

Auto-encoders with fibered latent spaces: A geometric approach to batch correction

Tariq Daouda*(1, 2, 3) Reda Chhaibi*(4) Prudencio Tossou(4, 6) Alexandra-Chloé Villani(1, 2, 3)

<https://arxiv.org/abs/2005.07852>

<https://github.com/tariqdaouda/fiberedAE/>



Fibered Auto-Encoders: Supervised geometric disentanglement

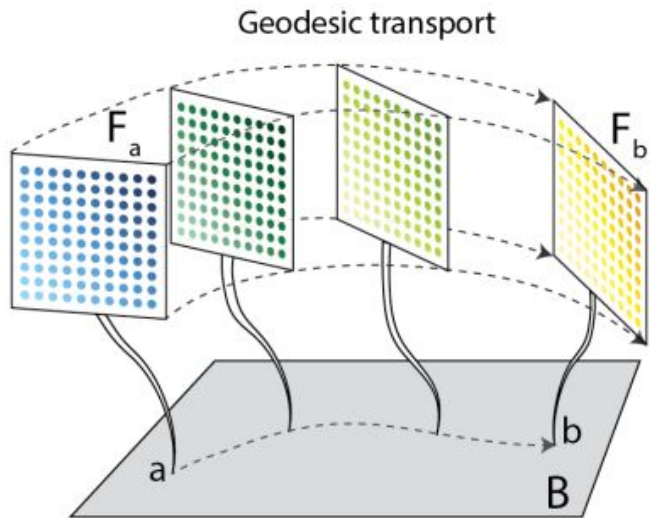


Figure 1.1. Conceptual sketch of the proposed method. We stratify the learned latent space M into a base space B and a fiber space F . Under this representation, samples can be formally transported from F_a to F_b , and geodesic interpolations between latent spaces F_a and F_b can be generated.

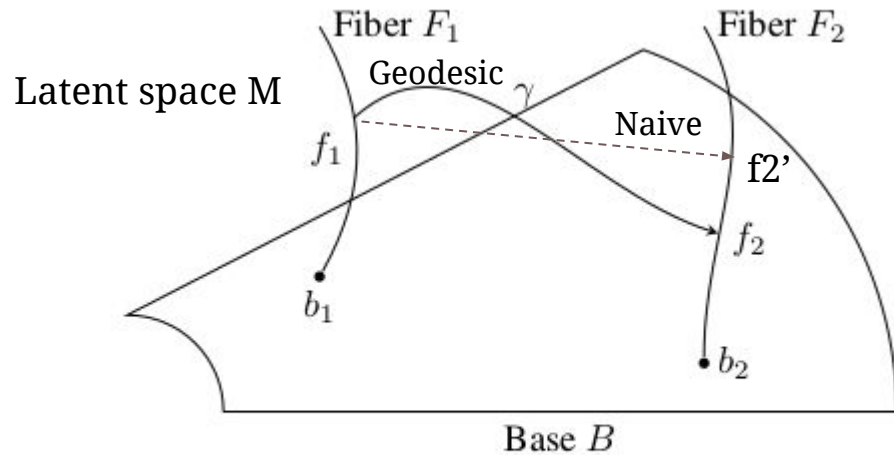


Figure 2.1. Illustration of geodesic transport with curve γ from the fiber F_1 to F_2 .

The geodesic is the shortest path w.r.t the natural metric respecting the auto-encoder formalism:

$$\mathcal{X} \xrightarrow{\text{Encoder}} M \xrightarrow{\text{Decoder}} \left(\mathcal{X}, \|\cdot\|^2 \right).$$

Fibered Auto-Encoders: Model and Training

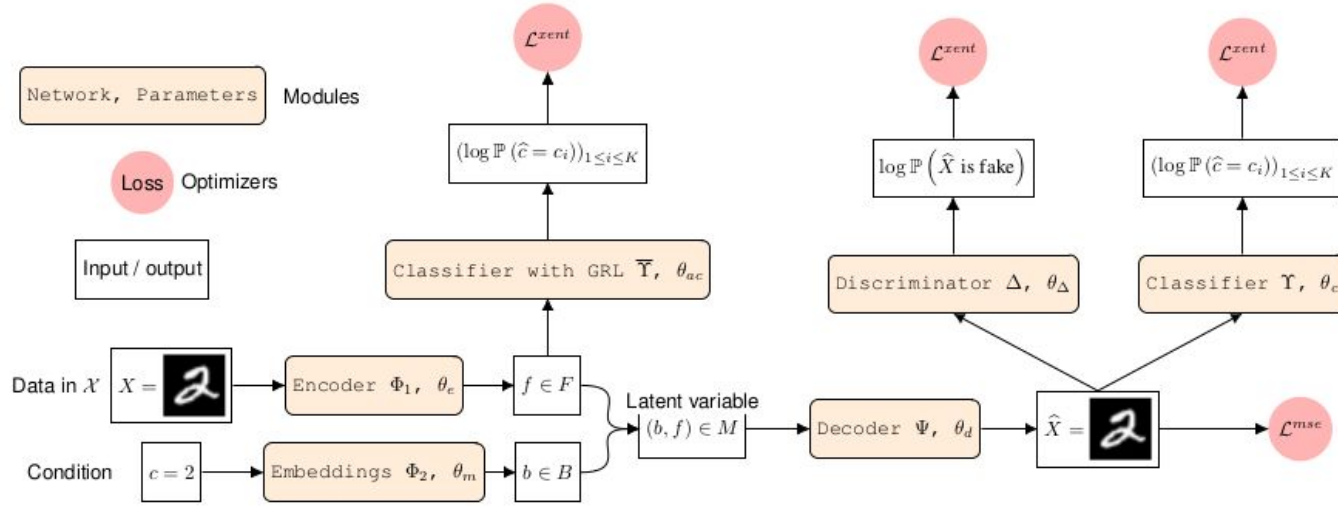


Figure 2.2. Network architecture. The general architecture is that of an auto-encoder receiving couples of samples and conditions (X, c) and outputting a reconstruction \hat{X} . The latent space is stratified into the fiber coordinate f (output of the bottleneck layer), and the base coordinate b encoding conditions. To the auto-encoder architecture we have added the classifier $\bar{\Upsilon}$ coupled with a GRL (Gradient Reversal Layer (Ganin & Lempitsky, 2014)) to disentangle f from b , the GAN discriminator Δ to ensure reconstruction realism, and the condition classifier Υ to prevent mode collapses.

Fibered Auto-Encoders: Manifold learning MNIST

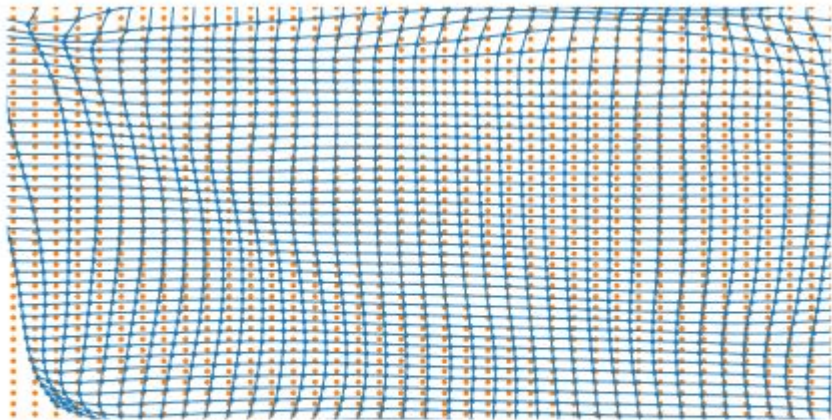


Figure 3.3. Diffeomorphism between F_4 and F_9 . The orange dots represent the original coordinates in F_4 , the blue dots are their corresponding images in F_9 computed through geodesic transport.

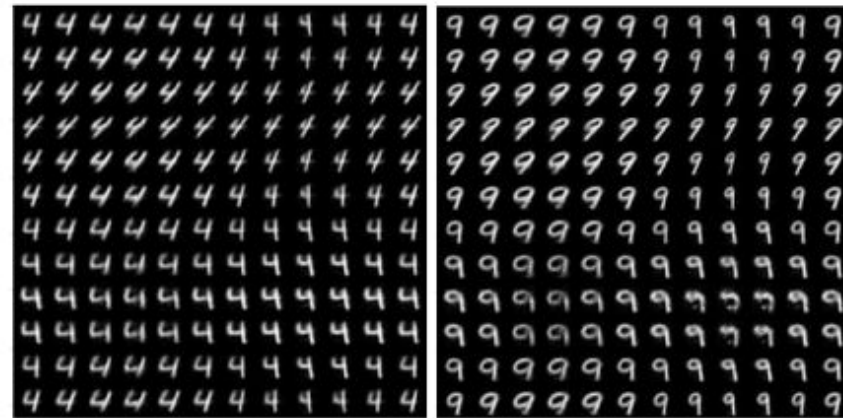


Figure 3.1. Manifold plots for the fibers F_4 and F_9 in MNIST. Images were generated using an evenly spaced grid in the standard fiber space $F = [-1, 1]^2$.

Fibered Auto-Encoders: Manifold learning MNIST (fiber space)

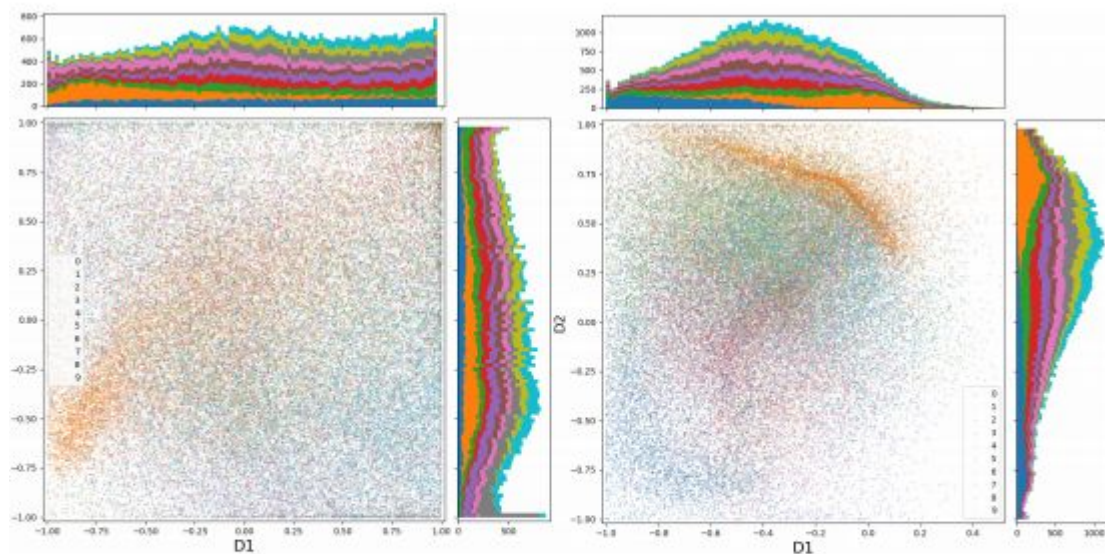


FIGURE 3.2. MNIST fiber space for all conditions. Every digit is a condition represented by a color. Fiber space is $[-1, 1]^2$. Left: network trained with condition adversarial training, right: without. Empirical distributions on fibers are closer to uniform, making them less distinguishable.

Fibered space
With and without
disentanglement
via a Gradient
Reversal Layer

Fibered Auto-Encoders: Manifold learning Olivetti (10 ex / class)

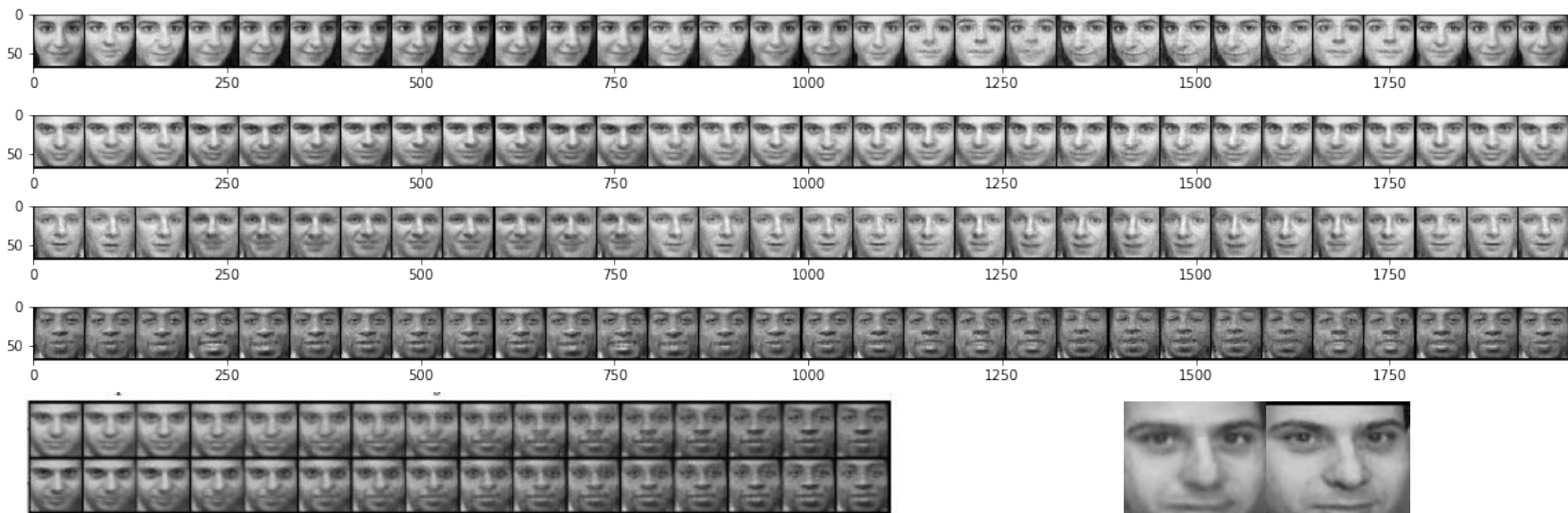


FIGURE 4.6. Geodesic interpolation from F_8 to F_{21} . Individual 8 gradually morphs into individual 21.

Smooth interpolation
100 image from 10 training examples

Fibered Auto-Encoders: Single cell integration

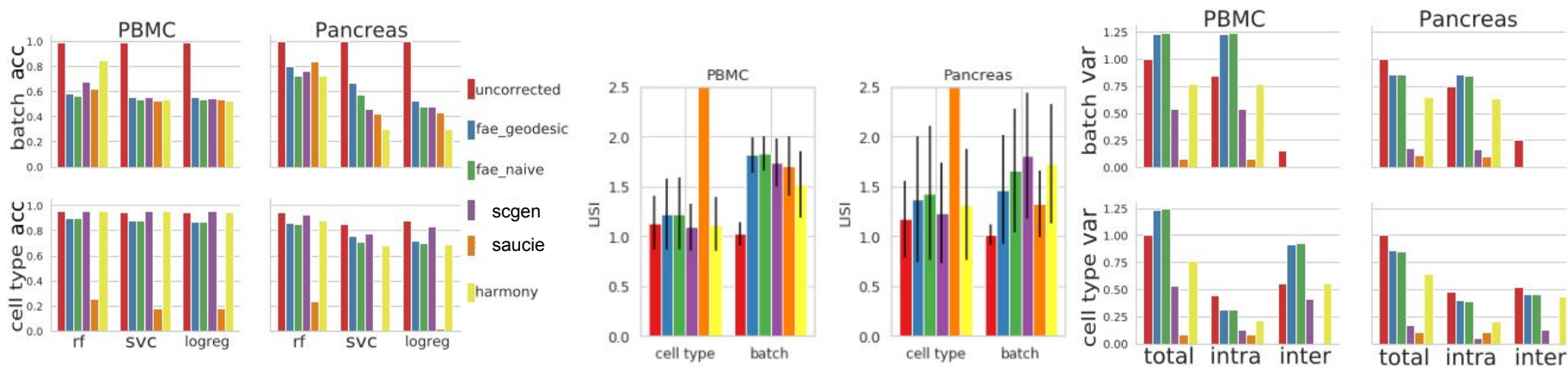


Figure 3.2. Left: Accuracies for uncorrected and batch correction methods on both dataset. Batch accuracy (lower is better) and, cell type accuracy (higher is better) are reported for random forest (rf), support vector classifier (svc) and logistic regression (logreg). Middle: LISI scores for uncorrected and batch correction methods on both datasets. LISI on cell type (closer to 1 is better), LISI on batch (higher is better). Error bars display the standard deviation. Right: Total variance and Ward's variance decomposition, for uncorrected data and batch correction methods. Total variance has been normalized to 1, on the uncorrected dataset.

Fibered Auto-Encoders: Single cell integration

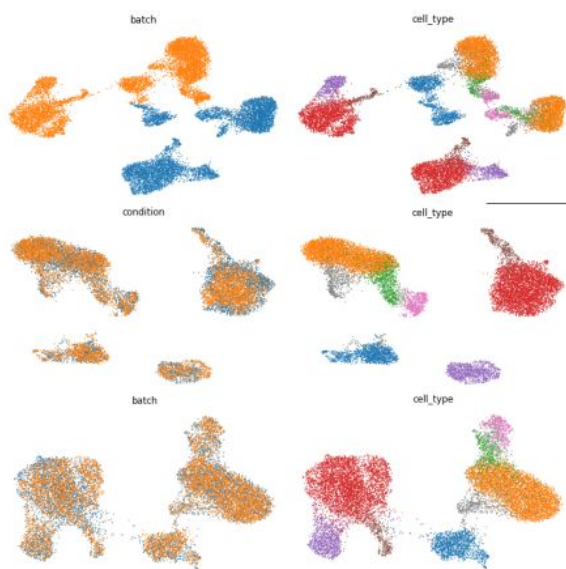


Figure 3.4. UMAP visualization of PBMC cells. Left column: cells colored by batch, Right: colored by cell types. From top to bottom: uncorrected data, scGen, geodesic transport (the plot of naive transport is very close to the naked eye). Transport conserves cell types relationships by keeping purple cells (i.e., CD16 monocytes) close to red cells (i.e., CD14 monocytes), which are related to each other.

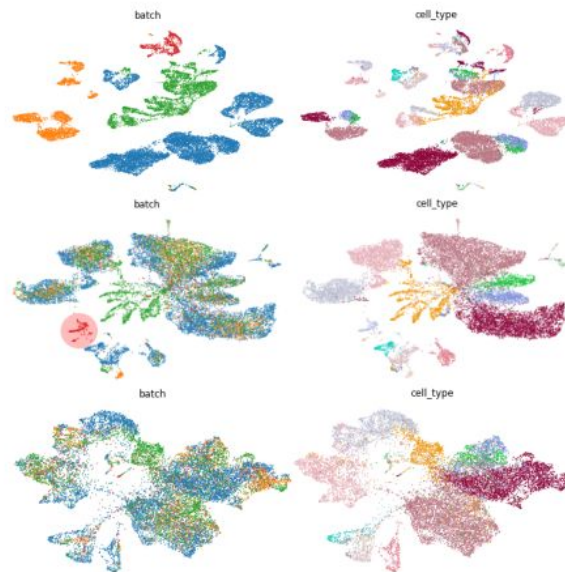
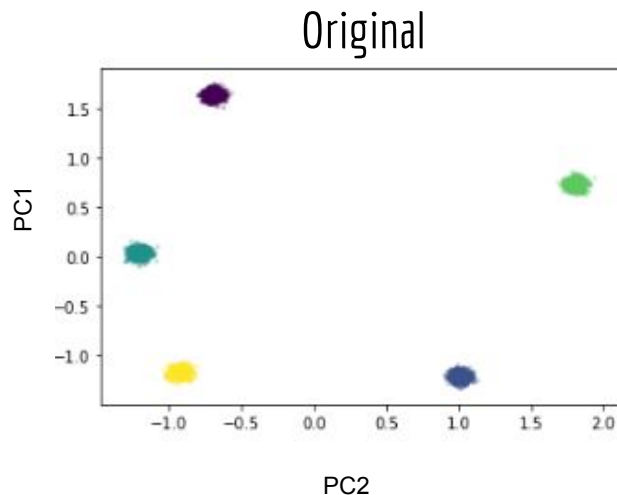
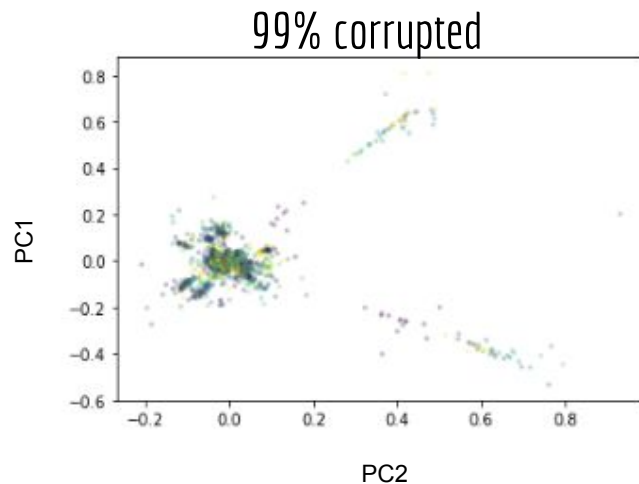


Figure 3.5. UMAP visualization of pancreas cells. Left column: cells colored by batch, right: colored by cell types. From top to bottom: uncorrected data, scGen (bottleneck size: 100), naive transport (bottleneck size: 16). Contrary to scGen, naive transport was able to integrate cells from the red batch despite the small sample size. This suggests that FAE are better at integrating datasets of small sample sizes

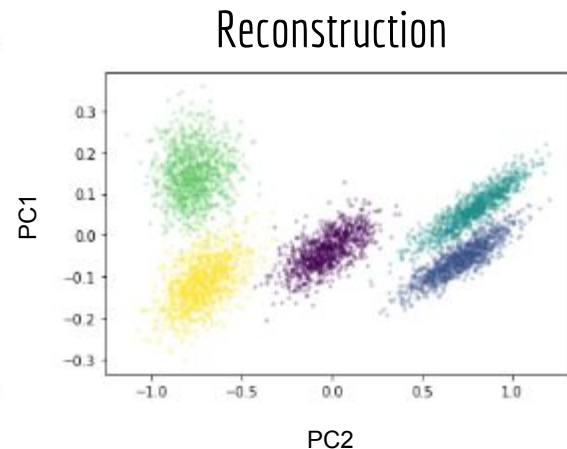
Fibered Auto-Encoders: Noise reduction with no teacher (5000x128x5)



Original signal: 5000 samples of a Gaussian in dim 128 (5 different means)



99% of values are dropped out following a bernoulli distribution



After training on corrupted inputs and targets, classes are retrieved