

The Surprising Effectiveness of Linear Unsupervised Image-to-Image Translation

Eitan Richardson

School of Computer Science and Engineering
The Hebrew University of Jerusalem
Jerusalem, Israel
Email: eitanrich@cs.huji.ac.il

Yair Weiss

School of Computer Science and Engineering
The Hebrew University of Jerusalem
Jerusalem, Israel
Email: yweiss@cs.huji.ac.il

Abstract—Unsupervised image-to-image translation is an inherently ill-posed problem. Recent methods based on deep encoder-decoder architectures have shown impressive results, but we show that they only succeed due to a strong locality bias, and they fail to learn very simple nonlocal transformations (e.g. mapping upside down faces to upright faces). When the locality bias is removed, the methods are too powerful and may fail to learn simple local transformations. In this paper we introduce *linear* encoder-decoder architectures for unsupervised image to image translation. We show that learning is much easier and faster with these architectures and yet the results are surprisingly effective. In particular, we show a number of local problems for which the results of the linear methods are comparable to those of state-of-the-art architectures but with a fraction of the training time, and a number of nonlocal problems for which the state-of-the-art fails while linear methods succeed.

I. INTRODUCTION

In unsupervised image-to-image translation we are given a set of images from domain A (e.g. black and white images of faces) and a set of images from domain B (e.g. color images of faces). We do not know the correspondence between images in the two sets (in fact, such a correspondence might not exist), and we nevertheless seek to learn a mapping from domain A to B . In a probabilistic view of the same problem, there exists some joint distribution $P_{A,B}$ over the two domains. We are given iid samples from the two *marginal* distributions P_A and P_B and we want to infer $P_{B|A}$.

This problem is inherently ill-posed. We can always define an arbitrary correspondence of the images in the two sets and learn a mapping from each image in A to its corresponding image in B . This is a perfectly valid mapping from A to B . Figure 1 illustrates the ill-posedness when both A and B are one dimensional.

Despite this inherent ambiguity, significant progress has been achieved in recent years on this problem using deep encoder-decoder architectures. Perhaps the most successful recent method is CycleGAN [1] which uses a deep encoder-decoder architecture (figure 2) together with adversarial and *cycle-consistency* loss terms. The method is demonstrated to succeed and generate high-quality outputs in a variety of tasks such as transforming black and white images to color images, turning horses into zebras and transforming edge images to real (e.g. shoes). Many methods followed

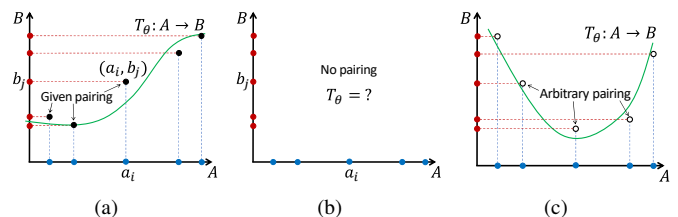


Fig. 1. (a) When pairing is given as supervision, domain translation is a regression problem – fitting some parametric transformation T_θ . (b) Without the pairs correspondence, the problem is ill-posed: (c) any arbitrary correspondence can be chosen, resulting in an arbitrary learned transformation.

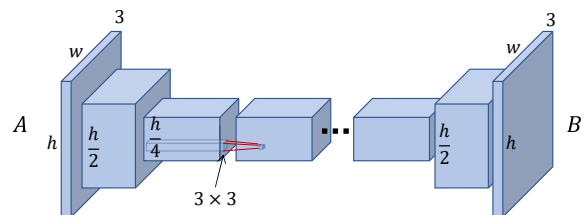


Fig. 2. The strong locality bias in CycleGAN and other *im2im* methods is mainly due to the large spatial dimension of the encoder-decoder bottleneck, typically $1/4$ of the input resolution, combined with a small convolution kernel size (e.g. 3×3).

CycleGAN, modifying the architecture [2], [3] or training objective [4] and improving different aspects of the problem, such as generating more diverse images [5], [6] and learning from fewer examples [7].

Although recent unsupervised domain translation (UDT) methods such as CycleGAN and MUNIT [5] have been successful on many unsupervised image-to-image translation problems, it is worth noting that all the problems they are demonstrated on are essentially *local*: each pixel in the output image depends only on nearby pixels in the input image and the general image structure is preserved. When this locality assumption does not hold, these methods fail to learn even very simple transformations. Figure 3 shows an example. Here the domain A is a set of vertically flipped faces and the domain B is a set of upright faces. All the algorithms have to do is learn to perform a vertical flip. As figure 3 shows, CycleGAN and MUNIT fail to learn this very simple transformation. Both

methods learn to map an upside-down face to an upright face, but the generated face is distorted and its resemblance to the input face is poor.

Presumably the strong locality bias in modern methods is due to the large spatial dimension of the encoder-decoder bottleneck that is used in both algorithms (typically 1/4 of the input resolution) (figure 2). Some of the methods (e.g. [1]) even have an optional loss term that is simply the $L1$ pixel-to-pixel distance between the input and generated images. Note that several of the UDT papers (e.g. [1], [2]) refer to the ill-posedness of the problem. Nevertheless, to the best of our knowledge, all successful methods solve the ill-posedness by using an architecture that is biased towards locality.

One way to remove the locality bias is to have an encoder-decoder bottleneck that has *no spatial dimension*. Figure 4 shows such an architecture based on the very recent ALAE method [8], which uses a StyleGAN-based [9] encoder-decoder¹. The bottleneck here is a vector of length 512 but with no spatial dimension and hence there is no particular bias towards local transformations. As shown in figure 5 when the bias towards local transformations is removed, the method learns an arbitrary mapping between the two domains, even for simple, local transformations. The figure shows the example of colorization. The deep architecture without a locality bias learns to map a gray level face image to a color image *of a different face*, even though cycle-consistency holds.

What is needed therefore is a method that can learn unsupervised image-to-image transformations but without the locality constraint. In this paper we present such a method. It is based on the assumption that the mapping from A to B is a *linear, orthogonal transformation*. Although this assumption is clearly restrictive, we show that the method is surprisingly effective – learning is much easier and faster with these architectures and yet the linear transformations are surprisingly expressive. In particular, we show a number of local problems for which the results of the linear methods are comparable to those of state-of-the-art deep architectures but with a fraction of the training time, and a number of nonlocal problems for which the state-of-the-art fails while linear methods succeed. Code will be made publicly available at <https://github.com/eitanrich/lin-im2im>.

II. OUR APPROACH

A. Linear Image-to-Image Translation

We wish to solve the following problem: given a set of images \mathcal{D}_A from domain A and a set of images \mathcal{D}_B from domain B find an *orthogonal linear transformation* T that best maps the set \mathcal{D}_A to the set \mathcal{D}_B .

How restrictive is the assumption that the transformation is linear and orthogonal? Note that any permutation of the pixels (e.g. flipping an image vertically or horizontally) is an orthogonal linear transformation. Less obvious transformations

¹Simply flattening the bottleneck in the CycleGAN architecture is not possible due to the large number of required parameters. Applying global average-pooling results in strongly distorted outputs images.

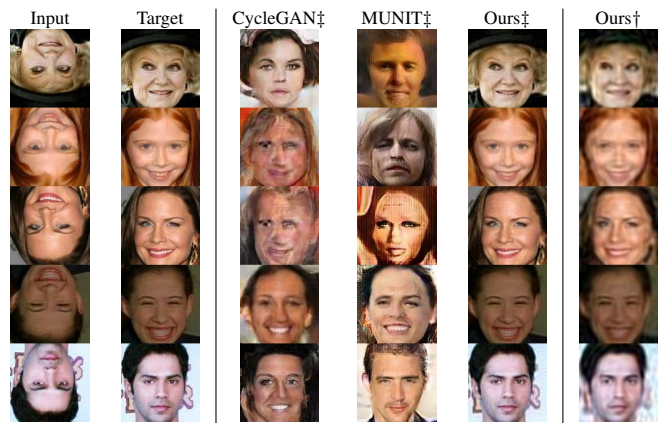


Fig. 3. Deep image-to-image translation methods are biased towards local changes and fail when the true transformation is not local (like *vertical-flip* shown here). Our proposed orthogonal transformation does not suffer from this bias and succeeds in learning non-local transformations. Domains pairing: ‡=Matching pairs exist (shuffled), †=Domains contain no matching pairs.

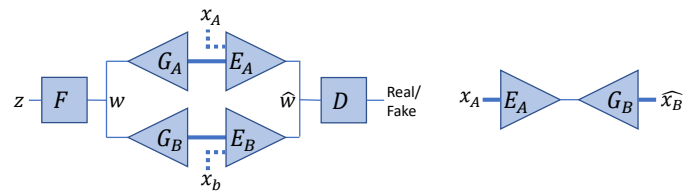


Fig. 4. ALAE-based bias-free UDT. Left: To train domain-translation, a second Generator-Encoder pair is added to the base ALAE. The combined loss ensures proper auto-encoding of each domain separately and also cycle-consistency across domains. Right: At inference time, the encoder and generator of the two domains are mixed. Tensors with a spatial dimension are shown as thick lines.



Fig. 5. A deep encoder-decoder architecture without a locality bias (the flat-bottleneck ALAE) converges to an arbitrary solution in the UDT problem ($A \rightarrow B$ has no resemblance to the reconstruction $A \rightarrow A$) even though cycle-consistency is maintained ($A \rightarrow B \rightarrow A$ resembles $A \rightarrow A$).

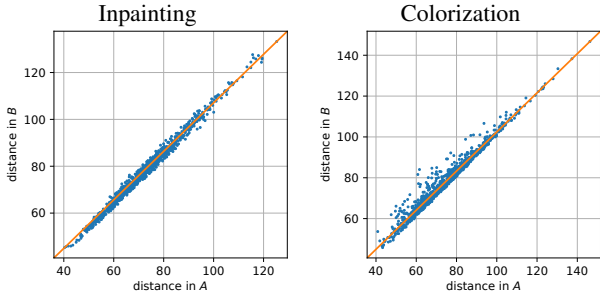


Fig. 6. A scatter plot of the L_2 distance between two images in domain A and the distance between the corresponding two images in domain B for 1000 randomly selected image pairs from FFHQ. Even though both colorization and inpainting are not invertible transformations, the distances are preserved, suggesting they can be approximated by an orthogonal transformation on a subspace of the original space.

such as inpainting or colorization are also well approximated by an orthogonal transformation. Figure 6 shows a scatter plot of the distance between two images in domain A and the distance between the corresponding two images in domain B for 1000 randomly selected image pairs. Even though both colorization and inpainting are not invertible transformations, the distances are approximately preserved, suggesting they can be approximated by an orthogonal transformation on a subspace of the original space. Specifically, in the colorization example, the gray level images occupy a subspace that is of $1/3$ the dimension of the original RGB images, and yet figure 6 suggests that if we restrict both RGB and gray level images to lie in a subspace of the same size, then an orthogonal transformation can approximate the mapping.

Even with the restriction to linear orthogonal transformations, the number of such transformations on full images is huge. For a simple example, consider 128×128 pixels color images. A linear transformation that maps a set A of such images to a set B of such images can be represented by a matrix of size 49152×49152 so that the number of free parameters is over 2.4 Billion. How can we learn such a matrix from finite training data? The following observation, shows that we can restrict ourselves to much smaller numbers of free parameters.

Observation 1: If every image in A can be approximated as $x_A = W_A z_A$ and each image in B can be approximated as $x_B = W_B z_B$ where z_A, z_B are vectors of length r ($r < d$ the image dimension) and W_A, W_B are orthogonal rectangular matrices (with orthonormal column vectors), then any linear, orthogonal transformation T from A to B can be approximated as:

$$T = W_B Q W_A^T \quad (1)$$

where Q is an $r \times r$ orthogonal matrix.²

Proof: This follows from the fact that we can write the transformation from A to B ($x_B = T x_A$) as a mapping from

²For simplicity, we assume here that the data has zero mean. In practice we subtract the mean before processing each dataset and add the mean back to produce the final output.

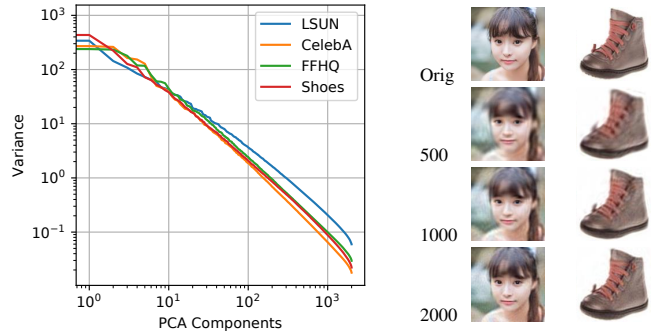


Fig. 7. Left: The PCA spectra of different image datasets behave similarly – eigenvalues decrease exponentially (with slightly lower rate in the less-structured LSUN dataset). Right: Reconstruction of FFHQ and Shoes test samples (128×128 pixels) using increasing numbers of PCA components.

z_A to z_B ($z_B = (W_B^T T W_A) z_A$). Defining $Q = W_B^T T W_A$ then it is easy to verify that $Q^T Q = I$.

B. Linear Transformation in PCA Subspace

Equation 1 has a simple interpretation as a *linear encoder decoder* architecture. The $r \times d$ matrix W_A^T encodes the input image x_A using a vector of length r , the matrix Q transforms this vector into the matching encoding of domain B , while the $d \times r$ matrix W_B decodes the vector into an image. This is similar to the architecture of modern deep image to image translation methods (e.g. figure 2), but with the important difference that the mapping from image to the latent vector is linear. In particular, this means that we can find W_A, W_B easily using Principal Component Analysis (PCA). Orthogonal transformations have been used in the past to model unsupervised domain translation by [10], [11], [12], but they use a deep encoder and decoder rather than the linear one that we use here.

How many PCA coefficients are required? Figure 7 shows that the PCA spectrum of commonly used image datasets falls off as a power law (linear in a log-log plot). This means that images of size 128×128 pixels can be very well approximated with a *linear* encoding and decoding using as few as 2000 PCA coefficients ($r = 2000$). Thus we can reduce our problem to that of finding an orthogonal matrix Q of size 2000×2000 .

C. Solving for Q

Finding an orthogonal transformation that maps a set of points in R^n to another set of points in R^n is a well-studied problem in computer graphics and computer vision [13], [14]. It is well known that in the infinite data setting, this problem can be solved efficiently up to a sign ambiguity. We briefly review this solution and then present our algorithm for the specific case of image datasets.

Observation 2: Denote by Σ_A, Σ_B the covariance matrices of the datasets $\mathcal{D}_A = \{x_1^A, \dots, x_n^A\}$, $\mathcal{D}_B = \{x_1^B, \dots, x_m^B\}$. Assume that the true relation between domains A and B is an orthogonal linear transformation T^* . If v is an eigenvector of Σ_A with eigenvalue λ then $T^* v$ is an eigenvector of the covariance Σ_B with the same eigenvalue.

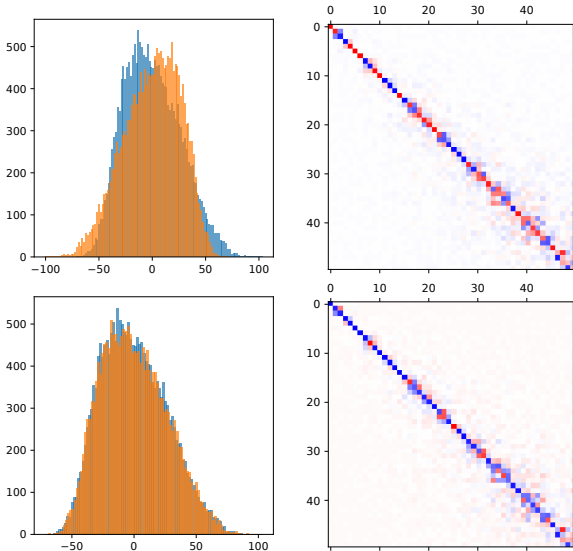


Fig. 8. Using the skewness to resolve PCA sign ambiguity. Left: Histograms of $W_0^A X_A$ and $W_0^B X_B$ (first coordinate in the data PCA embedding) before aligning the skewness (top) and after aligning (bottom). Right: The top-left block (first 50 coordinates) of Q before and after aligning the skewness (red=negative). Aligning the skewness makes Q much closer to I .

Observation 2 implies that if we had infinite data (so that we could reliably estimate Σ_A, Σ_B) and if the eigenvalues of each covariance matrix are unique, then we can recover Q exactly by sorting the eigenvectors in the two domains and mapping the k th eigenvector of Σ_A to the k th eigenvector of Σ_B . This solution still allows for a *sign ambiguity* since if v is an eigenvector of Σ_A so is $-v$. Furthermore, for any finite dataset the ordering of the eigenvectors may be changed as a result of sampling different datapoints from each dataset. Therefore the common practice is to initialize an iterative algorithm such as Iterated Closest Points (ICP) [13] from multiple initial conditions (each with a different sign chosen for the top eigenvectors) [10].

We make two modifications to the standard approach. First, we have found that for image datasets, the sign ambiguity can be mostly resolved using the *skewness* of the distribution (see figure 8). Specifically, given an eigenvector v we calculate the distribution of $v^T x$ on the dataset and calculate the skewness of that distribution:

$$skew(v) = \sum_i (v^T x_i - \mu)^3$$

If the skewness is negative, we replace v with $-v$. Second, in order to improve the accuracy of ICP, we use the “best buddy” heuristic [15]: if x_B is the closest point to $T x_A$ we also require that x_A be the closest point $T^T x_B$. See algorithm 1 for details. Once a set of corresponding pairs was found (e.g. using ICP), the orthogonal transformation Q is estimated using the *Procrustes* method (lines 12 – 13).

Note that our algorithm can be applied even when the two sets do not contain matching pairs. The existence of matching pairs improves the ICP convergence speed and the data efficiency (in the nonmatching case, the size of the

training sets needs to be sufficiently large so that cross-domain nearest-neighbors will be reasonably similar).

Algorithm 1: Orthogonal UDT in PCA subspace

Input: $\mathcal{D}_A = \{x_1^A, \dots, x_n^A\}$, $\mathcal{D}_B = \{x_1^B, \dots, x_m^B\}$, r
Result: Orthogonal transformation $T: A \rightarrow B$

- 1 Compute W_A, W_B : r principal components of $\mathcal{D}_A, \mathcal{D}_B$
- 2 Fix eigenvectors sign for positive skew
- 3 Compute PCA embedding $\{z_1^A, \dots, z_n^A\}, \{z_1^B, \dots, z_m^B\}$
- 4 Initialize $Q \leftarrow I$
- 5 **while not converged do**
- 6 $A \leftarrow \emptyset, B \leftarrow \emptyset$
- 7 **for** $i \leftarrow 1$ **to** n **do**
- 8 $j \leftarrow \arg \min_{j'} \|z_i^A Q - z_{j'}^B\|$
- 9 $k \leftarrow \arg \min_{k'} \|z_{k'}^A Q - z_j^B\|$
- 10 **if** $k = i$ **then**
- 11 $A.insert_row(z_i^A), B.insert_row(z_j^B)$
- 12 $U, S, V \leftarrow \text{SVD}(A^T B)$
- 13 $Q \leftarrow UV$
- 14 **return** $T \leftarrow W_A^T Q W_B$

In practice we find this algorithm to be very fast (training times of a few seconds for commonly used datasets, whereas training CycleGAN can take more than a day).

III. RESULTS

A. Datasets and Transformations

We conduct our experiments using two popular face datasets – CelebA [16] and FFHQ [9] and one dataset of non-face images – Shoes [17]. We trained all models on images resized to 128×128 pixels. To train our method we used 20K images from each domain. For the SOTA deep methods we used 50K images from CelebA and the entire 60K FFHQ train images. All results are shown on (unseen) test images.

We tested the following synthetic transformations:

Task Name	Domain A	Domain B
<i>vflip</i>	Vertically flipped	Original
<i>rotate</i>	90 degrees rotated left	Original
<i>colorize</i>	Grayscale image	Original
<i>inpaint</i>	Set to 0 in a mask	Original
<i>edges-to-real</i>	Sobel edges	Original
<i>super-res</i>	Reduced size by 8	Original

In most tests, sets A and B contained *matching* images – set A was generated by applying the selected transformation to the original images and then rearranging them in random order (shuffling). We also tested the *nonmatching* case in which the original dataset was first randomly split in half and then the selected transformation was applied to one half.

B. Qualitative Results

Figure 9 shows the orthogonal transformation learned by our proposed method (algorithm 1) on several transformations applied to the FFHQ dataset. As can be seen, our method can learn a variety of relations between A and B and the learned transformation can be applied successfully to unseen

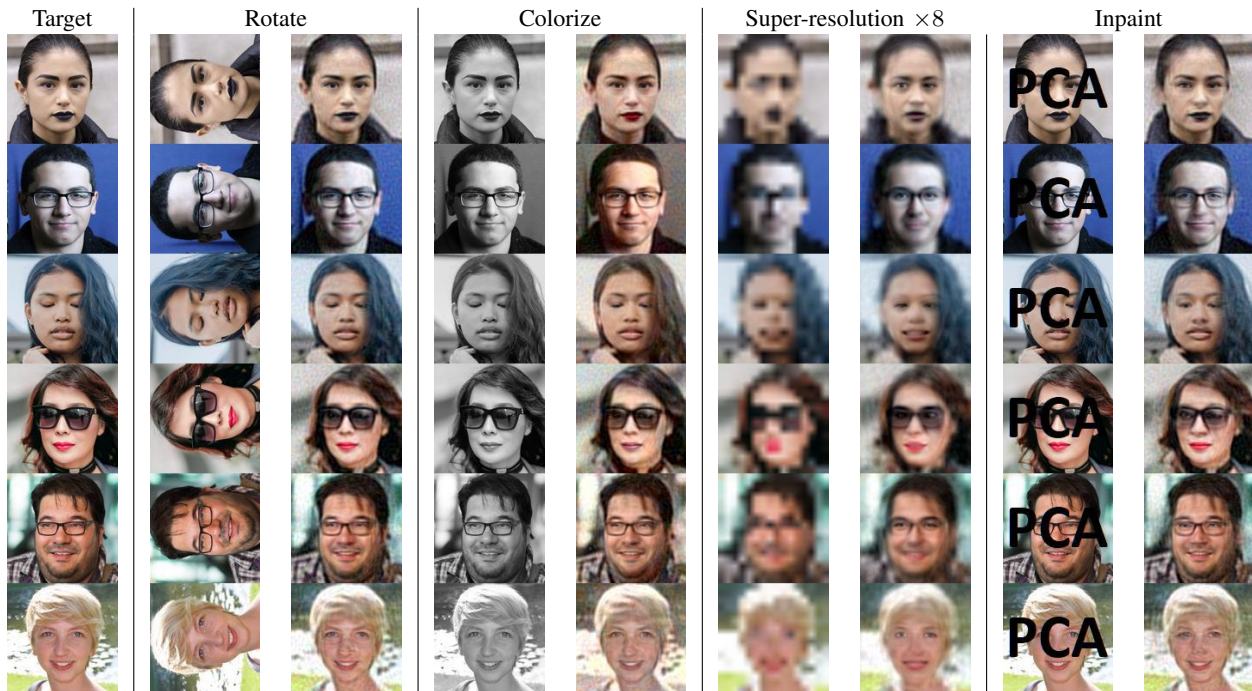


Fig. 9. Domain translation using our proposed orthogonal transformation in PCA domain, demonstrated on different FFHQ tasks. Training on 20K samples of each domain with unpaired matching pairs. Results shown on unseen test images.



Fig. 10. Results of our paired PCA-based linear transformation on non-face images – Shoes.

test images. Similar results are shown in figure 10 for the *Shoes* dataset. Results for CelebA are shown in figure 3.

Figure 11 shows the learned transformation in the *paired* setting for more challenging tasks – converting edge images to real images and in-painting where a large part of the original image is hidden. For the paired (supervised) case, the correspondence finding step (lines 6–11 of algorithm 1) is not required and the transformation is found in a single iteration using the two paired sets.

C. Quantitative Evaluation

Since we use synthetic true transformations, we can measure the quality of the learned transformation by comparing the transformed images to the “ground truth” target images. We used two commonly used image similarity measures – the mean squared error (MSE) and the structural similarity (SSIM)

[18]. Table I lists the transformation quality achieved by our method compared to the SOTA deep neural-network methods, as well as the training times. We compare both the *local* case, which is suitable to the architectural bias in the deep methods and the *nonlocal* case. For *colorization*, which is a *local* transformation, our method achieves similar results to CycleGAN and MUNIT, but at less than 1/1000 of the training time. For the two *nonlocal* tasks, our method achieves significantly better results than the deep methods.

IV. LIMITATIONS AND POSSIBLE EXTENSIONS

Figure 11 shows some limitations of our method. By definition, the main limitation of the proposed linear transformation is modeling true transformations that are very *non-linear*. An example is *edges to real-images*, in which our results are inferior to those of deep encoder-decoder methods that can model non-linear transformations. For such a setting, even a fully supervised application of our method (when the pairing between x_A and x_B are given to the algorithm) does not give good results, even if we allow arbitrary linear transformations.

Another limitation is that very complex image domains may require a large number of PCA coefficients to represent properly and still, very fine details may not be well reconstructed when passing through the low-dimensional PCA subspace.

A possible solution (figure 12) to the above limitations is combining the orthogonal transformation, which is easy to compute and not limited to local-changes with a bias-free deep encoder-decoder architecture, which excels at generating realistic and sharp images. The orthogonal transformation can guide the convergence of the over-expressive deep architecture towards the desirable solution. This approach can

TABLE I
QUANTITATIVE EVALUATION OF UNSUPERVISED IMAGE-TO-IMAGE TRANSLATION METHODS

Task	CycleGAN [‡]			MUNIT [‡]			Ours [‡]			Ours [†]		
	MSE	SSIM	T[h]	MSE	SSIM	T[h]	MSE	SSIM	T[h]	MSE	SSIM	T[h]
CelebA-colorize	0.0066	0.914	49	0.0256	0.750	52	0.0043	0.883	0.04	0.0071	0.761	0.04
CelebA-vflip	0.1167	0.358	43	0.1084	0.333	48	0.0012	0.917	0.04	0.0041	0.780	0.04
FFHQ-rot90	0.1267	0.302	39	0.1220	0.268	39	0.0023	0.870	0.05	0.0335	0.381	0.05

MSE is the mean-squared error between the input and target images, SSIM [18] estimates the perceptual similarity (higher the better) and T is the total training time in hours (including data loading).

Domains pairing: ‡=Matching pairs exist (shuffled), †=No matching pairs.

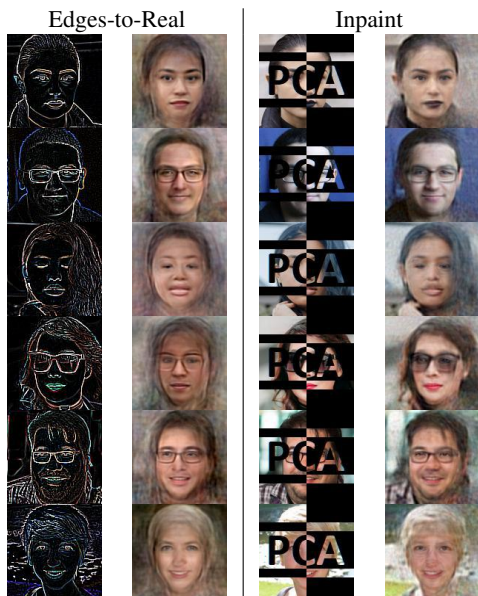


Fig. 11. Our proposed linear translation handling more challenging tasks in the supervised (paired) setting. Note that for the *edges-to-real* task, an unrestricted linear translation was chosen instead of an orthogonal one.

be thought of as a generalization of the “identity $L1$ ” loss term used by CycleGAN and other methods, replacing the naive assumption that $T(x_A) \approx x_A$ (e.g. zebras look like horses) with the more general assumption that the relation between A and B is approximately linear. Considering figure 4, the suggested loss term is therefore: $\mathcal{L}_{orthogonal} = \|G_A(w)W_A^T Q - G_B(w)W_B^T\|_2$, where W_A, W_B and Q are pre-computed using algorithm 1. As can be seen in figure 12, the deep encoder-decoder learns a $A \rightarrow B$ transformation that is more accurate than the linear transformation used for regularization (we used only 300 PCA coefficients in this test, making the linear orthogonal transformation blurry). Note that the deep encoder-decoder reconstruction itself (Input A vs $A \rightarrow A$) is not perfect. This is a current limitation of ALAE and other StyleGAN based auto-encoders.

V. DISCUSSION

Many recent deep encoder-decoder based methods demonstrate success on the inherently ill-posed unsupervised image-to-image translation problem. These methods employ sophisticated training methods such as adversarial and cycle-consistency loss. Despite their success, however, we show

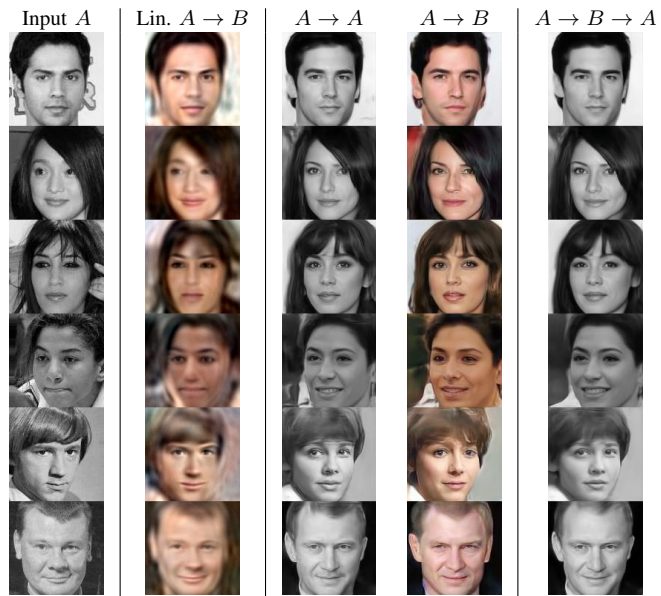


Fig. 12. ALAE (similar to figure 5, but with linear orthogonal transformation as a regularization term). Now the transformation $A \rightarrow B$ is very similar to the reconstruction $A \rightarrow A$, but with the added color information required for domain B – the transformation is no longer arbitrary.

that they rely heavily on a particular locality bias that is embodied in the architectures they use – assuming the local image structure is preserved between the two domains. We show that these methods fail when the true transformation is nonlocal and when the architecture bias is removed.

As an alternative, we presented a very different bias – *orthogonal linear transformations*. We show that such transformations can approximate a wide range of true domain relationships and solve different tasks such as in-painting and colorization, in addition to any geometric transformation. We suggest a highly efficient algorithm that expresses the orthogonal transformation in PCA subspace, requiring only a small fraction of the number of parameters compared to linear transformation in the full image space. The learning time for the linear transformations is a few seconds, compared to GPU-days for the deep methods.

We are not suggesting that linear methods should replace the deep encoder-decoder methods, however, we believe that presenting the effectiveness and surprising versatility of the orthogonal transformations can be of benefit to image-to-image research community and can help expanding the range

of solved problems beyond local transformations.

REFERENCES

- [1] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [2] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," in *Advances in neural information processing systems*, 2017, pp. 700–708.
- [3] J. Kim, M. Kim, H. Kang, and K. Lee, "U-gat-it: unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation," *arXiv preprint arXiv:1907.10830*, 2019.
- [4] Y. Hoshen and L. Wolf, "Identifying analogies across domains," in *International Conference on Learning Representations*, 2018.
- [5] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz, "Multimodal unsupervised image-to-image translation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 172–189.
- [6] H.-Y. Lee, H.-Y. Tseng, J.-B. Huang, M. Singh, and M.-H. Yang, "Diverse image-to-image translation via disentangled representations," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 35–51.
- [7] M.-Y. Liu, X. Huang, A. Mallya, T. Karras, T. Aila, J. Lehtinen, and J. Kautz, "Few-shot unsupervised image-to-image translation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 10 551–10 560.
- [8] S. Pidrorskyi, D. A. Adjeroh, and G. Doretto, "Adversarial latent autoencoders," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 104–14 113.
- [9] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 4401–4410.
- [10] Y. Hoshen and L. Wolf, "Non-adversarial unsupervised word translation," *arXiv preprint arXiv:1801.06126*, 2018.
- [11] C. Xing, D. Wang, C. Liu, and Y. Lin, "Normalized word embedding and orthogonal transform for bilingual word translation," in *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Denver, Colorado: Association for Computational Linguistics, May–Jun. 2015, pp. 1006–1011. [Online]. Available: <https://www.aclweb.org/anthology/N15-1104>
- [12] Y. Hoshen and L. Wolf, "Unsupervised correlation analysis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3319–3328.
- [13] P. J. Besl and N. D. McKay, "Method for registration of 3-d shapes," in *Sensor fusion IV: control paradigms and data structures*, vol. 1611. International Society for Optics and Photonics, 1992, pp. 586–606.
- [14] A. Myronenko and X. Song, "Point set registration: Coherent point drift," *IEEE transactions on pattern analysis and machine intelligence*, vol. 32, no. 12, pp. 2262–2275, 2010.
- [15] T. Dekel, S. Oron, M. Rubinstein, S. Avidan, and W. T. Freeman, "Best-buddies similarity for robust template matching," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2021–2029.
- [16] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- [17] A. Yu and K. Grauman, "Fine-grained visual comparisons with local learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 192–199.
- [18] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.