# LEARNING ROTATION INVARIANT FEATURES FOR CRYOGENIC ELECTRON MICROSCOPY IMAGE RECONSTRUCTION

*Koby Bibas\*, Gili Weiss-Dicker\*, Dana Cohen, Noa Cahan, Hayit Greenspan*

Faculty of Engineering, Tel Aviv University, Israel

## ABSTRACT

Cryo-Electron Microscopy (Cryo-EM) is a Nobel prize-winning technology for determining the 3D structure of particles at near-atomic resolution. A fundamental step in the recovering of the 3D single-particle structure is to align its 2D projections; thus, the construction of a canonical representation with a fixed rotation angle is required. Most approaches use discrete clustering which fails to capture the continuous nature of image rotation, others suffer from low-quality image reconstruction. We propose a novel method that leverages the recent development in the generative adversarial networks. We introduce an encoder-decoder with a rotation angle classifier. In addition, we utilize a discriminator on the decoder output to minimize the reconstruction error. We demonstrate our approach with the Cryo-EM 5HDB and the rotated MNIST datasets showing substantial improvement over recent methods.

***Index Terms***— Cryo-EM, 5HDB, Rotated MNIST, Deep learning, Image synthesis, Generative adversarial networks

## 1. INTRODUCTION

Unsupervised feature learning algorithms have emerged as a promising tool for learning representations from data [1, 2, 3]. Learning invariant image representation enables machine learning algorithms to achieve good generalization performance while using a small number of labeled training examples. An invariant representation is particularly valuable for the Cryo-EM, where the goal is to determine the 3D electron density of a particle from many noisy and randomly oriented 2D projections. Having a model that aligns the 2D particle pose to a canonical predefined posture could significantly improved the 3D reconstruction of the particle [4].

Existing classic methods for the problem of determining the 3D structure of a particle use a Gaussian mixture model to group these 2D views. However, this assumes a discrete set of projections where it is known that particle conformations are continuous. To face this issue, more recent machine learning-based disentanglement approaches do not impose a specific structure on the learned latent representations [1]. These methods, however, use a variational autoencoder (VAE). Using VAE

induces blurriness to the reconstructed image which might eliminate key components in the particle structure.

In this work, we propose a different approach. We use an encoder-decoder architecture and a discriminator. The discriminator penalizes the decoder generated images with rotated content. In this way, we ensure that the generated content has a fixed orientation, which later can be utilized to reconstruct the 3D shape of the particle. We demonstrate the effectiveness of our method on Cryo-EM and rotated MNIST datasets [1]. We improve the current leading approach mean squared error (MSE) by an order of magnitude on both the Cryo-EM 5HDB and the rotated MNIST datasets [1].

## 2. RELATED WORK

In Cryo-EM, the main challenge is to determine the structure of a protein or a particle. In this section, we describe related work that tackles this problem.

Classic statistical methods assume that the many 2D projection observations of the particle arise from either a single structure or from a discrete mixture of structures. Assuming there are a finite number of possible projections, the particle views are grouped into a discrete number of clusters [4]. These conformations are confounded by orientation in the collected images. Thus, their goal is to cluster each projection image into one of the possible finite sets of projections.

Modern approaches tackle the problem of disentangling latent variables in an unconstrained setting [5]. Others constrained the manifold of latent values to be homeomorphic to some known underlying data manifold to capture useful latent representations [6].

The recently suggested *spatial-VAE* [1] addresses this problem by formulating the generative model as a function of the spatial coordinates. This makes the reconstruction error differentiable with respect to latent rotation parameters which creates a representation that is independent of the content pose.

A similar approach to ours was taken by the *AttGAN* [7]. It uses an encoder-decoder architecture for facial attribute editing by conditioning the decoding of a given face latent representation on the desired attributes. Notice that in our case the rotation angle is a continuous variable as opposed to face attributes and this imposes an additional challenge.
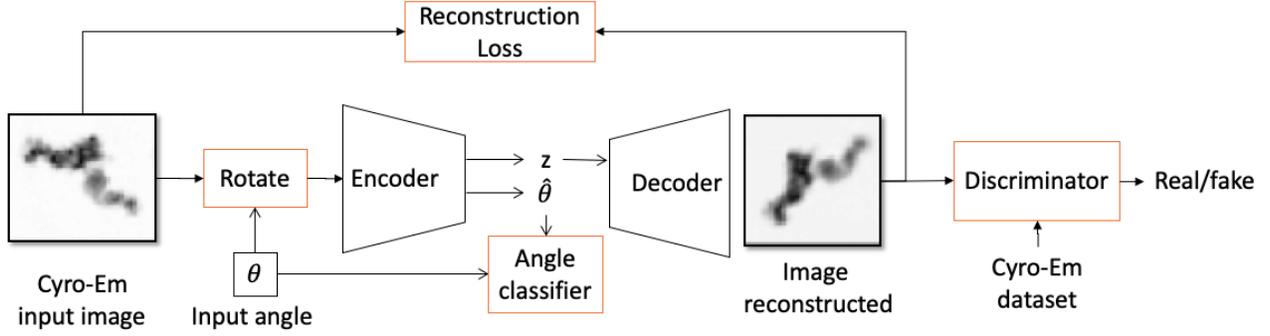
---

\* Equal contribution

**Fig. 1**: Our proposed scheme. We use encoder-decoder architecture with an angle classifier and a discriminator.

## 3. THE PROPOSED APPROACH

In this section, we describe our proposed method of disentangling the image content from the object pose in order to obtain a rotation-invariant image representation.

We employ an encoder-decoder architecture with a discriminator. A block diagram of the proposed system in shown in Figure 1. It contains the following stages: First, we apply a random rotation $\theta$ on the input image and propagate it through the encoder. We utilize one scalar from the latent vector to represent the input image rotation. Denote $\hat{\theta}$ as the predicted rotation (the latent variable value), we use the following loss for the angle classifier model

$$\mathcal{L}_{angle}(\theta, \hat{\theta}) = \exp\left\{|\theta - \hat{\theta}|\right\} - 1. \tag{1}$$

When an accurate prediction is achieved, the $\exp$ value is 1 and the total loss is 0. Next, we decode the latent vector with a decoder and compute the reconstruction loss with the original unrotated image. Denote $x$ and $\hat{x}$ as the input image and the reconstructed image respectively, the reconstruction loss is

$$\mathcal{L}_{rec}(x, \hat{x}) = ||x - \hat{x}||_2 + ||x - \hat{x}||_1. \tag{2}$$

Finally, we use a discriminator that gets as input the reconstructed images and the given dataset. We use the Wasserstein loss [8]. Denote $G$ as the generator (the decoder in our case), $D$ as the discriminator and $z$ as the latent vector without the rotation variable, the loss functions of the discriminator and of the decoder are

$$\begin{aligned}\mathcal{L}_{adv\text{-}disc}(x, D, G) &= D(x) - D(G(z)), \\ \mathcal{L}_{adv\text{-}decoder}(x, D, G) &= D(G(z)).\end{aligned} \tag{3}$$

The adversarial loss is

$$\mathcal{L}_{adv}(D, G) = \mathcal{L}_{adv\text{-}disc}(D, G) + \mathcal{L}_{adv\text{-}decoder}(D, G). \tag{4}$$

The final training loss function is the combination of the above loss functions

$$\mathcal{L}(\theta, \hat{\theta}, x, \hat{x}, D, G) = \mathcal{L}_{angle}(\theta, \hat{\theta}) + \mathcal{L}_{rec}(x, \hat{x}) + \mathcal{L}_{adv}(D, G). \tag{5}$$

In the following sections, we show that using our approach we manage to better reconstruct the unrotated images.

## 4. DATASETS

In order to evaluate performance, we conducted experiments using the following datasets.

**5HDB dataset [1, 9].** A Cryo-EM dataset that contains simulated 2D projections with random rotations and additive random noise. The dataset includes 20K simulated projections of integrin $\alpha$-IIb with integrin $\beta$-3. The image size is 40x40. We used 16K and 4K images for training and testing respectively.

**Rotated MNIST [1].** Each image from the MNIST dataset is rotated by a random angle sampled from $\mathcal{N}(0, \frac{\pi^2}{16})$. Training and testing sets consist of 60K and 10K images respectively.

For both datasets, we normalized the pixel values such that their values is between 0 and 1. We did not pre-process these datasets, besides the random rotation and normalization that were mentioned.

## 5. EXPERIMENTS

In this section, we present experiments[1] that test our proposed encoder-decoder with a discriminator scheme as a method to disentangle the image content from the pose. We compare our approach and the spatial-VAE method, which is considered a leading method in learning rotation-invariant features for Cryo-EM datasets.

For both 5HDB and rotated MNIST datasets, we trained our suggested model for 300 epochs with a learning rate value of $10^{-4}$ with a decrease by 0.1 after 200 epochs. For every 4 steps of the decoder, the discriminator was updated once. We used also a weight decay value that equals $10^{-5}$.

---

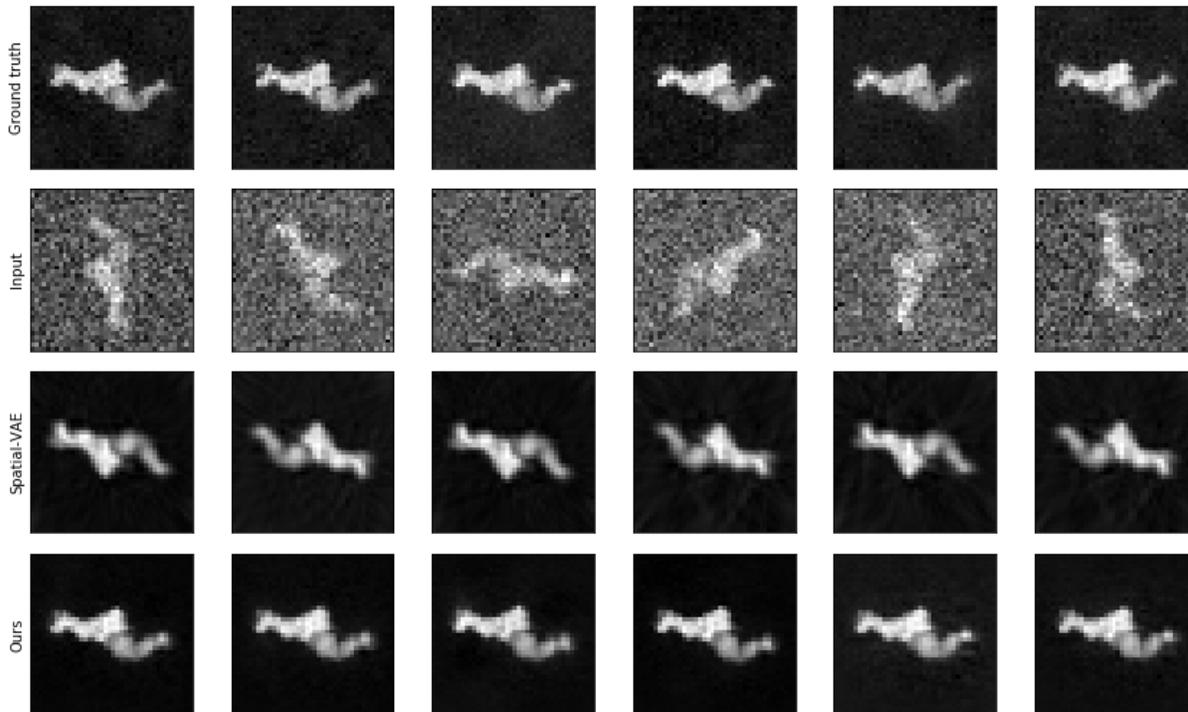[1]Code is available in `https://github.com/kobybibas/CryoEM_rotation_invariant`

**Fig. 2**: Proposed method results compared to baseline spatial-VAE for different rotation in the input. The first row is the ground true. The second row presents the model inputs. The 3rd and 4th rows show the spatial-VAE and our method results respectively.

**Table 1**: Performance of our rotation invariant auto-encoder (AE) and the spatial-VAE model on the test set

| Dataset | Method | Average MSE | Worst MSE |
|---------|--------|-------------|-----------|
| 5HDB | Spatial-VAE | 2.1 | 3.29 |
| | Rotation invariant AE | 0.3 | 0.81 |
| MNIST | Spatial-VAE | 66.07 | 121.82 |
| | Rotation invariant AE | 0.02 | 0.18 |

### 5.1. 5HDB dataset

In order to evaluate the performance of our method, we measured the MSE between the ground truth image and the decoder output. The average MSE of the 5HDB test set is described in Table 1. Our model outperforms the baseline in terms of average MSE compared to the original unrotated image by an order of magnitude.

In Figure 2 we present the outputs of the compared models. One can see that the object rotation of our model outputs is similar to the ground truth, where the rotation of the spatial-VAE model is different and is also changed based on the protein rotation in the input images: In the first, third, and fifth columns the orientation of the protein using the spatial-VAE model is flipped with respect to the ground truth.

We also evaluate the average MSE of the predicted angle by our method $|\theta - \hat{\theta}|^2$ on the 5HDB dataset. The result is an average MSE of 0.17 radians.

The worst-case image is the image with the highest MSE between the ground truth image (the unrotated image with no noise) and the output of the model. A visualization of the worst-case 5HDB images of our model and the spatial-VAE model is presented in Figure 3. The spatial-VAE suffers from blurriness and did not reconstruct the image correctly. On the contrary, our model worst-case image can be considered as a successful reconstruction.
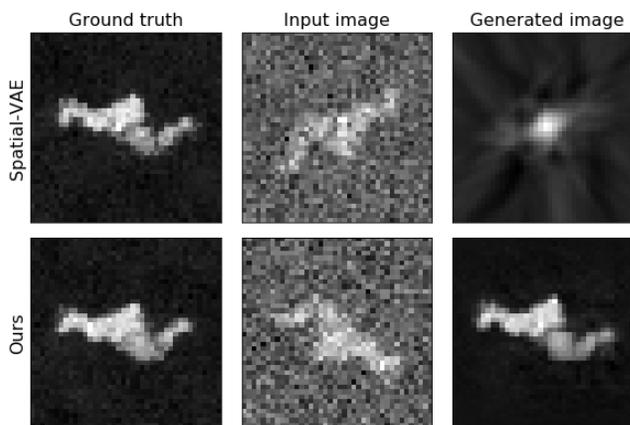


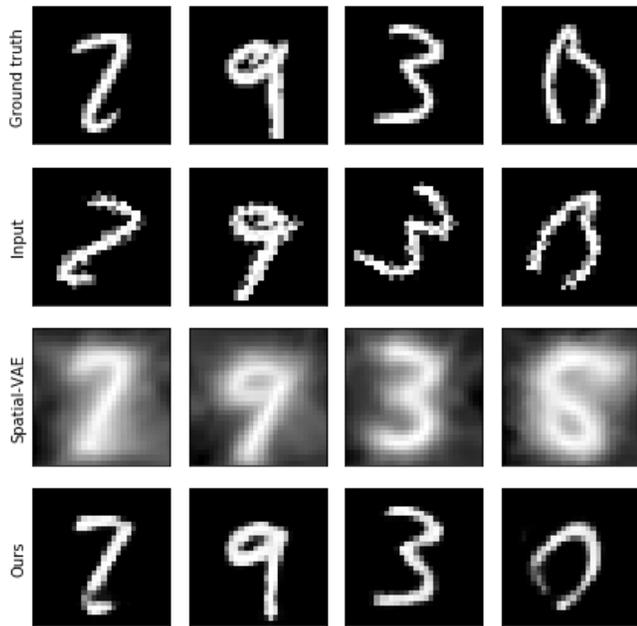**Fig. 3**: 5HDB dataset worst case MSE image.

**Fig. 4**: Rotated MNIST dataset. The first and second rows are the ground truth and the input images respectively. The 3rd and 4th rows show the spatial-VAE and our method results.

## 5.2. Rotated MNIST dataset

We use the same metrics in the evaluation of the rotated MNIST dataset as in the 5HDB dataset.

The average MSE of the rotated MNIST test set is shown in Table 1. Our method attains an average MSE of 0.02 which is two orders of magnitude better than the spatial-VAE which has an average MSE of 66.07.

We show in Figure 4 qualitative results of the rotated MNIST dataset. As shown in the second the fourth columns, our method reconstructs the unrotated image with greater accuracy than the spatial-VAE method. In addition, the spatial-VAE outputs are blurred which explains the high MSE values of this method.

## 6. CONCLUSION

In this work, we suggested a novel encoder-decoder architecture with a discriminator to produce a canonical representation of cryogenic electron microscopy images. Our suggested method offers an improvement on the 5HDB single-particle electron microscopy and rotated MNIST datasets. This is evident in both the quantitative and qualitative results.

We are currently exploring additional variations to the proposed architecture, and its generalization to additional attributes. In the future, our method can be extended to additional modalities, such as CT and MRI imaging, and can help generate canonical representations and invariant reconstruction in various tasks.

## 7. COMPLIANCE WITH ETHICAL STANDARDS

This is a numerical simulation study for which no ethical approval was required.

## 9. REFERENCES

[1] Tristan Bepler, Ellen Zhong, Kotaro Kelley, Edward Brignole, and Bonnie Berger, "Explicitly disentangling image content from translation and rotation with spatial-vae," in *Advances in Neural Information Processing Systems*, 2019, pp. 15409–15419.

[2] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze, "Deep clustering for unsupervised learning of visual features," in *ECCV*, 2018, pp. 132–149.

[3] Dotan Kaufman, Koby Bibas, Eran Borenstein, Michael Chertok, and Tal Hassner, "Balancing specialization, generalization, and compression for detection and tracking," *Proceedings of the British Machine Vision Conference (BMVC)*, 2019.

[4] Amit Singer and Fred J Sigworth, "Computational methods for single-particle cryo-em," *Annual review of biomedical data science*, vol. 3, 2020.

[5] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," in *Advances in neural information processing systems*, 2016, pp. 2172–2180.

[6] Luca Falorsi, Pim de Haan, Tim R Davidson, Nicola De Cao, Maurice Weiler, Patrick Forré, and Taco S Cohen, "Explorations in homeomorphic variational auto-encoding," *arXiv preprint arXiv:1807.04689*, 2018.

[7] Zhenliang He, Wangmeng Zuo, Meina Kan, Shiguang Shan, and Xilin Chen, "Attgan: Facial attribute editing by only changing what you want," *IEEE Transactions on Image Processing*, vol. 28, no. 11, pp. 5464–5478, 2019.

[8] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville, "Improved training of wasserstein gans," in *Advances in neural information processing systems*, 2017, pp. 5767–5777.

[9] Fu-Yang Lin, Jianghai Zhu, Edward T Eng, Nathan E Hudson, and Timothy A Springer, "$\beta$-subunit binding is sufficient for ligands to open the integrin $\alpha iib\beta 3$ headpiece," *Journal of Biological Chemistry*, vol. 291, no. 9, pp. 4537–4546, 2016.