore chose to use a fixed number of 10 iterations throughout the experiments.

<i>method</i>	iterations	5way-1shot	5way-5shot
ProtoNet-SOT _p	1	70.71	83.79
ProtoNet-SOT _p	2	71.10	84.01
ProtoNet-SOT _p	4	71.18	84.08
ProtoNet-SOT _p	8	71.20	84.10
ProtoNet-SOT _p	16	71.20	84.10

Table 4: **Sinkhorn iterations ablation study:** See text for details.

A.2 OT entropy regularization parameter λ

We measured the impact of using different values of the optimal-transport entropy regularization parameter λ (the main parameter of the Sinkhorn algorithm) on a variety of configurations (ways and shots) in Few-Shot-Classification (FSC) on *MiniImagenet* [39] in Fig. 6 as well as on the Person-Re-Identification (RE-ID) experiment on Market-1501 [47] in Fig. 7. In both cases, the ablation was executed on the validation set.

For FSC, in Fig. 6, the **left** plot shows that the effect of the choice of λ is similar across tasks with a varying number of ways. The **right** plot shows the behavior as a function of λ across multiple shot-values, where the optimal value of λ can be seen to have a certain dependence on the number of shots. Recall that we chose to use a fixed value of $\lambda = 0.1$, which gives an overall good accuracy trade-off. Note that a further improvement could be achieved by picking the best values for the particular cases. Notice also the log-scale of the x-axes to see that performance is rather stable around the chosen value.

For Re-ID, in Fig. 7, we experiment with a range of λ values on the validation set of the Market-1501 dataset. The results (shown both for mAP and rank-1 measures) reveal a strong resemblance to those of the FSC experiment in Fig. 6, however, the optimal choices for λ are slightly higher, which is consistent with the dependence on the shots number, since the re-ID tasks are typically

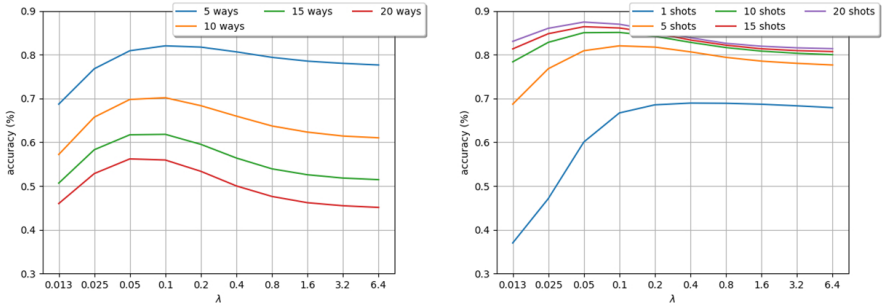
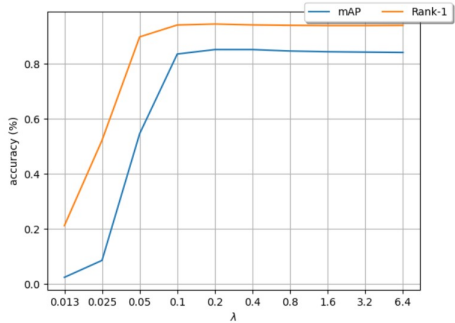


Fig. 6: **Ablation study on λ in *Few-Shot-Classification (FSC)***: Considering different ‘ways’ (left), and different ‘shots’ (right). See text for details.

large ones. In this re-ID ablation, we found that a value of $\lambda = 0.25$ gives good results across both datasets. We ask to note that in the paper we mistakenly reported that we used $\lambda = 1.0$, while in practice all our results were obtained using $\lambda = 0.25$.

Fig. 7: **Ablation study on λ in *Person-Re-Identification (Re-ID)***: Using the validation set of the Market-1501 dataset and considering both mAP and Rank-1 measures. See text for details.



B Unsupervised Clustering - further details

In this section we provide further details (due to lack of space in main paper) on the experiment on unsupervised clustering on the unit sphere (Exp. 5.1).

B.1 Separation between inter- and intra-class features

Fig. 8 depicts the average percentile of the in-class and out-class distances computed by the original and the SOT points. Each panel presents the distributions of both types of distances, for instances of a different level of noise. We compute the mean (and plus-minus half-std) percentiles, with respect to the entire set of pair-wise distances, for a fixed level of in-class noise (increasing from top-left to bottom-right panels), for a range of data dimensionality (x-axis). Naturally, the

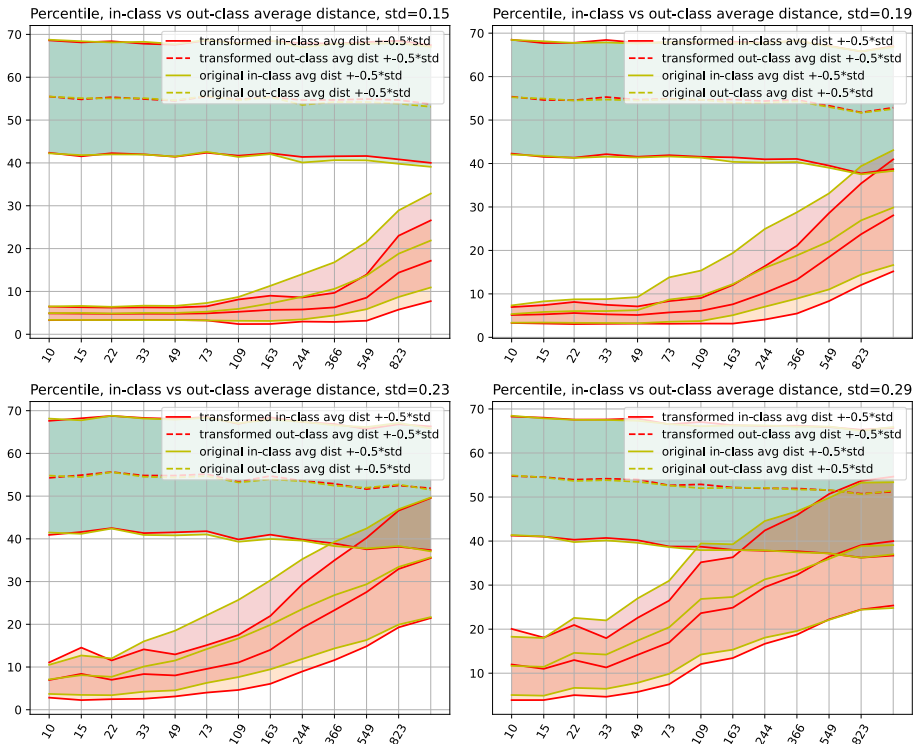


Fig. 8: **intra (in) vs. inter (out) class distances before and after SOT.** A strong indicative property of an embedding that works on class (cluster) objects is its ability to reduce embedded intra-class (pink shaded) pairwise feature distances compared to inter-class (green shaded) ones. SOT (red lines) consistently improves this separation compared to the baseline (brown lines) - leading to better downstream clustering and classification. **x-axis** represents data dimensionality; **y-axis** represents percentiles of pair-wise distances; The four panels present results for the noise standard deviations levels in $\{0.15, 0.19, 0.23, 0.29\}$

overlap between in-class and between-class distances increases both with dimensionality and with in-class noise. Nevertheless, across almost all sampled points, the situation is far better after SOT application (in red), compared to prior to SOT application (in brown). This can explain, in part, the effectiveness of using SOT in Euclidean-based downstream methods, like k -means and ProtoNet [34].

B.2 Evaluation on an extended set of measures

In Fig. 9 we evaluate the performance on additional popular clustering metrics, NMI and ARI (in addition to the accuracy measure we reported on in Figure 5 of the paper). The results shows the same trend as with accuracy, perhaps even stronger for NMI, where SOT significantly improves the clustering performance.

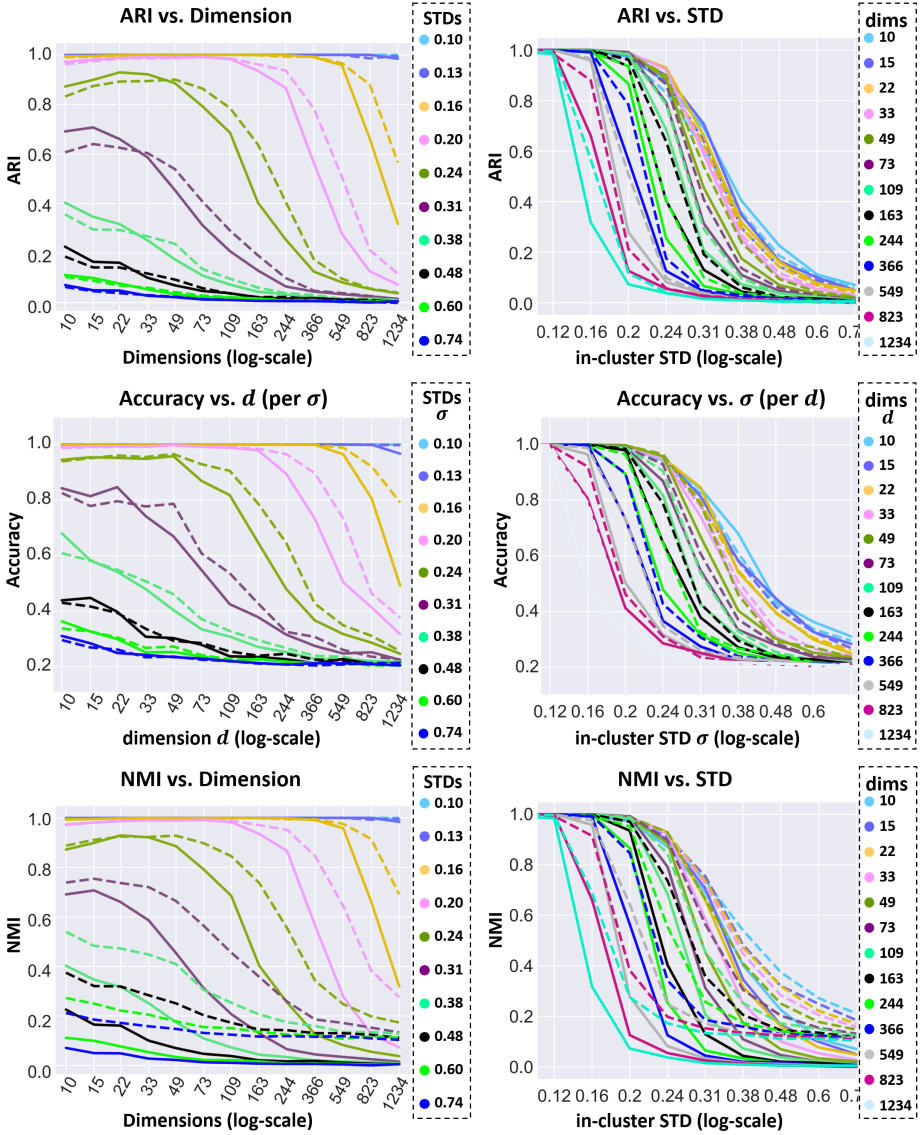


Fig.9: A controlled clustering experiment on the d -dimensional sphere - Extension of results from Figure 5 of the paper, with 2 additional measures: It can be seen that the SOT (dashed - -) shows superior results in all aspects (see text for explanations and interpretation). Clustering accuracy across different noise levels σ and dimensions d . **Note:** For each configuration, SOT is shown by a *dashed* line while the baseline features are shown by a *solid* line. For all 3 measures - the higher the better.