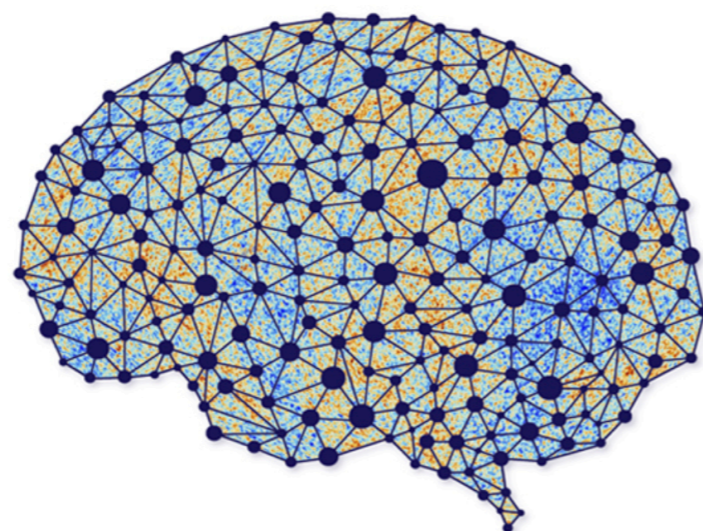


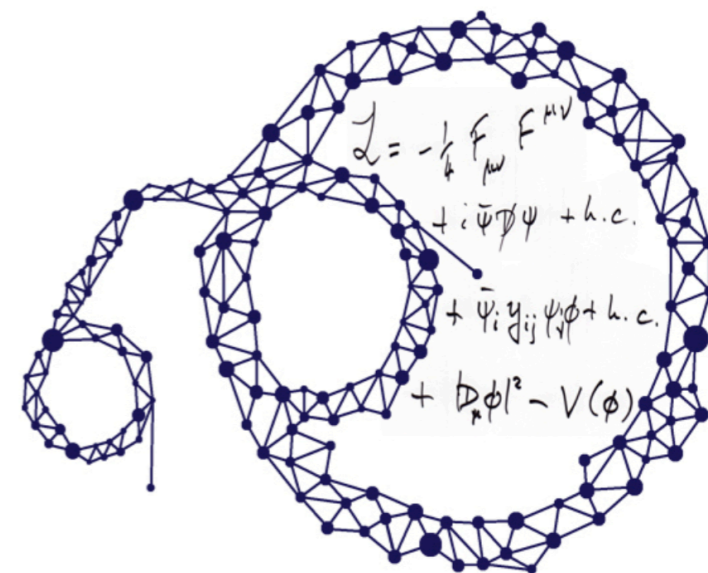
# Physics 361 - Machine Learning in Physics

## Lecture 24 – Normalizing flows

April 18<sup>th</sup> 2024



AI  
∩  
Universe



Moritz Münchmeyer

# Normalizing Flows

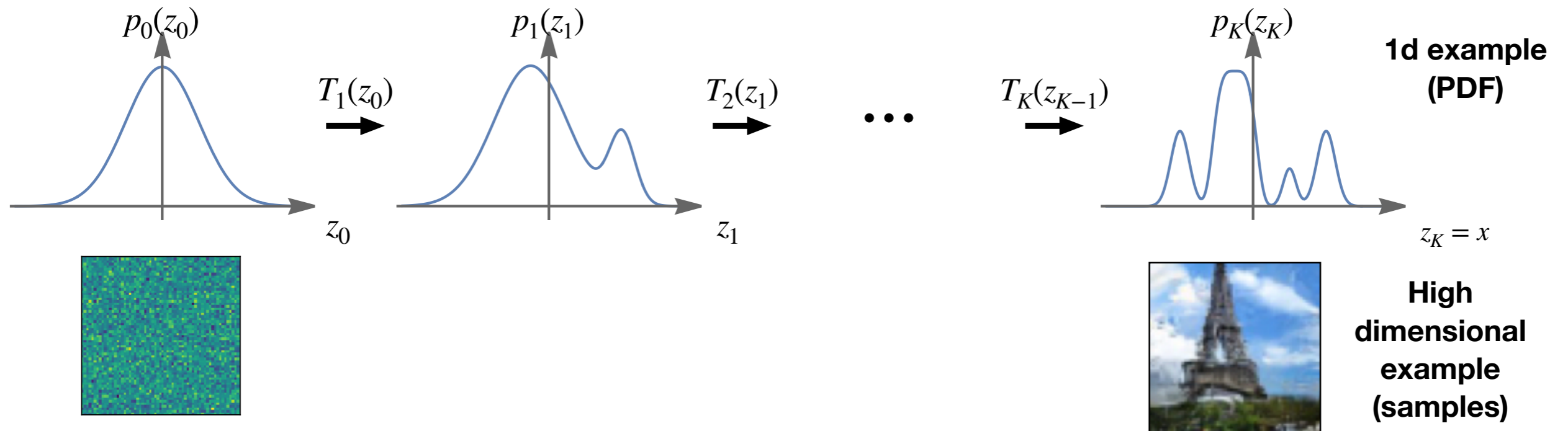
**Intro to Normalizing flows**

# Basics of Normalizing Flows

# Normalizing flows

- Normalizing flow: Series of learned transformations that **deform a simple base distribution into a complicated target distribution.**

Learned transformations  $T_i$   
e.g. parametrized by a neural network



- Difference with most other ML methods: We learn a probability distribution, rather than an arbitrary input->output mapping.
- Review: <https://arxiv.org/abs/1912.02762>. Widely used in physics e.g. in QFT, likelihood-free inference and cosmology

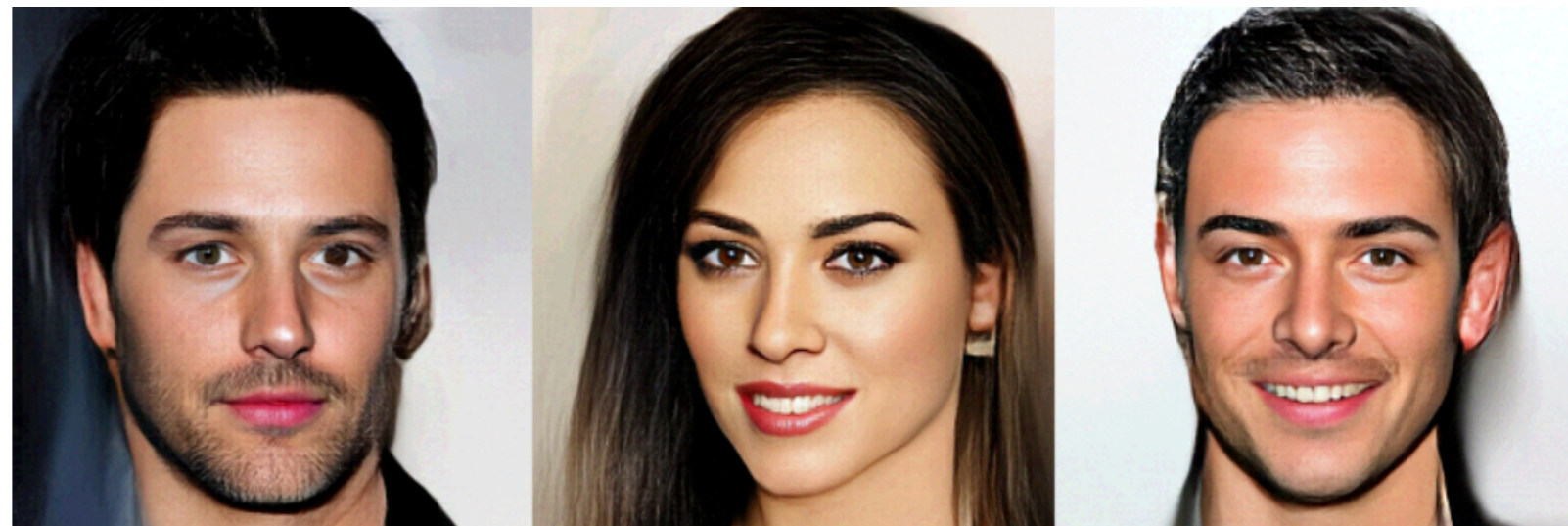
# Normalizing flows are generative models

- Like GANs and diffusion models, normalizing flows are generative models.
- They can be used to generate images too. However they are not currently as good at that as these other models.
- But they can do something other models cannot: give a normalized probability density for the sample. They are real PDFs.

Real-NVP flow



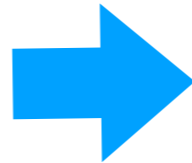
Glow flow



# Normalizing flows

- Transformation  $T$  (the “flow”)

$$\mathbf{x} = T(\mathbf{u})$$



- Change of variables of PDF

$$p_x(\mathbf{x}) = p_u(\mathbf{u}) |\det J_T(\mathbf{u})|^{-1}$$

- Chain many “simple” transformations together to make a complicated distribution:

$$T = T_K \circ \dots \circ T_1$$

$$\mathbf{z}_k = T_k(\mathbf{z}_{k-1})$$

- After training **two basic operations can be performed:**

- Exact density evaluation** (backward mode)

**Sample  $\mathbf{x}$    $p(\mathbf{x})$**

- Sampling from the distribution** (forward mode)

**Base distribution sample  $\mathbf{u}$   Target sample  $\mathbf{x}$**

# Basics and definitions

- data distribution:  $\vec{x}$ : D-dimensional real vector  
 $\vec{x} \sim P(\vec{x})$

E.g.:  $\vec{x} = \begin{pmatrix} \text{height} \\ \text{age} \\ \text{weight} \\ \text{strength} \end{pmatrix}$  of people

- Main idea of normalizing flows:  
Express  $\vec{x}$  as a transformation  $T$  of a real vector  $\vec{u}$  sampled from a simple base distribution.

$$\vec{x} = T(\vec{u}) \quad \text{where} \quad \vec{u} \sim P(\vec{u})$$

E.g.:  $\vec{u}$  is a Gaussian random variable



# Properties of the transformation

- To specify  $p(x)$  we need to specify
    - The transformation  $T(x; \theta)$  with "learned" parameters  $\theta$
    - The base distribution  $p_u(u; \phi)$  with "learned" parameters  $\phi$  } often kept static
  - Defining properties of transformation  $T$ :
    - $T$  must be invertible
    - Both  $T$  and  $T^{-1}$  must be differentiable
- Thus  $T$  is a diffeomorphism.
- $\Rightarrow u = T^{-1}(x)$  must also be  $D$ -dimensional.



# Properties of the transformation

- The flow transformation is a change of variables.

$$P_x(x) = P_u(u) \left| \det J_T(u) \right|^{-1}$$

$$\text{where } u = T^{-1}(x)$$

with Jacobian

$$J_T(u) = \begin{bmatrix} \frac{\partial T_1}{\partial u_1} & \dots & \frac{\partial T_1}{\partial u_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial T_D}{\partial u_1} & \dots & \frac{\partial T_D}{\partial u_D} \end{bmatrix}$$

$|\det J_T(u)|$  quantifies the local change of volume that "molds"  $P(u)$  into  $P(x)$ .

# Composition of transformations

- Invertible and differentiable transformations are *composable*. In practice we stack many simple base transformations:

$$T = T_K \circ T_{K-1} \circ \dots \circ T_1$$

$$\det J_{T_2 \circ T_1}(u) = \det J_{T_2}(T_1(u)) \det J_{T_1}(u) \text{ etc.}$$

- Example: Flow from Gaussian to cross shape with 4 transf.

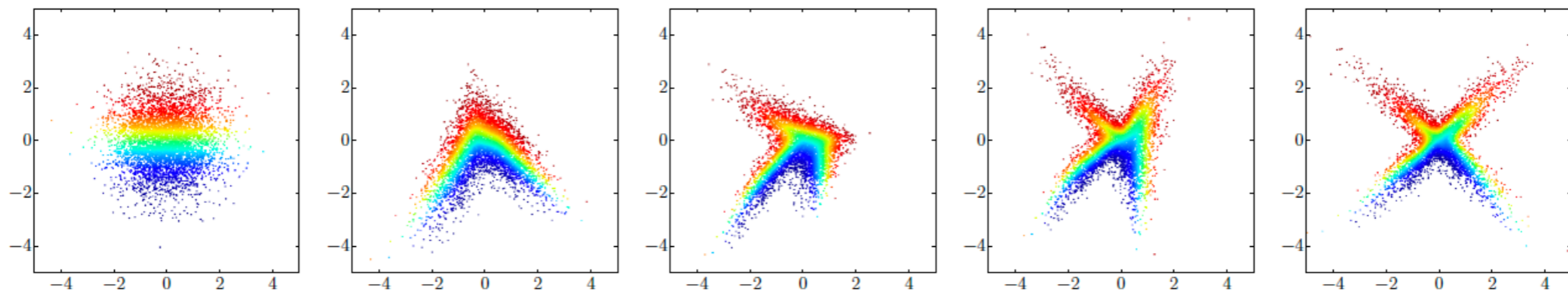


Figure 1: Example of a 4-step flow transforming samples from a standard-normal base density to a cross-shaped target density.

# Forward and backward use of the flow

- The flow (after training) has two fundamental operations:

- **Sampling**: Draw samples from  $p(u)$  and transform

$$x = T(u)$$

to get samples from  $x$ .

Other generative models (GANs, VAE) can also sample.

- **Density evaluation**: From a sample  $x$ , we can calculate  $p(x)$  as

$$p_x(x) = p_u(u) |\det J(u)|^{-1}$$

# Computational tradeoffs

- Sampling and density evaluation have different computational requirements.
- Density evaluation requires calculating  $T^{-1}$  and  $\det J_{T^{-1}}$ .
- For large  $D$ , these calculations can easily be forbiddingly expensive.

$\Rightarrow$  Need flow architectures that are both flexible (expressive) and fast.

- $T$  must be easily invertible. Most functions are not.
- Some flows are universal approximators, some aren't, and sometimes it is unknown.

# Training the flow

- Goal: Fit a flow  $p_x(x; \theta)$  to target distribution  $p_x^*(x)$   
↑  
parameters of transformation  $T(x; \theta)$   
and sometimes prior  $p(u; \theta)$
- Usually we have a collection of samples from  $p_x^*(x)$ , the training data.
- There are different measures of similarities between  $p_x$  and  $p_x^*$ .
- The most popular choice: Kullback-Leibler (KL) divergence

# Training the flow: KL divergence

- The Loss is given by

(see Lecture 4)

$$\mathcal{L}(\theta) = D_{\text{KL}} [P_x^* \parallel P_x(x; \theta)]$$

$$= - \mathbb{E}_{P_x^*} [\log p_x(x; \theta)] + \text{const}$$

make samples from training data  
Likely under flow pdf

- using  $p_x(x) = p_u(T^{-1}(x)) |\det J_{T^{-1}}(x)|$

$$\Rightarrow \mathcal{L}(\theta) = - \frac{1}{N} \sum_{n=1}^N \log p_u(T^{-1}(x_n; \theta)) + \log |\det J_{T^{-1}}(x_n; \theta)|$$

- Training: SGD for  $\nabla_{\theta} \mathcal{L}$

# Designing normalizing flows

# Designing efficient flows

- We stack many simple building blocks.

$$T = T_{\bar{K}} \circ \circ \circ T_1$$

$$z_k = T_k(z_{k-1})$$

$$\log |J_T(z)| = \sum_{k=1}^{\bar{K}} \log |J_{T_k}(z_{k-1})|$$

more depth  $\rightarrow O(k)$  growth in cost  $\ddot{\smile}$

- Note: making  $T_k$  invertible in theory and actually inverting it are very different requirements  
 $\implies$  want easily invertible functions



# Designing efficient flows

- We want a tractable Jacobian determinant.

For a general map  $D$  inputs  $\rightarrow D$  outputs  
calculating  $\det J$  has cost  $\mathcal{O}(D^3)$ , often  
intractable for large  $D$ .

Most flows are using forms for which  
 $\det J$  is  $\mathcal{O}(D)$ .

- We now discuss building blocks  $T_K$  that have  
this feature.

$$\dot{\vec{z}} = f(\vec{z})$$

# Planar flows

- Building block:

elementwise non-linearity

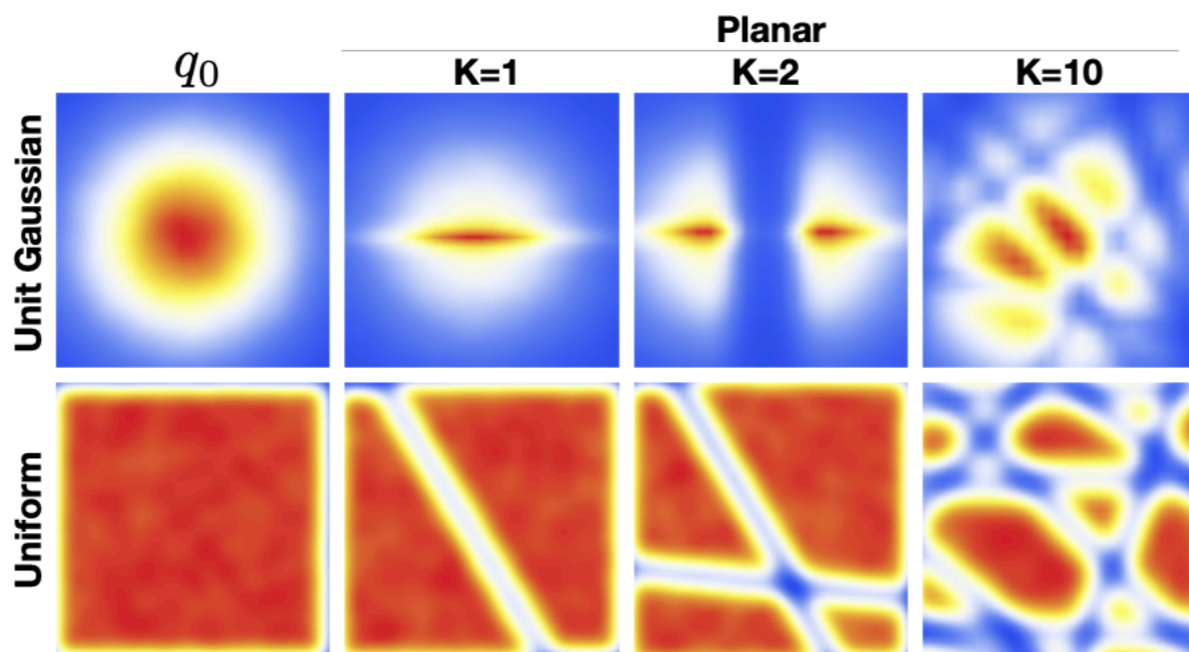
1505.05770

$$f(\vec{z}) = \vec{z} + \vec{u} h(\vec{w}^T \vec{z} + b)$$

vectors dim D

- can calculate  $\det J$  in  $\mathcal{O}(D)$  time.
- stack many of these transformations

- Example  
 $z : 2 \text{ dim.}$



Not suitable for large dimensions, since transformations are too local (i.e. affect a small volume)

# Autoregressive flows

- Autoregressive property:  $\vec{z} = (z_1, z_2, z_3 \dots z_i \dots z_D)$

variable  $z_i$  depends on  $z_{1:i}$  only

$$z'_i = \tau(z_i, \vec{h}_i)$$



„transformer“

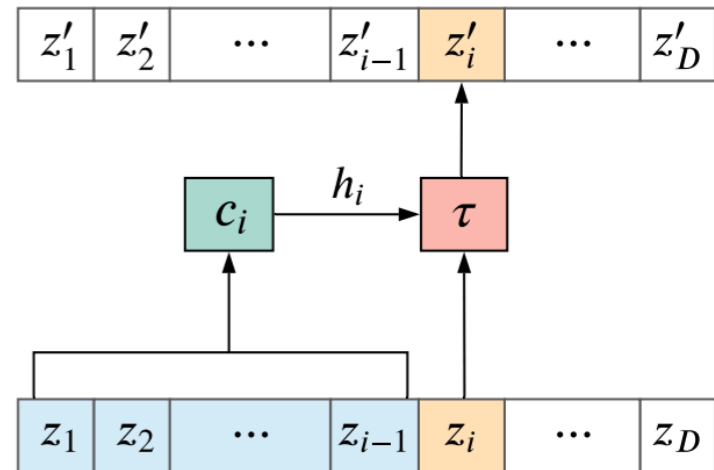
must be invertible.  
i.e. monotonic in  $z_i$

$$\vec{h}_i = c_i(\vec{z}_{<i})$$



„conditioner“

modulates the transformer



- This leads to a triangular jacobian:

$$J_T(z) = \begin{bmatrix} & & & & 0 \\ & & & & \approx \\ & & & & \approx \\ & & & & \approx \\ \approx & & & & \end{bmatrix} \Rightarrow \det J = \prod z_{ii} \quad \mathcal{O}(d)$$

# Affine Auto regressive flows

- Particularly simple Linear transformer

$$\tau(z_i) = \alpha_i z_i + \beta_i \quad \text{shift and scale } z_i$$

invertible for  $\alpha \neq 0$       "affine transformation"

or:  $\exp^{\alpha_i} z_i + \beta_i$

- The conditioner here is:

$\alpha, \beta$  depend on  $\{z_{<i}\}$  by some neural network parametrisation



NN: e.g. MLP  
or CNN

So the conditioner has the Learned parameters  $\Theta$ .

# Different autoregressive flows

- Again, autoregressive flows are

$$z'_i = T(z_i, \vec{h}_i)$$

↑  
"transformer", must be invertible.  
i.e. monotonic in  $z_i$

$$\vec{h}_i = C_i(\vec{z}_{<i})$$

↑  
"conditioner", modulates the  
transformer, not a bijection.

- Many different transformers and conditioners have been proposed. Popular:

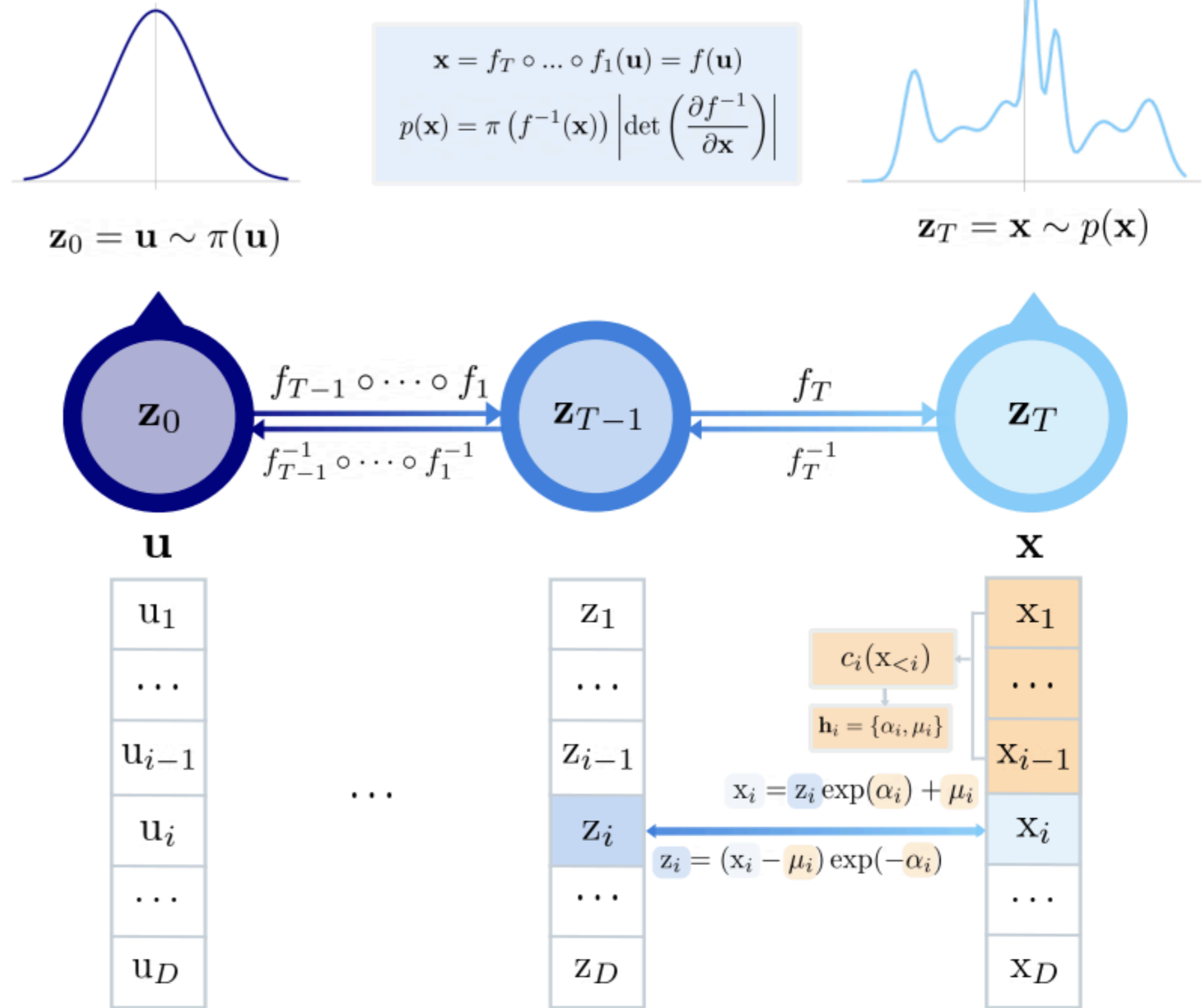
- Masked autoregressive flows (MAF)

- Inverse autoregressive flows (IAF)

They have different cost and expressivity.

# Illustration of the MAF flow

Source: 2310.03741



**Figure 2.** Diagram of how normalizing flows work, with the specific example of Masked Autoregressive Flows. The samples from the vector  $\mathbf{z}_0 = \mathbf{u}$ , sampled from the simple distribution  $\pi(\mathbf{u})$ , are deformed through the sequence of transformations  $f = f_T \circ \dots \circ f_1$  into those of  $\mathbf{z}_T = \mathbf{x}$ , which follow a more complex distribution  $p(\mathbf{x})$ . In the lower panel, we illustrate the conditioner that “masks out” the connections between  $z_i$  and  $\mathbf{h}_{<i}$ , as well as the affine functions applied to the vector components.

# Real-NVP

↙ 2d spatial array

- Baseline flow for many "image" applications  
1605.08803
- Special case of affine autoregressive flow:  
Split variables into two halves:

$$z'_{1:k} = z_{1:k}$$

$$z'_{k+1:d} = z_{k+1:d} \alpha(z_{1:k}) + \beta(z_{1:k})$$

⇒ Jacobian  $\begin{pmatrix} 1 & & | & 0 \\ & \ddots & & \\ \hline & & & \ddots \\ & & & 0 \end{pmatrix}$

Allows parallel computation of  $z'$  since all inputs are available.

- Stack many such layers with different orderings of the variables.

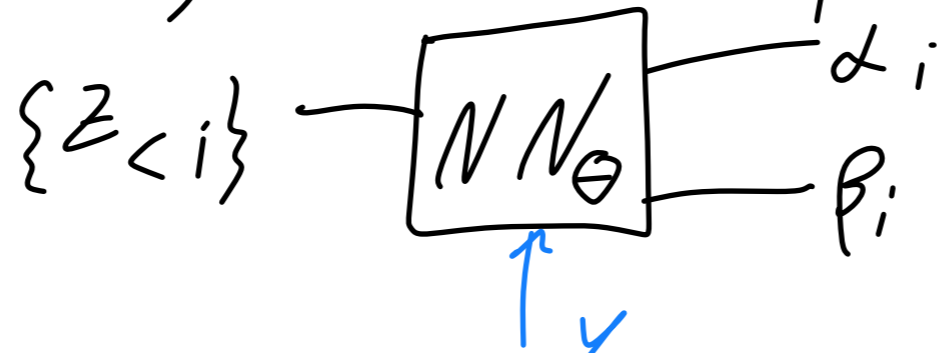
# Conditional flows

- We want to Learn not just PDFs  $p(x)$  but also conditional PDFs

$$p(x|y)$$

- This is achieved by making the flow transformations  $T$  dependent on the condition  $y$ .

- For example, in a autoregressive flow, the neural network which parametrizes  $z$  now also gets an input of  $y$ :



- Training works the same as before (KL-div.).

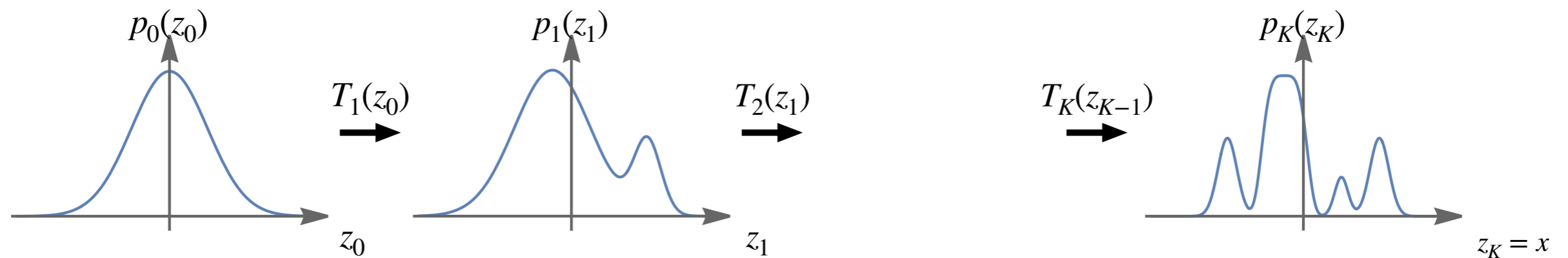


# **Research Example: Normalizing flows to model the matter distribution in cosmology**

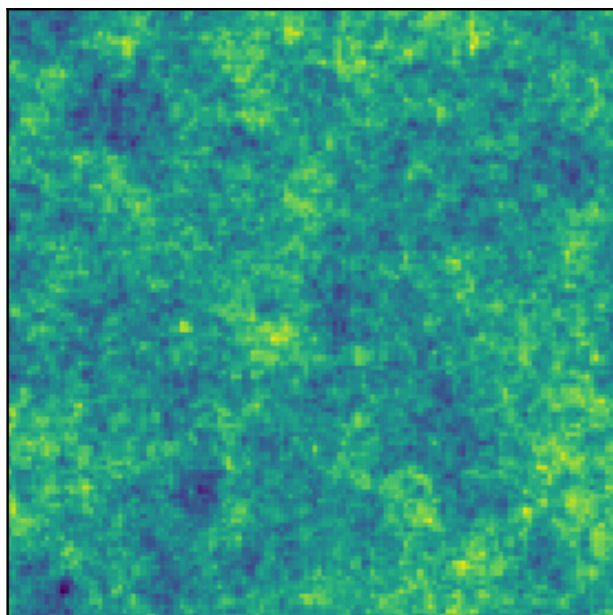
**(work from my group)**

# NFs vs structure formation

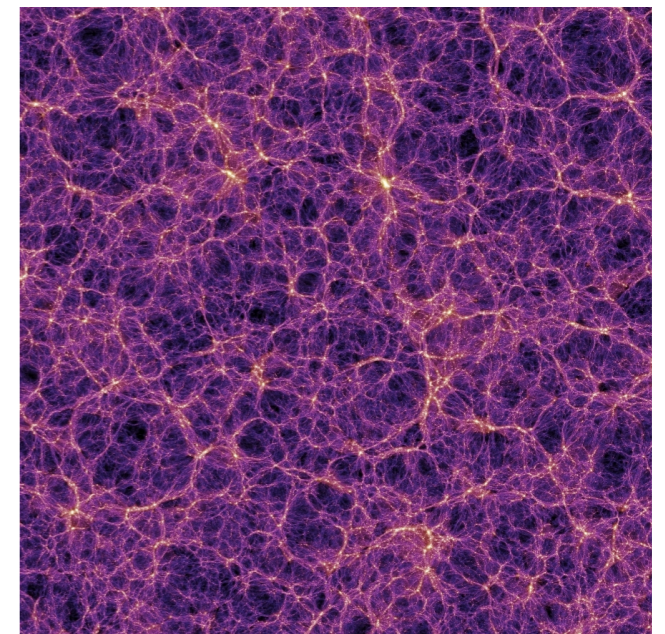
Gaussian initial conditions PDF morphs into complicated late-time matter distribution.



**Gaussian primordial matter perturbations**

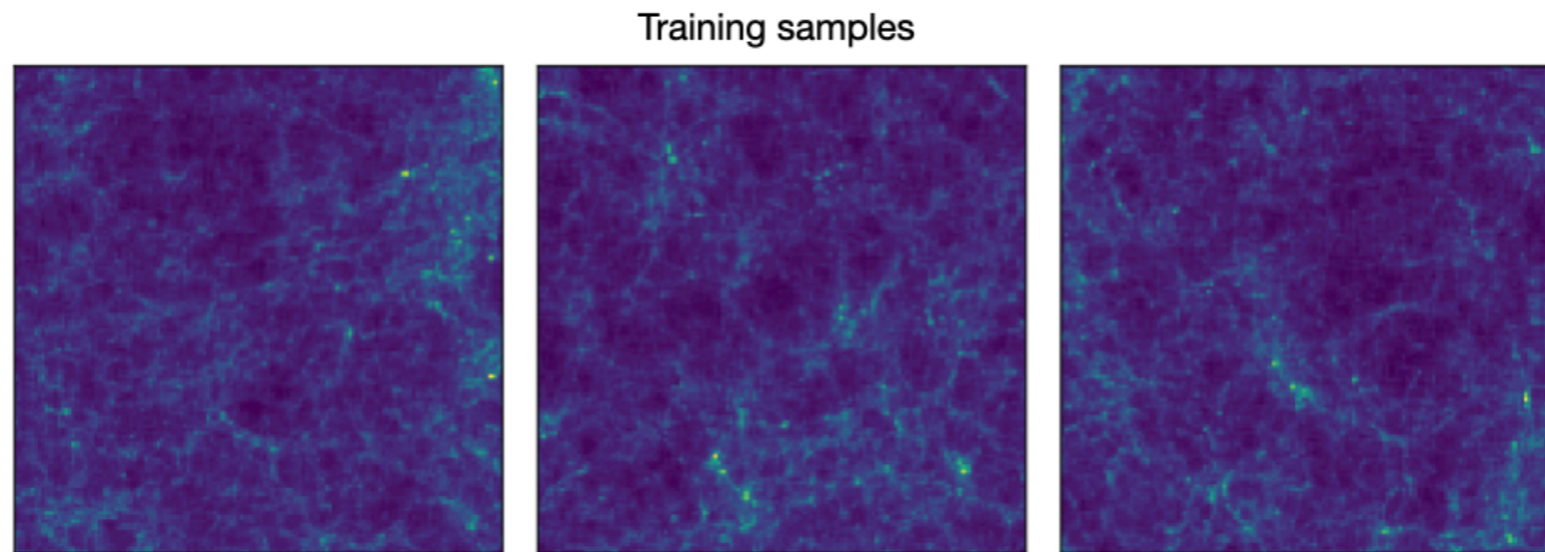
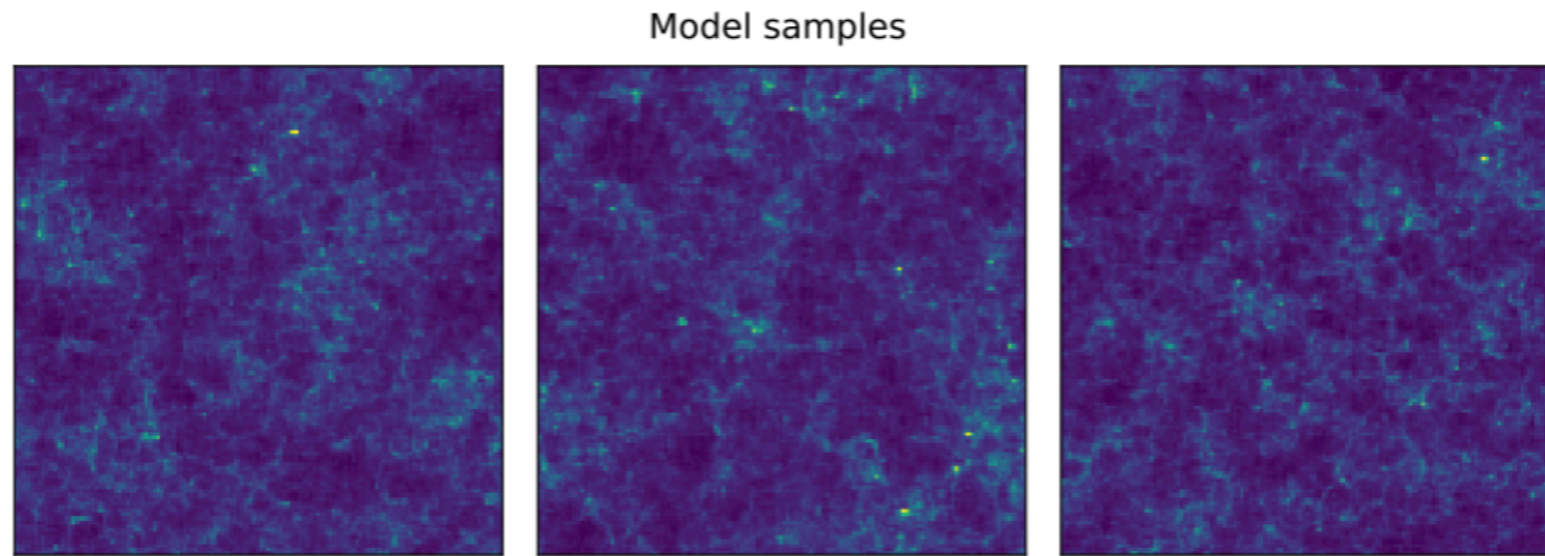
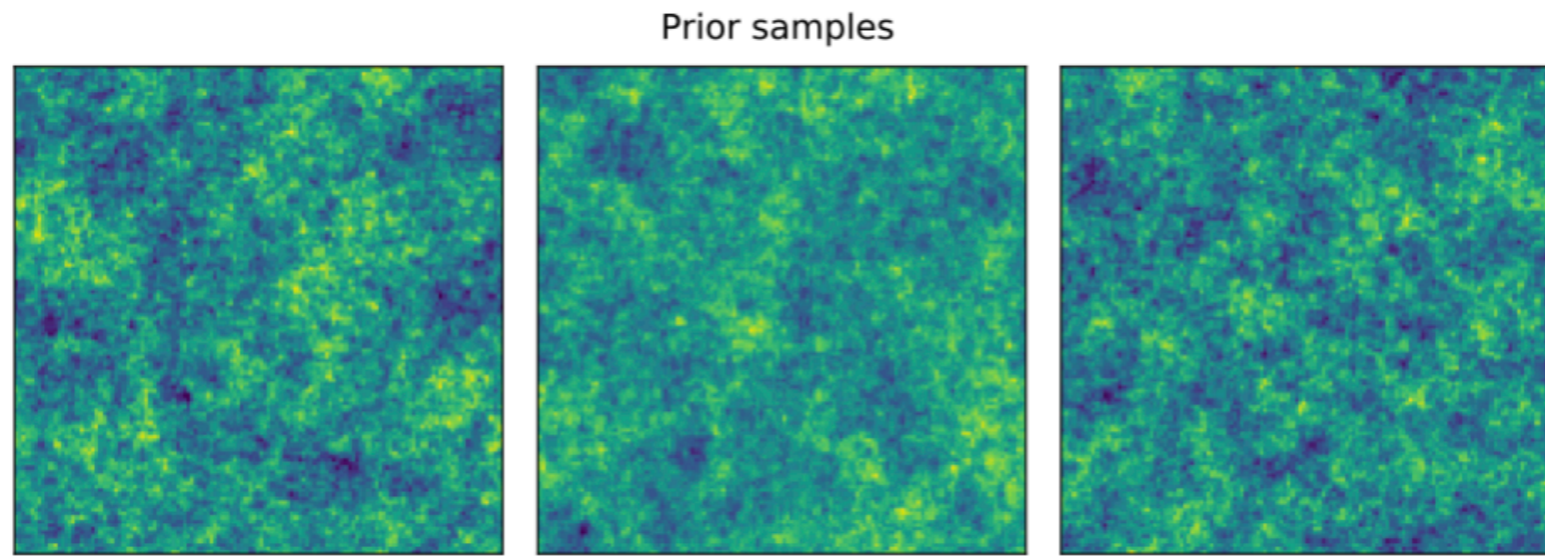


**Non-gaussian matter/galaxy distribution today**



**Cosmological time evolution**

# Flowing from a correlated Gaussian to today's matter distribution

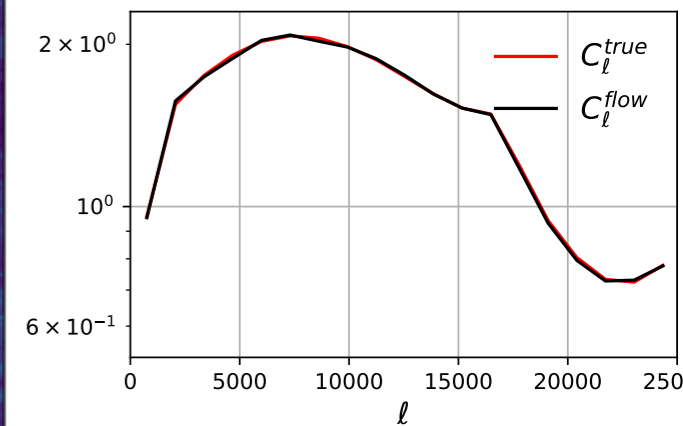


Flow:  
RealNVP

Density peaks  
match, as in  
physical structure  
formation.

The flow learned  
to deform a  
Gaussian PDF into  
a highly non-  
Gaussian PDF

Power spectrum matches



$125 h^{-1}$  Mpc

About  $128^3$   
particles  
per field

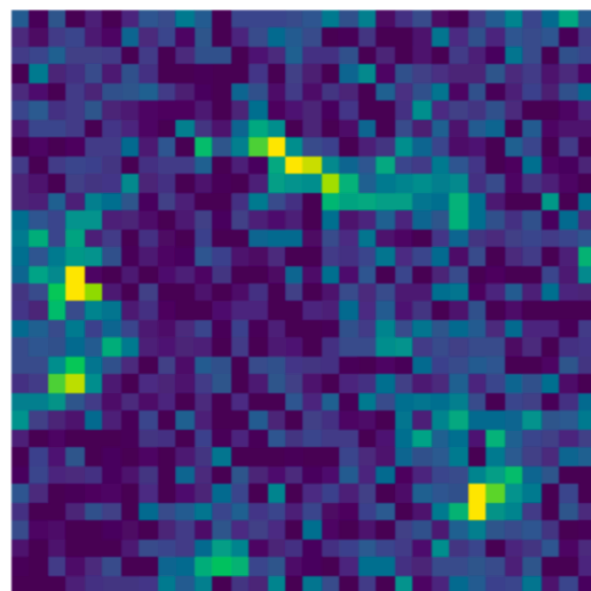
# De-noising with a Generative Prior

In data analysis in cosmology we often make use of **Gaussian priors (Wiener Filter)**. This is no longer justified for very high resolution observations. Using the trained normalizing flow we can **include non-Gaussian priors**:

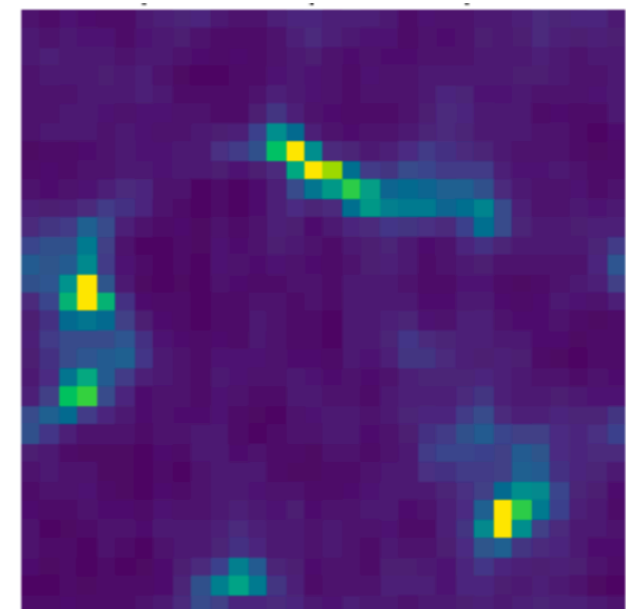
$$\ln p(y | d) = -\frac{1}{2}(y - d)^T N^{-1}(y - d) - \ln p_{\text{flow}}(y)$$

True matter field      Noisy observation

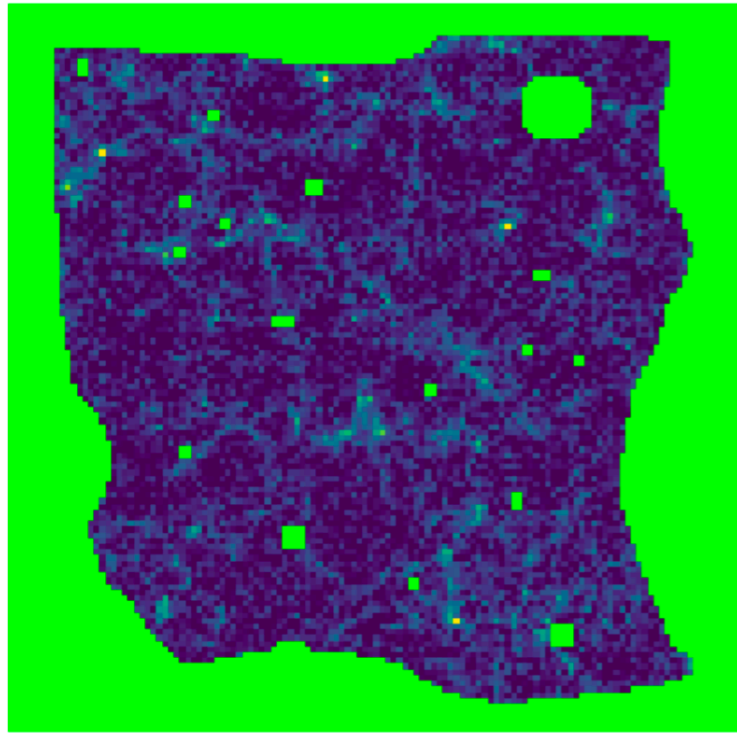
We use a flow trained on simulations of the matter distribution. Then we use this knowledge of the matter PDF to **de-noise an observation of the matter field by maximizing the posterior**.



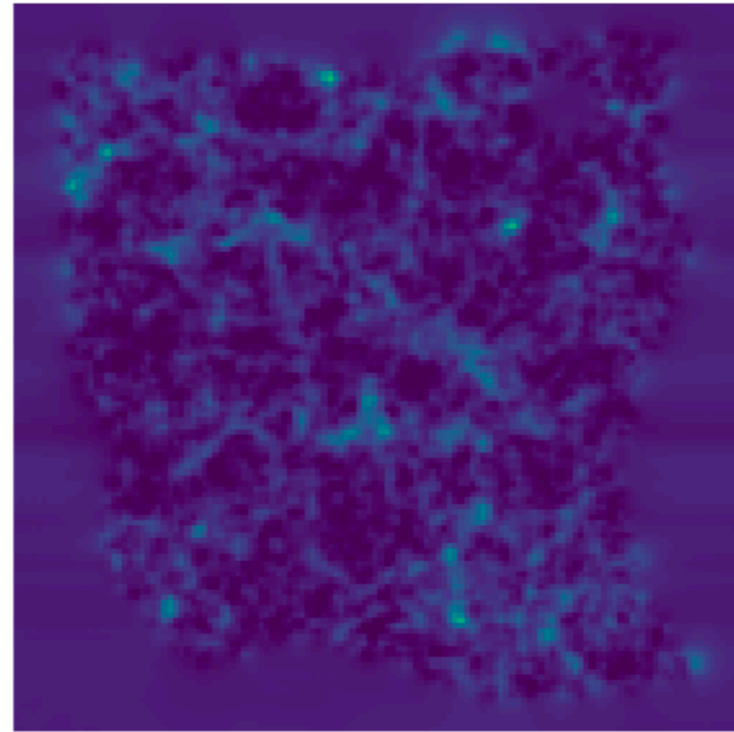
Flow based de-noising



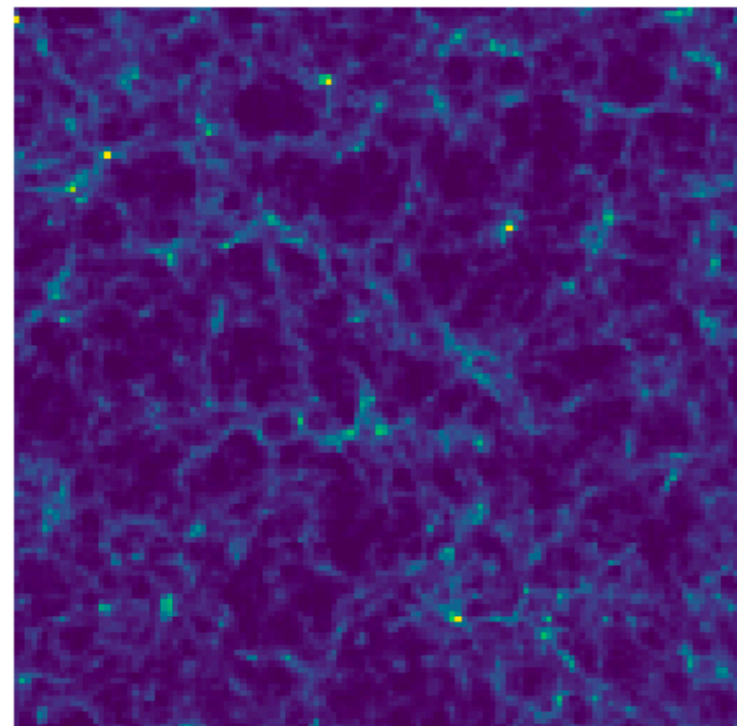
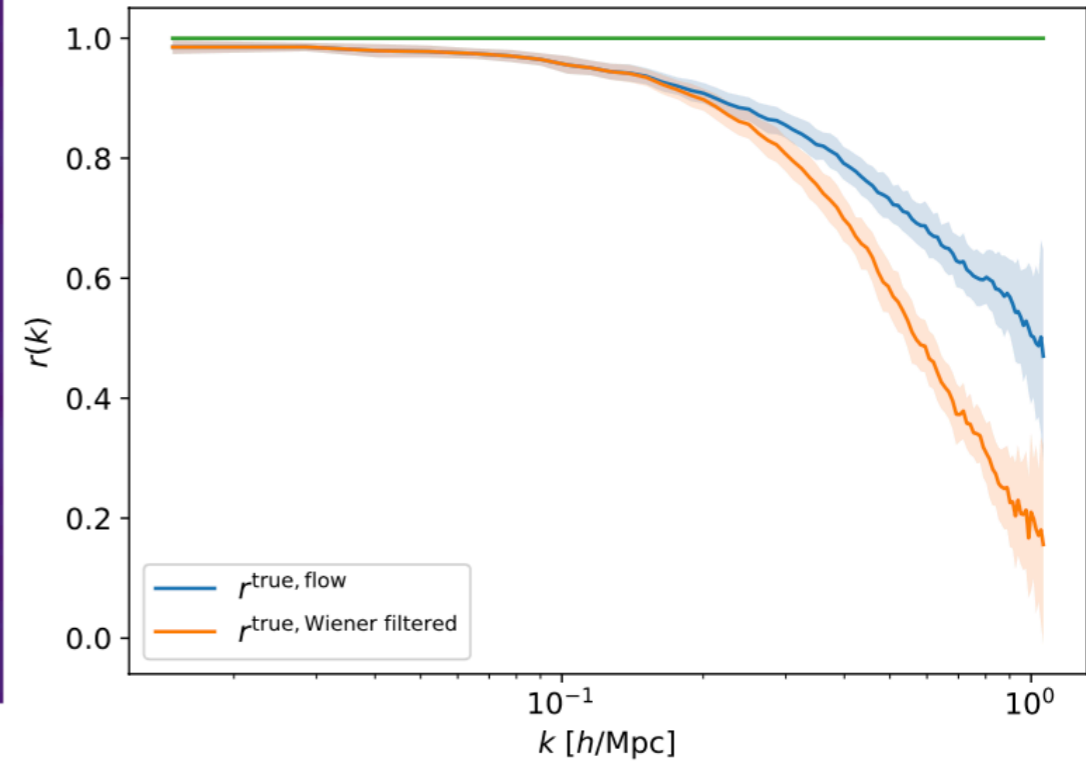
# De-noising the observed matter field



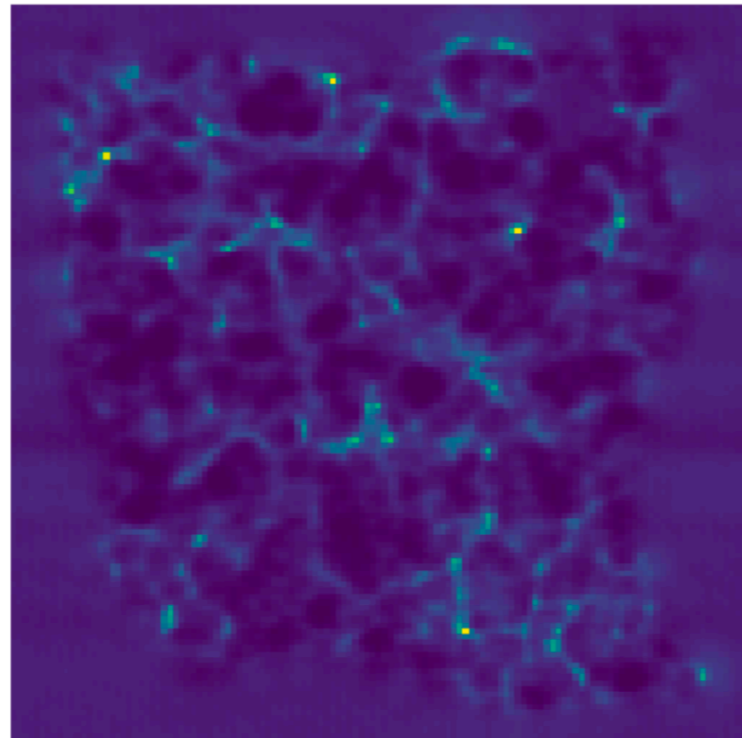
Observed (noisy, masked)



Wiener filtered



Truth



Flow MAP

As expected, the NF lowers the reconstruction noise on non-linear scales compared to the Wiener filter.

Generative de-noising is useful in many other domains.

Rouhiainen, MM: [arXiv:2211.15161](https://arxiv.org/abs/2211.15161) De-noising non-Gaussian fields in cosmology with normalizing flows