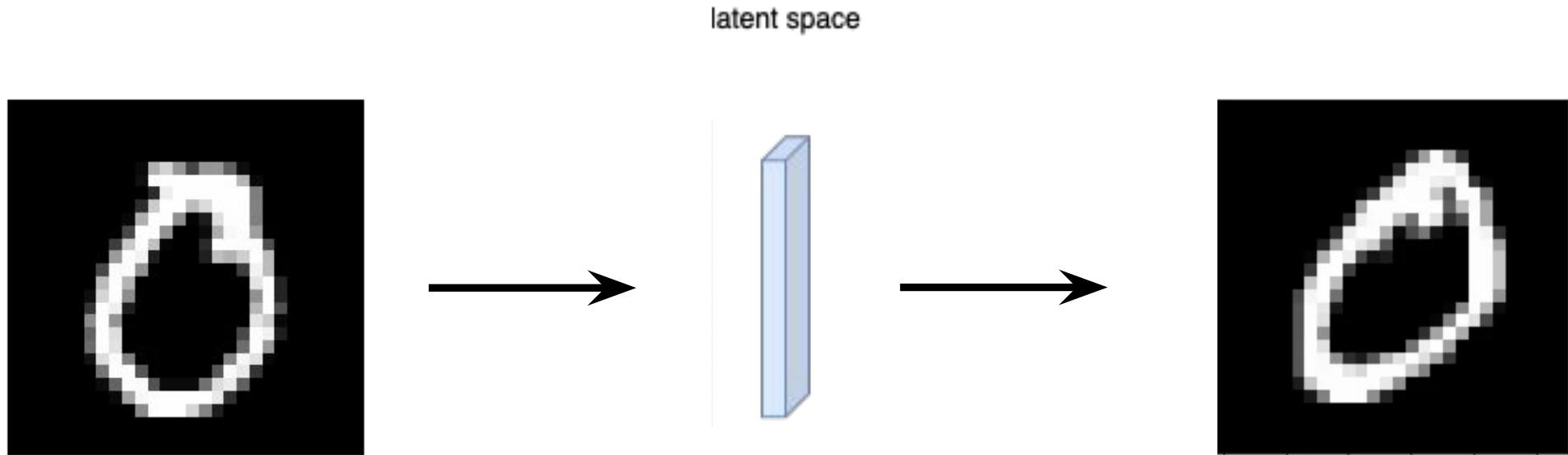


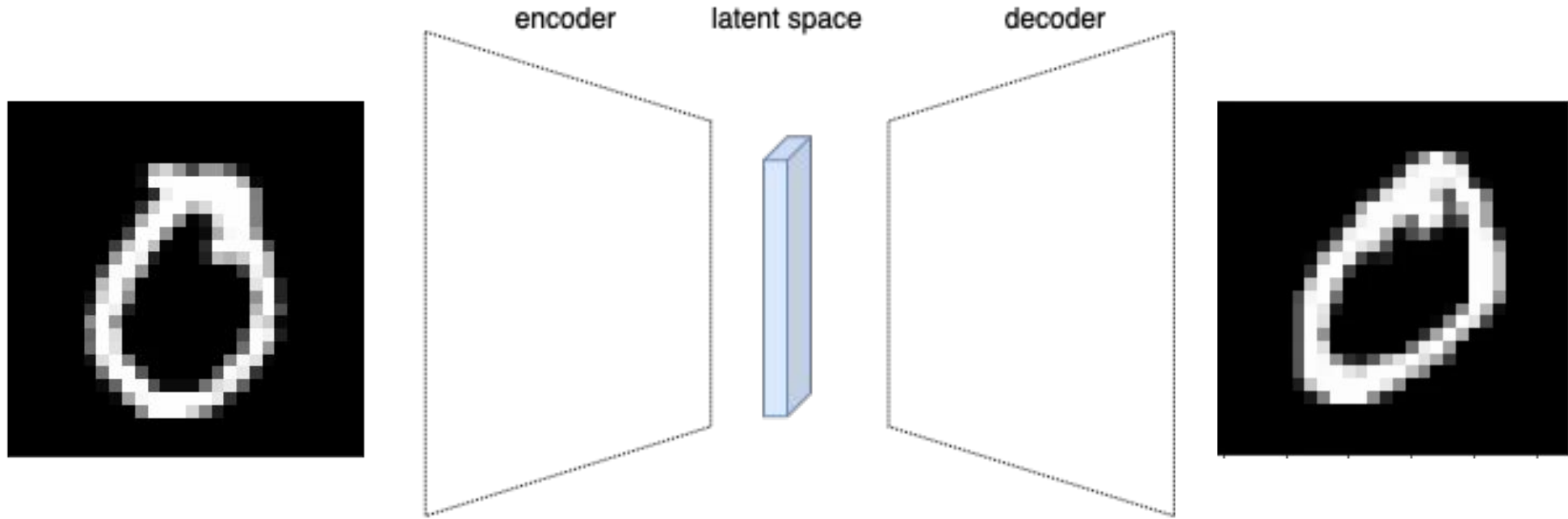
# Flow Factorized Representation Learning

Yue Song, T. Anderson Keller, Nicu Sebe, Max Welling

# Representation learning

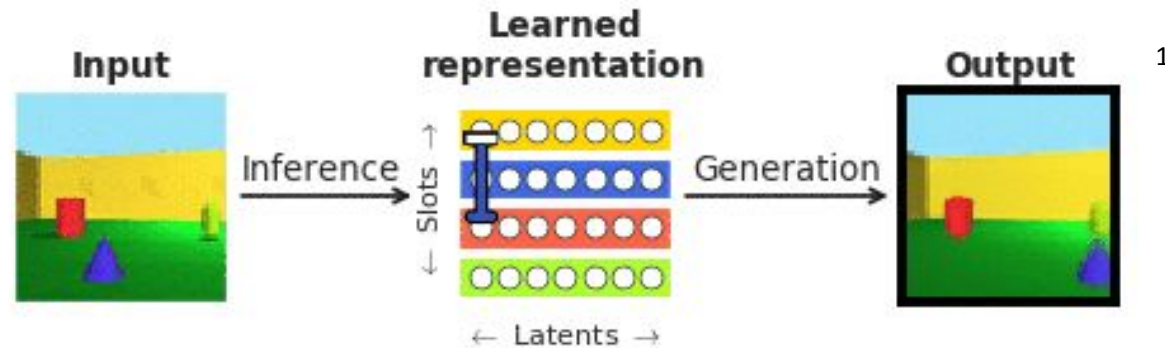


# Representation learning



$$\text{ELBO}_i(\lambda) = \mathbb{E}_{q_\lambda(z|x_i)}[\log p(x_i|z)] - D_{KL}(q_\lambda(z|x_i)||p(z))$$

# Disentanglement



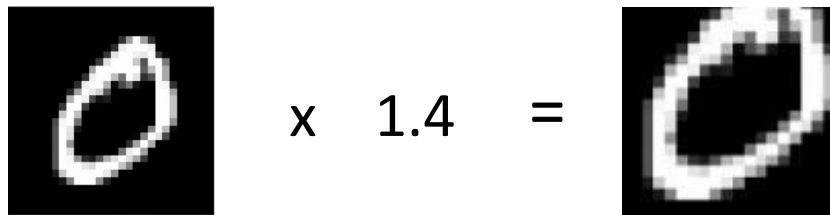
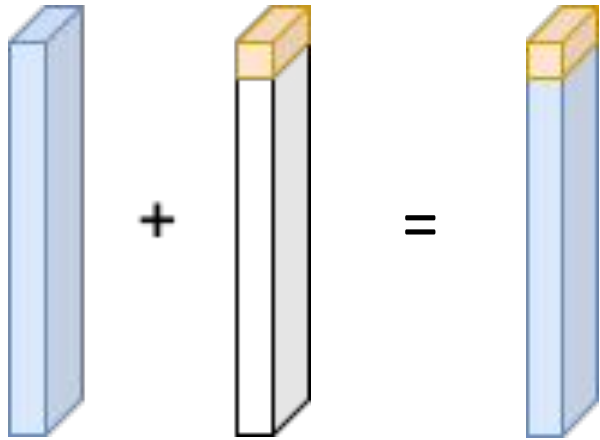
Related work:

- $\beta$ -VAE: restructures ELBO with variable  $\beta$  to enforce latent independence
- FactorVAE: uses discriminator to enhance loss for badly factorized latents
- SlowVAE: regularises for temporal sparsity using prior hyperparameters

<sup>1</sup> Unsupervised deep learning identifies semantic disentanglement in single inferotemporal face patch neurons, DeepMind

# Equivariance

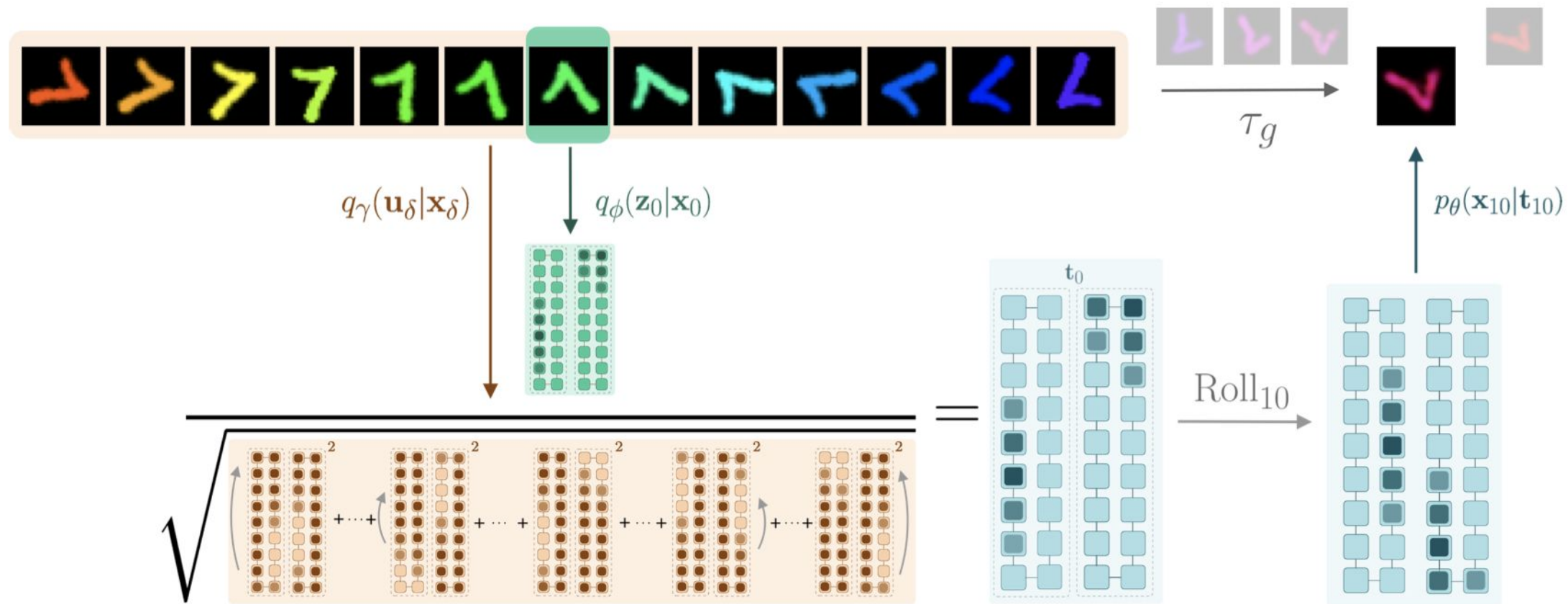
$$T'[f(x)] = f(T[x])$$



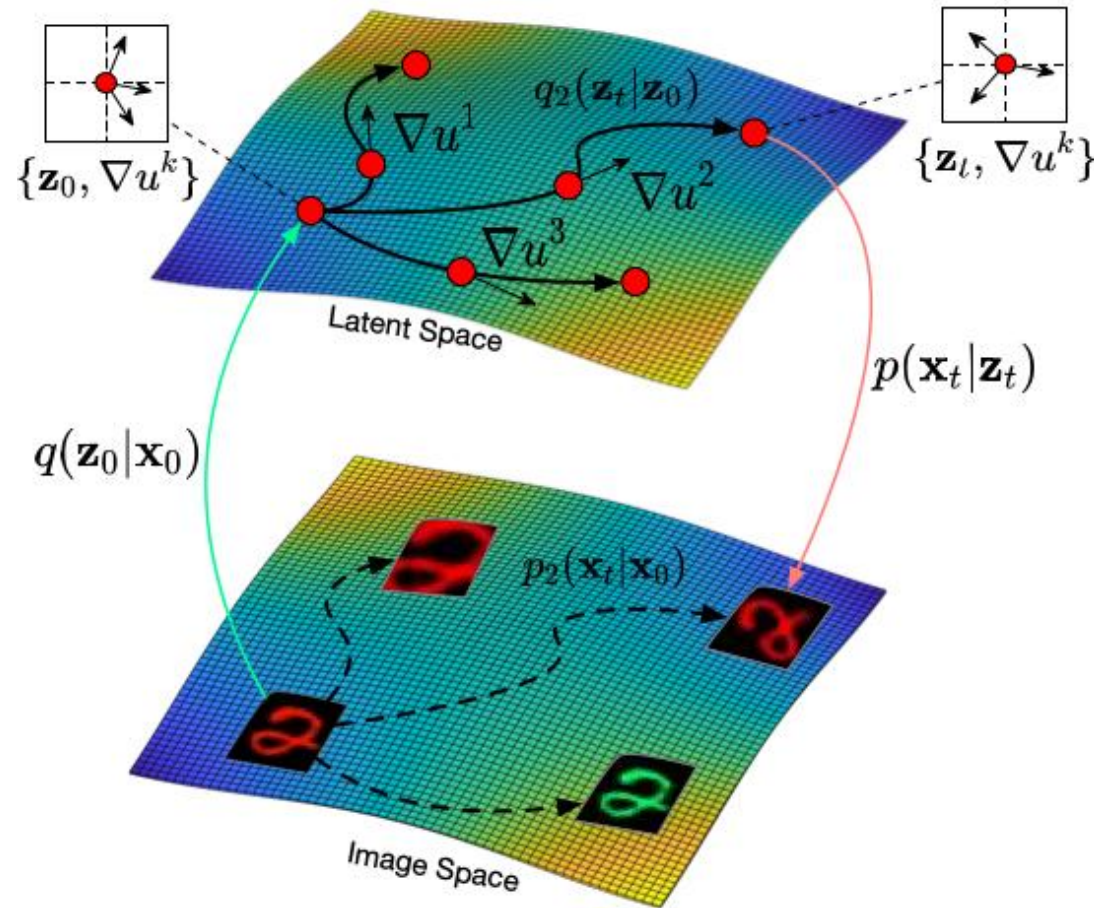
Related work:

- Topographic VAE

# Equivariance



# Alternative latent perspective



$$\int_{\mathbf{z}_0, \mathbf{z}_t} q(\mathbf{z}_0 | \mathbf{x}_0) q_k(\mathbf{z}_t | \mathbf{z}_0) p(\mathbf{x}_t | \mathbf{z}_t)$$

=

$$p_k(\mathbf{x}_t | \mathbf{x}_0)$$

# Alternative latent perspective

Supervised case

$$p(\bar{\mathbf{x}}, \bar{\mathbf{z}}, k) = p(k)p(\mathbf{z}_0)p(\mathbf{x}_0|\mathbf{z}_0) \prod_{t=1}^T p(\mathbf{z}_t|\mathbf{z}_{t-1}, k)p(\mathbf{x}_t|\mathbf{z}_t)$$

$$q(\bar{\mathbf{z}}|\bar{\mathbf{x}}, k) = q(\mathbf{z}_0|\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{z}_t|\mathbf{z}_{t-1}, k)$$

$$\log p(\bar{\mathbf{x}}|k) \geq \mathbb{E}_{q_\theta(\bar{\mathbf{z}}|\bar{\mathbf{x}}, k)} [\log p(\bar{\mathbf{x}}|\bar{\mathbf{z}}, k)] + \mathbb{E}_{q_\theta(\bar{\mathbf{z}}|\bar{\mathbf{x}}, k)} \left[ \log \frac{p(\bar{\mathbf{z}}|k)}{q(\bar{\mathbf{z}}|\bar{\mathbf{x}}, k)} \right]$$

$$= \sum_{t=0}^T \mathbb{E}_{q_\theta(\bar{\mathbf{z}}|k)} [\log p(\mathbf{x}_t|\mathbf{z}_t, k)] - \mathbb{E}_{q_\theta(\bar{\mathbf{z}}|k)} [\mathbf{D}_{\text{KL}} [q_\theta(\mathbf{z}_0|\mathbf{x}_0)||p(\mathbf{z}_0)]] - \sum_{t=1}^T \mathbb{E}_{q_\theta(\bar{\mathbf{z}}|k)} [\mathbf{D}_{\text{KL}} [q_\theta(\mathbf{z}_t|\mathbf{z}_{t-1}, k)||p(\mathbf{z}_t|\mathbf{z}_{t-1}, k)]]$$



# Alternative latent perspective

Weakly supervised  
case

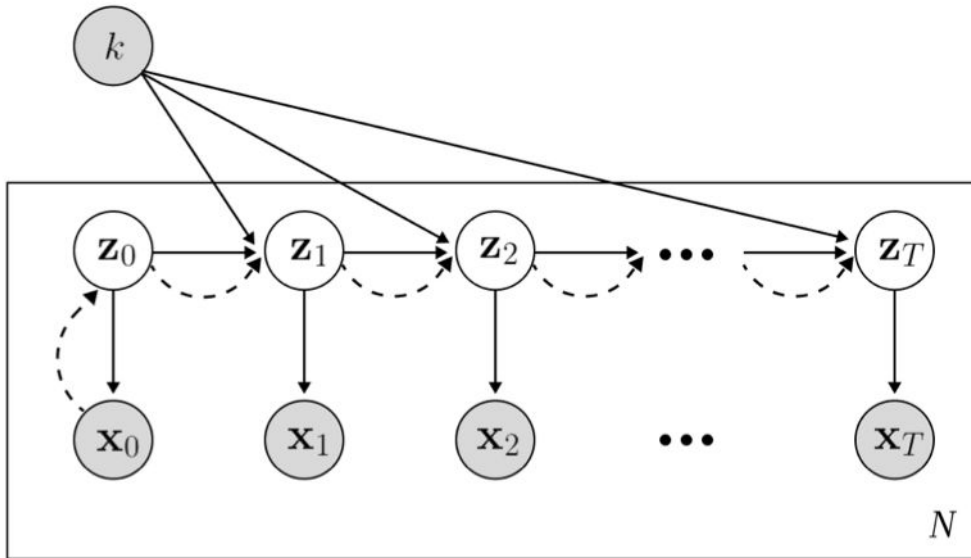
$$p(\bar{\mathbf{x}}, \bar{\mathbf{z}}, k) = p(k)p(\mathbf{z}_0)p(\mathbf{x}_0|\mathbf{z}_0) \prod_{t=1}^T p(\mathbf{z}_t|\mathbf{z}_{t-1}, k)p(\mathbf{x}_t|\mathbf{z}_t)$$

$$q(\bar{\mathbf{z}}, k|\bar{\mathbf{x}}) = q(k|\bar{\mathbf{x}})q(\mathbf{z}_0|\bar{\mathbf{x}}) \prod_{t=1}^T q(\mathbf{z}_t|\mathbf{z}_{t-1}, k)$$

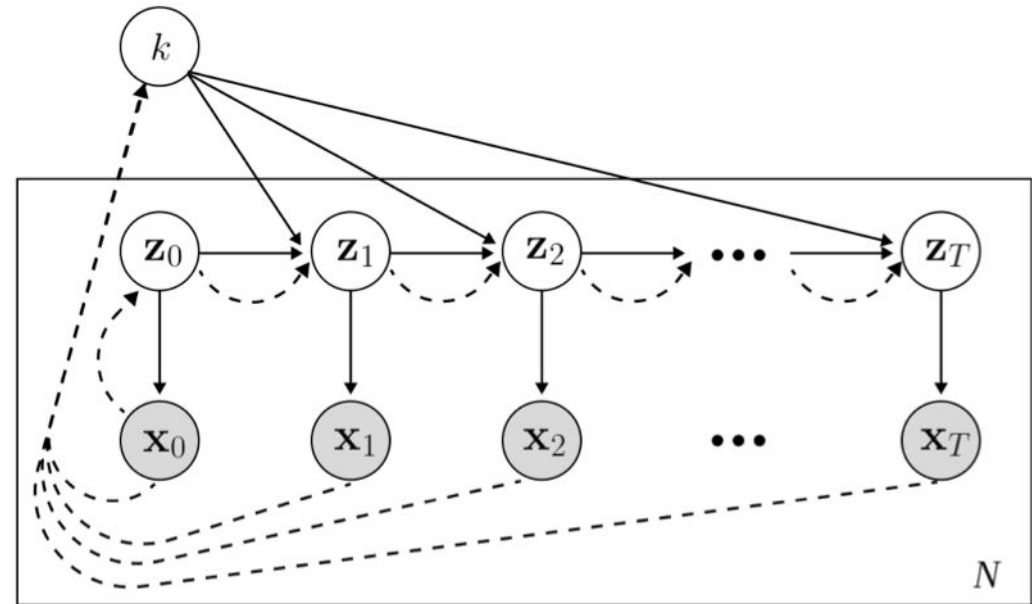
$$\log p(\bar{\mathbf{x}}) \geq \mathbb{E}_{q_\theta(\bar{\mathbf{z}}, k|\bar{\mathbf{x}})} [\log p(\bar{\mathbf{x}}|\bar{\mathbf{z}}, k)] + \mathbb{E}_{q_\theta(\bar{\mathbf{z}}, k|\bar{\mathbf{x}})} \left[ \log \frac{p(\bar{\mathbf{z}}|k)}{q(\bar{\mathbf{z}}|\bar{\mathbf{x}}, k)} \right] + \mathbb{E}_{q_\gamma(k|\bar{\mathbf{x}})} \left[ \log \frac{p(k)}{q(k|\bar{\mathbf{x}})} \right]$$

# Alternative latent perspective

Supervised case



Weakly supervised case



White nodes: latent variables

Grey nodes: observed data

Solid lines: generative model

Dashed lines: approximate posterior

# Time evolving priors & posteriors

$$D_{\text{KL}} [q_{\theta}(\mathbf{z}_t | \mathbf{z}_{t-1}, k) || p(\mathbf{z}_t | \mathbf{z}_{t-1}, k)] ?$$

- Continuity equation:  $\partial_t p(\mathbf{z}) = -\nabla \cdot (p(\mathbf{z}) \nabla u^k(\mathbf{z}))$
- Discrete particle evolution:  $\mathbf{z}_t = g(\mathbf{z}_{t-1}, k) = \mathbf{z}_{t-1} + \nabla_{\mathbf{z}} u^k$

Prior:

$$u^k = -D_k \log p(\mathbf{z}_t)$$

$$\partial_t p(\mathbf{z}_t) = D_k \nabla^2 p(\mathbf{z}_t)$$

$$p(\mathbf{z}_0) = \mathcal{N}(0, 1)$$

Posterior:

$$\log q(\mathbf{z}_t | \mathbf{z}_{t-1}, k) = \log q(\mathbf{z}_{t-1}) - \log |1 + \nabla_{\mathbf{z}}^2 u^k|$$

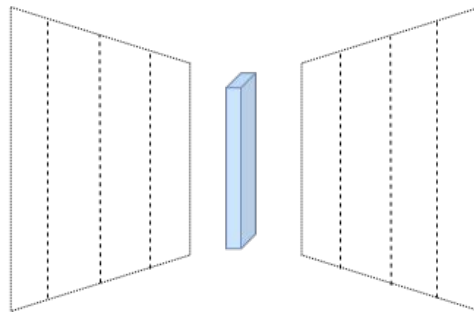
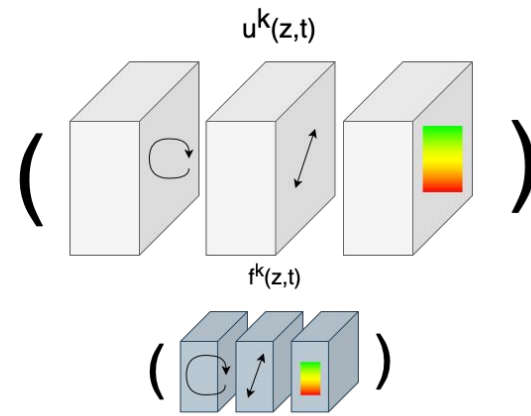
# Optimal Transport for posterior flow

$$\frac{\partial}{\partial t} u^k(\mathbf{z}, t) + \frac{1}{2} \|\nabla_{\mathbf{z}} u^k(\mathbf{z}, t)\|^2 = f(\mathbf{z}, t) \quad \text{subject to } f(\mathbf{z}, t) \leq 0$$

$$\rightarrow \mathcal{L}_{HJ} = \frac{1}{T} \sum_{t=1}^T \left( \frac{\partial}{\partial t} u^k(\mathbf{z}, t) + \frac{1}{2} \|\nabla_{\mathbf{z}} u^k(\mathbf{z}, t)\|^2 - f(\mathbf{z}, t) \right)^2 + \|\nabla u^k(\mathbf{z}_0, 0)\|^2$$

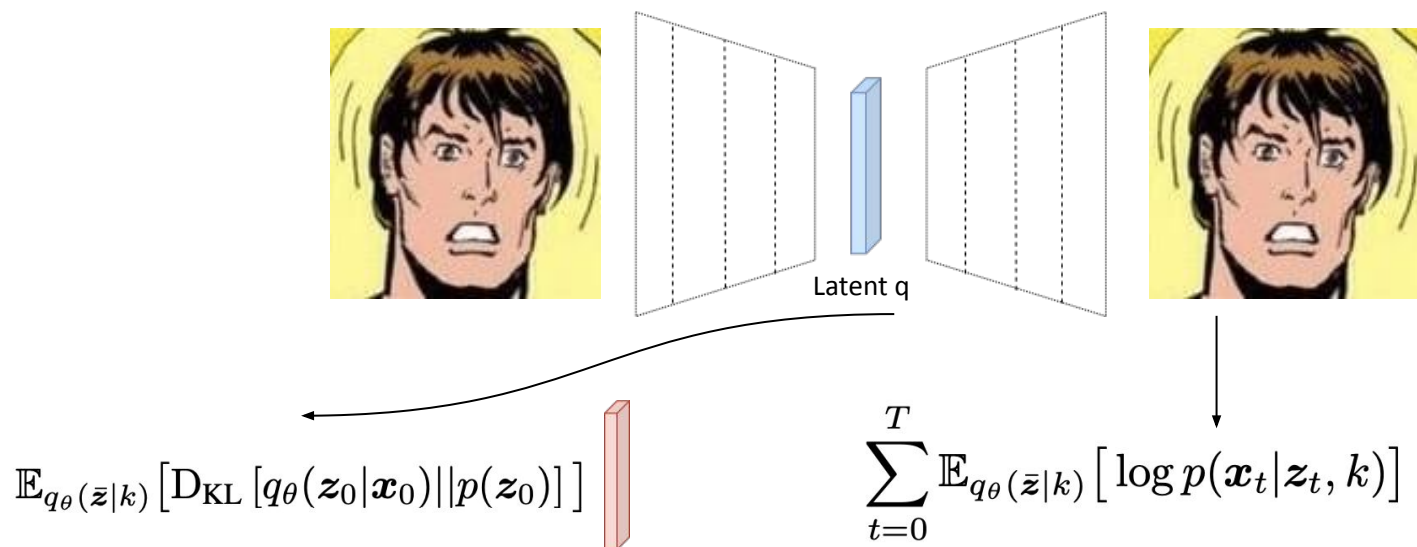
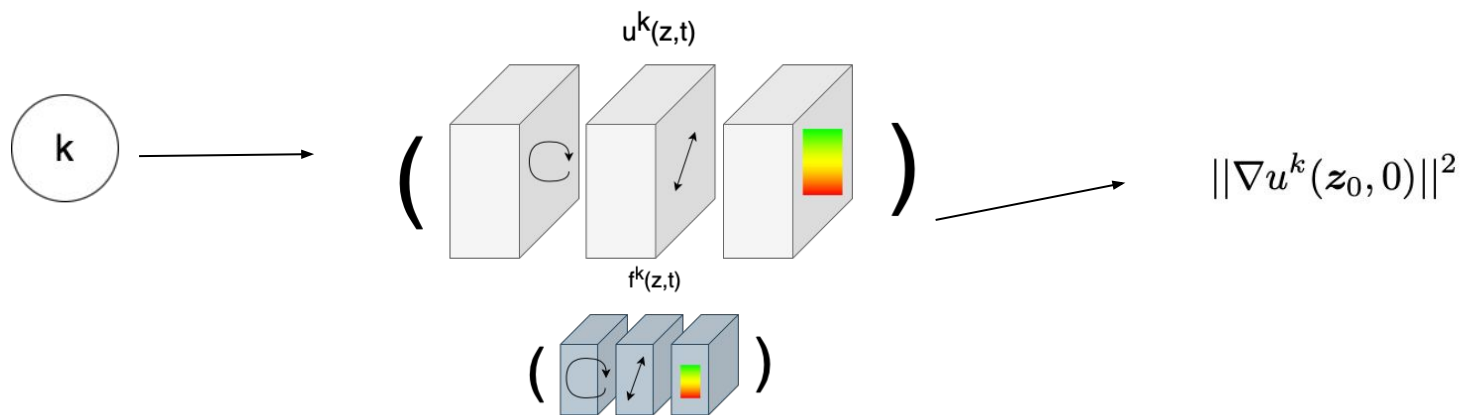


# Implementation



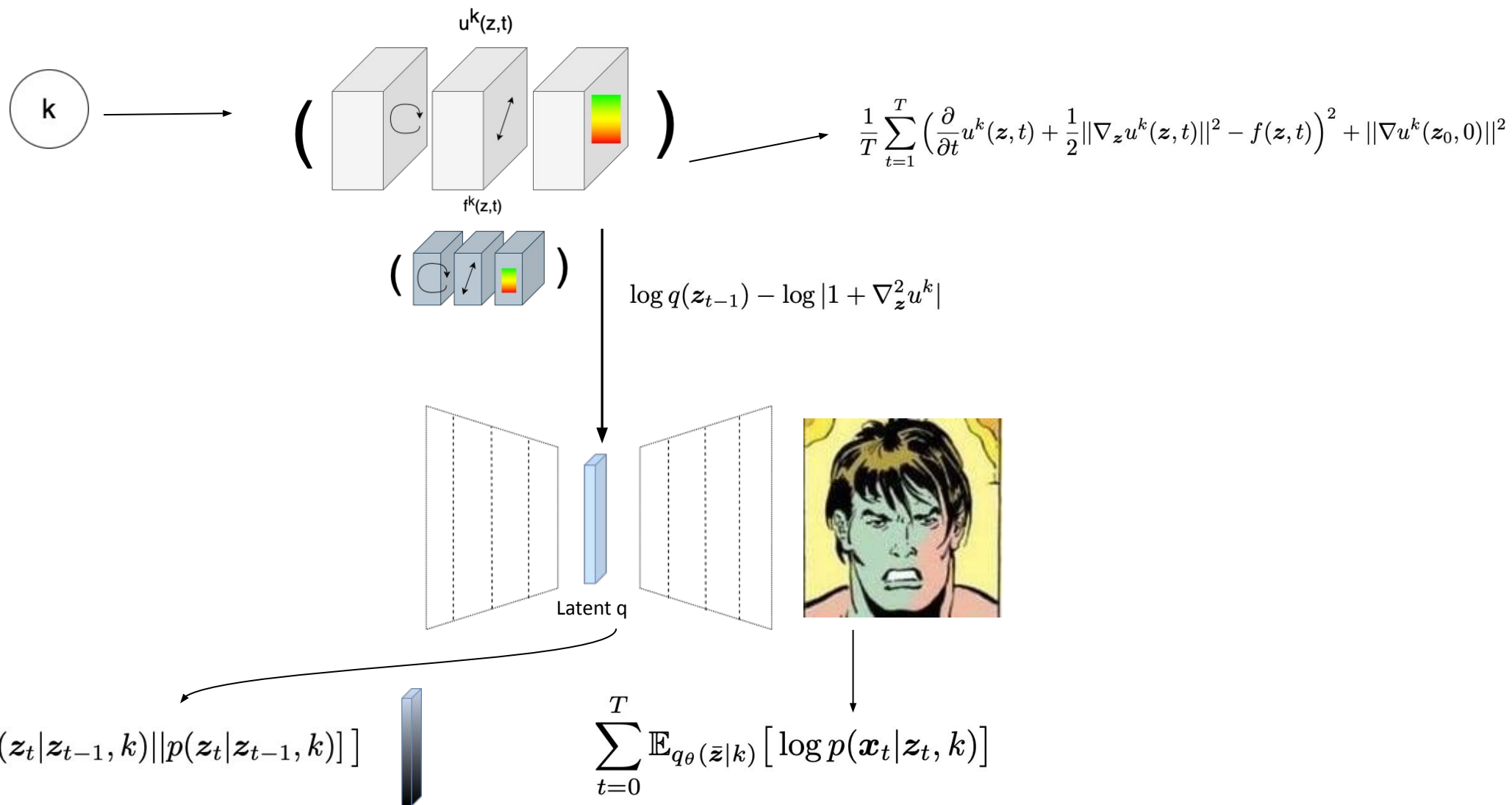
# Implementation

t= 0



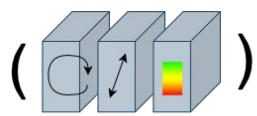
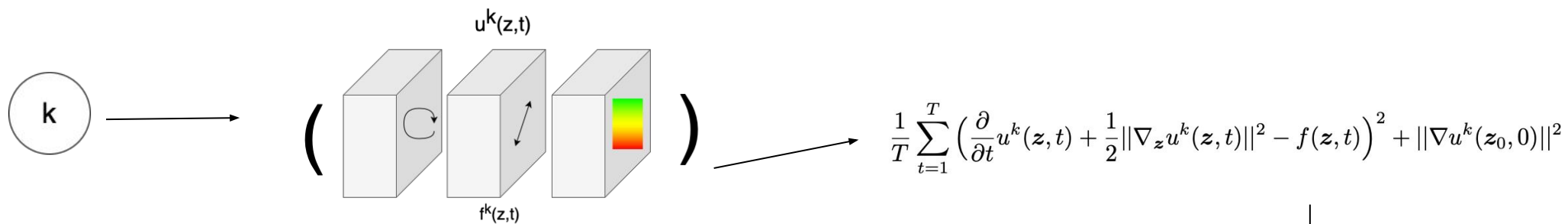
# Implementation

t= 1

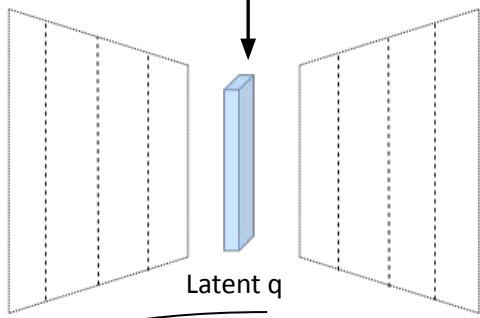


# Implementation

$t = 2 = T$



$$\log q(z_{t-1}) - \log |1 + \nabla_z^2 u^k|$$



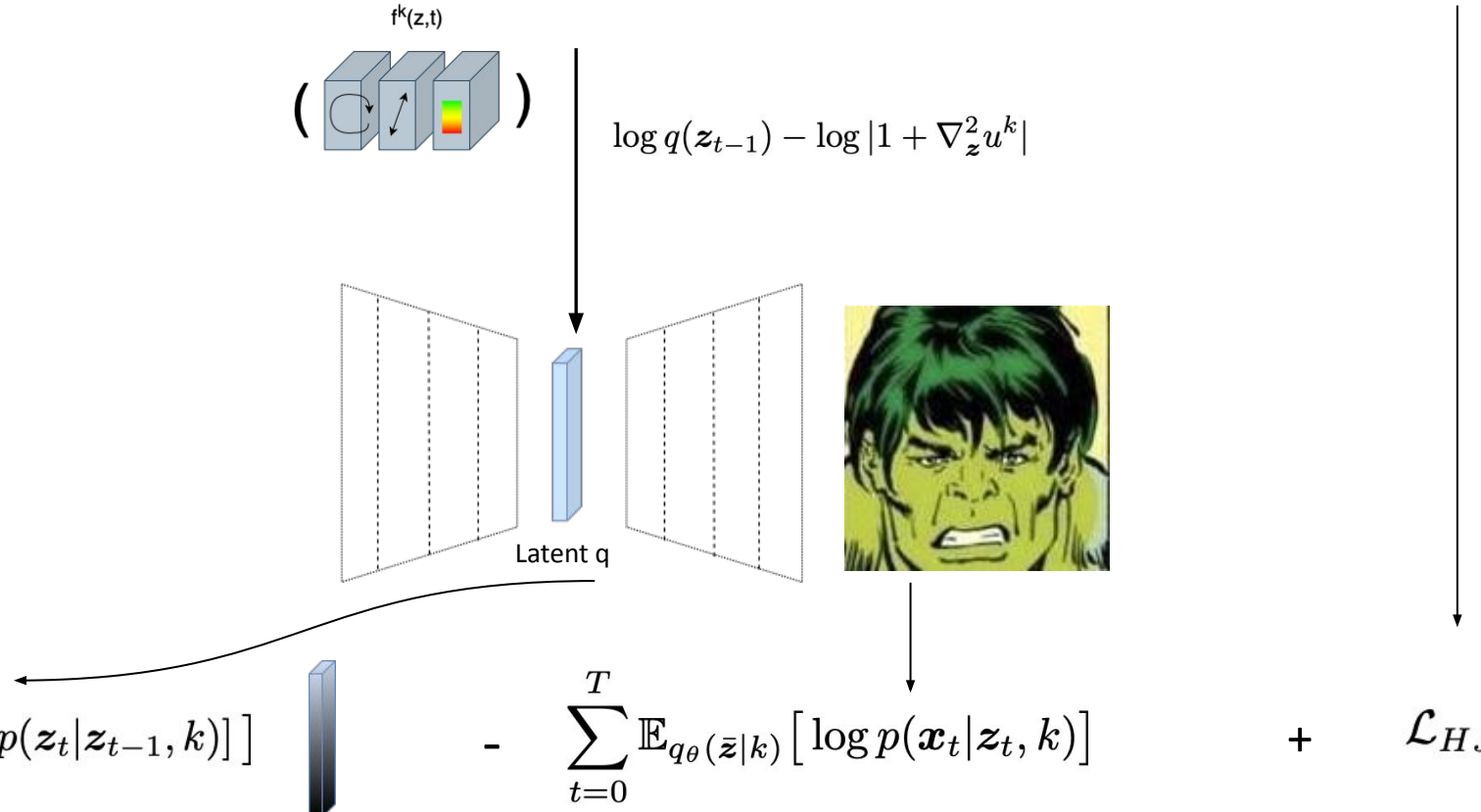
$\mathcal{L} =$

$$\mathbb{E}_{q_\theta(\bar{z}|k)} [\text{D}_{\text{KL}} [q_\theta(z_0|\mathbf{x}_0) || p(z_0)]]$$

$$+ \sum_{t=1}^T \mathbb{E}_{q_\theta(\bar{z}|k)} [\text{D}_{\text{KL}} [q_\theta(z_t|z_{t-1}, k) || p(z_t|z_{t-1}, k)]]$$


$$- \sum_{t=0}^T \mathbb{E}_{q_\theta(\bar{z}|k)} [\log p(\mathbf{x}_t|z_t, k)]$$

$$+ \mathcal{L}_{HJ}$$





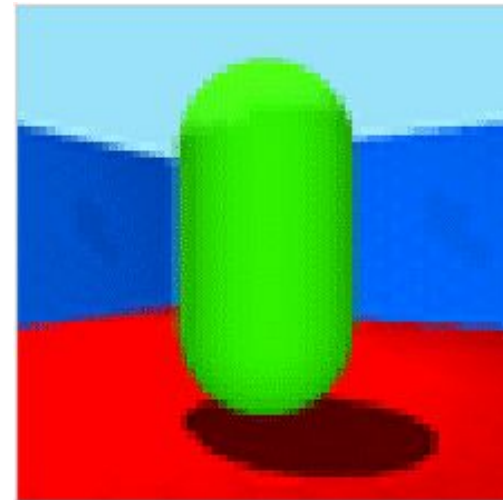
# Implementation

 : one latent dimension, corresponding to a normal distribution


Color in latent corresponds to distinct ground truth factor



Latent representation  $q$



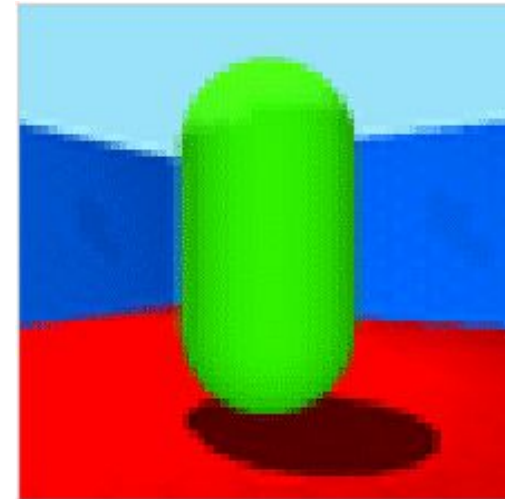
# Implementation

 : one latent dimension, corresponding to a time-evolving normal distribution

Color in latent corresponds to distinct ground truth factor

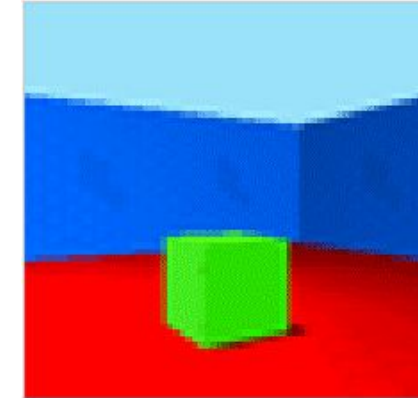


Latent representation  $q(z_t | z_{t-1}, k_{\text{color}})$



# Experiment setting

- Datasets:
  - MNIST: color, rotation, scale
  - Shapes3D: wall-, floor- & object hue, scale





lighting intensity



lighting x-dir



lighting y-dir



lighting z-dir



camera x-pos



camera y-pos



camera z-pos



object shape



robot x-movement



robot y-movement



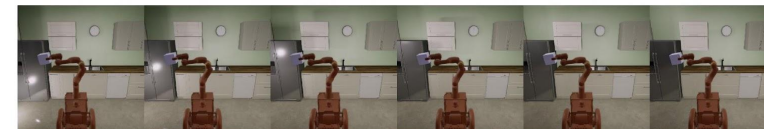
camera height



object scale



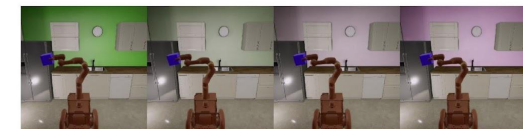
lighting intensity



lighting y-dir



object color



wall color

# Experiment setting

- Metrics:

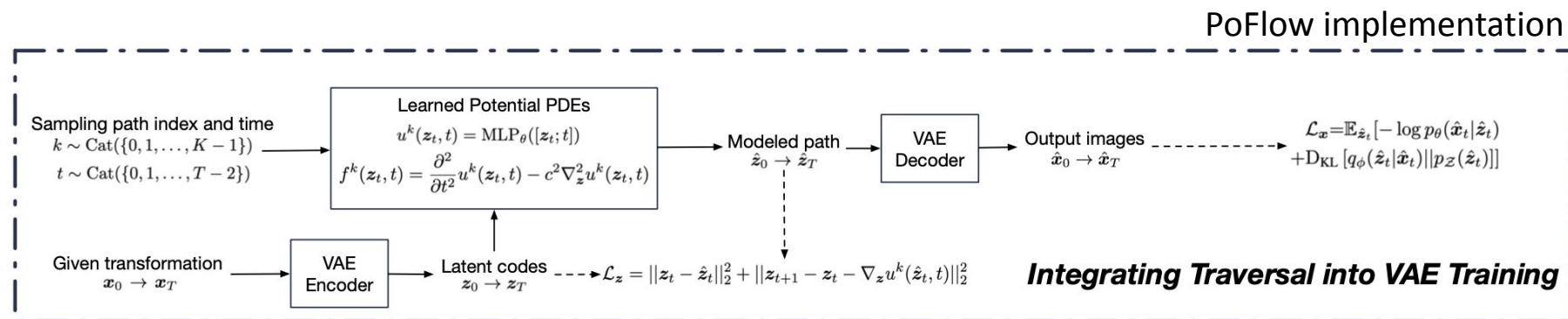
- Equivariance error:  $\mathcal{E}_k = \sum_{t=1}^T |\mathbf{x}_t - \text{Decode}(\mathbf{z}_t)|$ ,  $\mathbf{z}_t = \mathbf{z}_0 + \sum_{t=1}^T \nabla_{\mathbf{z}} u^k$

- Log-likelihood:  $\log p(\mathbf{x}_t)$

- VP score: (%) of correctly classified pairs  $[\mathbf{x}_0, \mathbf{x}_T]$  by a lightweight NN

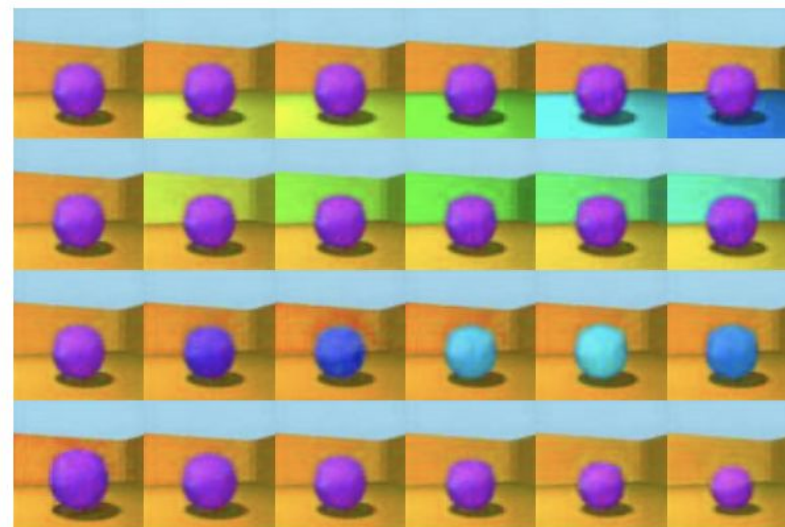
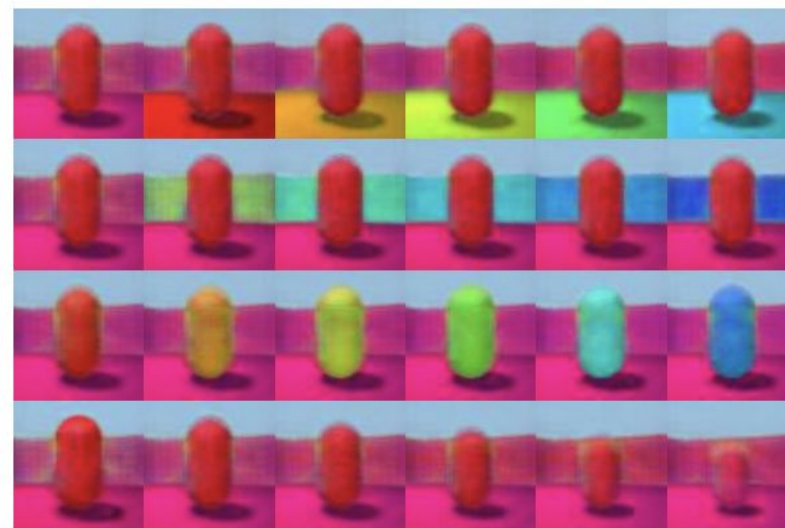
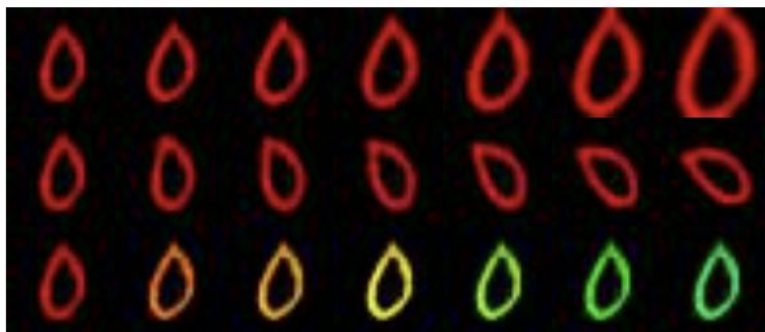
- Baselines:

- PoFlow
- Topographic VAE
- SlowVAE
- $\beta$ -VAE
- FactorVAE
- VAE





# Qualitative results

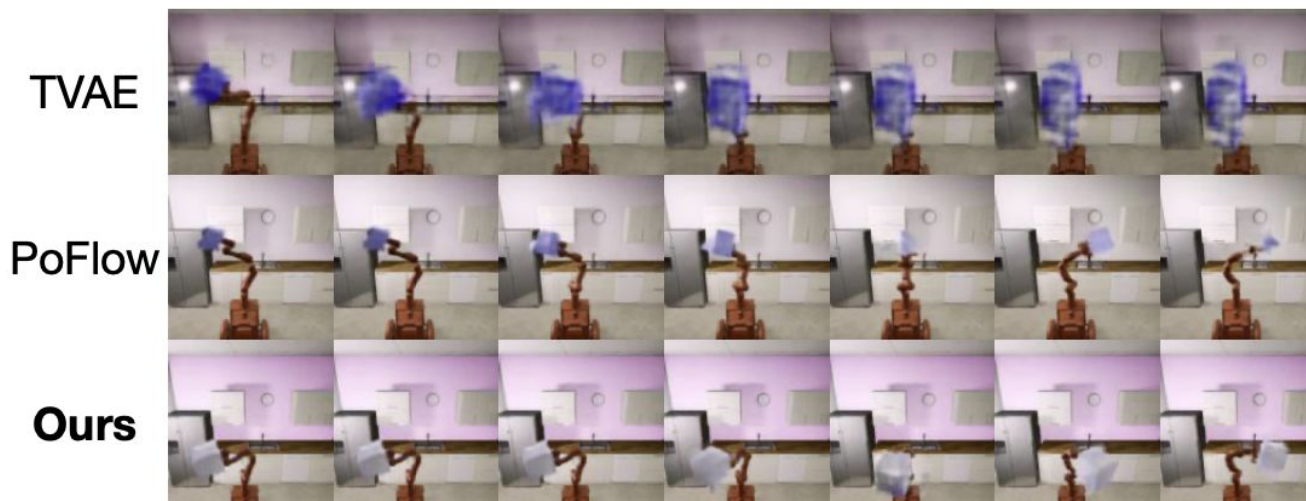


# Qualitative results

*Camera X-Pos*



*Robot X-move*





# Quantitative results

Methods	Supervision?	Equivariance Error ( $\downarrow$ )			Log-likelihood ( $\uparrow$ )
		Scaling	Rotation	Coloring	
VAE [47]	No (✗)	1275.31 $\pm$ 1.89	1310.72 $\pm$ 2.19	1368.92 $\pm$ 2.33	-2206.17 $\pm$ 1.83
$\beta$ -VAE [35]	No (✗)	741.58 $\pm$ 4.57	751.32 $\pm$ 5.22	808.16 $\pm$ 5.03	-2224.67 $\pm$ 2.35
FactorVAE [46]	No (✗)	659.71 $\pm$ 4.89	632.44 $\pm$ 5.76	662.18 $\pm$ 5.26	-2209.33 $\pm$ 2.47
SlowVAE [49]	Weak (✓)	461.59 $\pm$ 5.37	447.46 $\pm$ 5.46	398.12 $\pm$ 4.83	-2197.68 $\pm$ 2.39
TVAE [45]	Yes (✓)	505.19 $\pm$ 2.77	493.28 $\pm$ 3.37	451.25 $\pm$ 2.76	-2181.13 $\pm$ 1.87
PoFlow [79]	Yes (✓)	234.78 $\pm$ 2.91	231.42 $\pm$ 2.98	240.57 $\pm$ 2.58	-2145.03 $\pm$ 2.01
<b>Ours</b>	Yes (✓)	<b>185.42<math>\pm</math>2.35</b>	<b>153.54<math>\pm</math>3.10</b>	<b>158.57<math>\pm</math>2.95</b>	<b>-2112.45<math>\pm</math>1.57</b>
<b>Ours</b>	Weak (✓)	<b>193.84<math>\pm</math>2.47</b>	<b>157.16<math>\pm</math>3.24</b>	<b>165.19<math>\pm</math>2.78</b>	<b>-2119.94<math>\pm</math>1.76</b>

Table 1: Equivariance error  $\mathcal{E}_k$  and log-likelihood  $\log p(\mathbf{x}_t)$  on MNIST

Methods	Supervision?	Equivariance Error ( $\downarrow$ )				Log-likelihood ( $\uparrow$ )
		Floor Hue	Wall Hue	Object Hue	Scale	
VAE [47]	No (✗)	6924.63 $\pm$ 8.92	7746.37 $\pm$ 8.77	4383.54 $\pm$ 9.26	2609.59 $\pm$ 7.41	-11784.69 $\pm$ 4.87
$\beta$ -VAE [35]	No (✗)	2243.95 $\pm$ 12.48	2279.23 $\pm$ 13.97	2188.73 $\pm$ 12.61	2037.94 $\pm$ 11.72	-11924.83 $\pm$ 5.64
FactorVAE [46]	No (✗)	1985.75 $\pm$ 13.26	1876.41 $\pm$ 11.93	1902.83 $\pm$ 12.27	1657.32 $\pm$ 11.05	-11802.17 $\pm$ 5.69
SlowVAE [49]	Weak (✓)	1247.36 $\pm$ 12.49	1314.86 $\pm$ 11.41	1102.28 $\pm$ 12.17	1058.74 $\pm$ 10.96	-11674.89 $\pm$ 5.74
TVAE [45]	Yes (✓)	1225.47 $\pm$ 9.82	1246.32 $\pm$ 9.54	1261.79 $\pm$ 9.86	1142.01 $\pm$ 9.37	-11475.48 $\pm$ 5.18
PoFlow [79]	Yes (✓)	885.46 $\pm$ 10.37	916.71 $\pm$ 10.49	912.48 $\pm$ 9.86	924.39 $\pm$ 10.05	-11335.84 $\pm$ 4.95
<b>Ours</b>	Yes (✓)	<b>613.29<math>\pm</math>8.93</b>	<b>653.45<math>\pm</math>9.48</b>	<b>605.79<math>\pm</math>8.63</b>	<b>599.71<math>\pm</math>9.34</b>	<b>-11215.42<math>\pm</math>5.71</b>
<b>Ours</b>	Weak (✓)	<b>690.84<math>\pm</math>9.57</b>	<b>717.74<math>\pm</math>10.65</b>	<b>681.59<math>\pm</math>9.02</b>	<b>653.58<math>\pm</math>9.57</b>	<b>-11279.61<math>\pm</math>5.89</b>

Table 2: Equivariance error  $\mathcal{E}_k$  and log-likelihood  $\log p(\mathbf{x}_t)$  on Shapes3D



# Quantitative results

Methods	Lighting Intensity	Lighting X-dir	Lighting Y-dir	Lighting Z-dir	Camera X-pos	Camera Y-pos	Camera Y-pos
TVAE [45]	11477.81	12568.32	11807.34	11829.33	11539.69	11736.78	11951.45
PoFlow [79]	8312.97	7956.18	8519.39	8871.62	8116.82	8534.91	8994.63
Ours	<b>5798.42</b>	<b>6145.09</b>	<b>6334.87</b>	<b>6782.84</b>	<b>6312.95</b>	<b>6513.68</b>	<b>6614.27</b>

Table 3: Equivariance error ( $\downarrow$ ) on Falcol3D

Methods	Robot X-move	Robot Y-move	Camera Height	Object Scale	Lighting Intensity	Lighting Y-dir	Object Color	Wall Color
TVAE [45]	8441.65	8348.23	8495.31	8251.34	8291.70	8741.07	8456.78	8512.09
PoFlow [79]	6572.19	6489.35	6319.82	6188.59	6517.40	6712.06	7056.98	6343.76
Ours	<b>3659.72</b>	<b>3993.33</b>	<b>4170.27</b>	<b>4359.78</b>	<b>4225.34</b>	<b>4019.84</b>	<b>5514.97</b>	<b>3876.01</b>

Table 4: Equivariance error ( $\downarrow$ ) on Isaac3D

Table 5: VP Scores (%) on MNIST.

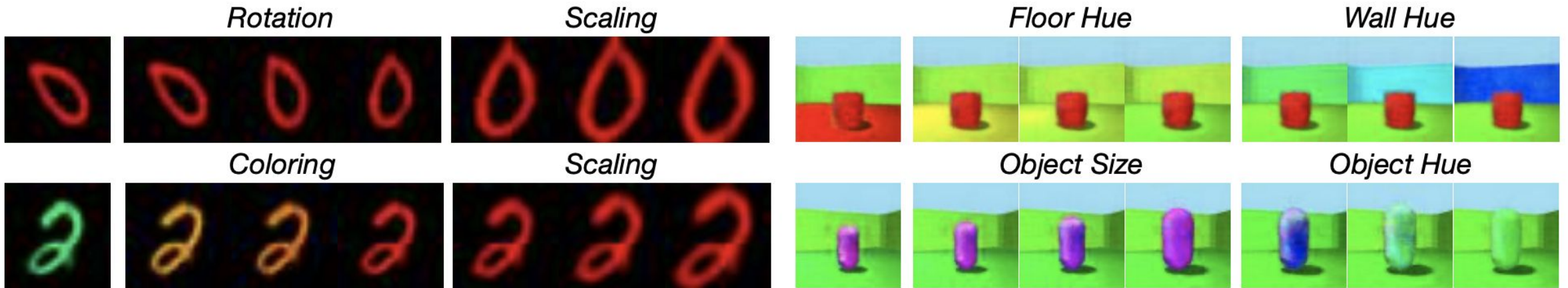
Training Set Split	Ours	PoFlow	TVAE	FactorVAE	$\beta$ -VAE
10%	<b>95.69</b>	93.05	89.91	85.92	87.31
1%	<b>92.71</b>	91.27	88.15	84.46	85.25

Table 6: VP Scores (%) on Shapes3D.

Training Set Split	Ours	PoFlow	TVAE	FactorVAE	$\beta$ -VAE
10%	<b>95.92</b>	91.48	88.27	84.49	85.91
1%	<b>77.03</b>	72.32	68.39	63.83	65.78

# Additional evaluations

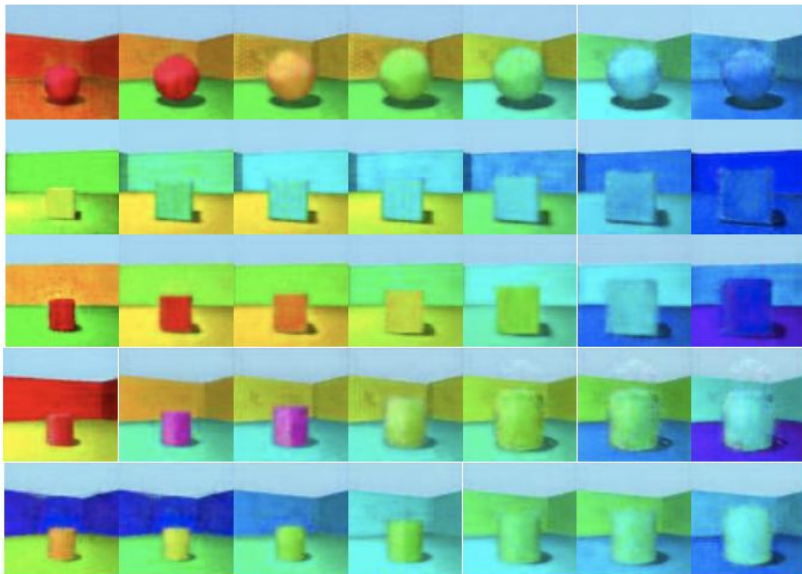
- 1: switching transformations



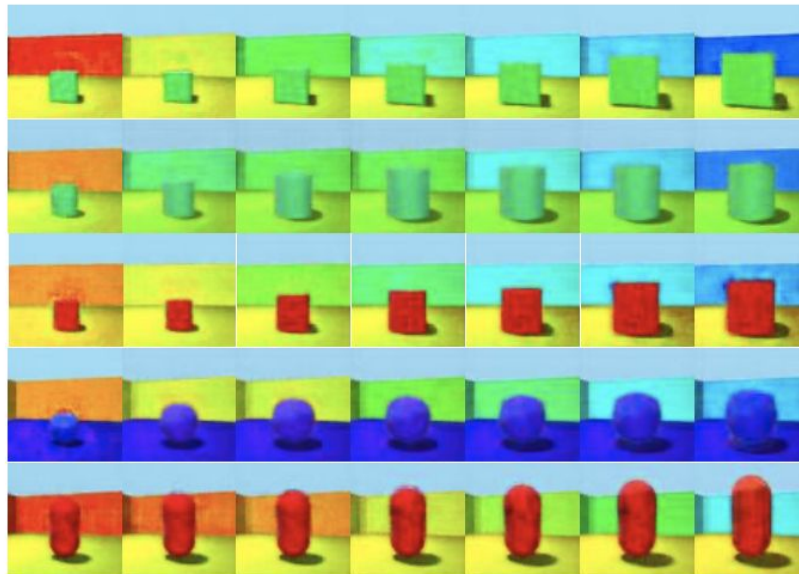
# Additional evaluations

- 2: superposing transformations

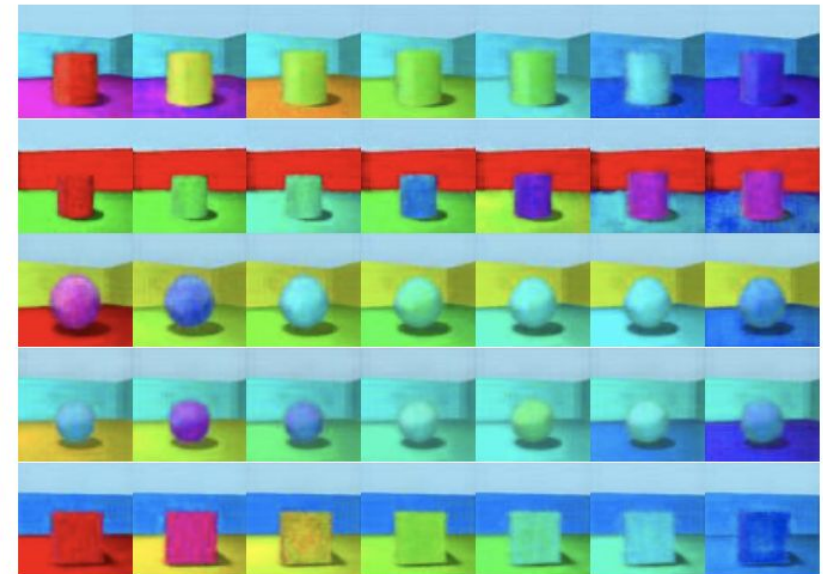
*All Transformations*



*Wall Hue + Scale*

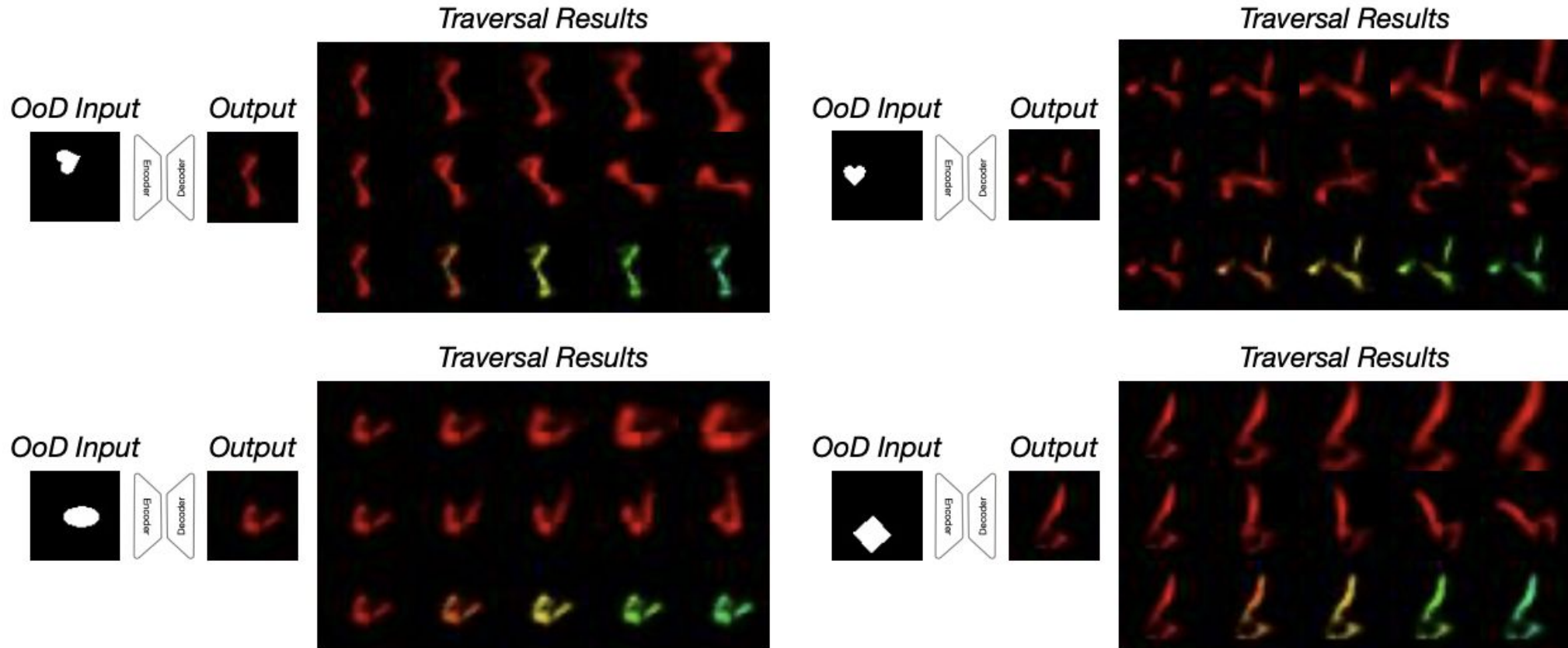


*Floor Hue + Object Hue*



# Additional evaluations

- 3: equivariance generalisation to new data





# Ablation studies

Table 7: Equivariance error of different priors.

Prior	Scaling	Rotation	Coloring
<b>SG</b>	190.24±2.18	158.93±3.25	164.18±2.77
<b>MoG</b>	188.23±2.45	157.79±2.86	161.49±2.62
<b>VAMP</b>	192.81±3.67	161.47±4.12	162.97±3.89
<b>Diffusion</b>	<b>185.42±2.35</b>	<b>153.54±3.10</b>	<b>158.57±2.95</b>

Table 8: Equivariance error of different PDEs.

Prior	Scaling	Rotation	Coloring
<b>Heat</b>	223.95±3.38	212.47±3.85	207.66±2.91
<b>FP</b>	211.54±3.17	188.59±3.92	194.73±3.09
<b>OHJ</b>	190.43±2.48	163.87±3.03	162.38±2.86
<b>GHJ</b>	<b>185.42±2.35</b>	<b>153.54±3.10</b>	<b>158.57±2.95</b>

# Ablation studies

Table 9: Equivariance error on MNIST of a different number of transformations ( $K$ ).

$K$	Scaling	Rotation	Coloring
1	$185.27 \pm 2.59$	–	–
2	$185.78 \pm 2.21$	$154.29 \pm 2.87$	–
3	$185.42 \pm 2.45$	$153.54 \pm 3.10$	$158.57 \pm 2.95$

Table 10: Equivariance error on MNIST of different sequence lengths ( $T$ ).

Sequence Length ( $T$ )	Scaling	Rotation
9	$185.42 \pm 2.35$	$153.54 \pm 3.10$
12	$214.47 \pm 2.59$	$198.72 \pm 2.89$

# Conclusion: positives

- Differentiated approach
- Quantitatively better model
- Intuitive & strong qualitative results

# Conclusion: negatives

- Code doesn't work
- More expensive to train
- Short on generalisation results
- Short on theoretical intuition



Thank you for your attention!

Questions?

```

1 import torch
2
3 #Randomly sample a transformation at each iteration
4 index = torch.randint(0, potential_number)
5 x_bar = sequence_generation(index)
6
7 #Generating index according to the supervision setting
8 if training_mode = "supervised":
9     index_potential = index
10 elif training_mode = "weakly-supervised":
11     index_potential = q_k(x_bar)
12
13 #initial element of the sequence
14 z, rho_z = flow_vae(x_bar[0])
15
16 #Future elements of the sequence obtained by latent flow
17 for t in range(0,T)
18     PDE_loss, delta_z, delat_rho_z = HJ_PDE(index_potential,z,t)
19
20     #Updates in the sample and probability space
21     z = z + delta_z
22     rho_z = rho_z + delat_rho_z
23
24     #Inference at every intermediate step
25     hat_xt = flow_vae.inference(z)
26
27     #Loss: PDE loss + reconstruction loss + KL div
28     loss += PDE_loss + CE(hat_xt,x_bar[t]) + KL(rho_z, prior_rho_z)
29
30 #KL div for index prediction (weakly-supervised setting)
31 if training_mode = "weakly-supervised":
32     loss += KL(index_potential,index)
33
34 loss.backward()
35 optimizer.step()

```

# Gumbel-Softmax trick

**Weakly-supervised setting.** For the Gumbel-Softmax trick, we re-parameterize  $q_\gamma(k|\bar{\mathbf{x}})$  by

$$y_i = \frac{e^{\frac{x_i + g_i}{\tau}}}{\sum_i e^{\frac{x_i + g_i}{\tau}}} \quad (15)$$

where  $x_i$  is the category prediction,  $g_i$  is the sample drawn from Gumbel distributions, and  $\tau$  is the small temperature to make softmax behave like argmax. We take the ‘hard’ binary prediction in the forward pass and use the straight-through gradient estimator [6] during backpropagation. The temperature  $\tau$  is initialized with 1 and is gradually reduced to 0.05 with the annealing rate  $3e-5$ .

We now turn to introduce why HJ equations could minimize the Wasserstein distance. As stated in [4], the  $L_2$  Wasserstein distance can be re-formulated in the fluid mechanical interpretation as

$$W^2 = \inf \int_D \int_0^1 \frac{1}{2} \rho(x, t) v(x, t)^2 dx dt \quad (16)$$

where the density satisfies the continuity equation ( $\partial_t \rho = -\nabla \cdot (\rho(x, t)v(x, t))$ ). If we introduce the momentum  $m(x, t) = \rho(x, t)v(x, t)$  and two Lagrange multipliers  $u$  and  $\lambda$ , the Lagrangian function of the Wasserstein distance would be:

$$L(\rho, m, \phi) = \int_D \int_0^1 \frac{\|m\|^2}{2\rho} + u(\partial_t \rho + \nabla \cdot m) - \lambda(\rho - s^2) \quad (17)$$

where the second term is the equality constraint, and the third term is an equality constraint with a slack variable  $s$ . Using integration by parts formula, the above equation can be re-written as

$$L(\rho, m, \phi) = \int_D \int_0^1 \frac{\|m\|^2}{2\rho} + \int_D u \rho|_0^1 - \int_D \int_0^1 (\partial_t u \rho + \nabla u \cdot m) - \lambda(\rho - s^2) \quad (18)$$

Based on the set of Karush–Kuhn–Tucker (KKT) conditions ( $\partial_m L = 0$ ,  $\partial_u L = 0$ ,  $\partial_\rho L = 0$ , and  $\lambda \geq 0$ ), we would have:

$$\begin{cases} \partial_m L = \frac{m}{\rho} - \nabla u = v - \nabla u = 0 \\ \partial_u L = \partial_t \rho + \nabla \cdot m = 0 \\ \partial_\rho L = -\frac{\|m\|^2}{2\rho^2} - \partial_t u - \lambda = -\frac{1}{2}\|v\|^2 - \partial_t u - \lambda = 0 \end{cases} \quad (19)$$

where the first condition indicates that the gradient  $\nabla u$  acts as the velocity field, and the third condition implies the optimal solution is given by the generalized HJ equation:

$$\partial_t u + \frac{1}{2}\|\nabla u\|^2 = -\lambda \leq 0 \quad (20)$$

We thus apply the generalized HJ equation (*i.e.*,  $\partial_t u + \frac{1}{2}\|\nabla u\|^2 \leq 0$ ) as the constraints. We further use an extra negative force because this would give more dynamics for modeling the posterior flow.

# Time complexity

	<b>Ours (weakly-supervised)</b>	<b>Ours (supervised)</b>	<b>PoFlow</b>	<b>TVAE</b>
Time (s)	1.076	0.723	0.638	0.591

# Semi-supervised setting

## Q6. Semi-supervised Learning

Thanks for the insightful advice. Of course, our method can be extended to the semi-supervised setting. One straightforward implementation could be switching between inferring the transformation index from sequences and directly enforcing the index during the training process. We test such a possibility on MNIST and present the results below:

Equivariance error on MNIST under different supervision settings.

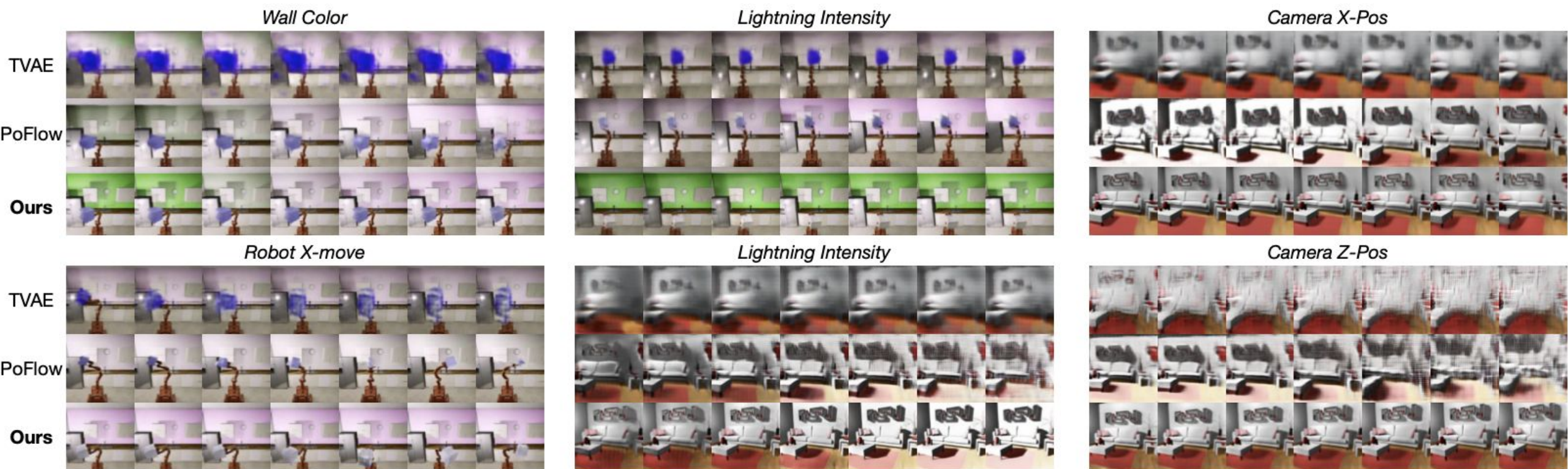
Supervision	Scaling	Rotation	Coloring
Weakly	193.84±2.47	157.16±3.24	165.19±2.78
Semi	186.07±2.58	152.87±3.36	160.13±2.82
Full	185.42±2.35	153.54±3.10	158.57±2.95

As can be seen, this semi-supervised setting has very close and even competitive performance against the supervised setting. We also observed accelerated convergence of this setting compared with the weakly-supervised fashion.

Another interesting semi-supervised learning setting is directly enforcing the index for some of the transformations while inferring the index for others. This setting might be useful when the number of transformations to model is very large. We would add a paragraph in the revised paper to discuss the many intriguing possibilities!



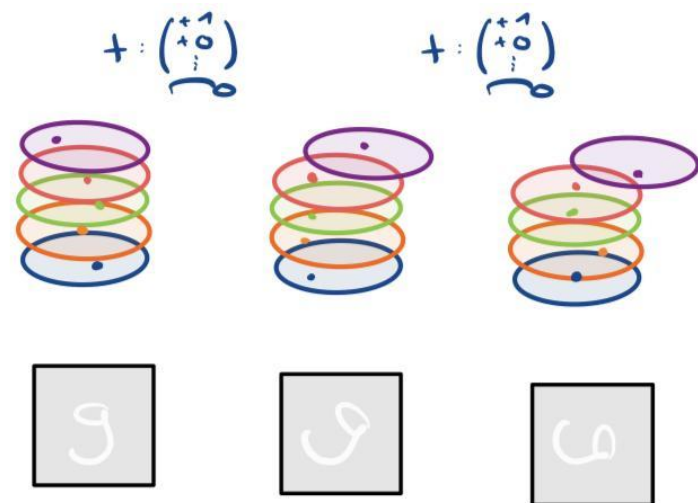
# Full qualitative results (Falcor3D, Isaac3D)





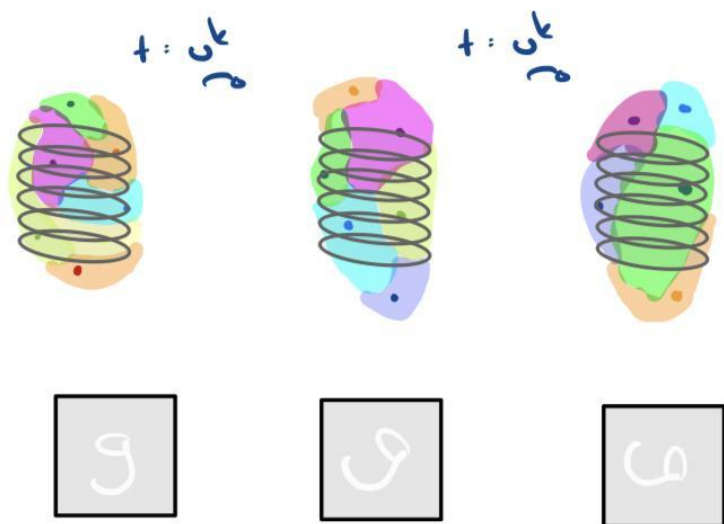


# Implementation



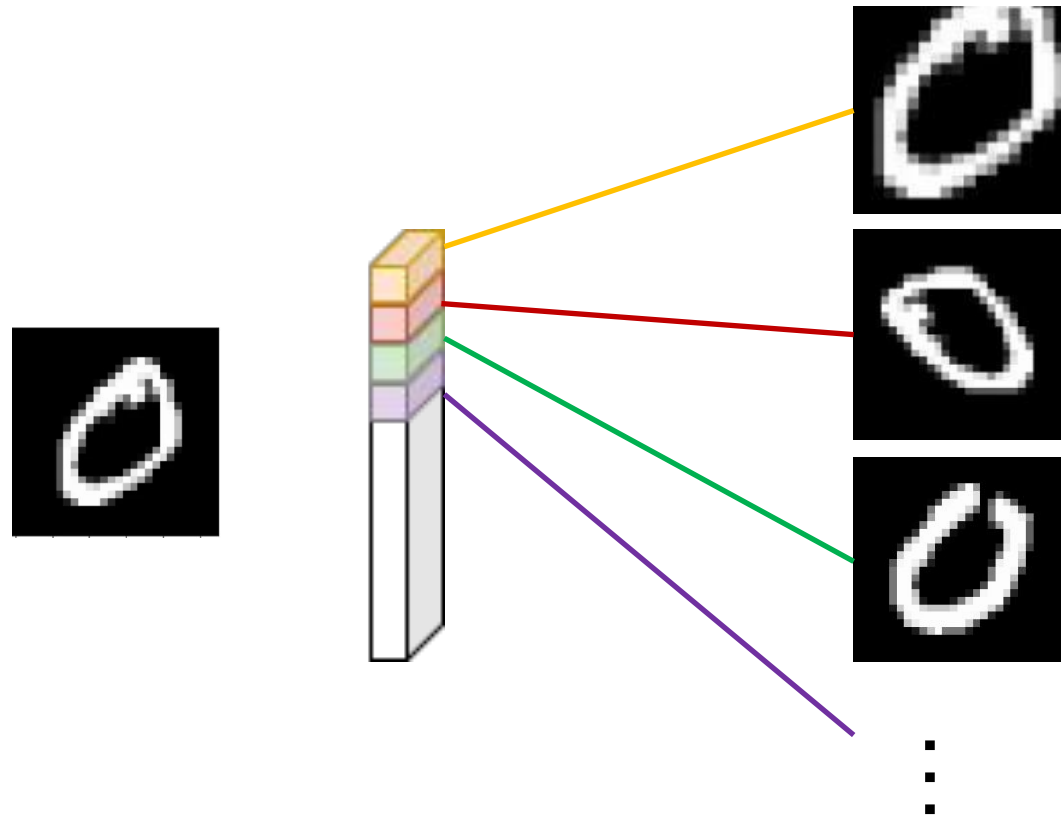
conventional  
disentangled  
VAE latent

○ gaussian distribution "dist",  
corresponds to one dimension



FFRL  
VAE latent

# Disentanglement & Equivariance



$$T'[f(x)] = f(T[x])$$

