# NOVEL TRACKING APPROACH BASED ON FULLY-UNSUPERVISED DISENTANGLEMENT OF THE GEOMETRICAL FACTORS OF VARIATION

### A PREPRINT

**Mykhailo Vladymyrov**[*]
Albert Einstein Center for Fundamental Physics
Laboratory for High Energy Physics
University of Bern
Bern 3012, Switzerland
mykhailo.vladymyrov@tki.unibe.ch

**Akitaka Ariga**
Albert Einstein Center for Fundamental Physics
Laboratory for High Energy Physics
University of Bern
Bern 3012, Switzerland
akitaka.ariga@lhep.unibe.ch

February 14, 2020

### ABSTRACT

Efficient tracking algorithms are a crucial part of particle tracking detectors. While a lot of work has been done in designing a plethora of algorithms, these usually require tedious tuning for each use case. (Weakly) supervised Machine Learning-based approaches can leverage the actual raw data for maximal performance. Yet in realistic scenarios, sufficient high-quality labeled data is not available. While training might be performed on simulated data, the reproduction of realistic signal and noise in the detector requires substantial effort, compromising this approach.

Here we propose a novel, fully unsupervised, approach to track reconstruction. The introduced model for learning to disentangle the factors of variation in a geometrically meaningful way employs geometrical space invariances. We train it through constraints on the equivariance between the image space and the latent representation in a Deep Convolutional Autoencoder. Using experimental results on synthetic data we show that a combination of different space transformations is required for meaningful disentanglement of factors of variation. We also demonstrate the performance of our model on real data from tracking detectors.

***Keywords*** Unsupervised Learning, Disentangling factors of variation, particle tracking

## 1 Introduction

Particle tracking detectors allow us to study elementary particle interactions by visualizing particle trajectories. Robust tracking algorithms are nowadays a fundamental component of all tracking detector techniques. Tracking techniques in particle physics have evolved along with technological developments, from implementations on hardware logical elements, computer data processing, GPU-accelerated algorithms, to modern Deep-Learning based approaches [1]–[3]. The advanced implementation of tracking algorithms can be seen for example in emulsion detector data reconstruction.

Nuclear photoemulsion (referred to as *emulsion* in further text) detectors are tracking detectors that allow the detection of charged particles with high spatial (50 nm) and angular (<1 mrad) resolution. They do not require a power supply during the experimental run. These properties enable fundamental physical experiments searching for short lived particles [4]–[7] and large scale experiments in remote regions [8]. The emulsion gel consists of small silver bromide crystals dispersed in a gelatin frame. When a charged particle passes through the emulsion gel, the crystals along its trajectory create latent image centers, which become visible under optical microscopes after chemical development (Figure 1).

---

[*]Presently at: Theodor Kocher Institute, University of Bern and Science IT Support, University of Bern
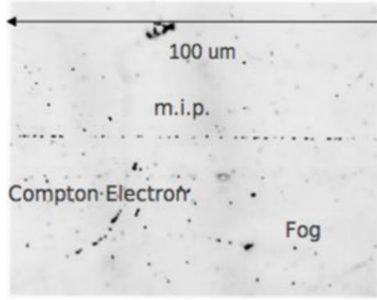
Figure 1: An emulsion detector viewed under a microscope. Tracks of charged particles are visible as sequences of silver grains. The detector is sensitive to particles with minimal ionization (m.i.p.).

Conventional track reconstruction is performed in three steps. First, 3D tomographic images of the emulsion detector are acquired using automated scanning microscopes. Next, the positions of silver grains ("hits") are located in the 3D image volumes, and finally tracks in the detector volume are reconstructed as a sequence of hits along straight lines [2], [9]. For this detector, the typical track curvature radius is significantly larger than the track length in a single emulsion film, and thus the local curvature is ignored.

Several tracking algorithms were developed during the evolution of the scanning systems, allowing for efficient track reconstruction in real-time during data acquisition [10]–[12], as well as for particle identification and energy measurements [13]–[15]. While satisfying the needs of many experiments they have several drawbacks. Their adaptation to different experimental condition, e.g. high track density or high background level requires tedious calibration ranging from extensive parameter tuning to performing dedicated test runs using e.g. accelerator beams. In addition, when the procedure of extracting the hits is separated from track reconstruction, the tracking algorithm cannot fully exploit the information available in the raw image data, compromising performance especially in the high background/track density cases.

Incorporating tracking based on classical Deep Learning, where the track parameters are predicted from the raw image data, would naturally address the latter issue. Yet, to train such a model in a supervised manner either one would need to provide massive amounts of labeled 3D raw image data for each experimental case, or training would need to be performed largely on simulated datasets. While suitable for some similar cases [3] this approach requires perfect knowledge of the optical microscope parameters, grain size distributions, detector noise, etc., which are not always available directly.

Similarly to recent works where, for example, the underlying factors of variation in the images of faces, such as eyes or hair color, glasses, or head tilt are disentangled [16], it is possible to identify geometrical factors of variation of tracks. Training such models in an unsupervised manner, i.e. where no track parameters (labels) are provided during the training can address the issues mentioned above simultaneously, by both leveraging raw image data for efficient track reconstruction and allowing simple adaptation to new configurations requiring the raw image dataset only.

In this work we aim at studying an unsupervised learning approach for extracting track parameters solely from the raw image data. Here we introduce a tracking approach based on the Deep Convolutional Autoencoder [17], [18] model that learns to disentangle the geometrical factors of variation (coordinates and angles of each track) in a fully unsupervised manner by imposing equivariance of the space transformation. While the reconstruction constraint alone fails to disentangle the factors of variation in a meaningful way, we show that adding a simple constraint on the translational invariance along the track line also does not lead to the desired disentanglement. We demonstrate that incorporating more sophisticated transformations in the latent representation is necessary to avoid the reference ambiguity.

The remaining of the paper is structured as follows. In Section 2. the details of the proposed equivariance constraints, latent representation interpretation and implementation details will be given. In Section 3. we will show how different constraints affect the performance and carry out an in-depth study of encoder and decoder performance separately to better understand the learned representation. Also, performance on real emulsion data will be shown. Finally, in Section 4 the applicability of the approach and future prospects will be discussed.
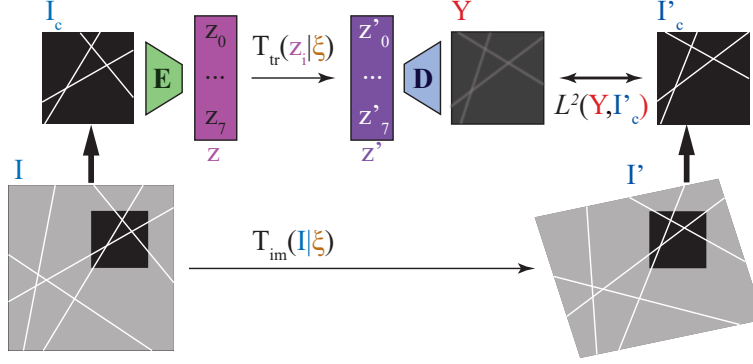
Figure 2: Schematic view of Deep Convolutional Autoencoder consisting of an encoder $E$ and a decoder $D$. Geometrical space transformations are applied to images $I$ and each of the latent representation blocks $z_i$.

## 2 Methods

### 2.1 Equivariance constraint

In this work, we will simplify the problem to the 2D case. We use synthetic image data of the emulsion detector tracks to perform a study of the proposed approach. Also, we demonstrate the performance of the trained model on real emulsion detector data.

We use a Deep Convolutional Autoencoder consisting of an encoder $E$ and a decoder $D$ as illustrated in Fig. 2, that are trained in an end-to-end manner. The encoder $E$ acts on $32 \times 32$ pixel images $I_c$, which are obtained from the full images $I$ using the cropping operation $C$: $I_c = C(I)$, producing the latent representation, $z$. In our setup, $z$ is used to estimate the geometrical parameters, namely position and angle, of tracks present in the image crop.

We then define a set of geometrical transformations acting both in the image representation space and in the latent representation space parametrized by the same parameter set $\xi$. In the image space $I' = T_{im}(I|\xi)$ and in the latent space of track parameters $z'_t = T_{tr}(z_t|\xi)$.

Given the encoder and decoder functions

$$z = E(I_c),$$
$$Y = D(z),$$

we then demand equivariance of both encoder and decoder under these transformations, i.e. the commutation of the encoder and decoder functions $E$ and $D$ with the transformations $T$ in corresponding domain:

$$T_{tr}(E(I)|\xi) = (E(T_{im}(I|\xi)); \forall \xi$$

$$D(T_{tr}(z|\xi)) = T_{im}(D(z)|\xi); \forall \xi$$

From which, assuming $I = D(E(I))$, it follows that

$$D(T_{tr}(E(I)|\xi)) = T_{im}(I|\xi); \forall \xi,$$

where cropping operations are omitted for brevity. This allows us to formulate the optimization problem in an end-to-end manner, primarily through the minimization of the $L^2$ loss between the cropped transformed image $I'_c = C(I')$ and the decoder output $Y$ (Figure 2):

$$E, D = \arg\min_{E,D} \mathbb{E}_{I,\xi} L^2(Y, I'_c) = \arg\min_{E,D} \mathbb{E}_{I,\xi} L^2(D(T_{tr}(E(C(I))|\xi)), C(T_{im}(I|\xi))).$$

We will show that with a sufficient set of transformations $T$, the model is able to learn a geometrically meaningful latent representation $E(I)$.

### 2.2 Interpretable latent representation

We limit the number of tracks potentially detected on each cropped image to $n = 8$, slightly above the maximum possible number of tracks per crop (five) in our data set (see section 2.5 for details). We will parametrize a track with
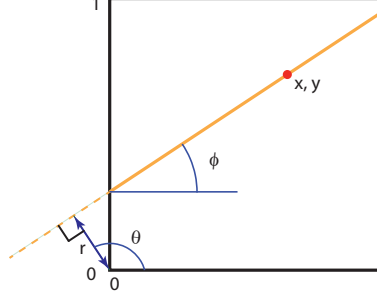
Figure 3: Coordinate range on the image is chosen to be $(0,0)$ at bottom left to $(1,1)$ at right top. Track line is parametrized either by a point on the track $(x,y)$ and sine and cosine of its slope angle $\phi$ as $c = \kappa \cos(\phi)$, $c = \kappa \sin(\phi)$, or by distance $r$ to the track from the origin and angle $\theta$ as $c = \kappa \cos(\theta)$, $c = \kappa \sin(\theta)$.

$n_p$ parameters. Thus, the encoder is designed to output a vector $z$ of length $n \cdot n_p$. This vector is then partitioned into $n$ chunks of length $n_p$. We further refer to these chunks as "track feature containers" $z_i$, each corresponding to one of the $n$ tracks.

We attribute an *a priori* meaning to each of the $n_p$ elements in track features $z_i$. Here we have explored three parametrization options:

1. $n_p = 4$, $z_i = (x_i, y_i, c_i, s_i)$ – the track's geometrical parameters, where $x_i$, $y_i$ are the coordinates of a point on the $i$-th track line and $c_i = \kappa \cos \phi_i$, $s_i = \kappa \sin \phi_i$ are proportional to the cosine and sine of the track inclination angle $\phi_i$. Coordinates on the image crop are $(0,0)$ in the lower left corner and $(1,1)$ in the upper right, $\phi \in [0, 2\pi)$, so that $\sin \phi, \cos \phi \in [-1,1]$ (Figure 3). Such a parametrization is chosen because it is continuous and confined, unlike e.g. $\phi$ itself or $\tan \phi$.

2. $n_p = 5$, $z_i = (x_i, y_i, c_i, s_i, a_i)$, where the first 4 elements are the track's geometrical parameters as in 1), and the parameter $a_i \in [0,1]$ shows the confidence of the encoder in the track presence. A value of $a_i \leq 0$ is defined as a disabled track, and $a_i > 0$ as enabled. By disabling some tracks, degenerate outputs, in which multiple containers predict the same track, can be prevented.

3. $n_p = 3$, $z_i = (r_i, c_i, s_i)$ – the track's geometrical parameters in the rho-theta parametrization [19]. $c_i = \kappa \cos \theta_i$, $s_i = \kappa \sin \theta_i$ – are proportional to the cosine and sine of the angle $\theta_i = \phi_i + \frac{\pi}{2}$, and $r_i$ is the distance from the origin $(0,0)$ to the track (Fig. 3).

The first two parametrizations are overparametrized yet more naturally occurring in the image representation. Importantly, these enable an explicit implementation of the translation invariance. The last one is the most common parametrization for 2D tracking (e.g. Hough transformation [19]).

Since by implementation the values are $z_i \in [-1,1]$, the parameter $r_i$ is scaled linearly into the $[0, \sqrt{2}]$ range.

## 2.3 Representation transformations

Image transformations are often used for image augmentation during model training [20] and in some cases, models are trained to recover original images from transformed ones [21]. Instead, here we explicitly apply transformations coherently in the image and latent spaces, as required for the equivariance constraint. As space transformations we employ affine transformations as a combination of rotation, scaling, skew, and translation. Under these affine transformations, straight lines are transformed to straight lines. We implemented these transformations coherently in the image and latent representation spaces. Details on the transformation's implementation are given in the Appendix A.

The main property of a line is the translational invariance along it. In the present work we tried to see if this translation transformation can be sufficient for disentangling the geometrical parameters of a line in the latent representation. Also, we have studied the effect of such a transformation when incorporated in addition to the affine transformations.

In this work we included the five model configurations of constraints based on the equivariance between the image space and the latent representation of the track line parameters (Table 1).

In the models, which output the track activation parameter $a$, when track container $z_i$ is marked by the encoder as not active ($a_i \leq 0$), we reset the geometrical parameters of this track to random values. This operation forces the decoder to learn to ignore the disabled track.

Table 1: Model configurations.

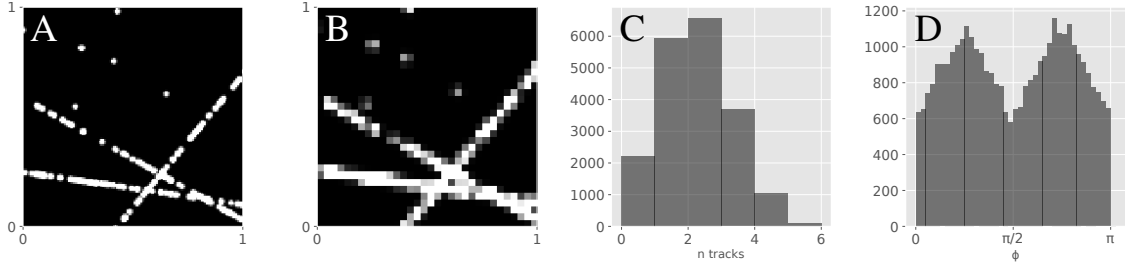| | Name | Description | $n_p$ | Parametrization |
|---|---|---|---|---|
| 1 | **AT+A** | Affine Transformations with track activation parameter $a$ | 5 | $(x_i, y_i, c_i, s_i, a_i)$ |
| 2 | **AT+TI+A** | Affine Transformations + Translation Invariance with track activation parameter $a$ | 5 | $(x_i, y_i, c_i, s_i, a_i)$ |
| 3 | **RT+TI+A** | Translation Invariance + Rotation transformation only with track activation parameter $a$ | 5 | $(x_i, y_i, c_i, s_i, a_i)$ |
| 4 | **AT+TI** | Affine Transformations + Translation Invariance | 4 | $(x_i, y_i, c_i, s_i)$ |
| 5 | **AT, rcs** | Affine Transformations in the (r,c,s) parametrization | 3 | $(r_i, c_i, s_i)$ |



Figure 4: A. Generated dataset sample. B. Downscaled sample used for model training. C, D. Distributions of number of tracks per image and track angles.

## 2.4 Loss function

The model training is performed by minimizing the loss function. In the models without the track activation parameter $a$, the loss function contains only the image term $L_{im}$, which describes the dissimilarity between the decoder output $Y$ and the transformed image $I'$ with the $L^2$ measure:

$$L = L_{im} = \underset{pixels}{\mathbb{E}} \left( \lambda\, \alpha_{sig}(c_0 + I'_c)^2 + (1-\lambda)\, \alpha_{L2} \right)\ (Y - I'_c)^2$$

Here $\mathbb{E}$ denotes averaging over all image pixels. The $(\lambda\, \alpha_{sig}(c_0 + I'_c)^2 + (1-\lambda)\, \alpha_{L2})$ term scales the loss in the image regions with high signal intensity $I'_c$ in the beginning of training ($\lambda = 1$): $\mathbb{E}\, \alpha_{sig}(c_0 + I'_c)^2 (Y - I'_c)^2$. The coefficient $c_0$ prevents the loss from dropping to zero in low intensity image regions. As the training progresses, $\lambda$ exponentially decreases, and the loss is relaxed to pure $L^2$: $\mathbb{E}\, \alpha_{L2}(Y - I'_c)^2$.

In the models with the track activation parameter $a$, the loss function contains three terms:

$$L = L_{im} + L_{unif\_act} + L_{bin\_act}$$

The first term describes the $L^2$ pixel value measure as described above. The next two terms address the information flow problem (also referred to as shortcut problem in ref. [22]), in which the geometrical parameters $(x, y, c, s)_i$ of multiple tracks describe the same track or some of them are ignored. This is achieved in two steps. First, we demand that each track $t_i$ is found (and marked as active by the parameter $a_i > 0$) on average at the same rate as the others:

$$L_{unif\_act} = \alpha_{unif} \underset{i}{\mathbb{E}} \left( \bar{a}_i - \bar{a} \right)^2,$$

where the mean activation of all tracks $\bar{a} = \mathbb{E}_{mb,i}\, a_i$ and the mean activation of a particular track $\bar{a}_i = \mathbb{E}_{mb}\, a_i$ are calculated over the minibatch. The final term $L_{bin\_act} = \alpha_{bin}\, \mathbb{E}\left( 1 - a^2 \right)$ forces the activation parameter $a$ to cluster at values $-1$ or $1$, enforcing the encoder decision on whether a track is enabled or disabled.

## 2.5 Training data

For model training and evaluation, we generate synthetic images resembling noisy emulsion data. They contain of two types of objects: "particle tracks" and noise, so called "fog". Tracks are chains of bright Gaussian spots with the spot density per unit length sampled from a Poisson distribution with mean $\mu$ located randomly along the straight lines with

deviation $d \in N(0, \sigma_d)$. Fog is represented by Gaussian spots uniformly distributed in the area with density $\rho$. The track density as well as the $\mu$, $\sigma_d$, and $\rho$ parameters approximately correspond to usual experimental conditions [9] and remain fixed throughout this study. While the generated dataset resolution matches the usual imaging resolution, we downscale the images by a factor of four to facilitate this study (Figure 4A,B).

### 2.6 Model implementation and training procedure

In both the encoder and decoder, we incorporated the CoordConv approach [23] in the first layer, which conceptually fits in our study. In this approach, two additional channels, containing $x$ and $y$ coordinates of pixels in the range of $[0, 1]$ correspondingly, are concatenated with input data channels before the first convolutional layer. In practice, we observe that CoordConv improves the performance.

Since we deal with several objects of the same nature, we found it reasonable to apply the decoder $D$ to the track feature containers $z_i$ of each $i$-th track separately and then merge the outputs. To this end we first process each of the containers $z_i$ with the same decoder network that outputs single channel map $Y_i = D(z_i')$ corresponding to the track $t_i$. Then these images are merged into the final output $Y$ such that for each pixel at coordinates $[k, l]$ in the image $Y_{[k,l]} = \sigma(\max_i Y_{i \, [k,l]})$. Here the sigmoid activation function $\sigma(x) = (1 + e^{-x})^{-1}$ is ensuring that the output pixel values are in the range $[0, 1]$. This not only reduces the number of parameters in the decoder, but also simplifies the study of encoder performance, as each $z_i$ has the same structure. Alternatively, shuffling the containers $z_i$ within each sample could be employed. Details on the encoder and decoder architectures are given in Table 2.

Table 2: Description of the encoder and decoder used in our models. conv(kernel size, dilation, # of channels) - convolution; c_conv - CoordConv, concatenation with 2 channels of $x$ and $y$ coordinates and convolution; AP - average pooling; FC - fully connected layer; c_tconv - transposed CoordConv: tiling input up to the target size, concatenation with 2 coordinate channels, and convolution. All convolutions and FC layers are followed by batch normalization [24] and ReLU [25] activation, unless otherwise stated.

| Encoder | c_conv(3x3, 1, 16), conv(3x3, 2, 16), conv(3x3, 2, 64), conv(3x3, 2, 128), AP(2x2), conv(1x1, 1, 128), conv(3x3, 2, 512), AP(2x2), conv(1x1, 1, 256), FC(1024), FC(512), FC(128), FC($n \cdot n_p$, tanh) |
|---|---|
| Decoder | *Single block:*<br>FC(16), c_tconv(1x1, 32), conv(1x1, 1, 32), conv(1x1, 1, 32), conv(1x1, 1, 16), conv(3x3, 1, 16), conv(3x3, 1, no activation)<br><br>*Blocks merging:*<br>Max(*blocks*), sigmoid activation |

The coefficients in the loss function were chosen empirically to balance the values of its terms: $c_0 = 0.3$, $\alpha_{sig} = 10,000$, $\alpha_{L2} = 300$, $\alpha_{unif} = 2.5 \times 10^5$, and $\alpha_{bin} = 55$. Training starts with $\lambda = 1$, and after $l = 100k$ iterations $\lambda$ is exponentially decreased every 5k iterations by a factor of 0.9, so that $L_{im}$ is relaxed to pure $L^2$ after about 200k iterations:

$$\lambda = \begin{cases} 1, & l < 100k \\ 0.9^{\frac{l-100k}{5k}} & l > 100k \end{cases}$$

We perform the training on 32x32 random crops from 40,000 images of 128x128 pixels until convergence for 500,000 iterations with minibatch of 128 images. All experiments were carried out using TensorFlow 1.12 [26]. Models were trained using the Adam optimizer [27] with initial learning rate of $6 \times 10^{-5}$ to allow for stabilization, rising to $1 \times 10^{-3}$ after 2k iteration was used. Afterwards the rate is decreased by a factor of 0.88 every 90,000 iterations. The loss function over the course of training for e.g. the **AT+TI** model is shown in Figure 5. The loss function evaluated on test dataset (green curve in Fig. 5; see section 3.3 for test dataset details) at training checkpoints confirms that model did not overfit to the training data. The training of each model took about 50 hours on a single GeForce GTX 1080 GPU.

## 3 Results

### 3.1 Autoencoder performance

First, we evaluate whether our models have learned to properly capture the content of presented images in both latent and image spaces. In Figure 6, the comparison of the outputs of the five models and the lines drawn according to the
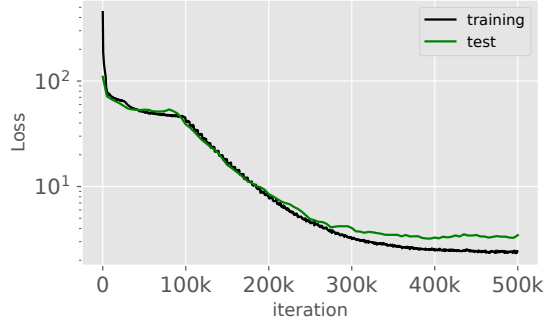
Figure 5: Training and test loss over the course of training of the **AT+TI** model. At iteration 100k the $\lambda$ parameter starts to decrease.

latent representation $z$ predicted for the image, interpreted as described above, are shown. It is clear that both **AT+A** and **AT+TI+A** models were able to build the geometrically meaningful latent representation $z$ in most cases. For the **AT+A** model, which does not employ the translational invariance, the output contains more false detection both in the image output and in the drawn track lines. Also, the image output is significantly less sharp in the beginning of the training for this model. **RT+TI+A** on the other hand did not manage to separate the factors of variation in the desired way, and it took much longer to converge, even just to mimic the desired image output. One can see inconsistency between the image output of the autoencoder and the track lines obtained by the latent representation, meaning that overall it did not grasp the concept of the geometrical space in the desired manner. None of the models properly learned the ability to "switch off" the tracks using the confidence parameter $a_i$. While this parameter is not completely ignored (blue lines in the tracks column in Fig. 6), in most cases the models have learned other ways to disable track parameter containers $z_i$, which are not used to encode lines in the image. These containers simply have geometrical parameters corresponding to lines outside of the image crop range, or have $x, y$ coordinates far away from the image crop center (e.g. **AT+A** and **AT+TI+A** in Fig. 6A). Performance of the **AT+TI** model was comparable or sometimes even better than that of the **AT+TI+A** model. On the downside, without the parameter $a_i$ acting as a regularizer, this model tends to attribute close parameters to several lines (e.g. **AT+TI** in Fig. 6E). The overparametrized models used the $x, y$ position to encode confidence in the track presence by placing them closer or further from the image along the track line (Figure 6E,F). Performance of the **AT, rcs** model is slightly worse than of the **AT+TI** model. The performance of the models clearly degrades when the number of tracks in the image crop is $\geq 4$. We assume that the main reason for this is that these cases were rather underrepresented in the training set.

### 3.2 Disentanglement of the geometrical variational factors

To better understand the learned representation, we have performed a careful dissection into both decoder and encoder in this and the following sections, that was possible since the latent representation was designed to be fully interpretable. We start with the visual analysis of the learned representation by verifying the output of the decoder for given values of $z_i$. To this end, we have run the decoder on the entire range of meaningful values of $z_i$. In addition, for this study we performed the prediction on the image coordinate area $(-1, -1) - (2, 2)$, i.e. 9 times bigger than the range of the original cropped image $(0, 0) - (1, 1)$. This way we can empirically see how well the decoder generalizes to a wider coordinate range. This is possible thanks to the CoordConv nature of the decoder: by changing the values in the coordinate channels we can perform the prediction at any position. In Figure 7 and Supplementary Figure 1 we show comparisons of decoder's predictions for a set of representative values of $z_i$ between all models.

It is clear that all models, which employ the activation parameter, have learned to suppress the output when the values of the parameters $a_i$ are small. The output of the **RT+TI+A** model does not correlate with the expectation at all, and while the output resembles lines, the learned representation is clearly not the desired one. It would be interesting to investigate which representation was found but this is outside the scope of this paper. Models employing translational invariance produced more elongated lines that fade out slower compared to the **AT+A** model (compare, for example, **AT+A** vs **AT+TI+A** and **AT+TI** in Fig. 7, rows 2-5). **AT+TI** shows even more pronounced and fine lines. All overparametrized models (i.e. all except **AT, rcs**) also suppress tracks with an $x, y$ position lying further from the image range (Fig. 7, top row).
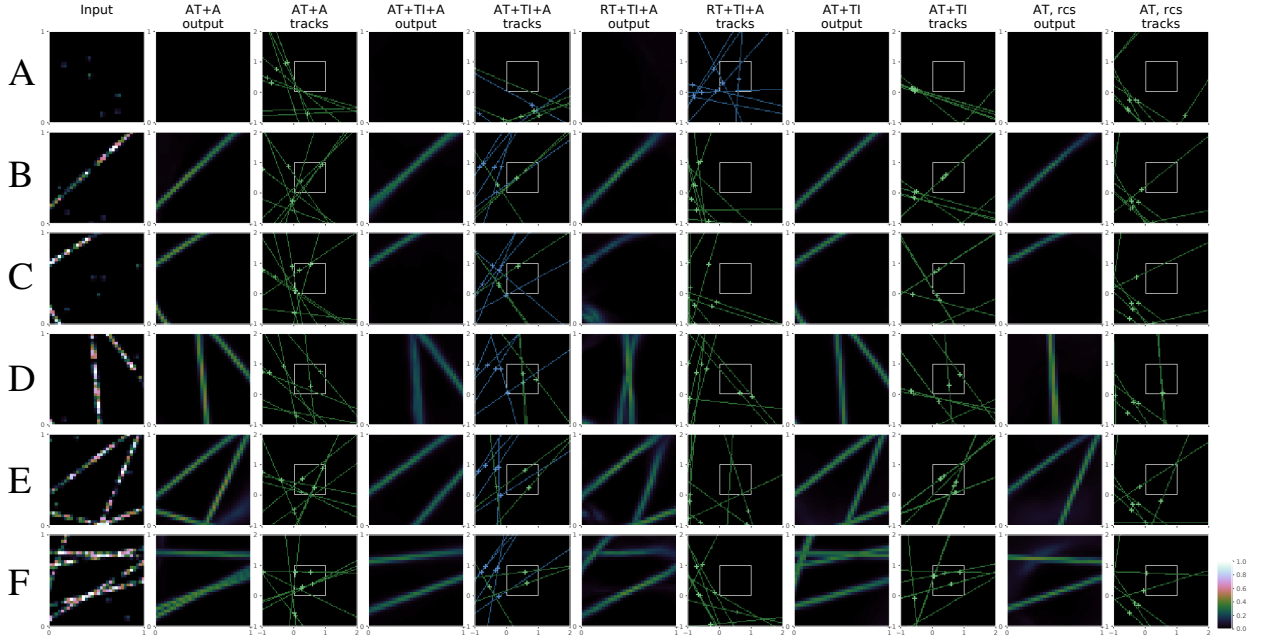
7

Figure 6: Autoencoder performance. A. No tracks in the image. B. One track in the image. C. Two tracks in the image. D. Three tracks in the image. Notice third track in left bottom corner, successfully detected by the **AT+TI** model. E. Four tracks in the image. Models have learned to represent confidence in track presence as distance from view center rather than the latent variable $a$. F. With $\geq 4$ tracks in the image the models start to partially fail.

Leftmost column shows input images. For each model the image output of the decoder and the lines drawn according to the latent representation $z$ are shown. In the input and output columns color depicts brightness. In track columns: green and blue lines correspond to enabled ($a > 0$) and disabled ($a \leq 0$) tracks; for models without parameter $a$ all tracks are shown in green; crosses of corresponding color show the predicted $x, y$ position; white frame shows span of the input image. (Best seen in electronic version.)
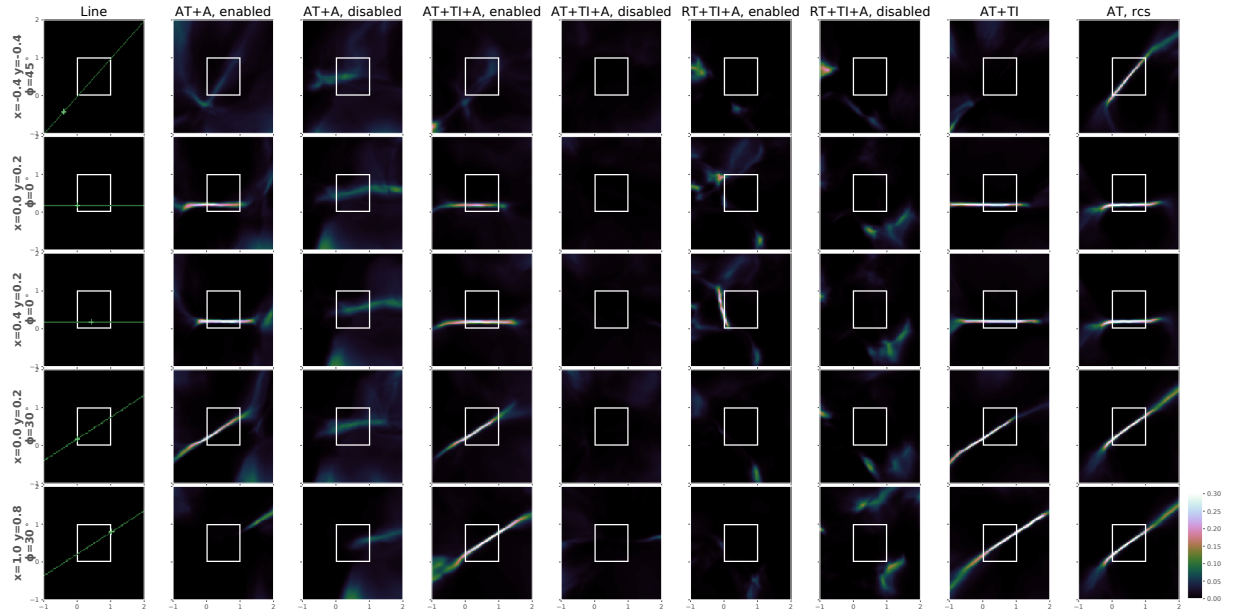


Figure 7: Comparison of decoder outputs for all models for given latent variable values. Tracks corresponding to the given values are shown in leftmost column. For disabled and enabled tracks, the activation parameter is set to $a = 0$ and $a = 1$ correspondingly. (Best seen in electronic version.)
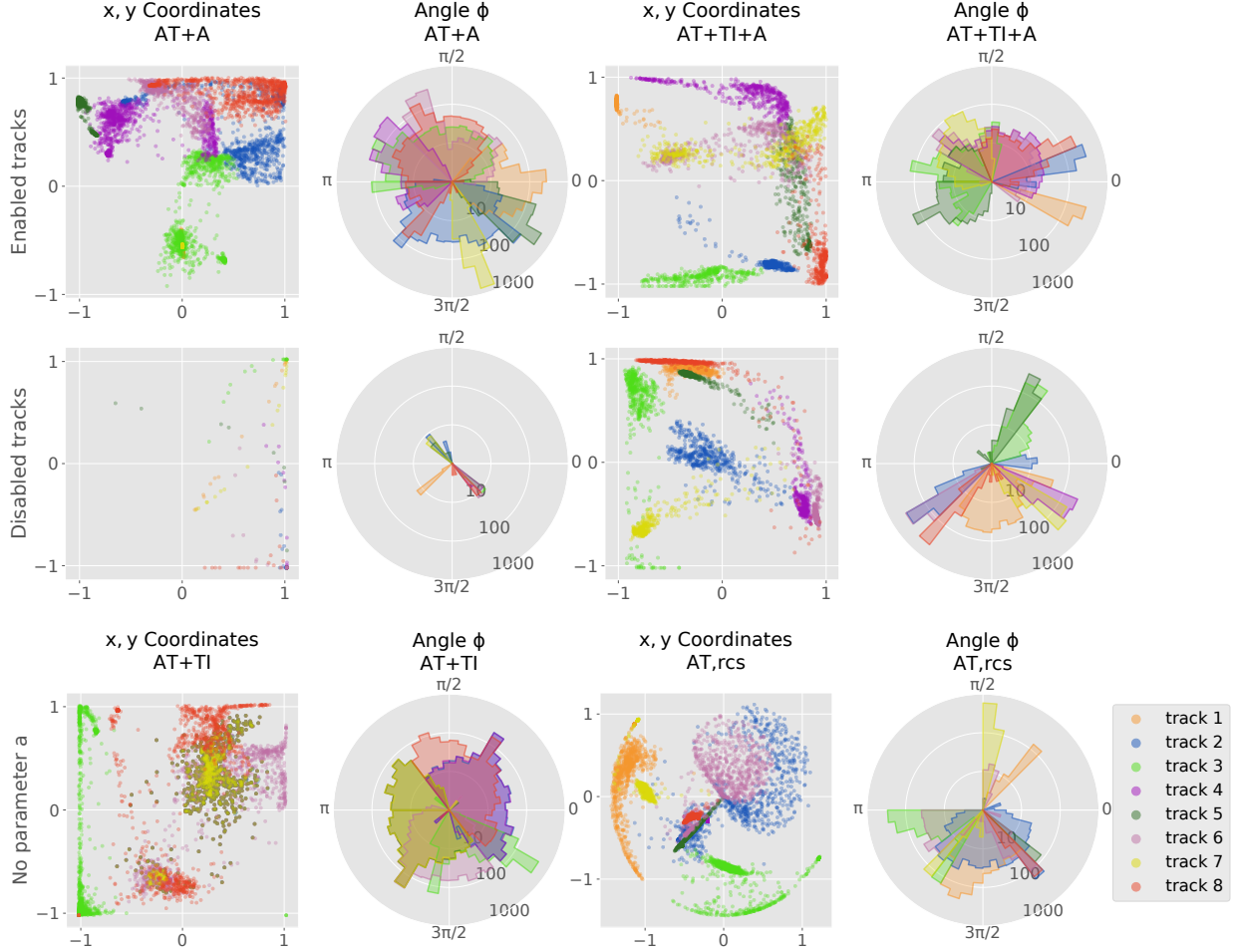
Figure 8: Distribution of the $x, y$ coordinates and inclination angle $\phi$ for the enabled (top row) and disabled (middle row) tracks for the **AT+A** and **AT+TI+A** models and the models without activation parameter **AT+TI**, **AT, rcs** (bottom row). Interpreted geometrical values for the eight track parameter containers are shown in different colors. Values for each container tend to cluster, covering a particular space region.

## 3.3 Performance of the track parameters' measurement

To study the performance of the model for tracking, we evaluate the distribution and resolution of the encoder outputs $z_i$. In Figure 8, the distributions of the predicted $x_i, y_i$ positions and the angle $\phi_i$ obtained from the $c_i, s_i$ parameters for each of the eight track feature containers, $z_i$, is shown for **AT+A**, **AT+TI+A**, **AT+TI**, and **AT, rcs** models. Since the latter has a different representation, for consistency we obtain the values of $x_i, y_i, \phi_i$ using the two-argument arctangent function as follows:

$$\theta_i = \arctan2(s_i, c_i)$$

$$x_i = r_i \cos\theta_i; y_i = r_i \sin\theta_i; \phi_i = \theta_i - \pi/2$$

We skip further studies of the **RT+TI+A** model, since the latent variables do not have the desired meaning, driving the geometrical analysis meaningless. We show these distributions separately for "enabled" and "disabled" tracks, according to the latent activation parameter $a_i$, where applicable. One can clearly see again that the models did not learn to use the parameter $a_i$, and e.g. the **AT+A** assigned almost all tracks the "enabled" value of the activation parameter. Instead, the reconstructed parameters for containers, which do not correspond to any tracks in the image, have rather localized $x, y$ positions and angle (see the peaks in the angular distribution, present for each of the eight track containers). The positions for existing tracks lie within or close to the coordinate region of the image $(0, 0) - (1, 1)$. While they tend to cluster for each track container, they rather uniformly cover a wide band in the $x, y$ space. Notably,

combined with a wide angular distribution, this localization does not limit the sensitivity region of the model (see, for example, the parameter distributions of the 8-th container for the enabled tracks in the **AT+A** model in Fig. 8).

In the angular space all directions are covered, leaving no blind spots. Each of the eight parameter containers covers a subspace with some overlap for models **AT+A** and **AT+TI+A**. This means, that only a fraction of the track containers is sensitive to any chosen direction. E.g. the **AT+TI+A** model would fail to detect >3 parallel lines with $70°$ inclination. Angular overlap in the **AT+A** model is very poor leading to poor detection of several parallel tracks in an image crop. In fact, three out of eight output containers have geometrical parameters corresponding to lines outside of the image (Supplementary Figure 2). In **AT+TI**, on the other hand, each container covers almost $\pi$ in angular space, and the overall distribution is rather uniform (See Supplementary Figure 2). This would allow to detect several parallel lines in a view (e.g. in Fig. 6F several tracks have similar angles).

The **AT, rcs** model did not learn to utilize most of the containers. Practically only the containers 2 and 6 learned to encode track lines, as seen in Figure 8. Nevertheless, even these two containers do not cover the whole angular range. For the remaining containers, the tracks lie outside of the image crop region, as seen on the $x, y$ distribution. This distribution is easy to interpret for this model since the track angle can be observed directly from the coordinates. For a circle with center at the origin and passing through some point $x, y$, the track line would be the tangent to the circle at this point. Arguably, the lack of flexibility due to minimal parametrization did not allow this model to efficiently switch off the tracks, leaving 2 almost always enabled and 6 always disabled.

To quantitatively evaluate the performance of the models we have processed the test dataset. This dataset was generated similarly to the training dataset with additional information on ground truth (GT) track positions and angles. It consists of 30,000 images with 5,000 sample images for each of 0, 1, ..., 5 tracks/image conditions.

First, we assign the reconstructed tracks (and active, i.e. $a > 0$ for models with activation parameter $a$) to the GT ones or mark them as fake. We evaluate the distance $\Delta r$ from image center between a predicted track and a GT track, and the difference in angle $\Delta \phi$. Then we build the $\chi^2$ as $\chi^2 = (\frac{\Delta r}{\sigma_r})^2 + (\frac{\Delta \phi}{\sigma_\phi})^2$. Here we use the theoretical position resolution which is defined by pixel size $\sigma_r = \frac{1px}{\sqrt{12}} = \frac{1}{32}\frac{1}{2\sqrt{3}} \approx 0.009 \approx 0.3$ px, and angular resolution defined by pixel size and track length $l$ within the image crop $\sigma_\phi = \frac{\sqrt{2}}{\sqrt{12}l} = \frac{1}{\sqrt{6}l}$, ($\sigma_\phi \approx 13$ mrad for $l = 32$ px).

The assignment is then performed sequentially, by selecting available prediction-GT pairs according to the minimum value of $\chi^2$, if $\chi^2 \leq 11.83$ ($3\sigma$ statistical significance, number of degrees of freedom $ndf = 2$). The remaining predicted tracks are split into two categories, fakes and duplicates. A track is considered as a duplicate if its $\chi^2$ to any of the used GT tracks or assigned prediction tracks is $\chi^2 < 2.3$ (i.e. within $1\sigma$), and as a fake otherwise. For the assigned tracks, we then evaluate the offsets $\Delta r$, $\Delta \phi$ as a function of number of tracks on the original image, as well as the fraction of assigned tracks (efficiency), number of fake tracks, and number of duplicate tracks (Figure 9). The actual coordinate and angular resolutions $\sigma_{r,mod}$, $\sigma_{\phi,mod}$ of the models can be estimated from these data as mean values of $\Delta r$ and $\Delta \phi$. For example, for the **AT+TI** model $\sigma_{r,mod} \approx 0.013 \approx 0.42$ px, $\sigma_{\phi,mod} \approx 15$ mrad for $l = 32$ px).

We then use mean resolutions for each model to show the $\chi^2 = (\frac{\Delta r}{\sigma_{r,mod}})^2 + (\frac{\Delta \phi}{\sigma_{\phi,mod}})^2$ distribution for these models for different numbers of tracks per image crop in Supplementary Figure 3. While for the **AT+A**, **AT+TI**, and **AT, rcs** models the distributions are consistent with a $\chi^2$ distribution with $ndf = 2$, for the **AT+TI+A** model, the peak is smeared and shifted towards higher values, consistent with higher resolution variance especially in the high track density region.

The resolution of the models is stable as the number of tracks grows. The **AT+TI** model has consistently higher resolution, as well as higher efficiency. Efficiency significantly decreases in all models with increasing number of tracks per image. We argue that it is caused by the fact that images with high track number were underrepresented in the training set (Figure 4C), and the efficiency would improve if a training set with high track multiplicity images were used. To support this claim, we have generated a training dataset of 60,000 images with a uniform distribution of track density, i.e. 10,000 images for 0, 1, ..., 5 tracks/image crop. We have then retrained the **AT+TI** model on this dataset. This has significantly (>20%) improved the efficiency for high track density (**AT+TI/U** model in Fig. 9). While the number of fake tracks is similar in the **AT+TI** model, since it lacks the regularization based on the latent parameter $a$, it tends to assign all of the available containers to tracks. This leads in turn to a larger number of duplicates. Nevertheless, this effect is suppressed for the **AT+TI/U** model, trained on the dataset with uniform track number representation. Another model without the activation parameter $a$ – the **AT, rcs** model does not produce many duplicates, most likely since each track is sensitive only to a narrow angular range, as shown above.

Finally, we processed a real emulsion dataset to qualitatively observe the performance of the **AT+TI/U** model (trained on synthetic data) in processing real experimental data. We used a single image out of a 3D tomographic image stack of size $640 \times 512$ pixels corresponding to $190 \times 150\,\mu m$ of emulsion detector area, irradiated with 400 GeV protons at
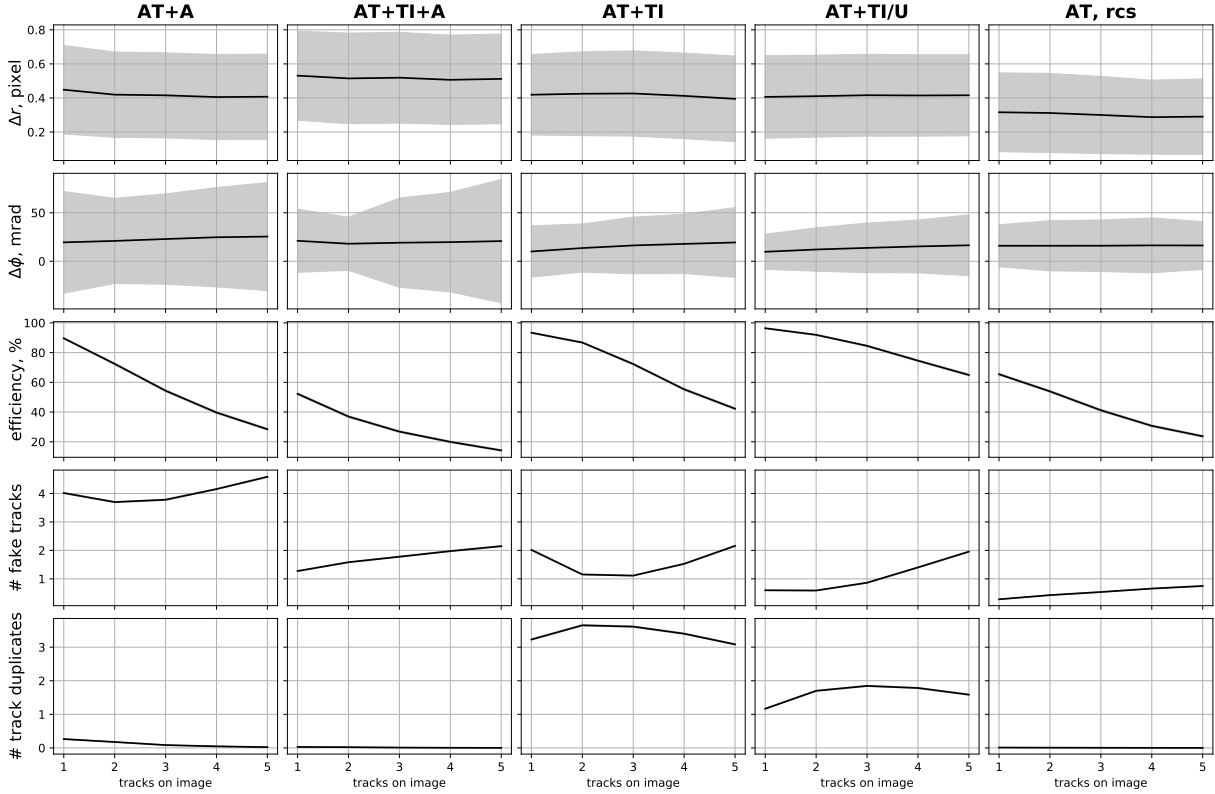
Figure 9: Tracking precision and efficiency as function of the number of tracks in the image for the five model configurations. From top to bottom: distance between GT track and the reconstructed one (less is better); angular difference (less is better); fraction of reconstructed tracks (more is better); number of reconstructed tracks that have no corresponding GT track, per image (less is better); number of reconstructed tracks that duplicate another reconstructed track, per image (less is better). Gray fill shows error range of 1 standard deviation of distribution for resolution parameters and of mean value for efficiency, fake tracks, and duplicates.

different angles at the SPS accelerator beam at CERN [28]. We preprocessed the image (Figure 10A) by downscaling it by a factor of four, inverting the image, and normalizing the color scale to match the training data properties (Figure 10B). We then divided the image into $5 \times 4$ non-overlapping $32 \times 32$ pixel crops and processed them independently. The resulting tracks are then assembled into the full image size and shown as 32 pixels long segments with highlighted $x, y$ position (Figure 10B, overlay). Even though the models were trained on synthetic data with different signal and noise distributions from experimental data, one can appreciate the agreement between real tracks and predicted ones (e.g. tracks pointed to by arrowheads), confirming the effectiveness of our method. For real-life applications, the model must be trained on the raw experimental data, to learn the true signal and noise distributions from it.

## 4 Discussion and outlook

Disentangling factors of variation remains a hot topic for several years in representation learning research. Moreover, developing models that are capable of abstracting high-level concepts from raw data can lead to plenty of direct practical applications. Many works approached this problem by employing variational autoencoders with regularizations in the latent representation that enforce disentangling [16] or autoencoders combined with adversarial training [22]. In most cases, after disentangling the factors of variation, a few labeled samples can be used to associate the factors with interpretable measures in a quantitative way. Previously it has been shown [29] how a few labeled samples can improve disentanglement itself. In another work [30] it was shown that applying an equivariance constraint, i.e. changing one factor of variation, corresponding to a change of one dimension of the disentangled representation in a predictable manner, leads to disentangled variables. Nevertheless, to the best of our knowledge no previous works have tried to extract meaningful quantitative information in a fully unsupervised manner.
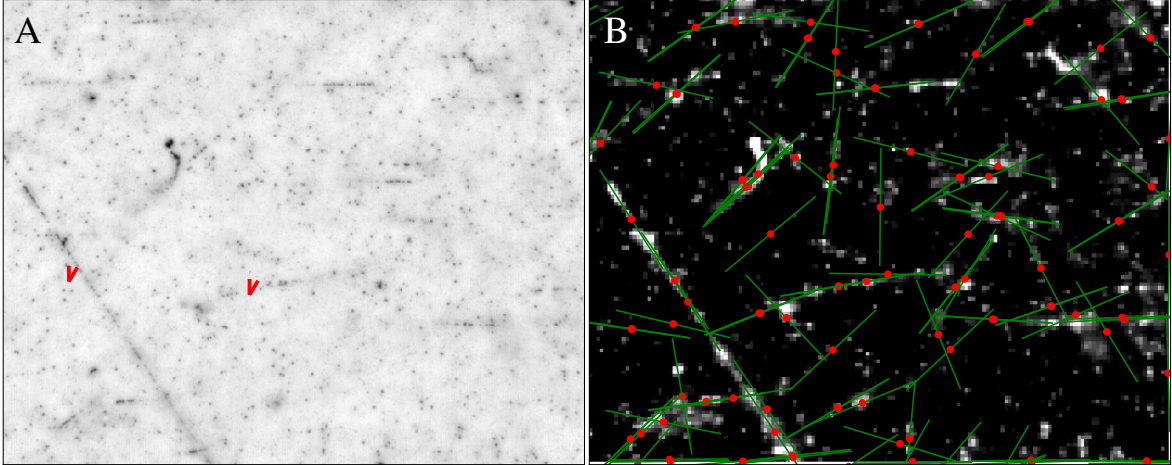
11

Figure 10: Tracking particles in emulsion data with the **AT+TI/U** model. A. Wide field microscopy image of an emulsion detector. The image corresponds to a $380\mu m \times 300\mu m$ detector surface. B. A downscaled and inverted image, overlaid with track segments (32 pixels long) reconstructed with the **AT+TI/U** model. Examples of tracks reconstructed consistently with our model are highlighted by arrowheads.

In this work, we have demonstrated that imposing equivariance constraints on the autoencoder under geometrical transformations in the image and latent representation domains enables the model to "discover" the existence of multiple lines in the presented images in a fully unsupervised manner. Incorporating simple affine transformation such as translation, rotation, scaling and skew as equivariances between the image and latent spaces allows the models to successfully disentangle the factors of variation in the image data into geometrically meaningful parameters (coordinates and angles of lines). Adding the possibility to "switch off" a predicted track with an additional activation parameter $a$ does not drastically change the results (models **AT+A** and **AT+TI+A**). While it can help to prevent the shortcut problem and reduces the number of track duplicates, these models did not learn to exploit it.

Incorporating the translation along the track line in addition to the whole set of affine transformations enforces the line detection. However, employing only the translational invariance, even together with rotational transformations (**RT+TI+A** model), leads to reference ambiguity so that latent parameters do not correspond to the desired geometrical variables. We believe that with a few calibration measurements it would be possible to find a mapping from these latent parameters to the desired geometrical variables, but this lies out of the scope of this work. As we have shown, a larger set of transformations allows the model to immediately learn the latent representation in an unambiguous, geometrically meaningful way. The minimal subset of affine transformations, sufficient for disentangling the factors of variation without reference ambiguity, will be explored in future works.

In addition to the coordinate-angle parametrization, which gives the models more freedom in sample exploration, we have studied the classical rho-theta parametrization. Under the set of affine transformations, this model is also able to learn a meaningful parametrization in a fully unsupervised manner. Yet the model performance is slightly worse than, for example, the **AT+TI** model, most likely because the incorporated CoordConv approach prefers to have a natural $x - y$ coordinate representation.

The main weak point of our current implementation is that neither the background nor the grain distribution along the lines are in any way represented by the current models, which may impair line detection in the case of high background rate. In addition, the employed transformations in the image domain affect the image parameter distributions, such as brightness (corresponding to the $dE/dx$ energy loss in emulsion detector) and sharpness. Models with additional global and per-track latent parameters without an *a priori* assigned meaning would naturally overcome these hurdles. Training them would require dealing with the shortcut problem in these parameters, and thus would benefit from employing the adversarial framework [22], [31]. While the aim of this work was to carefully study the proposed approach in general, we leave this aspect to further studies.

We expect this approach to have a large potential in the analysis and extraction of geometrical properties from image data. In further work we plan to adapt this technique to the location of tracks in full resolution 3D tomographic microscopy data, or data from Liquid Argon Time Projection Chamber detectors [32], [33], which would be a direct extension of this approach. Also adding more samples with a higher track number in the training dataset is expected to improve the efficiency and resolution in cases with high track density.

While designed to detect simple line structures, this technique has the potential to be used for locating and parametrizing other objects, such as splines. This would enable the tracking of particles in magnetic fields and pave the way to novel automated image vectorization techniques. Being fully unsupervised, this approach can leverage all available raw dataset with no extra work required.

## Acknowledgement

## References

[1]   Niwa K, K. Hoshino, and K. Niu, "Auto scanning and measuring system for the emulsion chamber," in *International Cosmic ray Symposium of High Energy Phenomena*, Tokyo, 1974, p. 149.

[2]   A. Alexandrov, A. Buonaura, L. Consiglio, *et al.*, "The Continuous Motion Technique for a New Generation of Scanning Systems," *Scientific Reports*, vol. 7, no. 1, p. 7310, Dec. 2017, ISSN: 2045-2322. DOI: 10.1038/s41598-017-07869-3. [Online]. Available: http://www.nature.com/articles/s41598-017-07869-3.

[3]   R. Acciarri, C. Adams, R. An, *et al.*, "Convolutional neural networks applied to neutrino events in a liquid argon time projection chamber," *Journal of Instrumentation*, vol. 12, no. 03, P03011–P03011, Mar. 2017, ISSN: 1748-0221. DOI: 10.1088/1748-0221/12/03/P03011. [Online]. Available: http://stacks.iop.org/1748-0221/12/i=03/a=P03011?key=crossref.b0d7d7efe193e427b920796770b371c9.

[4]   K. Kodama, N. Ushida, C. Andreopoulos, *et al.*, "Observation of tau neutrino interactions," *Physics Letters B*, vol. 504, no. 3, pp. 218–224, Apr. 2001, ISSN: 03702693. DOI: 10.1016/S0370-2693(01)00307-0. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0370269301003070?via%7B%5C%%7D3Dihub%20https://linkinghub.elsevier.com/retrieve/pii/S0370269301003070.

[5]   N. Agafonova, A. Aleksandrov, A. Anokhina, *et al.*, "Discovery of $\tau$ Neutrino Appearance in the CNGS Neutrino Beam with the OPERA Experiment," *Physical Review Letters*, vol. 115, no. 12, p. 121 802, Sep. 2015, ISSN: 0031-9007. DOI: 10.1103/PhysRevLett.115.121802. [Online]. Available: https://link.aps.org/doi/10.1103/PhysRevLett.115.121802.

[6]   C. Ahdida, R. Albanese, A. Alexandrov, *et al.*, "Sensitivity of the SHiP experiment to Heavy Neutral Leptons," *Journal of High Energy Physics*, vol. 2019, no. 4, p. 77, Apr. 2019, ISSN: 1029-8479. DOI: 10.1007/JHEP04(2019)077. arXiv: 1811.00930. [Online]. Available: http://arxiv.org/abs/1811.00930%20http://link.springer.com/10.1007/JHEP04(2019)077.

[7]   S. Aghion, A. Ariga, M. Bollani, *et al.*, "Nuclear emulsions for the detection of micrometric-scale fringe patterns: an application to positron interferometry," *Journal of Instrumentation*, vol. 13, no. 05, P05013–P05013, May 2018, ISSN: 1748-0221. DOI: 10.1088/1748-0221/13/05/P05013. [Online]. Available: http://stacks.iop.org/1748-0221/13/i=05/a=P05013?key=crossref.0ead6fcf419f49c586e07afa1d88567a.

[8]   R. Nishiyama, A. Ariga, T. Ariga, *et al.*, "First measurement of ice-bedrock interface of alpine glaciers by cosmic muon radiography," *Geophysical Research Letters*, vol. 44, no. 12, pp. 6244–6251, Jun. 2017, ISSN: 00948276. DOI: 10.1002/2017GL073599. [Online]. Available: http://doi.wiley.com/10.1002/2017GL073599.

[9]   A. Ariga, T. Ariga, A. Ereditato, *et al.*, "A Nuclear Emulsion Detector for the Muon Radiography of a Glacier Structure," *Instruments*, vol. 2, no. 2, p. 7, Apr. 2018, ISSN: 2410-390X. DOI: 10.3390/instruments2020007. [Online]. Available: http://www.mdpi.com/2410-390X/2/2/7.

[10]  N. Armenise, M. De Serio, M. Ieva, *et al.*, "High-speed particle tracking in nuclear emulsion by last-generation automatic microscopes," *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, vol. 551, no. 2-3, pp. 261–270, Oct. 2005, ISSN: 01689002. DOI: 10.1016/j.nima.2005.06.072. [Online]. Available: http://linkinghub.elsevier.com/retrieve/pii/S0168900205013240%20https://linkinghub.elsevier.com/retrieve/pii/S0168900205013240.

[11]  A. Ariga and T. Ariga, "Fast $4\pi$ track reconstruction in nuclear emulsion detectors based on GPU technology," *Journal of Instrumentation*, vol. 9, no. 04, P04002–P04002, Apr. 2014, ISSN: 1748-0221. DOI: 10.1088/1748-0221/9/04/P04002. [Online]. Available: http://stacks.iop.org/1748-0221/9/i=04/a=P04002?key=crossref.51c05f33a6f743cf8adde0c7dc2f3253.

[12]  M. Yoshimoto, T. Nakano, R. Komatani, and H. Kawahara, "Hyper-track selector nuclear emulsion readout system aimed at scanning an area of one thousand square meters," *Progress of Theoretical and Experimental Physics*, vol. 2017, no. 10, Oct. 2017, ISSN: 2050-3911. DOI: 10.1093/ptep/ptx131. [Online]. Available: http://academic.oup.com/ptep/article/doi/10.1093/ptep/ptx131/4554499/Hypertrack-selector-nuclear-emulsion-readout.

[13] A. Radovic, M. Williams, D. Rousseau, *et al.*, "Machine learning at the energy and intensity frontiers of particle physics," *Nature*, vol. 560, no. 7716, pp. 41–48, Aug. 2018, ISSN: 0028-0836. DOI: 10.1038/s41586-018-0361-2. [Online]. Available: http://www.nature.com/articles/s41586-018-0361-2.

[14] C. Adams, M. Alrashed, R. An, *et al.*, "Deep neural network for pixel-level electromagnetic particle identification in the MicroBooNE liquid argon time projection chamber," *Physical Review D*, vol. 99, no. 9, p. 092 001, May 2019, ISSN: 2470-0010. DOI: 10.1103/PhysRevD.99.092001. arXiv: 1903.05663. [Online]. Available: http://arxiv.org/abs/1903.05663%20https://link.aps.org/doi/10.1103/PhysRevD.99.092001.

[15] L. Dominé and K. Terao, "Scalable Deep Convolutional Neural Networks for Sparse, Locally Dense Liquid Argon Time Projection Chamber Data," Mar. 2019. arXiv: 1903.05663. [Online]. Available: http://arxiv.org/abs/1903.05663.

[16] A. Kumar, P. Sattigeri, and A. Balakrishnan, "Variational Inference of Disentangled Latent Concepts from Unlabeled Observations," Nov. 2017. arXiv: 1711.00848. [Online]. Available: http://arxiv.org/abs/1711.00848.

[17] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, Mar. 1998, ISSN: 00189219. DOI: 10.1109/5.726791. arXiv: 1703.07980. [Online]. Available: http://ieeexplore.ieee.org/document/726791/%20http://arxiv.org/abs/1703.07980.

[18] F. Li, H. Qiao, and B. Zhang, "Discriminatively boosted image clustering with fully convolutional auto-encoders," *Pattern Recognition*, vol. 83, pp. 161–173, Nov. 2018, ISSN: 00313203. DOI: 10.1016/j.patcog.2018.05.019. arXiv: 1703.07980. [Online]. Available: https://linkinghub.elsevier.com/retrieve/pii/S0031320318301936.

[19] R. O. Duda and P. E. Hart, "Use of the Hough transformation to detect lines and curves in pictures," *Communications of the ACM*, vol. 15, no. 1, pp. 11–15, Jan. 1972, ISSN: 00010782. DOI: 10.1145/361237.361242. [Online]. Available: http://portal.acm.org/citation.cfm?doid=361237.361242.

[20] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," May 2015. arXiv: 1505.04597. [Online]. Available: http://arxiv.org/abs/1505.04597.

[21] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial Transformer Networks," *Advances in Neural Information Processing Systems*, vol. 2015-Janua, pp. 2017–2025, Jun. 2015, ISSN: 10495258. arXiv: 1506.02025. [Online]. Available: http://arxiv.org/abs/1506.02025.

[22] A. Szabó, Q. Hu, T. Portenier, M. Zwicker, and P. Favaro, "Understanding Degeneracies and Ambiguities in Attribute Transfer," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11209 LNCS, Springer Verlag, 2018, pp. 721–736, ISBN: 9783030012274. DOI: 10.1007/978-3-030-01228-1_43. [Online]. Available: http://link.springer.com/10.1007/978-3-030-01228-1%7B%5C_%7D43.

[23] R. Liu, J. Lehman, P. Molino, *et al.*, "An Intriguing Failing of Convolutional Neural Networks and the CoordConv Solution," *Advances in Neural Information Processing Systems*, vol. 2018-Decem, pp. 9605–9616, Jul. 2018, ISSN: 10495258. arXiv: 1807.03247. [Online]. Available: http://arxiv.org/abs/1807.03247.

[24] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," Feb. 2015. arXiv: 1502.03167. [Online]. Available: http://arxiv.org/abs/1502.03167.

[25] V. Nair and G. E. Hinton, "Rectified Linear Units Improve Restricted Boltzmann Machines," in *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ser. ICML'10, USA: Omnipress, 2010, pp. 807–814, ISBN: 978-1-60558-907-7. [Online]. Available: http://dl.acm.org/citation.cfm?id=3104322.3104425.

[26] Martín Abadi, Ashish Agarwal, Paul Barham, *et al.*, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, Software available from tensorflow.org. [Online]. Available: http://tensorflow.org/.

[27] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*, Dec. 2014. arXiv: 1412.6980. [Online]. Available: http://arxiv.org/abs/1412.6980.

[28] S. Aoki, A. Ariga, T. Ariga, *et al.*, "DsTau: Study of tau neutrino production with 400 GeV protons from the CERN-SPS," Jun. 2019. arXiv: 1906.03487. [Online]. Available: http://arxiv.org/abs/1906.03487.

[29] F. Locatello, M. Tschannen, S. Bauer, *et al.*, "Disentangling Factors of Variation Using Few Labels," May 2019. arXiv: 1905.01258. [Online]. Available: http://arxiv.org/abs/1905.01258.

[30] H. Kim and A. Mnih, "Disentangling by Factorising," *35th International Conference on Machine Learning, ICML 2018*, vol. 6, pp. 4153–4171, Feb. 2018. arXiv: 1802.05983. [Online]. Available: http://arxiv.org/abs/1802.05983.

[31] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, *et al.*, "Generative Adversarial Networks," Jun. 2014. arXiv: `1406.2661`. [Online]. Available: `http://arxiv.org/abs/1406.2661`.

[32] C. Adams, M. Alrashed, R. An, *et al.*, "First measurement of $\nu_\mu$ charged-current $\pi^0$ production on argon with the MicroBooNE detector," *Physical Review D*, vol. 99, no. 9, p. 091 102, May 2019, ISSN: 2470-0010. DOI: `10.1103/PhysRevD.99.091102`. arXiv: `1811.02700`. [Online]. Available: `http://arxiv.org/abs/1811.02700%20https://link.aps.org/doi/10.1103/PhysRevD.99.091102%20http://dx.doi.org/10.1103/PhysRevD.99.091102`.

[33] D. Brailsford, "DUNE: Status and Perspectives," Apr. 2018. arXiv: `1804.04979`. [Online]. Available: `http://arxiv.org/abs/1804.04979`.

## A  Affine transformations in latent representations

In this section, we describe the implementation of the affine transformations for the employed $(x, y, c, s)$ and $(r, c, s)$ representations. Additionally, we provide the parametrization of the transformations themselves, as well as the range of parameters used during training in this study. We produce the transformation functions in the latent and image spaces as a combination of rotation, scaling, skew, and shift. First, we define the corresponding transformation operations of the track line parameters in the following way (index $i$ is omitted for brevity). Rotation:

$$T_{rot}(z_t|\xi_r) = (x, y, \cos(\phi + \xi_r), \sin(\phi + \xi_r)); \phi = \arctan2(s, c); \xi_r \in (-\frac{\pi}{4}, \frac{\pi}{4})$$

where $\arctan2(y, x)$ is the two-argument arctangent function. Scaling:

$$T_{scale_x}(z_t|\xi_{sc_x}) = (\xi_{sc_x} x, y, \frac{\xi_{sc_x} c}{\epsilon + \sqrt{(\xi_{sc_x} c)^2 + s^2}}, \frac{s}{\epsilon + \sqrt{(\xi_{sc_x} c)^2 + s^2}}); \xi_{sc_x} \in (0.7, 1.3); \epsilon = 10^{-6}$$

$$T_{scale_y}(z_t|\xi_{sc_y}) = (x, \xi_{(sc_y)} y, \frac{c}{\epsilon + \sqrt{c^2 + (\xi_{sc_y} s)^2}}, \frac{\xi_{sc_y} s}{\epsilon + \sqrt{c^2 + (\xi_{sc_y} s)^2}}); \xi_{sc_y} \in (0.7, 1.3); \epsilon = 10^{-6}$$

Skew:
$$T_{skew_x}(z_t|\xi_{sk_x}) = (x + \xi_{sk_x} y, y, \cos\phi, \sin\phi); \phi = \arctan2(s, c + \xi_{sk_x} s); \xi_{sk_x} \in (-0.4, 0.4)$$
$$T_{skew_y}(z_t|\xi_{sk_y}) = (x, y + \xi_{sk_y} x, \cos\phi, sin\phi); \phi = \arctan2(s + \xi_{sk_y} s, c); \xi_{sk_y} \in (-0.4, 0.4)$$

Translation:
$$T_{trans}(z_t|\xi_{t_x}, \xi_{t_y}) = (x + \xi_{t_x}, y + \xi_{t_y}, c, s), \xi_{t_x}), \xi_{t_x} \in (-0.4, 0.4)$$

Rotation is performed around the coordinate origin; scaling $x$ and scaling $y$ preserves the $y$ and $x$ coordinates intact; skew $x$ and skew $y$ preserves points on the $x$ and $y$ axes correspondingly. The employed range of transformations is a trade-off between urging the model to learn the desired representation and preserving most of the original tracks in the image crop after the transformation.

The combined space transformation is then produced by consecutively applying these transformations:

$$z_t' = T_{tr}(z_t|\xi) = T_{trans}(T_{skew_y}(T_{skew_x}(T_{scale_y}(T_{scale_x}(T_{rot}(z_t|\xi_r)|\xi_{sc_x})|\xi_{sc_y})|\xi_{sk_x})|\xi_{sk_y})|\xi_{t_x}, \xi_{t_y}),$$

$$\xi = (\xi_r, \xi_{sc_x}, \xi_{sc_y}, \xi_{sk_x}, \xi_{sk_y}, \xi_{t_x}, \xi_{t_y})$$

In some models, we add an additional transformation of the track parameters $z_t$, corresponding to translational invariance (TI) along the line. This is implemented as a random shift of the $x, y$ parameters along the line:

$$T_{t.i.}(z_t|r) = (x + rc, y + rs, c, s); \ r = rand(-0.5, 0.5).$$

We apply these transformations only to the enabled tracks according to the value of $a$. For the disabled tracks, the parameters are set to random values in the $(-1, 1)$ range, enforcing the decoder to learn to ignore disabled tracks:
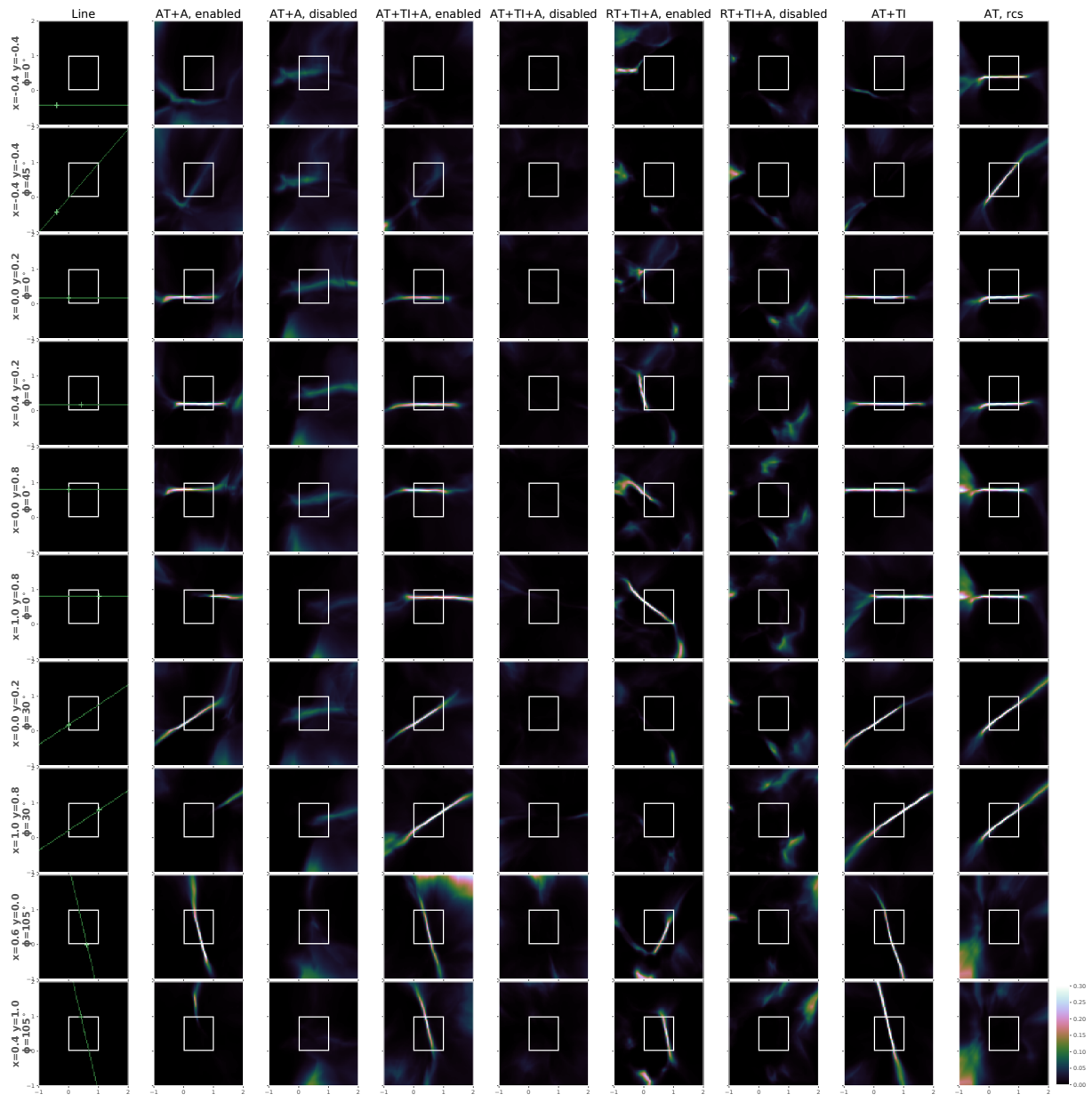
$$z' = \begin{cases} (z_t', \sigma(\gamma a)), & a > 0 \\ (r_1, r_2, r_3, r_4, \sigma(\gamma a)); \ r_i = rand(-1, 1), & a \leq 0 \end{cases}$$

where $\gamma = 20$ and the sigmoid function $\sigma(a) = \frac{1}{1+e^{-a}}$ is applied to the activation parameter $a$ for implementation reasons.
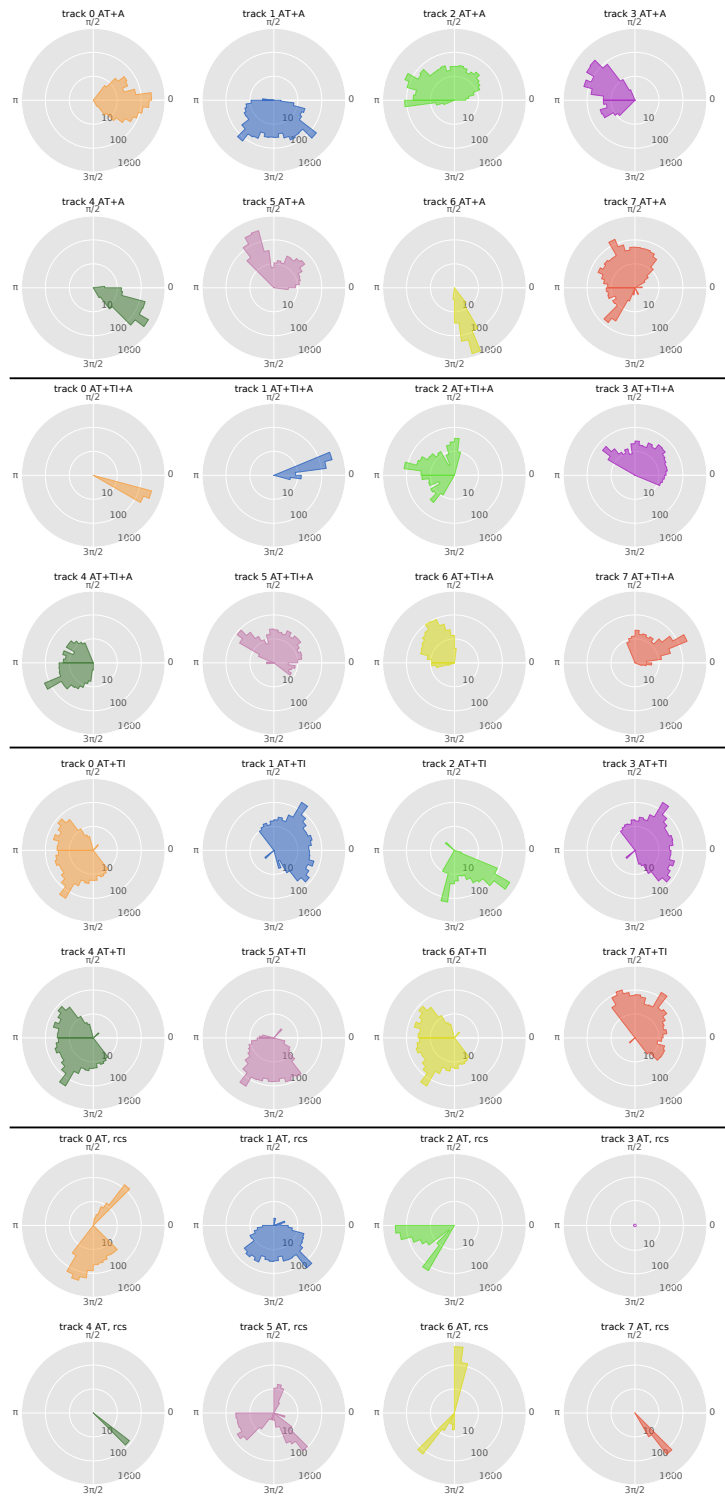
For the model employing the $(r, c, s)$ representation, we first transform the parameters to the $(x, y, c, s)$ representation as: $\theta = \arctan2(s, c)$, $\phi = \theta - \pi/2$, $x = r \cos\theta$; $y = r \sin\theta$; $c = \cos\phi$; $s = \sin\phi$, and then apply the shown above transformations. Afterwards, the inverse transformation to the $(r, c, s)$ representation is applied.

The parameter set $\xi$ is drawn from a uniform random distribution for each training sample on each iteration and is fed into the network along with the images. The input images $I$ of size 96×96 pixels are cropped to 32×32 as shown in Figure 2 and fed as the network input $I_c$. The same input images $I$ are then elastically transformed with the transformation function $I' = T_{im}(I|\xi)$ using the same parameter set $\xi$. For $T_{im}$ we have employed the `tf.contrib.image.transform` function from the TensorFlow library [26]. The origin of the transformations corresponds to pixel coordinates (48,48) in the input image, i.e. the lower bottom corner of the crop. After being cropped to 32×32 pixels, the images are used as network output target $I_c'$. We use larger 96×96 input images to ensure that the final crop after transformation does not contain regions outside of the input image. Examples of the space transformations are shown in the Supplementary Figure 4.
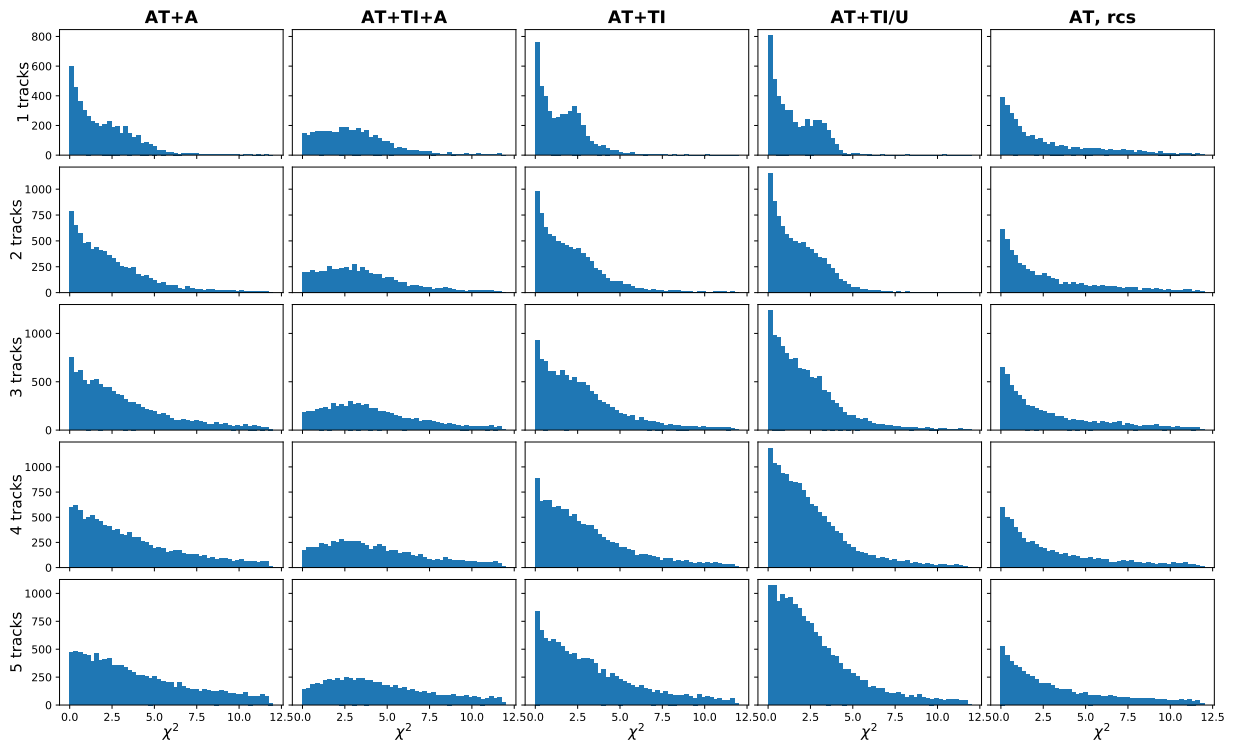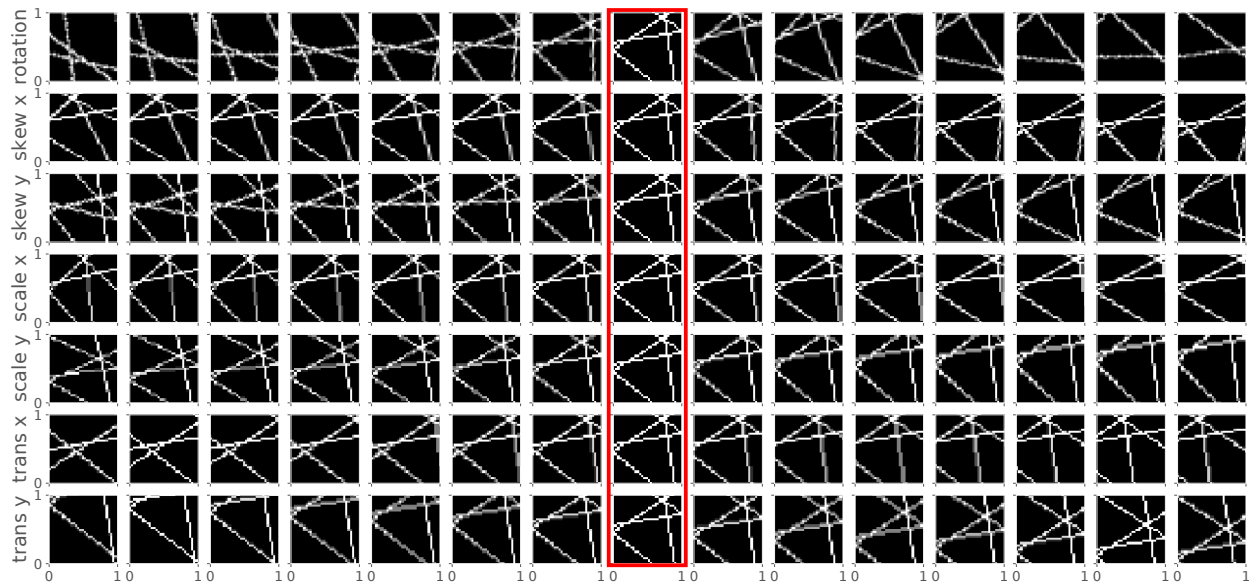
## Supplementary Material



Supplementary Figure 1: Comparison of decoder outputs for all models for several selected latent variable values. True tracks corresponding to the given values are shown in leftmost column.

Supplementary Figure 2: Distributions of the predicted track angles by each of the track feature containers in the latent representation. From top to bottom: **AT+A**, **AT+TI+A**, **AT+TI**, and **AT, rcs** models. Color coding is the same as in Figure 8.

Supplementary Figure 3: $\chi^2$ distributions for the **AT+A**, **AT+TI+A**, **AT+TI**, **AT+TI/U**, and **AT, rcs** models for 1, 2, 3, 4, and 5 tracks per image crop.



Supplementary Figure 4: Examples of image transformations (rows): rotation and skew, scale, and translation along x and y.
The transformations performed with 15 parameters (columns) are presented, linearly distributed in the employed parameter range. Transformations in the middle column correspond to identity transformations (highlighted).