# Spatial Transformer $K$-Means

**Romain Cosentino** [1]  **Randall Balestriero** [1]  **Yanis Bahroun** [2]  **Anirvan Sengupta** [2 3]  **Richard Baraniuk** [1]
**Behnaam Aazhang** [1]

## Abstract

$K$-means defines one of the most employed centroid-based clustering algorithms with performances tied to the data's embedding. Intricate data embeddings have been designed to push $K$-means performances at the cost of reduced theoretical guarantees and interpretability of the results. Instead, we propose preserving the intrinsic data space and augment $K$-means with a similarity measure invariant to non-rigid transformations. This enables (i) the reduction of intrinsic nuisances associated with the data, reducing the complexity of the clustering task and increasing performances and producing state-of-the-art results, (ii) clustering in the input space of the data, leading to a fully interpretable clustering algorithm, and (iii) the benefit of convergence guarantees.

## 1. Introduction

Clustering algorithms aim at discovering patterns in the data that enable their characterization, identification, and separation. The development of such a framework without any prior information regarding the data remains one of the milestones of machine learning that would assist clinicians, physicists, and data scientists, among others, with a better pattern discovery tool (Bertsimas et al., 2020; Greene & Cunningham, 2005).

While supervised learning has been converging toward the almost exclusive use of Deep Neural Networks (DNN), avoiding the development of handcrafted features to provide the desired linearly separable embedding map, unsupervised clustering algorithms take various forms depending on the application at hand (Estivill-Castro, 2002; Ma et al., 2019; Wagstaff et al., 2001). For instance, the usage of SIFT features combined with clustering algorithm for medical imaging (Nam et al., 2009), the extraction of DNNs em-

bedding used as the input of the $K$-means algorithm for computer vision tasks (Xie et al., 2016), and the combination of signal-processing features extractors combined with Gaussian mixture model to understand the nature of the various seismic activities (Seydoux et al., 2020). The important role of clustering algorithms in assisting medical diagnoses as well as scientific discoveries highlight the importance of the development of an *interpretable* and *theoretically guaranteed* tool (Dolnicar, 2003; Xu & Wunsch, 2010).

In this work, we focus our attention on the $K$-means clustering algorithm (MacQueen, 1967) and its application to 2-dimensional signals, such as images or time-frequency representations. Well-known for its simplicity, efficiency, and interpretability, the $K$-means algorithm partitions the data space into $K$ disjoint regions. Each region is represented by a centroid, and each datum is assigned to the closest centroid's region. The integral part in the design of a clustering algorithm is the choice of an appropriate distance, and the number of clusters (Frey & Jojic, 2002; He et al., 2013; Raytchev & Murase, 2001). While the Euclidean distance makes the design of the algorithm straightforward, this measure of similarity might omit the geometrical relationships between data points (Steinbach et al., 2004). In fact, a small rigid perturbation of an image, such as rotation or translation, is enough to change the cluster assignment.

There are two major difficulties in constructing a distance for a clustering algorithm; on the one hand, the metric should take into account the geometry of the data, e.g., be invariant to rigid transformations for images, and on the other hand, the metric should be interpretable as it is tied to the interpretability of the algorithm (Steinbach et al., 2004).

In this work, we tackle these two difficulties by introducing in our similarity measure the spatial transformations inherent to the geometry of the data at hand. In particular, we: $(i)$ formulate an interpretable and theoretically guaranteed $K$-means framework capable of exploiting the symmetry within the data, $(ii)$ extend prior work on metrics invariant to rigid transformations to non-rigid transformations, thus taking into account a more realistic set of nuisances and $(iii)$ allow the learnability of the symmetry underlying the data at hand, therefore enabling the exploration of data where the equivalence classes are yet to be determined.

---

[1]Rice University [2]Flatiron Institute [3]Rutgers University. Correspondence to: Romain Cosentino <rom.cosentino@gmail.com>.

,

To learn the symmetry in the data and perform their transformations, we will use the spatial transformer framework, which was successfully introduced in Jaderberg et al. (2015). This allows us to provide a learnable metric invariant to non-rigid transformations that is used as the $K$-means distortion error.

While many approaches to learn and estimate non-rigid transformations have been proposed, we will follow one of nowadays mainstream approaches developed in Jaderberg et al. (2015) where the Thin Plate Spline is used as a differentiable deformation model. Our attempt is, in fact, not to compare among deformation models but to consider a way to approach the learnability of invariances in an unsupervised setting such that it is effective, tractable, and interpretable.
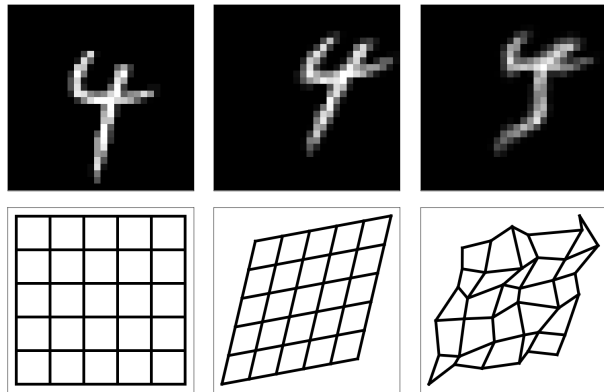
Our contributions can be summarized as follows:

- We propose a novel approach to tackle clustering using a novel adaptive similarity measure within the $K$-means framework that considers non-rigid transformations, Sec. 3.1.
- We derive an appropriate update rule for the centroids that drastically improves both the interpretability of the centroids and their quality, Sec. 3.2.
- We provide convergence guarantees and geometrical interpretations of our approach, Sec. 3.3, 3.4.
- Finally, we show numerically that our unsupervised algorithm competes with state-of-the-art methods on various datasets while benefiting from interpretable results, Sec. 5.

## 2. Background

### 2.1. Invariant Metrics

The development of measures invariant to specific deformations has been under investigation in the computer vision community for decades (Fitzgibbon & Zisserman, 2002; Lim et al., 2004; Simard et al., 2012). By considering affine transformations such as shearing, translation, and rotation of the data as being nuisances, these approaches propose a distance that reduces the variability intrinsic to high-dimensional images. These works are considered as appearance manifold-based framework; that is, the distance are quantified by taking into account geometric proximity (Basri et al., 1998; Ho et al., 2003; Murase & Nayar, 1993; Su & Chou, 2001).

While the development of affine invariant metrics is pretty standard, their extension to more general non-rigid transformations requires more attention. Recently, various deep learning methods proposed ways to learn diffeomorphic transformations (Balakrishnan et al., 2018; Dalca et al., 2019; Detlefsen et al., 2018; Lohit et al., 2019; Shapira Weber et al., 2019). Others adopt a more theoretically grounded



*Figure 1.* **Spatial Transformations -** Visualizations of a sample taken from the MNIST dataset and its transformed versions. Each image results from the application of the spatial transformer that take as input the original signal (top left), and the grid displayed below its transformed version. (*Left*) we observe the original image and its associated original transformation grid, which corresponds to the identity transform. (*Middle*) the image has been transformed by the affine transformation induced by the associated grid. (*Right*) the image transformed by the non-rigid transformation using the TPS induced by the grid below it.

approach based on group theory as in Allassonnière et al. (2015); Durrleman et al. (2013); Freifeld et al. (2015); Zhang & Fletcher (2015) as well as the statistical "pattern theory" approach developed in Dupuis et al. (1998); Grenander & Grenander (1993).

### 2.2. Spatial Transformer

The transformer operator, denoted by $\mathcal{T}$, allows for non-rigid image transformations. It is based on the composition of two mappings; a deformation map and a sampling function. The deformation maps a uniform grid of 2-dimensionak coordinates to provide its transformed version. The sampling function samples the signal with respect to a given grid of 2-dimensional coordinates.

The mapping we select to enable the learnability of the transformation in the coordinate space of the 2-dimensional signal is the Thin-Plate-Spline (TPS) interpolation technique (Bookstein, 1989; Duchon, 1976; Nejati et al., 2010) which produces smooth surfaces from $\mathbb{R}^2$ to $\mathbb{R}^2$ (Morse et al., 2005). We refer the reader to Appendix E for details regarding this method. We consider as learnable parameters of the TPS a set of 2-dimensional coordinates, called landmarks, and denoted by $\nu$. Given a set of landmarks, the TPS provides the transformation map of a 2-dimensional grid. That is, the euclidean plane is bent according to the learned landmarks.

In Fig. 1, we show on the bottom right the grid associated with the $\ell = 6^2$ landmarks. Each grid corresponds to the

spatial transformation applied to the hand-written digit 4. The transformation of the signal based on these new coordinates is produced by performing bilinear interpolation using the original signal (top left) and the new coordinates; the details are provided in Appendix E.

The spatial transformer is the composition of these two maps and is defined as

$$\mathcal{T}(x, \nu), \qquad (1)$$

where $x \in \mathbb{R}^n$ is the original 2-dimensional signal, $\nu \in \mathbb{R}^{2\ell}$ is the set of 2-dimensional transformed coordinate to be learned. Note that $2\ell$ can be smaller than the dimension of the image as the TPS interpolates to re-scale the transformation to any size.

Such a framework composing the TPS and bilinear interpolation has been defined as spatial transformer in Jaderberg et al. (2015). However, in their work, the inference of the non-rigid transformations is performed using each datum as the input of a "localisation network"; instead, we directly learn the transformation parameters.

## 3. Spatial Transformer $K$-means

We now introduce the spatial transformer $K$-means, ST $K$-means, our proposed solution that composes the spatial transformer and the $K$-means algorithm.

### 3.1. Formalism

We recall that in this work we will consider 2-dimensional signals defined by their width and height, such as images and time-frequency representation of time-series. Given a set of 2-dimensional signals $,\{x_i\}_{i=1}^N$, with $x_i \in \mathbb{R}^n$, the $K$-means algorithm aims at grouping the data into $K$ distinct clusters defining the partition $\mathcal{C} = \{C_k\}_{k=1}^K$, with $\cup_k C_k = \{x_i\}_{i=1}^N$ and $C_i \cap C_j = \emptyset, \forall i \neq j$. Each cluster $C_k$ of the partition is represented by a centroid $\mu_k \in \mathbb{R}^n, \forall k \in \{1, \ldots, K\}$.

As for the $K$-means algorithm, the goal of the ST $K$-means is to find centroids minimizing the following distortion error

$$\min_{\mathcal{C}, \mu_1, \ldots, \mu_K} \sum_{k=1}^K \sum_{i:x_i \in C_k} d(x_i, \mu_k) . \qquad (2)$$

The assignment of a signal $x_i$ to a cluster $C_k$ is achieved through the evaluation of the similarity measure, $d$, between the signal and each centroid. A signal $x_i$ belongs to the cluster $C_l$ if and only if $l = \arg\min_k d(x_i, \mu_k)$. While the standard $K$-means algorithm makes use of the Euclidean distance, i.e., $d(x_i, \mu_k) = \|x_i - \mu_k\|_2^2$, we instead propose to use the following deformation invariant similarity measure

$$d(x_i, \mu_k) := \min_{\nu \in \mathbb{R}^{2\ell}} \|\mathcal{T}(x_i, \nu) - \mu_k\|_2^2 , \qquad (3)$$

which is a *Quasipseudosemimetric*, see Appendix A for details and proof.

This similarity measure represents the least-square distance between the centroids and the datum that has been fit to the centroid via the spatial transformer operator. Once this fitting is done for each centroid, the cluster assignment is done based on the argmin of those distances, i.e., the data $x_i$ is assigned to $\arg\min_k d(x_i, \mu_k)$. Therefore, the underlying assumption of our approach is that the distance between the optimal transformation of a signal into a centroid belonging to the same "class" should be smaller than the distance between its optimal transformation into a centroid that does not. That is, let $x_i$ be geometrically near $\mu_k$, then $\min_{\nu \in \mathbb{R}^{2\ell}} \|\mathcal{T}(x_i, \nu) - \mu_k\|_2^2 < \min_{\nu \in \mathbb{R}^{2\ell}} \|\mathcal{T}(x_i, \nu) - \mu'_k\|_2^2$.

This measure requires solving a non-convex optimization problem. It can be achieved in practice by exploiting the spatial transformer's differentiability with respect to the landmarks $\nu$. As a result, we can learn the transformation by performing gradient-descent based optimization (Kingma & Ba, 2014); further details regarding this optimization are given in Appendix B as well as solutions to facilitate the optimization of the non-convex objective by exploiting the manifold geometry.

The crucial property of the measure we propose is its invariance to deformations that are spanned by the spatial transformer; formal proofs and definitions are proposed in Appendix A.3. This means that evaluating Eq. 3 with any datum that is transformed from the spatial transformer will produce the same value, as long as no information is lost.

### 3.2. Learning the Spatial Transformer $K$-means

Solving the optimization problem in Eq. 2, similarly to $K$-means, is an NP-hard problem. A popular tractable solution nonetheless exists and is known as the two-step Lloyd algorithm (Lloyd, 1982).

In the ST $K$-means, the first step of the Lloyd algorithm consists of assigning the data to the clusters using the newly defined measure of similarity in Eq. 3 . The second step is the update of the centroids using the previously determined cluster assignment. It corresponds to the result of the optimization problem: $\arg\min_{\mu_k} \sum_{i:x_i \in C_k} d(x_i, \mu_k)$, provided in following Proposition 1.

**Proposition 1.** *The centroids update of the ST $K$-means algorithm are given by*

$$\mu_k^\star := \frac{1}{|C_k|} \sum_{i:x_i \in C_k} \mathcal{T}(x_i, \nu_{i,k}^\star), \ \forall k \qquad (4)$$

*where $|C_k|$ denotes the cardinal of the set $C_k$, $\nu_{i,k}^\star$ is the set of parameters of the TPS that best transforms*

*the signal $x_i$ into the centroid $\mu_k$, that is, $\nu_{i,k}^\star = \arg\min_{\nu \in \mathbb{R}^{2l}} \|\mathcal{T}(x_i, \nu) - \mu_k\|_2^2$ (proof in Appendix A.2).*

The averaging in Eq. 4 is performed on the transformed version of the signals. The ST $K$-means thus considers the topology of the signal's space. A pseudo-code of the centroid update Eq. 4 is presented in Algo. 1.

---

**Algorithm 1** Centroids Updates of ST $K$-means

---

**Input:** Cluster $C_k$, TPS parameters $\left\{ \nu_{i,k}^\star \right\}_{i:x_i \in C_k}$

**Output:** Centroids update $\mu_k^\star$
  1: Initialize $\mu_k = 0$
  2: **for** $i : x_i \in C_k$ **do**
  3:    Compute $\mu_k = \mu_k + \mathcal{T}_\ell(x_i; \nu_{i,k}^\star)$,  Eq. 4
  4: $\mu_k^\star = \frac{\mu_k}{|C_k|}$

---

The ST $K$-means, which aims to minimize the distortion error Eq. 2 is done by alternating between the two steps detailed above until convergence, as summarized in Algo. 2.

---

**Algorithm 2** Spatial Transformer $K$-means

---

**Input:** Initial centroids $\mu_k$, dataset $\{x_i\}_{i=1}^N$

**Output:** Cluster partition $\{C_k\}_{k=1}^K$
  1: **repeat**
  2:    **for** $i = 1$ to $N$ **do**
  3:       **for** $k = 1$ to $K$ **do**
  4:          Compute and store $d(x_i, \mu_k)$ by solving Eq. 3
  5:          Assign $x_i$ to $C_l$ where $l = \arg\min_k d(x_i, \mu_k)$
  6:    Update the centroid $\mu_k$ using Algo. 1
  7: **until** Convergence

---

The update in Eq. 4, induced by our similarity measure, alleviates a fundamental limitation of the standard $K$-means. In fact, in the standard $K$-means, the average of the data belonging to a cluster $C_k$, $\frac{1}{|C_k|} \sum_{i:x_i \in C_k} x_i$, consists of an averaging of the signals without deforming them, which, as a result, does not account for the non-euclidean geometry of the signals (Klassen et al., 2004; Srivastava et al., 2005).

### 3.3. Convergence of the Spatial Transformer $K$-means

As we mentioned, our development is motivated by the interest in proposing a novel way to think about invariance in an unsupervised fashion while conserving the interpretability and theoretical guarantees of the $K$-means algorithm. We propose here to prove the convergence of the ST $K$-means algorithm following the generalization of clustering algorithms via the Bregman divergence as developed in Banerjee et al. (2005). In their work, they provide the class of distortion function that admits an iterative relocation scheme

where a global objective function, such as the one in Eq. 2, is progressively decreased. We, therefore, prove that Algo. 2 monotonically decreases the distortion error of the ST $K$-means in Eq. 2 which in turn implies that Algo. 2 converges to a local optimal.

**Proposition 2.** *Under the assumption that the spatial transformation optimization problem in Eq. 3, reaches a unique global minimum, the ST $K$-means algorithm described in Algo. 2 terminates in a finite number of step at a partition that is locally optimal (Proof in Appendix A).*

### 3.4. Geometrical Interpretation of the Similarity Measure

One of the great benefit of the $K$-means algorithm is the interpretability of the regions composing its partitioning. In particular, they are related to Voronoi diagrams which are well studied partitioning techniques (Aurenhammer, 1991; Aurenhammer et al., 2013). Following this framework, we propose now to highlight the regions defined by the ST $K$-means algorithm. This is achieved by analysing the following sets $\forall k \in \{1, \ldots, K\}$

$$R_k = \{x \in \mathbb{R}^n | d(x, \mu_k) \le d(x, \mu_j), \ \forall j \ne k\}, \quad (5)$$

where we recall $d(x, \mu_k) = \min_{\nu \in \mathbb{R}^{2\ell}} \|\mathcal{T}(x, \nu) - \mu_k\|_2^2$. Such a partitioning falls in the framework of a special type of Voronoi diagram.

**Proposition 3.** *The partitioning induced by the ST $K$-means corresponds to a weighted Voronoi diagram where each region's size depends on the per data spatial transformations (proof and details in Appendix A.5).*

While the Euclidean $K$-means induces a Voronoi diagram where each region is a polytope, the ST $K$-means does not impose such a constraint of its geometry. The similarity measure we propose adapts the geometry of each data to each centroid and thus induces a specific metric space for each data-centroid pair. In particular, for each data-centroid pair, the ST $K$-means has a particular metric that induces the boundary of the regions. In a more general setting, each region is defined as the orbit of the centroid with respect to the transformations induced by the spatial transformer, thus defining regions that depend on the orbit's shape instead of polytopal ones.

This geometric observation can lead to efficient initializations for the ST $K$-means (Arthur & Vassilvitskii, 2006), as well as the evaluation of its optimality (Bhattacharya et al., 2016). Besides, one can perform in depth study to understand the shape of the regions spanned by our approach to understand the fail cases of the algorithm for a particular application Har-Peled & Raichel (2014); Xia & Aïssa (2018). One can also compare the partitioning achieved in our approach with the one of DNN as in Balestriero et al. (2019) to gain more insights into both models.

## 3.5. Computational Complexity & Parameters

The time complexity of ST $K$-means is $O(NK(\ell^3 + \ell n))$. In fact, the ST $K$-means computes a TPS of computational complexity $O(\ell^3 + \ell n)$ for each sample of the $N$ samples and each of the $K$ centroids, as in Eq. 3. In practice, $\ell$ is of the order $2^6$. The number of parameters of the model is $2\ell \times N \times K$; it depends on the number of samples, clusters, and landmarks.

To speed up the computation, we $(i)$ pre-compute the matrix inverse responsible for the dominating cubic term, see Appendix E for implementation details regarding the TPS, and $(ii)$ implement ST $K$-means on GPU with SymJAX (Balestriero, 2020) where high parallelization renders the practical computation time near constant with respect to the number of landmarks as we depict in Fig. 2.
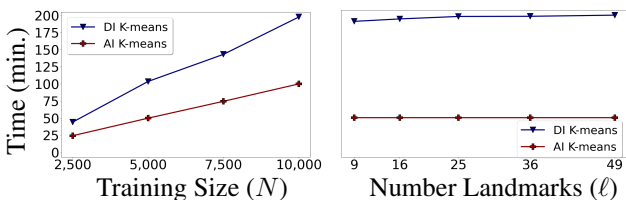


*Figure 2.* **Computational Training Time** - Comparison between our ST $K$-means and the Affine Invariant (AI) $K$-means computational times on the Arabic Characters dataset. The input pixel size is $n = 1024$. (*Left*) shows the computational time for varying training set sizes and $\ell = 7^2$. (*Right*) shows the computational time as a function of the number of landmarks, $\ell$, for $N = 10,000$. Since the AI $K$-means does not use the TPS algorithm, its computational time is constant as a function of the number of landmarks. We can observe that our process to speed up the computation enables the tractability of the ST $K$-means.

# 4. Experimental Setup

In this section, we detail the experimental settings followed to evaluate the performances of our model. For all the experiments, the number of clusters is set to be the number of classes the dataset contains for all clustering algorithms. The various datasets and their train-test split to optimize the model's parameters and update the centroids of the different models are described in Appendix F.

## 4.1. Evaluation Metrics

For all the experiments, the accuracy is calculated using the metric proposed in Yang et al. (2010) and defined as

$$\text{Accuracy} = \max_m \frac{1}{N} \sum_{i=1}^{N} 1_{\{l_i = m(\hat{l}_i)\}} \ , \qquad (6)$$

where $l_i$ is the ground-truth label, $\hat{l}_i$ the cluster assignment and $m$ all the possible one-to-one mappings between clusters and labels. The results in Table 1 are taken as the best

score on the test set based on the ground truth labels among 10 runs as in Xie et al. (2016). We also provide on the same run the normalized mutual information (NMI) (Romano et al., 2014), and adjusted rand index (ARI) (Hubert & Arabie, 1985).

## 4.2. Cross Validation Settings

We provide in Appendix. C the details regarding the benchmark models and their cross-validation settings.

Our model requires the cross-validation of hyper-parameters: the number of landmarks and the learning rate to learn the similarity measure in Eq. 3. However, the clustering framework does not allow the use of label information to perform the cross-validation of the parameters. We thus need to find a proxy for it to determine the optimal model parameters. Interestingly, the distortion error related used in the ST $K$-means, Eq. 2, appears to be negatively correlated to the accuracy, as displayed in Fig. 5. Note that the use of the distortion error is commonly used as a fitness measure in $K$-means, for example, when cross-validating the number of clusters.

We cross-validate the number of landmarks, $\ell$, which defines the resolution of the transformation, which we optimize over the following grid, $[3^2, 4^2, 5^2, 6^2, 7^2, 8^2]$. Then, the learning of the landmarks, $\nu$, is done via Adam optimizer. The learning rate is picked according to $[10^{-4}, 5 \times 10^{-4}, 10^{-3}, 5 \times 10^{-3}, 10^{-2}, 5 \times 10^{-2}]$. We train our method for 150 epochs for all the datasets, with batches of size 64. As for $K$-means and AI $K$-means, the centroids' initialization of the ST $K$-means is performed by the $K$-means++ algorithm. Importantly, the same procedure is applied to all datasets.

Note that during the training, both the similarity measure in Eq. 3 and the clustering update are performed, Eq. 6. During the algorithm's testing phase, the centroids remain fixed, and only the similarity measure is performed to assign each testing datum to a cluster.

# 5. Results and Interpretations

In this section, we report and interpret the results obtained by our ST $K$-means and competing models.

## 5.1. Clustering Accuracy

We report in Table 1 the accuracy of the different models considered on the different datasets. Our approach shows to outperform existing models on most datasets. Our model equals the performance of AI $K$-means on Affine MNIST and is only outperformed by VaDE (MLP) on MNIST.

Whereas the various deep learning approaches perform well on datasets for which their architectures were developed,
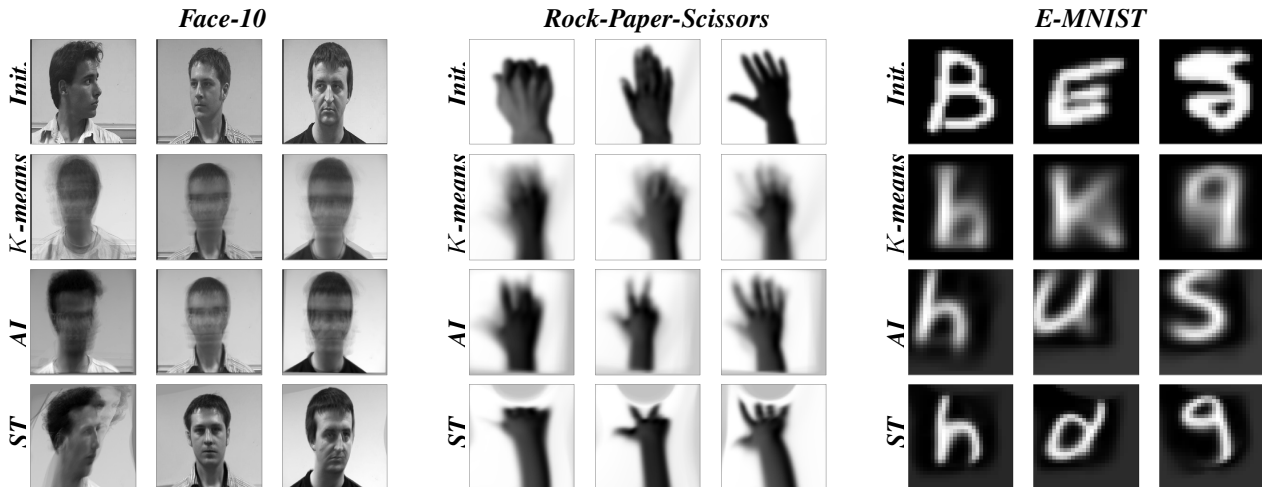
*Figure 3.* **Centroids** - We depict some centroids for the different $K$-means algorithms. The centroid at initialization are displayed in the *nth1* row. The centroids learned by $K$-means are shown in the $2^{nd}$ *row*, by the Affine invariant $K$-means in the $3^{rd}$ *row*, and by our ST $K$-means in the $4^{th}$ *row*. By comparing the results of the AI $K$-means ($3^{rd}$ *row*) with the standard $K$-means ($2^{nd}$ *row*), we can see that using only affine transformations slightly improves the $K$-means centroids and reduces the superposition issue that $K$-means suffers from. By comparing the results of our ST $K$-means ($4^{th}$ *row*) with the other methods, it is clear that using non-rigid transformations significantly improves the quality of the centroids, making them sharper and removing the issue related to the non-additiveness of images. Note that $K$-means iteratively updates the centroids and cluster assignments, as such, the class associated to a specific centroid usually changes during training (additional centroid vizualisations are proposed in Appendix G.2).

e.g., MNIST and its derivatives: E-MNIST, Arabic Characters, they show limited performance on higher resolution datasets with a small number of samples, such as Rock-Paper-Scissors, Face-10 as well as the two toy examples. In fact, they are composed of only 700 training data and 300 testing data. In the following sections, we interpret various visualizations of the $K$-means variants used in this work.

### 5.2. Interpretability: Centroids Visualization

We propose in Fig. 3 to visualize the centroids obtained via $K$-means, AI $K$-means, and our ST $K$-means. Supplementary visualizations are provided in Appendix G.2. For each dataset, the first row shows the clusters after initialization from $K$-means++. The three following rows show the centroids obtained via the $K$-means, AI $K$-means, and ST $K$-means algorithms, respectively.

We observe that, for all datasets, the $K$-means centroids are not lying on the data manifold as they are unrealistic images that could not occur naturally in the dataset. Besides, they appear to be blurry and hardly interpretable. These drawbacks are due to the update rule that consists in the average of the data belonging to each cluster in the pixel space. The AI $K$-means algorithm drastically reduces the centroids' blurriness induced by such an averaging as it considers the average of affinely transformed data. However, our ST $K$-means produces the crispest centroids and does not introduce any ambiguity in between the different clusters. In fact, the update of our method, Eq. 4, takes

into account the non-linear structure of the manifold by taking the average over data transformed using a non-rigid transformation.

Interestingly, Fig. 3 shows that even if at initialization multiple centroids assigned to the same class are attributed to different clusters, the ST $K$-means is able to recover this poor initialization thanks to its explicit manifold modeling and centroid averaging technique. For instance, in the Rock-Paper-Scissors dataset, although at initialization, two centroids correspond to the class paper, the ST $K$-means learns centroids of each of the three classes within this dataset. In the Face-10 dataset, some centroids learned correspond to the rotation of the initialization; even in such extreme change of pose, the centroids remain crisp in most cases.

### 5.3. Interpretability: Embedding Visualization

To get further insights into the disentangling capability of the ST $K$-means, we compare the 2-dimensional projections of the data using t-SNE (Maaten & Hinton, 2008), of the $K$-means, AI $K$-means and ST $K$-means. Supplementary visualizations are provided in Appendix G.1.

The t-SNE visualizations, for both the AI and ST $K$-means, are obtained by extracting the optimal transformation that led to the assignment. Precisely, for each image $x_i$, we compute $l = \arg\min_k d(x_i, \mu_k)$ and extract the optimal parameter $\nu_{i,l}^{\star}$ which is then used to obtain the transformed image fed as the input of the t-SNE.
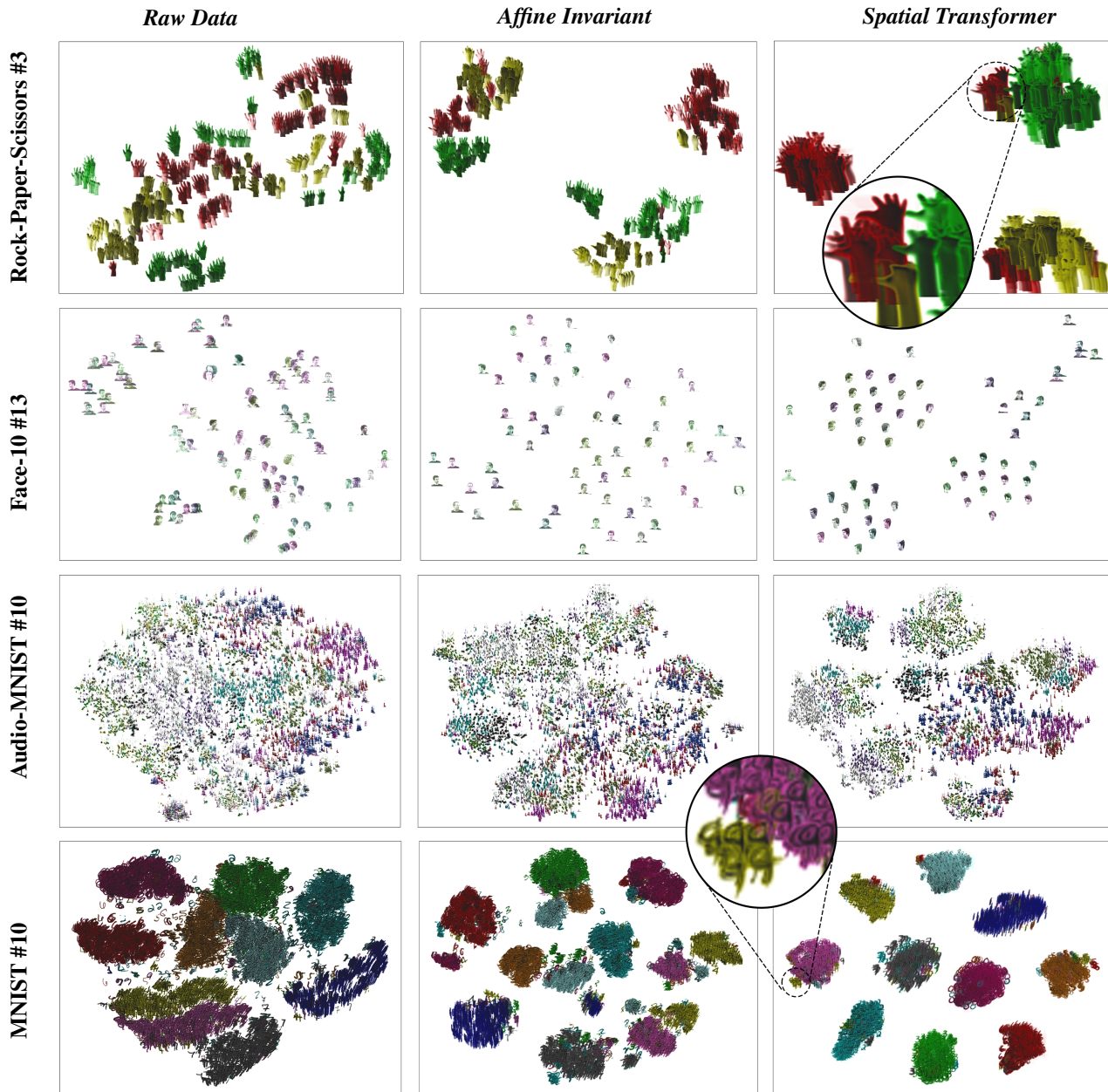
*Figure 4.* **2-dimensional t-SNE** - (# denotes the number of clusters) - We suggest the reader to zoom in the plots to best appreciate the visualizations. - The raw data (*left column*), the affinely transformed data using the AI distance, i.e., we extract the best affine transformation of the data that corresponds to the centroid it was assigned and perform the t-SNE on these affinely transformed data, (*middle column*), the data transformed with respect to the TPS as per Eq. 3, i.e., the same process as previously mentioned but we consider the spatial transformer instead, and then perform the dimension reduction on these transformed data, (*right column*). Each row corresponds to a different datasets: Rock-Paper-Scissors, Face-10, AudioMNIST, and MNIST are depicted from the top to bottom row. For all the figures, the colors of the data represent their ground truth labels. We observe that across datasets, both the affine transformations learned on the data and the non-rigid transformations help to define more localized clusters. One can observe that for the Face-10 dataset, while the dataset contains 13 clusters, we can see that the ST $K$-means induced transformations lead to a 2-dimensional space where the faces are clustered 3 majors orientations. The top left cluster corresponds to faces pointing left, the bottom one face pointing right, and the bottom right one face pointing front. We also propose to zoom-in two locations where the ambiguity in the transformation induced by the spatial transformer is noticeable. In particular, we show two cases where the non-rigid transformations are too large for certain samples leading to an erroneous clustering assignment, e.g., in the MNIST dataset, the yellow samples in the lense are initial instance of the class 4 that have been transformed into digit that geometrically ressemble the centroid of the cluster 9, thus being assigned to the 9's cluster. The same concept is shown in the Rock-Paper-Scissors lense where some instance of the classes rock and paper are assigned to the class scissors (additional t-SNE vizualisations are proposed in Appendix G.1).

Table 1. Clustering results in % of the test set accuracy Eq. 6 - Following the benchmarks evaluation method, the best accuracy (ACC) over 10 runs are displayed - We also provide the associated normalized mutual information (NMI) and adjusted rand index (ARI) - the number of clusters is denoted by # next to the dataset name and where (†): Xie et al. (2016) and (‡): Jiang et al. (2016).

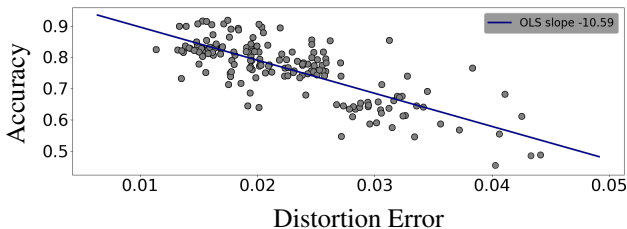| | *Deep Learning* | Aff. MNIST #10 | Diffeo. MNIST #10 | MNIST #10 | Audio MNIST #10 | E-MNIST #26 | Rock-Paper-Sci. #3 | Face-10 #13 | Arabic Char. #28 | Aff. MNIST #10 | Diffeo. MNIST #10 | MNIST #10 | Audio MNIST #10 | E-MNIST #26 | Rock-Paper-Sci. #3 | Face-10 #13 | Arabic Char. #28 | Aff. MNIST #10 | Diffeo. MNIST #10 | MNIST #10 | Audio MNIST #10 | E-MNIST #26 | Rock-Paper-Sci. #3 | Face-10 #13 | Arabic Char. #28 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | ACC | | | | | | | | NMI | | | | | | | | ARI | | | |
| $K$-means | ✗ | 68 | 61 | 53 | 10 | 39 | 40 | 20 | 19 | - | - | 50 | 1 | 39 | 5 | 18 | 27 | - | - | 39 | 0 | 21 | 4 | 0 | 1 |
| AI $K$-means | ✗ | **100** | 91 | 75 | 29 | 48 | 72 | 31 | 30 | - | - | 62 | 18 | 45 | 30 | 30 | 37 | - | - | 54 | 10 | 26 | 24 | 3 | 17 |
| **ST $K$-means** | ✗ | **100** | 99 | 92 | **41** | 65 | 86 | 45 | 51 | - | - | 82 | **26** | 63 | 63 | 53 | 61 | - | - | 83 | **15** | 46 | 63 | 20 | 38 |
| AE + $K$-means | ✓ | 72 | 60 | 66 | 13 | 41 | 48 | 37 | 23 | - | - | 64 | 1 | 40 | 9 | 27 | 33 | - | - | 59 | 0 | 28 | 6 | 26 | 15 |
| DEC (MLP) (†) | ✓ | 84 | 77 | 84 | 10 | 55 | 46 | 33 | 24 | - | - | 83 | 1 | 51 | 12 | 20 | 32 | - | - | 80 | 0 | 31 | 8 | 3 | 13 |
| DEC (Conv) | ✓ | 70 | 68 | 78 | 15 | 60 | 54 | 38 | 29 | - | - | 74 | 3 | 56 | 18 | 31 | 39 | - | - | 69 | 1 | 37 | 13 | 17 | 16 |
| VaDE (MLP) (‡) | ✓ | 68 | 65 | **94** | 11 | 20 | 50 | 36 | 26 | - | - | **89** | 1 | 12 | 16 | 27 | 30 | - | - | **85** | 0 | 8 | 11 | 14 | 10 |
| VaDE (Conv) | ✓ | 65 | 59 | 81 | 14 | 58 | 55 | 40 | 46 | - | - | 78 | 2 | 55 | 20 | 35 | 53 | - | - | 80 | 0 | 38 | 15 | 18 | 29 |



*Figure 5.* **Accuracy vs Distortion Error** - Clustering accuracy, Eq. 6, of ST $K$-means algorithm on the MNIST dataset as a function of the distortion error, Eq. 2, using the similarity measure, Eq. 3. Each gray dot is associated with a specific set of hyperparameters, e.g., the learning rate and the number of landmarks for the spatial transformer. The accuracy is negatively correlated to the distortion error (see the blue line corresponding to the ordinary least square fit), indicating that the distortion error is an appropriate metric to cross-validate the hyper-parameters of the ST $K$-means algorithm, which is crucial in an unsupervised setting as the labels are not available.

We can observe in Fig. 4 that the affine transformations ease the data separation in this 2-dimensional space. The ST $K$-means also drastically enhances the separability of the different clusters. When using ST $K$-means, the data are clustered based on macroscopically meaningful and interpretable parameters, making the model's performance possible to understand. For instance, for the Face-10 dataset, the t-SNE representation of the ST $K$-means clusters' shows that faces are grouped according to three significant orientations, left, right, and front. These three clusters are more easily observed in our ST $K$-means than in the affine invariant model. However, the 13 different orientations present

in the dataset remain too subtle to be captured by the ST $K$-means.

For the MNIST dataset, the last row and column of Fig. 4, we observe that most of the incorrectly clustered images are almost indistinguishable from samples of the cluster they have been attributed. In particular, we highlight this by proposing to zoom-in into the cluster of hand-written 9 in Fig. 4. We can see that the yellow instances are samples from the class 4 that have been transformed such that they resemble the 9's centroid in Fig. 3. We also provide a zoom-in on one of the clusters obtained on the rock-paper-scissors dataset, first row and last column of Fig. 4. The incorrectly clustered data are the ones that, when transformed, easily fit the scissors shape.

## 6. Conclusion

Designing an unsupervised algorithm that is robust to non-rigid transformations remains challenging, despite the tremendous breakthrough in machine learning. The problem lies in appropriately limiting the size of the transformations. We showed that the spatial transformer could achieve this as the number of landmarks allows the learnability of a coarse to fine grid of transformation. However, such a parameter controlling the size of the transformation should be designed as well as be learned per-cluster or per-sample. Besides this difficulty, we showed that we could conserve the interpretability of the $K$-means algorithm applied in the input data space while drastically improving its performances. Such a framework should be favored in clustering applications where the explainability of the decision is critical.

# References

Allassonnière, S., Durrleman, S., and Kuhn, E. Bayesian mixed effect atlas estimation with a diffeomorphic deformation model. *SIAM Journal on Imaging Sciences*, 8(3): 1367–1395, 2015.

Altwaijry, N. and Al-Turaiki, I. Arabic handwriting recognition system using convolutional neural network. *Neural Computing and Applications*, pp. 1–13, 2020.

Arthur, D. and Vassilvitskii, S. k-means++: The advantages of careful seeding. Technical report, 2006.

Aurenhammer, F. Voronoi diagrams—a survey of a fundamental geometric data structure. *ACM Computing Surveys (CSUR)*, 23(3):345–405, 1991.

Aurenhammer, F., Klein, R., and Lee, D. *Voronoi diagrams and Delaunay triangulations*. World Scientific Publishing Company, 2013.

Balakrishnan, G., Zhao, A., Sabuncu, M. R., Guttag, J., and Dalca, A. V. An unsupervised learning model for deformable medical image registration. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9252–9260, 2018.

Balestriero, R. Symjax: symbolic cpu/gpu/tpu programming. *arXiv preprint arXiv:2005.10635*, 2020.

Balestriero, R., Cosentino, R., Aazhang, B., and Baraniuk, R. The geometry of deep networks: Power diagram subdivision. In *Advances in Neural Information Processing Systems*, volume 32, 2019.

Banerjee, A., Merugu, S., Dhillon, I. S., Ghosh, J., and Lafferty, J. Clustering with bregman divergences. *Journal of machine learning research*, 6(10), 2005.

Basri, R., Roth, D., and Jacobs, D. Clustering appearances of 3d objects. In *Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 414–420. IEEE, 1998.

Becker, S., Ackermann, M., Lapuschkin, S., Müller, K.-R., and Samek, W. Interpreting and explaining deep neural networks for classification of audio signals. *arXiv preprint arXiv:1807.03418*, 2018.

Bertsimas, D., Orfanoudaki, A., and Wiberg, H. Interpretable clustering: an optimization approach. *Machine Learning*, pp. 1–50, 2020.

Bhattacharya, A., Jaiswal, R., and Ailon, N. Tight lower bound instances for k-means++ in two dimensions. *Theoretical Computer Science*, 634:55–66, 2016.

Bookstein, F. L. Principal warps: Thin-plate splines and the decomposition of deformations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(6):567–585, 1989.

Cohen, G., Afshar, S., Tapson, J., and Van Schaik, A. Emnist: Extending mnist to handwritten letters. In *International Joint Conference on Neural Networks*, pp. 2921–2926. IEEE, 2017.

Cosentino, R. and Aazhang, B. Learnable group transform for time-series. In *International Conference on Machine Learning*, 2020.

Dalca, A. V., Rakic, M., Guttag, J., and Sabuncu, M. R. Learning conditional deformable templates with convolutional networks. *arXiv preprint arXiv:1908.02738*, 2019.

Deng, L. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.

Detlefsen, N. S., Freifeld, O., and Hauberg, S. Deep diffeomorphic transformer networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4403–4412, 2018.

Dolnicar, S. Using cluster analysis for market segmentation-typical misconceptions, established methodological weaknesses and some recommendations for improvement. *Australasian Journal of Market Research*, 11(2):5–12, 2003.

Duchon, J. Interpolation des fonctions de deux variables suivant le principe de la flexion des plaques minces. *Revue Française d'Automatique, Informatique, Recherche Opérationnelle. Analyse Numérique*, 10(R3):5–12, 1976.

Dupuis, P., Grenander, U., and Miller, M. I. Variational problems on flows of diffeomorphisms for image matching. *Quarterly of applied mathematics*, pp. 587–600, 1998.

Durrleman, S., Allassonnière, S., and Joshi, S. Sparse adaptive parameterization of variability in image ensembles. *International Journal of Computer Vision*, 101(1):161–183, 2013.

Estivill-Castro, V. Why so many clustering algorithms: a position paper. *ACM SIGKDD explorations newsletter*, 4 (1):65–75, 2002.

Fitzgibbon, A. and Zisserman, A. On affine invariant clustering and automatic cast listing in movies. In *European Conference on Computer Vision*, pp. 304–320. Springer, 2002.

Freifeld, O., Hauberg, S., Batmanghelich, K., and Fisher, J. W. Highly-expressive spaces of well-behaved transformations: Keeping it simple. In *Proceedings of the*

*IEEE International Conference on Computer Vision*, pp. 2911–2919, 2015.

Frey, B. J. and Jojic, N. Fast, large-scale transformation-invariant clustering. In *Advances in Neural Information Processing Systems*, pp. 721–727, 2002.

Gourier, N., Hall, D., and Crowley, L. J. Estimating face orientation from robust detection of salient facial structures. *International Workshop on Visual Observation of Deictic Gestures*, 2004.

Greene, D. and Cunningham, P. Producing accurate interpretable clusters from high-dimensional data. In *European Conference on Principles of Data Mining and Knowledge Discovery*, pp. 486–494. Springer, 2005.

Grenander, U. and Grenander, E. *General Pattern Theory: A Mathematical Study of Regular Structures*. Oxford Mathematical Monographs. Clarendon Press, 1993. ISBN 9780198536710. URL https://books.google.com/books?id=Z-8YAQAAIAAJ.

Har-Peled, S. and Raichel, B. On the complexity of randomly weighted voronoi diagrams. In *Proceedings of the thirtieth annual symposium on Computational geometry*, pp. 232–241, 2014.

He, K., Wen, F., and Sun, J. K-means hashing: An affinity-preserving quantization method for learning binary compact codes. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2938–2945, 2013.

Ho, J., Yang, M.-H., Lim, J., Lee, K.-C., and Kriegman, D. Clustering appearances of objects under varying illumination conditions. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition.*, volume 1, pp. I–I. IEEE, 2003.

Hubert, L. and Arabie, P. Comparing partitions. *Journal of classification*, 2(1):193–218, 1985.

Inaba, M., Katoh, N., and Imai, H. Applications of weighted voronoi diagrams and randomization to variance-based k-clustering. In *Proceedings of the tenth annual symposium on Computational geometry*, pp. 332–339, 1994.

Jaderberg, M., Simonyan, K., Zisserman, A., and Kavukcuoglu, K. Spatial transformer networks. In *Advances in Neural Information Processing Systems*, pp. 2017–2025, 2015.

Jiang, Z., Zheng, Y., Tan, H., Tang, B., and Zhou, H. Variational deep embedding: An unsupervised and generative approach to clustering. *arXiv preprint arXiv:1611.05148*, 2016.

Johnson, H. J. and Christensen, G. E. Landmark and intensity-based, consistent thin-plate spline image registration. In *Biennial International Conference on Information Processing in Medical Imaging*, pp. 329–343. Springer, 2001.

Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Klassen, E., Srivastava, A., Mio, M., and Joshi, S. H. Analysis of planar shapes using geodesic paths on shape spaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(3):372–383, 2004.

Letscher, D. Vector weighted voronoi diagrams and delaunay triangulations.

Lim, J., Ho, J., Yang, M.-H., Lee, K.-c., and Kriegman, D. Image clustering with metric, local linear structure, and affine symmetry. In *European Conference On Computer Vision*, pp. 456–468. Springer, 2004.

Lloyd, S. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 28(2):129–137, 1982.

Lohit, S., Wang, Q., and Turaga, P. Temporal transformer networks: Joint learning of invariant and discriminative time warping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12426–12435, 2019.

Ma, Q., Zheng, J., Li, S., and Cottrell, G. W. Learning representations for time series clustering. *Advances in neural information processing systems*, 32:3781–3791, 2019.

Maaten, L. v. d. and Hinton, G. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(Nov): 2579–2605, 2008.

MacQueen, J. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, pp. 281–297, Berkeley, Calif., 1967. University of California Press.

Moroney, L. Rock, paper, scissors dataset, feb 2019. URL http://laurencemoroney.com/rock-paper-scissors-dataset.

Morse, B. S., Yoo, T. S., Rheingans, P., Chen, D. T., and Subramanian, K. R. Interpolating implicit surfaces from scattered surface data using compactly supported radial basis functions. In *ACM SIGGRAPH*, pp. 78–es. 2005.

Murase, H. and Nayar, S. K. Learning and recognition of 3d objects from appearance. In *Proceedings IEEE Workshop on Qualitative Vision*, pp. 39–50, 1993. doi: 10.1109/WQV.1993.262951.

Nam, J. E., Maurer, M., and Mueller, K. A high-dimensional feature clustering approach to support knowledge-assisted visualization. *Computers & Graphics*, 33(5):607–615, 2009.

Nejati, M., Amirfattahi, R., and Sadri, S. A fast hybrid approach for approximating a thin-plate spline surface. In *2010 18th Iranian Conference on Electrical Engineering*, pp. 204–208. IEEE, 2010.

Raytchev, B. and Murase, H. Unsupervised face recognition from image sequences based on clustering with attraction and repulsion. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition.*, volume 2, pp. II–II. IEEE, 2001.

Romano, S., Bailey, J., Nguyen, V., and Verspoor, K. Standardized mutual information for clustering comparisons: one step further in adjustment for chance. In *International Conference on Machine Learning*, pp. 1143–1151. PMLR, 2014.

Seydoux, L., Balestriero, R., Poli, P., De Hoop, M., Campillo, M., and Baraniuk, R. Clustering earthquake signals and background noises in continuous seismic data with unsupervised deep learning. *Nature communications*, 11(1):1–12, 2020.

Shapira Weber, R. A., Eyal, M., Skafte, N., Shriki, O., and Freifeld, O. Diffeomorphic temporal alignment nets. In *Advances in Neural Information Processing Systems*, volume 32, pp. 6574–6585. Curran Associates, Inc., 2019.

Simard, P. Y., LeCun, Y. A., Denker, J. S., and Victorri, B. Transformation invariance in pattern recognition–tangent distance and tangent propagation. In *Neural Networks: Tricks of the Trade*, pp. 235–269. Springer, 2012.

Srivastava, A., Joshi, S. H., Mio, W., and Liu, X. Statistical shape analysis: Clustering, learning, and testing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(4):590–602, 2005.

Steinbach, M., Ertöz, L., and Kumar, V. The challenges of clustering high dimensional data. In *New Directions in Statistical Physics*, pp. 273–309. Springer, 2004.

Su, M.-C. and Chou, C.-H. A Modified Version of the K-means Algorithm with a Distance Based on Cluster Symmetry. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(6):674–680, 2001.

Wagstaff, K., Cardie, C., Rogers, S., and Schroedl, S. Constrained k-means clustering with background knowledge. 2001.

Wakin, M. B., Donoho, D. L., Choi, H., and Baraniuk, R. G. The Multiscale Structure of Non-differentiable Image Manifolds. In *Wavelets XI*, volume 5914, pp. 59141B. International Society for Optics and Photonics, 2005.

Xia, M. and Aïssa, S. Unified analytical volume distribution of poisson-delaunay simplex and its application to coordinated multi-point transmission. *IEEE Transactions on Wireless Communications*, 17(7):4912–4921, 2018.

Xie, J., Girshick, R., and Farhadi, A. Unsupervised deep embedding for clustering analysis. In *International Conference on Machine Learning*, pp. 478–487, 2016.

Xu, R. and Wunsch, D. C. Clustering algorithms in biomedical research: a review. *IEEE Reviews in Biomedical Engineering*, 3:120–154, 2010.

Yang, Y., Xu, D., Nie, F., Yan, S., and Zhuang, Y. Image clustering using local discriminant models and global integration. *IEEE Transactions on Image Processing*, 19 (10):2761–2773, 2010.

Zhang, M. and Fletcher, P. T. Finite-dimensional lie algebras for fast diffeomorphic image registration. In *International conference on information processing in medical imaging*, pp. 249–260. Springer, 2015.

# A. Properties of ST K-means and Proofs

## A.1. ST K-means Similarity Measure: a Quasipseudosemimetric

**Proposition 4.** *The similarity measure defined by $\min_{\nu \in \mathbb{R}^{2\ell}} \|\mathcal{T}(x, \nu) - \mu\|$ is a Quasipseudosemimetric.*

*Proof.* Let's first define the orbit of an image with respect to the TPS transformations. Note that, the TPS does not form a group as it is a piecewise mapping. However, we know that it approximate any diffeomorphism on $\mathbb{R}^2$. Therefore, for sake of simplicity, we will make a slight notation abuse by considering the orbit, equivariance, and others group specific properties as being induced by the spatial transformer $\mathcal{T}$.

**Definition 1.** *We define the orbit an image $x$ under the action the $\mathcal{T}$ by*
$$\mathcal{O}(x) = \left\{ \mathcal{T}(x, \nu) | \nu \in \mathbb{R}^{2\ell} \right\}. \tag{7}$$

Let's now consider each metric statement: 1) It is non-negative as per the use of a norm.

2) **Pseudo:** $\min_{\nu \in \mathbb{R}^{2\ell}} \|(\mathcal{T}(x, \nu) - \mu\| = 0 \Leftrightarrow \exists \nu \in \mathbb{R}^{2\ell}, s.t. \ x = \mathcal{T}(x, \nu) \Leftrightarrow x \sim_{\mathcal{T}} \mu$, that is, $x$ and $\mu$ are equivariant with respect to the transformations induced by $\mathcal{T}$. Thus, $d(x, \mu) = 0$ for possibly distinct values $x$ and $\mu$, however, these are not distinct when we consider the data as any possible point on their orbit with respect to the group of diffeomorphism. In fact, the distance is equal to $0$ if and only if, $\mu$ and $x$ are equivariant.

3) **Quasi:** The asymmetry of the distance is due to the non-volume preserving deformations considered. In fact, we do not consider the Haar measure of the associated diffeomorphism group and consider the $L_2$ distance with respect to the Lebesgue measure. Although the asymmetry of $d$ does not affect our algorithm or results, a symmetric metric can be built by normalizing the distance by the determinant of the Jacobian of the transformation. Such a normalization would make the metric volume-preserving and as a result make the distance symmetric.

4) **Semi:** If $x, x', x'' \in \mathcal{O}$, then $d(x, x'') = d(x, x') = d(x', x'') = 0$ as it exist a $\nu, \nu', \nu''$ such that the TPS maps each data onto the other as per definition of the orbit, thus the triangular inequality holds. If $x, x'' \in \mathcal{O}$ and $x' \notin \mathcal{O}$, we have $d(x, x'') = 0 \leq d(x, x') + d(x', x'')$. If $x, x' \in \mathcal{O}$ and $x'' \notin \mathcal{O}$, we have $d(x', x'') = d(x, x'')$, and since $0 \leq d(x, x')$, the inequality is respected. However, if $x, x', x''$ belong to three different orbits, then we do not have the guarantee that then triangular inequality holds. In fact, it will depend on the distance between the orbits which is specific to each dataset. $\square$

## A.2. ST K-means Updates: Proof of Proposition 1

We consider the Féchet mean of the centroid $k$ to be the solution of the following optimization problem, $\arg\min_{\mu_k} \sum_{i:x_i \in C_k} d(x_i, \mu_k)$. Using our similarity measure, we obtain the following.

*Proof.* The Fréchet mean for the cluster $C_k$ is defined as $\arg\min_{\mu_k} \sum_{i:x_i \in C_k} \|\mathcal{T}(x_i, \nu^\star) - \mu_k\|^2$ since the optimization problem is convex in $\mu_k$ (as the result of the composition of the identity map and a norm which are both convex) we have $\mu_k^\star : \nabla_\mu \sum_{i:x_i \in C_k} \|\mathcal{T}(x_i, \nu^\star) - \mu_k\|^2 = 0$. with,
$$\nabla_\mu \sum_{i:x_i \in C_k} \|\mathcal{T}(x_i, \nu^\star) - \mu_k\|^2 = 2(|C_k|) \times \mu_k + 2 \sum_{i:x_i \in C_k} \mathcal{T}(x_i, \nu^\star). \tag{8}$$
$\square$

## A.3. ST K-means Similarity Measure: Invariance Property

Motivated by the fact that small non-rigid transformations, usually, do not change nature of an image, we propose to exploit the invariance property of the similarity measure we proposed.

In this section, for sake of simplicity we will assume that the transformations belong to the group of diffeomorphism. In practice, the TPS can only approximate element of such group, and the constraint we impose on the transformation, e.g., number of landmark, also limit the type of diffeomorphism that can be approximated, therefore, we could instead consider that we approximate a subgroup of the diffeomorphism group.

Let's define an invariant similarity measure under the action of such group. That is, the similarity between two 2-dimensional signals remain the same under any diffeomorphic transformations. We propose to define the invariance in the framework of

centroid-based clustering algorithm as follows.

**Definition 2.** *An invariant similarity measure with respect to diff($\mathbb{R}^2$) is defined as $d : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}^+$ such that for all images $x \in \mathbb{R}^n$, all centroids $\mu \in \mathbb{R}^2$, and all group elements $\forall g \in \text{diff}(\mathbb{R}^2)$, we have*

$$d(x, \mu) = d(g \star x, \mu), \tag{9}$$

*where $g \star x$ denotes the action of the group element $g$ onto the image $x$.*

The similarity used in Eq. 3 of the optimization problem is diff($\mathbb{R}^2$)-invariant as per Definition 2.

**Proposition 5.** *The similarity $\min_{g \in \text{diff}(\mathbb{R}^2)} \|g \star x - \mu\|$ is diff($\mathbb{R}^2$)-invariant.*

*Proof.* Let consider $g^\star = \arg\min_{g \in \text{diff}(\mathbb{R}^2)} \|g \star x - \mu\|$, we have $\arg\min_{g \in \text{diff}(\mathbb{R}^2)} \|g \cdot g' \star x - \mu\| = g^\star \cdot g'^{-1}$, where $g'^{-1}$ is the inverse group element of $g'$. In fact, $\|g^\star \cdot g'^{-1} \cdot g' \star x - \mu\| = \|g^\star \star x - \mu\|$. Since for all $g' \in \text{diff}(\mathbb{R}^2)$, it exists an inverse element $g'^{-1}$, we have that $\forall g \in \text{diff}(\mathbb{R}^2)$, $d(g' \star x, \mu) = d(x, \mu)$.

That is, by definition of the group, there is always another element that minimizes the loss function by using the composition between the inverse element of the group that has just been added, $g'$, and the optimal element $g^\star$. $\square$

## A.4. ST $K$-means Convergence: Proof of Proposition 2

*Proof.* Following the notation of Sec. E, we can define the spatial transformer operator $\mathcal{T}$ as the composition of the TPS and bilinear interpolation map. That is, $\mathcal{T}(x, \nu) = \Gamma[F(\nu), x]$. Now the aim is to prove that $\min_{\nu \in \mathcal{R}^{2l}} \|\mathcal{T}(x, \nu) - \mu\|_2^2$ defines a Bregman divergence measure as in (Banerjee et al., 2005). In such a case, Algo. 2 defines a special case of the Bregman divergence hard-clustering algorithm again defined in (Banerjee et al., 2005) which is proven to converge.

Let's first start by making an assumption on the data $x$, we can without loss of generality assume that they are non-negative as we are dealing either with images or time-frequency representation where a modulus is applied to obtain the 2-dimensional real representation. Then, we also assume that the minimum over the transformation parameters $\nu$ reaches a global unique minimum, denoted by $\nu^\star$. Now,

$$\|\mathcal{T}(x, \nu^\star) - \mu\|_2^2 = \langle \mathcal{T}(x, \nu^\star), \mathcal{T}(x, \nu) \rangle + \langle \mu, \mu \rangle - 2\langle \mathcal{T}(x, \nu^\star), \mu \rangle$$
$$= \langle \mathcal{T}(x, \nu^\star), \mathcal{T}(x, \nu^\star) \rangle - \langle \mu, \mu \rangle - \langle \mathcal{T}(x, \nu^\star) - \mu, 2\mu \rangle \ ,$$

Now it is clear that $\mu = \mathcal{T}(\mu, 0)$ which consists in the identity transform of the centroid $\mu$. Then we denote by $\phi_{\nu^\star}(y) = \langle \mathcal{T}(y, \nu^\star), \mathcal{T}(y, \nu^\star) \rangle$ where $y \in \mathbb{R}^n$, and obtain that,

$$\|\mathcal{T}(x, \nu^\star) - \mu\|_2^2 = \phi_{\nu^\star}(x) - \phi_0(\mu) - \langle \mathcal{T}(x, \nu^\star) - \mathcal{T}(\mu, 0), \nabla\phi_0(\mu) \rangle \ .$$

Now, we know that $\langle x, x \rangle$ is non-decreasing w.r.t each dimension since the image or time-frequency representation are positive real valued, and the inner product defines a strictly convex map. Then, we also know that $\mathcal{T}(x, \nu^\star) = \Gamma[F(\nu^\star), x]$ is defined as the composition of the TPS for the coordinate and the bilinear map for the image, which can be formulated as a linear transformation with respect to the data $x : Ax$, where $A$ is a structured sparse matrix where each block denotes the dependency to nearby pixels. Therefore this mapping is convex. As a composition between a non-decreasing w.r.t each dimension and strictly convex function with a convex function, $\phi_{\nu^\star}$ is strictly convex, which complete the proof.

$\square$

## A.5. ST $K$-means: Weighted Voronoi Diagram

*Proof.* Let's start by re-writting the similarity measure as to analytically express a metric tensor that would be the weight in the weighted Voronoi diagram the ST $K$-means defines. Using App. E, we can re-write $d(x, \mu_k) = \min_{\nu \in \mathbb{R}^{2\ell}} \|\mathcal{T}(x, \nu) - \mu_k\|_2^2 = \min_{\nu \in \mathbb{R}^{2\ell}} \|A(\nu)x - \mu_k\|_2^2$, where $A(\nu)$ is bilinear in the coordinates that are induced by the TPS. In this formulation we can observe that $\nu$ defines the displacement vector w.r.t the original uniform grid of landmark. That is, if $\nu$ is the null vector, then $A(\nu)x = x$. Now, we assume that such linear operator is inversible, i.e., the TPS transformation is invertible (note that this is not always the case (Johnson & Christensen, 2001)). Then, we can re-write $d(x, \mu_k)$ as

$$\min_{\nu \in \mathbb{R}^{2\ell}} \left\|x - A(\nu_{x,k})^{-1}\mu_k\right\|_{A(\nu_{x,k})^T A(\nu_{x,k})}^2 \ ,$$

where $\|x\|_{A(\nu_{x,k})^T A(\nu_{x,k})} = x^T A(\nu_{x,k})^T A(\nu_{x,k}) x$, and $A(\nu_{x,k})^T A(\nu_{x,k})$ defines the metric tensor, and the notation $\nu_{x,k}$ indicates that the displacement vector $\nu$ depends on the centroid $\mu_k$ and the datum $x$. Also, note that while $A(\nu_{x,k})$ defines

the transformation operator to map $x$ onto $\mu_k$, $A(\nu_{x,k})^{-1}$ is the inverse operator mapping the centroid $\mu_k$ to the datum $x$.

Now the tuple of cells $\{R_k\}_{k=1}^{K}$ defines such as

$$R_k = \left\{ x \in \mathbb{R}^n \mid \left\| x - A(\nu_{x,k})^{-1} \mu_k \right\|_{A(\nu_{x,k})^T A(\nu_{x,k})} \leq \left\| x - A(\nu_{x,j})^{-1} \mu_j \right\|_{A(\nu_{x,j})^T A(\nu_{x,j})}, \ \forall j \neq k \right\},$$

defines a weighted Voronoi diagram (Inaba et al., 1994; Letscher), where we observe that the metric tensor is dependant on all the spatial transformations. $\square$

## B. Implementation Details

Note that the deformation invariant similarity measure we introduced in Eq. 3 differs from the affine invariant distances developed in (Fitzgibbon & Zisserman, 2002; Lim et al., 2004; Simard et al., 2012). All previously defined measures of error rely on the assumption that the manifold can be locally linearized and as a result the tangent space is used as a proxy to learn the optimal affine transformation. However, the work of (Wakin et al., 2005) suggests that tangent planes fitted to image manifold continually twist off into new dimensions as the parameters of the affine transformations vary due to a possible the intrinsic multiscale structure of the manifold. As such, the alignment of two images can be done by linearizing the manifold. To do so, (Wakin et al., 2005) propose to consider the multiscale structure of the manifold, we simplify their approach by applying a low-pass filter on the images and the centroid prior to learn the affine transformation best aligning them. Then, we optimize the remaining part of the TPS to account for diffeomorphic transformations. These two steps are similar to the one used in (Jaderberg et al., 2015).

## C. Alternative Methods

We compare our model with well-known clustering techniques using deep neural networks. We performed experiments for the VaDE (Jiang et al., 2016) and DEC (Xie et al., 2016) using the code made publicly available by the authors. We use the annotation (MLP) as a reference to the MLP architecture used in their experiments (see Appendix D for details). To fairly compare our model to the DEC and VaDE models, we proposed a convolutional architecture to the DEC and VaDE networks, denoted by DEC (Conv) and VaDE (Conv) (see Appendix D for details). Finally, we evaluate the performance of an augmented $K$-means algorithm trained using the features extracted by an Autoencoder, denoted by AE $+ K$-means in the following.

The parameters of the different models mentioned above are learned by stochastic gradient descent (Adam optimizer (Kingma & Ba, 2014)). In all the experiments, the learning rate are cross-validated following the approach in (Xie et al., 2016) according to $[10^{-4}, 5 \times 10^{-4}, 10^{-3}, 5 \times 10^{-3}, 10^{-2}, 5 \times 10^{-2}]$. The internal parameters that are model dependent, e.g., the number of pre-training epoch and the update intervals, are also cross-validated.

We also compare our ST $K$-means to the closely related $K$-means and affine invariant $K$-means, denoted by AI $K$-means. For each run, all three $K$-means algorithms start from the same initial centroids using the $K$-means++ algorithm developed by Arthur & Vassilvitskii (2006) to speed up the convergence of the $K$-means algorithm.

## D. Neural Network Architectures

For both architectures , the decoder architecture is symmetric to the encoder and the batch size is set to $64$.

**MLP:** The MLP architecture from input data to bottleneck hidden layer is composed of $4$ fully connected ReLU layers with dimensions $[500, 500, 2000, 10]$.

**Conv:** The CONV architecture is composed of $3$ $2d$-convolutional ReLU layers with $32$ filters of size $5 \times 5$, and $2$ fully connected ReLU layers with dimension $[400, 10]$. For each layer, a batch normalization is applied.

## E. Thin-Plate-Spline Interpolation

Let's consider two set of landmarks, the source ones $\nu_s = \{u_i, v_i\}_{i=1}^{\ell}$ and the transformed $\nu_t = \{u_i', v_i'\}_{i=1}^{\ell}$ where $\ell$ denotes the number of landmarks. The TPS aim at finding a mapping $F = (F_1, F_2)$, such that $F(u, v) = (F_1(u, v), F_2(u, v)) = (u', v')$, that is, the mapping between two set of landmarks. The particularity of the TPS is that it learns such a mapping by

minimizing the interpolation term, and a regularization that consists in penalizing the bending energy.

The TPS optimization problem is defined by

$$\min_F \sum_{i=1}^N \|(u_i', v_i') - F(u_i, v_i)\|^2 + \lambda \int \int \left[ (\frac{\partial^2 F}{\partial u^2})^2 + 2(\frac{\partial^2 F}{\partial u \partial v})^2 + (\frac{\partial^2 F}{\partial v^2})^2 \right] dudv. \tag{10}$$

In our model, the source landmarks are considered to be the coordinates of a uniform grid. Also note that both the source landmarks and transformed ones are usually a subset of the set of coordinate of the images. For instance, for the MNIST dataset of size $28 \times 28$, the landmarks would be a grid of size $\ell \times \ell$, where $\ell < 28$. While the mapping is based on the landmark, it is then applied to the entire image coordinate. In fact, $F = (F_1, F_2)$ is mapping $\mathbb{R}^2 \to \mathbb{R}^2$, where $F_1$ (resp. $F_2$) corresponds to the mapping from $(x, y)$ to the first dimension $x'$ (resp. the second dimension $y'$).

The solution of the TPS optimization problem, Eq. 10, provides the following analytical formula for $F$

$$F_1(u, v) = u' = a_1^{(1)} + a_u^{(1)} u + a_v^{(1)} v + \sum_{i=1}^\ell w_i^{(u)} U(|(u_i, v_i) - (u, v)|), \tag{11}$$

$$F_2(u, v) = v' = a_1^{(2)} + a_u^{(2)} u + a_v^{(2)} v + \sum_{i=1}^\ell w_i^{(v)} U(|(u_i, v_i) - (u, v)|), \tag{12}$$

where $|.|$ is the $L_1$-norm, $a_1, a_u, a_v$ are the parameters governing the affine transformation, and $w_i$ are parameters responsible for non-rigid transformations as they stand as a weight of the non-linear kernel $U$. The non-linear kernel $U$ is expressed by $U(r) = r^2 \log(r^2), \forall r \in \mathbb{R}_+$.

Based on the landmarks $\nu_s$ and $\nu_t$, we can obtain these parameters by solving a simple system of equation define by the following operations

$$\mathcal{L}^{-1} \mathcal{V} = \begin{bmatrix} (W^{(x)} | a_1^{(x)} a_x^{(x)} a_y^{(x)})^T \\ (W^{(x)} | a_1^{(y)} a_x^{(y)} a_y^{(y)})^T \end{bmatrix}. \tag{13}$$

where the matrix $\mathcal{L} \in \mathbb{R}^{(\ell+3) \times (\ell+3)}$, is defined as

$$\mathcal{L} = \left[ \begin{array}{c|c} \mathcal{K} & \mathcal{P} \\ \hline \mathcal{P}^T & \mathcal{O} \end{array} \right], \mathcal{K} = \begin{bmatrix} 0 & U(r_{12}) & \ldots & U(r_{1\ell}) \\ U(r_{21}) & 0 & \ldots & U(r_{2\ell}) \\ \ldots & \ldots & \ldots & \ldots \\ U(r_{\ell 1}) & \ldots & \ldots & 0 \end{bmatrix}, \mathcal{P} = \begin{bmatrix} 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \\ \ldots & \ldots & \ldots \\ 1 & x_\ell & y_\ell \end{bmatrix}$$

where $r_{ij} = |(u_i, v_i) - (u_j, v_j)|$, $\mathcal{K} \in \mathbb{R}_+^{\ell \times \ell}$, and $\mathcal{V} = \begin{bmatrix} x_1' & x_2' & \ldots & x_\ell' | 0 & 0 & 0 \\ y_1' & y_2' & \ldots & y_\ell' | 0 & 0 & 0 \end{bmatrix}$.

Note that, since the matrix $\mathcal{L}$ depends only on the source landmarks, and that in our case these are unchanged, its inverse can be computed only once. The only operation required to be computed for each data and each centroid is the matrix multiplication $\mathcal{L}^{-1} \mathcal{V}$ providing the parameters of the TPS transformation, as per Eq. 11, 12. Given these parameters, the mapping $F$ can be applied to to each coordinate of the image.

Now in order to render the image, one can perform bilinear interpolation as it is achieved in. Besides, the bilinear interpolation will allow the propagation of the gradient through any differentiable loss function.

Given an image $x_1 \in \mathbb{R}^n$ where $n = W \times H$, $W$ denotes the width and $H$ the height of the image, and two sets of landmarks $\nu_s = \{u_i, v_i\}_{i=1}^\ell$ ,uniform grid coordinate of $x_1$, and $\nu_t = \{u_i', v_i'\}_{i=1}^\ell$, the transformation of the uniform grid, which are a subset of the image coordinate, We are able to learn a mapping $F = (F_1, F_2)$ such that for each original pixel coordinate, we have their transformed coordinates. In fact, given any position $(u, v)$ on the original image, the mapping $F$ provides the new positions $(u', v')$ as per Eq. 11, Eq. 12.

Now, from this transformed the coordinates space, we can render an image $x_2 \in \mathbb{R}^n$ using, as in (Jaderberg et al., 2015), the bilinear interpolation function $\Gamma : \mathbb{R}^2 \times \mathbb{R}^n \to \mathbb{R}$ which takes as input the original image $x_1$ and the transformed pixel

coordinates $(u', v')$, and outputs the pixel value of the transformed image at a given pixel coordinate

$$
\begin{aligned}
x_2(k,l) &= \Gamma[F(u_k, v_l), x_1] \\
&= \Gamma[(u'_k, v'_l), x_1] \\
&= \sum_{t,h \in \{0,1\}} \sum_{i=1}^{W} \sum_{j=1}^{H} x_1(i,j)\delta(\lfloor u'_k + t \rfloor - i) \times \delta(\lfloor v'_l + h \rfloor - j)(u'_k - \lfloor u'_k \rfloor)^{\delta(t)}(v'_l - \lfloor v'_l \rfloor)^{\delta(h)} \\
&\qquad \times (1 - (v'_l - \lfloor v'_l \rfloor))^{\delta(t-1)}(1 - (u'_k - \lfloor u'_k \rfloor))^{\delta(h-1)},
\end{aligned}
$$

where $\delta$ is the Kronecker delta function and $\lfloor . \rfloor$ is the floor function rounding the real coordinate to the closest pixel coordinate.

## F. Datasets

**MNIST (Deng, 2012):** is a handwritten digit dataset containing $60.000$ training and $10.000$ test images of dimension $28 \times 28$ representing 10 classes.

**Rigid MNIST:** we randomly sample one instance of each MNIST class and generate 100 random affine transformations for each sample. A third of the data are used for testing.

**Non-rigid MNIST:** we randomly sample one instance of each MNIST class and generate 100 random affine transformations for each sample as well as random transformations using the TPS method. A third of the data are used for testing.

**Audio MNIST (Becker et al., 2018):** is composed of 30000 recordings of spoken digits by 60 different speaker and sampled at $48kHz$ of 1sec long. It consists of 10 classes. We use 10000 data for testing and 20000 for testing. This dataset will be transformed into a time-frequency representation. This representations, can be considered as images, are usually used as the common representation of audio recordings (Cosentino & Aazhang, 2020).

**E-MNIST (Cohen et al., 2017):** is a handwritten letters dataset merging a balanced set of the uppercase and lowercase letters into a single 26 classes dataset of dimension $28 \times 28$.

**Rock-Paper-Scissors (Moroney, 2019):** images of hands playing rock, paper, scissor game, that is, a 3 classes dataset. The dimension of each image is $300 \times 300$. The training set is composed of 2520 data and the testing set 372.

**Face-10 (Gourier et al., 2004):** images of the face of 15 people, wearing glasses or not and having various skin color. For each individual, different samples are obtained with different pose orientation varying from $-90$ degrees to $+90$ degrees vertical degrees. The dimension of each image is $288 \times 384$. The training set is composed of 273 data and testing set of 117 data.

**Arabic Char (Altwaijry & Al-Turaiki, 2020):** Handwritten Arabic characters written by 60 participants. The dataset is composed of $13,440$ images in the training set and 3360 in the test set. The dimension of each image is $32 \times 32$.

## G. Supplementary Visualization

### G.1. Additional t-SNE Visualisations
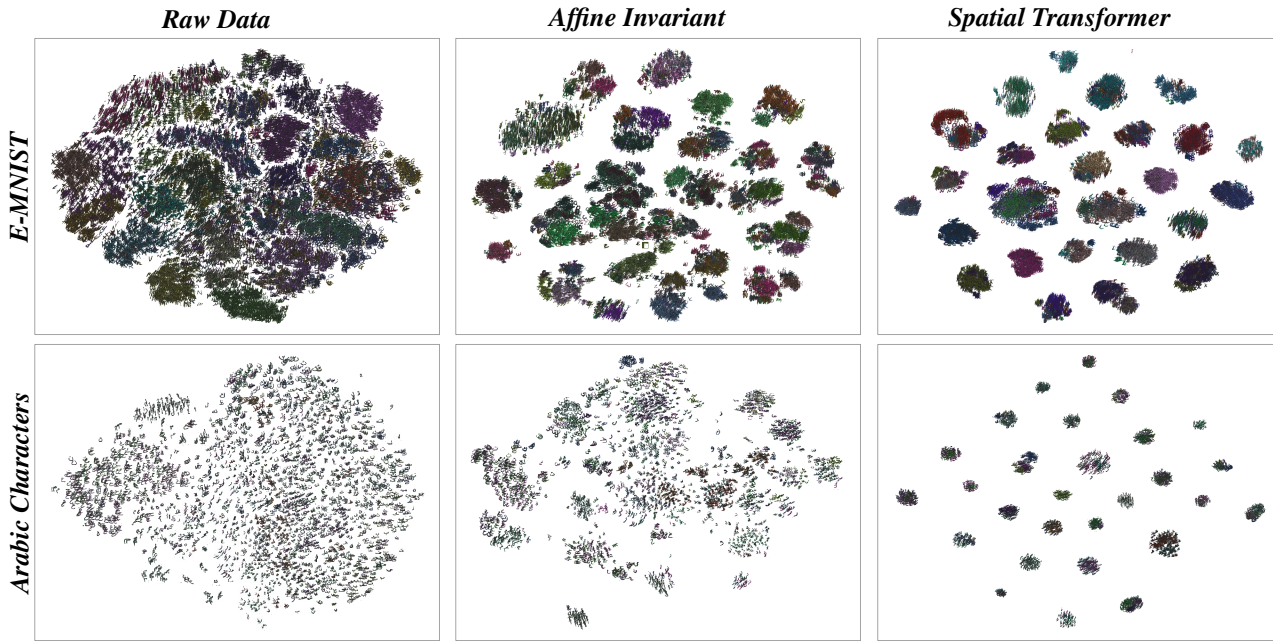
### G.2. Additional Centroid Visualisations

*Raw Data*      *Affine Invariant*      *Spatial Transformer*



Figure 6. 2-**dimensional t-SNE Vizualisation** - The raw data (*left column*), the affinely transformed data, i.e., we extract the transformation of the data that corresponds to the centroid it was assigned and perform the t-SNE on these affinely transformed data, (*middle column*), the data transformed with respect to non-rigid transformations as per Eq. 3, i.e., the same process as previously mention, but we consider the transformation induced by the TPS, and then perform the dimension reduction on these transformed data, (*right column*). Each row corresponds to a different dataset, E-MNIST, Arabic Characters, are depicted from the top to bottom row. For all the figures, the colors of the data represent their ground truth labels.
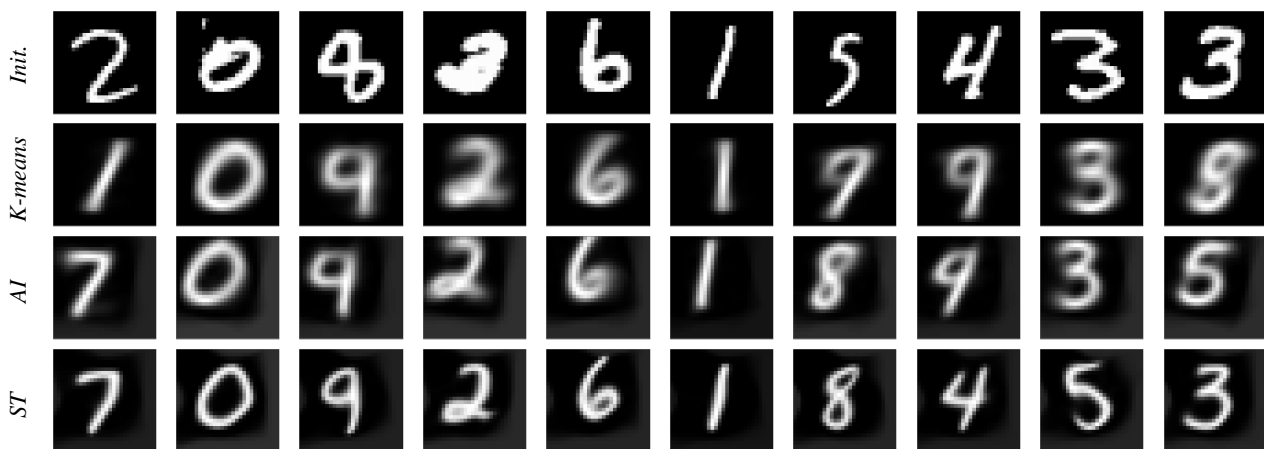


Figure 7. *Additional* - **MNIST Centroids Visualization (dim 28x28)** - Depiction of the initialization of the per-cluster centroids **top row** and the final per-cluster centroids of the K-means, Affine invariant K-means, and ST K-means (proposed) methods.
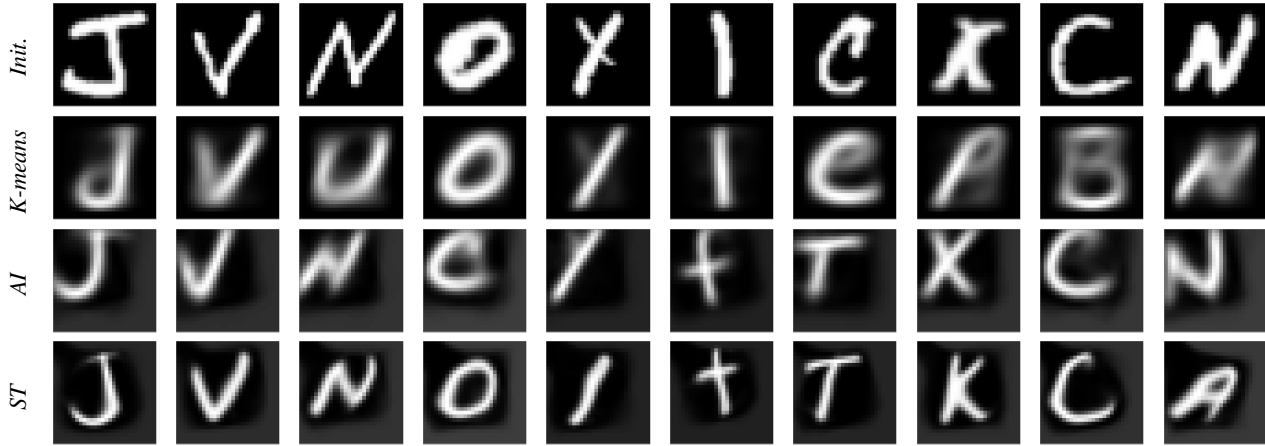
*Figure 8. Additional* - 10 **out of** 26 **E-MNIST Centroids Visualization (dim 28x28)** - Depiction of the initialization of the per-cluster centroids **top row** and the final per-cluster centroids of the K-means, AI K-means, and ST K-means (proposed) methods.
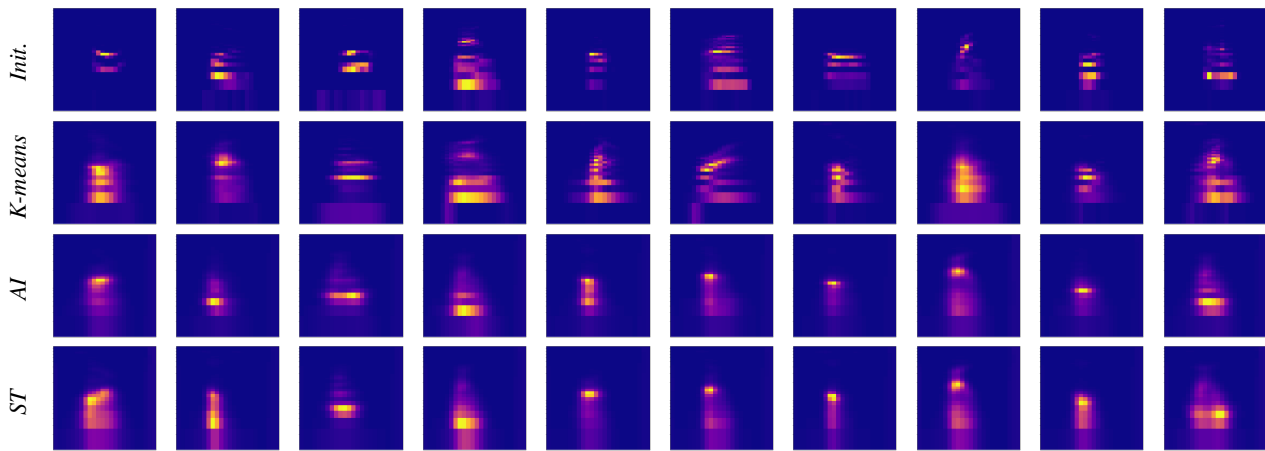


*Figure 9. Additional* - **Audio MNIST Centroids Visualization (dim 64x24)** - Depiction of the initialization of the per-cluster centroids **top row** and the final per-cluster centroids of the K-means, AI K-means, and ST K-means (proposed) methods.
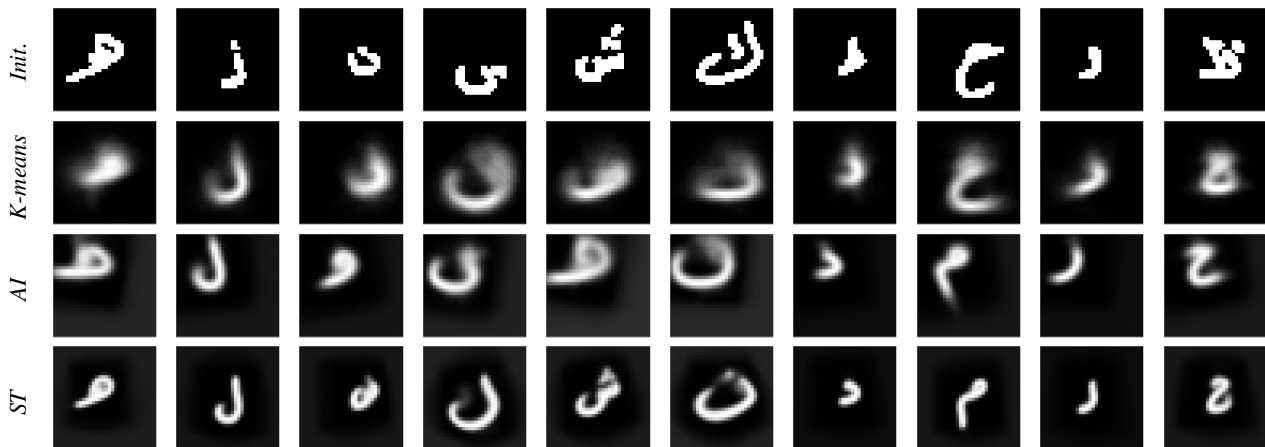


*Figure 10. Additional* - 10 **out of** 28 **Arab Characters Centroids Visualization (dim 32x32)** - Depiction of the initialization of the per-cluster centroids **top row** and the final per-cluster centroids of the K-means, AI K-means, and ST K-means (proposed) methods.
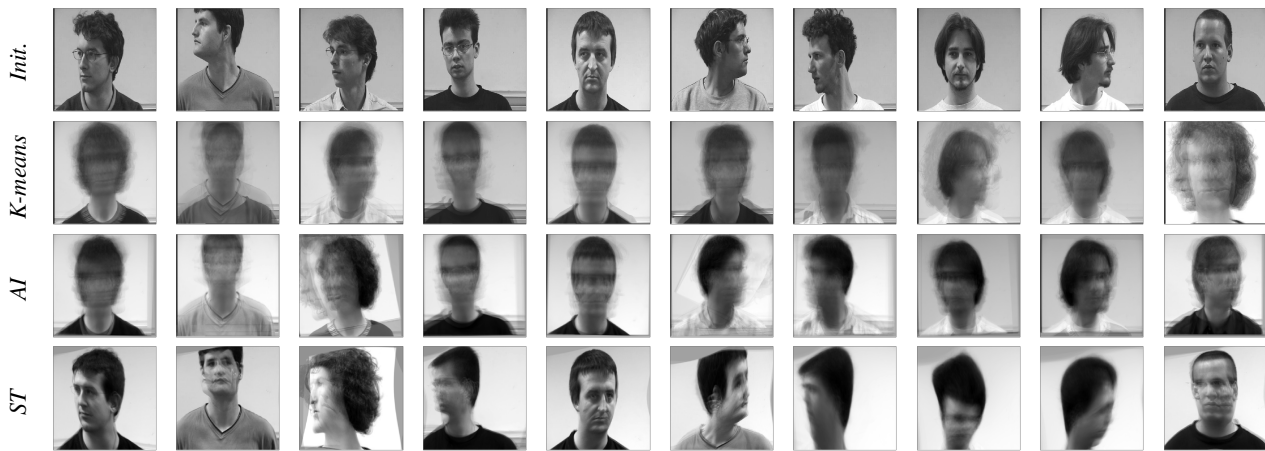
*Figure 11. Additional* - 10 **out of** 13 **Face Position Centroids Visualization (dim 288x384)** - Depiction of the initialization of the per-cluster centroids **top row** and the final per-cluster centroids of the K-means, AI K-means, and ST K-means (proposed) methods.