

Tutorial “Normalizing Flows”

Part 2

Ullrich Köthe

Visual Learning Lab, Heidelberg University

CVPR 2021, June 2021



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386



European Research Council
Established by the European Commission

Recap:

How do normalizing flows work? How do they differ from other generative models?

Ullrich Köthe

Visual Learning Lab, Heidelberg University

Tutorial „Normalizing Flows“ at CVPR 2021

June 2021



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386



European Research Council
Established by the European Commission



Generative Modelling

- Deep learning success story
 - Compute predictions y directly from complex data x
 - Point estimates: $\hat{y} \approx y^* = \operatorname{argmax} p(y | x)$, posteriors: $\hat{p}_\theta(y | x) \approx p(y | x)$
 - Relies on **discriminative / transductive machine learning**
(does not first build a “model of the world” as traditional sciences do)
 - Problem: discriminative models are hard to interpret, explain, validate
- ⇒ Generative modelling
- Turn the problem around: learn the data generation likelihood $p(x | y)$
 - More difficult: requires **insight** beyond mere prediction capability
 - Solve the original task via Bayes theorem

$$p(y | x) = \frac{p(x | y) p(y)}{p(x)}$$

Feynman: “What I cannot create, I do not understand.”

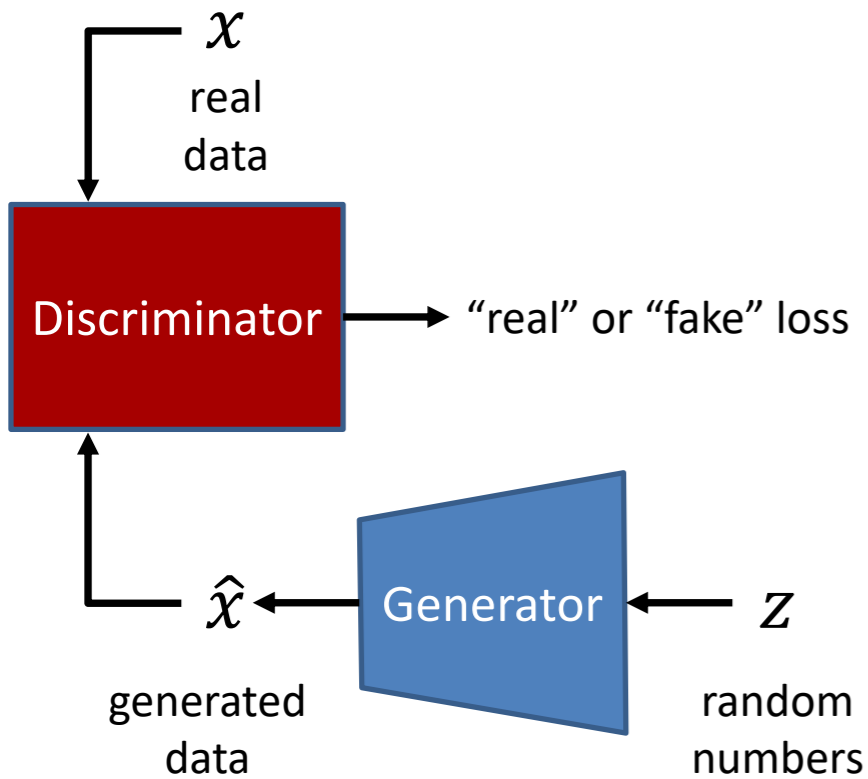




Generative Modelling as a Basis for Interpretable Deep Learning

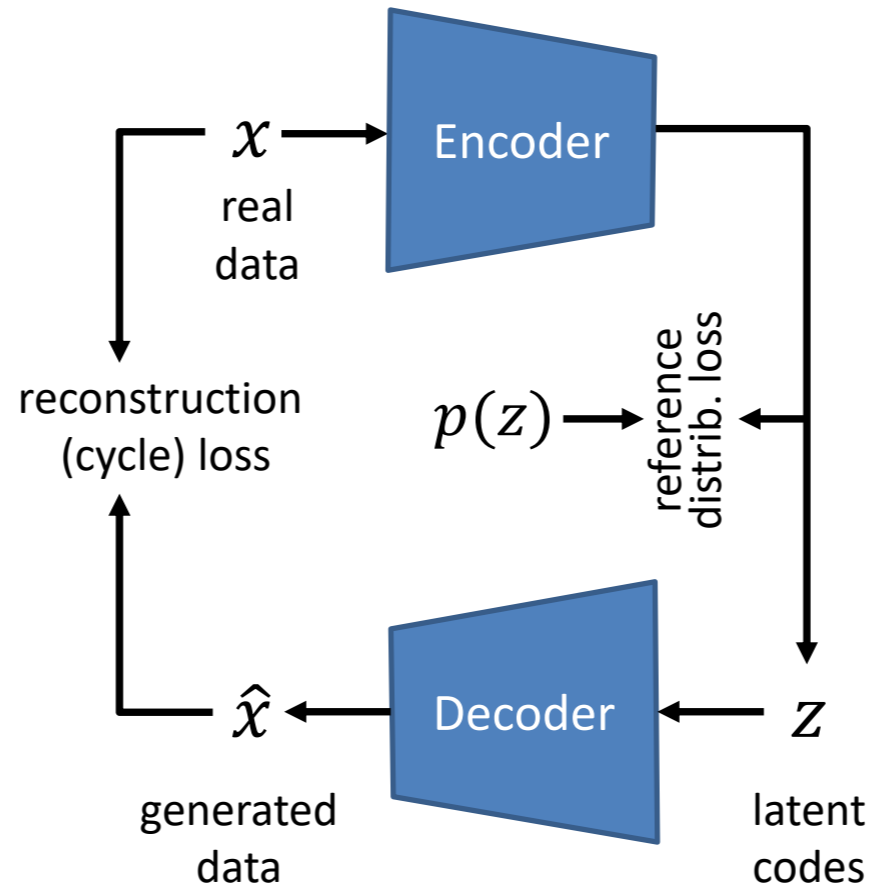
GANs

(Generative Adversarial Networks)



generation only

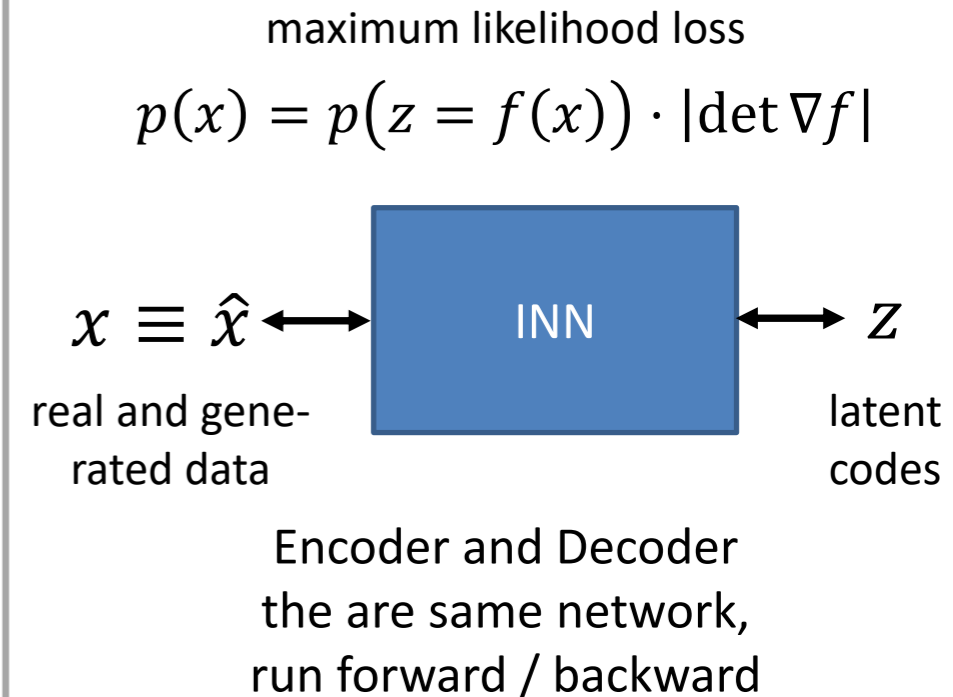
(Variational) Autoencoders



lossy encoding / decoding

Normalizing Flows

(Invertible Neural Networks, INNs)



lossless encoding / decoding

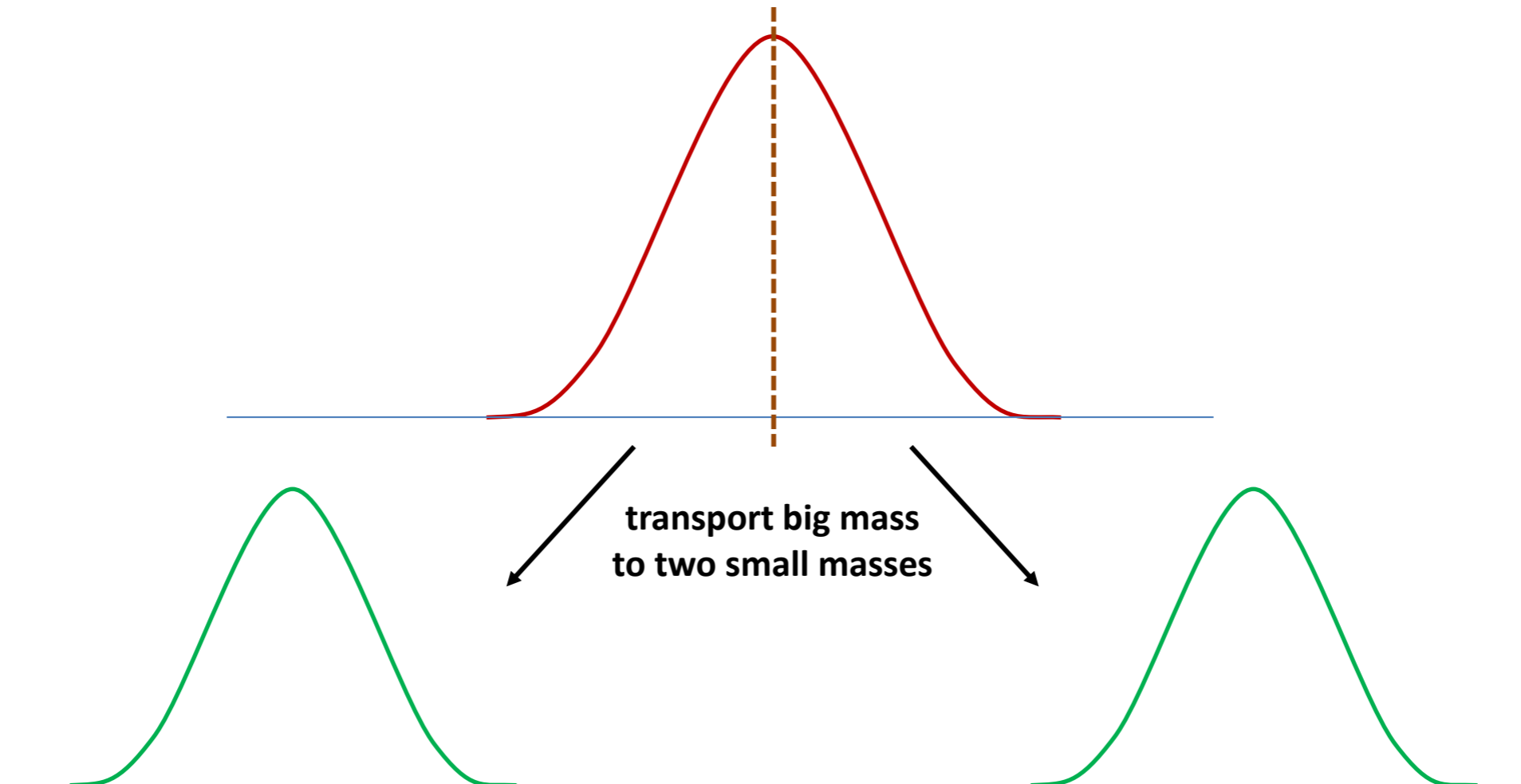




Normalizing flows

Model complicated probabilities as bijective mappings of simple ones

- Example: transport (“flow”) from simple “sand pile” to target „sand piles“

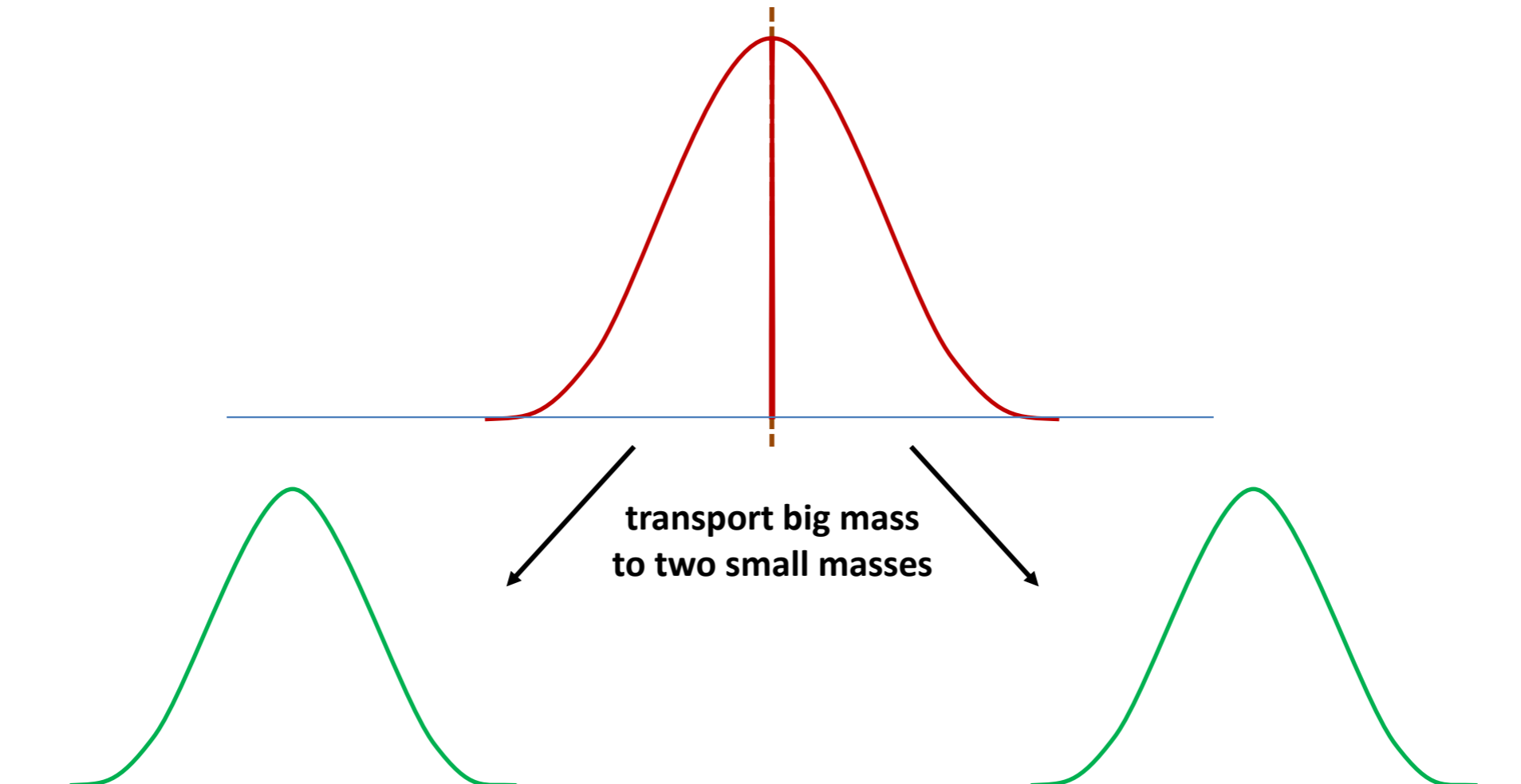




Normalizing flows

Model complicated probabilities as bijective mappings of simple ones

- Example: transport (“flow”) from simple “sand pile” to target „sand piles“

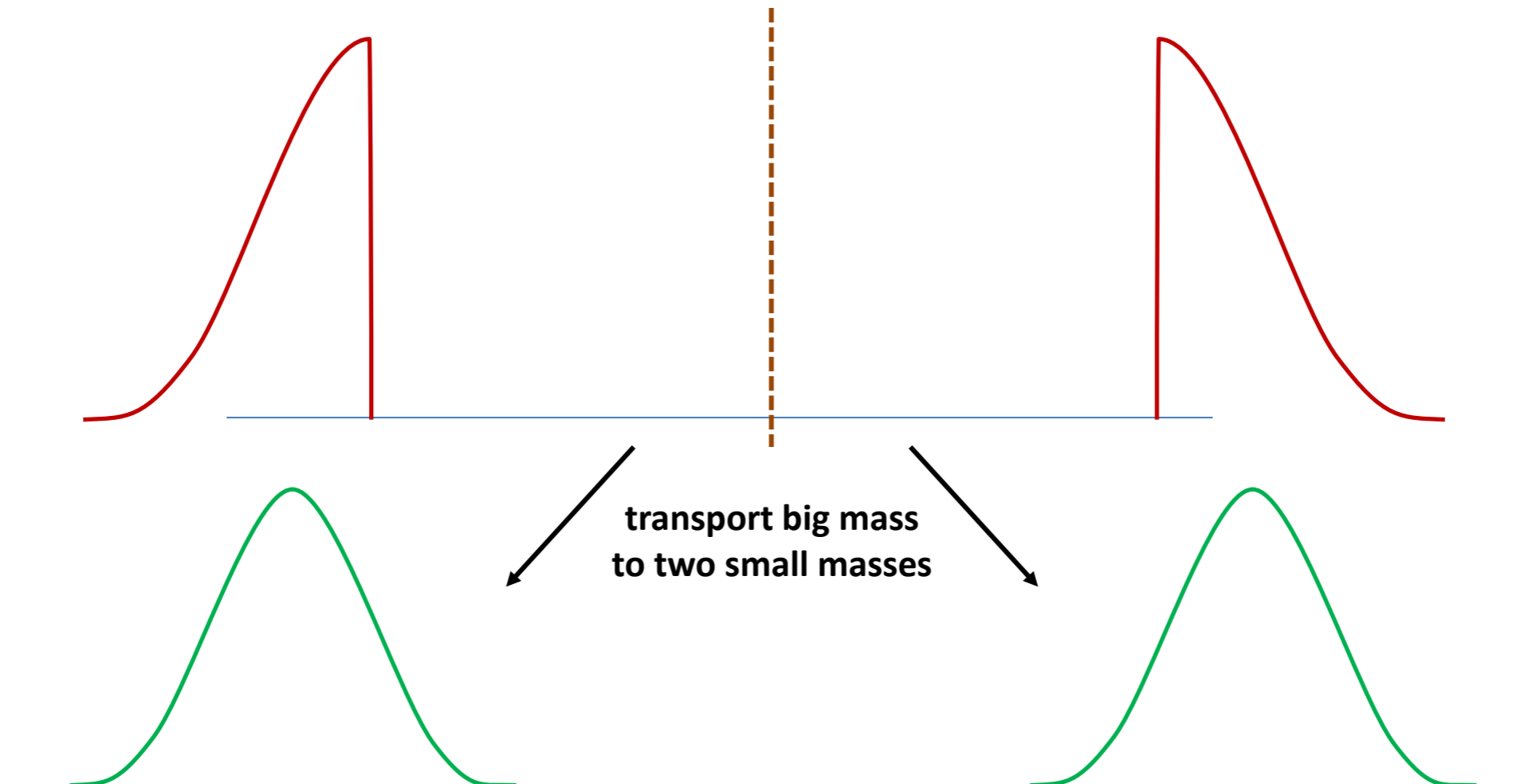




Normalizing flows

Model complicated probabilities as bijective mappings of simple ones

- Example: transport (“flow”) from simple “sand pile” to target „sand piles“

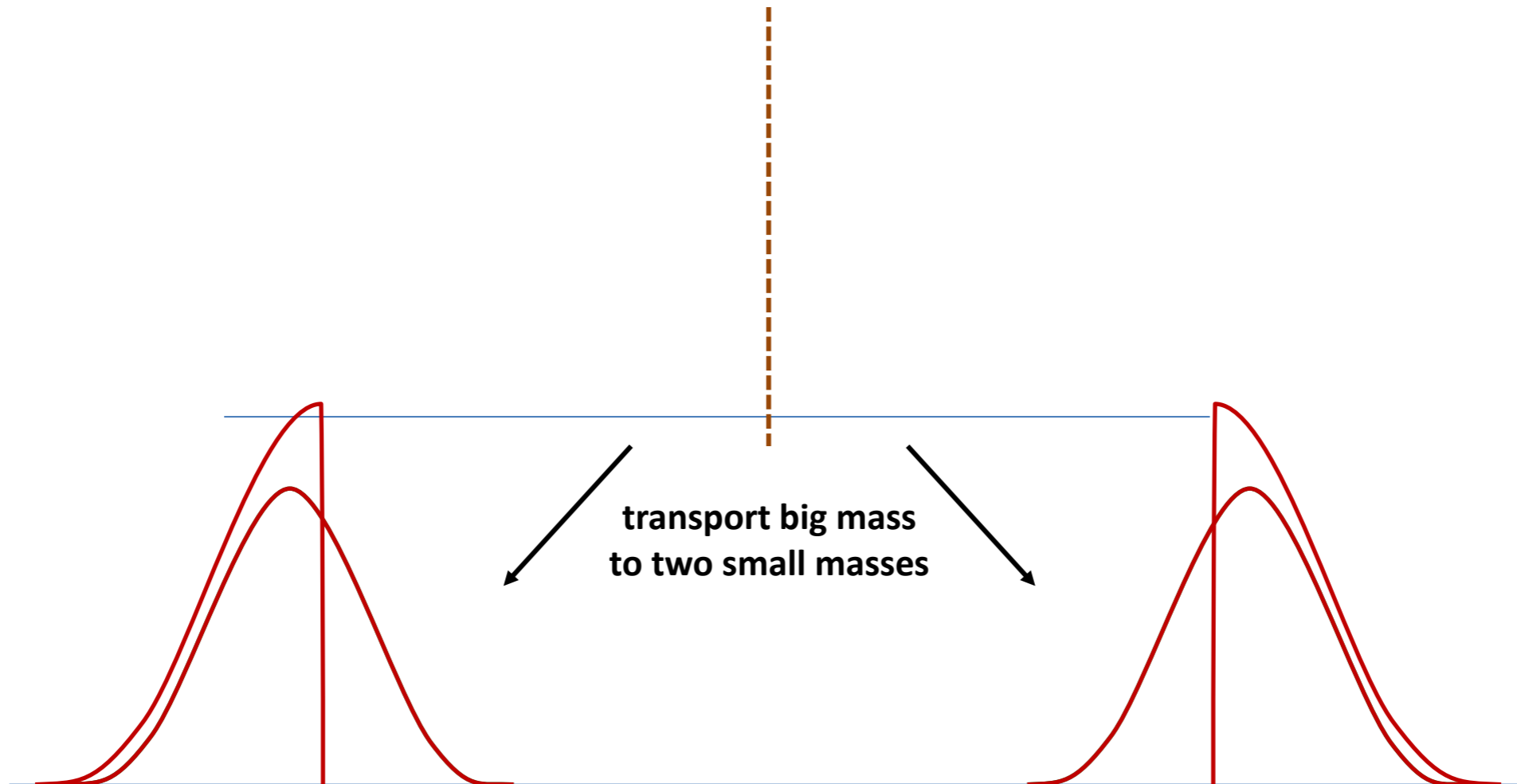




Normalizing flows

Model complicated probabilities as bijective mappings of simple ones

- Example: transport (“flow”) from simple “sand pile” to target „sand piles“

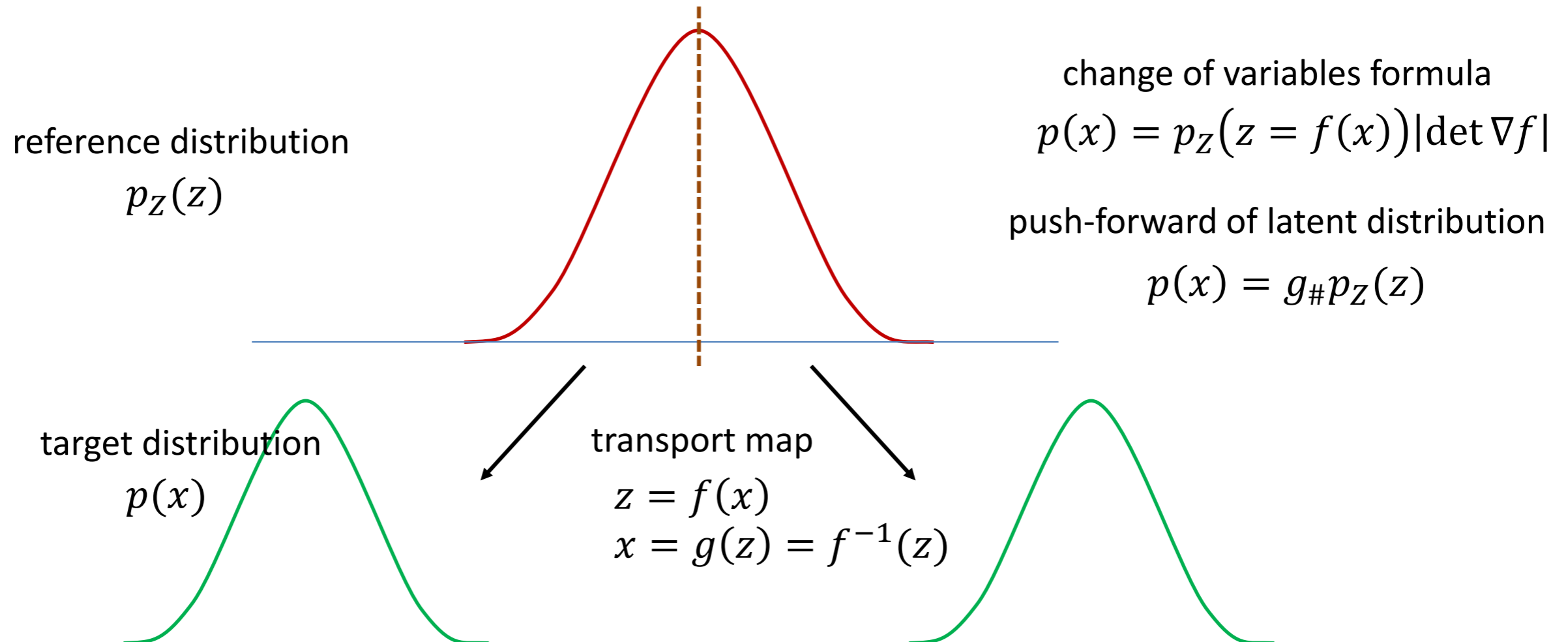




Normalizing flows

Model complicated probabilities as bijective mappings of simple ones

- Mathematically: target distribution is a push-forward of reference distribution





Multiple Possibilities for Normalizing Flows

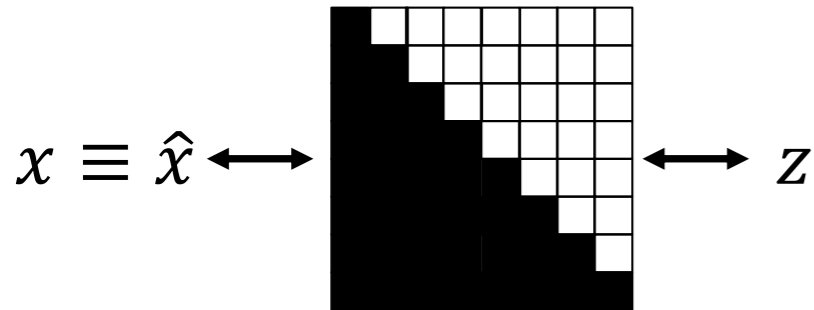
Autoregressive Models

Chain rule decomposition:

$$p(x_1, \dots, x_D) = \prod_i p_i(x_i | x_{<i})$$

triangular reparameterization:

$$\forall i: x_i = f_i(z_i, x_{<i}) \text{ monoton.}$$



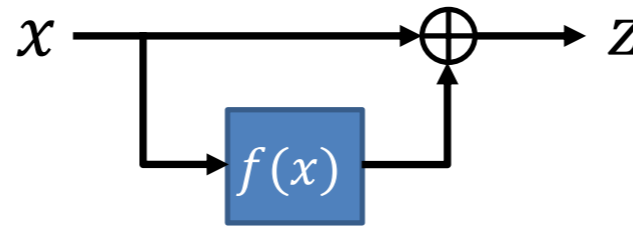
inverse direction inefficient

⇒ use two complementary nets

example: parallel WaveNet

iResNets (invertible residual networks)

Residual block:



$$z = x + f(x)$$

is invertible when

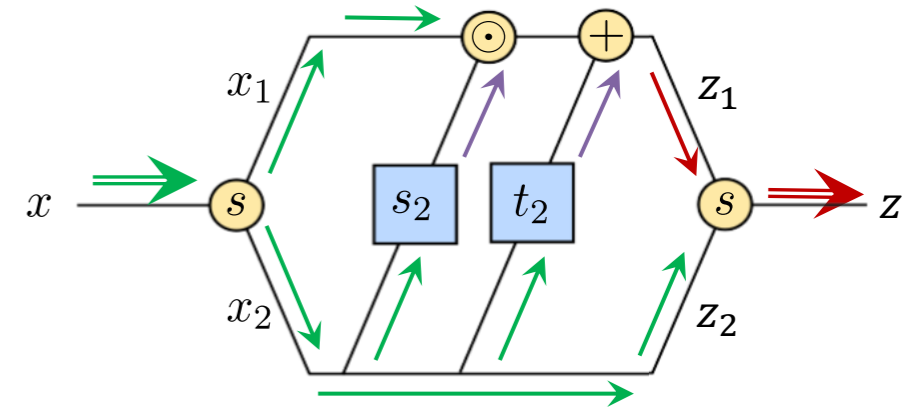
$$\|f(x)\|_{\text{Lipshitz}} < 1$$

inverse direction is reasonably efficient (fixpoint or Newton iterations)

example: Residual Flow Net

RealNVP

Affine coupling layer:



$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} x_1 \cdot s_2(x_2) + t_2(x_2) \\ x_2 \end{bmatrix}$$

inverse is equally efficient:

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} (z_1 - t_2(z_2)) / s(z_2) \\ z_2 \end{bmatrix}$$

example: GLOW

How do you make ResNets invertible and why would you care?

Ullrich Köthe

Visual Learning Lab, Heidelberg University

Tutorial „Normalizing Flows“ at CVPR 2021

June 2021



**UNIVERSITÄT
HEIDELBERG**
ZUKUNFT
SEIT 1386



European Research Council
Established by the European Commission



Recap: What is a ResNet?

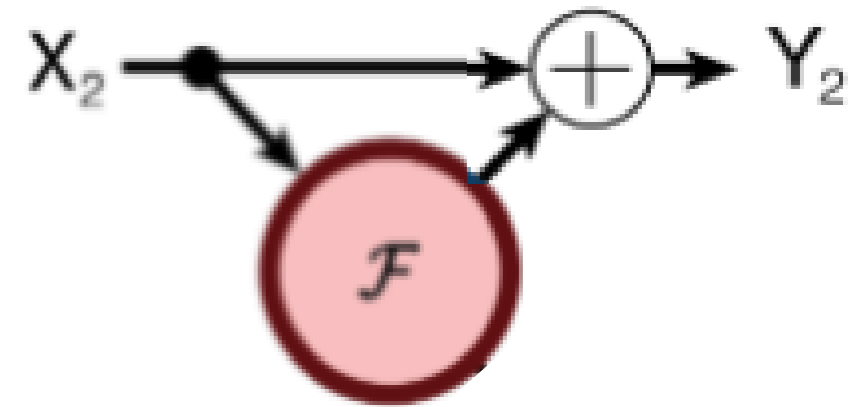
- Instead of modeling the transition from layer l to $l + 1$

$$z_{l+1} = \mathcal{F}_l(z_l)$$

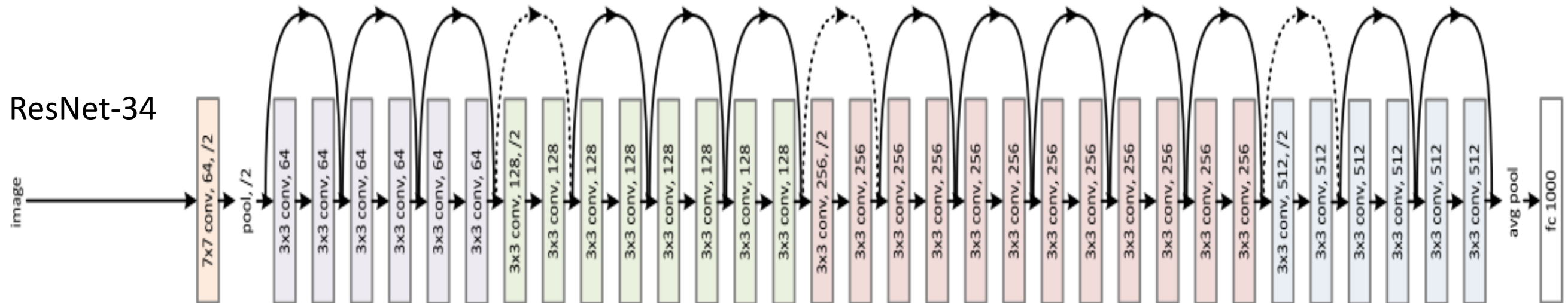
model the difference (residual) between consecutive layers

$$z_{l+1} - z_l = \mathcal{F}_l(z_l) \Leftrightarrow z_{l+1} = z_l + \mathcal{F}_l(z_l)$$

- Each layer (“residual block”) consists of a skip connection and a parallel feed-forward transformation
- Advantage: no vanishing gradients even for very deep networks



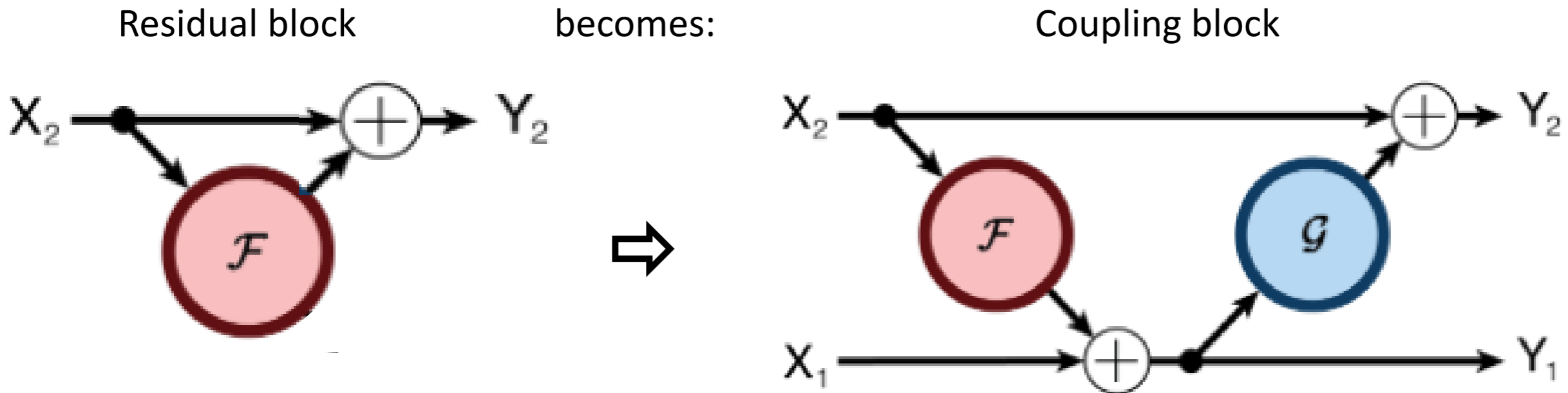
residual block





RevNets: Memory-efficient backpropagation

- Simple application of coupling layers: replace residual blocks with coupling blocks
 - Do not store activations during the forward pass of training
 - Recompute them on the fly during backpropagation, using the invertible architecture





RevNets: Memory-efficient backpropagation

- Simple application of coupling layers: replace residual blocks with coupling blocks
 - Do not store activations during the forward pass of training
 - Recompute them on the fly during backpropagation, using the invertible architecture

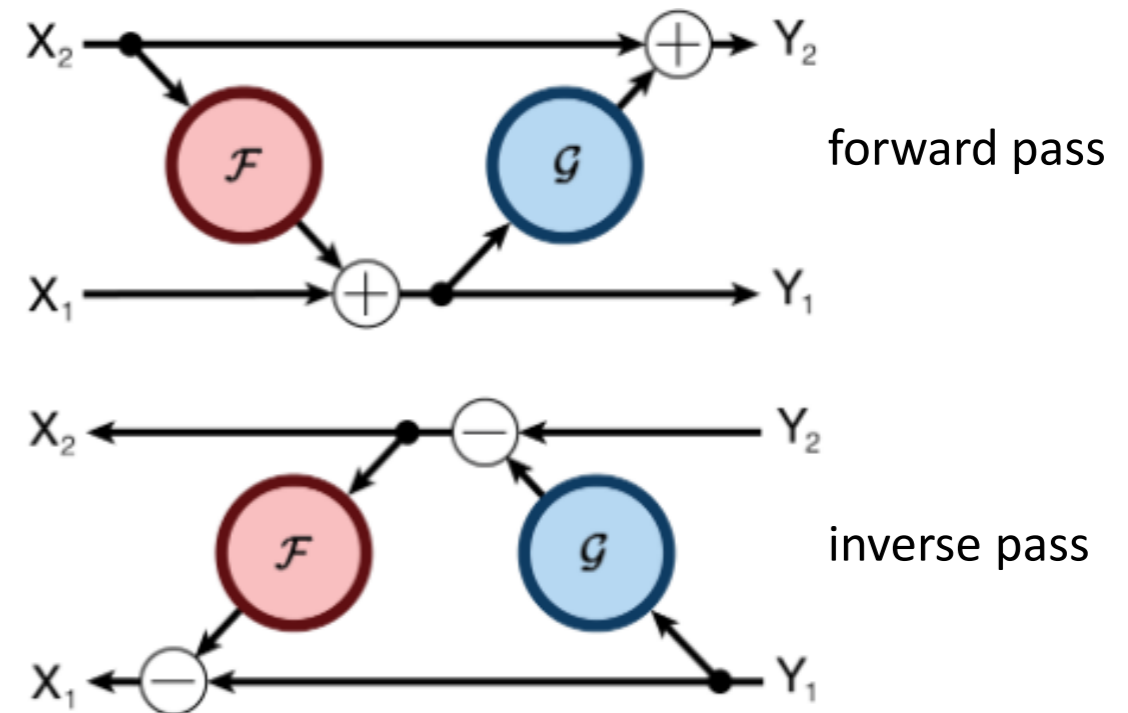
Algorithm 1 Reversible Residual Block Backprop

```

1: function BLOCKREVERSE( $(y_1, y_2), (\bar{y}_1, \bar{y}_2)$ )
2:    $z_1 \leftarrow y_1$ 
3:    $x_2 \leftarrow y_2 - \mathcal{G}(z_1)$ 
4:    $x_1 \leftarrow z_1 - \mathcal{F}(x_2)$ 
5:    $\bar{z}_1 \leftarrow \bar{y}_1 + \left(\frac{\partial \mathcal{G}}{\partial z_1}\right)^\top \bar{y}_2$ 
6:    $\bar{x}_2 \leftarrow \bar{y}_2 + \left(\frac{\partial \mathcal{F}}{\partial x_2}\right)^\top \bar{z}_1$ 
7:    $\bar{x}_1 \leftarrow \bar{z}_1$ 
8:    $\bar{w}_{\mathcal{F}} \leftarrow \left(\frac{\partial \mathcal{F}}{\partial w_{\mathcal{F}}}\right)^\top \bar{z}_1$ 
9:    $\bar{w}_{\mathcal{G}} \leftarrow \left(\frac{\partial \mathcal{G}}{\partial w_{\mathcal{G}}}\right)^\top \bar{y}_2$ 
10:  return  $(x_1, x_2)$  and  $(\bar{x}_1, \bar{x}_2)$  and  $(\bar{w}_{\mathcal{F}}, \bar{w}_{\mathcal{G}})$ 
11: end function
  
```

} inverse pass

“isolated” autodiffs of G and F





RevNets: Memory-efficient backpropagation

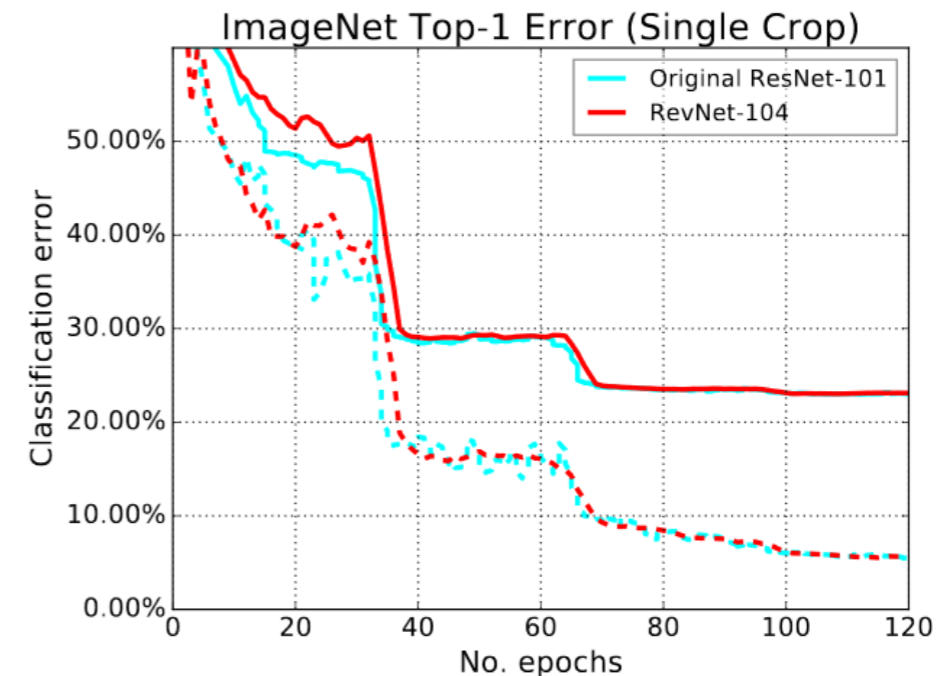
- Performance example: ResNet-101 vs. RevNet-104 on ImageNet

Dataset	Version	Params (M)	Units	Parameter Cost	Activation Cost
ImageNet	ResNet-101	44.5	3-4-23-3	~ 178MB	~ 5250MB
ImageNet	RevNet-104	45.2	2-2-11-2	~ 180MB	~ 1440MB

- Very similar behavior:

Top-1 classification error	
ResNet-101	RevNet-104
23.01%	23.10%

- Trade-off: greatly reduces memory consumption for 2-4 times the compute





Application: i-RIM 3D

- Allows training of very big nets: 3-dimensional convolutions, many layers
 - fastMRI Challenge: MRI reconstruction from 8x less raw data

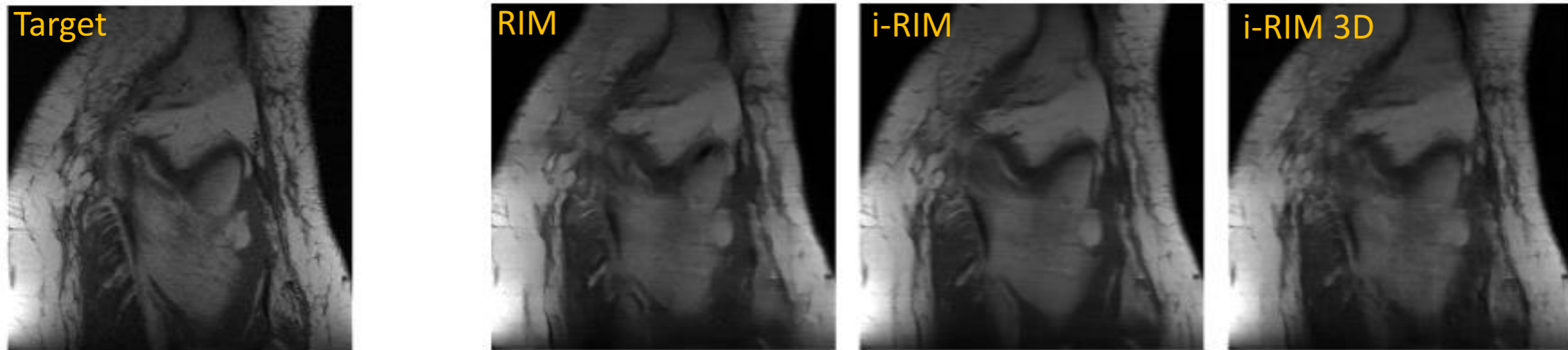


Table 1: Comparison of memory consumption during training and testing.

	RIM	i-RIM 2D	i-RIM 3D
Size Machine State (η, s) (CDHW)	$130 \times 1 \times 480 \times 320$	$64 \times 1 \times 480 \times 320$	$64 \times 32 \times 480 \times 320$
Memory Machine State (η, s) (in GB)	0.079	0.039	1.258
Number of steps T	1/4/8	1/4/8	1/4/8
Network Depth (#layers)	5/20/40	50/200/400	50/200/400
Memory during Testing (in GB)	0.60 / 0.65 / 0.65	0.20 / 0.24 / 0.31	5.87/6.03 / 6.25
Memory during Training (in GB)	2.65 / 6.01 / 10.49	2.47 / 2.49 / 2.51	11.51 / 11.76 / 11.89



Making ResNets Invertible: i-ResNets and Residual Flows

- Can one create an invertible network while keeping the original ResNet architecture?

$$\mathbf{x}_{t+1} = \mathbf{x}_t + g_{\theta_t}(\mathbf{x}_t) \quad \text{forward pass}$$

- How to ensure a bijective mapping?
- How to compute the inverse efficiently?
- How to perform maximum likelihood training?

- The mapping is guaranteed to be bijective if $\frac{\partial \mathbf{x}_{t+1}}{\partial \mathbf{x}_t} > 0$

- Sufficient condition: Lipschitz bound on g_{θ_t} : $\|g_{\theta_t}(x_t^{(1)}) - g_{\theta_t}(x_t^{(2)})\| \leq \lambda \|x_t^{(1)} - x_t^{(2)}\|$ with $\lambda < 1$

⇒ Expressive power of each block is limited, need more blocks

⇒ Blocks can be inverted using fixed point iterations or Newton iterations:

$$\begin{aligned} \mathbf{x}_t^0 &= \mathbf{x}_{t+1} \\ \mathbf{x}_t^{i+1} &= \mathbf{x}_{t+1} - g_{\theta_t}(\mathbf{x}_t^i) \end{aligned} \quad \text{backward pass}$$



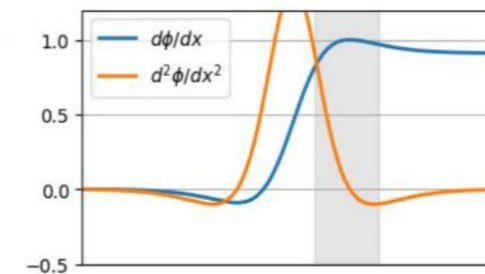
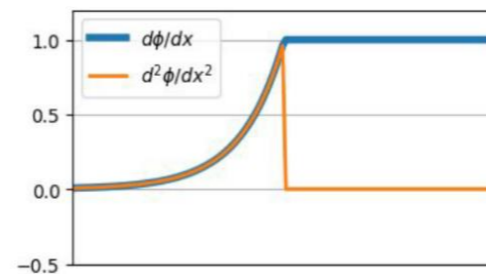
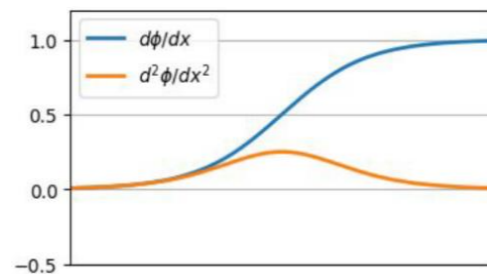
Making ResNets Invertible: i-ResNets and Residual Flows

- How to achieve the Lipschitz bound?
 - Concatenation is Lipschitz, when each transition is so
 - Linear/convolutional layers: **normalize weight matrices** with $c < 1$ and **largest singular** value $\tilde{\sigma}_i \leq \|W_i\|_2$ estimated by (one iteration of) power method

$$\tilde{W}_i = \begin{cases} cW_i/\tilde{\sigma}_i, & \text{if } c/\tilde{\sigma}_i < 1 \\ W_i, & \text{else} \end{cases}$$

Randomly initialise \vec{x}_0
 for $i = 1$ to n do
 $\vec{x}_i \leftarrow W^T W \vec{x}_{i-1}$
 end for
 $\sigma_{max} \leftarrow \frac{\|W\vec{x}_n\|_2}{\|\vec{x}_n\|_2}$

- Activation function: $\forall x: |\phi'(r)| \leq 1$ is fulfilled by many $\phi(r)$, but training involves the **gradient** of the log-determinant of the **Jacobian** (the **first** derivative), i.e. the **second** derivative $\phi''(r)$
- Many common $\phi(r)$ have $\phi'(r) \approx 1 \Rightarrow \phi''(r) \approx 0$, i.e. suffer from **vanishing gradients**
- \Rightarrow **Choose $\phi(r) = \text{LipSwish}(r) = 0.909 r / (1 + \exp(-\beta r))$**





Making ResNets Invertible: i-ResNets and Residual Flows

- How to perform maximum likelihood training?
 - Need the gradient of the log-determinant of the Jacobian
 - Approximate via truncated power series or unbiased log density estimator

$$\ln \left| \det \left(I + J_{g_{\theta_t}}(\mathbf{x}_t) \right) \right| \approx \sum_{k=1}^n (-1)^{k+1} \frac{\text{tr} \left(J_{g_{\theta_t}}(\mathbf{x}_t)^k \right)}{k} \approx \mathbb{E}_{n,v} \left[\sum_{k=1}^n \frac{(-1)^{k+1}}{k} \frac{v^T [J_g(x)]^k v}{\mathbb{P}(N \geq k)} \right]$$

- Very recent new possibility: use relative gradient (i.e. multiplicative instead of additive perturbation)
- ⇒ Gradient update calculation reduces to matrix-vector products (try on your own risk :-)

$$W_t \leftarrow W_t + \gamma (\mathbf{x}_{t-1} (\delta_t^T W_t^T) + \mathbb{I}) W_t$$

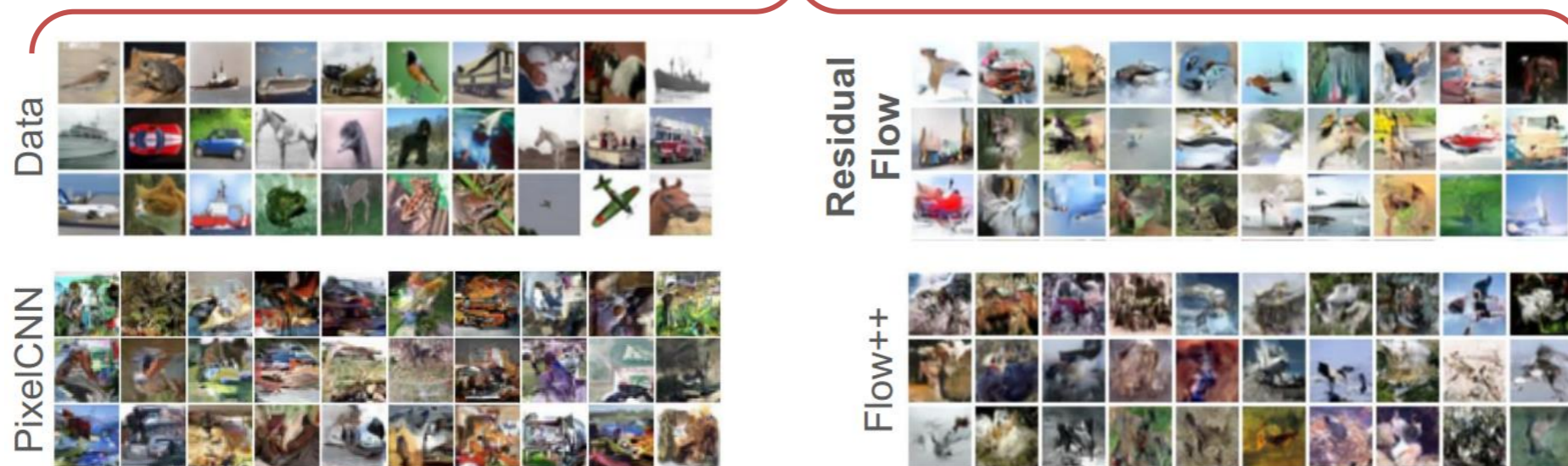




Making ResNets Invertible: i-ResNets and Residual Flows

- Improvements of Residual Flow over i-ResNet apparent visually and in the numbers

Model	MNIST	CIFAR-10	ImageNet 32×32	ImageNet 64×64
Real NVP (Dinh et al., 2017)	1.06	3.49	4.28	3.98
Glow (Kingma and Dhariwal, 2018)	1.05	3.35	4.09	3.81
FFJORD (Grathwohl et al., 2019)	0.99	3.40	—	—
Flow++ (Ho et al., 2019)	—	3.29 (3.09)	— (3.86)	— (3.69)
i-ResNet (Behrmann et al., 2019)	1.05	3.45	—	—
Residual Flow	0.97	3.29	4.02	3.78



RealNVP: Invertibility via Coupling Layers

Ullrich Köthe

Visual Learning Lab, Heidelberg University

Tutorial „Normalizing Flows“ at CVPR 2021

June 2021



**UNIVERSITÄT
HEIDELBERG**
ZUKUNFT
SEIT 1386



European Research Council
Established by the European Commission



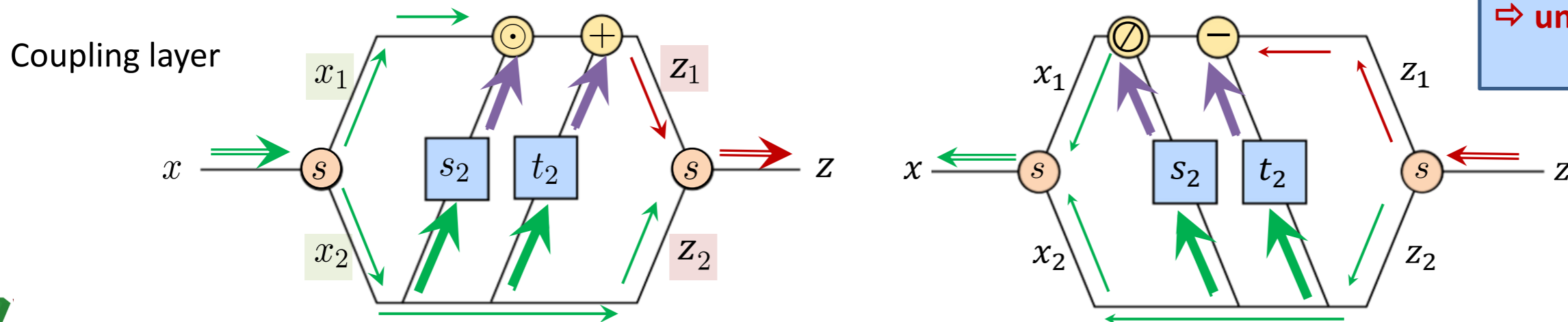
Invertible Neural Networks (INNs) with Coupling Layers

Powerful generative models: RealNVP („non-volume preserving“) [Dinh et al. 2017]

- Network is a sequence of *affine coupling layers*
 - Each coupling layer **splits** its input $x \in \mathbb{R}^D$ into two halves $x_1, x_2 \in \mathbb{R}^{D/2}$
 - Upper half is subjected to an affine transformation \Rightarrow outputs $z_1, z_2 \in \mathbb{R}^{D/2}$
 - Affine coefficients are computed by standard fully connected or convolutional networks
- $s_2 \in \mathbb{R}_+^{D/2}$ and $t_2 \in \mathbb{R}^{D/2}$ from the lower half's data

Forward computation: $z_1 = x_1 \odot s_2(x_2) + t_2(x_2), \quad z_2 = x_2$

Inverse computation: $x_1 = (z_1 - t_2(z_2)) \oslash s_2(z_2), \quad x_2 = z_2$

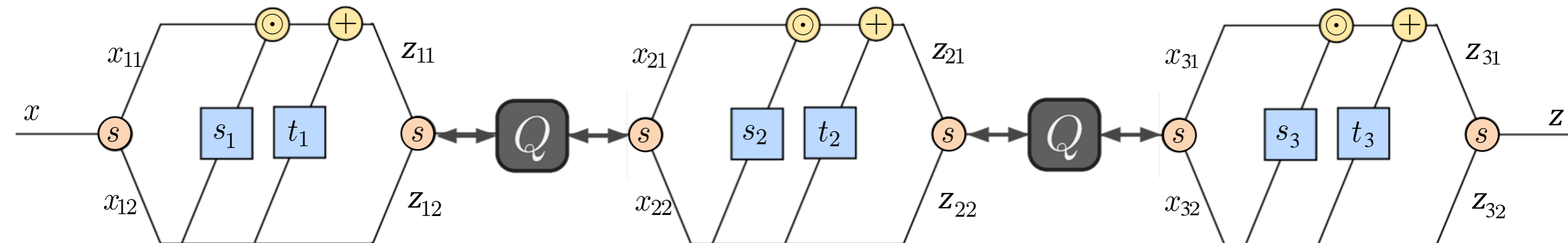


nested functions
 s_2 and t_2 are
 always executed in
 the same direction
 \Rightarrow unrestricted neural
 networks



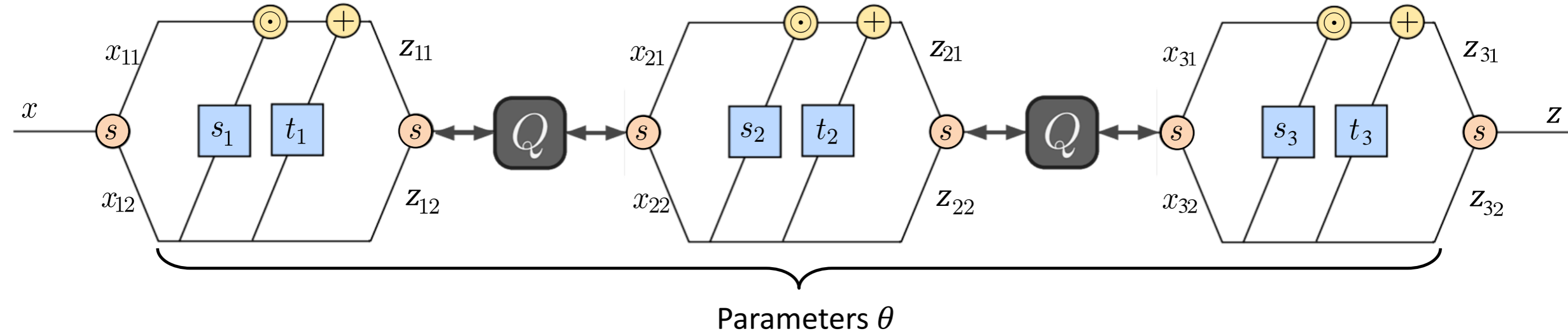
Deep INN

- Concatenate many coupling layers
- Alternate with orthogonal layers Q
 - ⇒ Active (upper lane) and passive (lower lane) dimensions change in each layer
 - Random permutations or projections are good enough, learning Q is not necessary
- Surprisingly powerful despite its simplicity
- Similar to autoencoder: forward mode = encoder, backward mode = decoder
 - Encoder and decoder are merged into a single network
 - Lossless encoding due to invertibility (no bottleneck)





Training Deep INNs with Maximum Likelihood Loss



Tractable data likelihood via change-of variables formula: $p_{\theta}(x) = p_Z(z = f_{\theta}(x)) \cdot |\det \nabla f_{\theta}(x)|$

⇒ Negative log-likelihood has especially simple form when $p_Z(z)$ is standard normal

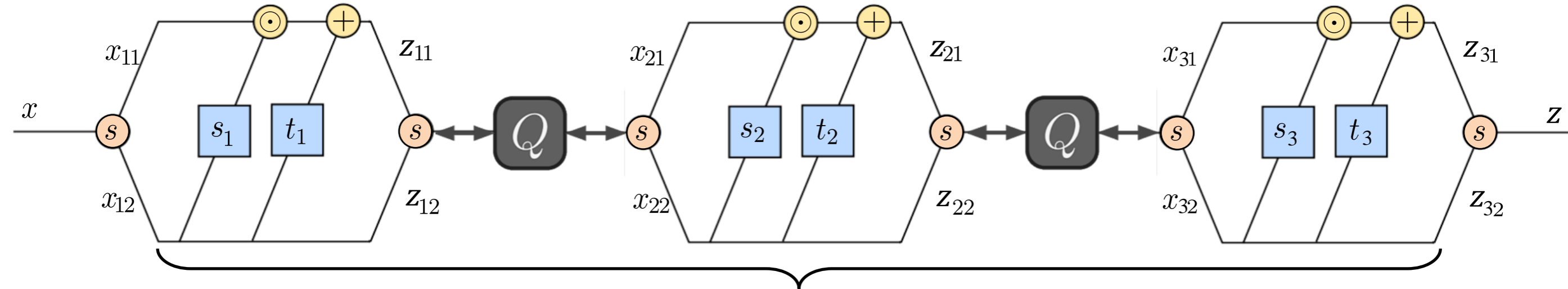
$$\begin{aligned}
 -\log p_{\theta}(x) &= -\log p_Z(z = f_{\theta}(x)) - \log |\det \nabla f_{\theta}(x)| \\
 &= \frac{D}{2} \log 2\pi + \frac{1}{2} \|f_{\theta}(x)\|_2^2 - \sum_l \text{sum}(\log s_{\theta,l}(x_{l2}))
 \end{aligned}$$

with $s_{\theta,l}(x_{l2})$ the multipliers at coupling layer l (note: $\log \det Q = 0$)





Training Deep INNs with Maximum Likelihood Loss



Parameters θ

Negative log-likelihood:
$$-\log p_{\theta}(x) = \frac{D}{2} \log 2\pi + \frac{1}{2} \|f_{\theta}(x)\|_2^2 - \sum_l \text{sum}(\log s_{\theta,l}(x_{l2}))$$

\Rightarrow Train by minimizing the NLL objective over training set $\{x^{(i)}\}_{i=1}^N$:

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} p_{\theta}(\{x^{(i)}\}_{i=1}^N) = \arg \max_{\theta} \prod_{i=1}^N p_{\theta}(x^{(i)}) = \arg \min_{\theta} \sum_{i=1}^N -\log p_{\theta}(x^{(i)}) \\ &= \arg \min_{\theta} \sum_{i=1}^N \left(\frac{1}{2} \|f_{\theta}(x^{(i)})\|_2^2 - \sum_l \text{sum}(\log s_{\theta,l}(x_{l2}^{(i)})) \right) \end{aligned}$$



Conditional Modeling with INNs

- In practice, we often need to model conditionals $p(x | y)$ or $p(y | x)$ rather than $p(x)$
- Example: Generative classification
 - x are features, y are class labels
 - determine posterior $p(y | x)$ using Bayes rule
$$p(y | x) \sim p(x | y) p(y)$$
 - ⇒ learn the likelihood $p(x | y)$ via a conditional INN (specifically, an IB-INN)
- Example: Solving inverse problems
 - x are hidden system parameters, y are observations of the system behavior
 - determine the posterior $p(x | y = \hat{y})$ to estimate parameters x from measured \hat{y}
 - ⇒ learn $p(x | y)$ using synthetic data from a simulation $y = g(x; \text{noise})$ of the forward process



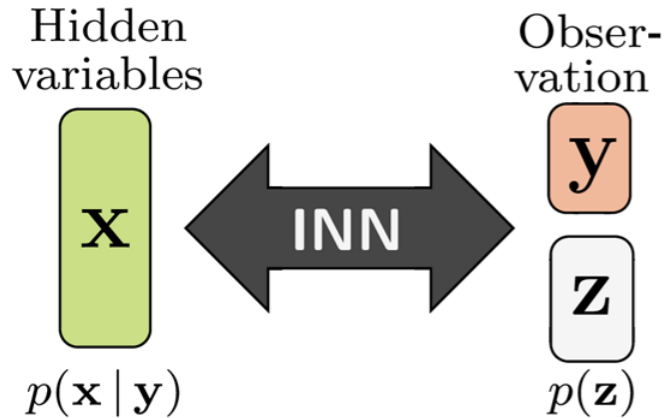


INN Architectures for Conditional Inference

Split latent space

training: $(y, z) = f_{\theta}(x)$
s.t. $p(z) = \mathcal{N}(0, \mathbb{I})$

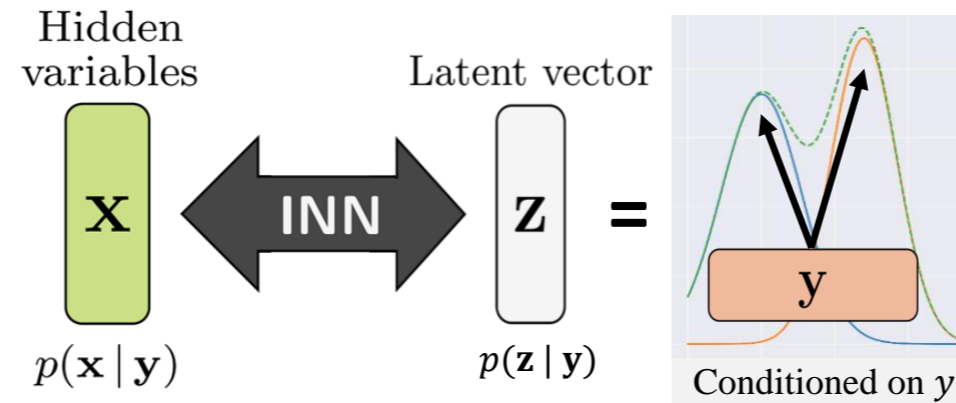
inference: sample $z \sim \mathcal{N}(0, \mathbb{I})$
compute $x = f_{\theta}^{-1}(\hat{y}, z)$
 $\Rightarrow x \sim p(x | \hat{y})$



Latent mixture INN

training: $z = f_{\theta}(x)$
s.t. $p(z) = \text{GMM}(z; y) = \sum_y \mathcal{N}(\mu_y, \Sigma_y)$

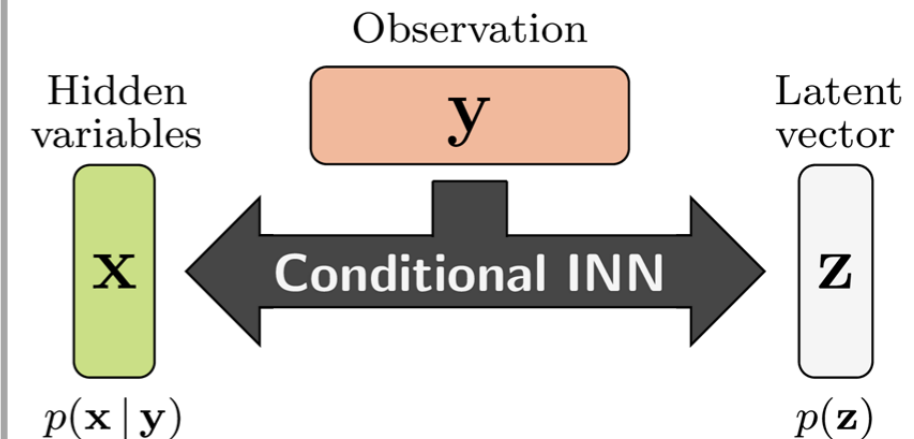
inference: sample $z \sim \mathcal{N}(\mu_{\hat{y}}, \Sigma_{\hat{y}})$
compute $x = f_{\theta}^{-1}(z)$
 $\Rightarrow x \sim p(x | \hat{y})$



Conditional INN

training: $z = f_{\theta}(x; y)$
s.t. $p(z) = \mathcal{N}(0, \mathbb{I})$

inference: sample $z \sim \mathcal{N}(0, \mathbb{I})$
compute $x = f_{\theta}^{-1}(z; \hat{y})$
 $\Rightarrow x \sim p(x | \hat{y})$



historically first

classification, disentanglement

inverse inference



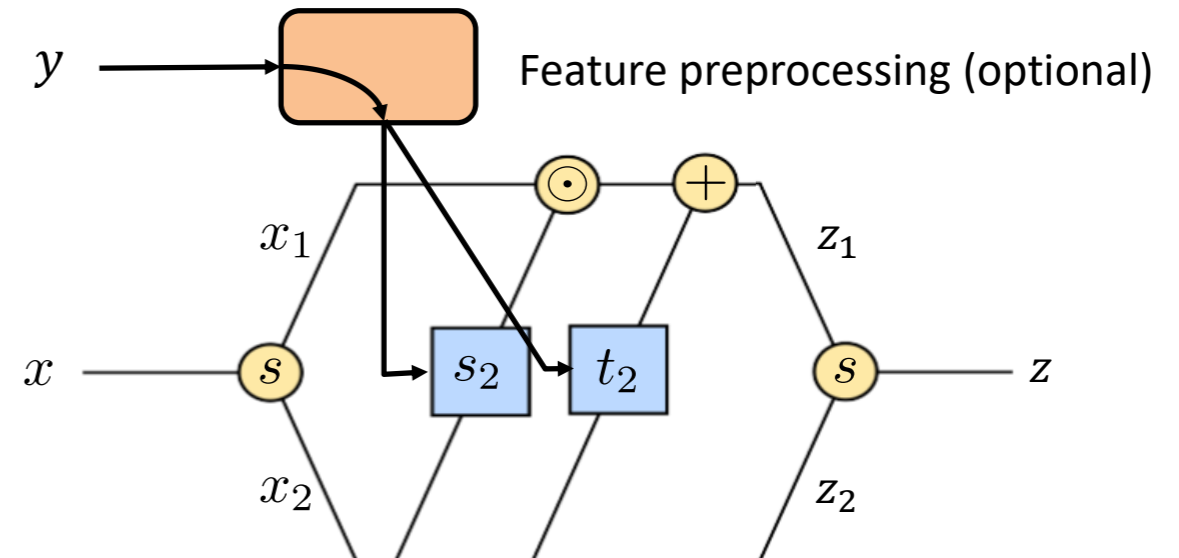
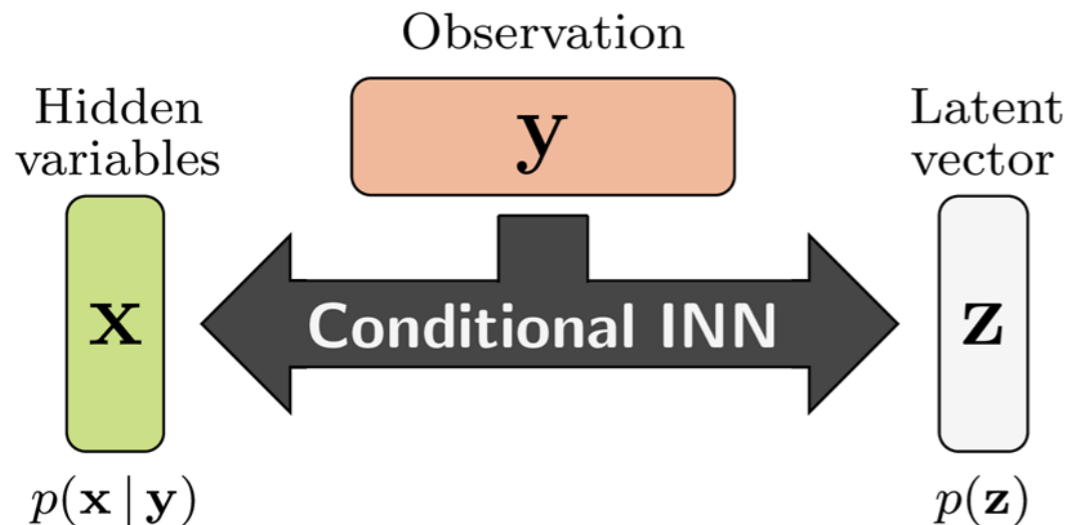
Conditional INN (cINN)

Conditional INN (cINN) adapts vanilla INN for conditional probabilities

- Reparametrize $x \sim p(x | y)$ as $x = g_\theta(z; y)$ with $z \sim p_z(z)$ and forward process $z = f_\theta(x; y) = g_\theta^{-1}(x; y)$
- Minimum log-likelihood loss becomes

$$\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^N \left(\frac{1}{2} \|f_\theta(x^{(i)}; y^{(i)})\|_2^2 - \sum_l \text{sum} \left(\log s_{\theta,l} \left(x_{l2}^{(i)}; y^{(i)} \right) \right) \right)$$

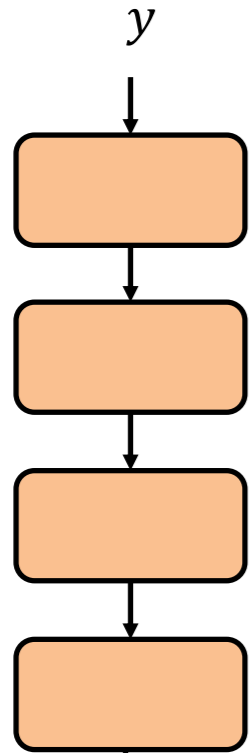
simple change of coupling layer architecture:
feed y as additional input to subnets s, t





cINNs Turn Deterministic Networks into Probabilistic Ones

Deterministic network



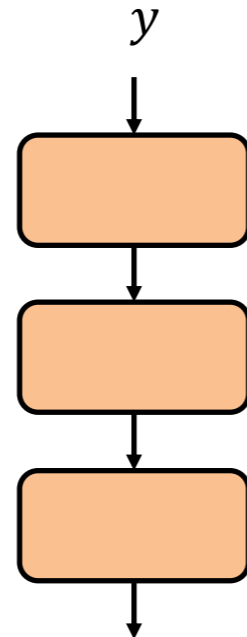
$x = \hat{h}(y)$ single output

$$\text{loss: } \hat{h} = \arg \min \sum_i (h(y_i) - x_i^*)^2$$

ground truth: x^*

⇒

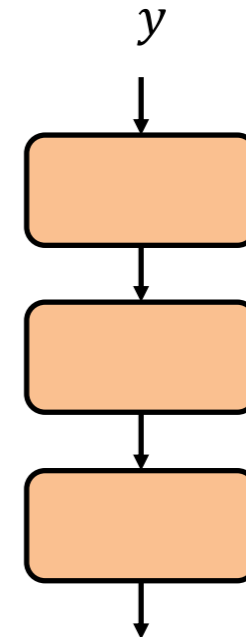
remove final layer(s)



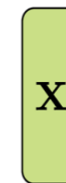
$c = h'(y)$
learned features

⇒

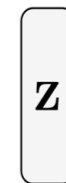
attach cINN



diverse outputs



Conditional INN



$$x \sim p(x | y) \Leftrightarrow x = \hat{g}(z; \hat{h}'(y)) \text{ with } z \sim p(z)$$

$$\text{loss: } \hat{g}, \hat{h}' = \arg \max \sum_i \log p(x_i^* | y_i)$$



Example: Image-to-Image Translation

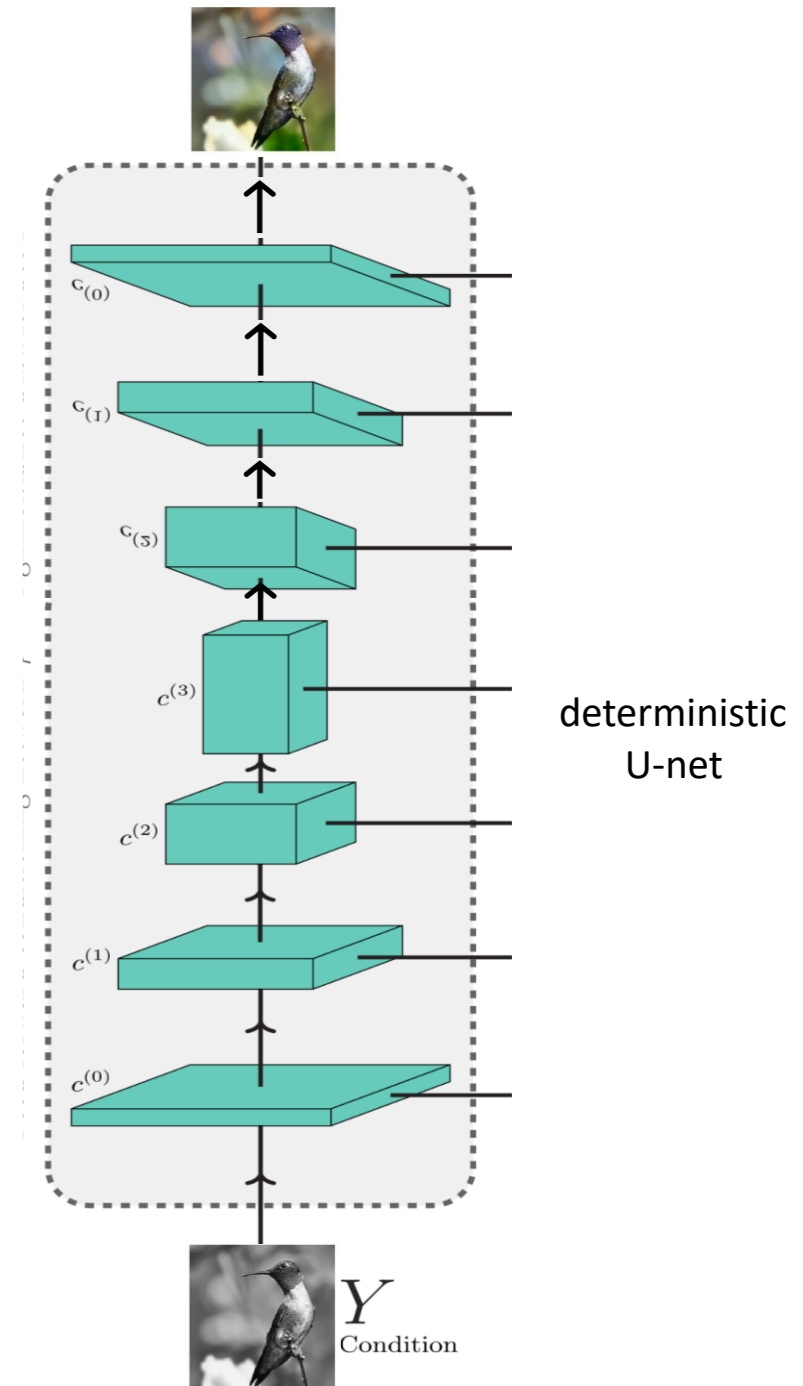
- Colorization as an inverse problem:
 - forward process: turn color image to grayscale by taking the L-channel in Lab color space
 - inverse problem: reconstruct **realistic** color channels
 $y = L \Rightarrow \hat{x} = [a, b]$



y



$[y, \hat{x}]$



- deterministic network: single result



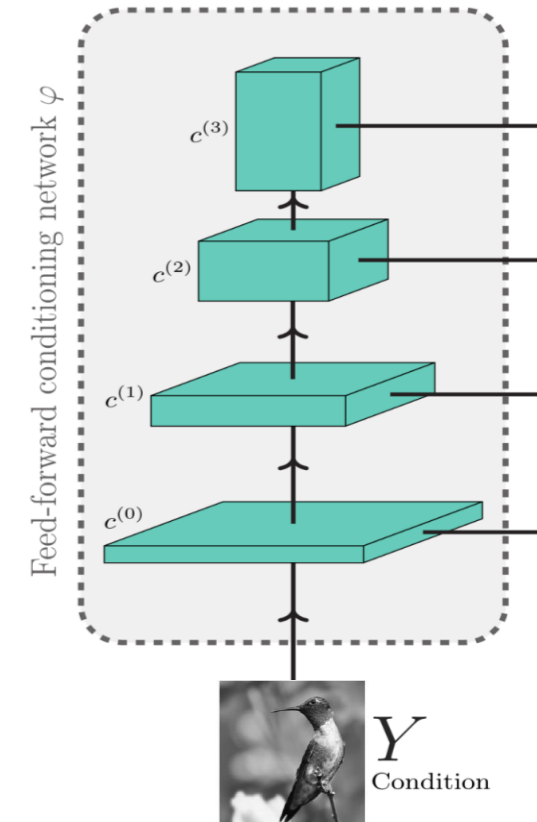
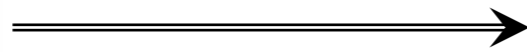


Example: Image-to-Image Translation

- Colorization as an inverse problem:
 - forward process: turn color image to grayscale by taking the L-channel in Lab color space
 - inverse problem: reconstruct **realistic** color channels
- $$y = L \Rightarrow \hat{x} = [a, b]$$



y





Example: Image-to-Image Translation

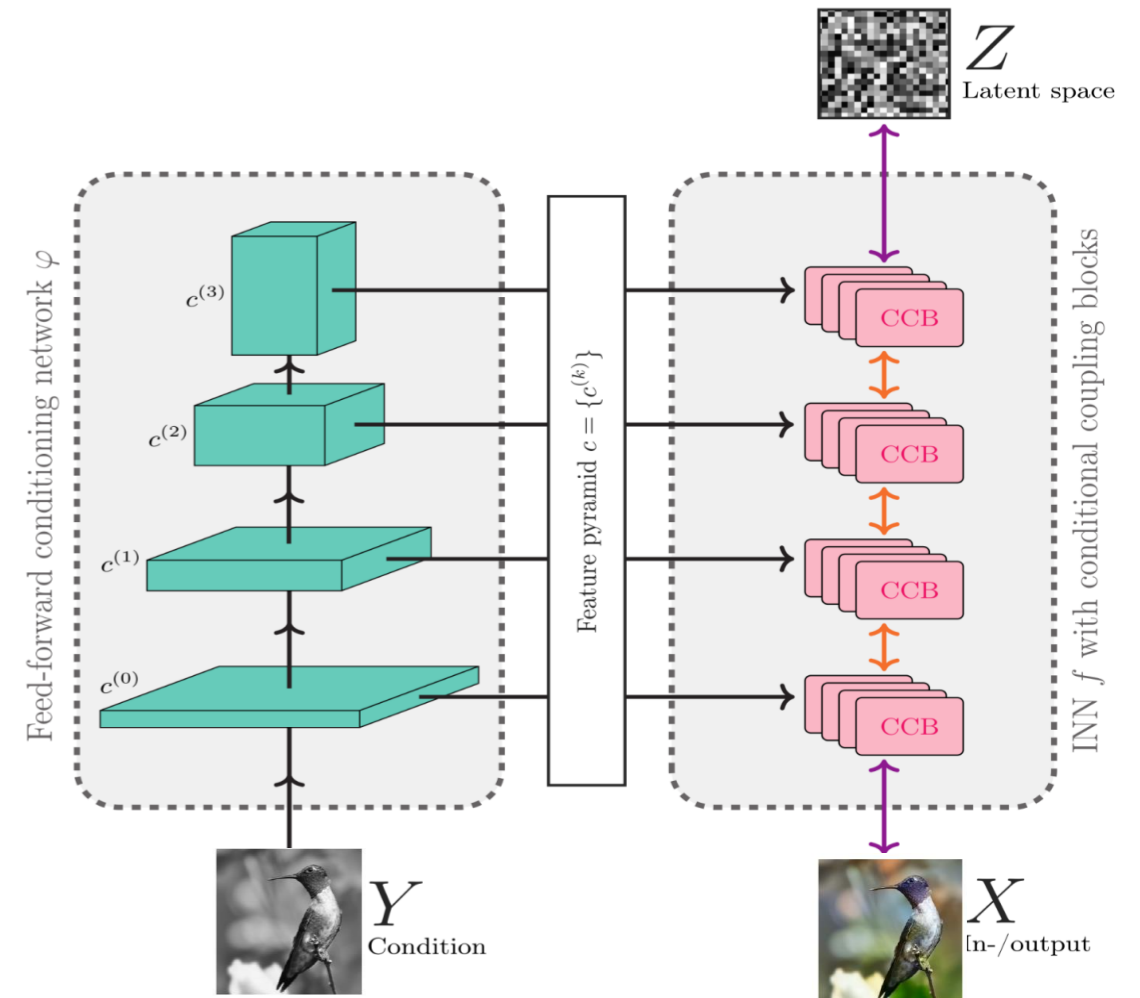
- Colorization as an inverse problem:
 - forward process: turn color image to grayscale by taking the L-channel in Lab color space
 - inverse problem: reconstruct **realistic** color channels
 $y = L \Rightarrow \hat{x} = [a, b]$



y

$\longrightarrow p(x|y)$

- cINN: diverse results





Example: Image-to-Image Translation

- Colorization as an inverse problem:
 - forward process: turn color image to grayscale by taking the L-channel in Lab color space
 - inverse problem: reconstruct **realistic** color channels
 $y = L \Rightarrow \hat{x} = [a, b]$



y

$$z \sim p_z(z)$$

$$x \sim \hat{p}(x | y)$$



- cINN: diverse results
- Quiz: Which color image is the ground-truth?





Example: Image-to-Image Translation

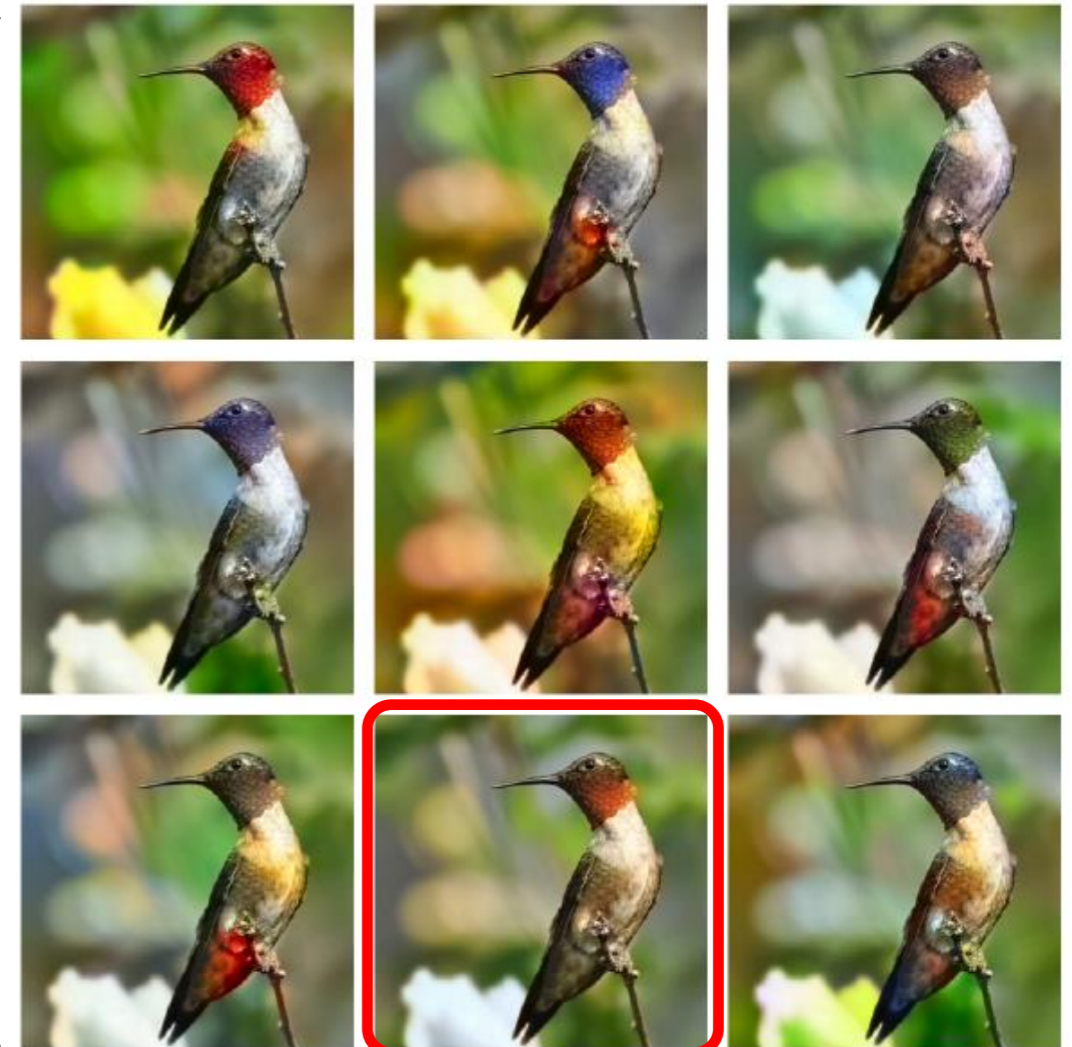
- Colorization as an inverse problem:
 - forward process: turn color image to grayscale by taking the L-channel in Lab color space
 - inverse problem: reconstruct **realistic** color channels
 $y = L \Rightarrow \hat{x} = [a, b]$



y

$$z \sim p_z(z)$$

$$x \sim \hat{p}(x | y)$$



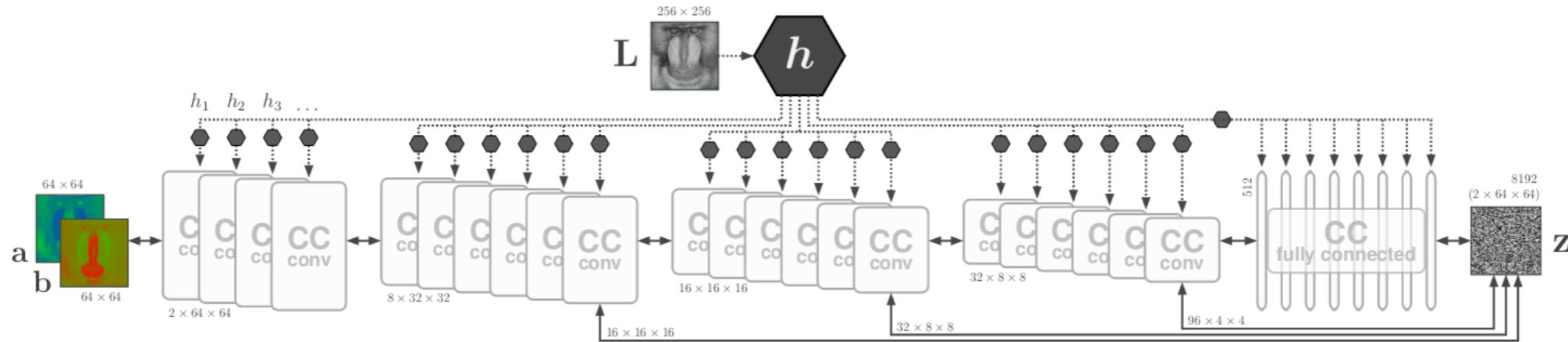
- cINN: diverse results
- Quiz: Which color image is the ground-truth?





cINN Architecture for Colorization

- Four convolutional stacks (with four to six coupling layers)
- Fully connected stack as backend (eight coupling layers)
- Coupling layers separated by random orthogonal matrices to mix channels
- Large feature detection network h (VGG), small conditioning networks h_l



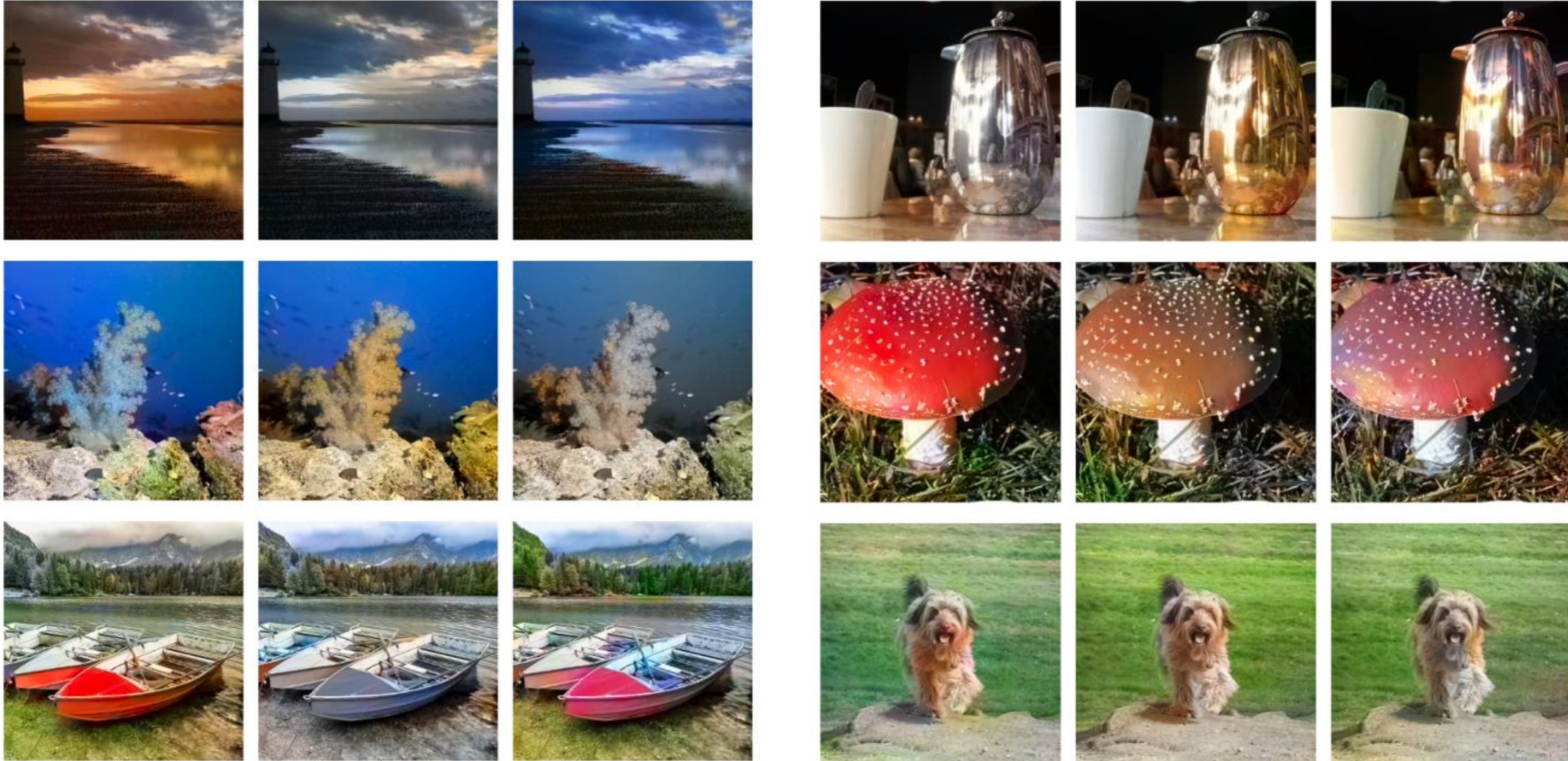
- Multi-scale decomposition via *Haar-Wavelet* down-sampling (standard max pooling not invertible)

$$\underbrace{\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}}_{c \times 2 \times 2} = \left(\underbrace{\begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}}_{\text{average}}, \underbrace{\begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{bmatrix}}_{\text{horizontal}}, \underbrace{\begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}}_{\text{vertical}}, \underbrace{\begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix}}_{\text{diagonal}} \right) \cdot \underbrace{\begin{bmatrix} a \\ h \\ v \\ d \end{bmatrix}}_{4 \cdot c \times 1 \times 1}$$





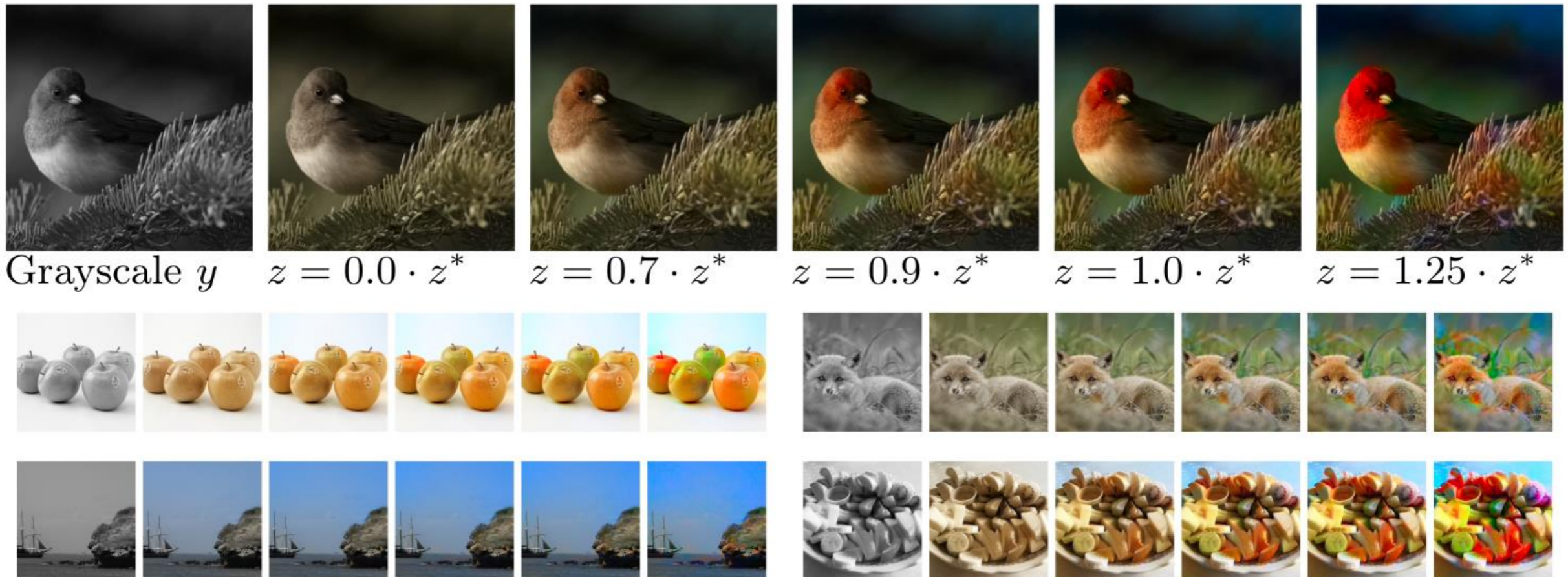
Colorization Examples





Colorization: Meaningful Latent Manipulations

- Magnitude of latent vector encodes color saturation
 - Linear interpolation from $z = 0$ outwards gradually increases saturation





Colorization: Meaningful Latent Manipulations

- Color transfer

- Encode color of input image $i = [L_i, a_i, b_i]$:

$$z_i = f(x = [a_i, b_i]; h'(y = L_i))$$

- Reconstruct color for a different grayscale image L_c : $\hat{x}_i = [\hat{a}_i, \hat{b}_i] = g(z_i; h'(y = L_c))$ with $g = f^{-1}$ while *keeping* the latent code z_i

Inputs
 $[L_i, a_i, b_i]$



New condition L_c

Outputs
 $[L_c, \hat{a}_i, \hat{b}_i]$





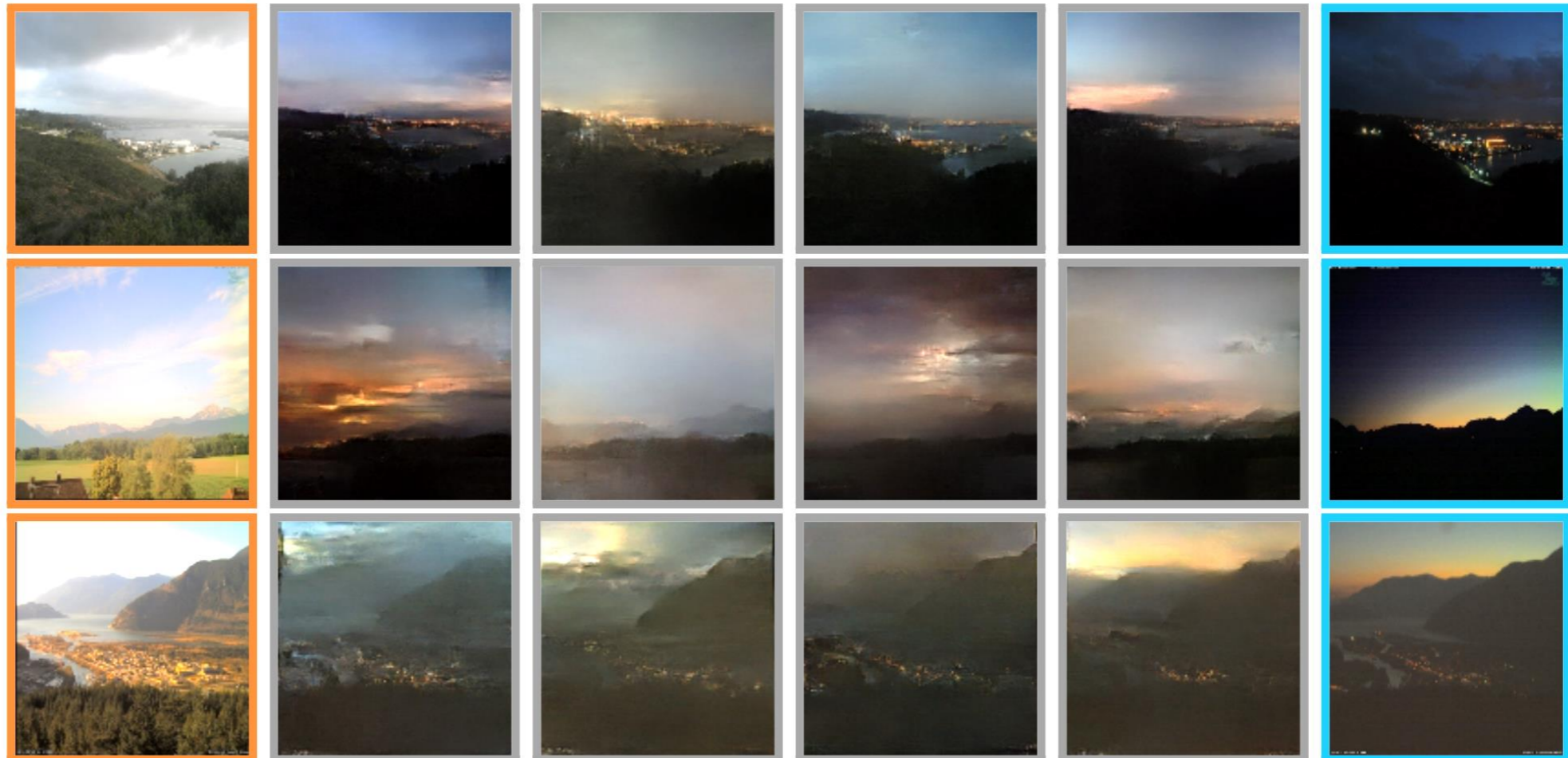
cINN for Image-to-Image Transformation

- Results:

Condition y
Day image






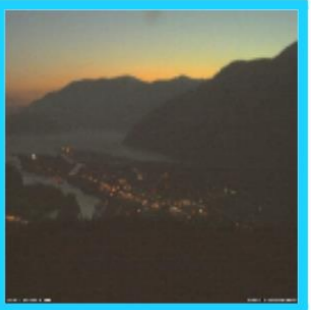
Generated x
Night images

Ground truth x
Night image



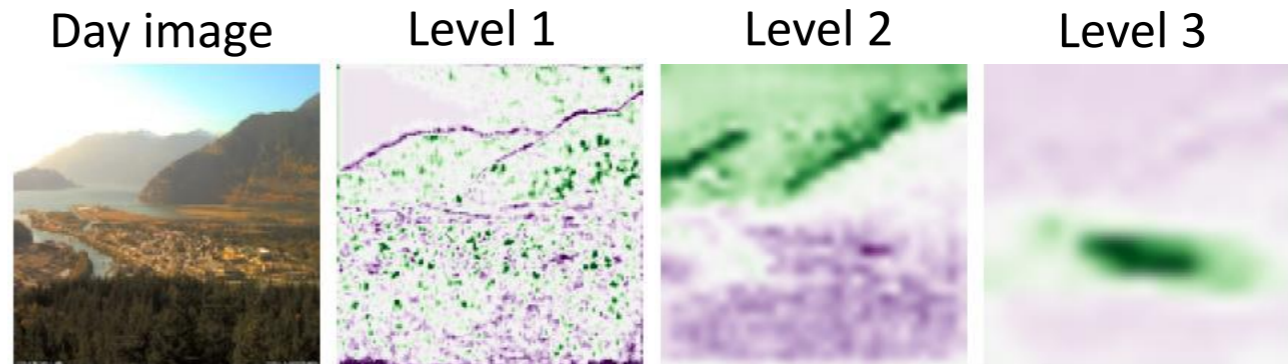


cINN for Image-to-Image Transformation

- Results: **Condition y**
Day image

Generated x
Night images

Ground truth x
Night image


- Multi-scale features learned by the conditioning network:

- Level 1: edges and texture
- Level 2: foreground / background
- Level 3: populated areas (lights!)



Solving Inverse Problems with Invertible Neural Networks

Ullrich Köthe

Visual Learning Lab, Heidelberg University

joint work with **Lynton Ardizzone, Stefan Radev, Jakob Kruse, Tim Adler,
Carsten Rother, Lena Maier-Hein**

Tutorial „Normalizing Flows“ at CVPR 2021

June 2021



**UNIVERSITÄT
HEIDELBERG**
ZUKUNFT
SEIT 1386



European Research Council
Established by the European Commission



Towards an INN-based solution: Linear Toy Example

- Forward process: given parameters $x_1, x_2 \sim \mathcal{N}(0,1)$, observation y arises according to
$$y = x_1 + x_2 = g(x_1, x_2)$$
- **Inverse** $(x_1, x_2) = g^{-1}(\hat{y})$ for given observation \hat{y} is **undefined**
 - Classical regularization: minimum norm solution $x_1 = x_2 = \frac{\hat{y}}{2}$ (disregards ambiguity!)
- Bayesian solution:
 - Introduce latent variable $z = x_1 - x_2 \quad \Leftrightarrow \quad (y, z) = g_{\text{aug}}(x_1, x_2) = (x_1 + x_2, x_1 - x_2)$ **is invertible!**
 - Reparametrize posterior $p(x_1, x_2 | y)$ as $(x_1, x_2) = g_{\text{aug}}^{-1}(y, z) = \left(\frac{y+z^{(t)}}{2}, \frac{y-z^{(t)}}{2}\right)$ with $z \sim \mathcal{N}(0,2)$
 - Given actual observation \hat{y} , repeat for $t \in 1, \dots, T$:
 - Sample $z^{(t)} \sim \mathcal{N}(0,2)$ and compute $x_1^{(t)} = \frac{\hat{y}+z^{(t)}}{2}$ and $x_2^{(t)} = \frac{\hat{y}-z^{(t)}}{2}$
 - Return $\left\{ \left(x_1^{(t)}, x_2^{(t)} \right) \right\}_{t=1}^T$ as a sample from the Bayesian posterior $p(x_1, x_2 | \hat{y})$

Generalize this to complex settings (non-linear g , noise, high dimensions) by INNs.

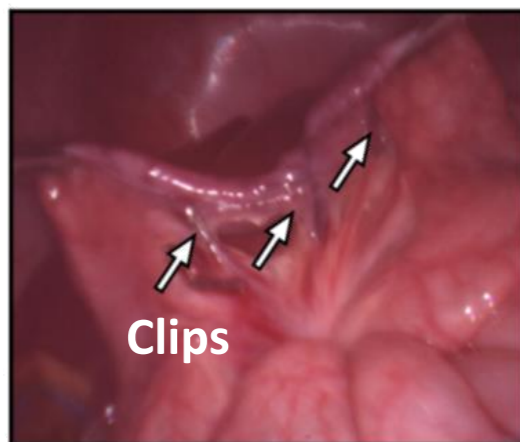




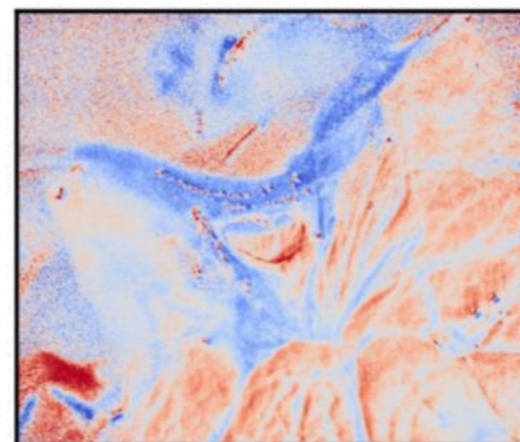
Application: Multispectral Endoscopy

Endoscopes for minimally invasive surgery

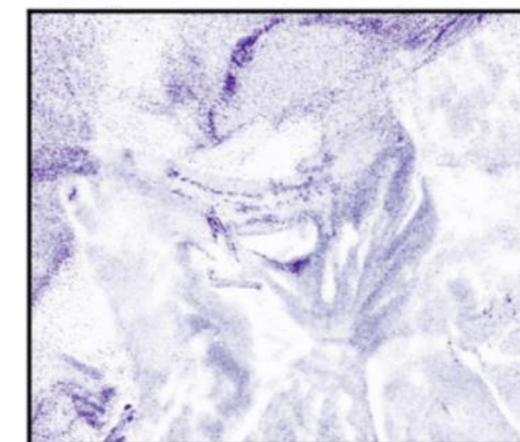
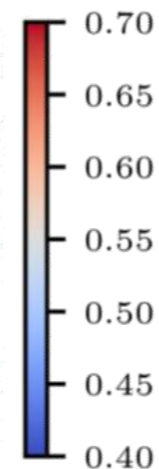
- can be equipped with a multispectral camera
- tissue state x (e.g. blood oxygenation) affects the observed color spectrum y
- **Task:** given spectrum, find posterior distribution of tissue state parameters
- Forward process $s(x)$ is implemented by Monte Carlo simulation



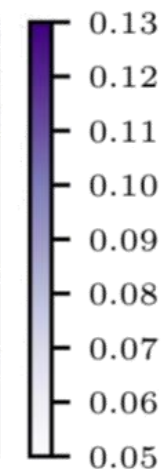
c) RGB image



a) Median sO_2



b) Est. uncertainty

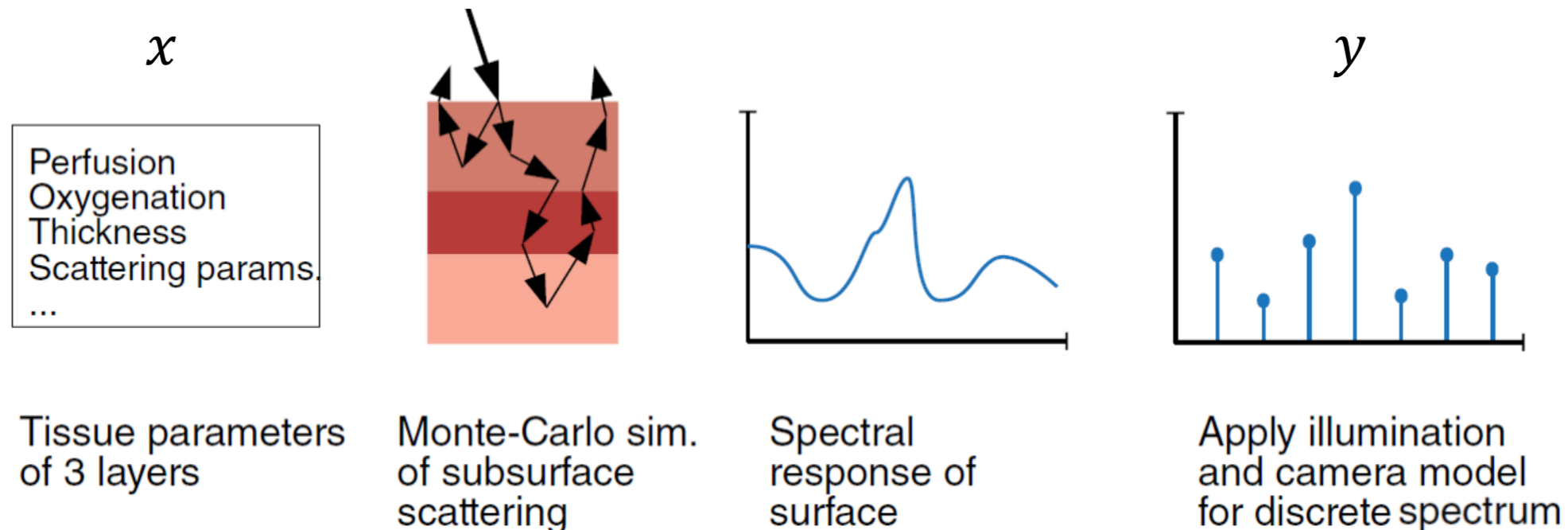




Application: Multispectral Endoscopy

Invert the forward process $s(x)$ implemented by Monte Carlo simulation:

- training: INN learns $[y, z] = f_{\theta}(x) \approx s_{aug}(x)$ with $p(z) \sim \mathcal{N}(0, \mathbb{I})$
- inference: given observed spectrum \hat{y} , sample $\{z_i \sim p(z)\}_{i=1}^M$ and compute posterior sample $\{x_i = f_{\theta}^{-1}(\hat{y}, z_i)\}_{i=1}^M$ (independently for every pixel)
- determine mean and variance from $\{x_i\}$ – works especially well for blood oxygenation

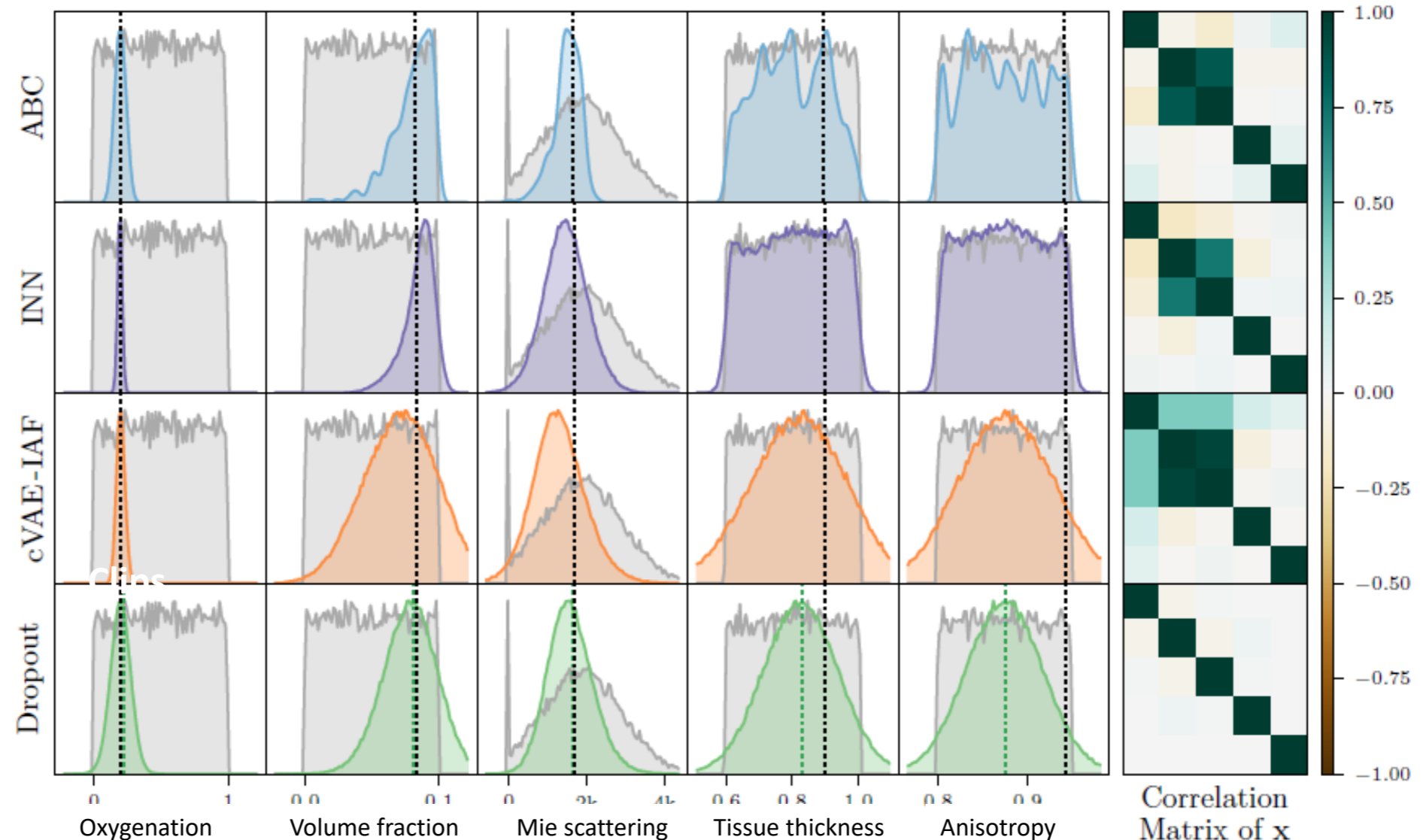




Application: Multispectral Endoscopy

Results

- INN performs well
- not all parameters are identifiable

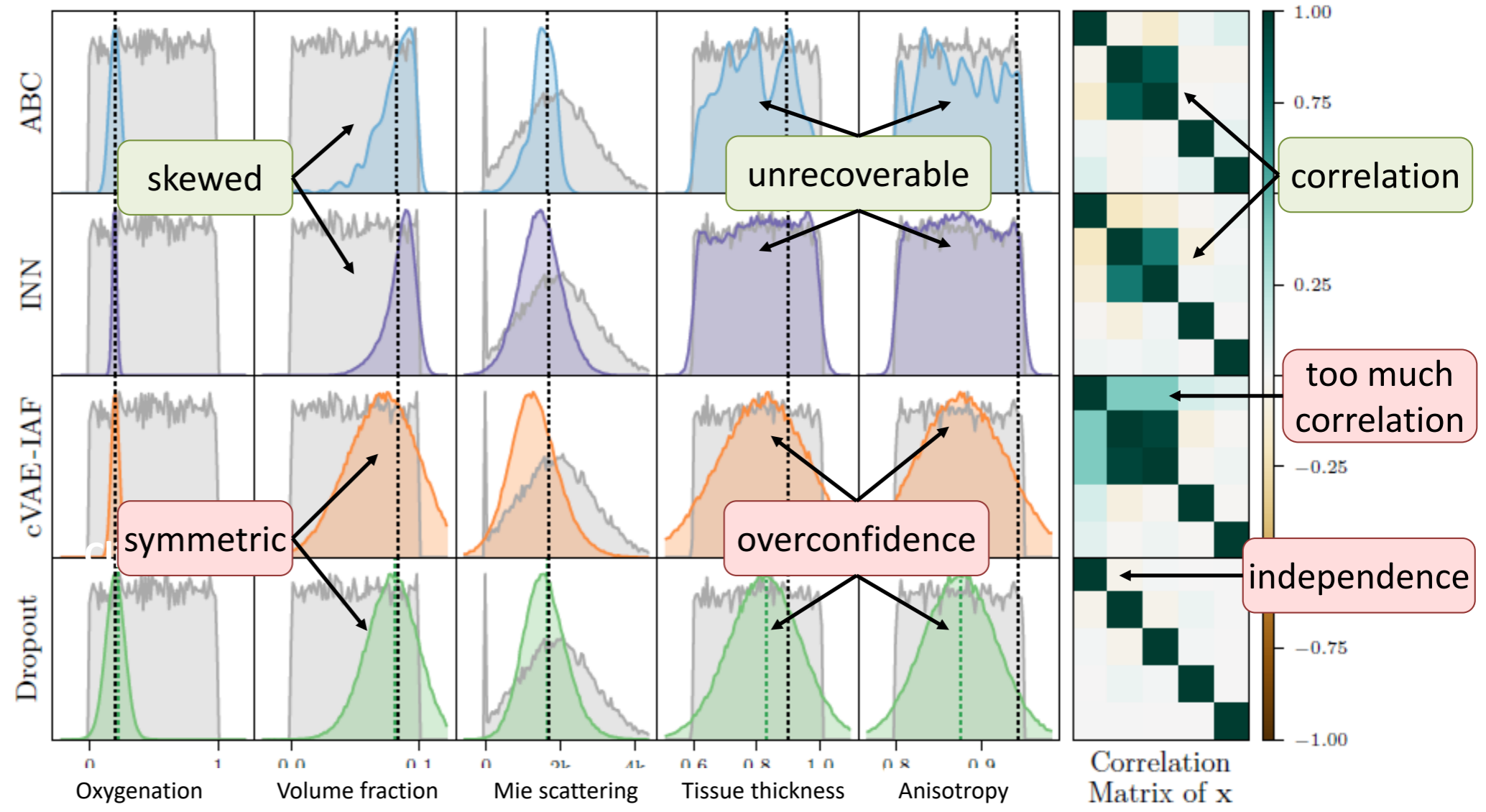




Application: Multispectral Endoscopy

Results

- INN performs well
- not all parameters are identifiable
- incorrect results for other methods
 - skewed distribut. appear symmetric
 - non-identifiable parameters have spurious mode
 - correlation is too weak or too strong



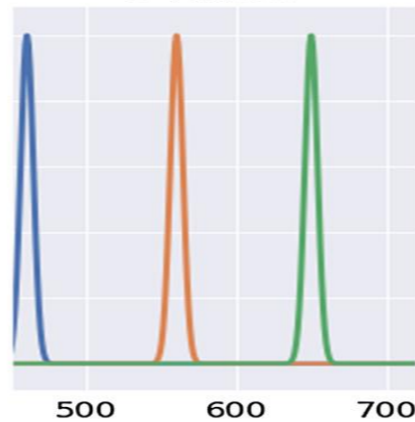


Experimental Design for Multispectral Endoscopy

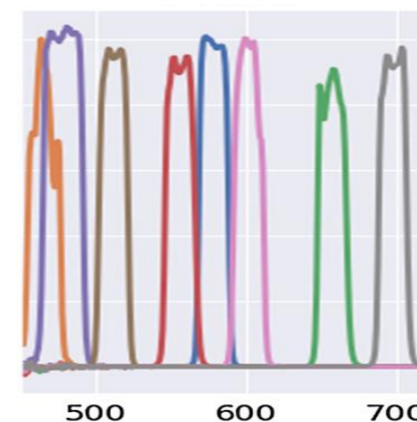
Analysis of posteriors: Which camera should be used?

- 3 to 27 spectral channels
- Which gives reliable results at best price and usability?

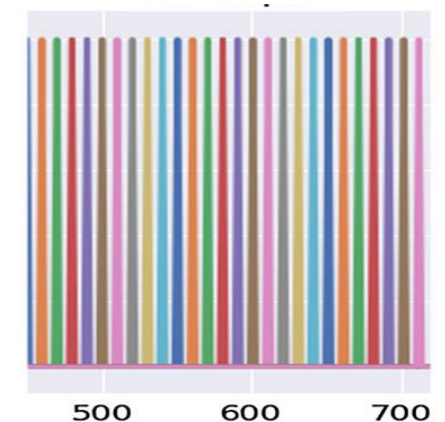
3 spectral channels



8 spectral channels

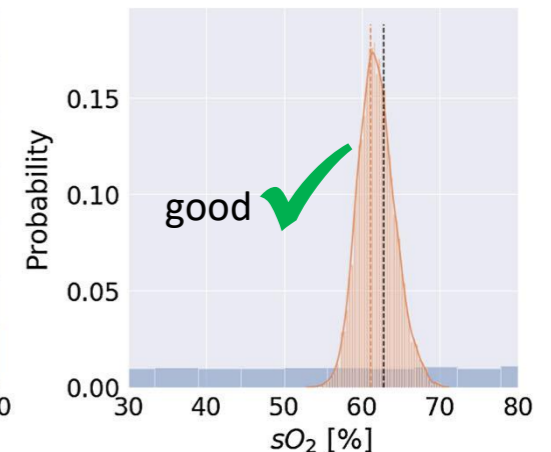
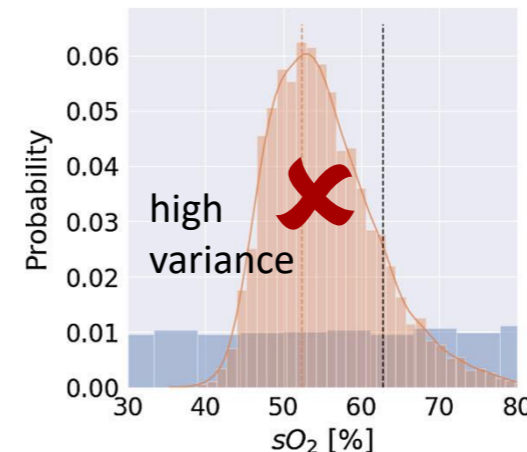
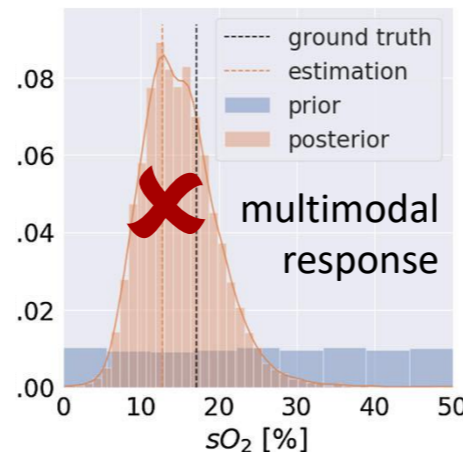


27 spectral channels



- posterior oxygen level histograms:

⇒ camera with 8 channels offers best trade-off between price and accuracy





INN Architecture for Endoscopy Application

- Forward process: given tissue parameters x , spectrum y arises from MC simulation g

$$y = g(x)$$

- Bayesian solution:

- Introduce **latent variables** z collecting the information about x that got lost in $y = g(x)$

$$y, z = g_{\text{aug}}(x)$$

- Train INN for $g_{\text{aug}}(x)$ with $p(z) = \mathcal{N}(0, \mathbb{I})$ and $y \perp z$, using synthetic training data from the simulation

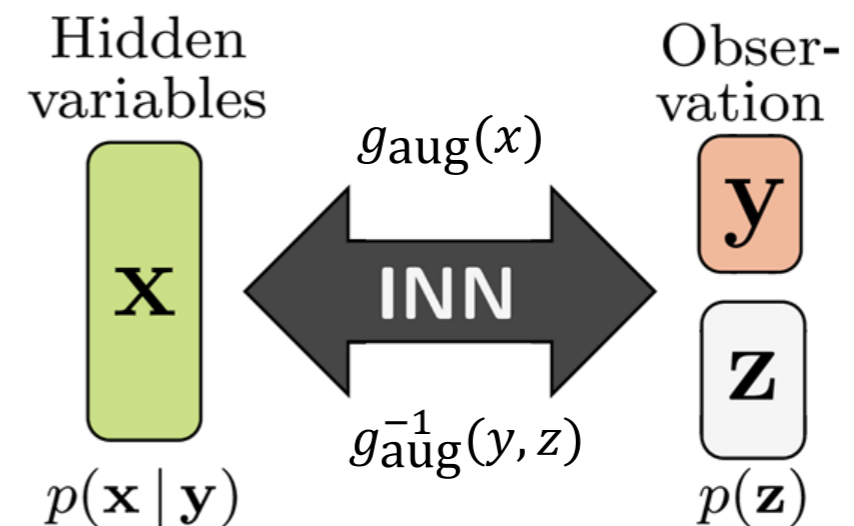
- Inference for real observation y_{obs} :

- For $t \in 1, \dots, T$:

- Sample $z^{(t)} \sim \mathcal{N}(0, \mathbb{I})$

- compute $x^{(t)} = g_{\text{aug}}^{-1}(y_{\text{obs}}, z^{(t)})$

- Return $\{(x^{(t)})\}_{t=1}^T$ as a sample from Bayesian posterior $p(x | y_{\text{obs}})$



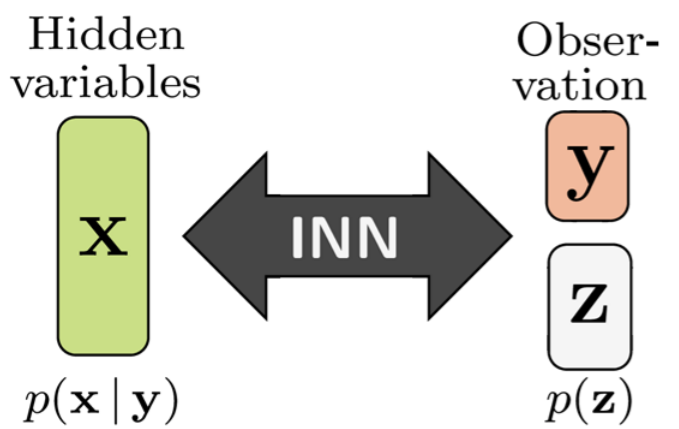


INN Architectures for Conditional Inference

Split latent space

training: $(y, z) = f_{\theta}(x)$
s.t. $p(z) = \mathcal{N}(0, \mathbb{I})$

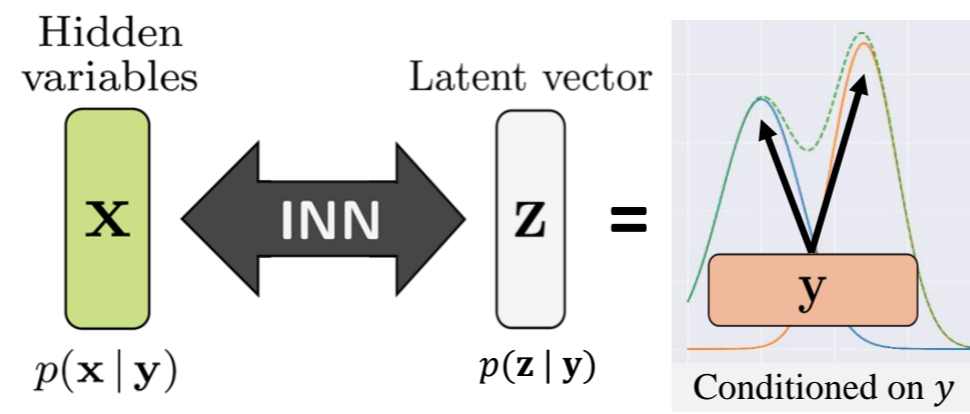
inference: sample $z \sim \mathcal{N}(0, \mathbb{I})$
compute $x = f_{\theta}^{-1}(\hat{y}, z)$
 $\Rightarrow x \sim p(x | \hat{y})$



Latent mixture INN

training: $z = f_{\theta}(x)$
s.t. $p(z) = \text{GMM}(z; y) = \sum_y \mathcal{N}(\mu_y, \Sigma_y)$

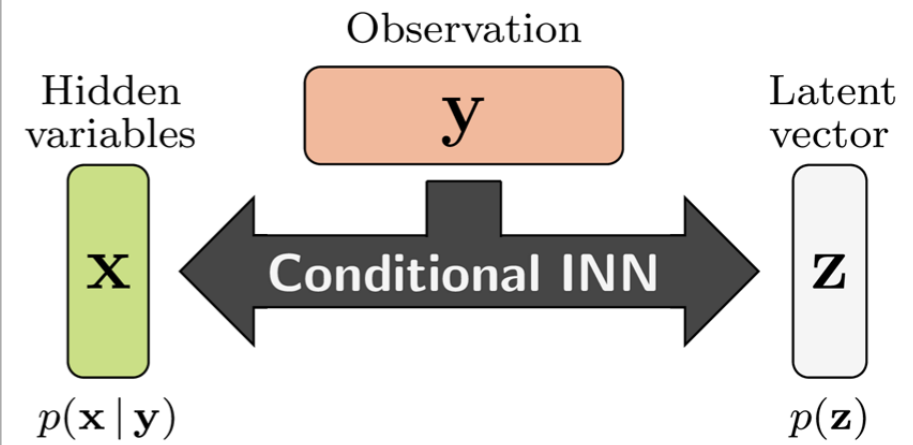
inference: sample $z \sim \mathcal{N}(\mu_{\hat{y}}, \Sigma_{\hat{y}})$
compute $x = f_{\theta}^{-1}(z)$
 $\Rightarrow x \sim p(x | \hat{y})$



Conditional INN

training: $z = f_{\theta}(x; y)$
s.t. $p(z) = \mathcal{N}(0, \mathbb{I})$

inference: sample $z \sim \mathcal{N}(0, \mathbb{I})$
compute $x = f_{\theta}^{-1}(z; \hat{y})$
 $\Rightarrow x \sim p(x | \hat{y})$



historically first

classification, disentanglement

inverse inference





BayesFlow:

Model-Based Inverse Inference with cINNs

Model-based inverse inference:

- system with intrinsic parameters x (hidden) and observations y (measurable)
- good scientific **understanding of the forward process: How does y arise from given x ?**
(e.g. differential equations, simulations)
- solve the **inverse problem: Which hidden parameters x explain some actual observations \hat{y} ?**
- usually no analytic solution
 - ambiguous outcomes due to information loss from x to $y \Rightarrow$ must estimate posterior $p(x | \hat{y})$
 - simplest approach: manually adjust x until outcomes match $\hat{y} \Rightarrow$ neglects uncertainty
 - traditional Bayesian inference: sampling methods (MCMC, HMC, ...) \Rightarrow very expensive
- standard ML methods are often not applicable
 - lack of training data with known ground truth x^*
 - only point estimates, no posteriors (i.e. no diverse outputs)
- cINN can elegantly solve the Bayesian inverse problem





BayesFlow:

Model-Based Inverse Inference with cINNs

cINNs make clever use of the known forward model to solve the inverse problem

- run cINN in **forward** mode for model-based **training**
 - use known forward model to create synthetic training data \Rightarrow cINN becomes a fast surrogate
 - train with diverse forward scenarios and noise \Rightarrow cINN learns the ambiguity and uncertainty
 - run cINN in **backward** mode for inverse **inference**
 - use actual observations \hat{y} as condition
 - sample many latents $\{z_k \sim p(z)\}_{k=1}^N$
 - run cINN backwards $\{x_k = g(z_k; \hat{y})\}_{k=1}^N$
- the $\{x_k\}$ are a sample of the Bayesian posterior $p(x | \hat{y})$
- “Train forward, get the inverse for free.”

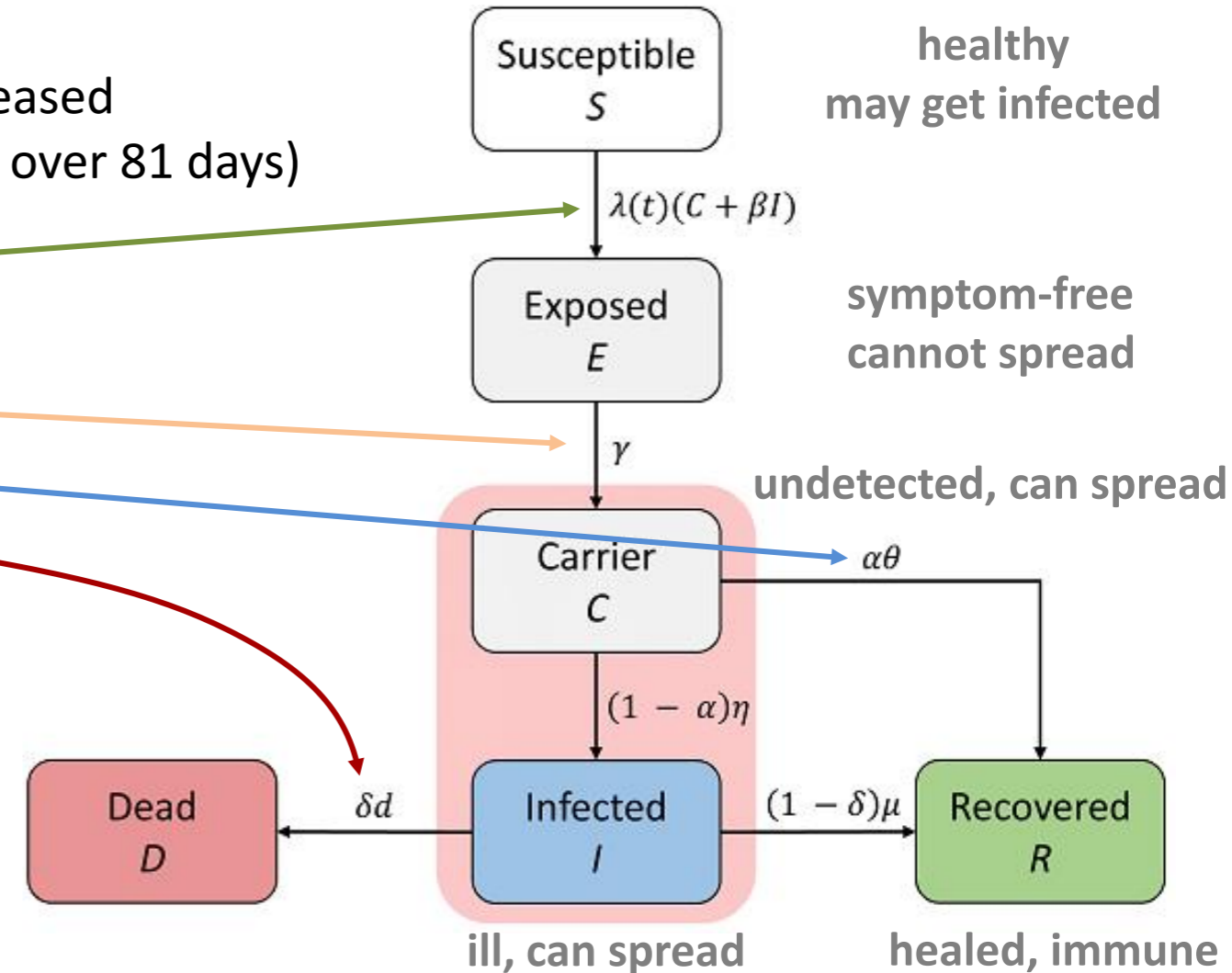




BayesFlow for Covid-19 Epidemiology

Epidemiology as a difficult inverse problem:

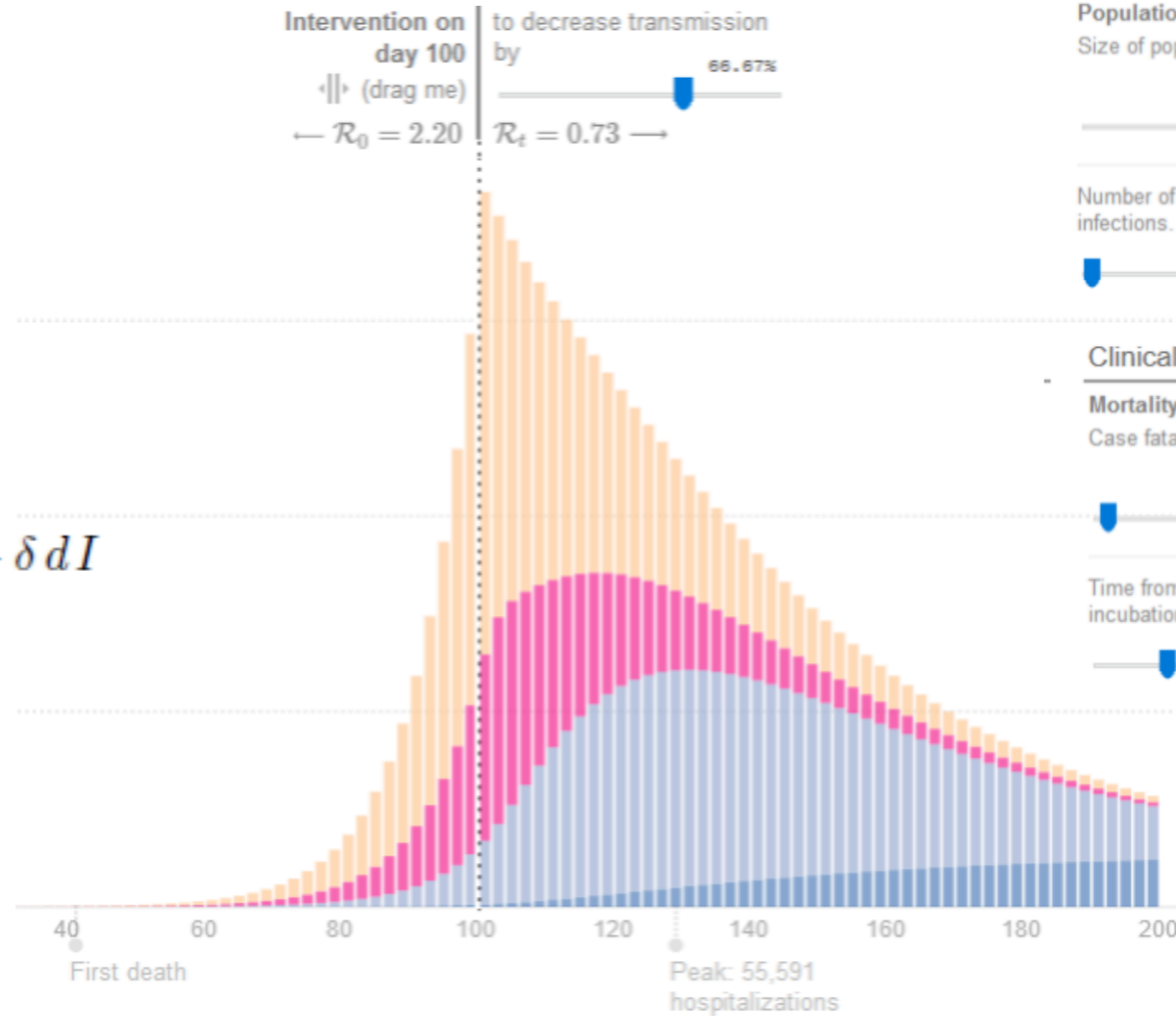
- observations: time series of infected, recovered, deceased (as of June 2020: 243 measurements = 3 observables over 81 days)
- 34 hidden parameters:
 - infection rate $\lambda(t)$
 - latent period $1/\gamma$
 - undetected fraction α
 - case fatality rate δ
 - ...
- prior knowledge:
 - SIR-type compartmental model (ODE system similar to Lotka-Volterra)
 - dates of government interventions
 - sources of reporting errors
 - ...





Forward Model: Epidemic Calculator

$$\begin{aligned} \frac{dS}{dt} &= -\lambda(t) \left(\frac{C + \beta I}{N} \right) S \\ \frac{dE}{dt} &= \lambda(t) \left(\frac{C + \beta I}{N} \right) S - \gamma E \\ \frac{dC}{dt} &= \gamma E - (1 - \alpha)\eta C - \alpha\theta C \\ \frac{dI}{dt} &= (1 - \alpha)\eta C - (1 - \delta)\mu I - \delta dI \\ \frac{dR}{dt} &= \alpha\theta C + (1 - \delta)\mu I \\ \frac{dD}{dt} &= \delta dI \end{aligned}$$



Transmission Dynamics

Population Inputs	Basic Reproduction Number \mathcal{R}_0	Transmission Times
Size of population. 7,000,000	Measure of contagiousness: the number of secondary infections each infected individual produces. 2.2	Length of incubation period, T_{inc} . 5.44 days
Number of initial infections. 1		Duration patient is infectious, T_{inf} . 2.9 Days

Clinical Dynamics

Mortality Statistics	Recovery Times	Care statistics
Case fatality rate. 2.00 %	Length of hospital stay 28.6 Days	Hospitalization rate. 20.00 %
Time from end of incubation to death. 32 Days	Recovery time for mild cases 11.1 Days	Time to hospitalization. 5 Days



BayesFlow for Epidemiology: The Networks

- Inference problem: observation sequence (IRD) \Rightarrow parameter posteriors
 - Solve with BayesFlow network: cINN with statistical preprocessing networks for y
 - Training: end-to-end optimization of maximum likelihood loss with 70000 simulations

Convolutional:

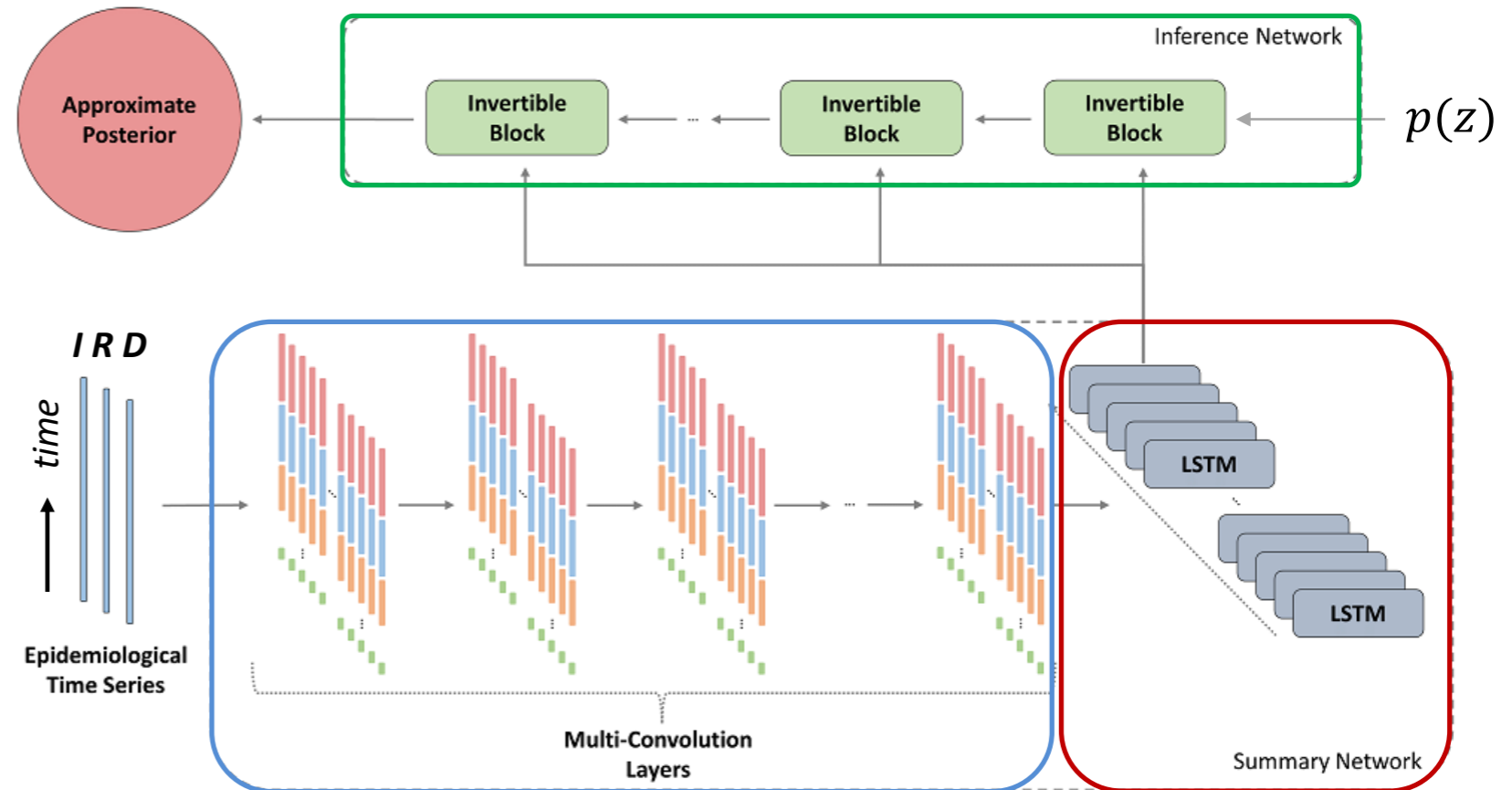
- noise reduction
- feature detection

Recurrent (LSTM):

- variable-length sequence to fixed size summary

Invertible (cINN):

- posterior inference

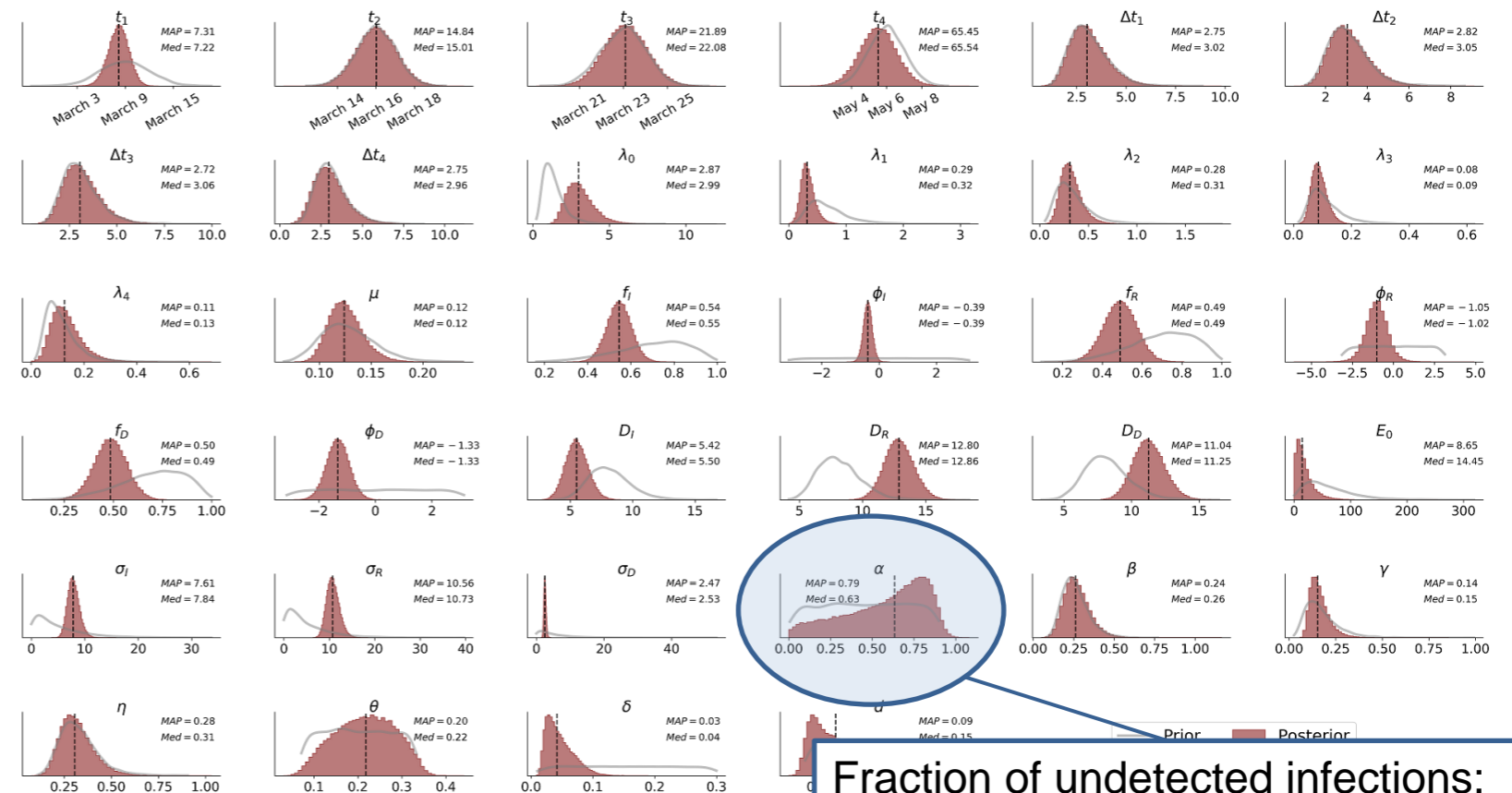




BayesFlow for Epidemiology: Covid-19 Marginal Posteriors

Results: marginal posteriors for first wave in Germany (March – June 2020, 81 time steps)

- High fraction of undetected infections:
63% (median), 79% (mode)
- Serial interval: 9-10 days
- High likelihood to transmit disease *before* diagnosis
- time to recovery:
4.6 days (undetected infections)
11.3 days (diagnosed cases)
(3.2 + 8.1 days before/after diagnosis)
- often non-Gaussian behavior



Fraction of undetected infections:
uniform prior \Rightarrow peaked posterior

Correspond well to clinical findings



BayesFlow for Epidemiology: Strengths

Well-calibrated **uncertainty quantification**

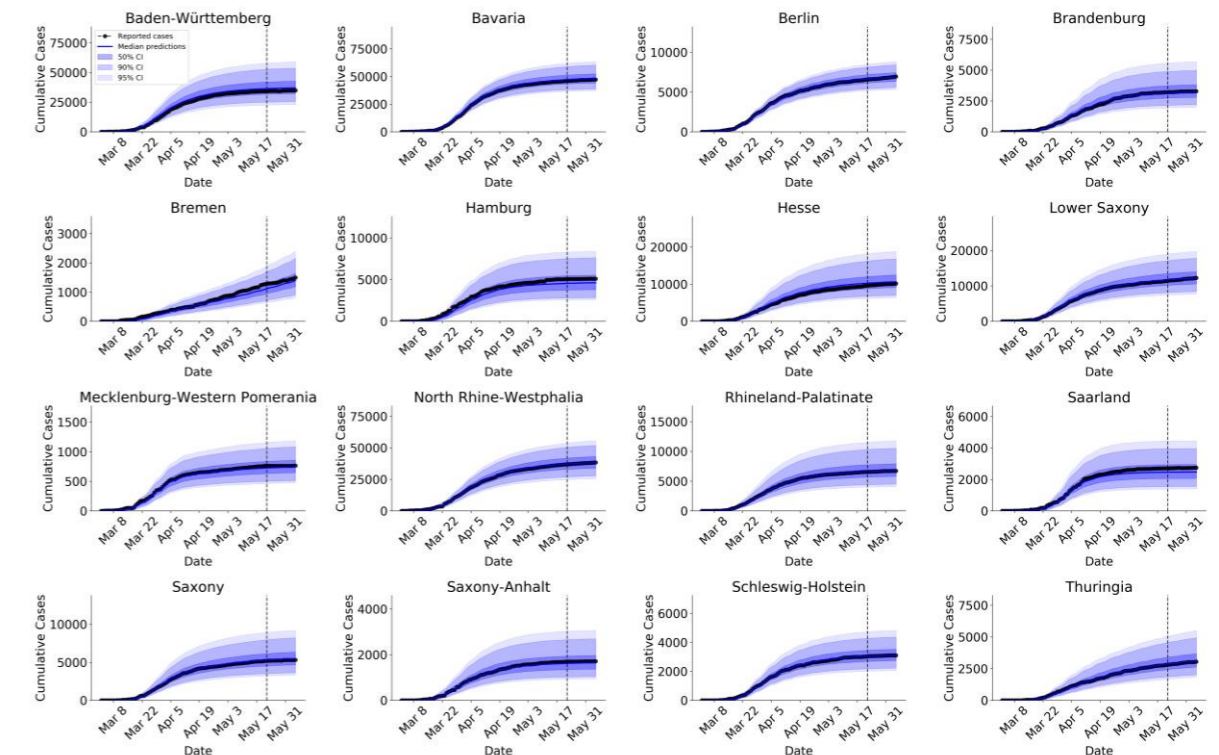
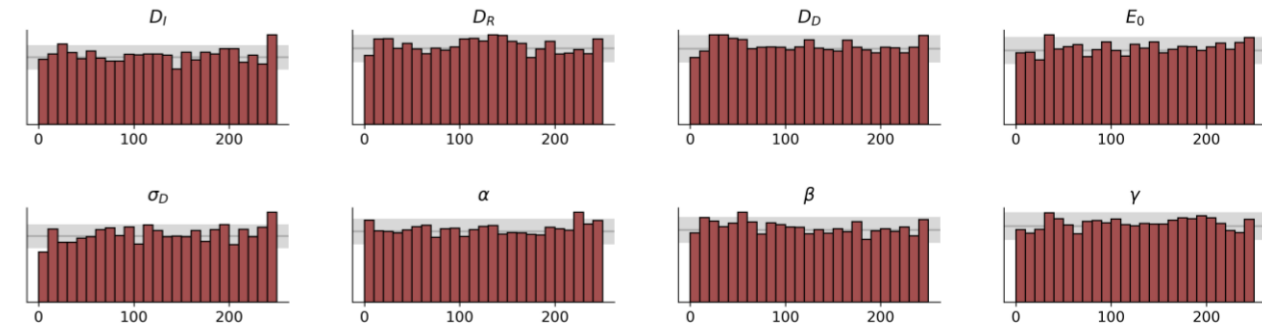
- $q\%$ confidence intervals are hit $\approx q\%$ of the time
- much better than classical estimators (e.g. least squares fitting, manual parameter tuning, ...)

Efficient backward operation \Rightarrow **fast inference**

- train once, predict often
- in contrast, MCMC runs from scratch for each \hat{y}
- Bayesflow upfront training effort $\approx 10 - 100x$ of single MCMC inference

\Rightarrow **training effort amortizes quickly**

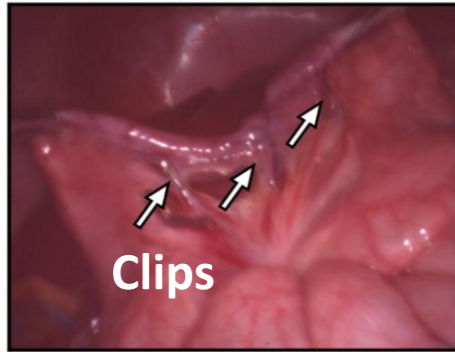
analysis of German states with identical network



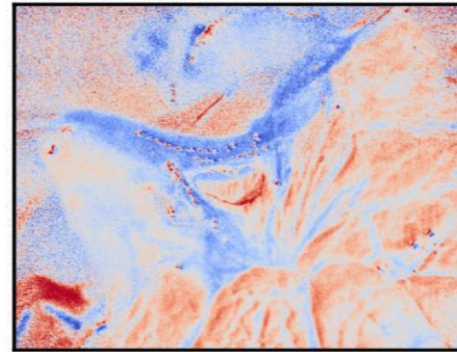


Very diverse inverse problems were solved with INNs/BayesFlow

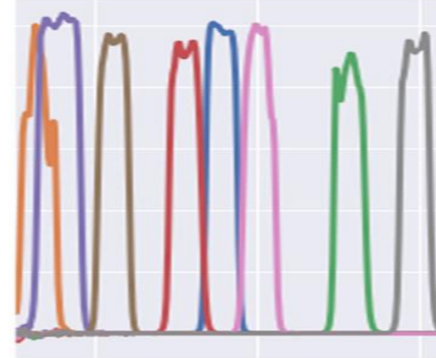
- Surgery:



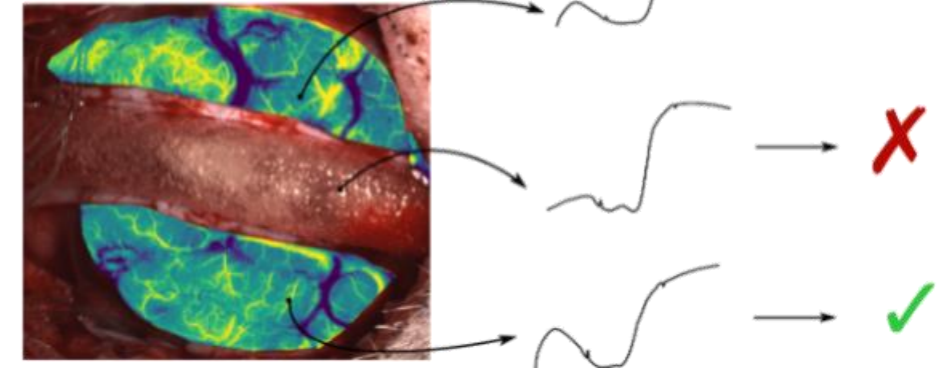
blood oxygenation



experimental design



outlier detection

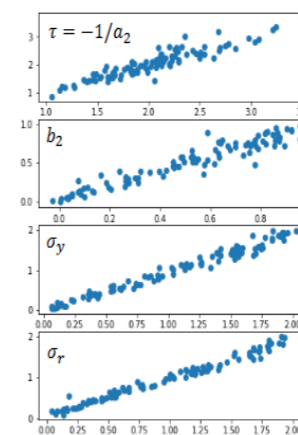
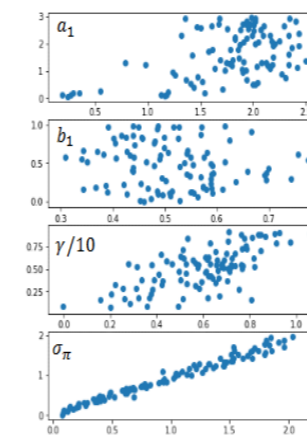


- Photo-acoustic imaging
- Particle physics
- Astrophysics
- Environmental physics
- Cognitive Science
- Inverse kinematics of robots
- Mechanical engineering

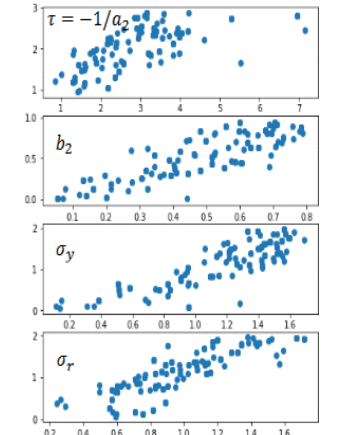
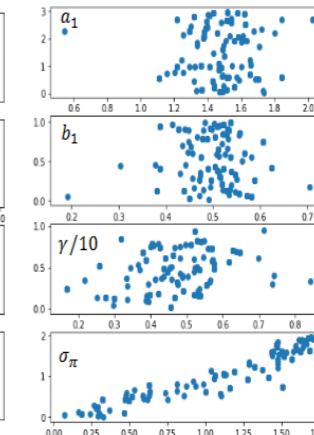
- Finance:

BayesFlow
beats MCMC

BayesFlow



KDE-MCMC



Ardizzone et al. "Analyzing inverse problems with invertible neural networks", ICLR 2019.
Adler et al. "Out of distribution detection for intra-operative functional imaging", UNSURE 2019.
Shiono "Estimation of agent-based models using BayesFlow", SSRN 2020.

Guaranteed disentanglement with Nonlinear ICA and Incompressible Flows

Ullrich Köthe

Visual Learning Lab, Heidelberg University
joint work with **Carsten Rother, Peter Sorrenson**

Tutorial „Normalizing Flows“ at CVPR 2021

June 2021



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386



European Research Council
Established by the European Commission

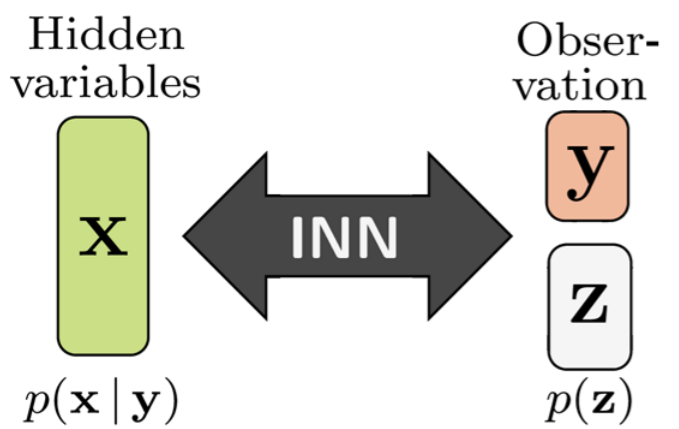


INN Architectures for Conditional Inference

Split latent space

training: $(y, z) = f_{\theta}(x)$
s.t. $p(z) = \mathcal{N}(0, \mathbb{I})$

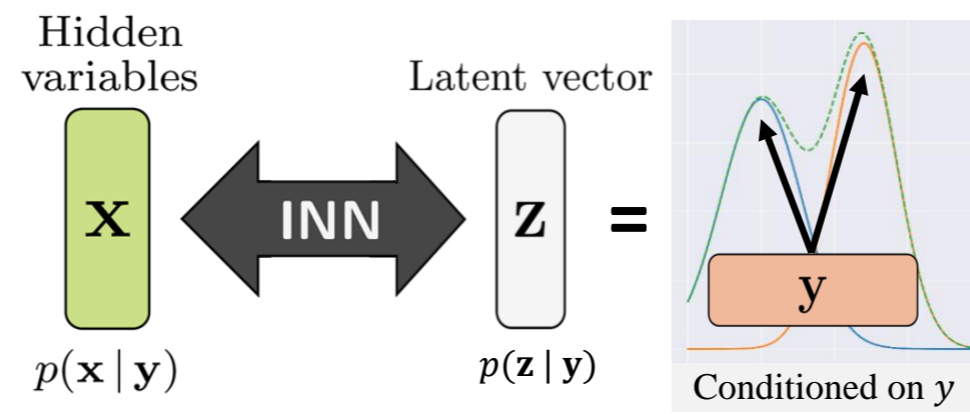
inference: sample $z \sim \mathcal{N}(0, \mathbb{I})$
compute $x = f_{\theta}^{-1}(\hat{y}, z)$
 $\Rightarrow x \sim p(x | \hat{y})$



Latent mixture INN

training: $z = f_{\theta}(x)$
s.t. $p(z) = \text{GMM}(z; y) = \sum_y \mathcal{N}(\mu_y, \Sigma_y)$

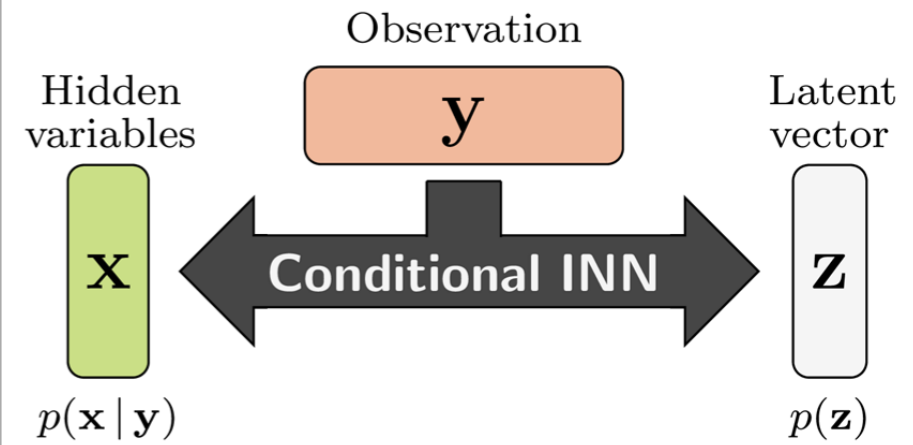
inference: sample $z \sim \mathcal{N}(\mu_{\hat{y}}, \Sigma_{\hat{y}})$
compute $x = f_{\theta}^{-1}(z)$
 $\Rightarrow x \sim p(x | \hat{y})$



Conditional INN

training: $z = f_{\theta}(x; y)$
s.t. $p(z) = \mathcal{N}(0, \mathbb{I})$

inference: sample $z \sim \mathcal{N}(0, \mathbb{I})$
compute $x = f_{\theta}^{-1}(z; \hat{y})$
 $\Rightarrow x \sim p(x | \hat{y})$



historically first

classification, disentanglement

inverse inference



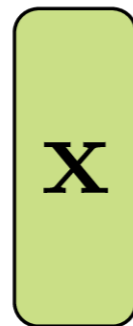


Interpretable Latent Spaces with Latent Mixture INNs (LM-INNs)

Interpretable latent spaces are a key to explainable machine learning

- Latent Mixture INNs are especially suitable for this task
 - Variation of cINNs: condition y acts on the latent space, not on the function g
 - cINN: $x \sim p(x | y) \Leftrightarrow z \sim p(z), \quad x = g(z; y)$
 - LM-INN: $x \sim p(x | y) \Leftrightarrow z \sim p(z | y), \quad x = g(z)$
 - Define $p(z | y)$ as a **mixture of Gaussians** instead of a single Gaussian
 - Especially simple when y is a class label: learn one mixture component $p(z | y = k)$ per label k

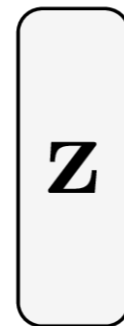
Hidden variables



$p(\mathbf{x} | \mathbf{y})$

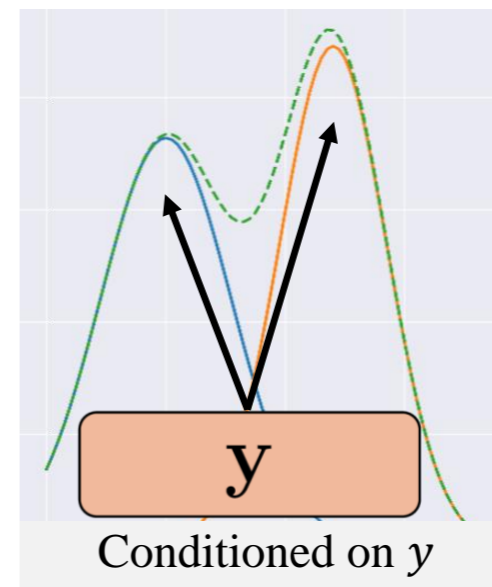


Latent vector



$p(\mathbf{z} | \mathbf{y})$

=





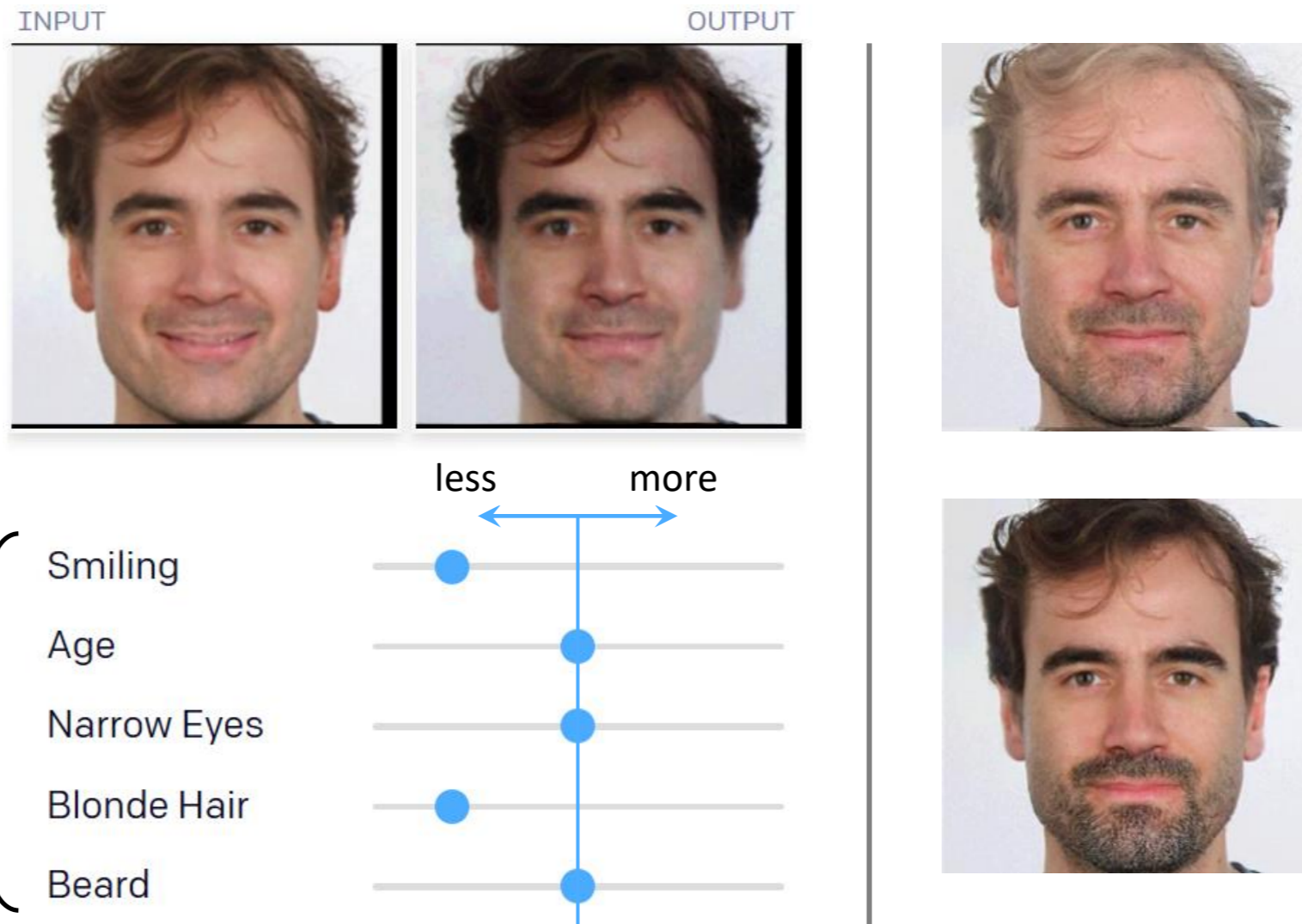
What is Disentanglement?

Train network so that each latent feature has a single interpretable effect

Example: GLOW INN

- Try your own face:
openai.com/blog/glow/

changing the value of a single latent feature has a coordinated and intuitive effect on many pixels simultaneously





Disentanglement: Definition

- Definition by Bengio et al.:
 - A disentangled representation has recovered the “*informative factors of variation*” in a dataset
 - Disentangled latent features separate different categories of information (e.g. identity, pose and background) into independent degrees of freedom
- Disentangled representations are **interpretable** by humans and **generalize well** for downstream tasks and transfer learning
- Methods so far empirically work well, but have no theoretical guarantees
- We apply the theory of **nonlinear ICA** to **INNs** to derive such guarantees





What is Disentanglement?

- Latent dimensions should have one and only one isolated effect on the data



- β -VAE disentangles azimuth whereas VAE entangles it with other variables

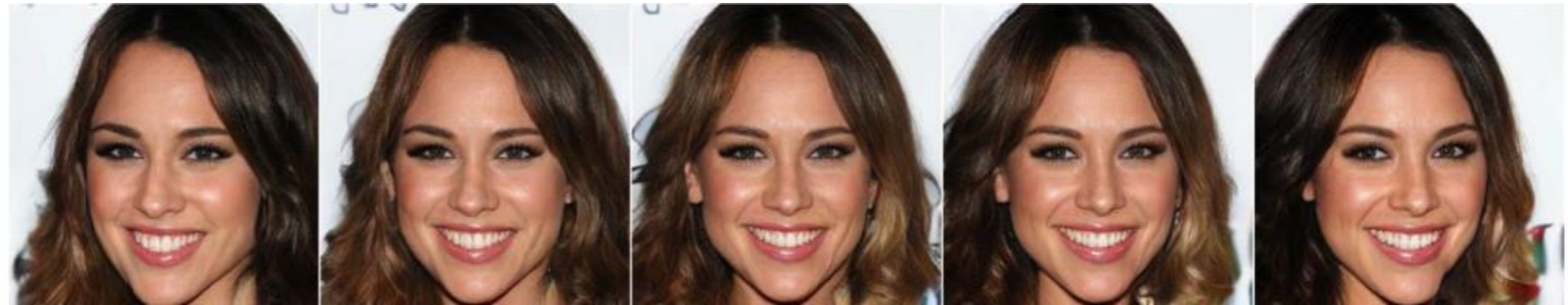




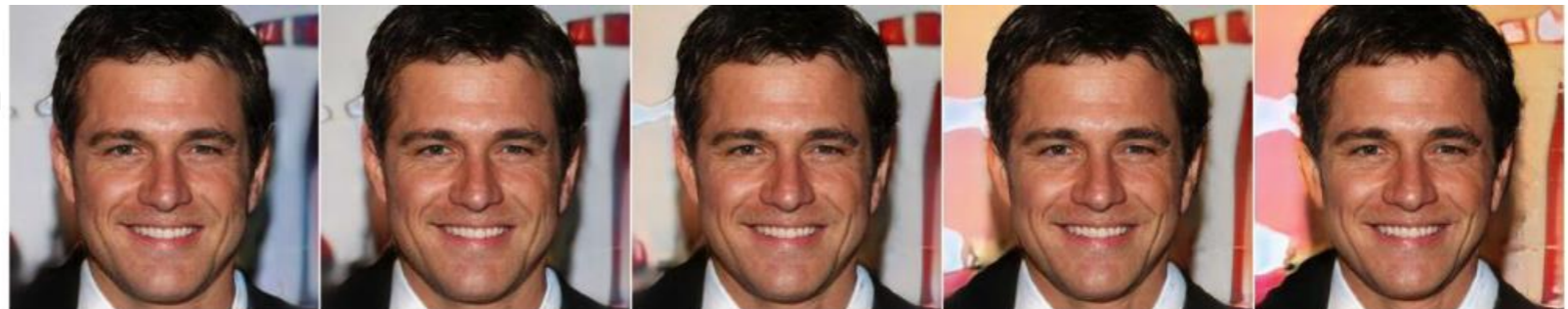
ID-GAN

- Information-Distillation Generative Adversarial Network is probably state-of-the-art
- Combine VAE encoder with conditional GAN generator

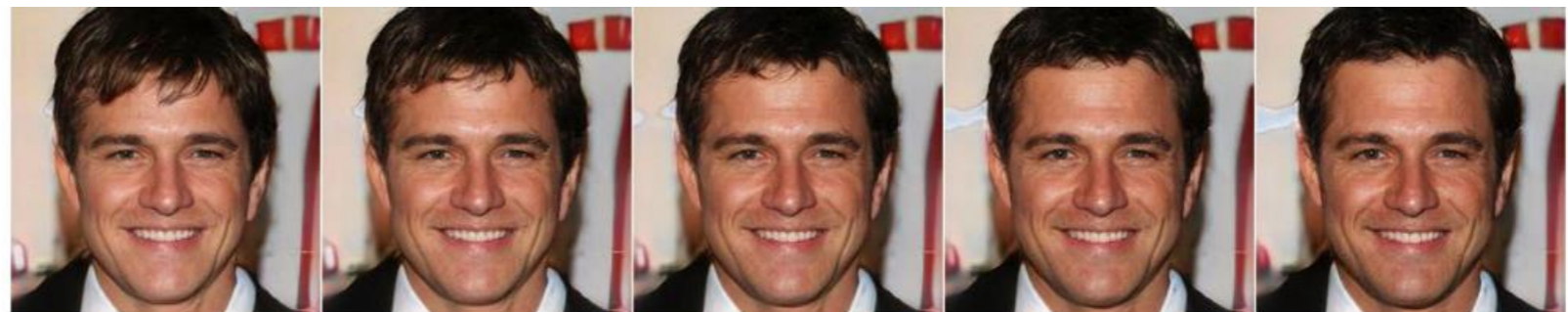
- Works well on large images (CelebA-HQ: 1024x1024)
- GAN conditioned on β -VAE latent code
- Additional cycle constraint: maximize mutual information between latent codes of real and fake images



azimuth



BG color



bang





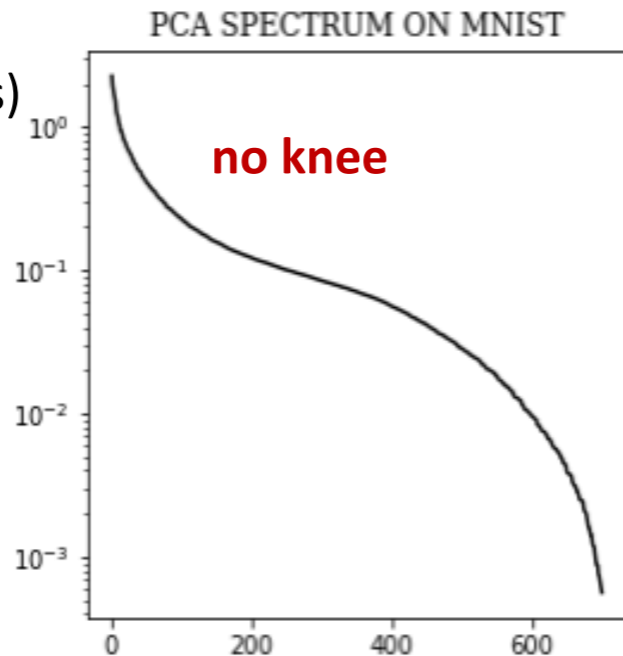
LM-INNs for Disentanglement

- **Fundamental disentanglement:** separate content from noise
 - number of content dimensions = intrinsic dimension of the dataset
 - similar to autoencoder bottleneck, but intrinsic dimension is learned (not chosen as a hyperparameter)
- **Content disentanglement:** disentangle content subspace into meaningful features

Classical: linear disentanglement by PCA

- do eigen decomposition
- sort features (eigenvectors) by energy (eigenvalues)

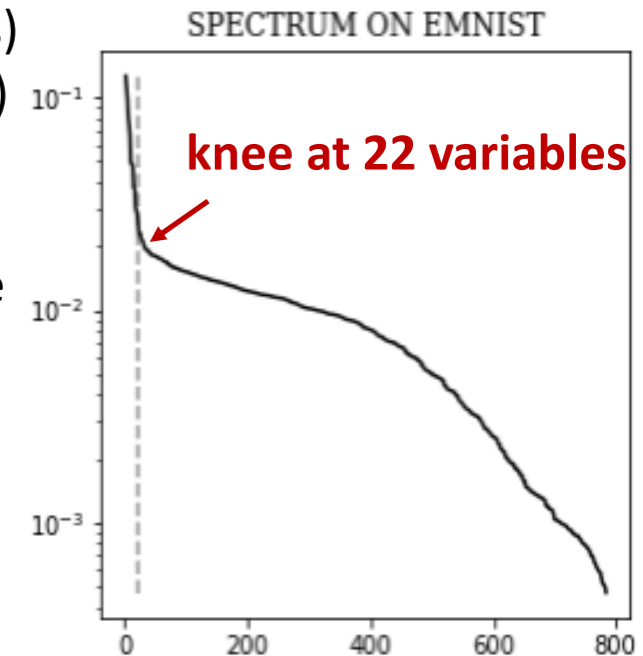
Spectrum is usually smooth
 ⇒ no clear choice for intrinsic dimension



New: non-linear ICA by LM-INN

- train LM-INN ($y = \text{class labels}$)
- sort features (latent variables) by energy (latent variance)

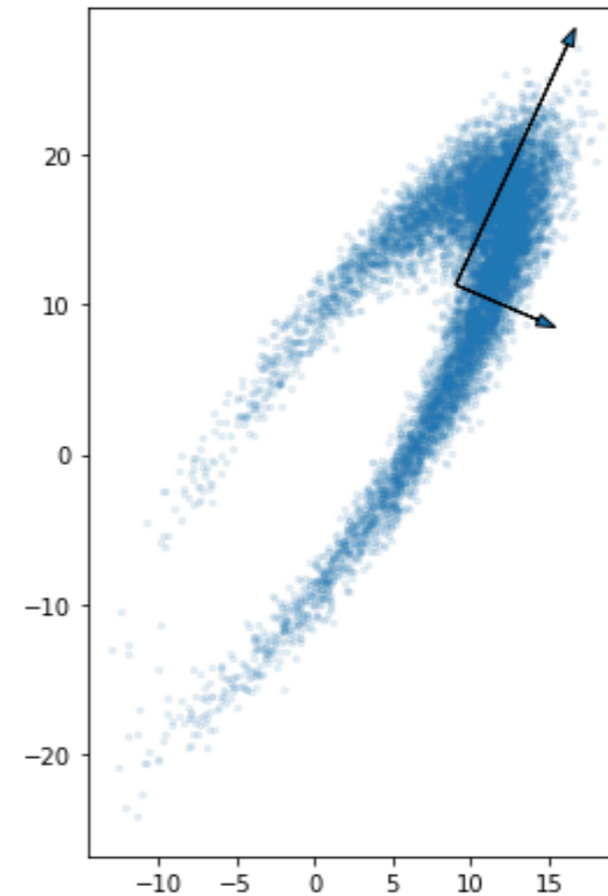
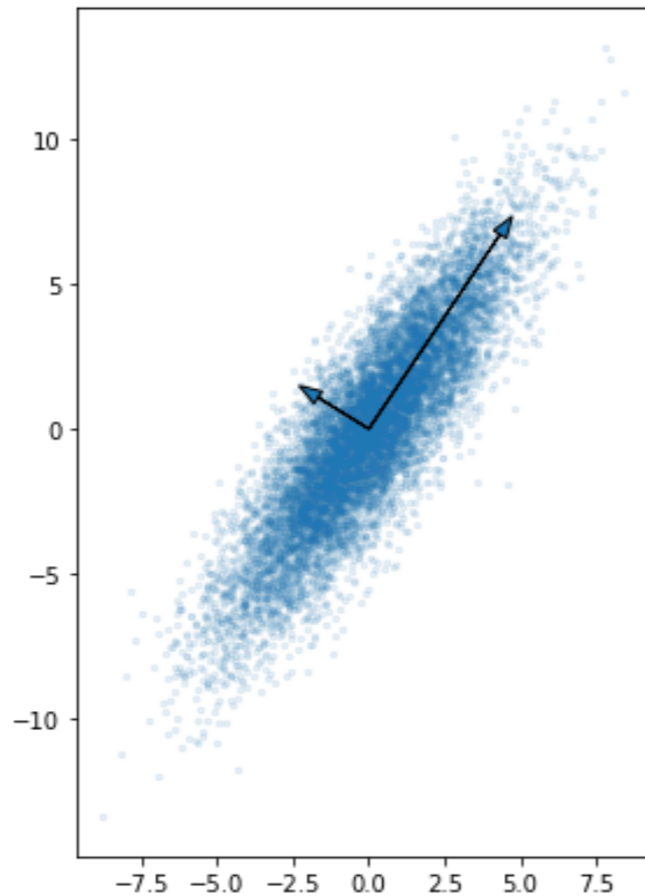
Spectrum may have marked knee
 ⇒ clear identification of the intrinsic dimension





Recap: PCA (Principal Component Analysis)

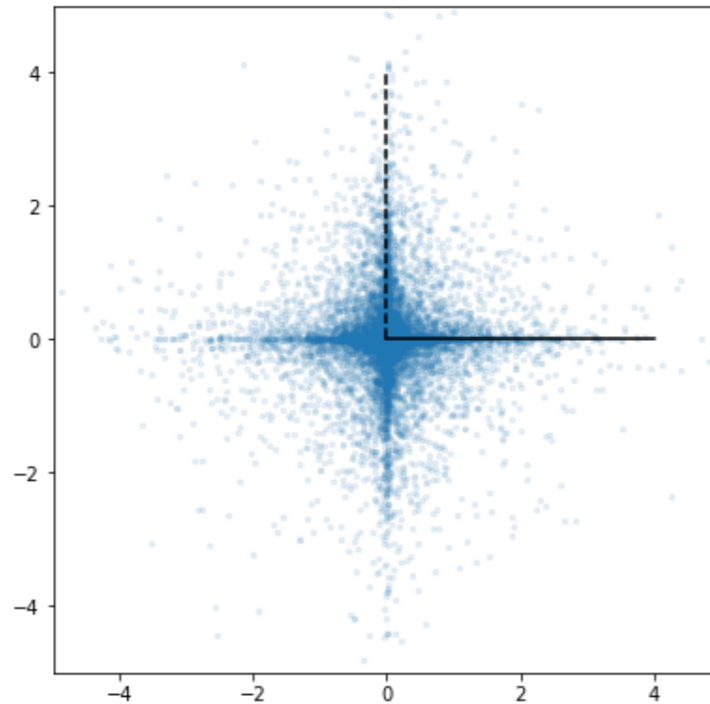
- Classical method for unsupervised disentanglement with a linear transformation
 - Finds **uncorrelated** basis vectors for multivariate Gaussian distributions
 - Can be applied to non-Gaussian data, but cannot fully disentangle them



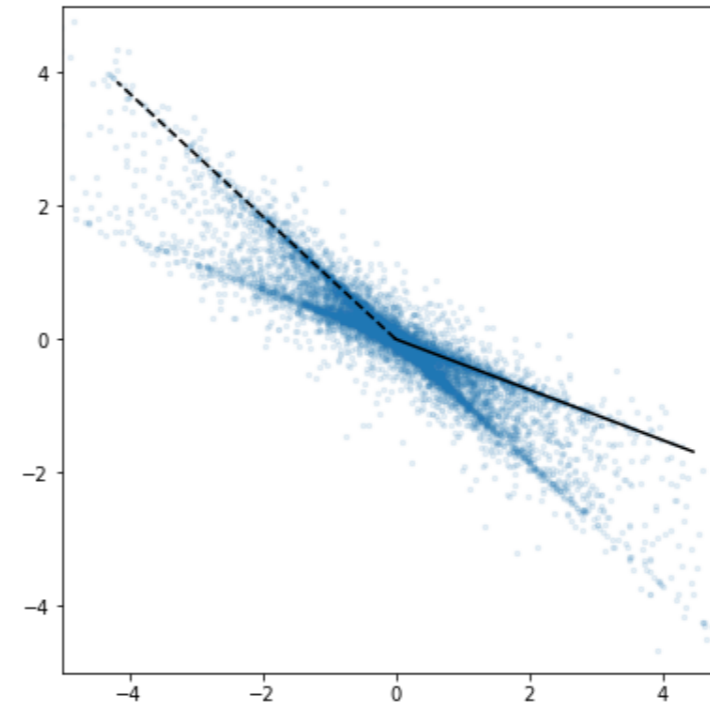


Recap: ICA (Independent Component Analysis)

- Roughly: Independent Component Analysis generalizes PCA to non-gaussian case
 - Apply arbitrary invertible linear transformation to factorial **non-Gaussian latent distribution**



Latent space
(non-gaussian,
independent dimensions)



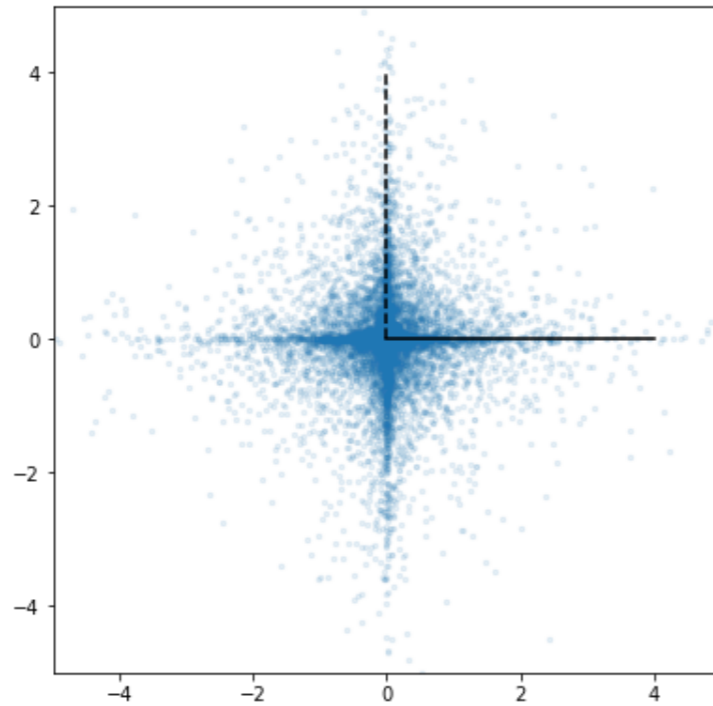
Data space

Invertible linear
transformation
→

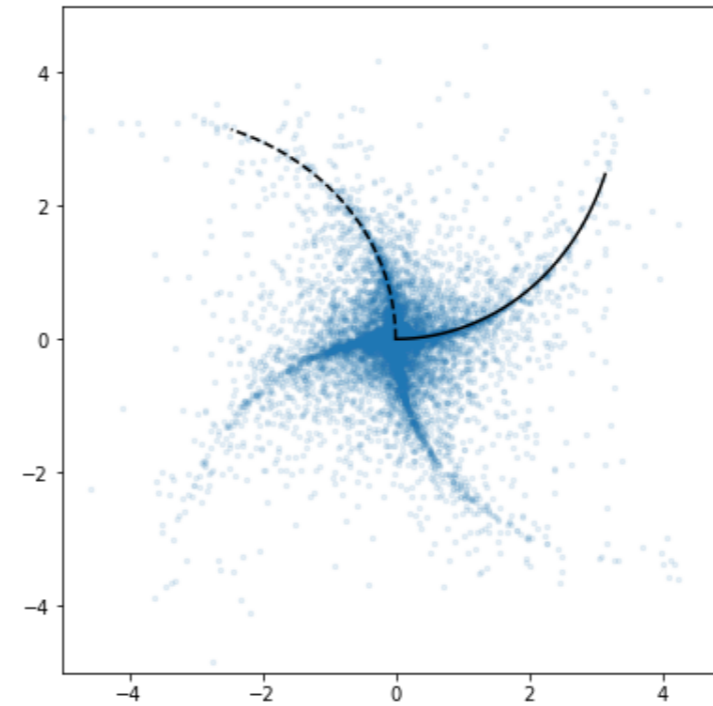


Nonlinear ICA

- Replace the linear transformation with an **invertible non-linear** transformation



Latent space



Invertible **non-linear**
transformation



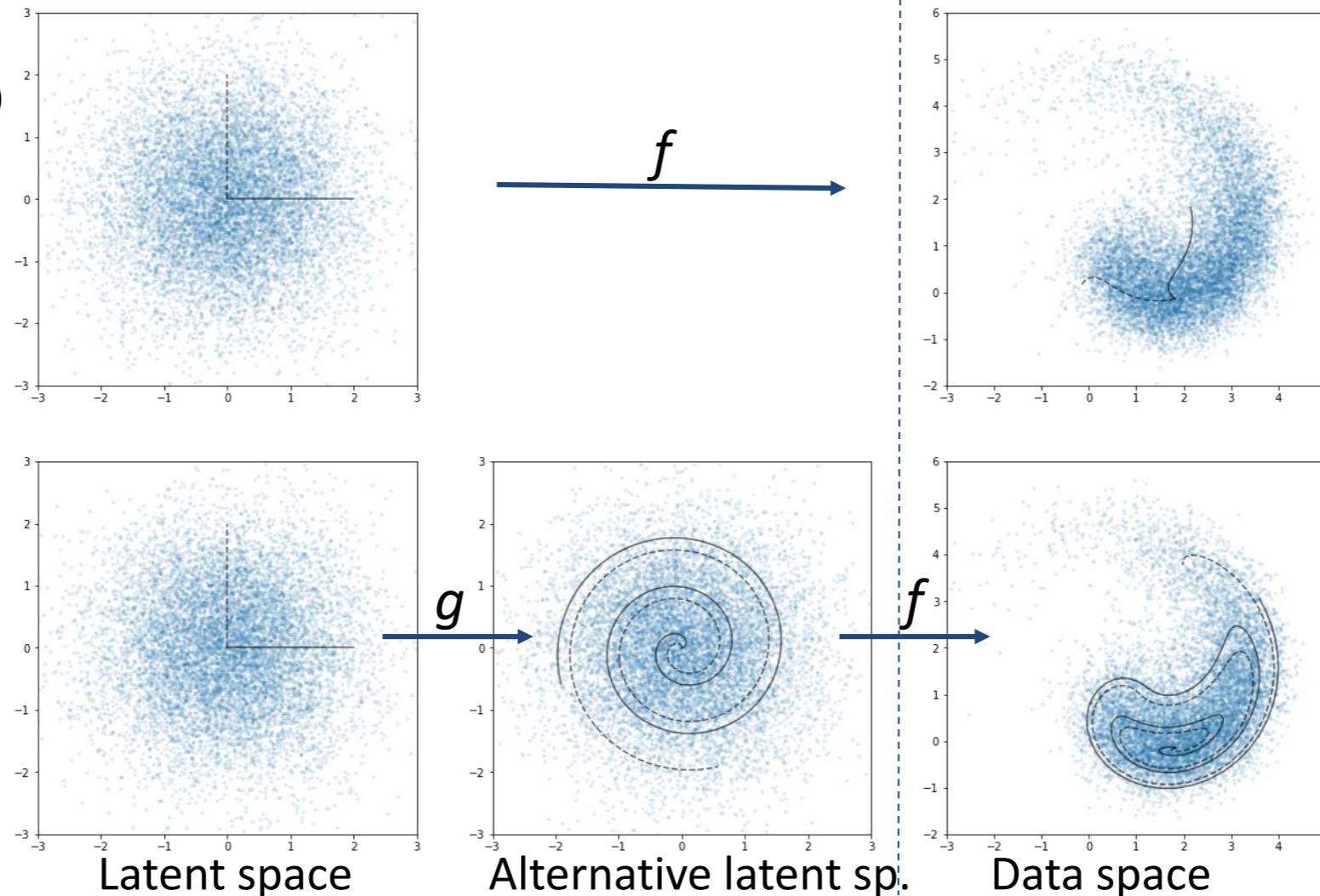
Data space



Non-linear ICA as a Disentanglement Method

- Disentanglement: undo the non-linear mixing of given data, recover latent space
- This in general impossible: non-linear mappings are **too flexible** \Rightarrow ambiguity unresolvable

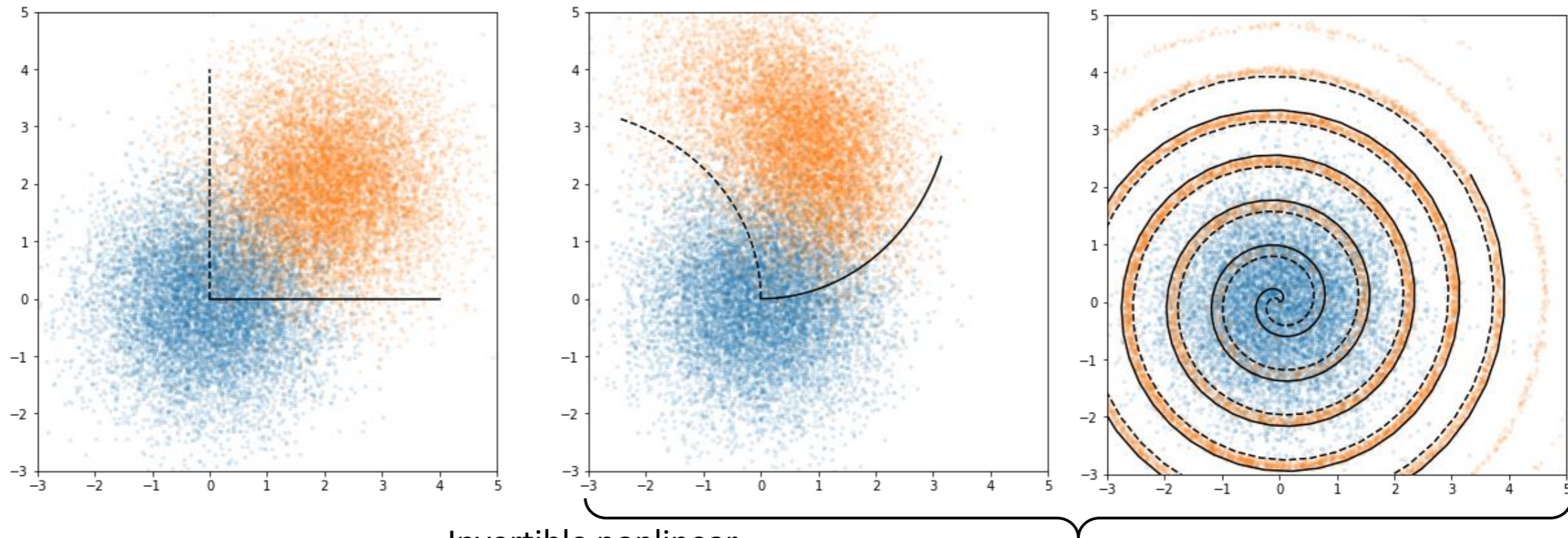
- Transformations $f(z)$ and $f(g(z))$ produce identical distributions
- \Rightarrow non-linear ICA is **unidentifiable**





ICA as a Disentanglement Method

- Fundamental insight: we need to **constrain** transformations g in the latent space
 - Constrain latent distributions by **conditioning**, e.g. by introducing a **mixture** distribution



Latent space
(mixture with **conditionally**
independent dimensions)

Invertible nonlinear
transformation

Alternative data spaces
(different solutions place the
colors differently)



LM-INNs for Disentanglement

LM-INNs fulfill the theoretical assumptions of **non-linear ICA**

Important negative result [Hyvärinen & Pajunen 1999]:

Fully unsupervised *non-linear* disentanglement is impossible

General non-linear transformations are too powerful – can fit everything

Recent positive results: non-linear disentanglement becomes identifiable with additional conditioning information, e.g.

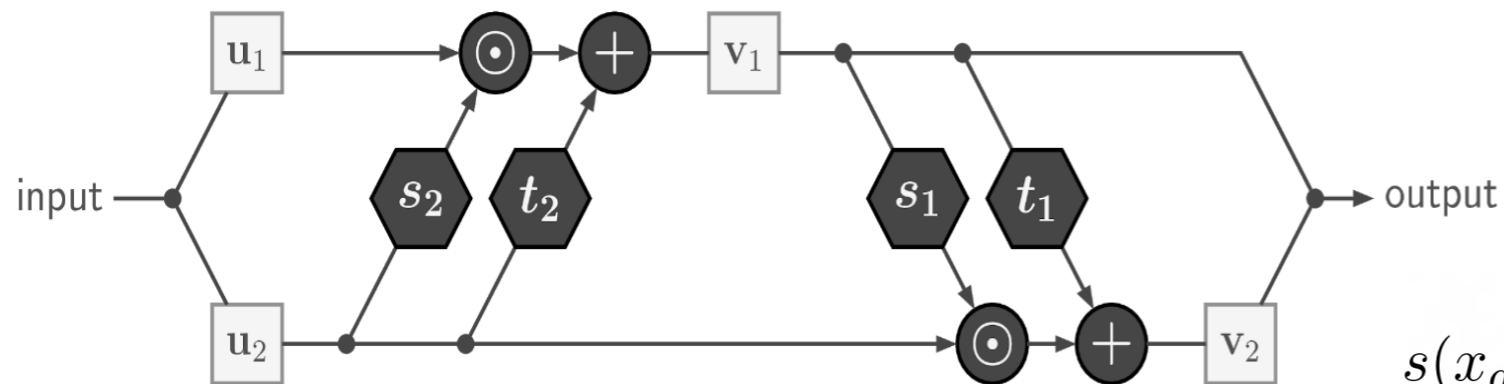
- Temporal relations [Hyvärinen & Morioka 2017, Hyvarinen, Sasaki, Turner 2018]
 - Multi-modal observations [Gresele, Rubenstein, Mehrjou, Locatello, Schölkopf 2020]
 - Class labels [our work]
- ⇒ **Mathematical guarantees that non-linear ICA finds the *true generative factors* and the *true intrinsic dimension* in certain situations** (generalizing this is a hot research topic).





General Incompressible-flow Networks (GIN)

- Modification of Real NVP coupling block architecture
 - Constrain the Jacobian to determinant 1
 - This differs from additive coupling (NICE): **space can be scaled** in some dimensions, when this is compensated for by a counter change in the remaining dimensions



$$\Rightarrow \sum_{i=1}^d s(x_i) = 0 \quad \Rightarrow \log |J| = 0$$

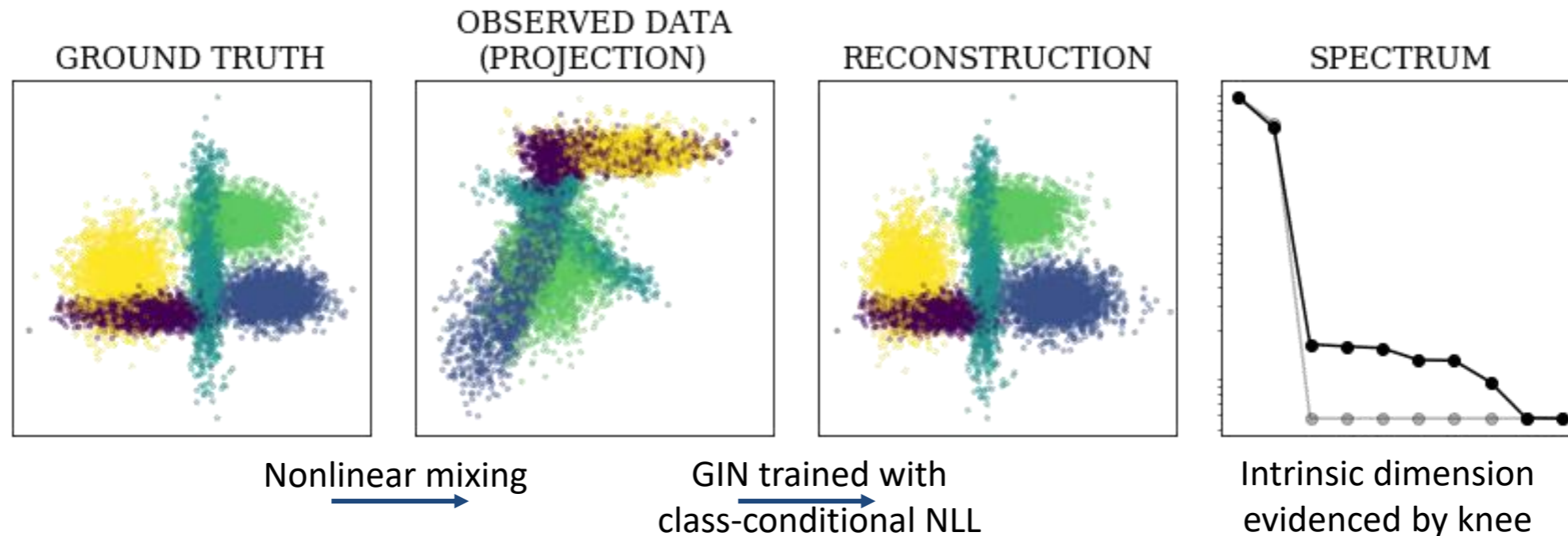
$$s(x_d) = - \sum_{i=1}^{d-1} s(x_i) \quad \text{or} \quad s(x_i) \leftarrow s(x_i) - \frac{1}{d} \sum_{i=1}^d s(x_i)$$

- Advantage:
 - Total “Variance” is preserved
 - ⇒ **Spectrum** of latent variables can be sorted and **interpreted as in PCA**



Artificial Data Experiments

- True generative process:
 - 5 Gaussian mixture components (“class labels”), 2 meaningful dimensions, 8 noise dimensions
 - Mapped to 10-dimensional data space using a random non-linear transformation
- Task of the INN:
 - Determine that intrinsic dimension is 2
 - Recover the GMM within the meaningful dimensions, given the class labels

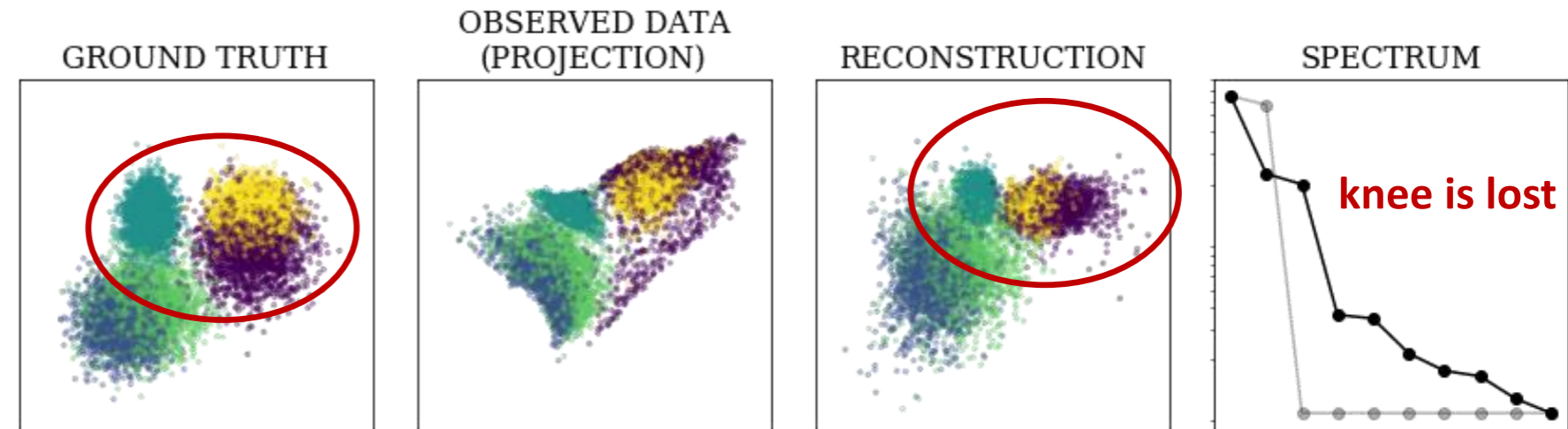




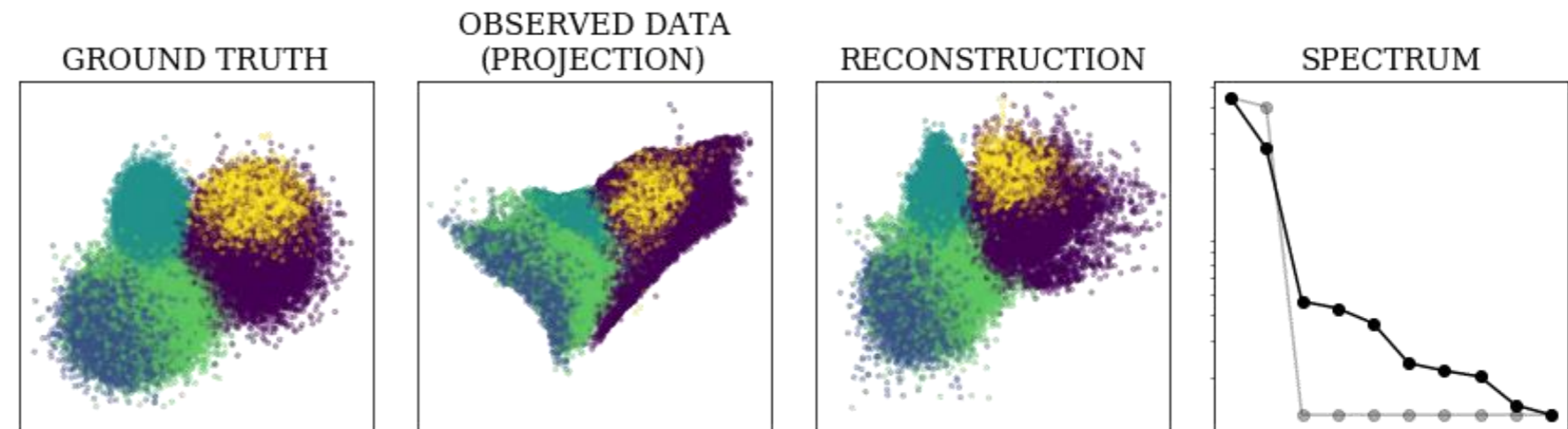
Artificial Data Experiments

- Intuition: works, because all clusters must be disentangled **simultaneously**
 - Breaks down, when clusters have **no overlap**: model transforms these clusters independently
 - Can be caused by lack of training data:

10^4 data points



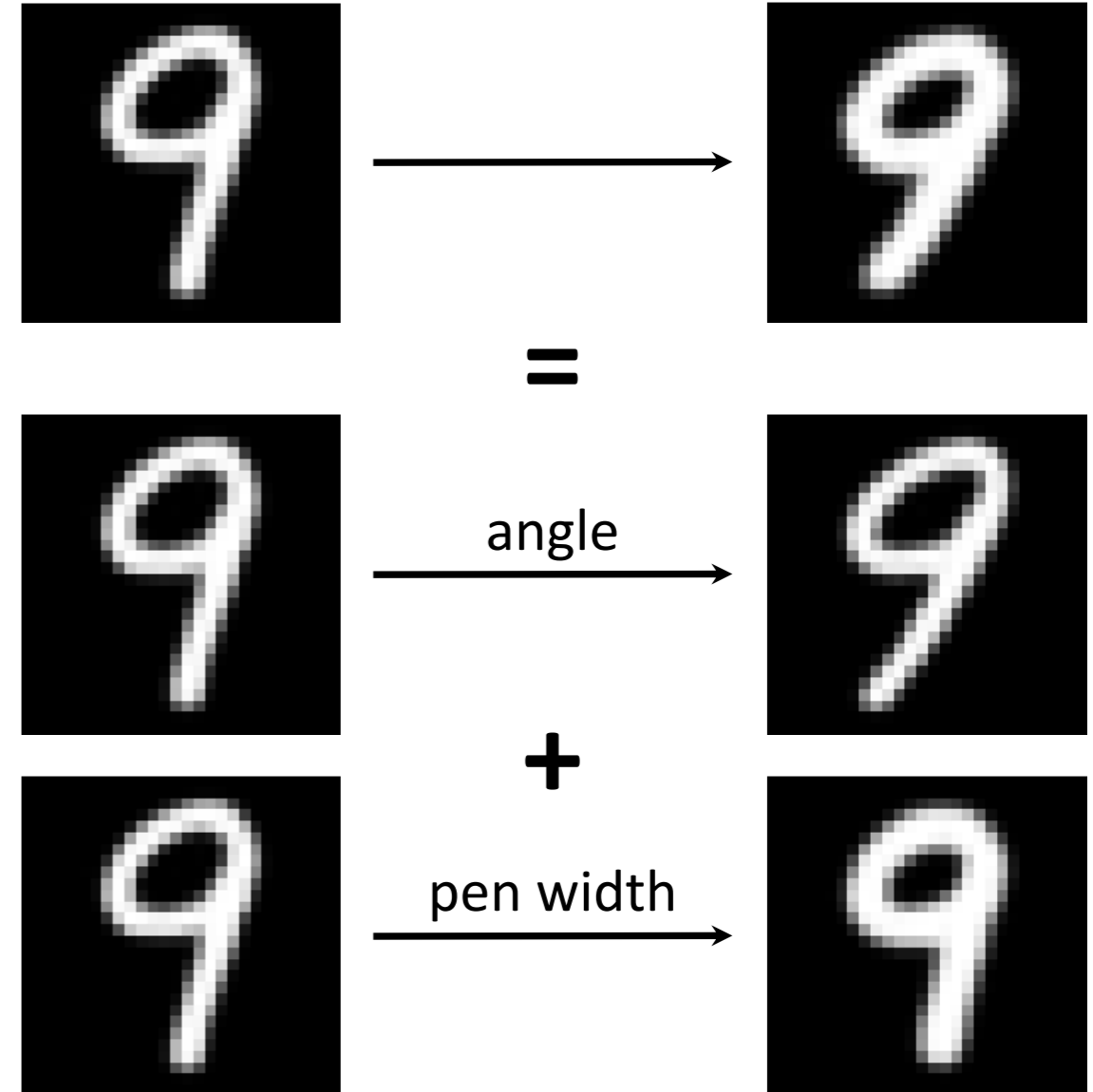
10^5 data points





LM-INNs for Disentanglement

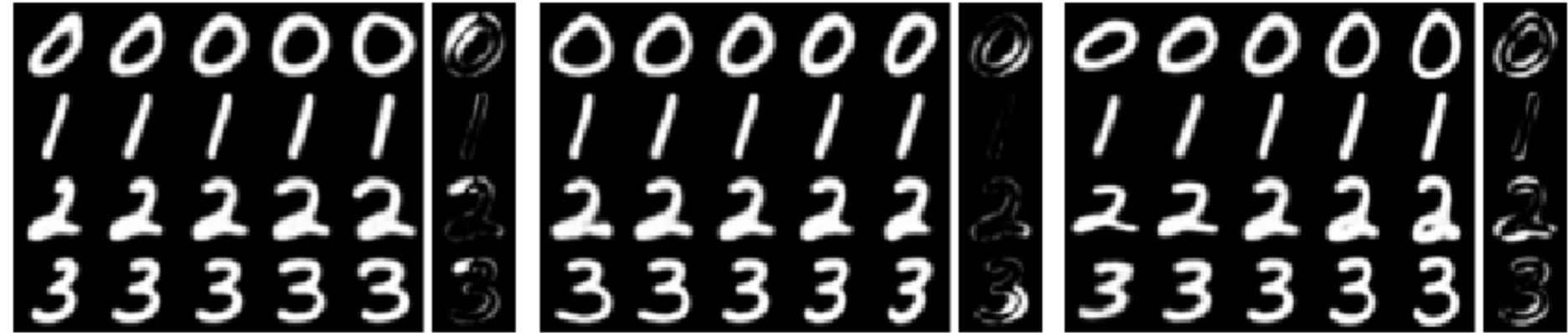
- Identify **intrinsic factors** of complicated data distributions, which intuitively explain variability
- Express complex/coordinated changes of the data as a **combination** of simple changes in the factors
- Example: EMNIST handwritten digits: latent factors are characteristics of handwriting styles





Application to EMNIST

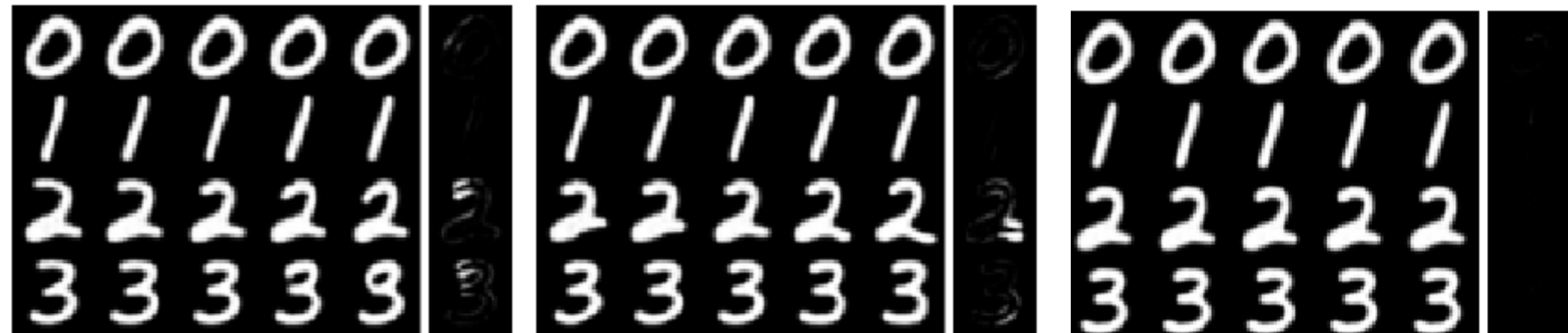
- First 8 latent variables control global properties
- Following 14 control local shape
- Remaining 762 have no visible effect



Var. 1: upper width

Var. 8: lower width

Var. 3: height



Var. 13: top left of 2,3,7

Var. 16: tail of 2

Var. 23: no effect

difference images



Application to EMNIST

Latent space
interpolation



Independent effect
of first 8 most
significant latent
dimensions



(animations not visible
in PDF version)



IB-INNs – Building (more) interpretable models with INN-based generative classifiers

Ullrich Köthe

Visual Learning Lab, Heidelberg University

joint work with **Lynton Ardizzone, Radek Mackowiak, Jakob Kruse, Carsten Rother**

Tutorial „Normalizing Flows“ at CVPR 2021

June 2021



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386



European Research Council
Established by the European Commission



IB-INNs: Generative Classifiers

- What is a generative classifier (GC)?
 - Classifier: given image x , predict label y of most salient object
 - A discriminative classifier (DC): learns the class **posterior probability** $p(y | x)$
 - Generative classifier: instead learns the **data likelihood** $p(x | y)$ and computes the posterior indirectly by **Bayes rule**:

$$p(y | x) = \frac{p(x | y)p(y)}{\sum_{y'} p(x | y')p(y')}$$

- GCs promise to foster **reliability and interpretability**
 - uncertainty quantification, outlier detection, robustness against distribution shifts
 - discovery of meaningful features
 - *but*: predictive performance of GCs used to be unconvincing \Rightarrow discriminative classifiers (DCs) prevailed
- Old idea, but so far discriminative classifiers have much better performance






The Information Bottleneck Principle

- Naively trained generative classifiers: **poor classification** accuracy in comparison to DCs
 - Tend to overfit
- Information bottleneck principle overcomes this problem
 - Introduce latent representation z , where all information flows through – “**bottleneck**”
 - Latent variables z should be: **highly informative for y** (= good classification)
keep only as much information about x as needed (= no overfitting)
- Minimize Information Bottleneck (IB) loss

$$\mathcal{L}_{\text{IB}} = I(x, z) - \beta \cdot I(y, z)$$

Generative aspect (minimize spurious information about x)
 
 Discriminative aspect (maximize information about class labels)

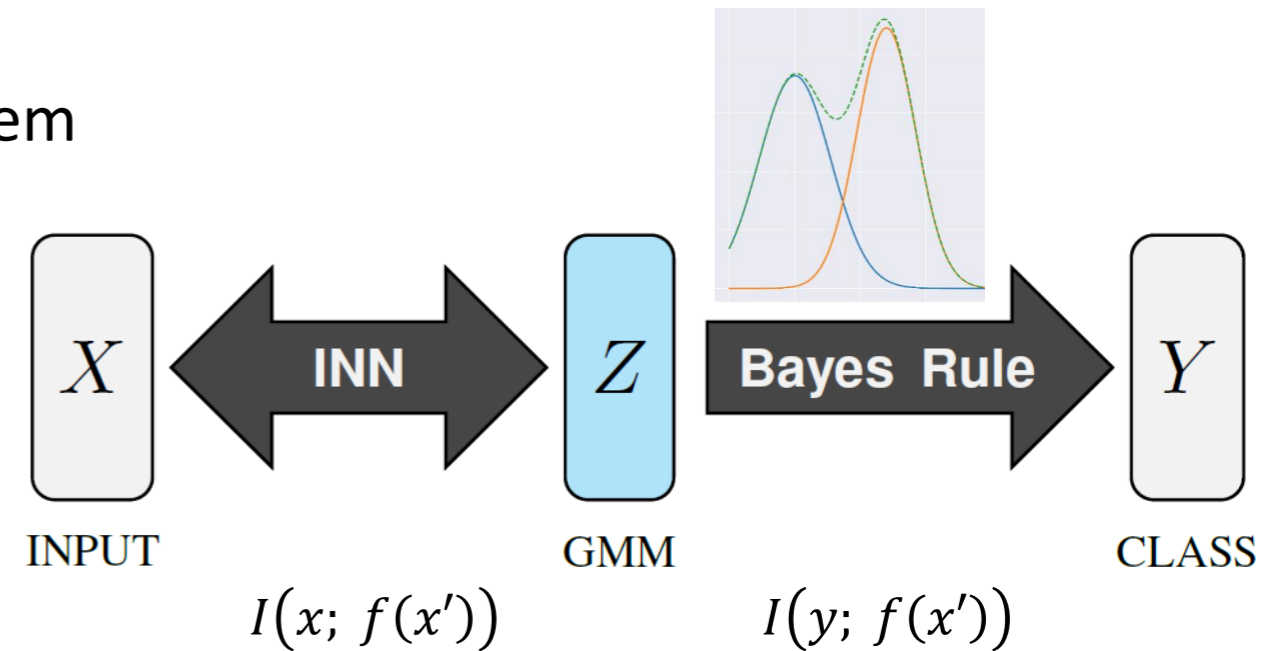
Trade-off parameter

with Mutual Information (MI) $I(y, z) = D_{\text{KL}}\left(p(y, z) \parallel p(y)p(z)\right)$



IB-INNs: Generative Classifiers

- Learning $p(x | y)$ is a **density estimation** problem
 - Normalizing flows are good at density estimation
 - We actually model $p(z | y)$ as a latent GMM and train an INN to transform this into $p(x | y)$
 - ⇒ this is a latent mixture INN
 - The model can be trained using the **Information Bottleneck Principle**



- Problem: INNs are lossless encoders – where is the bottleneck?
 - Train the INN mapping $z = f(x)$ with **noise augmented data** $x' = x + \epsilon$:

$$\mathcal{L}_{IB} = \lim_{\epsilon \rightarrow 0} I(x; z_{\epsilon} = f(x')) - \beta \cdot I(y; z_{\epsilon} = f(x'))$$

- Intuitively: noise ensures lossy encoding ⇒ prevents divergence of $I(x; z_{\epsilon} = f(x'))$
- Surprisingly, mutual information $I(x; z_{\epsilon})$ reduces to the usual maximum likelihood loss

$$I(x; z_{\epsilon}) = \mathbb{E}_{p(x), p(\epsilon)} [-\log p(x + \epsilon)] = \mathbb{E}_{p(x), p(\epsilon)} [-\log p_Z(z_{\epsilon}) - \log \det \nabla f]$$



IB-INN: Training an LM-INN as a Generative Classifier

LM-INN can approximate the IB loss arbitrarily well \Rightarrow IB-INN

- Successfully trained on CIFAR-10 (10 classes, 32^2 images) and ImageNet (1000 classes, 224^2 images)
- Depending on β , the IB-INN emphasizes generative or discriminative performance
 - at $\beta = 1$, bits/dimension (=generative performance) comparable to a purely generative model
 - at $\beta = 32$, test accuracy (=discriminative performance) comparable to a discriminative ResNet
 - ImageNet: **good trade-off at $\beta = 8$**



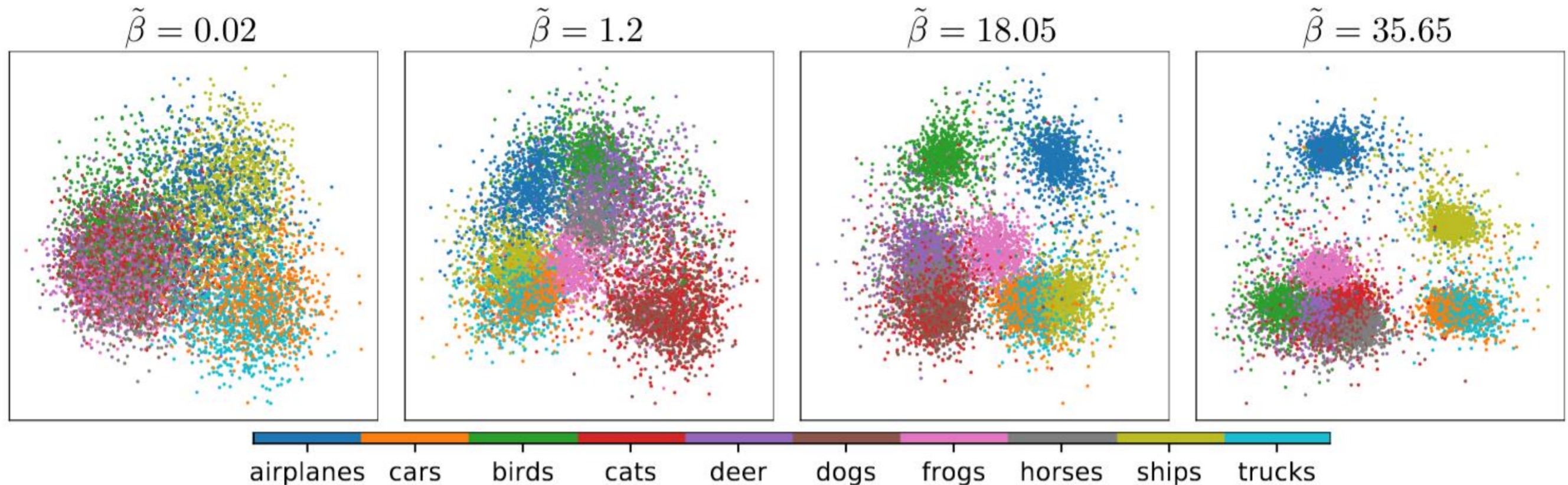
β	Bits/dim. (\downarrow)	Acc. (%) (\uparrow)
1	4.34	67.30
2	6.14	71.73
4	8.72	73.69
8	12.35	74.59
16	17.43	75.54
32	22.68	76.18
∞	47.01	76.27
0	2.59	–
ResNet	–	77.40



IB-INNs:

Benefits of GC (1): Interpretability

- Class separation improves as β (= importance of $I(y; z)$) increases
 - CIFAR-10 examples (PCA projection of latent space)





IB-INNs:

Benefits of GC (1): Interpretability

- Heatmaps for attention area of the most probable classes
 - Back-project relevant latent features to image space regions
 - Thanks to invertibility, the heat-maps **represent the true decision process**, not a post-hoc explanation

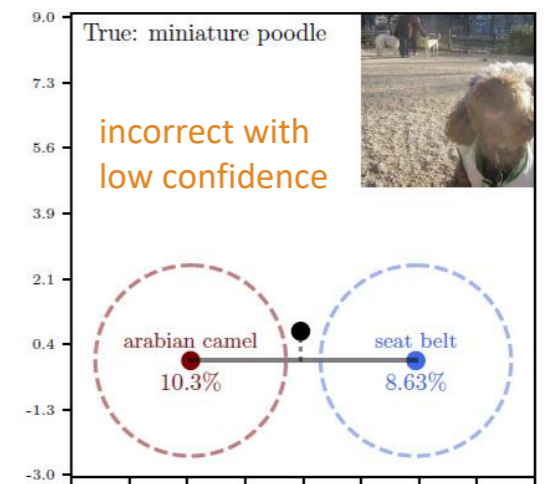
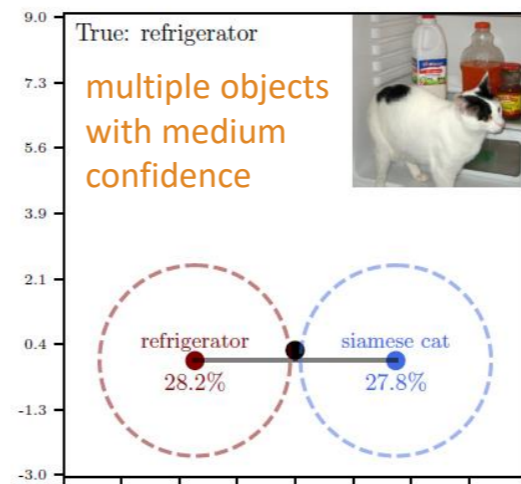
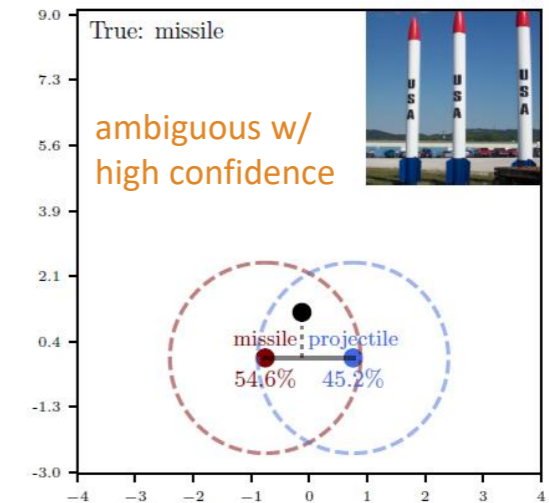
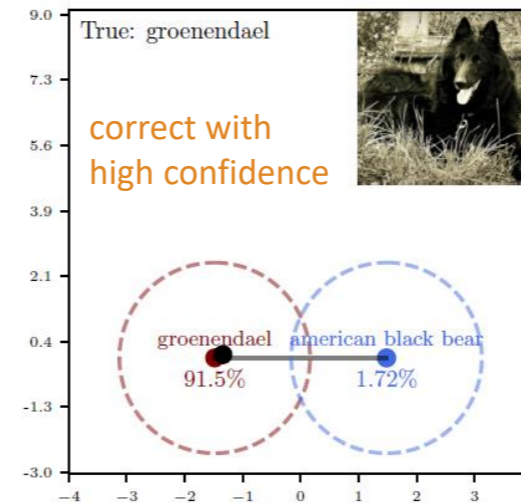
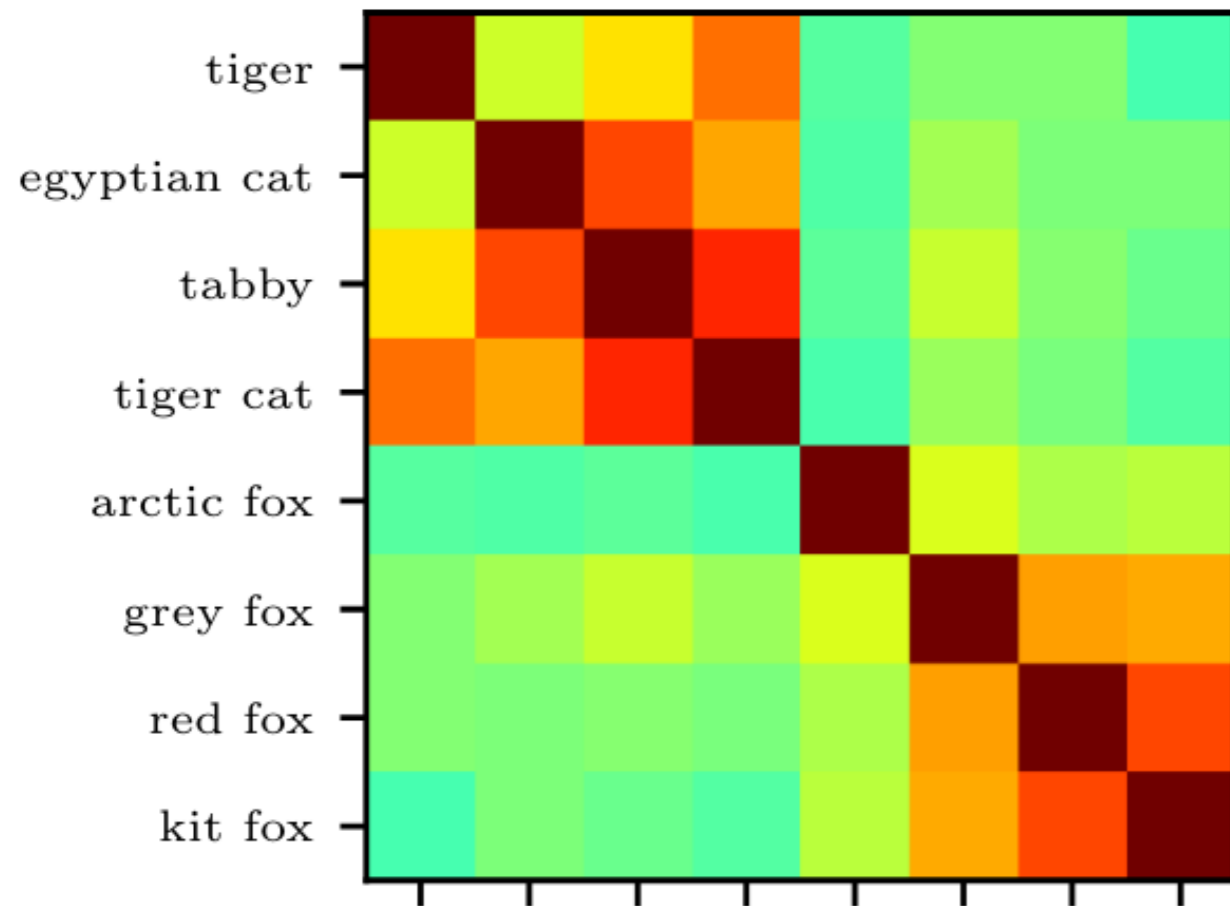




IB-INNs:

Benefits of GC (1): Interpretability

- Pairwise distances between class centers in z-space reflect class similarity and confidence
 - ImageNet examples

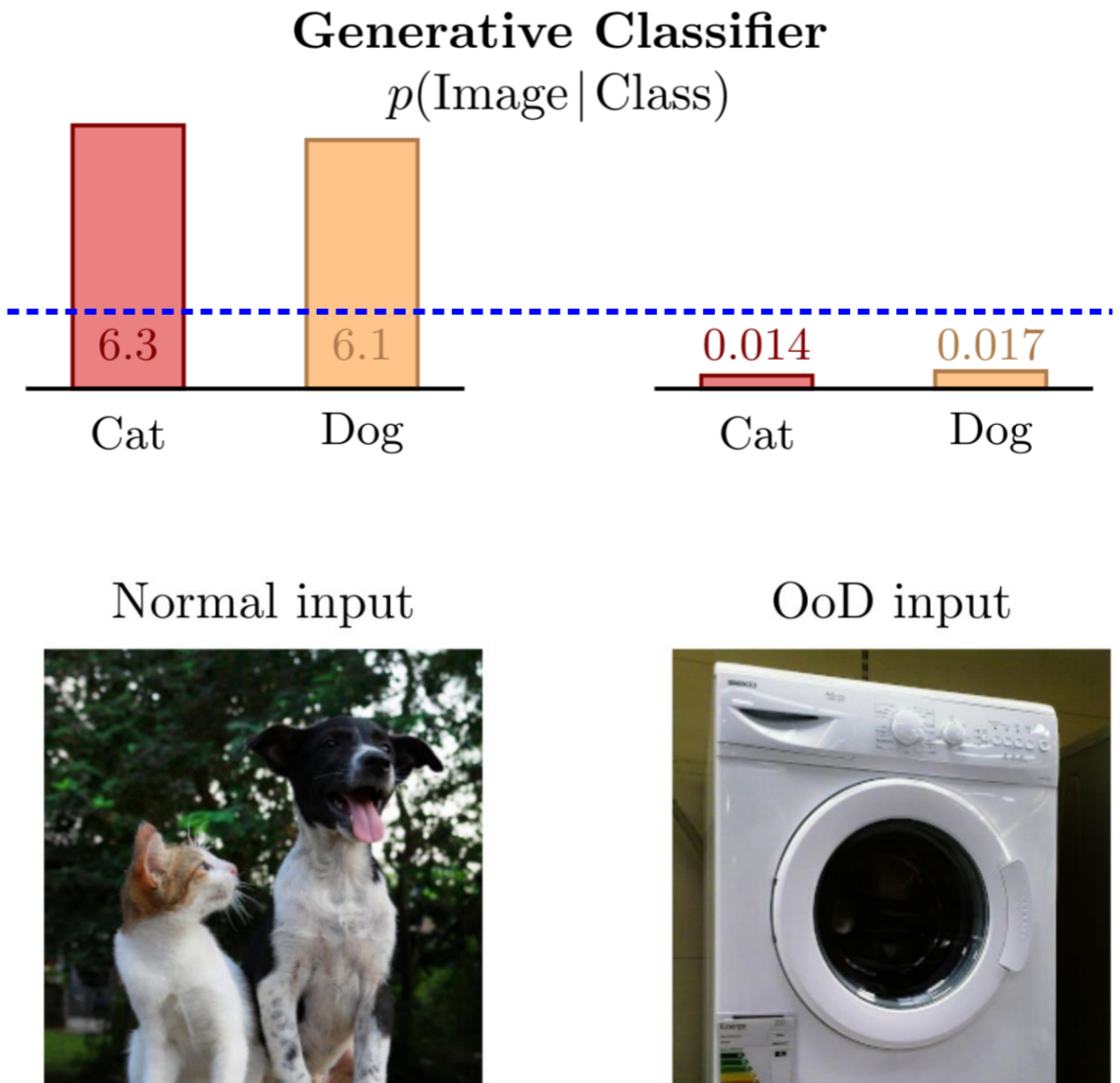




IB-INNs:

Benefits of GC (2): Out-of-Distribution Detection

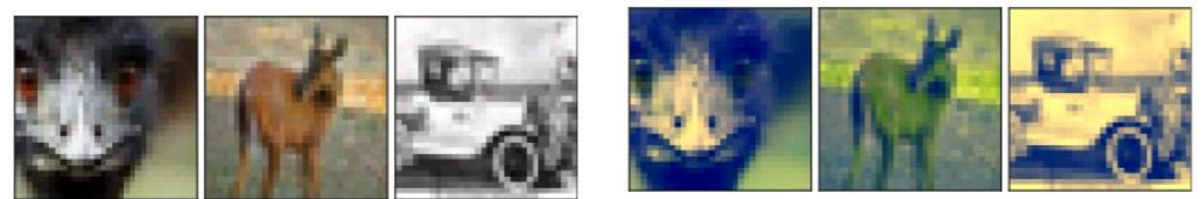
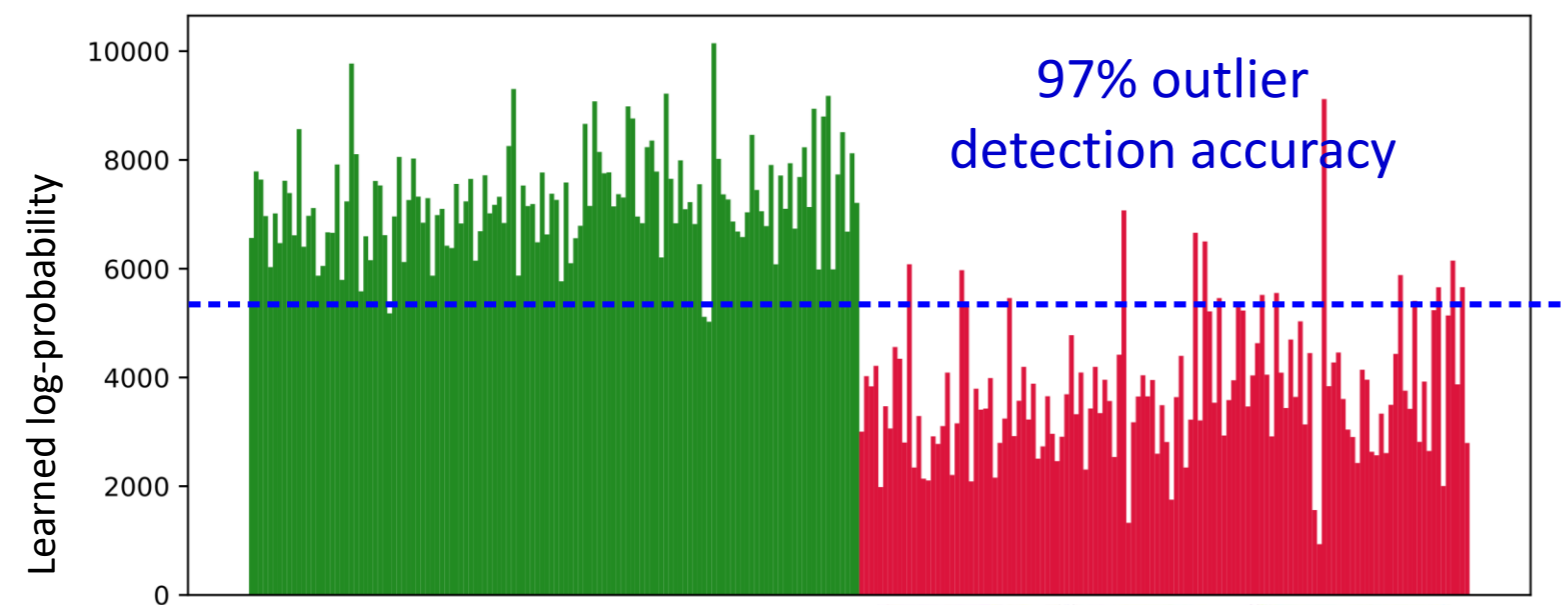
- Outliers have low likelihood for every class
 - Intuitively: can separate in-/outliers using threshold on likelihood
 - Many interesting open questions
 - Which outlier scores does IB-INN support? (e.g. typicality tests, WAIC, ...)
 - What does it mean for an instance to be an outlier in *high dimensions*? (much of our intuition is based on low dimensional case and thus misleading)
 - Which latent re-parameterizations are sensitive for which type of outlier?





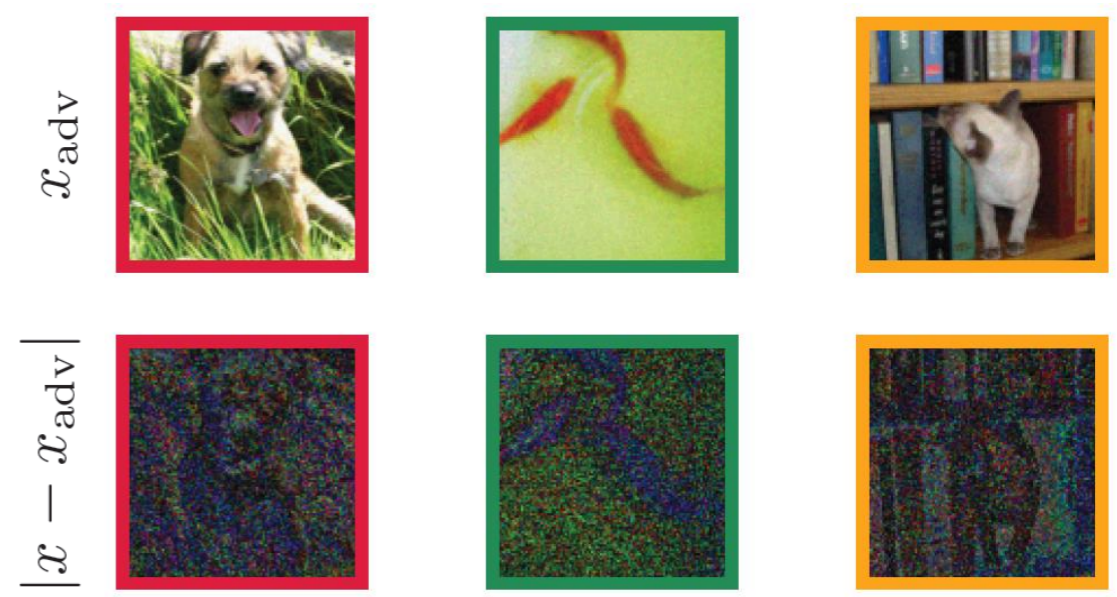
IB-INNs: Benefits of GC (2): Out-of-Distribution Detection

- Outliers have low likelihood for every class
 - Artificial outliers: scrambled colors (CIFAR-10)



Inlier images (test set) Same images, scrambled colors

Adversarial examples (ImageNet)



Minimal perturbations to get confident incorrect predictions

⇒ IB-INN improves adversarial robustness, but does not in itself solve the problem

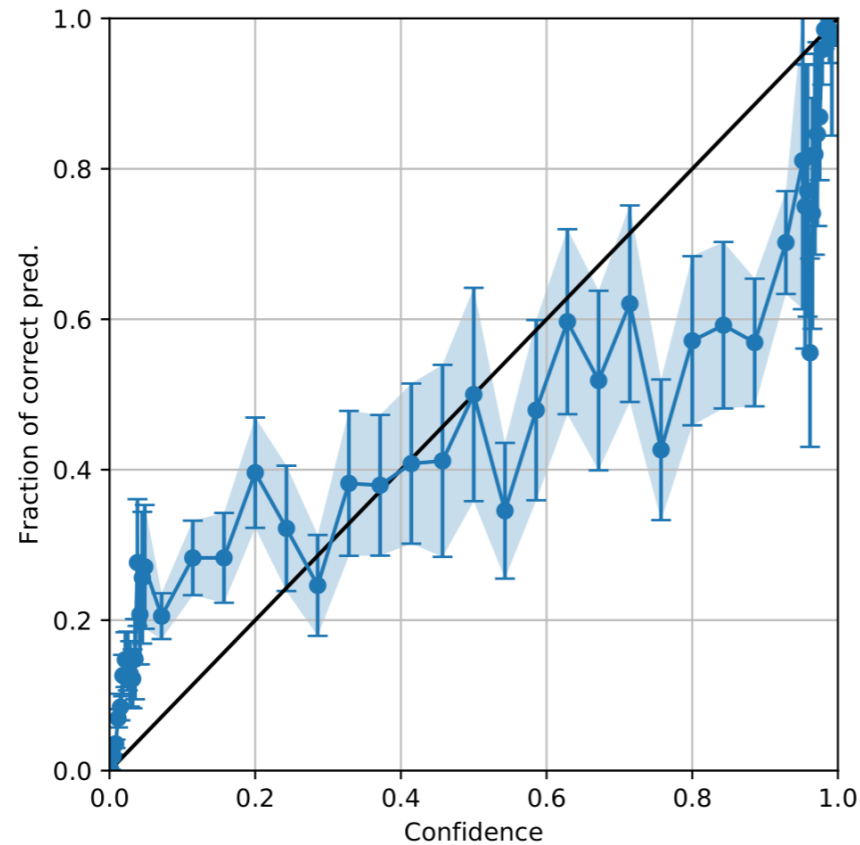


IB-INNs:

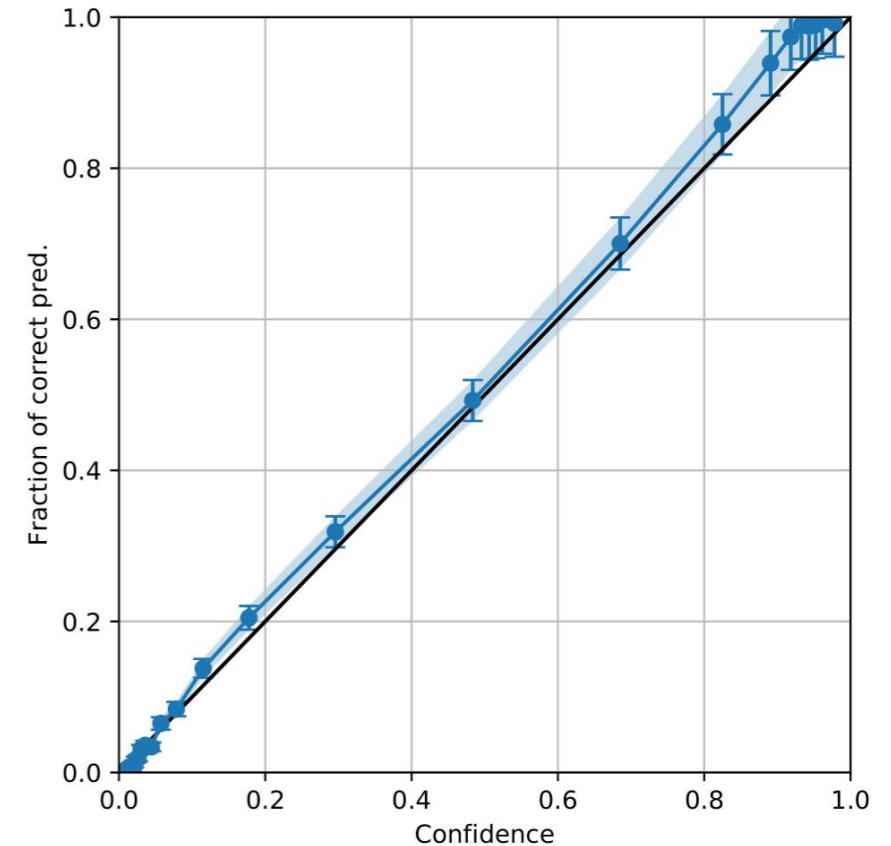
Benefits of GC (3): Uncertainty Calibration

- Calibration = consistency of confidence vs. actual performance
 - If classifier is 90% confident about class label, it should be right 90% of the time, neither less nor more
 - Problematic for discriminative classifiers [Guo et al. 2017] – IB-INNs are much better calibrated

DC: ResNet-18
(CIFAR-10)



GC: IB-INN $\beta=1$
(CIFAR-10)





Summary

Public code of our FrEIA library: <https://github.com/VLL-HD/FrEIA>

- INNs are very good density estimators:
 - Not yet quite as good as GANs (as trained by the Big Guys with 300 GPUs in parallel 😊)
 - But with much stronger mathematical interpretation and guarantees
- Three main approaches to incorporate additional information
 - Conditional INN: learn $p_{\mathbf{z}}(\mathbf{z} = f_{\text{INN}}(\mathbf{x}; \mathbf{y}))$
 - Latent mixture INN: learn $p_{\mathbf{z}}(\mathbf{z} = f_{\text{INN}}(\mathbf{x}) | \mathbf{y})$
 - Augmented latent space INN: learn $p_{\mathbf{y}, \mathbf{z}}(\mathbf{y}, \mathbf{z} = f_{\text{INN}}(\mathbf{x}))$
 - We get the full posterior $p(x | y)$, both exactly and through samples
- Future work:
 - Improve architectures and training
 - Strengthen validation and mathematical guarantees
 - Apply to various problems in natural and life sciences
 - Better incorporation of prior knowledge from the application domain



Thanks to our team and collaborators!



Visual Learning Lab, Uni Heidelberg:

Lynton Ardizzone

Jakob Kruse

Jens Müller

Felix Draxler

Radek Mackowiak

Peter Sorrenson

Carsten Rother

York University, Canada:

Marcus Brubaker

German Cancer Research Center, Heidelberg:

Tim Adler, Sebastian Wirkert, Lena Maier-Hein

Inst. for Theoretical Astrophysics, Uni Heidelberg:

Victor Ksoll, Ralf Klessen

Inst. for Environmental Physics, Uni Heidelberg:

André Butz, Florian Kleinicke

Psychologisches Institut, Uni Heidelberg:

Stefan Radev, Ulf Mertens, Andreas Voss



References

Public code of our INNs and papers in the FrEIA library: <https://github.com/VLL-HD/FrEIA>

Adler, T., et al.: “Uncertainty-aware performance assessment of optical imaging modalities with invertible neural networks”, Intl. J. Computer Assisted Radiology and Surgery 14(6):997–1007 (2019).

Adler, T. J., et al. „Out of distribution detection for intra-operative functional imaging“, In: *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging and Clinical Image-Based Procedures* (pp. 75-82), arXiv:1911.01877 (2019).

Ardizzone, L. et al.: “Analyzing inverse problems with invertible neural networks”, ICLR 2019, arXiv:1808.04730 (2018).

Ardizzone, L., Lüth, C., Kruse, J., Rother, C., Köthe, U.: “Conditional Invertible Neural Networks for Diverse Image-to-Image Translation”, GCPR 2020, arXiv:1907.02392 (2019).

Ardizzone, L., Mackowiak, R., Kruse, J., Rother, C., Köthe, U.: “Training Normalizing Flows with the Information Bottleneck for Competitive Generative Classification”, NeurIPS 2020 (oral), arXiv:2001.06448 (2020).

Behrmann, J., Duvenaud, D., & Jacobsen, J. H.: “Invertible residual networks”, ICML 2019, arXiv:1811.00995 (2018).

Bellagente, M. et al.: “Invertible Networks or Partons to Detector and Back Again”, SciPost Phys. 2020, arXiv:2006.06685 (2020).

Bieringer, S. et al.: “Measuring QCD Splittings with Invertible Networks“, arXiv:2012.09873 (2020).

Bengio, Y., et al.: “Representation Learning: A Review and New Perspectives“, IEEE PAMI 35(8):1798-1828 (2013).

Bloem-Reddy, B., & The, Y.W.: “Probabilistic symmetry and invariant neural networks“, JMLR 21(90):1–61(2020).

Draxler, F., Schwarz, J., Schnörr, C., Köthe, U.: “Characterizing the Role of a Single Coupling Layer in Affine Normalizing Flows“, GCPR (best paper finalist) (2020).

Dinh, L., Sohl-Dickstein, J., Bengio, S.: “Density estimation using Real NVP“, arXiv:1605.08803 (2016).

Guo, Ch., et al.: “On calibration of modern neural networks“, ICML (2017).

Khemakhem, I., et al.: “Variational autoencoders and nonlinear ICA: A unifying framework“, AISTAT (2020).





References

- Kingma, D., & Dhariwal, P.: “Glow: Generative flow with invertible 1x1 convolutions”, NIPS (2018).
- Kobyzev, S., Prince, JD., Brubaker, M.: “Normalizing Flows: An Introduction and Review of Current Methods”, arXiv:1908.09257 (2019).
- Kruse, J., et al.: “Benchmarking Invertible Architectures on Inverse Problems”, ICML Workshop INNF, arXiv:2101.10763 (2019).
- Kruse, J., Detommaso, G., Köthe, U., Scheichl, R. “HINT: Hierarchical Invertible Neural Transport for Density Estimation and Bayesian Inference”, AAAI 2021, arXiv:1905.10687 (2019)
- Ksoll, V. et al.: “Stellar Parameter Determination from Photometry using Invertible Neural Networks”, MNRAS, arXiv:2007.08391 (2020).
- Ksoll, V. et al.: “Measuring Young Stars in Space and Time”, part 1: arXiv:2012.00521 , part 2: arXiv:2012.00524 (2020).
- Lee, W., et al.: “High-Fidelity Synthesis with Disentangled Representation”, arXiv:2001.04296 (2020).
- Mackowiak, R., Ardizzone, L. et al.: “Generative Classifiers as a Basis for Trustworthy Computer Vision”, CVPR 2021 (oral), arXiv:2007.15036 (2020).
- Müller, J., Schmier, R., Ardizzone, L., Rother, C., Köthe, U. (2020). “Learning Robust Models Using The Principle of Independent Causal Mechanisms” arXiv:2010.07167 (2020).
- Nölke, JH. et al.: “Invertible Neural Networks for Uncertainty Quantification in Photoacoustic Imaging”, arXiv:2011.05110 (2020).
- Papamakarios, G. et al.: “Normalizing Flows for Probabilistic Modeling and Inference”, JMLR, 22(57):1-64, arXiv:1912.02762 (2021).
- Putzky, P., & Welling, M.: “Invert to Learn to Invert”, NeurIPS (2019).
- Radev, S., et al.: “BayesFlow: Learning Complex Stochastic Models with Invertible Neural Networks”, IEEE TNNLS, arXiv:2003.06281 (2020).
- Radev, S., et al.: “Model-based Bayesian inference – an application to the COVID-19 pandemics in Germany”, arXiv:2010.00300 (2020).
- Shiono, T.: “Estimation of Agent-Based Models Using Bayesian Deep Learning Approach of BayesFlow”, SSRN:3640351 (2020).
- Sorrenson, P., Rother, C., Köthe, U.: “Disentanglement by Nonlinear ICA with General Incompressible-flow Networks (GIN)”, ICLR 2020, arXiv:2001.04872 (2020).
- Tishby, N., et al.: “The Information Bottleneck Method”, Allerton Conf. Communic., Control, Computing, arXiv:0004057 (1999).

