

Machine learning for physicists

<https://github.com/wangleiphy/ml4p>

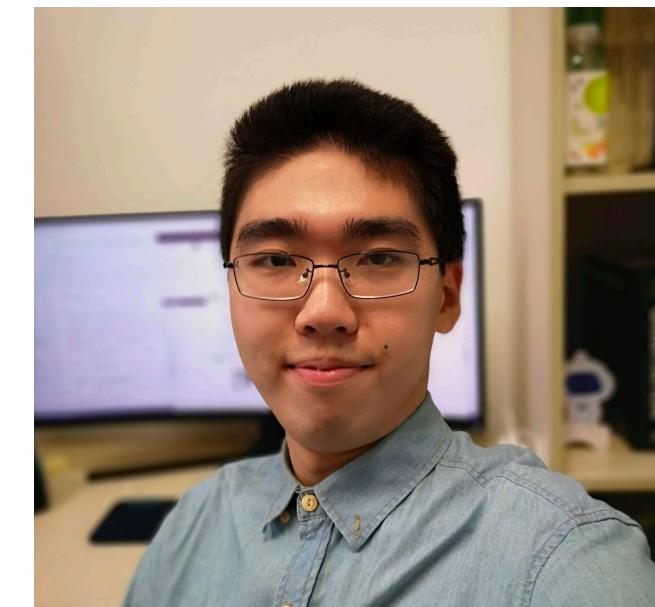
每周四上午10点

课程微信群

2.23	Overview
3.2	Machine learning practices
3.9	A hitchhiker's guide to deep learning
3.16	Research projects hands-on
3.23	Symmetries in machine learning
3.30	Differentiable programming
4.6	Generative models-I
4.13	Generative models-II
4.20	Research projects presentation
4.27	AI for science: why now ?



助教



李子航



李扬帆

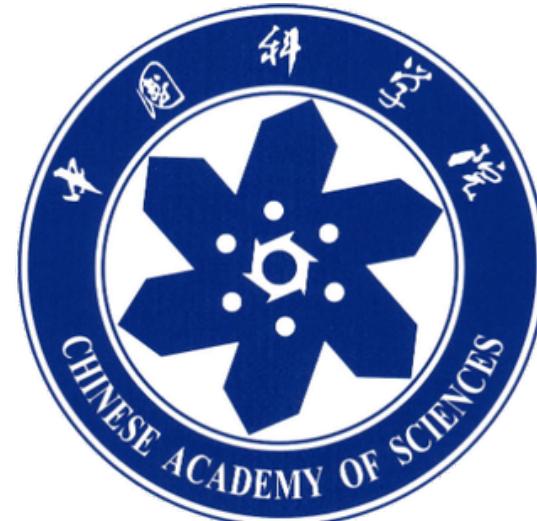
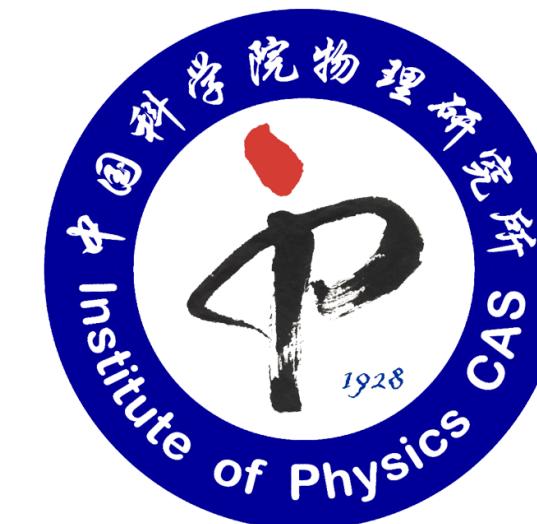
考核方式: project + presentation (1学分)

Generative AI for Science

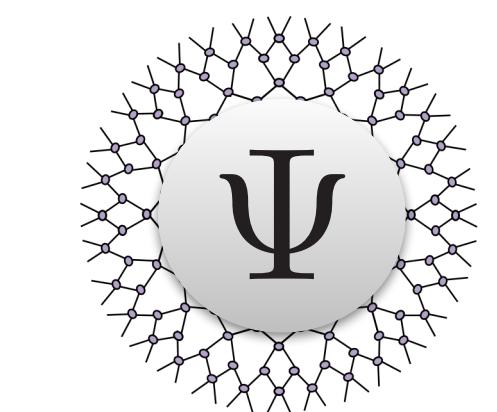
Lei Wang (王磊)

Institute of Physics, CAS

<https://wangleiphy.github.io>

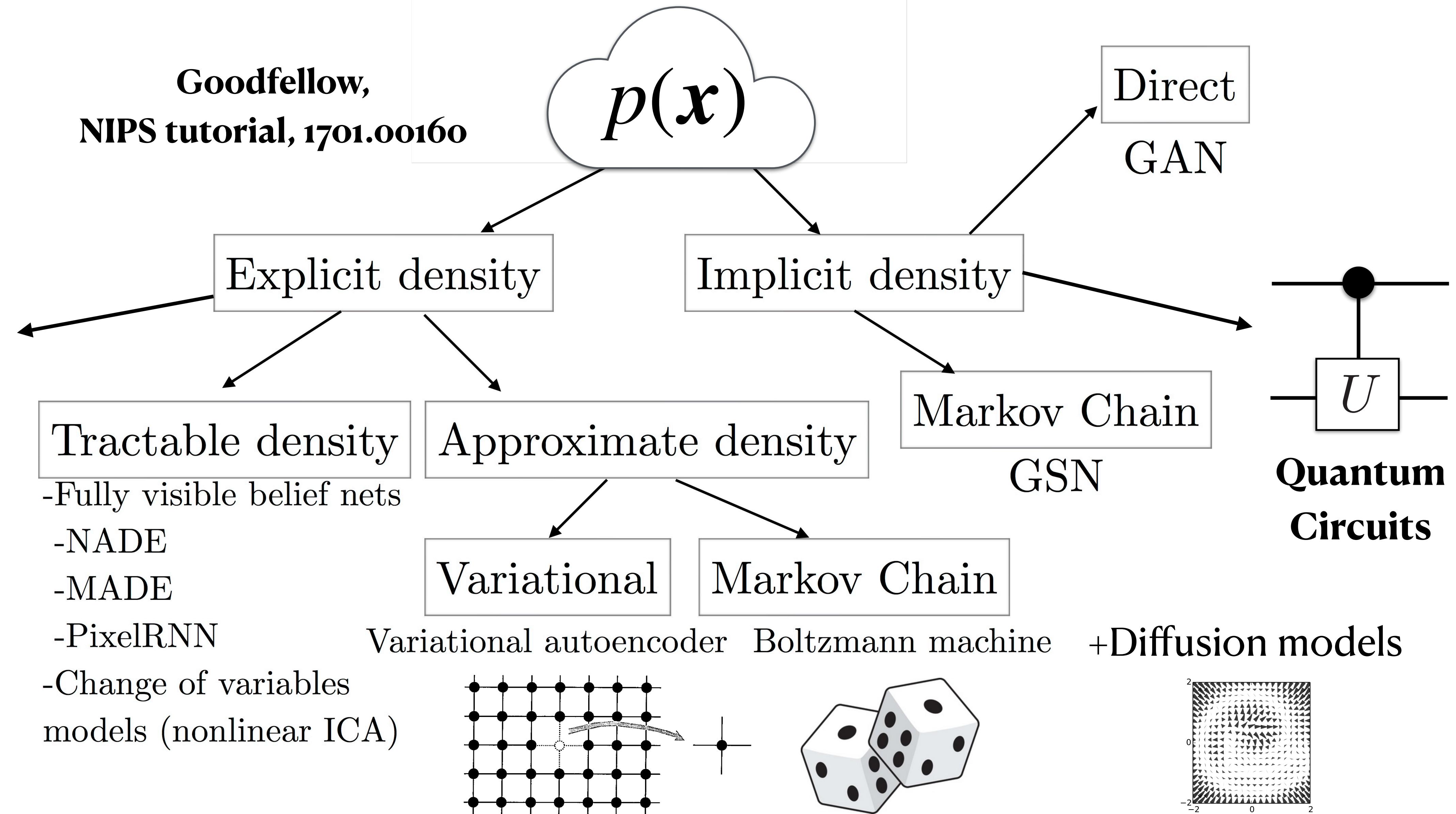


Generative models and their physics genes



**Tensor
Networks**

**Goodfellow,
NIPS tutorial, 1701.00160**



Variational autoencoders

Kingma, Welling, 1312.6114

Close connection to the variational calculus we have learned

$$p(\mathbf{x}) = \frac{e^{-\beta E(\mathbf{x})}}{Z}$$

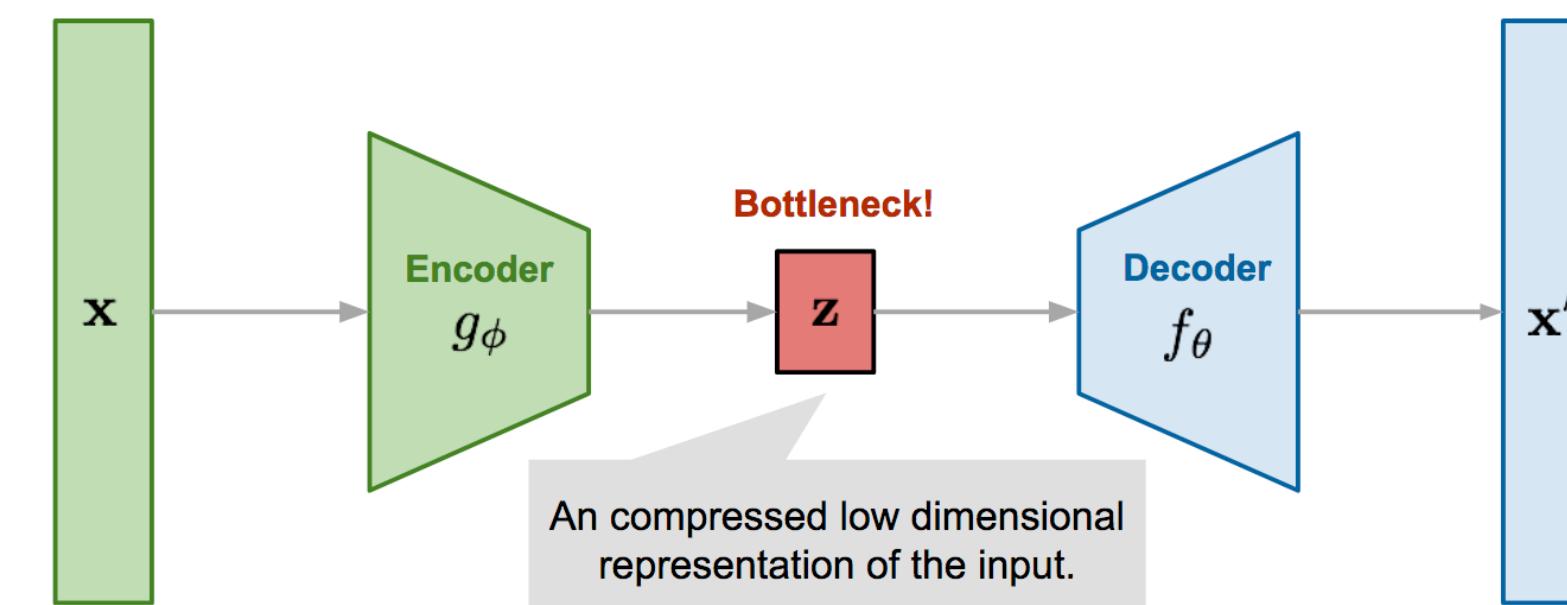
Variational free energy

$$\int d\mathbf{x} q(\mathbf{x}) [\ln q(\mathbf{x}) + \beta E(\mathbf{x})] \geq -\ln Z$$

$$p(z | \mathbf{x}) = \frac{p(\mathbf{x}, z)}{p(\mathbf{x})}$$

Variational Bayes/Variational inference

$$\int dz q(z | \mathbf{x}) [\ln q(z | \mathbf{x}) - \ln p(\mathbf{x}, z)] \geq -\ln p(\mathbf{x})$$



For each data we introduce

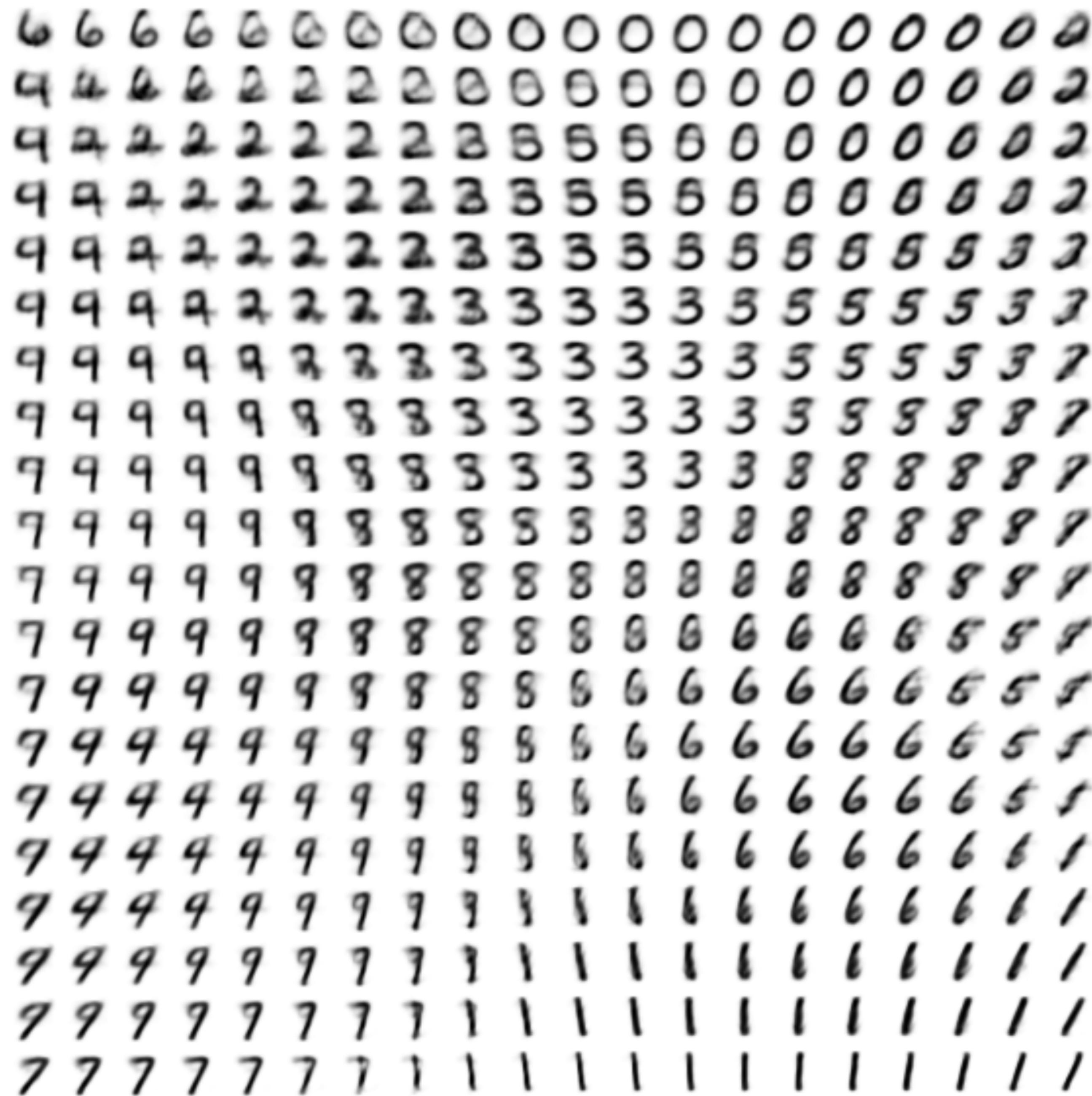
$$\mathcal{L}(\mathbf{x}) = \langle -\ln p(\mathbf{x}, \mathbf{z}) + \ln q(\mathbf{z}|\mathbf{x}) \rangle_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})}, \quad (53)$$

which is a variational upper bound of $-\ln p(\mathbf{x})$ since $\mathcal{L}(\mathbf{x}) + \ln p(\mathbf{x}) = \text{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) \geq 0$. We see that $q(\mathbf{z}|\mathbf{x})$ provides a variational approximation of the posterior $p(\mathbf{z}|\mathbf{x})$. By minimizing \mathcal{L} one effectively pushes the two distributions together. And the variational free energy becomes exact only when $q(\mathbf{z}|\mathbf{x})$ matches to $p(\mathbf{z}|\mathbf{x})$. In fact, $-\mathcal{L}$ is called evidence lower bound (ELBO) in variational inference.

We can obtain an alternative form of the variational free energy

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = -\langle \ln p_{\theta}(\mathbf{x}|\mathbf{z}) \rangle_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} + \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})). \quad (54)$$

The first term of Eq. (54) is the reconstruction negative log-likelihood, while the second term is the KL divergence between the approximate posterior distribution and the latent prior. We also be explicit about the network parameters θ, ϕ of the encoder and decoder.

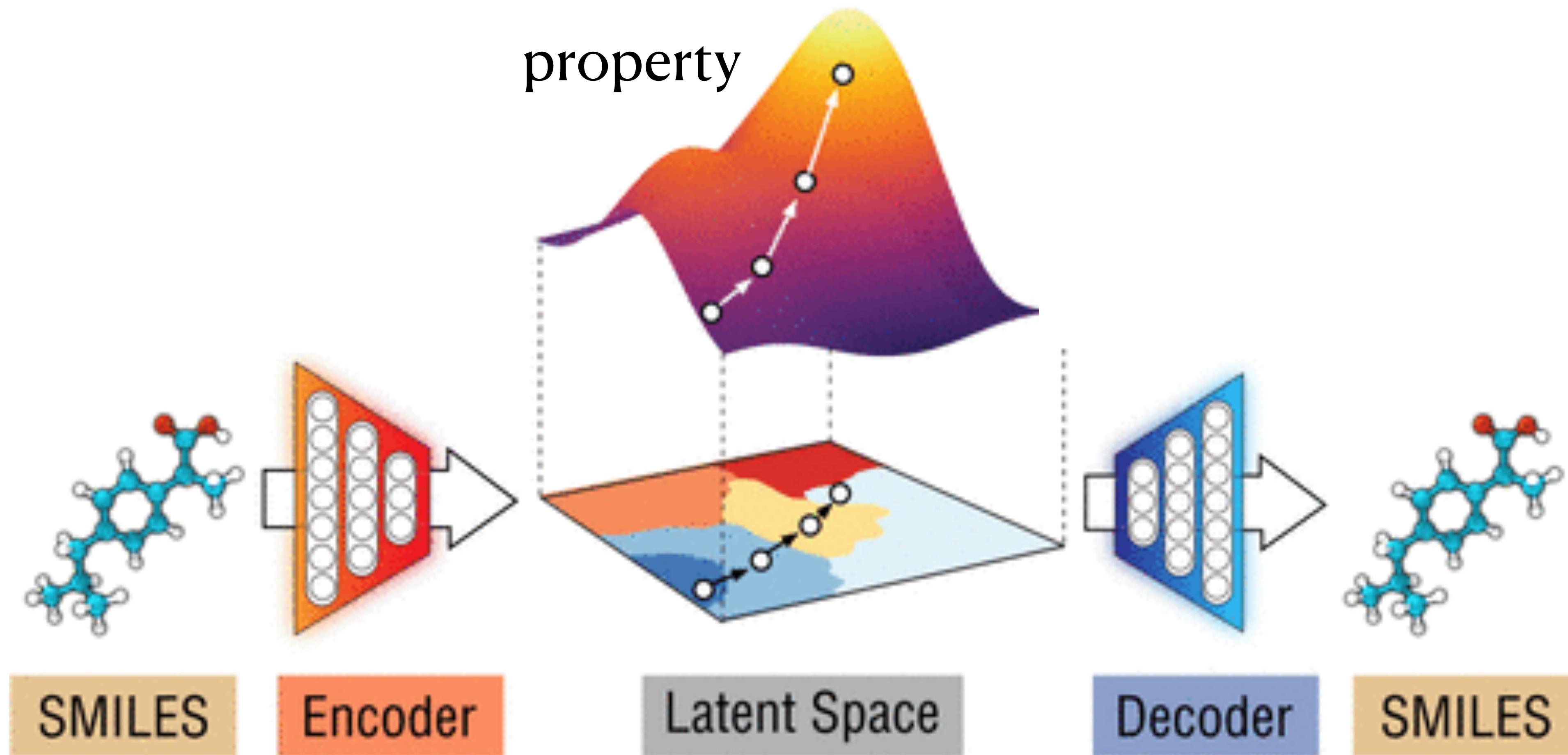


Learned MNIST latent space

Kingma, Welling, 1312.6114

Chemical design using continuous latent variables

Gomez-Bombarelli et al, 1610.02415



Normalizing flows



Parallel WaveNet 1711.10433

<https://deepmind.com/blog/high-fidelity-speech-synthesis-wavenet/>

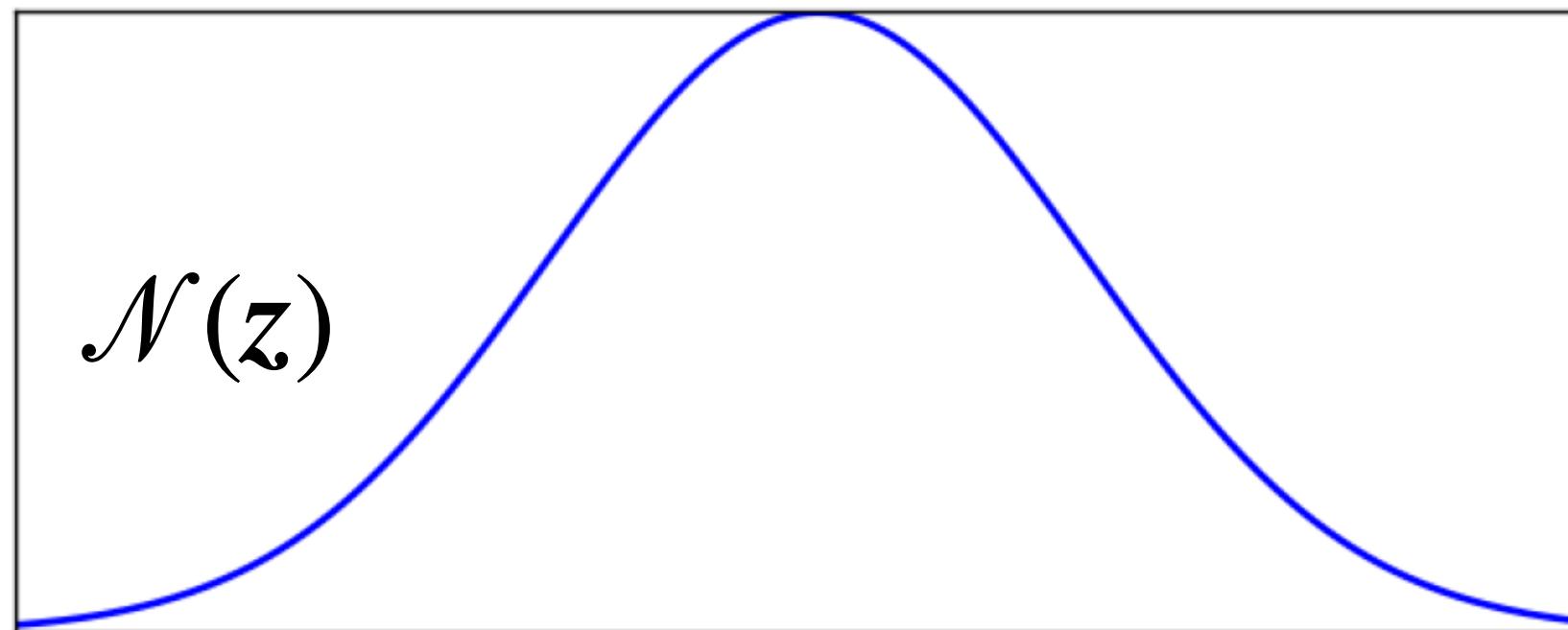


Glow 1807.03039

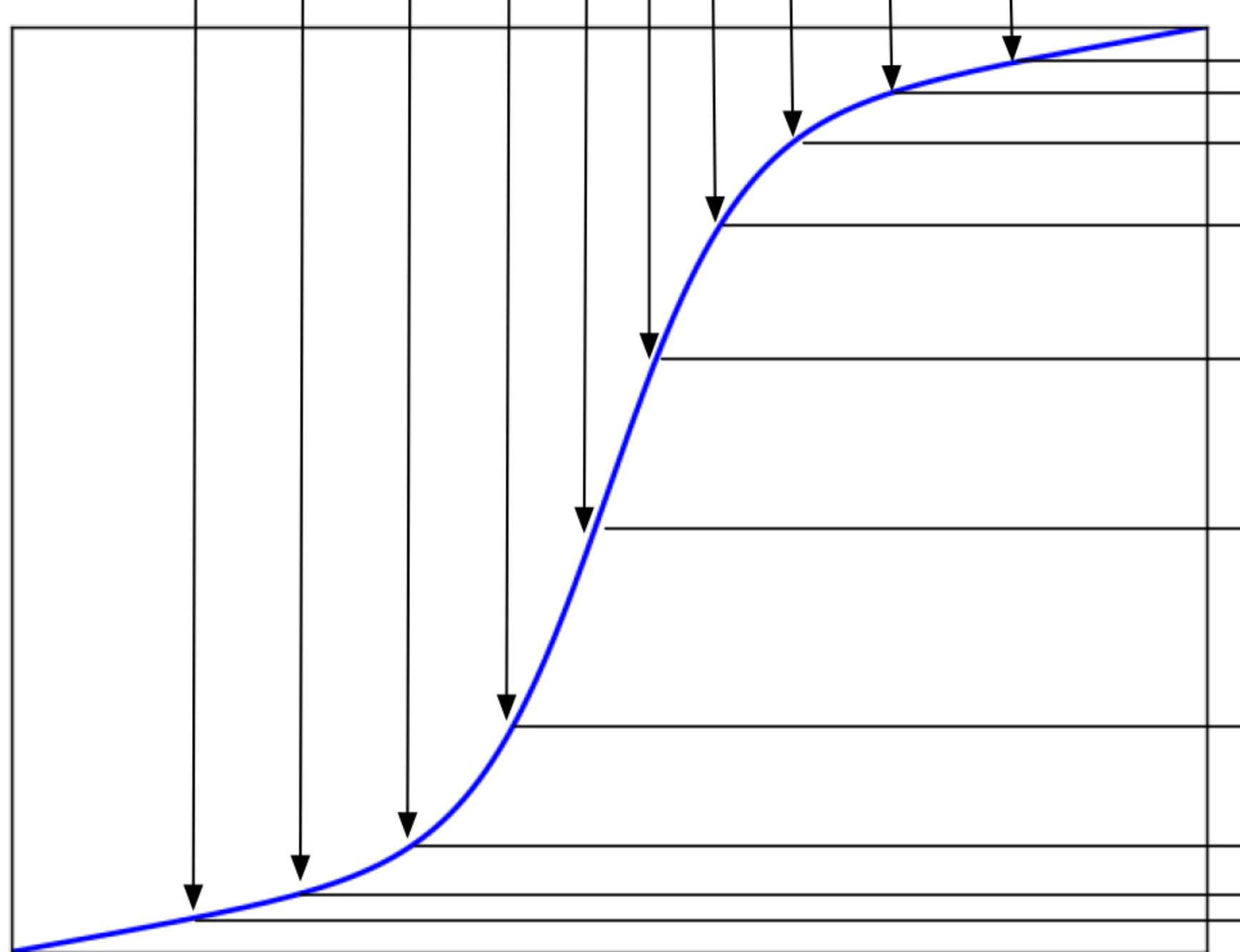
<https://blog.openai.com/glow/>

Normalizing flow in a nutshell

Base
density

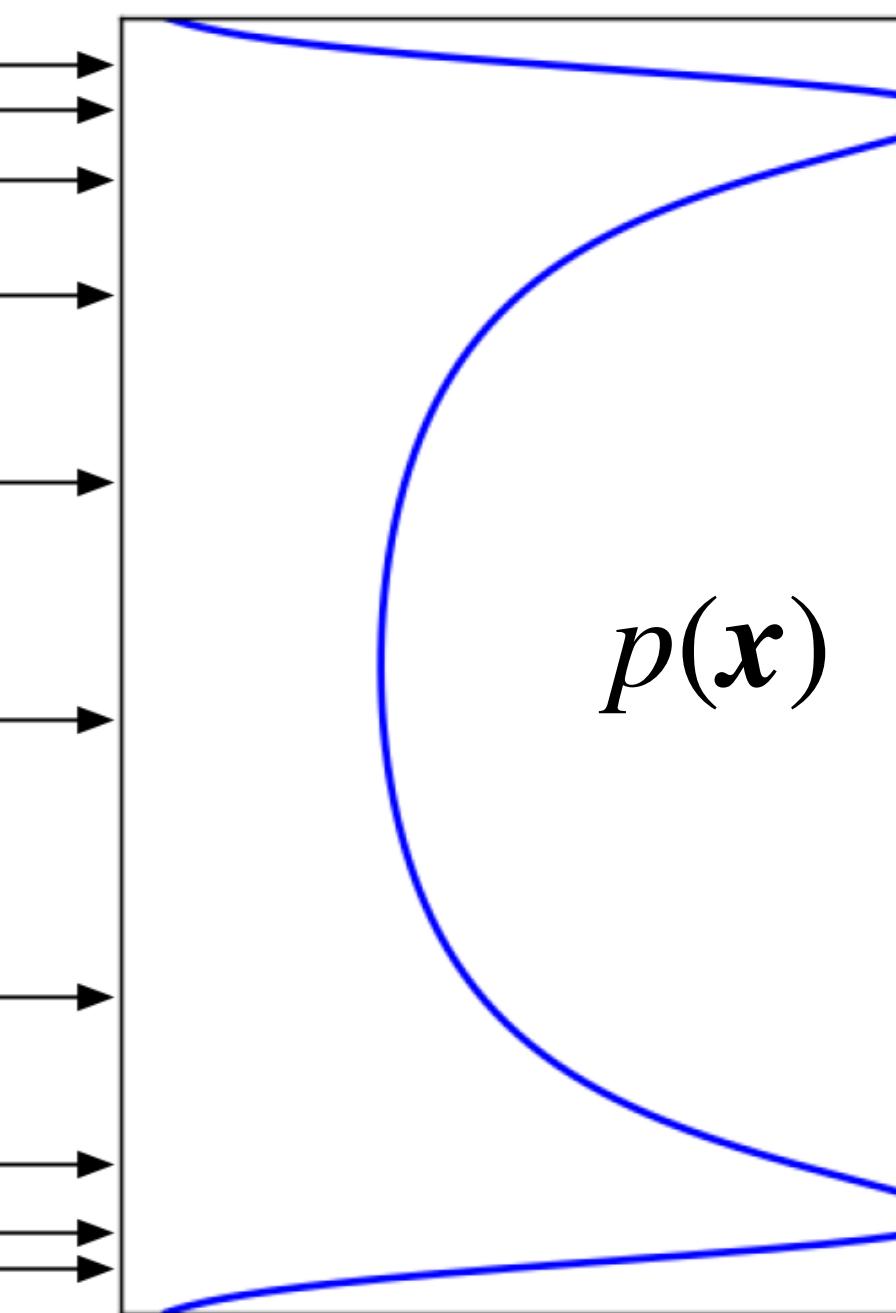


“neural net”
with 1 neuron



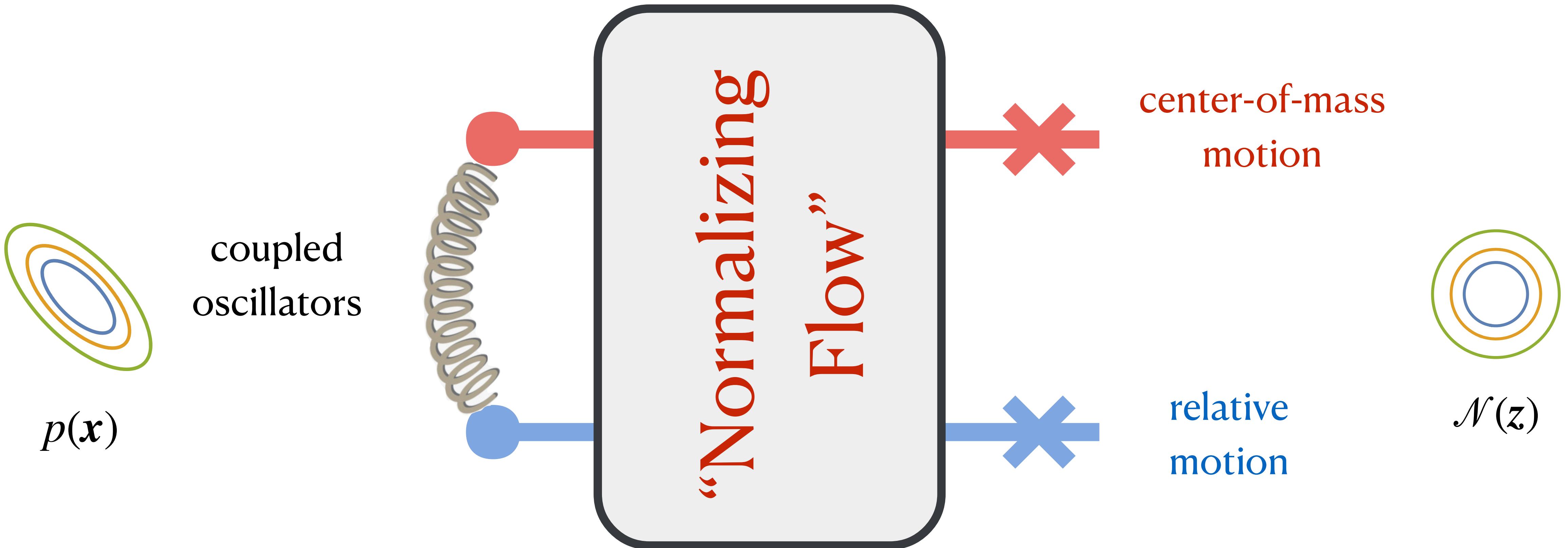
$$p(x) = \mathcal{N}(z) \left| \det \left(\frac{\partial z}{\partial x} \right) \right|$$

Review article 1912.02762



Target
density

Normalizing flow for physics: an intuition



High-dimensional, composable, learnable, nonlinear transformations

Example of a building block

Forward

$$\begin{cases} \mathbf{x}_< = \mathbf{z}_< \\ \mathbf{x}_> = \mathbf{z}_> \odot e^{s(\mathbf{z}_<)} + t(\mathbf{z}_<) \end{cases}$$

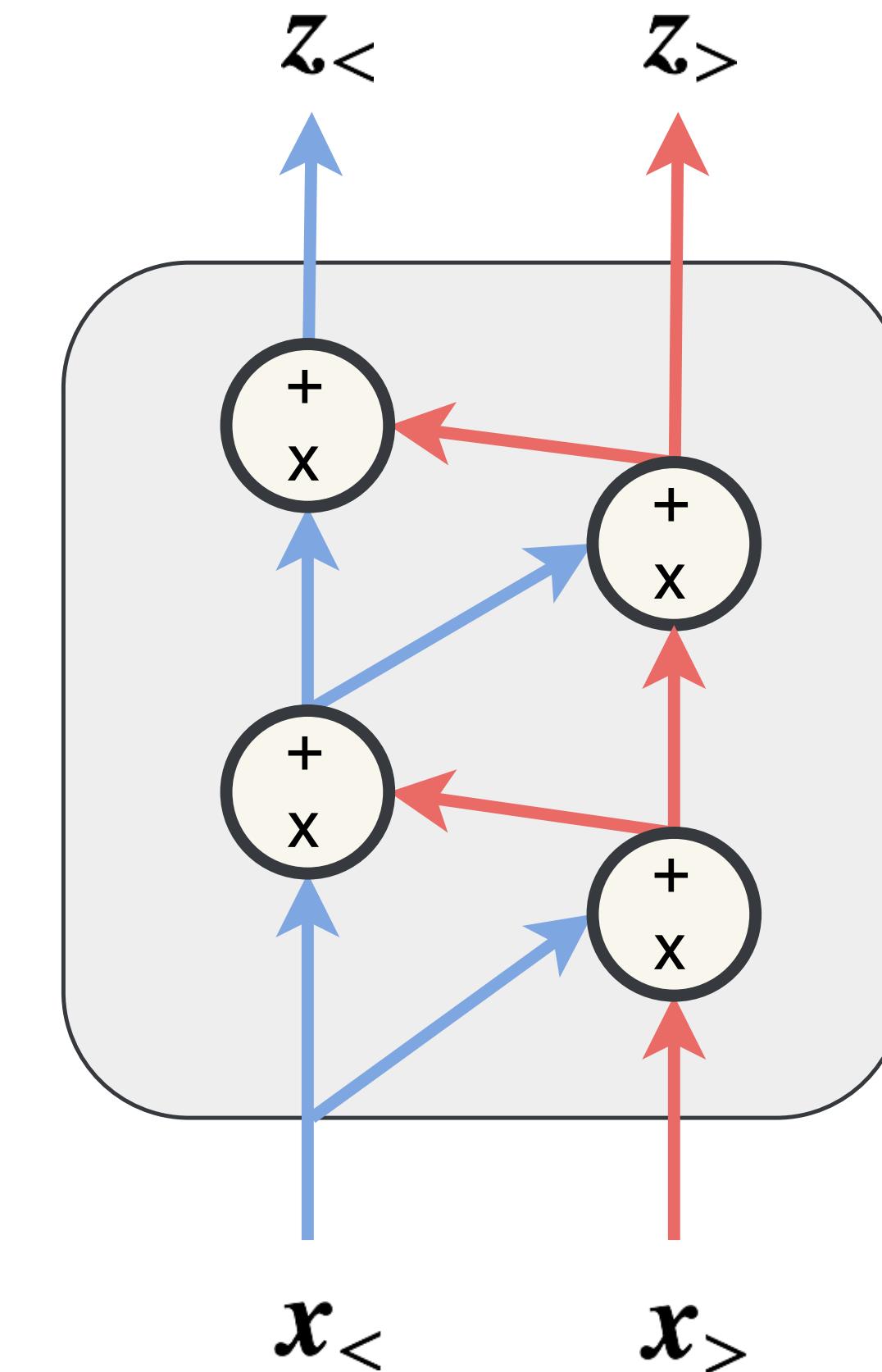
arbitrary
neural nets

Inverse

$$\begin{cases} \mathbf{z}_< = \mathbf{x}_< \\ \mathbf{z}_> = (\mathbf{x}_> - t(\mathbf{x}_<)) \odot e^{-s(\mathbf{x}_<)} \end{cases}$$

Log-Abs-Jacobian-Det

$$\ln \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right| = \sum_i [s(\mathbf{z}_<)]_i$$

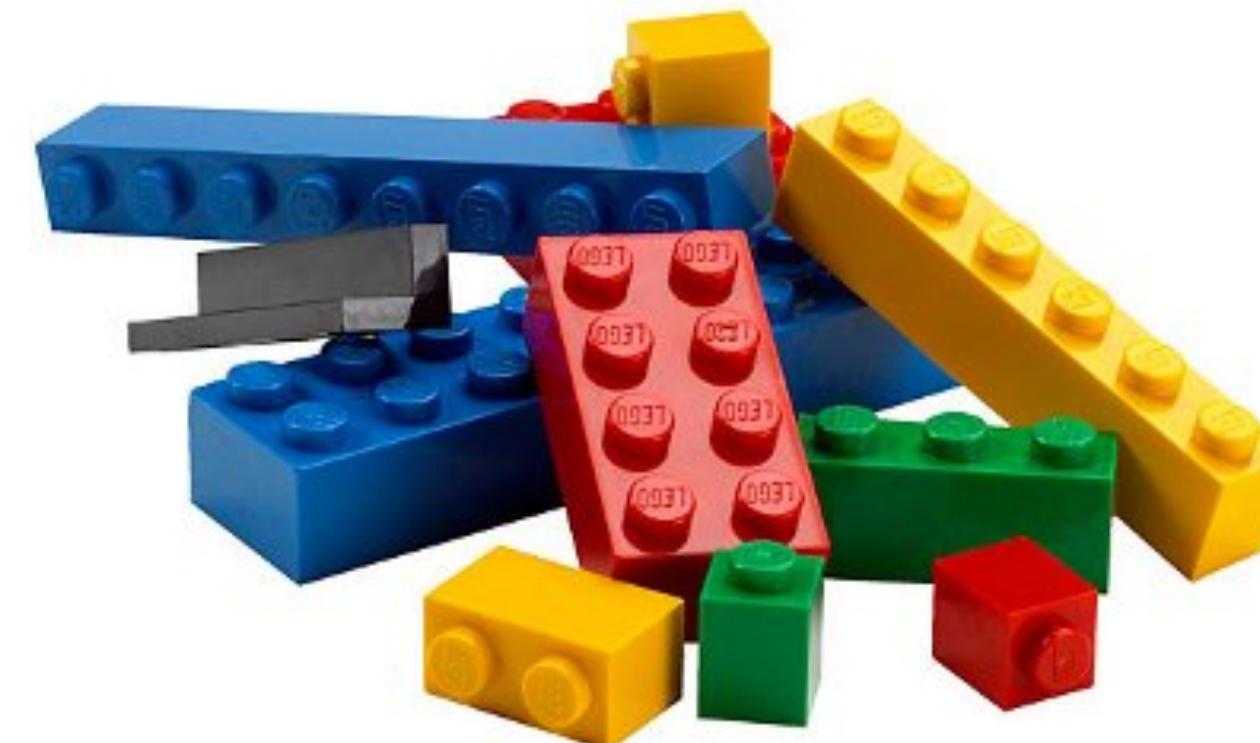


Real NVP, Dinh et al, 1605.08803

Turns out to have surprising connection Störmer–Verlet integration (later)

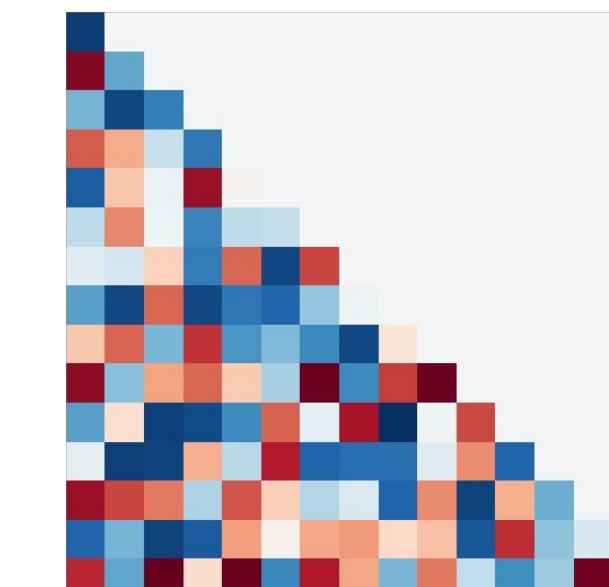
Flow architecture design

Composability

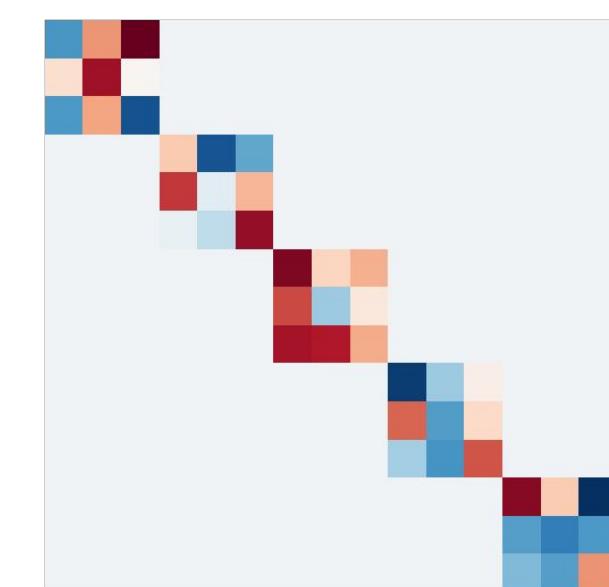


Balanced
efficiency &
inductive bias

$$\left| \det \left(\frac{\partial z}{\partial x} \right) \right|$$



Autoregressive



Neural RG

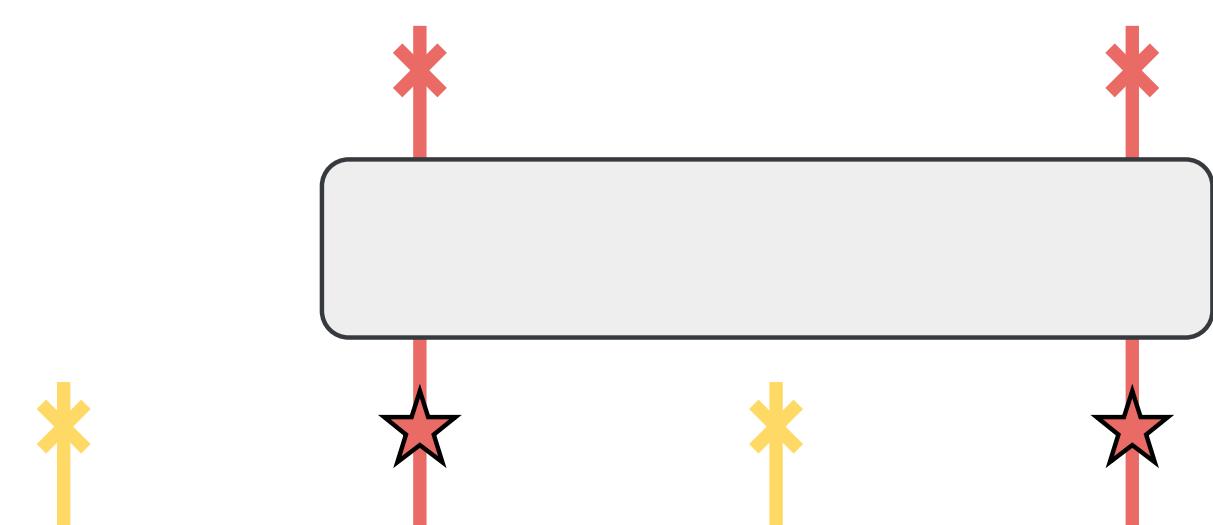
$$z = \mathcal{T}(x)$$
$$\mathcal{T} = \mathcal{T}_1 \circ \mathcal{T}_2 \circ \mathcal{T}_3 \circ \dots$$

$$\frac{\partial p(x, t)}{\partial t} + \nabla \cdot [p(x, t)v] = 0$$

Continuous flow

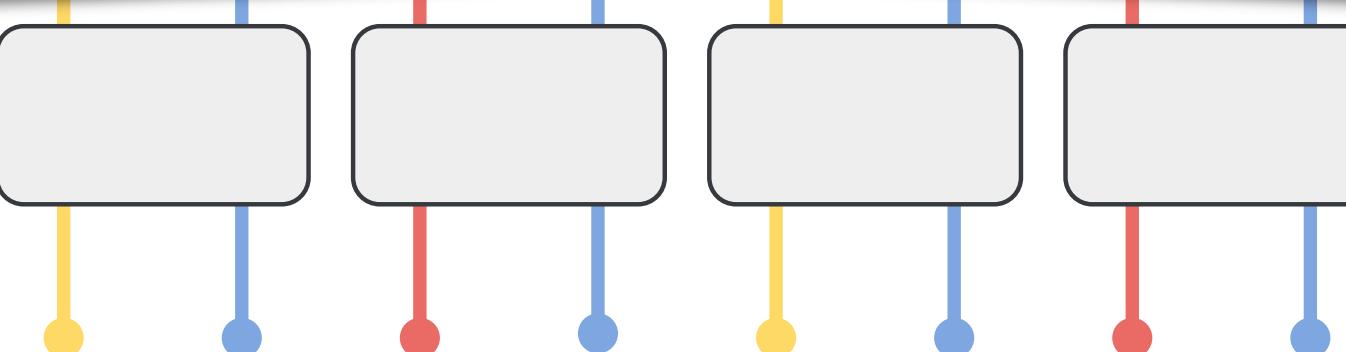
Neural network renormalization group

Collective variables



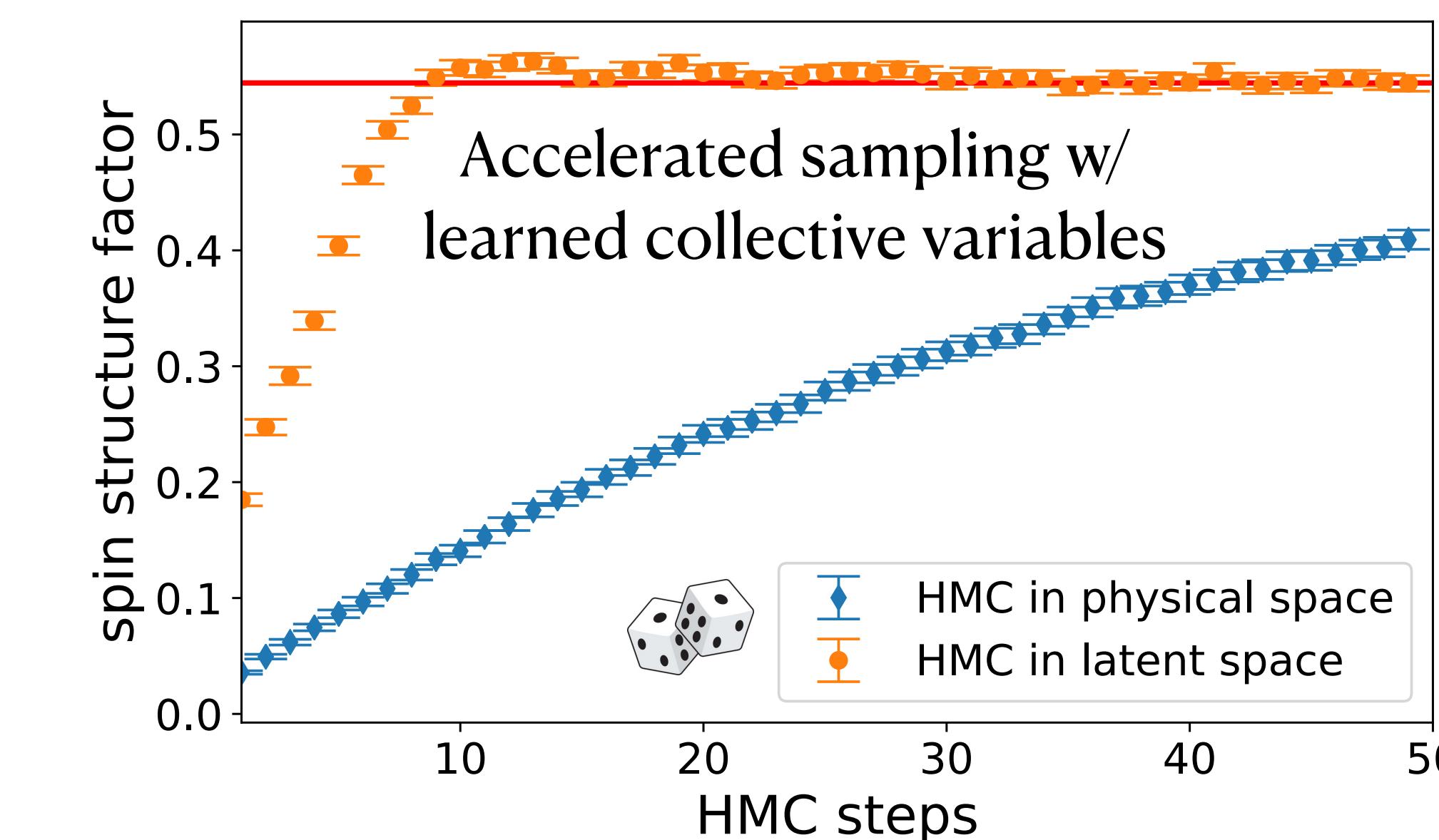
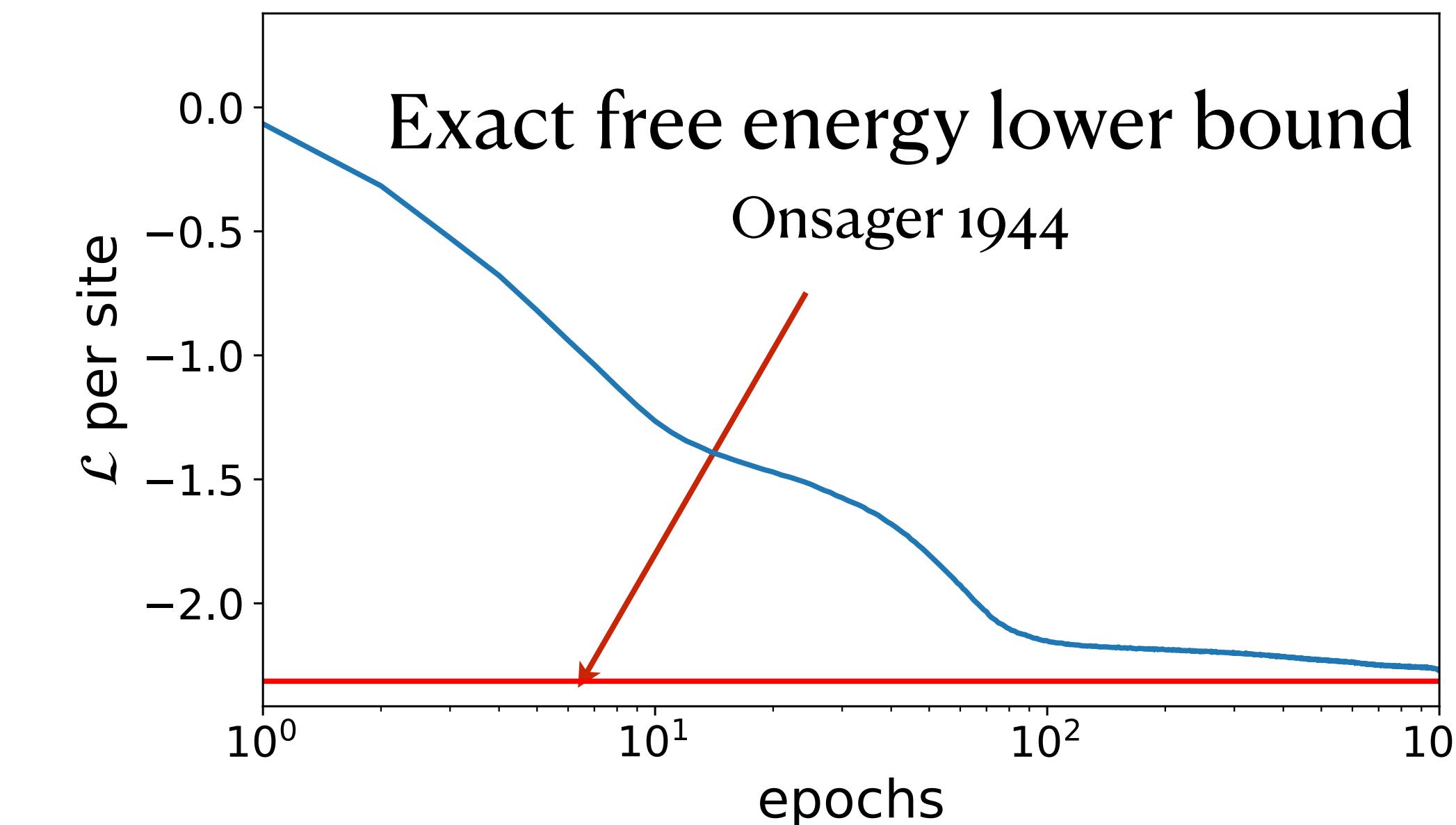
Probability Transformation

$$\ln p(\mathbf{x}) = \ln \mathcal{N}(\mathbf{z}) - \ln \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right|$$

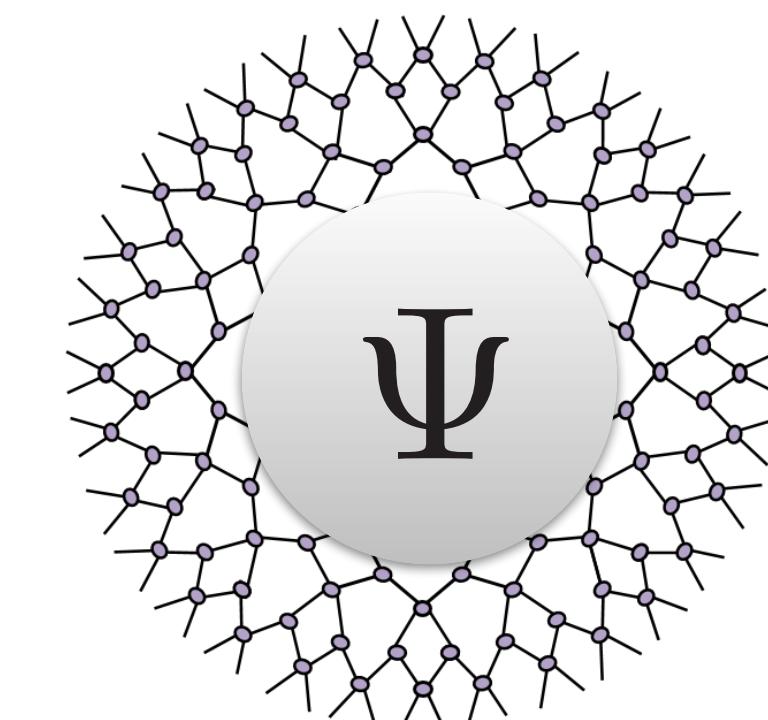
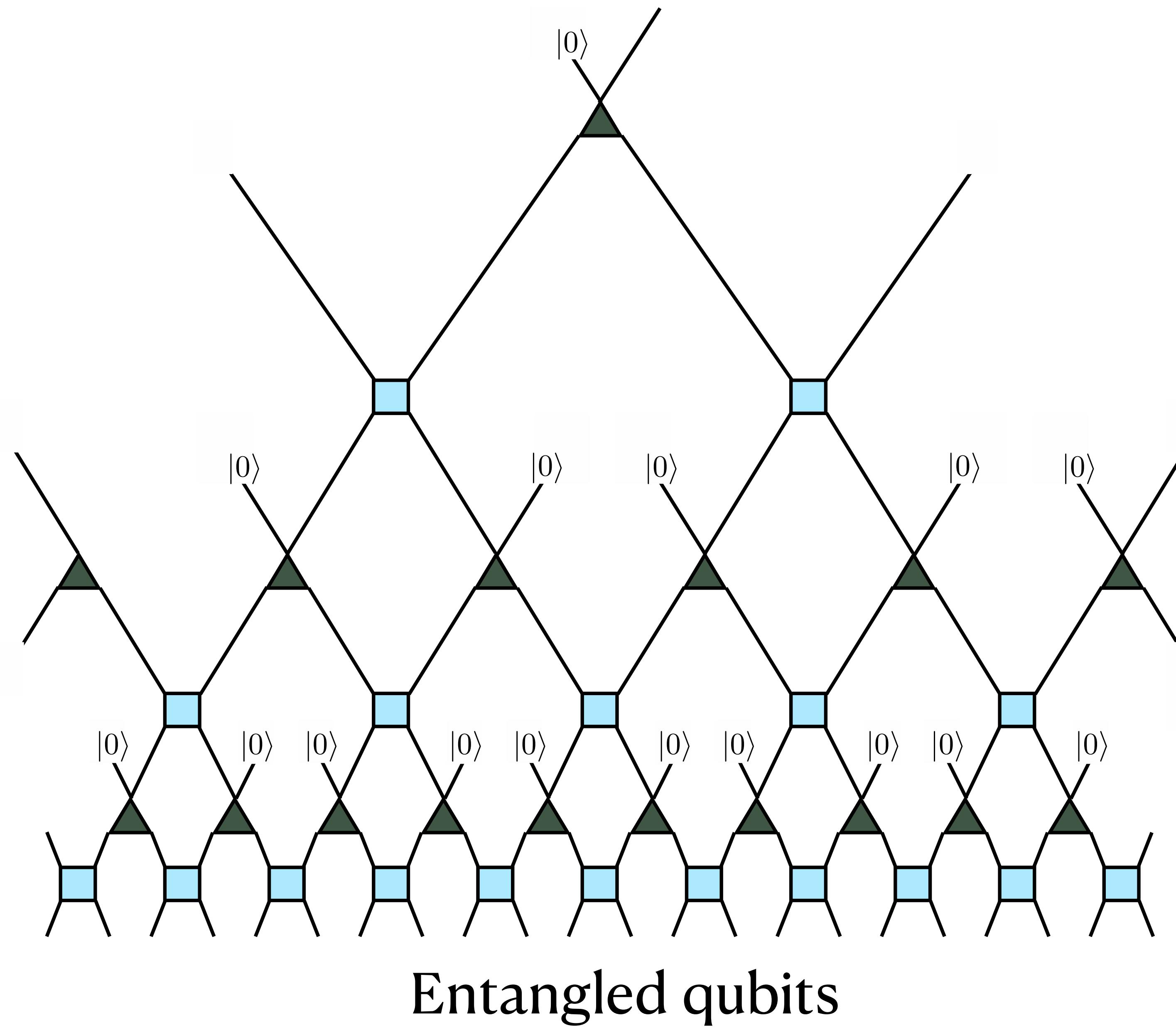


Physical variables

Li, LW, PRL '18 [lio12589/NeuralRG](#)

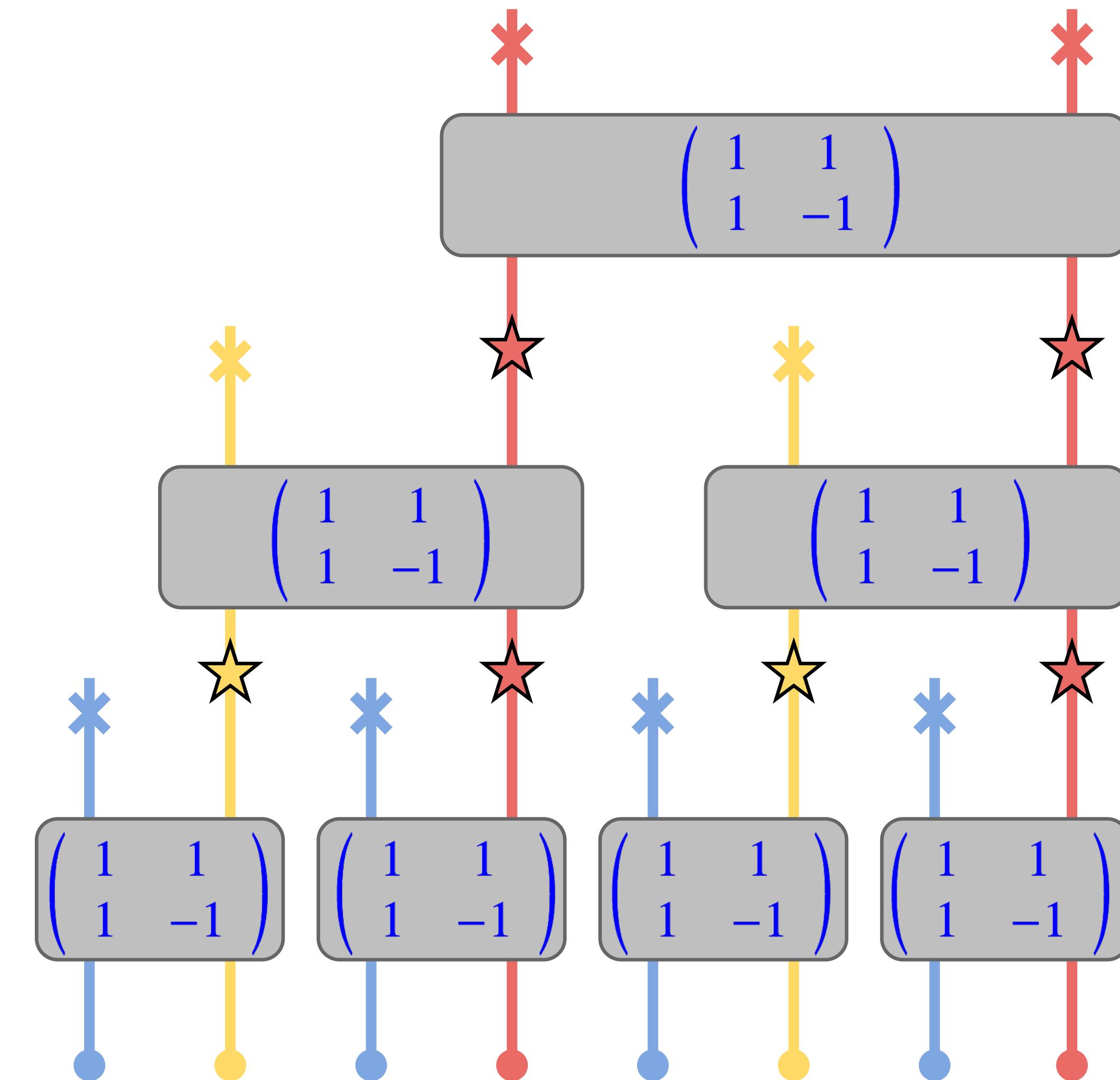
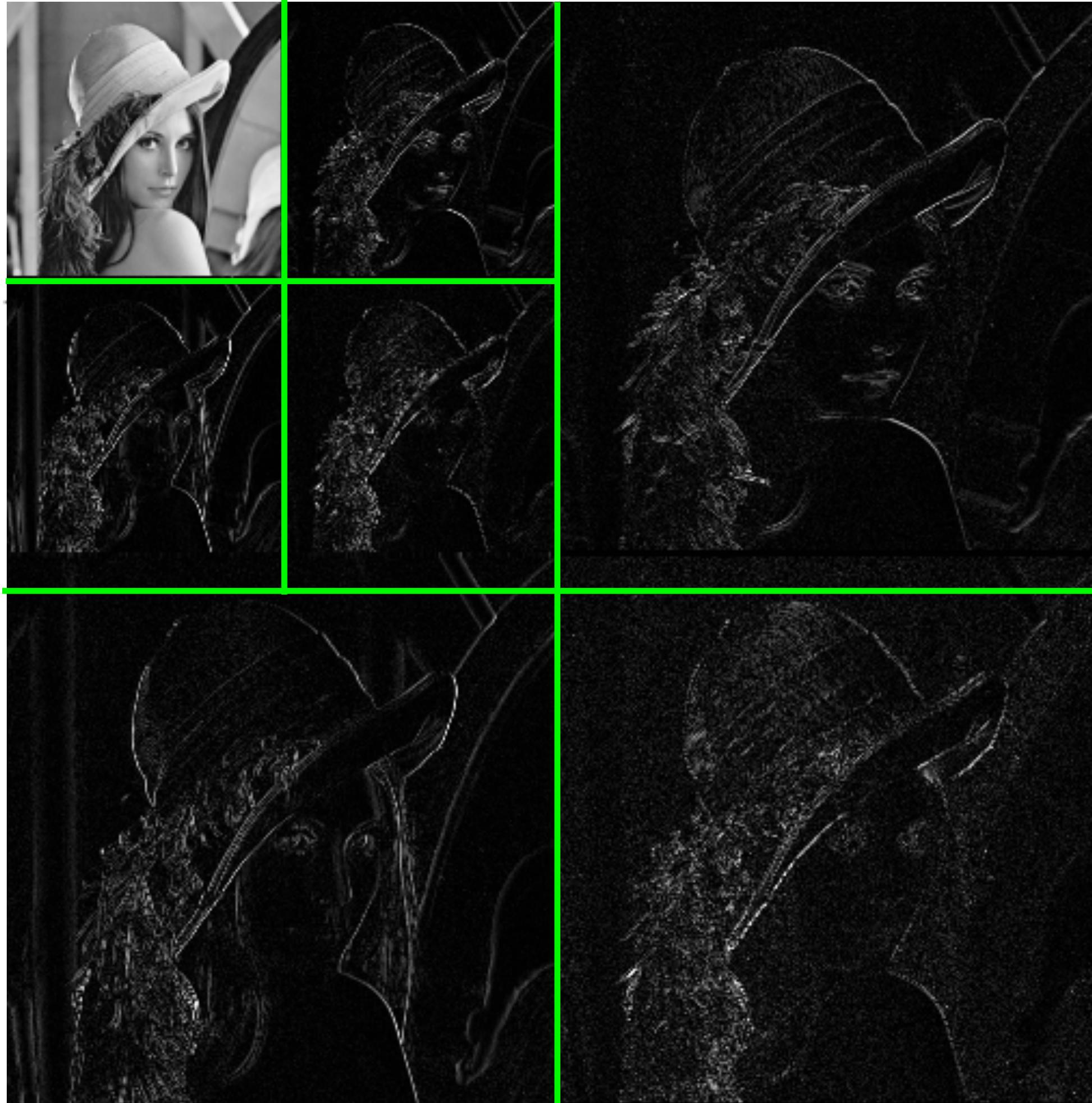


Quantum version of the architecture



**Multi-Scale
Entanglement
Renormalization
Ansatz**

Connection to wavelets



Nonlinear & adaptive generalizations of wavelets

Continuous normalizing flows

$$\ln p(\mathbf{x}) = \ln \mathcal{N}(\mathbf{z}) - \ln \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right|$$

Consider infinitesimal change-of-variables Chen et al 1806.07366

$$\mathbf{x} = \mathbf{z} + \varepsilon \boldsymbol{\nu}$$

$$\varepsilon \rightarrow 0$$

$$\frac{d\mathbf{x}}{dt} = \boldsymbol{\nu}$$

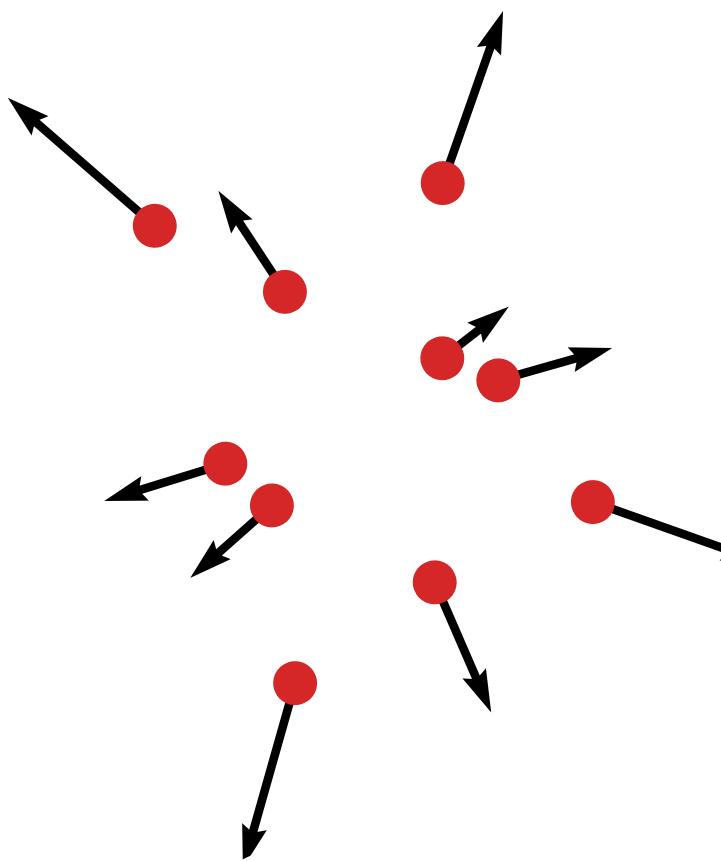
$$\ln p(\mathbf{x}) - \ln \mathcal{N}(\mathbf{z}) = - \ln \left| \det \left(1 + \varepsilon \frac{\partial \boldsymbol{\nu}}{\partial \mathbf{z}} \right) \right|$$

$t = 1$ $t = 0$

$$\frac{d \ln p(\mathbf{x}, t)}{dt} = - \nabla \cdot \boldsymbol{\nu}$$

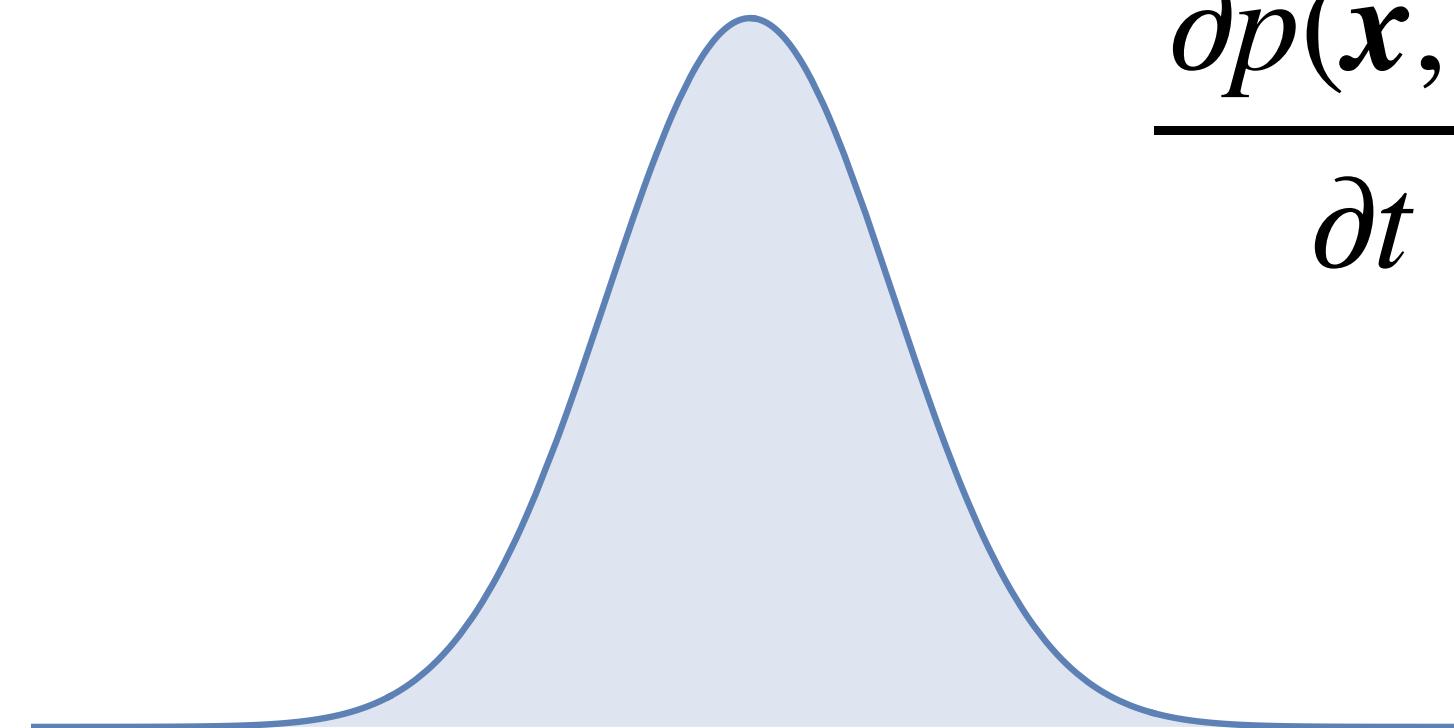
Fluid physics behind flows

$$\frac{dx}{dt} = v$$
$$\frac{d \ln p(x, t)}{dt} = - \nabla \cdot v$$



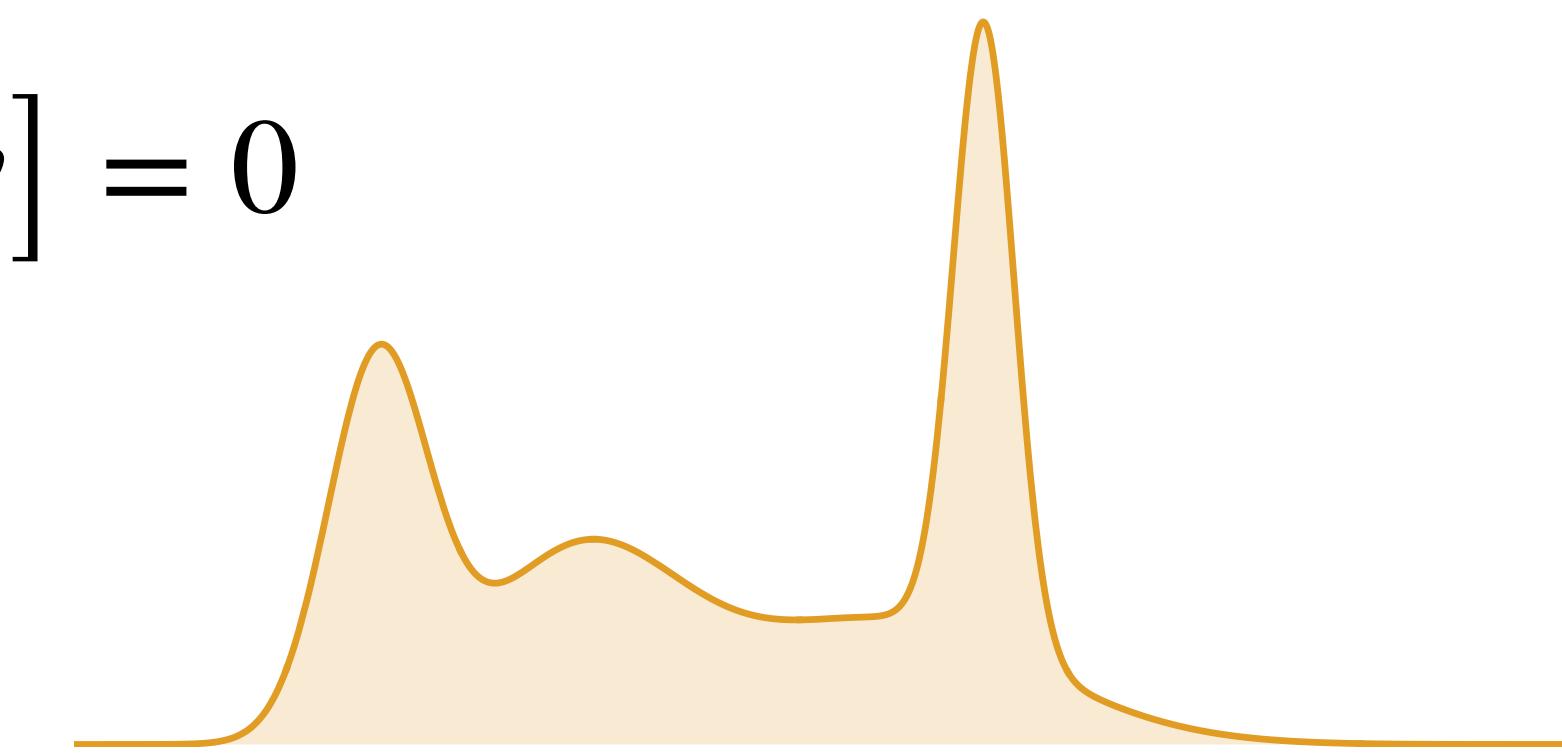
$$\frac{d}{dt} = \frac{\partial}{\partial t} + v \cdot \nabla$$

“material derivative”



Simple density

$$\frac{\partial p(x, t)}{\partial t} + \nabla \cdot [p(x, t)v] = 0$$



Complex density



Zhang, E, LW 1809.10188

[wangleiphy/MongeAmpereFlow](https://github.com/wangleiphy/MongeAmpereFlow)

Infinitesimal Flows Another way to reduce the computational overhead of normalizing flows is to use an ordinary differential equation to generate f [2]. In this case, the probability distribution changes over a finite time from $p(\mathbf{x}; 0)$ to $p(\mathbf{z}; T)$, where \mathbf{z} is the end point of a curve defined by the ODE $\dot{\mathbf{x}}(t) = \mathbf{v}(\mathbf{x}(t))$, $\mathbf{x}(0) = \mathbf{x}$. For a small time step dt , we can approximate $\mathbf{x}(t + dt)$ to first order as $\mathbf{x}(t + dt) = \mathbf{x}(t) + dt\mathbf{v}(\mathbf{x}(t)) + \mathcal{O}(dt^2)$. Plugging this into Eq. 2 yields:

$$\log p(\mathbf{x} + dt\mathbf{v}(\mathbf{x}) + \mathcal{O}(dt^2); t + dt) = \log p(\mathbf{x}; t) - \log |\mathbf{J}_f(\mathbf{x})| \quad (3)$$

$$= \log p(\mathbf{x}; t) - \log |\mathbf{I} + dt\mathbf{J}_v(\mathbf{x}) + \mathcal{O}(dt^2)| \quad (4)$$

Taking a Taylor series gives:

$$\log p(\mathbf{x}; t + dt) + dt\mathbf{v}(\mathbf{x})^T \nabla \log p(\mathbf{x}; t + dt) = \log p(\mathbf{x}; t) - dt \text{Tr}(\mathbf{J}_v(\mathbf{x})) + \mathcal{O}(dt^2) \quad (5)$$

which, in the limit as $dt \rightarrow 0$, becomes:

$$\frac{\partial \log p(\mathbf{x}; t)}{\partial t} = -\mathbf{v}(\mathbf{x})^T \nabla \log p(\mathbf{x}; t) - \text{Tr}(\mathbf{J}_v(\mathbf{x})) = -\mathbf{v}^T \nabla \log p(\mathbf{x}; t) - \nabla \cdot \mathbf{v} \quad (6)$$

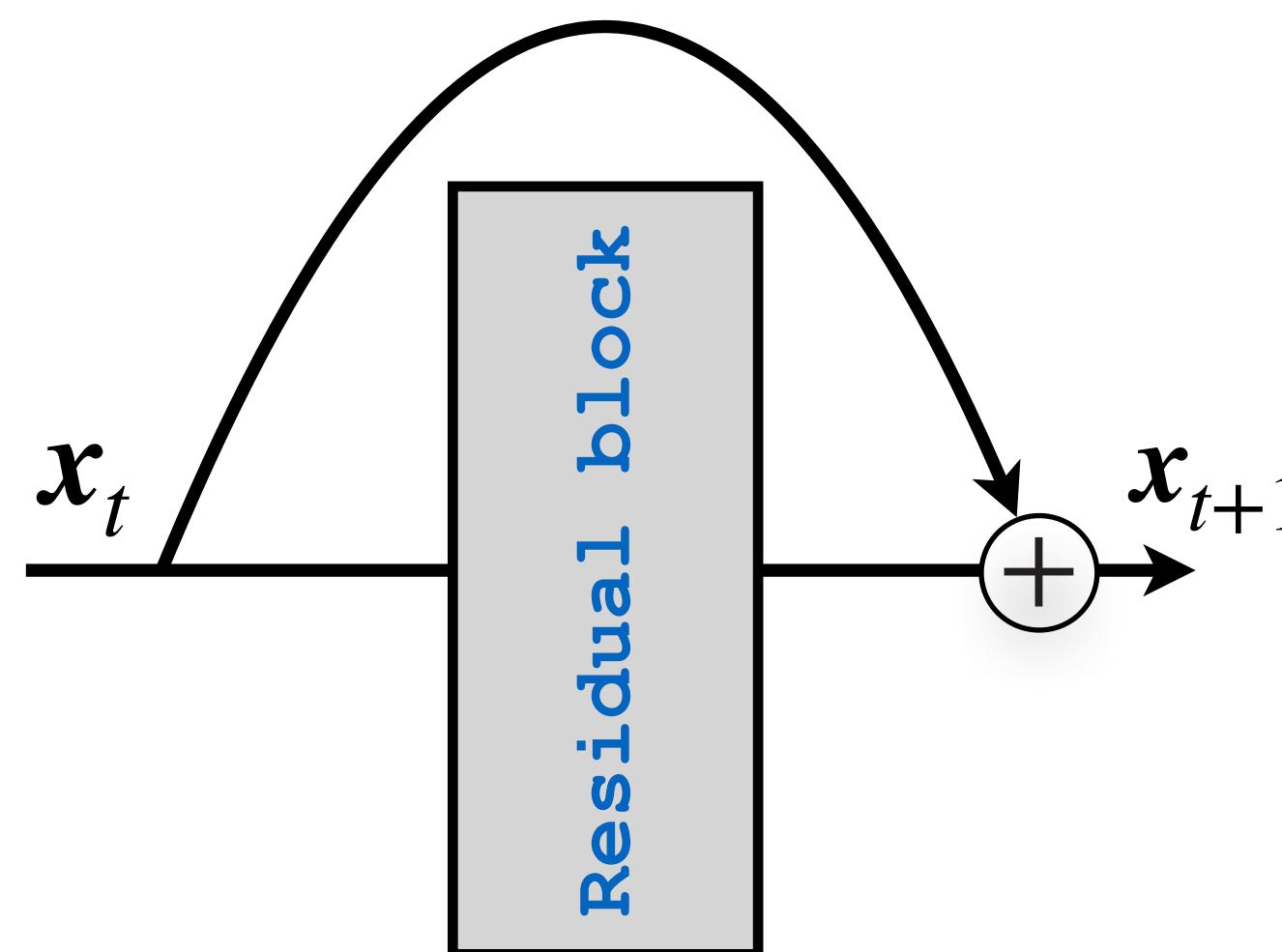
after some rearranging of terms. Here $\nabla \cdot$ is the divergence of a vector field, which is just another way of writing the trace of the Jacobian. The right-hand side of this equation is also the trace of the *Stein operator* of the distribution $p(\mathbf{x})$ applied to the function $\mathbf{v}(\mathbf{x})$, and plays a critical role in Stein variational gradient descent (SVGD) [13]. Switching from the log density to the density (and dropping the t for clarity), we find this expression can be simplified considerably:

$$\begin{aligned} \frac{1}{p(\mathbf{x})} \frac{\partial p(\mathbf{x})}{\partial t} &= -\mathbf{v}^T \frac{\nabla p(\mathbf{x})}{p(\mathbf{x})} - \nabla \cdot \mathbf{v} \\ \frac{\partial p(\mathbf{x})}{\partial t} &= -\mathbf{v}^T \nabla p(\mathbf{x}) - p(\mathbf{x}) \nabla \cdot \mathbf{v} \\ &= -\nabla \cdot (\mathbf{v}(\mathbf{x}) p(\mathbf{x})) \end{aligned} \quad (7)$$

This may also be familiar as the drift term of the Fokker-Planck equation [11, Eq. 6.48] or the continuity equation for conservation of mass in fluid mechanics. We will denote the change to a

Neural Ordinary Differential Equations

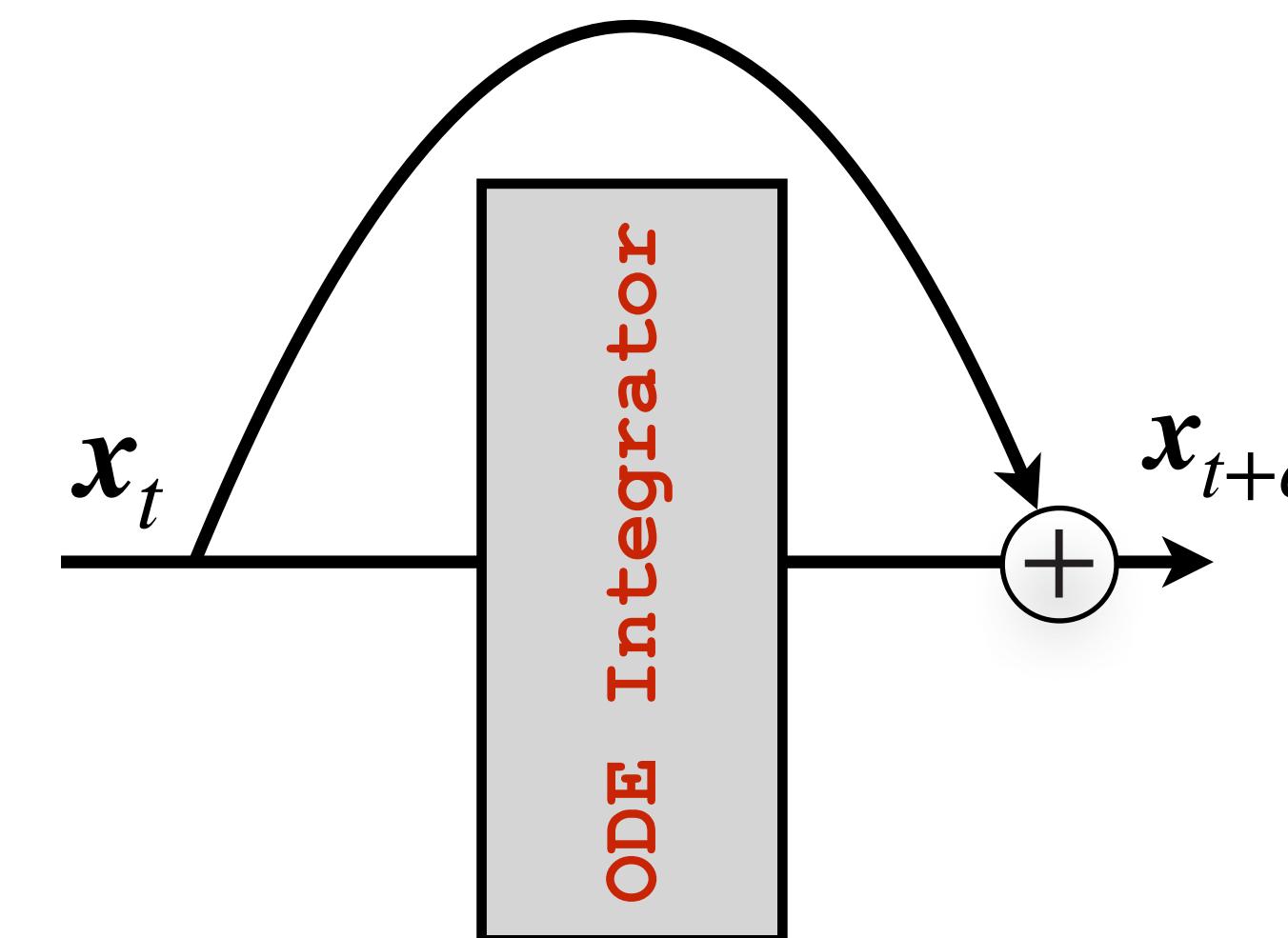
Residual network



$$x_{t+1} = x_t + v(x_t)$$

Chen et al, 1806.07366

ODE integration

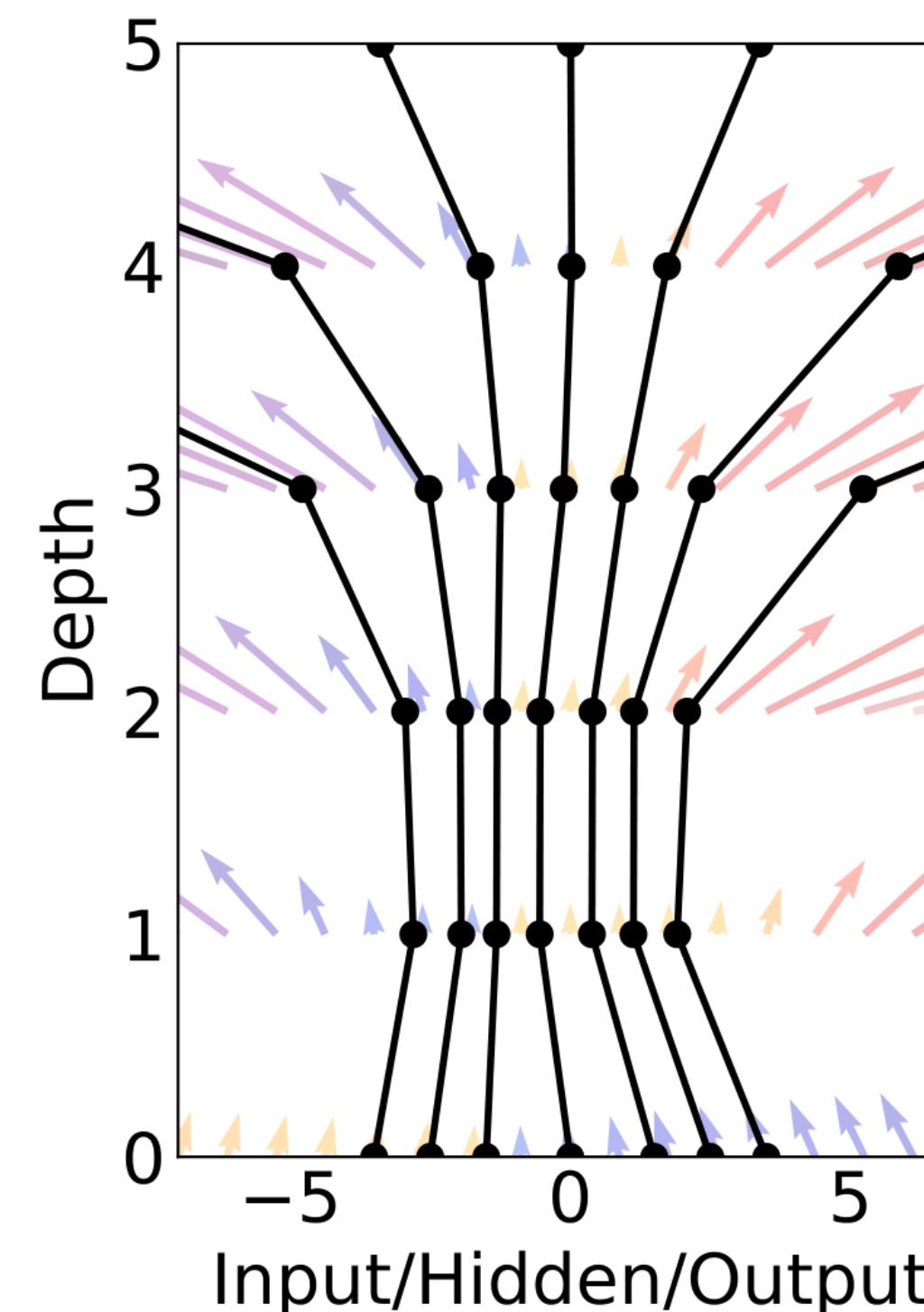


$$dx/dt = v(x)$$

Harbor et al 1705.03341
Lu et al 1710.10121,
E Commun. Math. Stat 17'

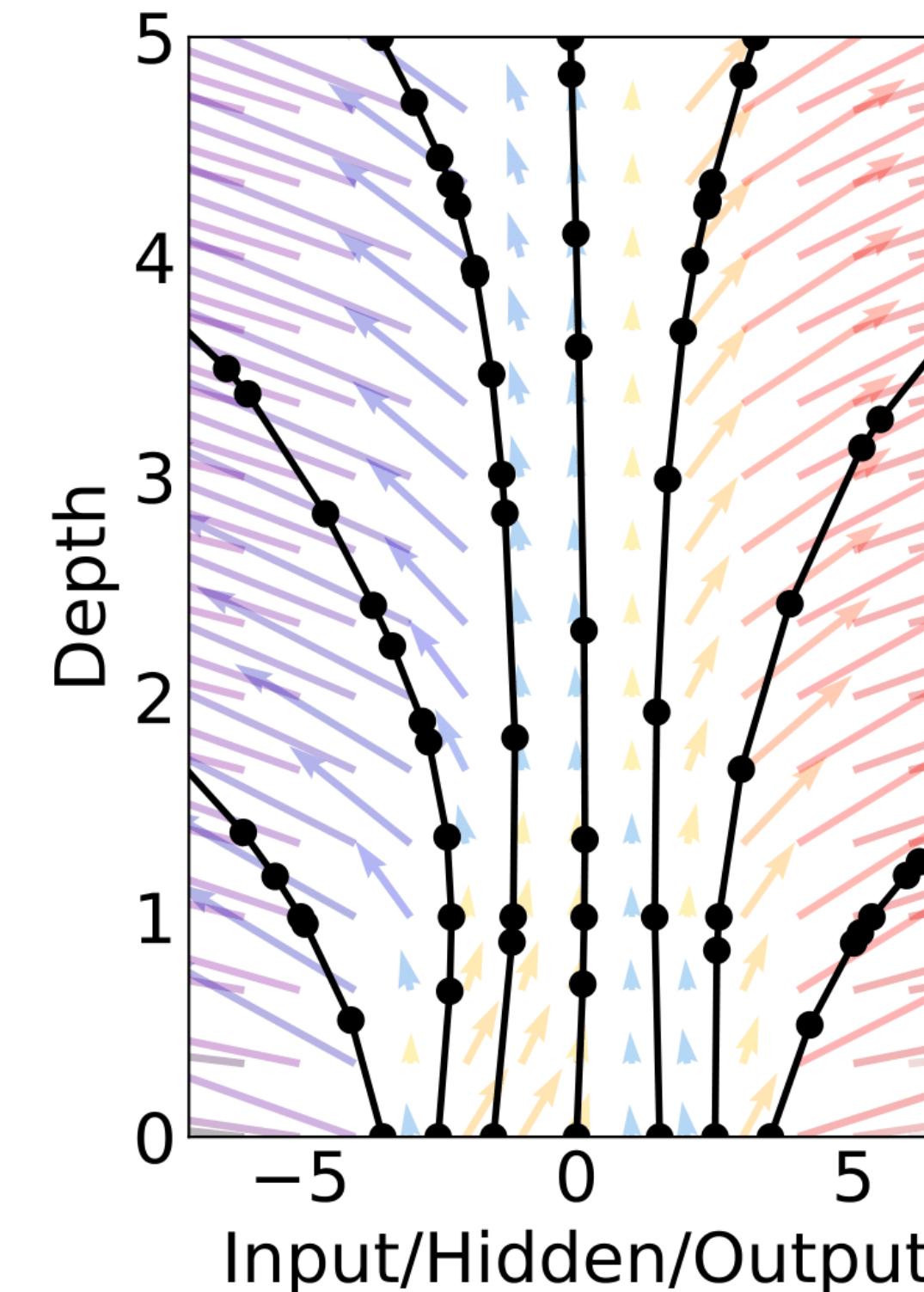
Neural Ordinary Differential Equations

Residual network



$$x_{t+1} = x_t + v(x_t)$$

ODE integration



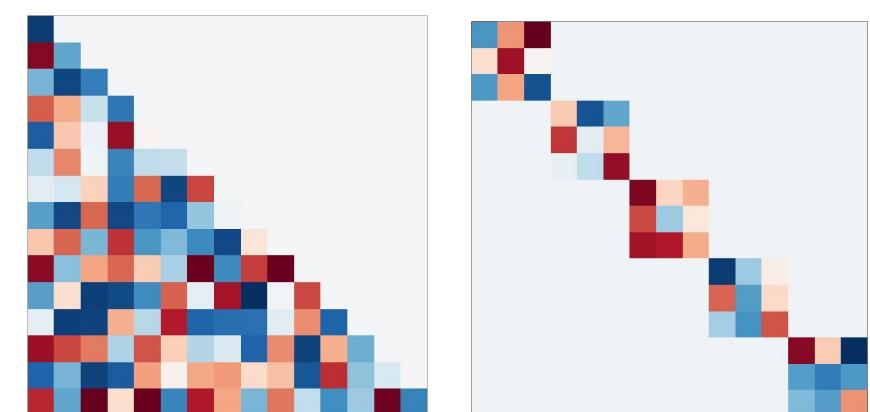
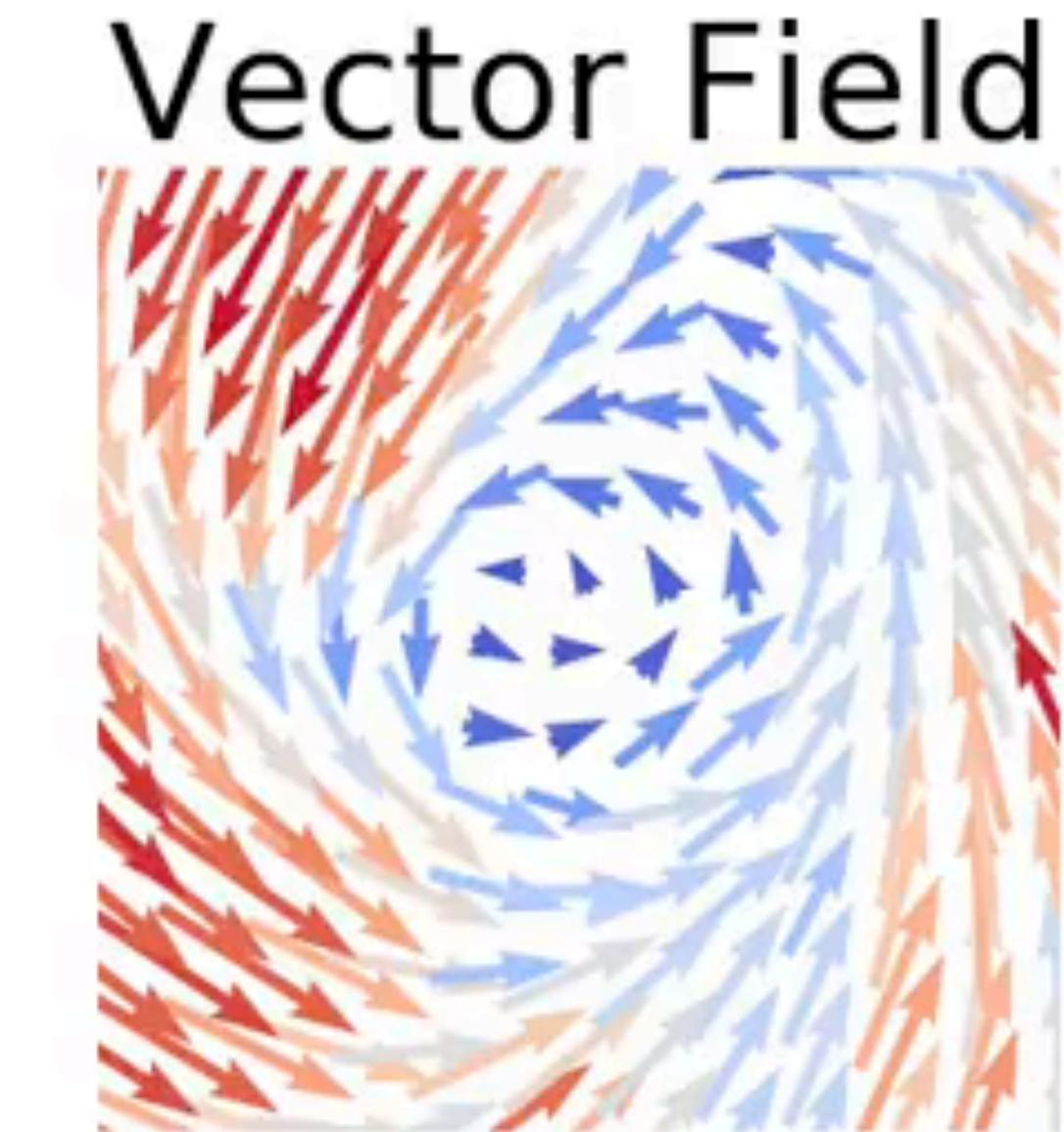
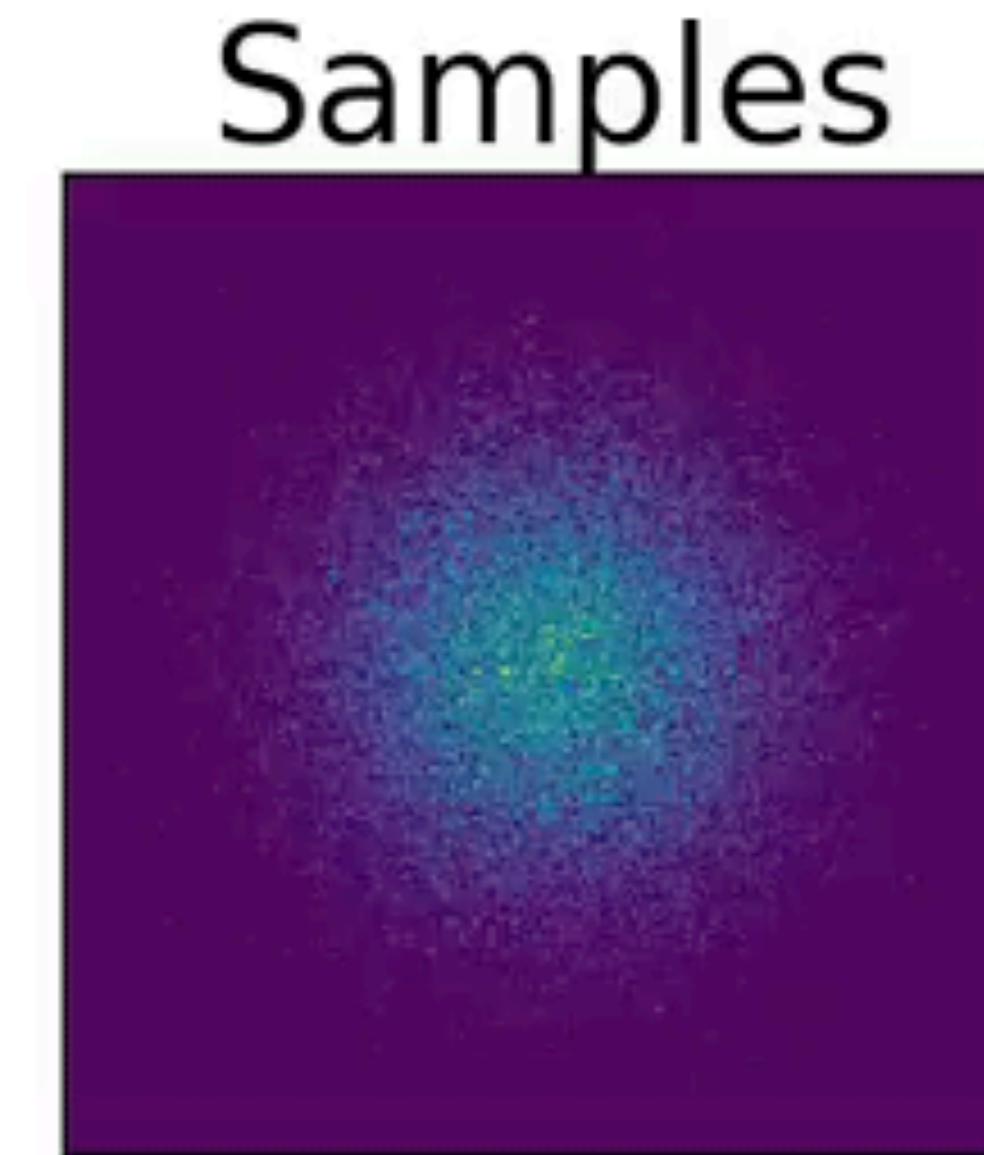
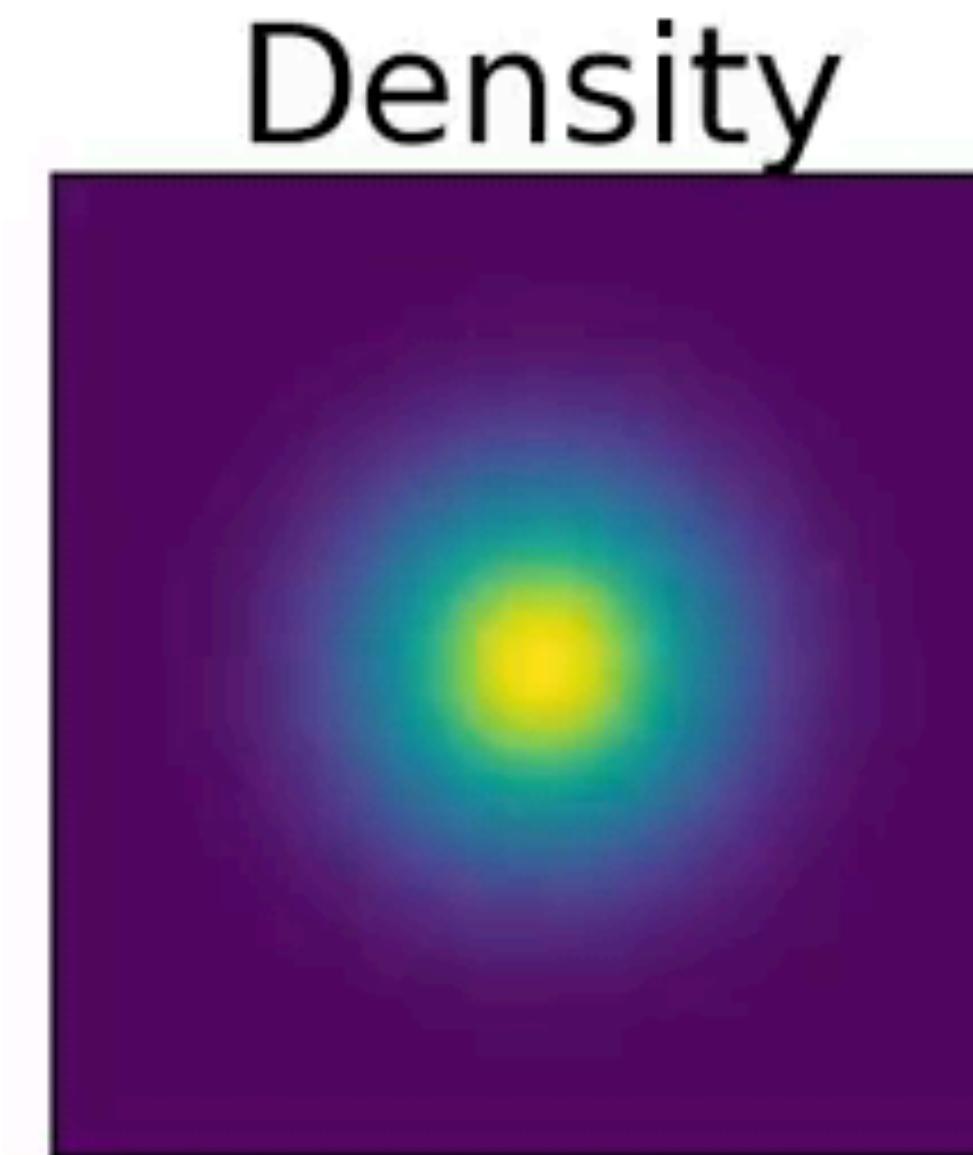
$$dx/dt = v(x)$$

Chen et al, 1806.07366

Harbor et al 1705.03341
Lu et al 1710.10121,
E Commun. Math. Stat 17'...

Continuous normalizing flows implemented with NeuralODE

Chen et al, 1806.07366, Grathwohl et al 1810.01367

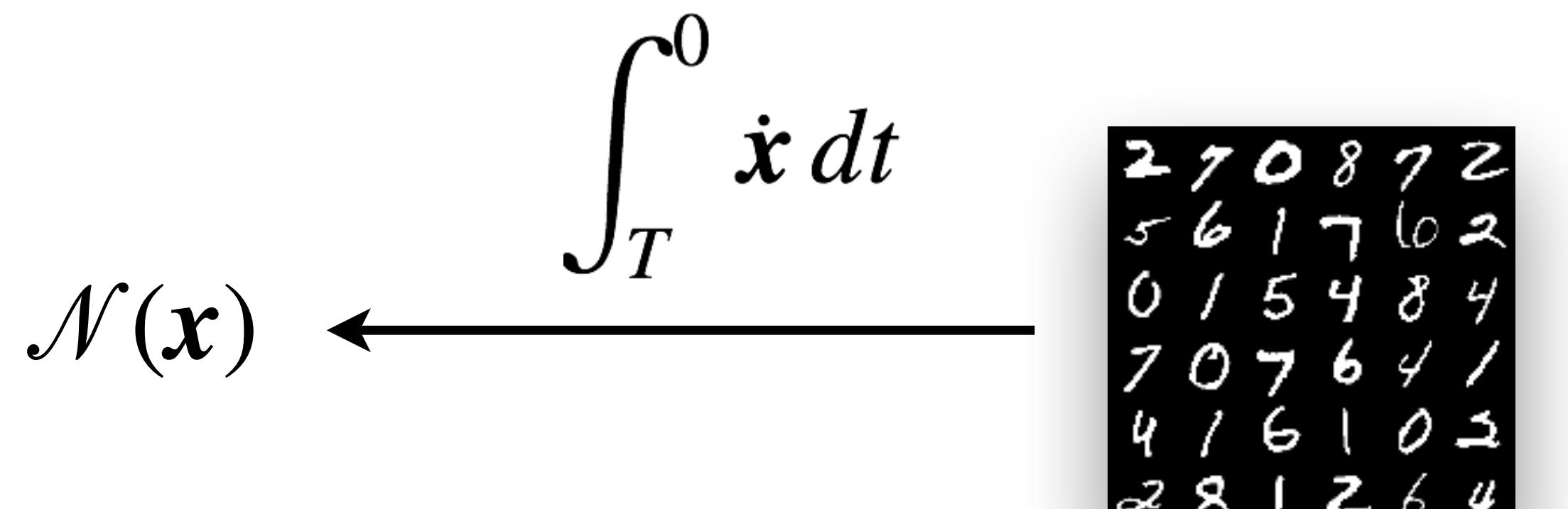


Continuous normalizing flow have no structural
constraints on the transformation Jacobian

The two use cases

Zhang, E, LW, 1809.10188

Maximum likelihood estimation



Variational free energy

$$\mathcal{N}(\mathbf{x}) \xrightarrow{\int_0^T \dot{\mathbf{x}} dt} \frac{e^{-E(\mathbf{x})}}{Z}$$

“learn from data”

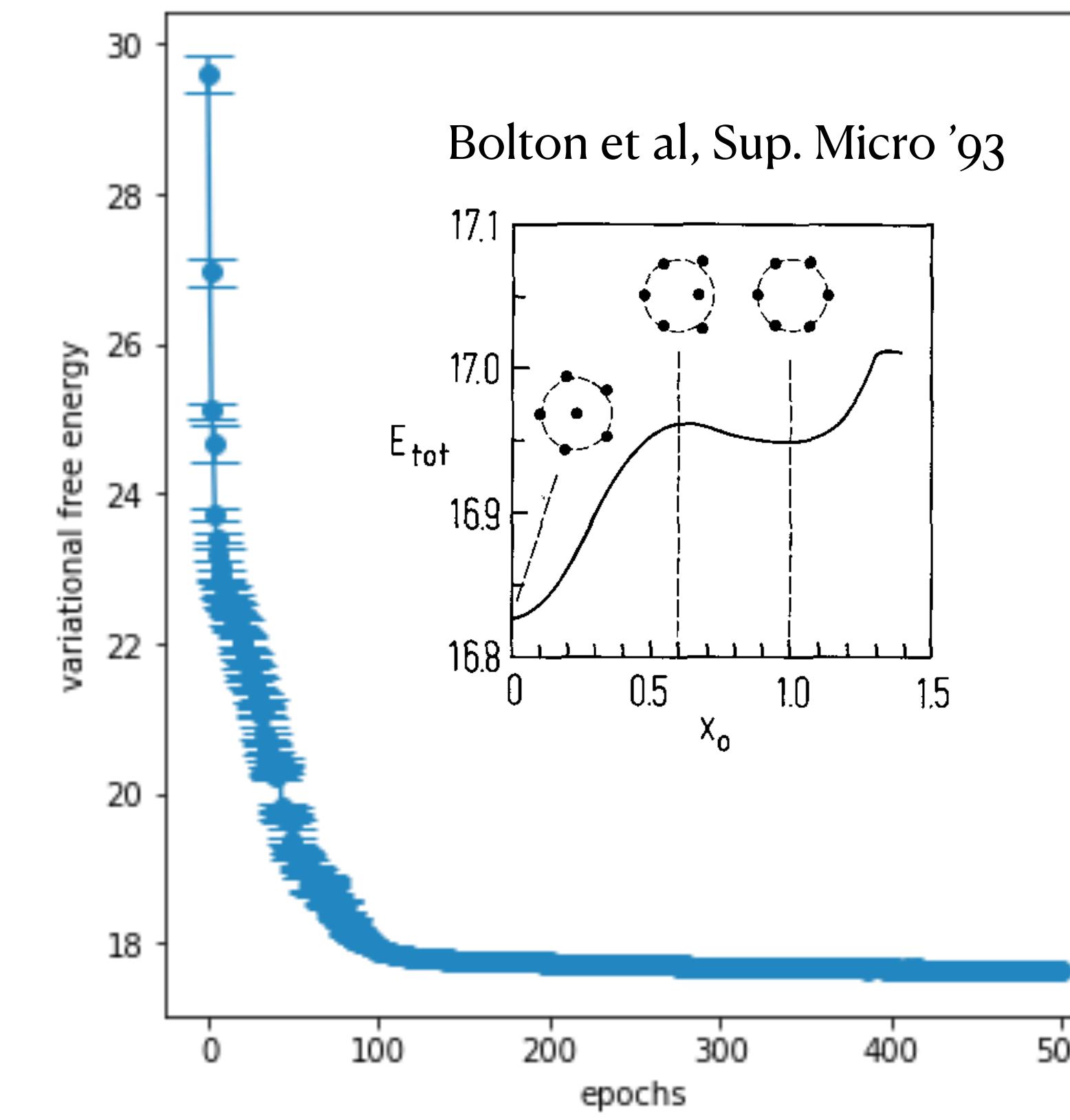
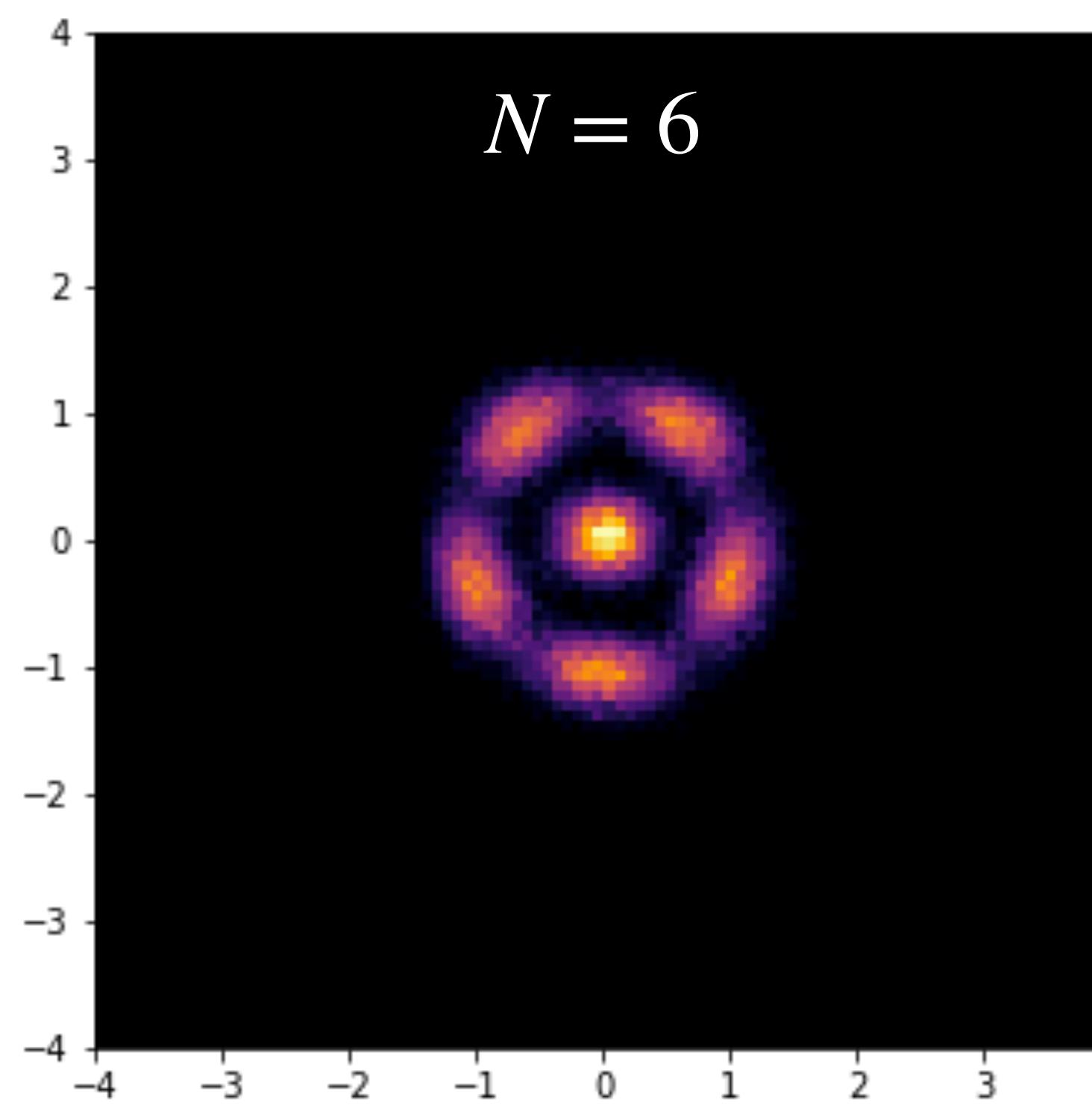
$$\mathcal{L} = -\mathbb{E}_{\mathbf{x} \sim \text{data}} [\ln p(\mathbf{x})]$$

“learn from Hamiltonian”

$$F = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [k_B T \ln p(\mathbf{x}) + H(\mathbf{x})]$$

Demo: Classical Coulomb gas in a harmonic trap

$$H = \sum_{i < j} \frac{1}{|x_i - x_j|} + \sum_i x_i^2$$





Fermi Flow

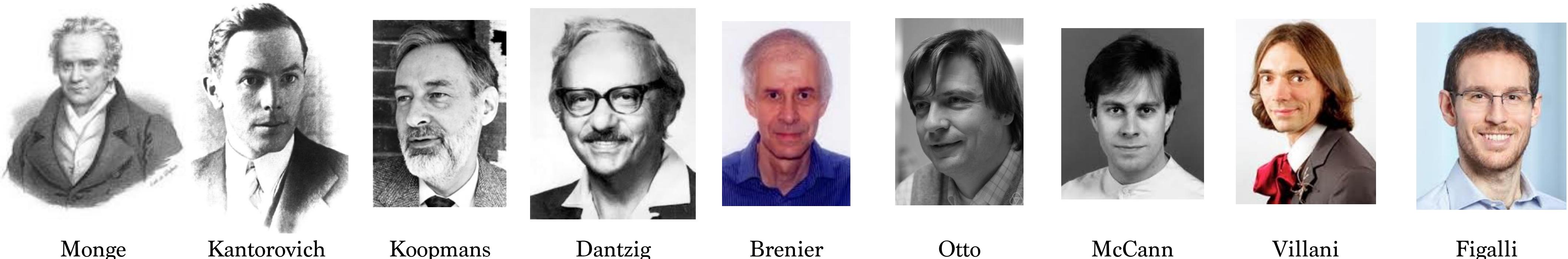
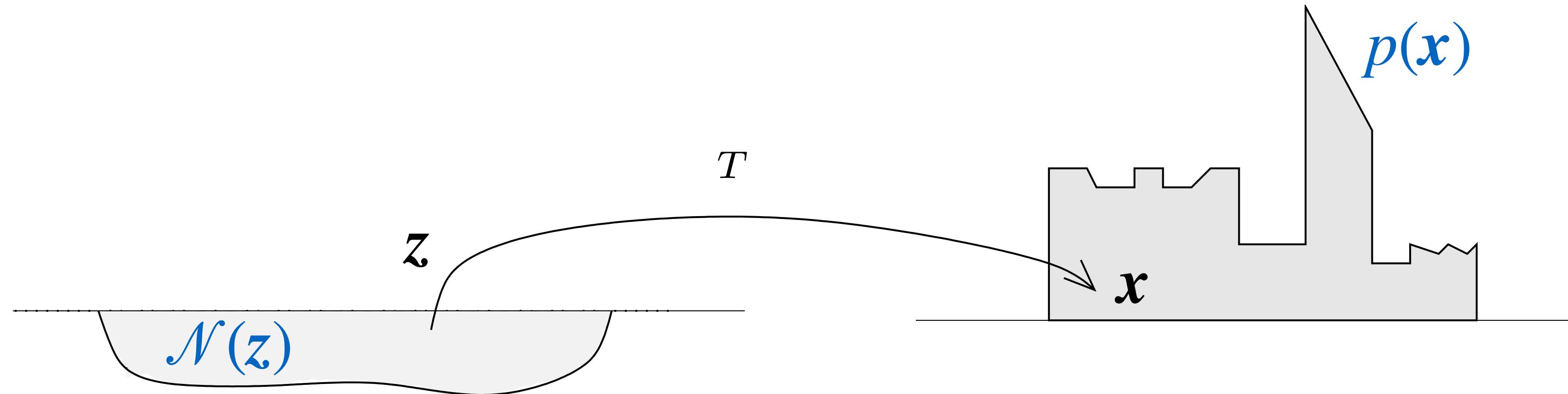
Xie, Zhang, LW, 2105.08644, JML '22

github.com/fermiflow

Continuous flow of electron density in a quantum dot

Optimal Transport Theory

Monge problem (1781): How to transport earth with optimal cost ?



Monge

Kantorovich

Koopmans

Dantzig

Brenier

Otto

McCann

Villani

Figalli

Nobel Prize in Economics '75

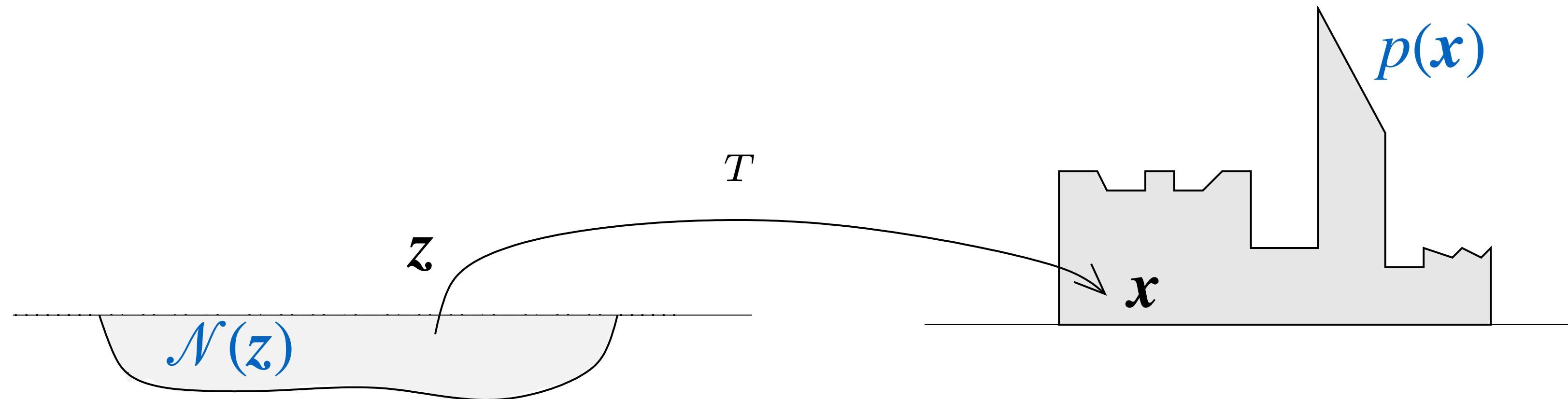
Fields Medal '10

Fields Medal '18

from Cuturi, Solomon NISP 2017 tutorial

Optimal Transport Theory

Monge problem (1781): How to transport earth with optimal cost ?



Brenier theorem (1991)

Under certain conditions
the optimal map is

$$z \mapsto x = \nabla u(z)$$

Monge-Ampère Equation

$$\frac{\mathcal{N}(z)}{p(\nabla u(z))} = \det \left(\frac{\partial^2 u}{\partial z_i \partial z_j} \right)$$

Monge-Ampère Flow

Zhang, E, LW 1809.10188



[wangleiphy/MongeAmpereFlow](#)

$$\frac{\partial \rho(\mathbf{x}, t)}{\partial t} + \nabla \cdot [\rho(\mathbf{x}, t) \nabla \varphi] = 0$$

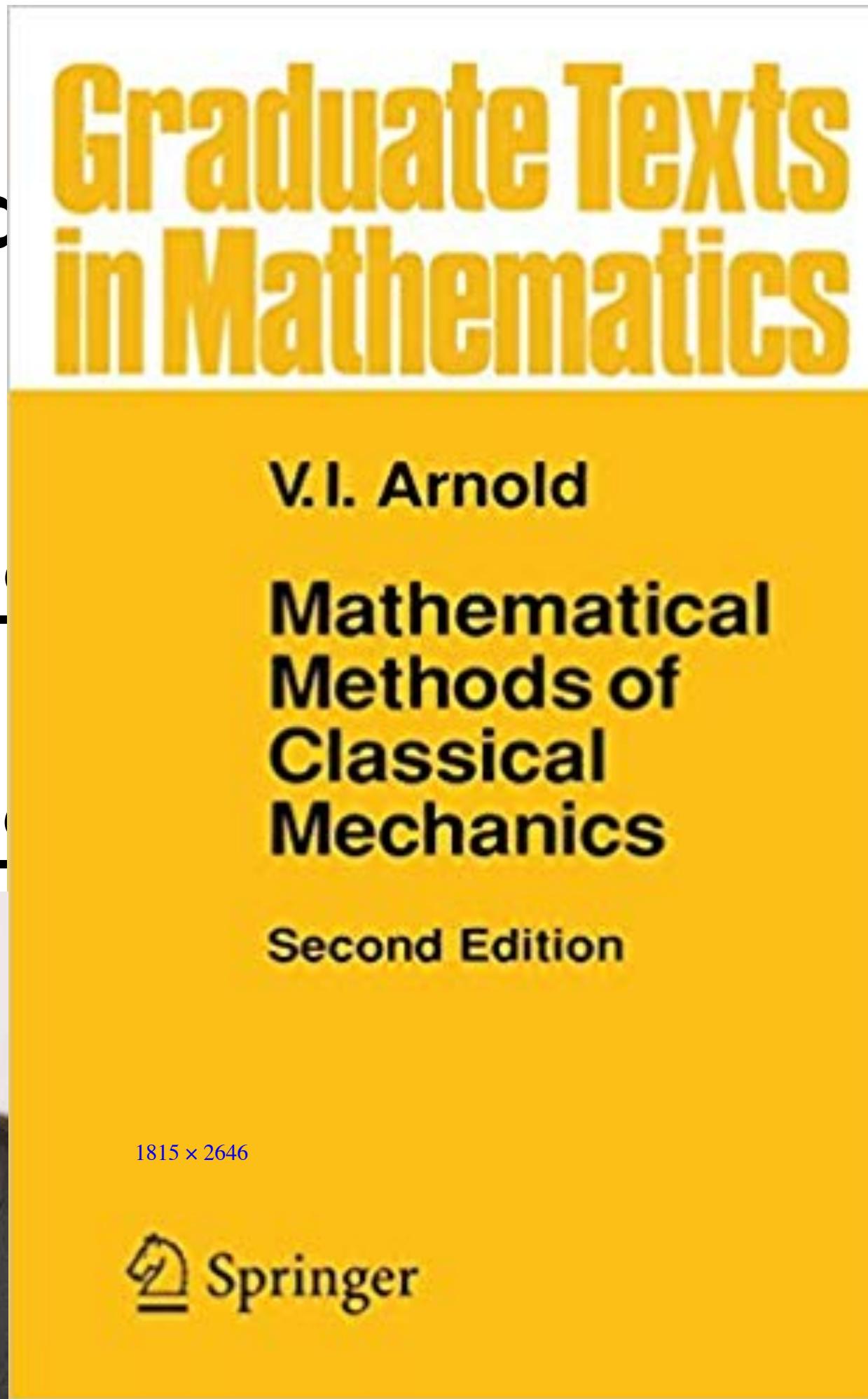
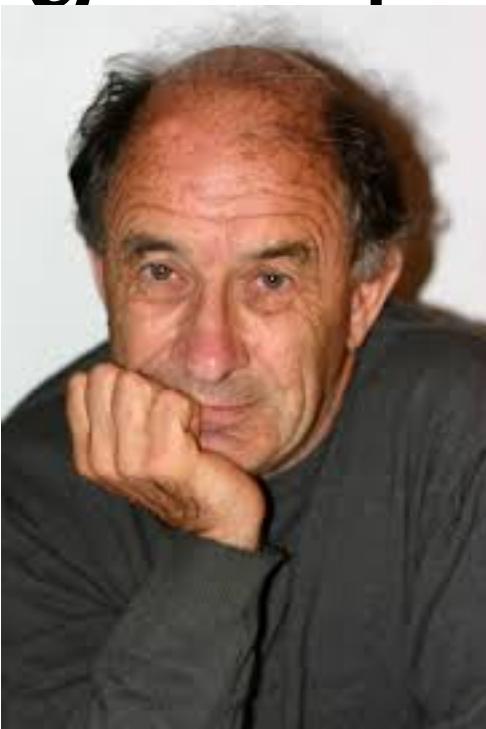
- ① Drive the flow with an “irrotational” velocity field
- ② Impose symmetry to the scalar valued potential for symmetric generative model

$$\varphi(g\mathbf{x}) = \varphi(\mathbf{x}) \implies \rho(g\mathbf{x}) = \rho(\mathbf{x})$$

Hamiltonian dynamics: phase space flow

Hamiltonian eq

$$\begin{cases} \dot{p} = -\frac{\partial H}{\partial q} \\ \dot{q} = +\frac{\partial H}{\partial p} \end{cases}$$



pace va

$$= (p, q)$$

lectic m

$$\begin{pmatrix} I \\ -I \end{pmatrix}$$

Hamilton

Feng Kang Qin Mengzhao
冯康 秦孟兆 著

浙江科学技术出版社

辛几何算法

哈密尔顿系统的

$$\begin{aligned} f^*\omega &= \omega \\ \omega &= \sum dp_i \wedge dq^i \end{aligned}$$

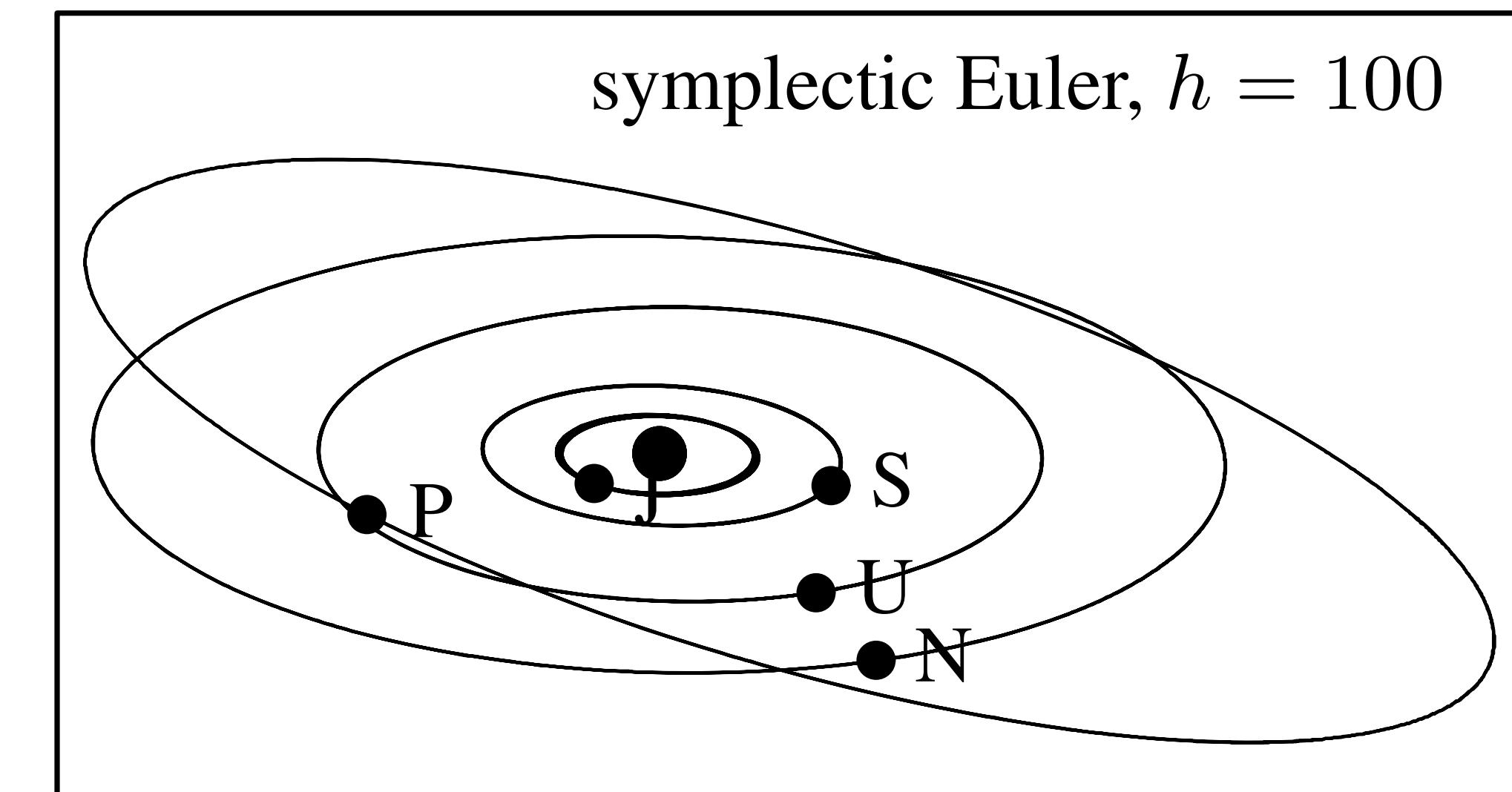
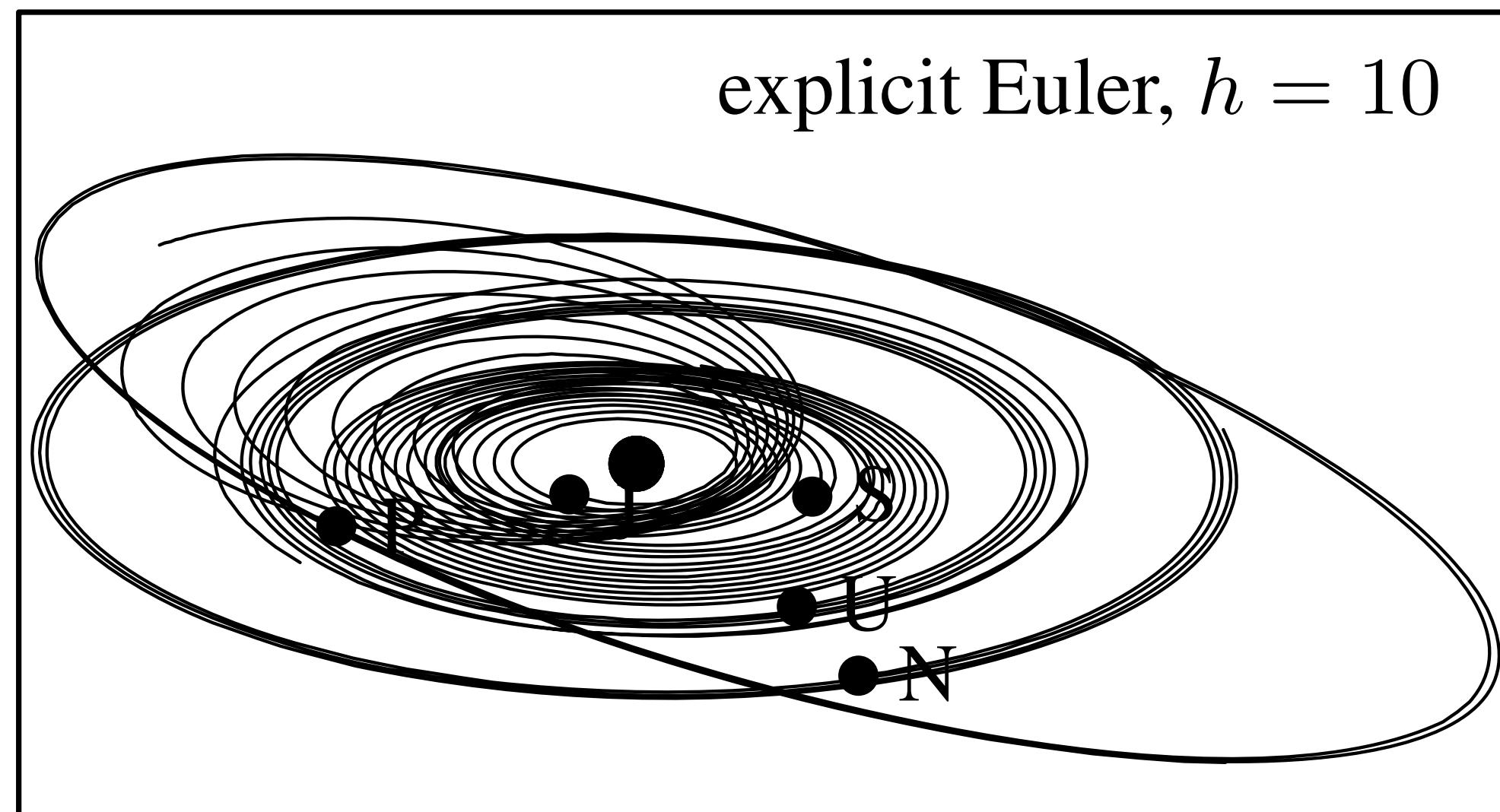
Symplectic Geometric
Algorithms for
Hamiltonian Systems



$$\nabla_x H(x) J$$

ic gradient flow

Symplectic Integrators



Canonical Transformations

$$x = (p, q) \quad \longleftrightarrow \quad z = (P, Q)$$

Change of variables

which satisfies

$$(\nabla_x z) J (\nabla_x z)^T = J$$

symplectic condition

one has

$$\dot{z} = \nabla_z K(z) J \quad \text{where} \quad K(z) = H \circ x(z)$$

Preserves Hamiltonian dynamics in the “latent phase space”

A probabilistic perspective

Canonical transformation deforms phase space density $p(x) = e^{-\beta H(x)}$

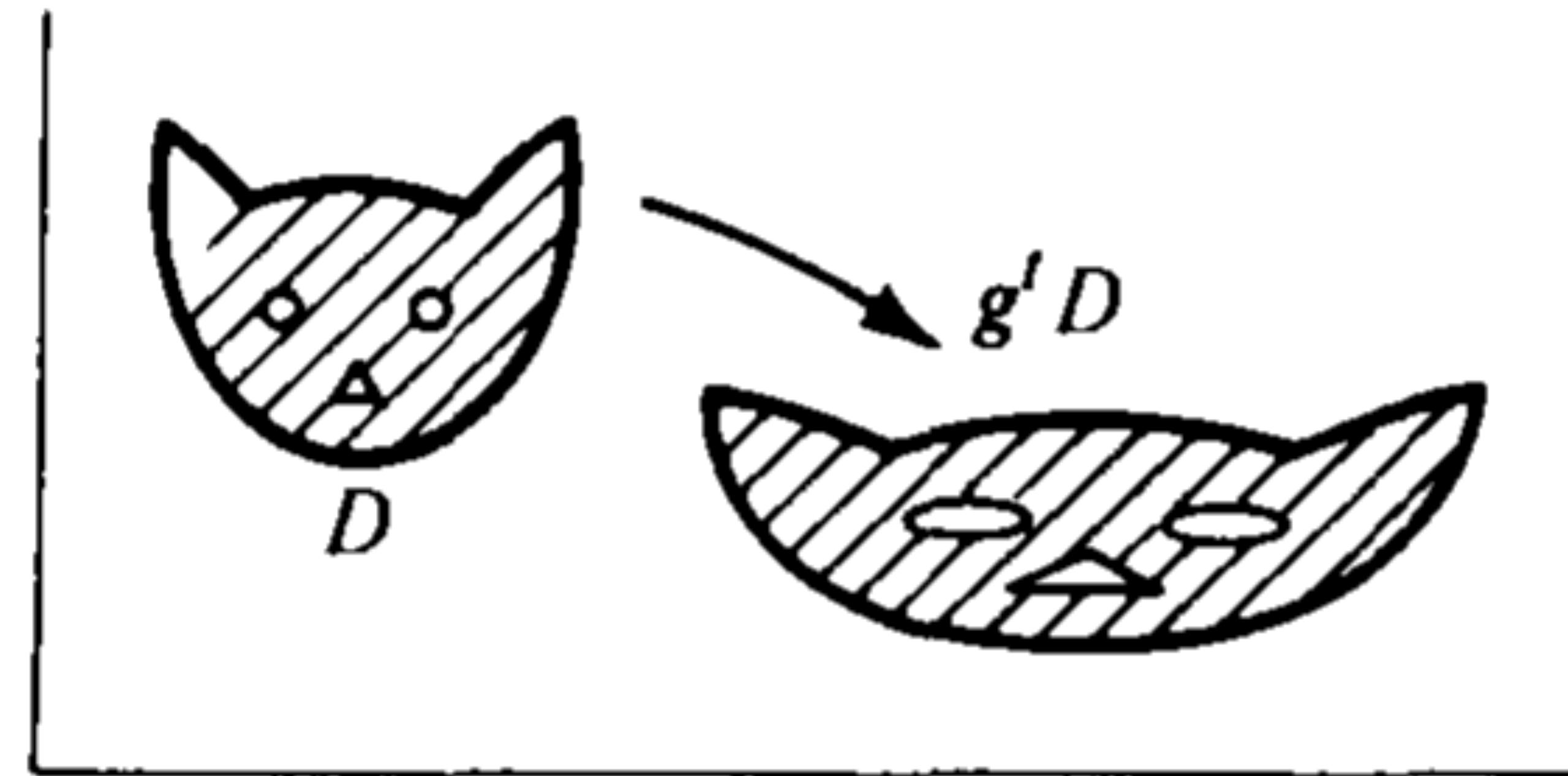
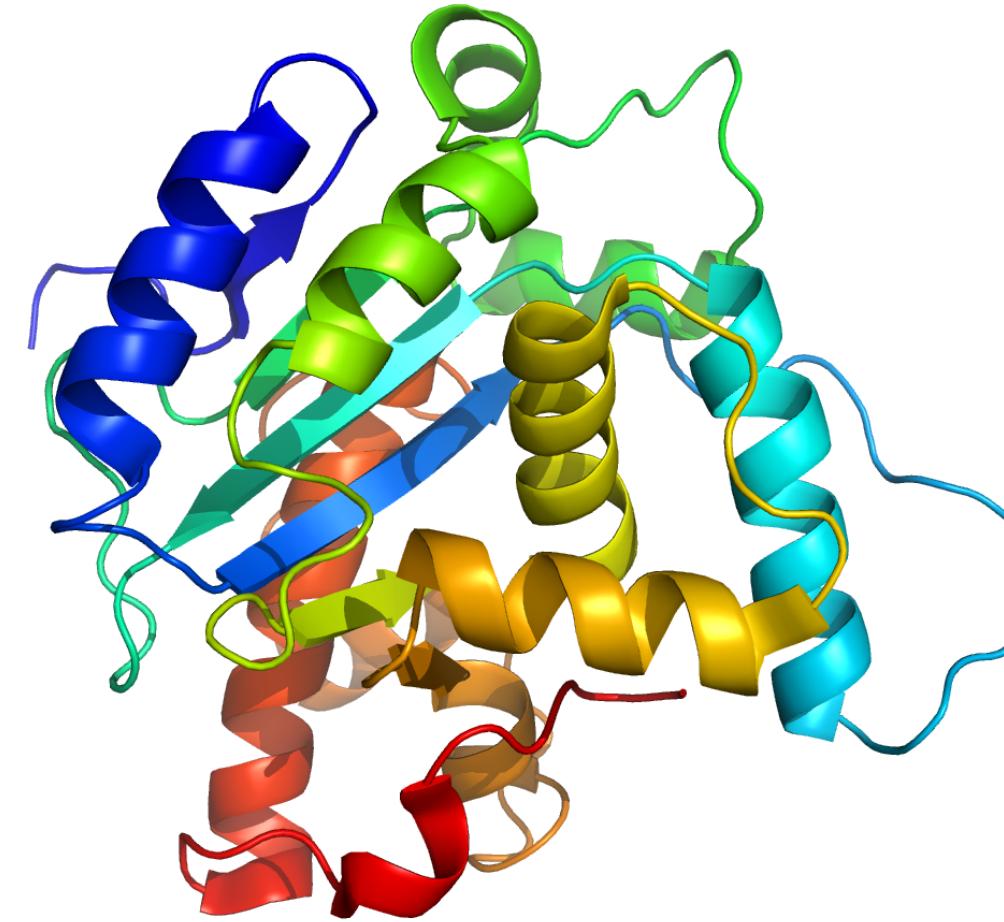


Image from Arnold '78

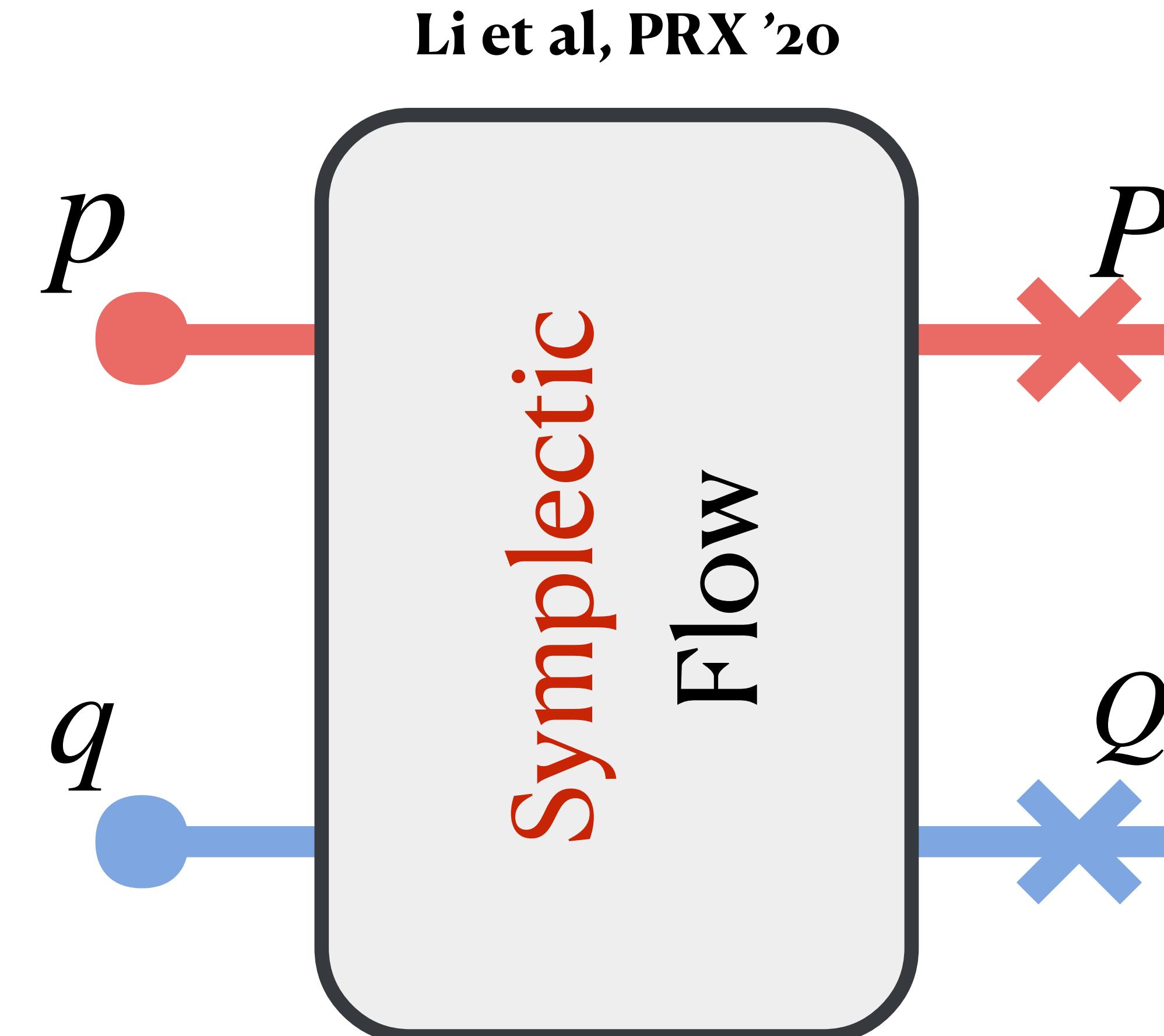
Use a symplectic flow to learn canonical transformations!

Neural Canonical Transformations

$$H(p, q)$$



physical space



$$K(P, Q) = \sum_k \frac{P_k^2 + \omega_k^2 Q_k^2}{2}$$



latent space

Learn the network parameter and the latent harmonic frequency

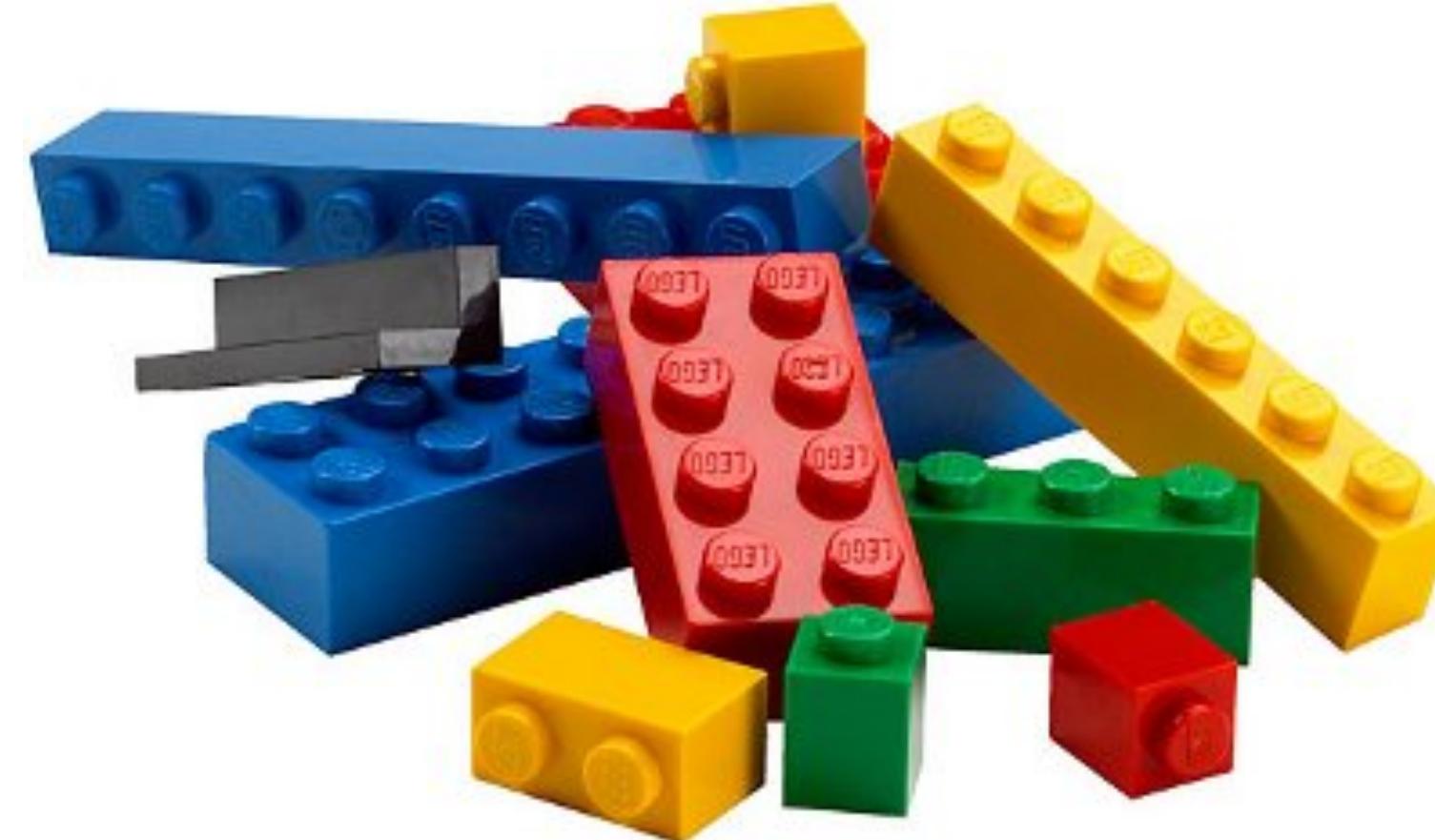
Modular design of the symplectic network

$$z = \mathcal{T}(x)$$

$$\mathcal{T} = \mathcal{T}_1 \circ \mathcal{T}_2 \circ \mathcal{T}_3 \circ \dots$$

$$(\nabla_x z) J (\nabla_x z)^T = J$$

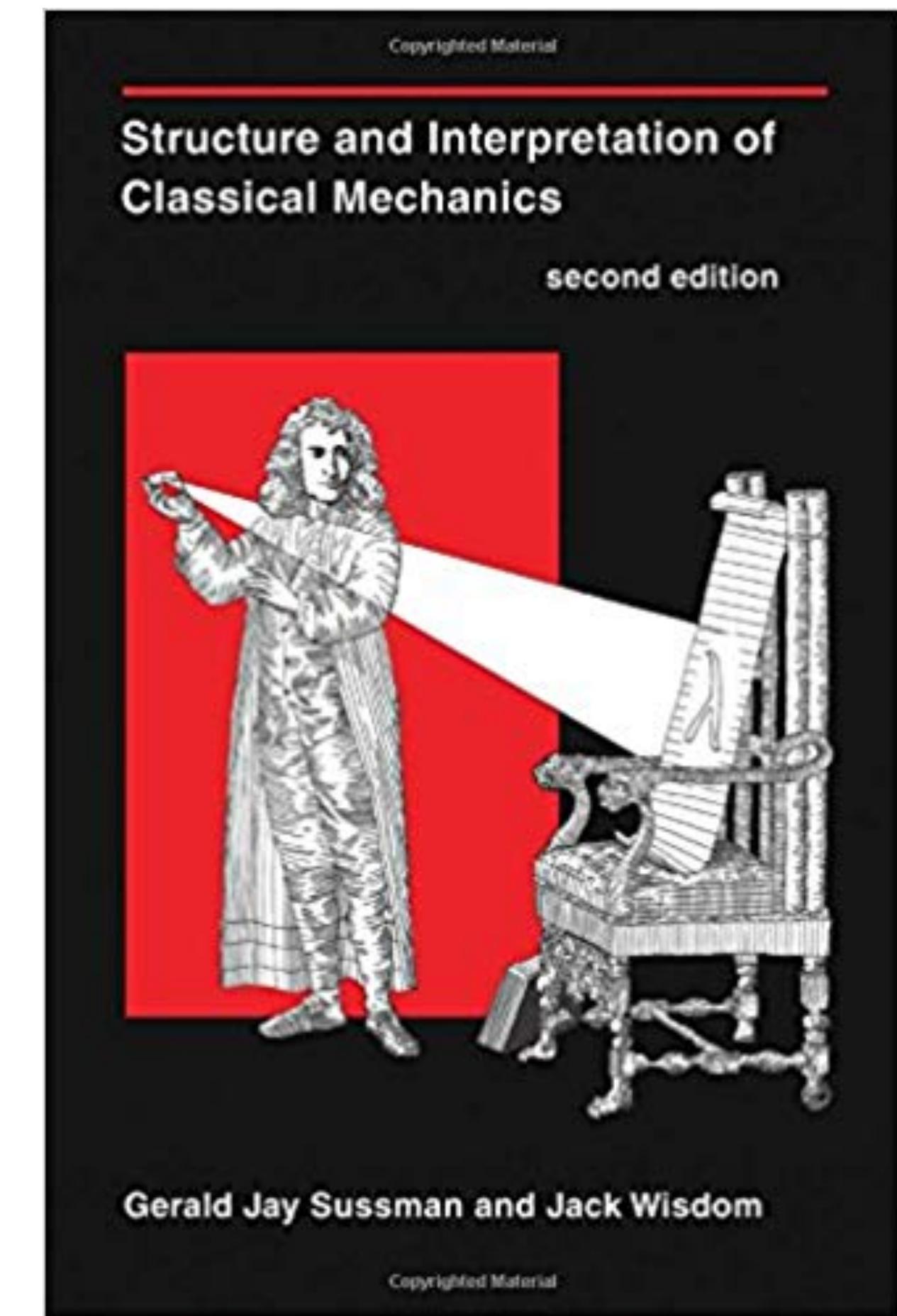
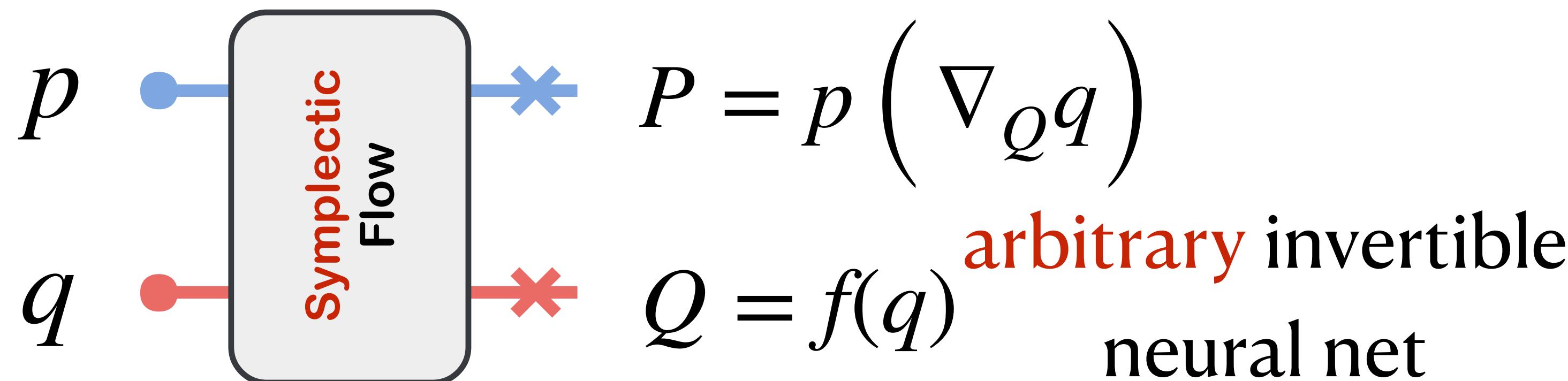
symplectic group



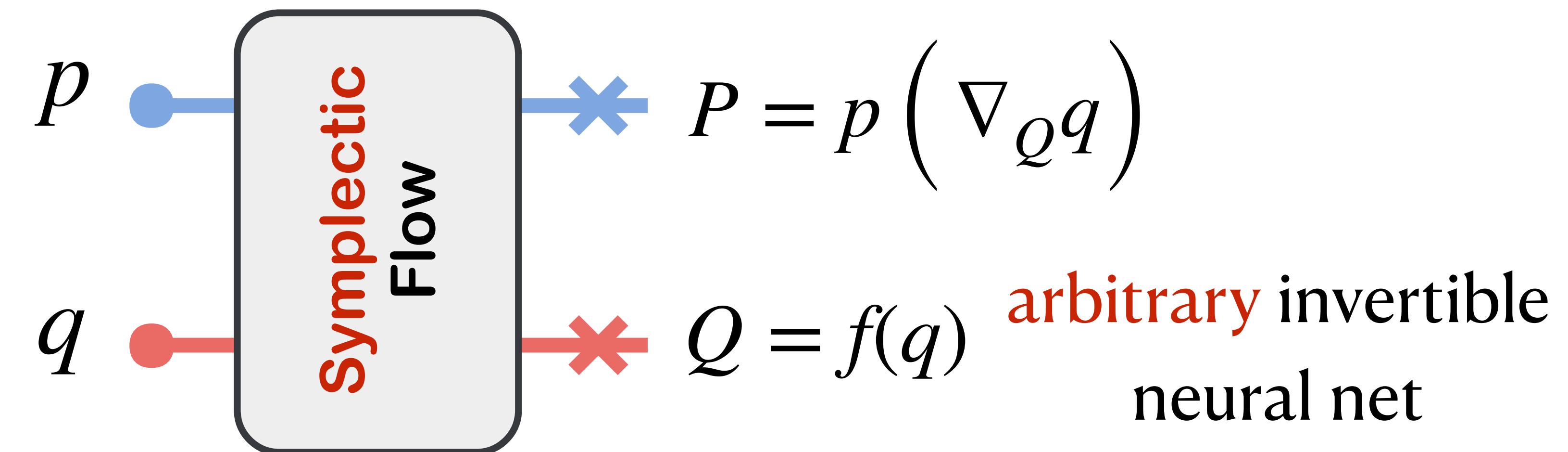
Compose symplectic primitives to form a deep neural network

Symplectic primitives

- Linear transformation: Symplectic Lie algebra
- Continuous-time flow: Symplectic generating functions
 - Symplectic integrator as a “layer” in deep net
 - Harbor et al 1705.03341
 - Lu et al 1710.10121
 - E, Commun. Math. Stat 17
 - Chen et al 1806.07366
 - Zhu et al, 2004.13830
 - ...
- Neural point transformation



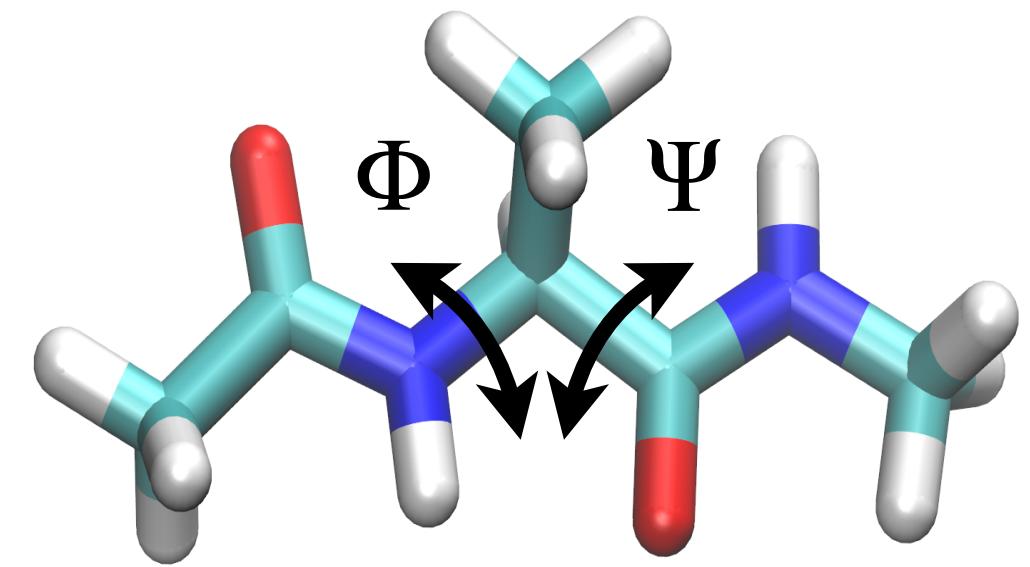
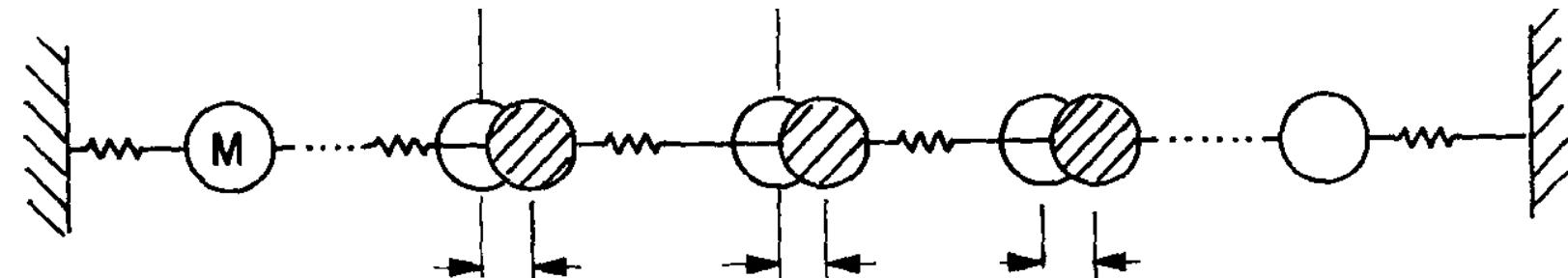
Neural point transformation



```
Q = forward(params, q)
_, vjp = jax.vjp(lambda Q: inverse(params, Q), Q)
P = vjp(p)[0]
```

How is this going to be useful?

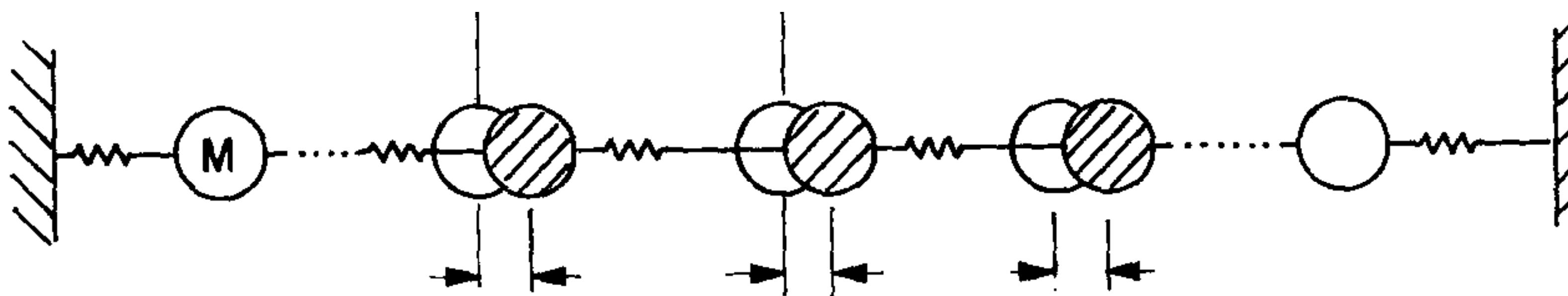
Let's play with examples!



3	4	2	1	9	5	6	2	1	8
8	9	1	2	5	0	0	6	6	4
6	7	0	1	6	3	6	3	7	0
3	7	7	9	4	6	6	1	8	2
2	9	3	4	3	9	8	7	2	5
1	5	9	8	3	6	5	7	2	3
9	3	1	9	1	5	8	0	8	4
5	6	2	6	8	5	8	8	9	9
3	7	7	0	9	4	8	5	4	3
7	9	6	4	7	0	6	9	2	3

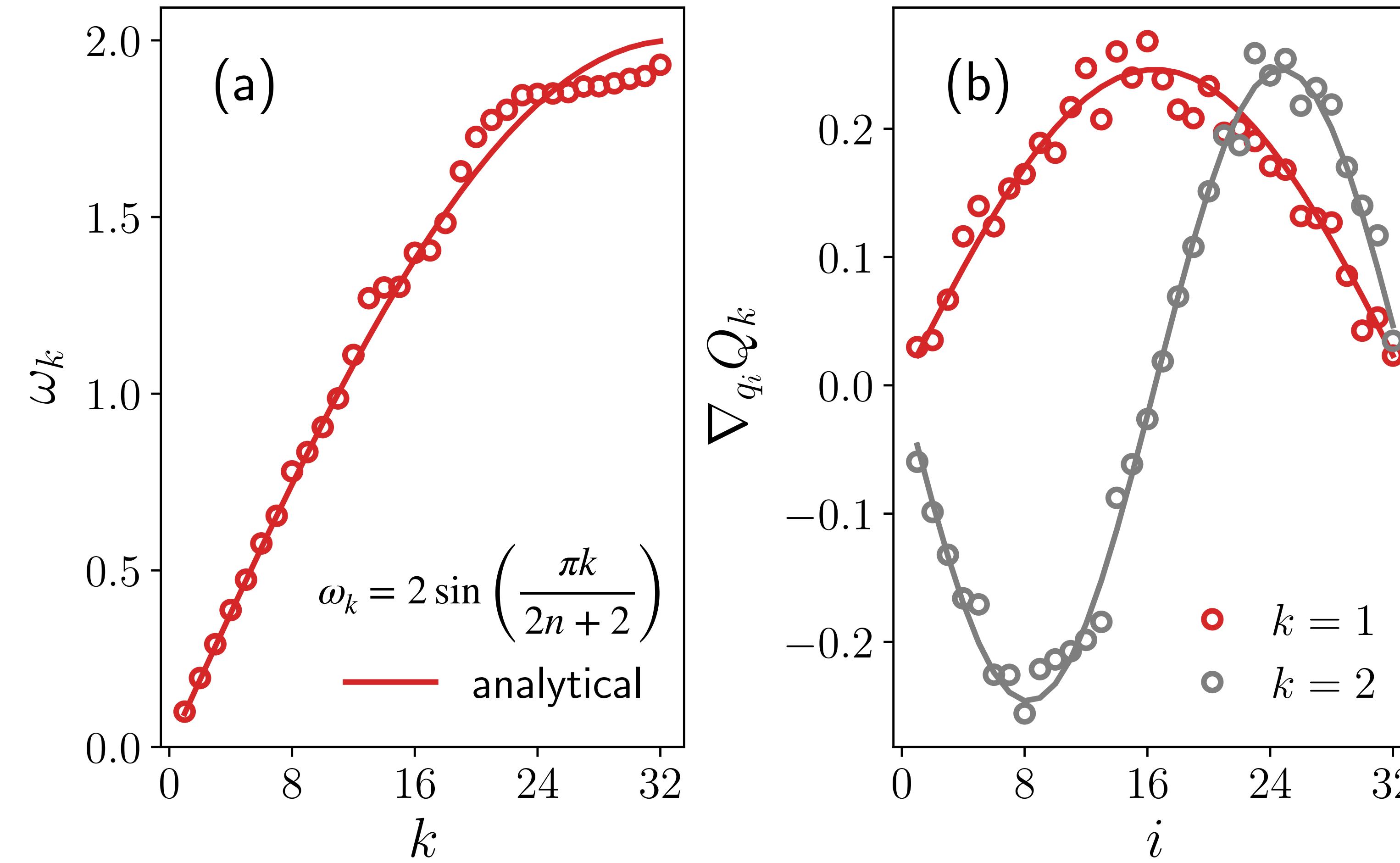
Warm up: Harmonic Chain

$$H = \frac{1}{2} \sum_{i=1}^n [p_i^2 + (q_i - q_{i-1})^2]$$



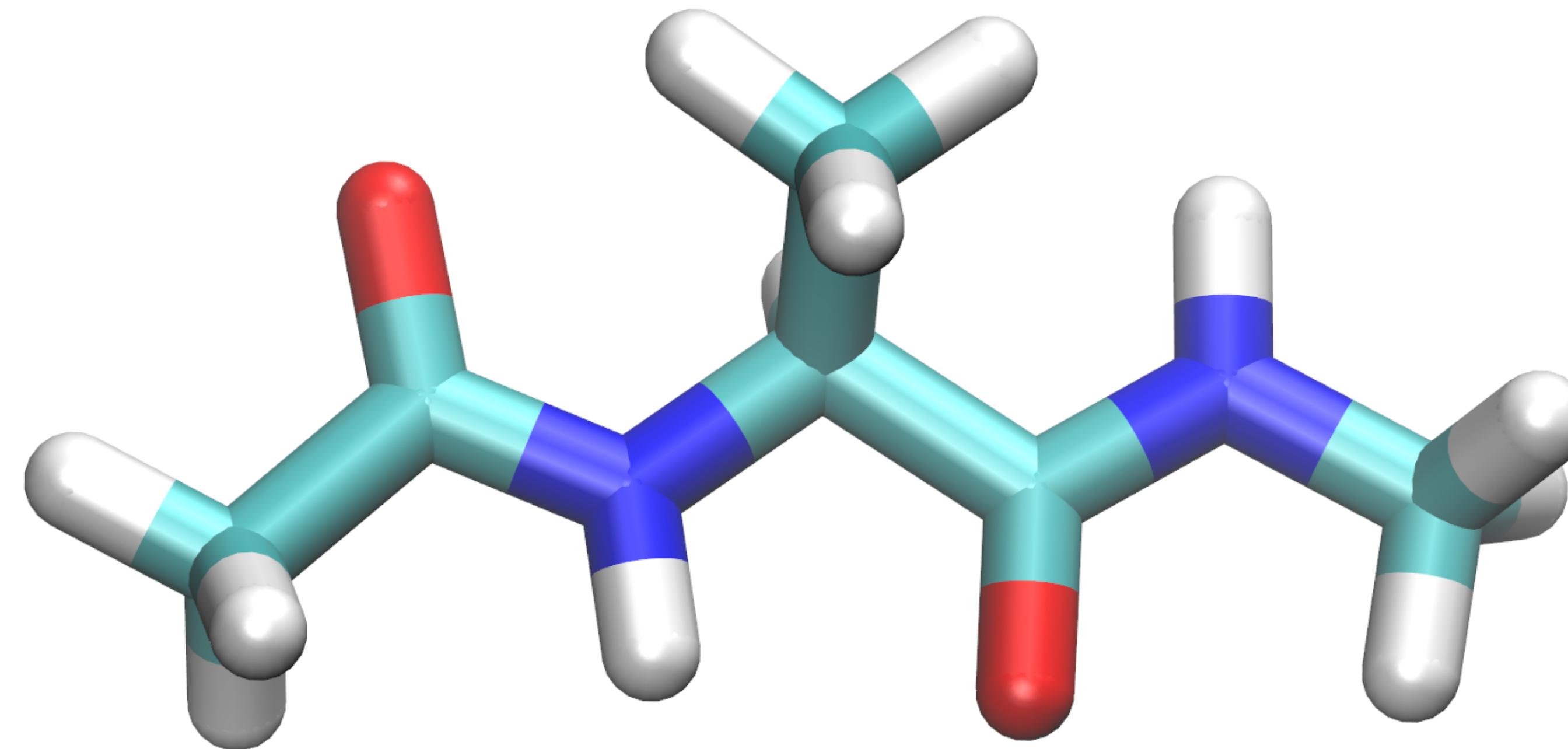
Fermi–Pasta–Ulam–Tsingou problem w/o nonlinearity

Learning the normal modes



Consistency check: neural nets can learn linear coordinate transformations

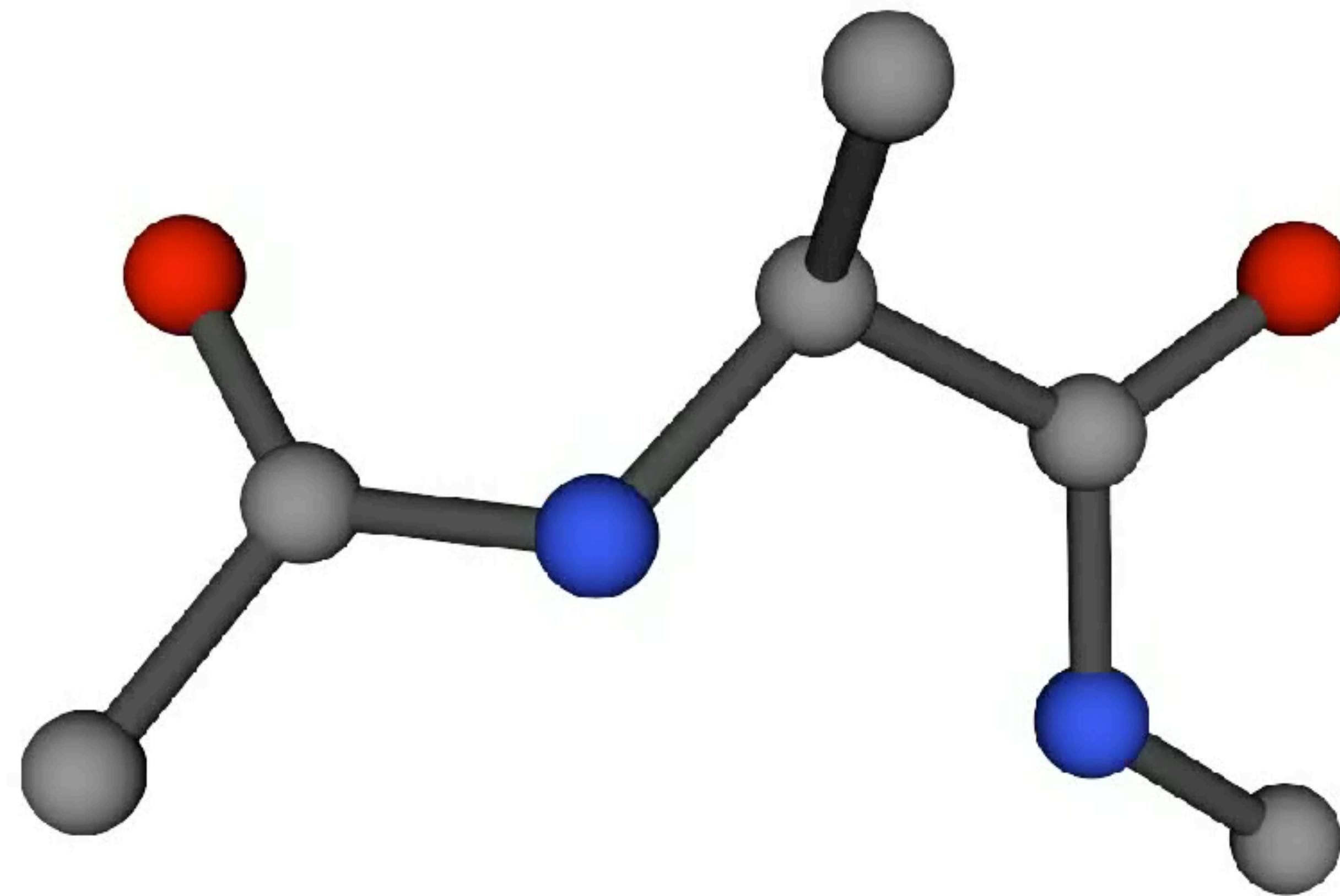
Alanine Dipeptide

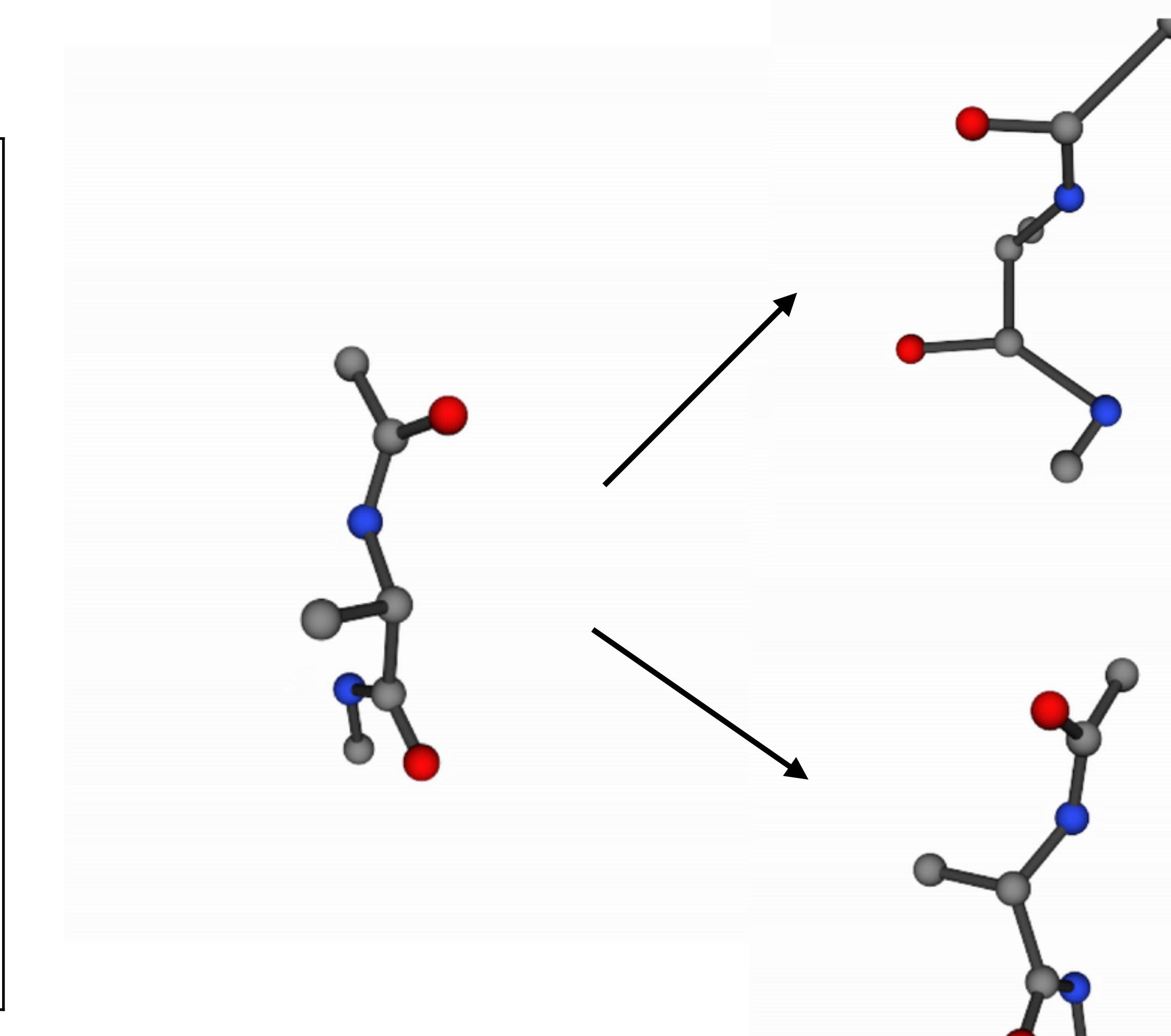
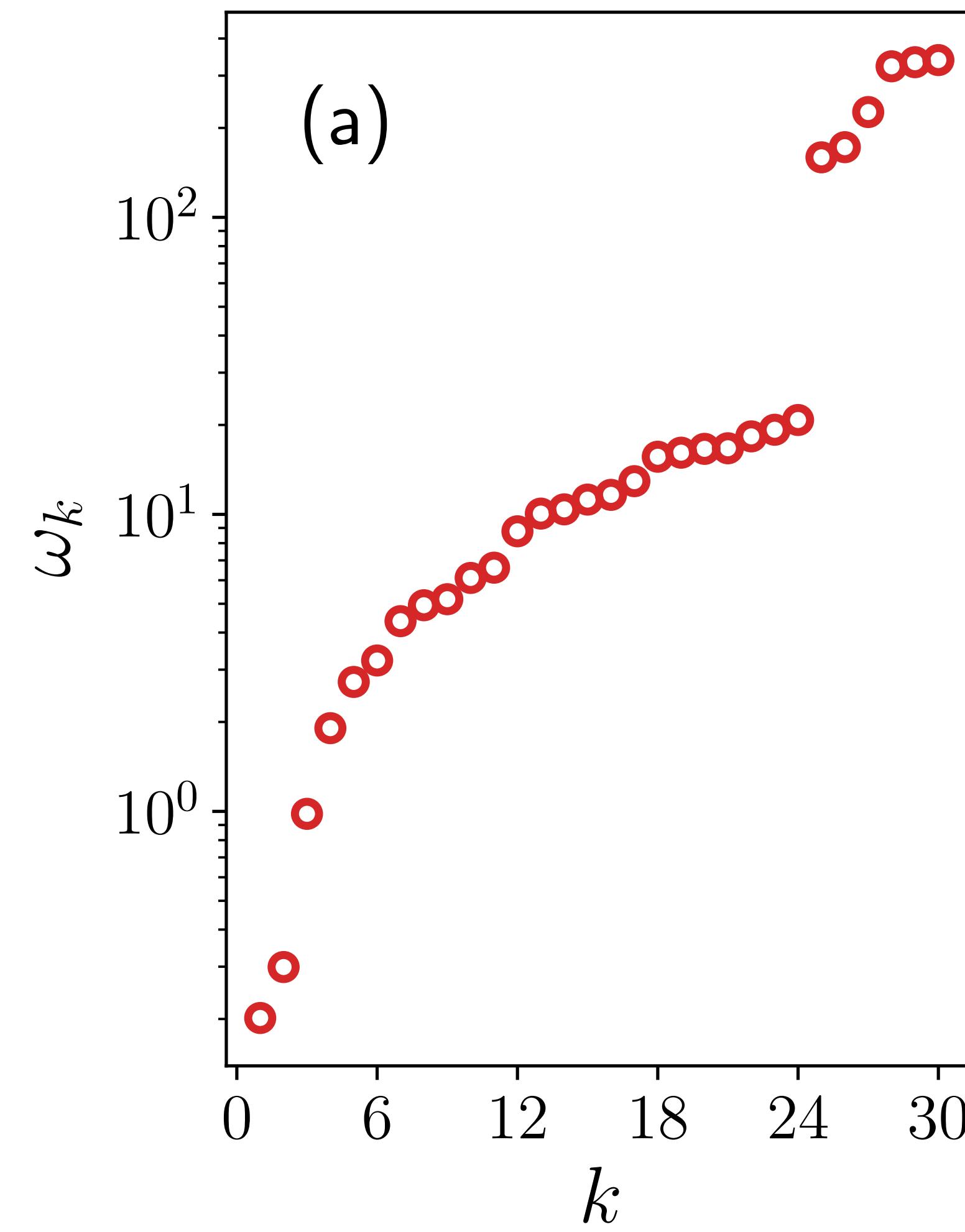


250 ns molecular dynamics simulation data at 300 K

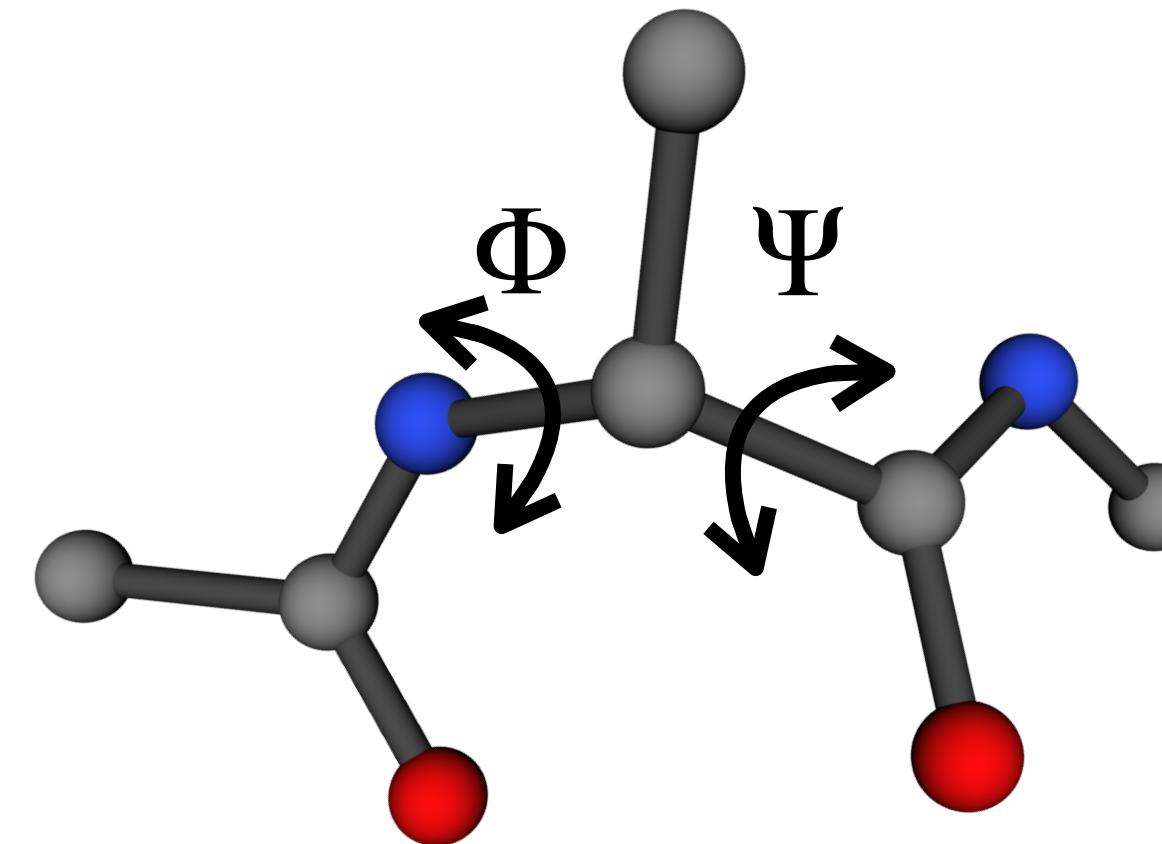
<https://markovmodel.github.io/mdshare/ALA2/#alanine-dipeptide>

More than 3 hours of video ...

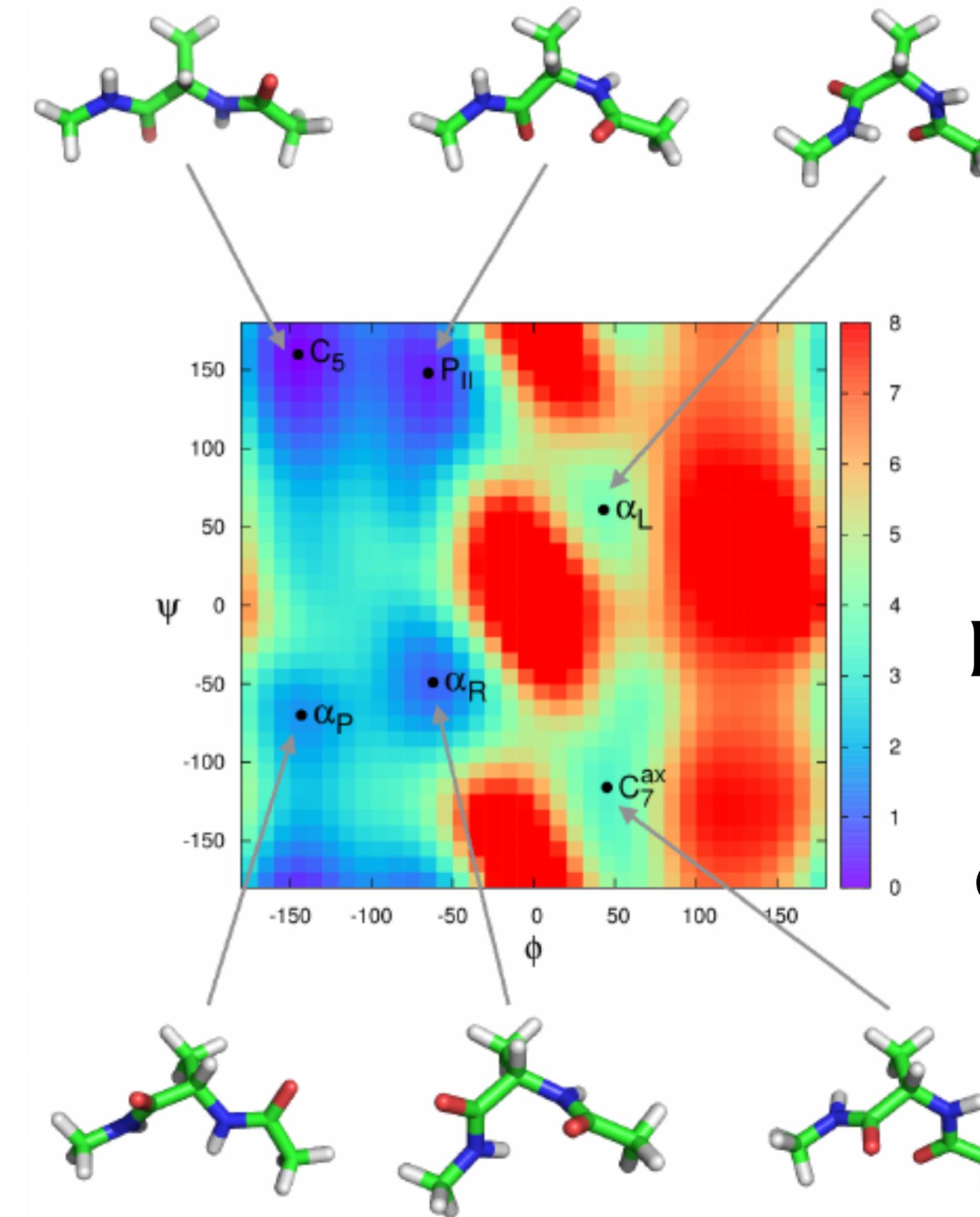
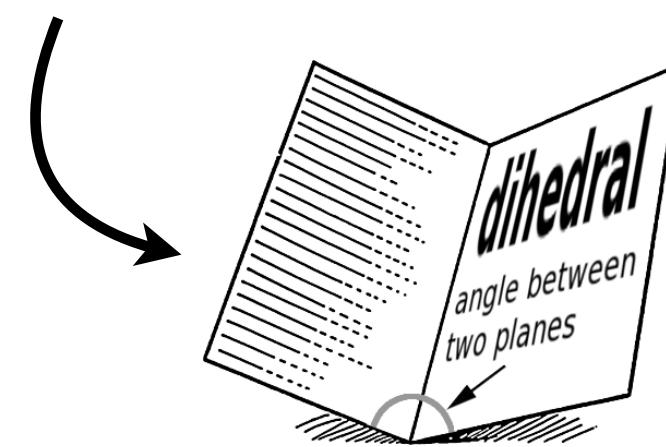




Neural canonical transformation identifies nonlinear slow modes!



slow motion of the
two torsion angles



Ramachandran
plot of stable
conformations

**Dimensional reduction to slow collective variables
useful for control, prediction, enhanced sampling...**

check the paper 1910.00024, PRX '20 for more examples & applications

MNIST handwritten digits



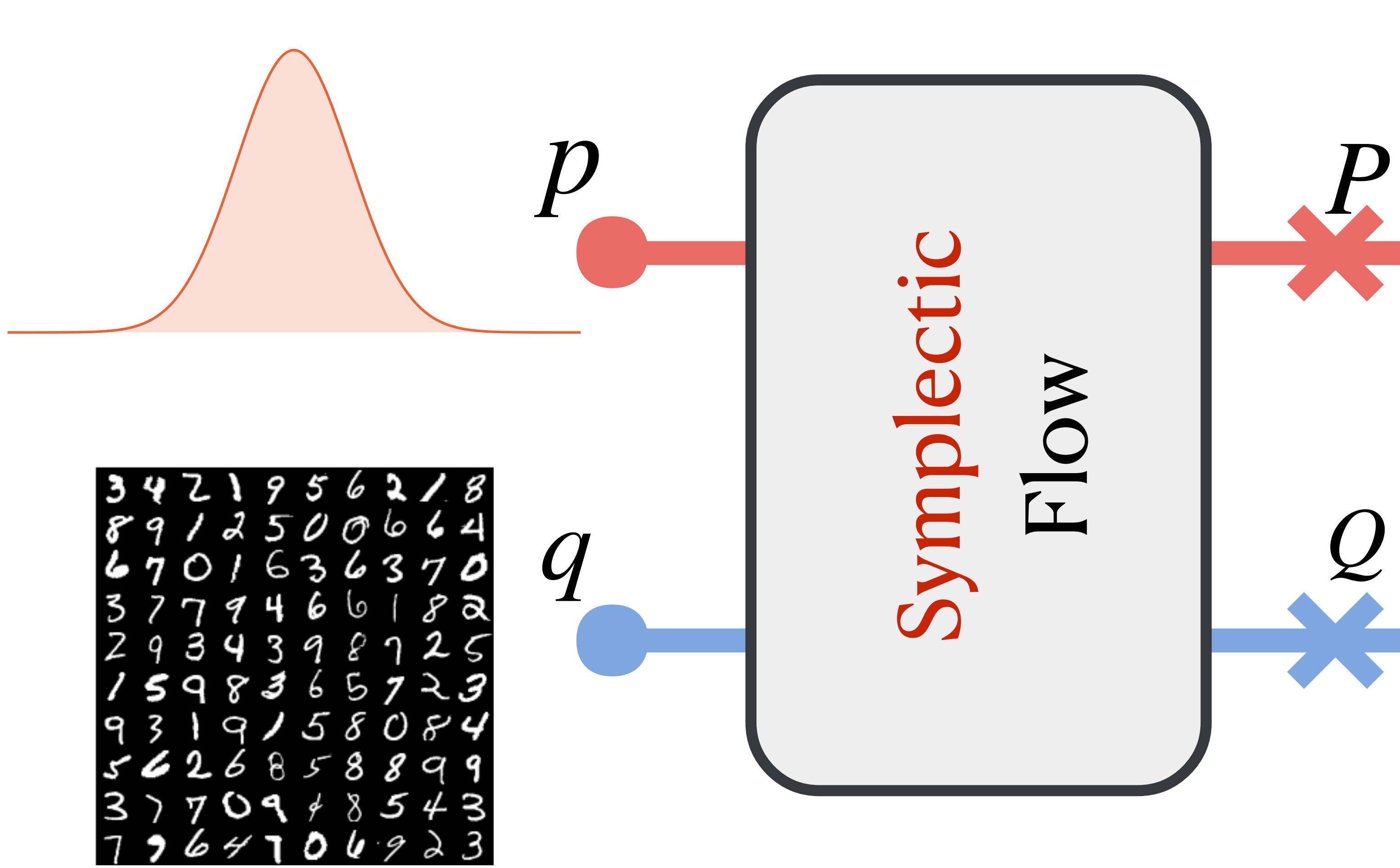
Data scientists:

“50,000 grayscale images
with 28x28 pixels”

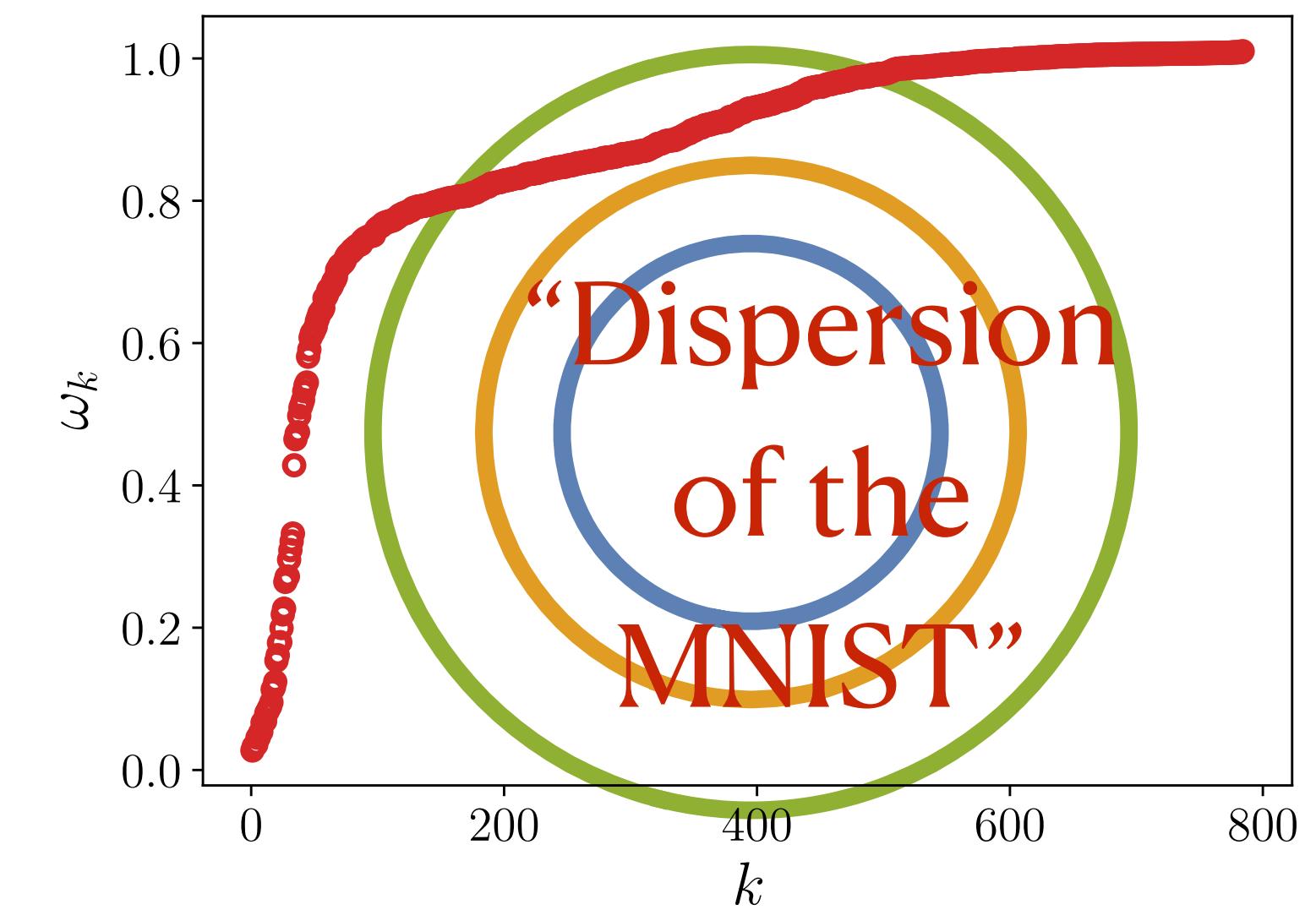
Physical Chemists:

“Stable conformations
of a molecule with 784
degrees of freedom”

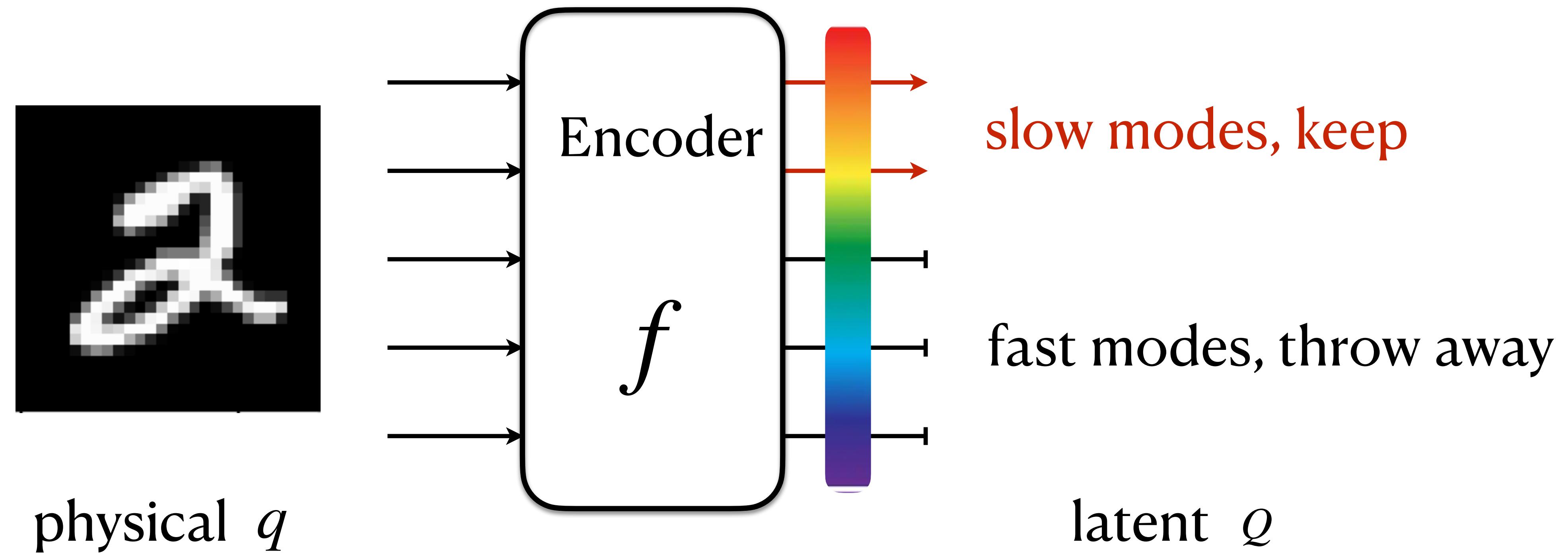
Learning slow variables of MNIST



$$K(P, Q) = \sum_k \frac{P_k^2 + \omega_k^2 Q_k^2}{2}$$



Conceptual compression



Compress by keeping slow collective variables

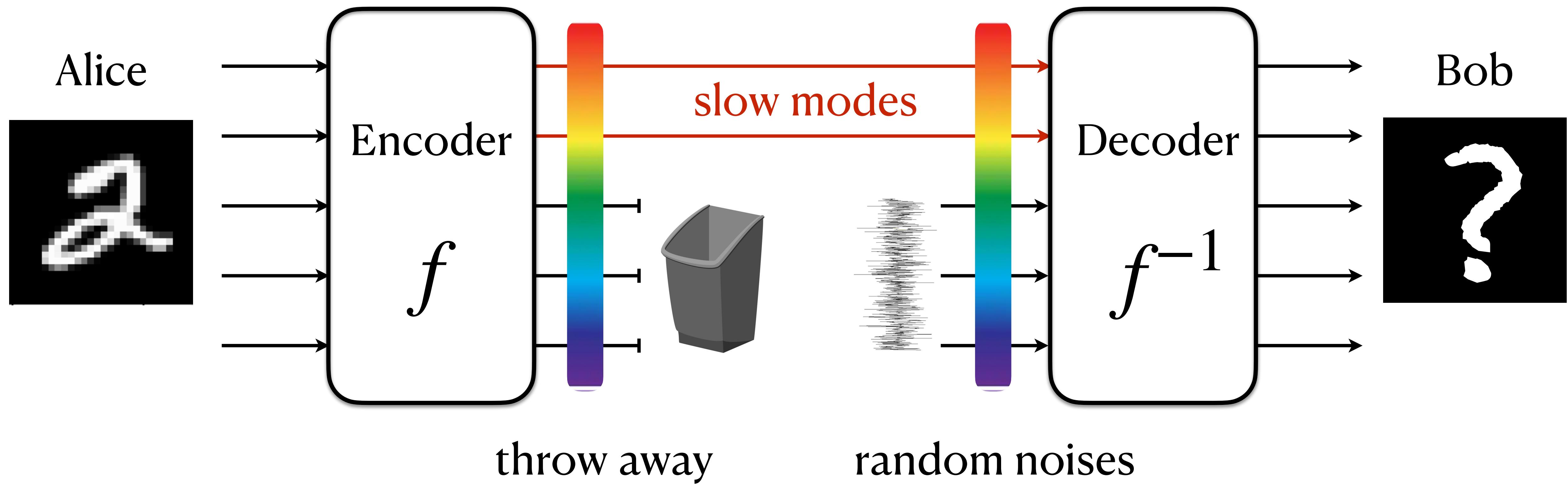
Kingma et al, 1312.6114

Gregor et al, 1604.08772

Dinh et al, 1605.08803

autoencoders/hierarchical network architecture/hyperbolic latent space...

Neural Alice-Bob game



Conceptual Compression of MNIST

Original

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

1/784 kept

2 3 5 1 0 2 3 9 0 3
5 9 3 2 2 5 9 9 3 3
4 9 6 3 5 4 8 7 0 5
2 5 3 2 0 9 9 2 4 1
7 8 0 2 3 2 9 3 5 9
8 3 4 9 6 0 8 3 2 5
9 2 5 4 8 9 3 5 9 6
3 4 9 4 7 9 8 7 2 1
3 1 6 3 4 5 9 5 5 3
4 3 4 6 0 0 0 0 0 4

Conceptual Compression of MNIST

Original

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

2/784 kept

6 9 3 4 8 3 3 1 0 3
3 9 0 0 3 6 9 7 7 8
6 2 3 4 9 6 0 3 4 5
2 5 3 3 5 4 6 6 7 9
4 2 2 0 4 0 1 4 6 6
0 1 2 9 2 0 2 3 3 9
2 0 3 3 0 2 8 0 3 8
8 3 7 9 2 4 2 8 6 9
5 1 5 3 5 2 9 9 4 0
8 3 1 3 4 3 3 8 9 9

Conceptual Compression of MNIST

Original

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

3/784 kept

6 9 6 4 7 5 5 9 4 3
4 9 0 3 0 8 8 9 7 0
2 0 6 5 4 3 0 5 3 0
9 6 5 3 3 3 0 8 4 5 9
4 3 2 3 9 8 9 7 5 4
6 0 5 9 3 9 7 9 4 8
5 3 3 8 4 5 6 6 3 6
2 9 7 5 4 3 2 8 3 7
3 9 0 7 7 0 4 3 5 2
3 0 5 6 3 2 1 2 3 7

Conceptual Compression of MNIST

Original

A 10x10 grid of handwritten digits, each digit being a 28x28 pixel image. The digits are arranged in a single row, with some digits appearing multiple times. The digits are rendered in white on a black background.

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

4/784 kept

A 10x10 grid of handwritten digits, similar to the original but with significantly fewer features. Only 4 out of the 784 features from the original dataset have been retained. These digits are rendered in white on a black background.

6 5 6 7 3 0 8 4 9 3
7 7 0 6 2 3 4 8 7 8
9 8 0 4 4 3 0 8 3 4
5 1 8 5 6 0 6 8 5 9
4 3 5 3 9 9 4 9 3 8
4 3 8 7 0 8 5 7 3 5
6 0 8 3 7 6 6 8 3 3
0 3 4 7 6 9 0 8 0 7
5 9 8 4 3 9 9 3 9 9
3 2 6 5 8 7 8 2 7

Conceptual Compression of MNIST

Original

A 10x10 grid of handwritten digits, each digit rendered in white against a black background. The digits are arranged in a single row.

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

5/784 kept

A 10x10 grid of handwritten digits, each digit rendered in white against a black background. The digits are arranged in a single row, showing only 5 of the 784 digits from the original set.

6 4 6 4 8 6 0 9 2 8
7 7 6 6 6 3 4 8 7 5
8 9 6 8 4 6 6 5 5 2
5 5 0 8 0 7 6 8 5 9
4 0 2 5 9 3 4 3 0 1
8 4 6 9 3 9 3 1 5 5
0 0 0 2 3 7 4 0 8 3 6
2 3 9 3 6 9 3 6 0 9
0 7 3 0 7 3 9 3 9 4
4 2 6 5 2 5 9 3 8 7

Conceptual Compression of MNIST

Original

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

10/784 kept

6 4 6 4 9 4 0 9 1 7
8 9 0 3 6 3 4 9 7 5
2 0 2 9 4 6 0 3 3 5
5 1 7 9 0 4 6 1 5 9
4 2 2 3 9 2 1 9 5 1
2 1 9 7 8 9 7 3 3 5
0 0 9 8 9 8 6 9 3 6
3 3 9 3 1 4 2 1 2 7
6 9 8 4 2 8 4 7 8 0
8 2 6 4 2 5 1 0 5 7

Conceptual Compression of MNIST

Original

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

15/784 kept

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 5 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

Conceptual compression of MNIST

Original

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 2 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

20/784 kept

6 7 6 4 3 5 0 9 1 7
3 9 0 6 3 3 4 9 7 5
2 0 2 5 4 3 6 0 5 5
5 1 9 8 0 4 6 1 8 9
4 2 2 3 9 2 1 9 3 1
2 1 6 9 3 9 7 1 7 5
0 0 3 8 9 8 0 8 5 6
3 8 9 3 1 4 2 1 2 7
6 9 8 4 3 5 4 7 9 0
8 2 6 5 2 5 1 0 8 7

“A Hamiltonian Extravaganza”

—Danilo J. Rezende@DeepMind

Symplectic ODE-Net, 1909.12077



SIEMENS

Hamiltonian Graph Networks with ODE Integrators, 1909.12790



Symplectic RNN, 1909.13334



Equivariant Hamiltonian Flows, 1909.13739



Hamiltonian Generative Network, 1909.13789



<http://tiny.cc/hgn>

Neural Canonical Transformation with Symplectic Flows, 1910.00024



See also Bondesan & Lamacraft, *Learning Symmetries of Classical Integrable Systems*, 1906.04645

Training: Monte Carlo Gradient Estimators

Review: 1906.10652

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} [f(x)]$$

Reinforcement learning

Variational inference

Variational Monte Carlo

Variational quantum algorithms

...

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} [f(x)] = \mathbb{E}_{x \sim p_{\theta}} [f(x) \nabla_{\theta} \ln p_{\theta}(x)]$$

Score function estimator (REINFORCE)

Pathwise estimator (Reparametrization trick) $x = g_{\theta}(z)$

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} [f(x)] = \mathbb{E}_{z \sim \mathcal{N}(z)} [\nabla_{\theta} f(g_{\theta}(z))]$$

10.1 Guidance in Choosing Gradient Estimators

With so many competing approaches, we offer our rules of thumb in choosing an estimator, which follow the intuition we developed throughout the paper:

- If our estimation problem involves continuous functions and measures that are continuous in the domain, then using the pathwise estimator is a good default. It is relatively easy to implement and a default implementation, one without other variance reduction, will typically have variance that is low enough so as not to interfere with the optimisation.
- If the cost function is not differentiable or a black-box function then the score-function or the measure-valued gradients are available. If the number of parameters is low, then the measure-valued gradient will typically have lower variance and would be preferred. But if we have a high-dimensional parameter set, then the score function estimator should be used.
- If we have no control over the number of times we can evaluate a black-box cost function, effectively only allowing a single evaluation of it, then the score function is the only estimator of the three we reviewed that is applicable.
- The score function estimator should, by default, always be implemented with at least a basic variance reduction. The simplest option is to use a baseline control variate estimated with a running average of the cost value.
- When using the score-function estimator, some attention should be paid to the dynamic range of the cost function and its variance, and to find ways to keep its value bounded within a reasonable range, e.g., transforming the cost so that it is zero mean, or using a baseline.
- For all estimators, track the variance of the gradients if possible and address high variance by using a larger number of samples from the measure, decreasing the learning rate, or clipping the gradient values. It may also be useful to restrict the range of some parameters to avoid extreme values, e.g., by clipping them to a desired interval.
- The measure-valued gradient should be used with some coupling method for variance reduction. Coupling strategies that exploit relationships between the positive and negative components of the density decomposition, and which have shared sampling paths, are known for the commonly-used distributions.
- If we have several unbiased gradient estimators, a convex combination of them might have lower variance than any of the individual estimators.
- If the measure is discrete on its domain then the score-function or measure-valued gradient are available. The choice will again depend on the dimensionality of the parameter space.
- In all cases, we strongly recommend having a broad set of tests to verify the unbiasedness of the gradient estimator when implemented.

Mohamed et al, 1906.10652

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}} [f(x)]$$

When to use which ?

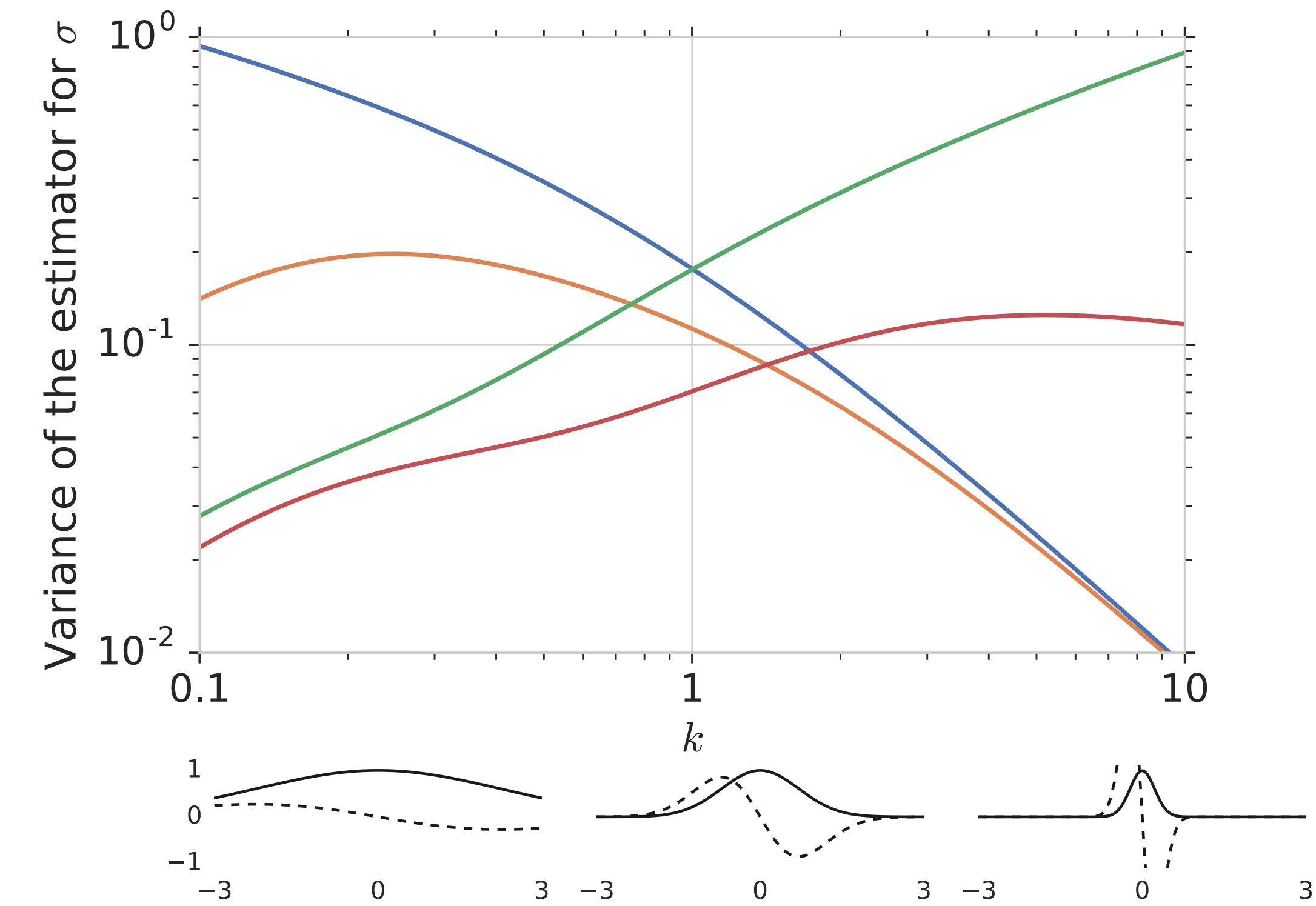
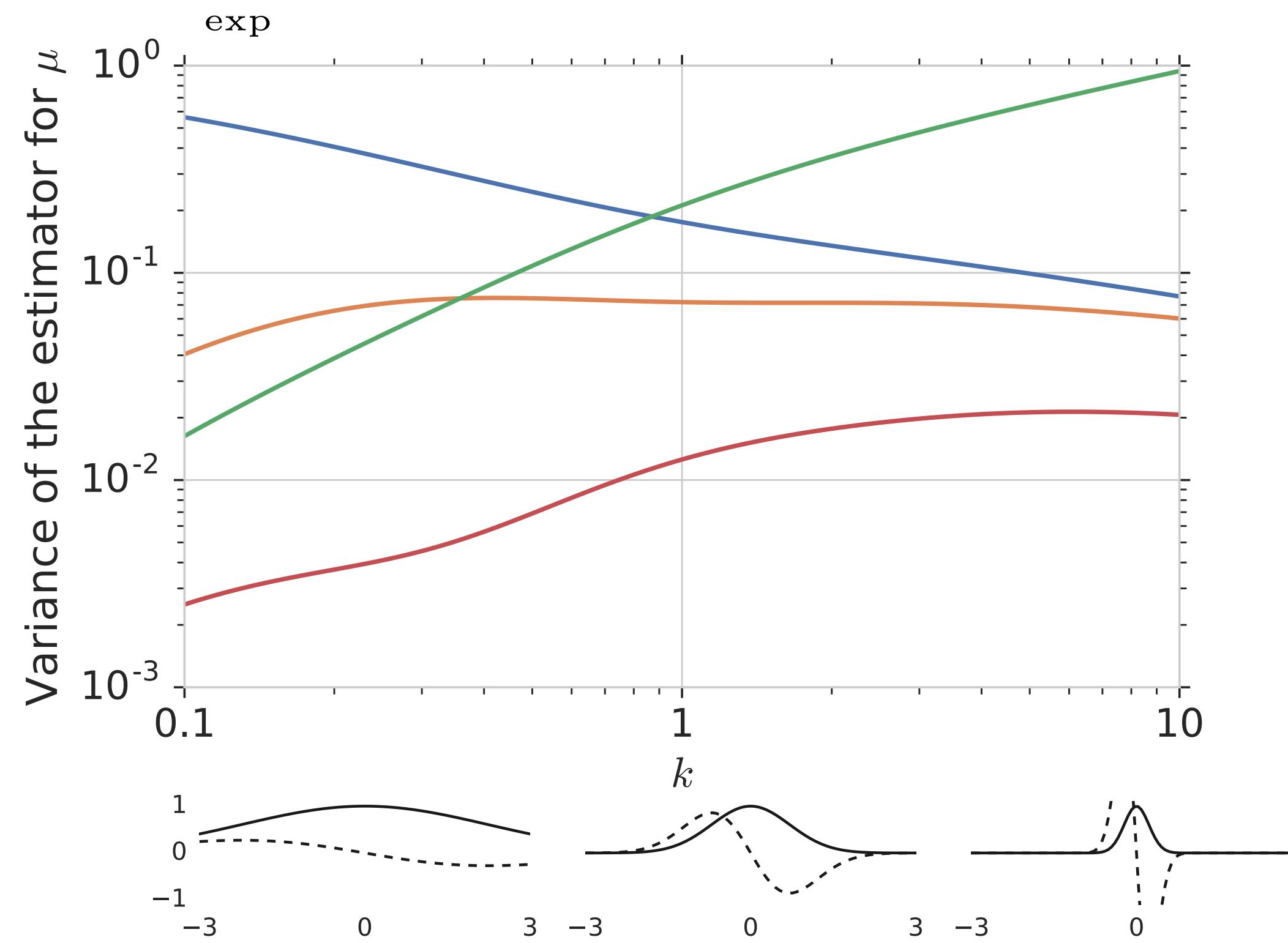
More discussions

Roeder et al, 1703.09194

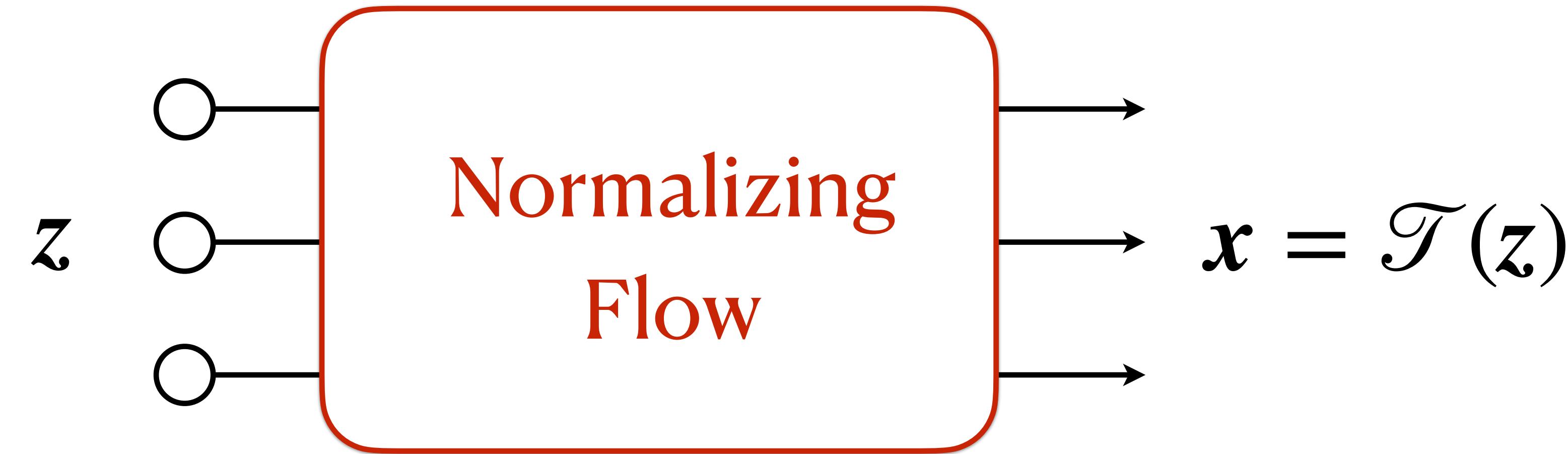
Vaitl et al 2206.09016, 2207.08219

$$\eta = \nabla_{\theta} \int \mathcal{N}(x|\mu, \sigma^2) f(x; k) dx; \quad \theta \in \{\mu, \sigma\}$$

— Score function — Score function + variance reduction — Pathwise — Measure-valued + variance reduction
— Value of the cost --- Derivative of the cost



Symmetries



Invariance

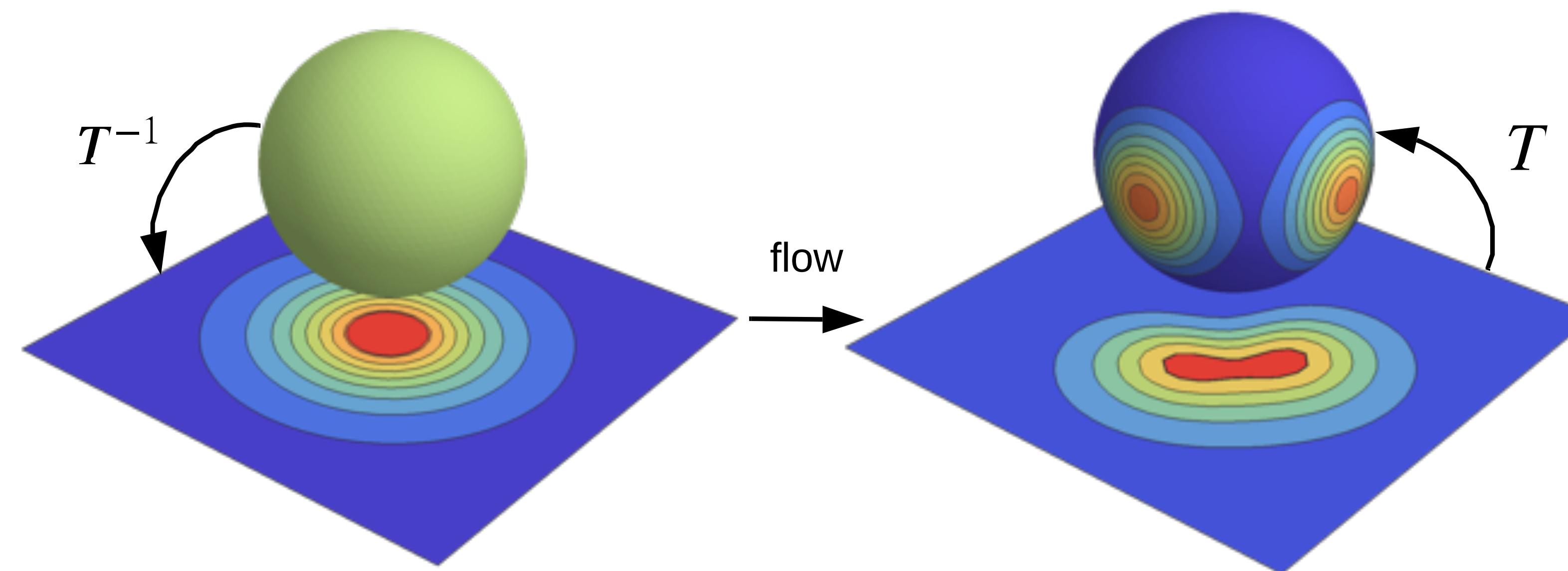
$$\rho(gx) = \rho(x)$$

Equivariance

$$\mathcal{T}(gz) = g\mathcal{T}(z)$$

Spatial symmetries, permutation symmetries, gauge symmetries...

Flow on manifolds

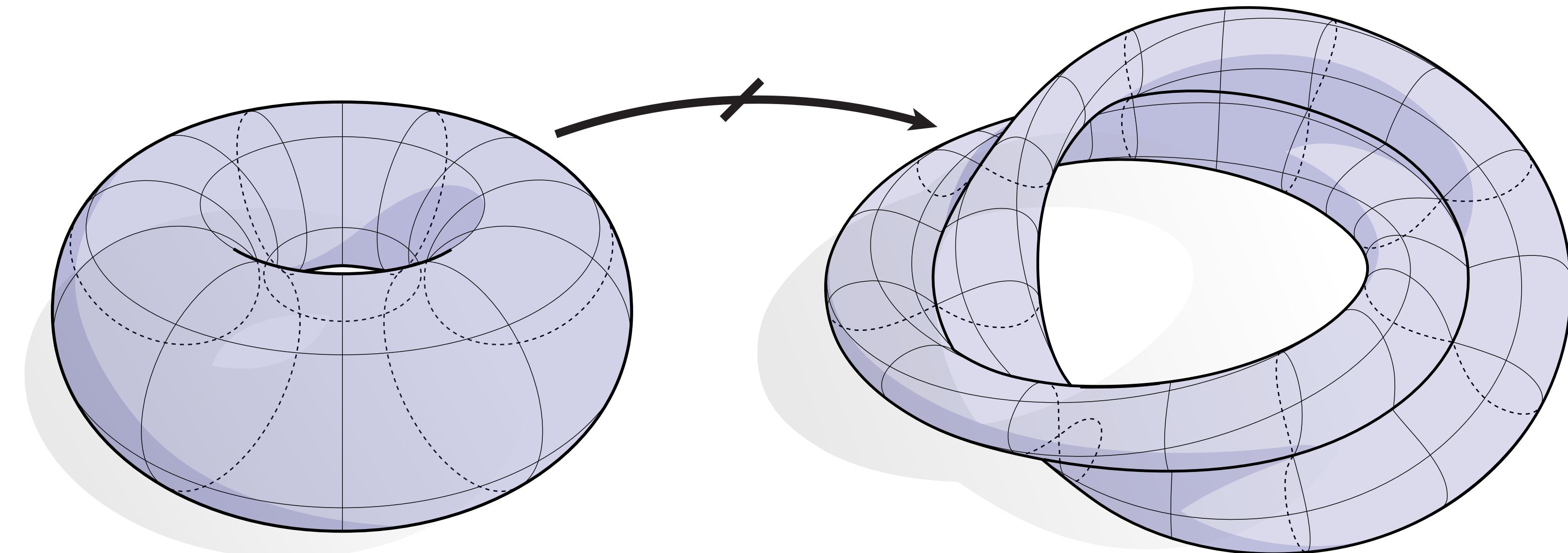


Periodic variables, gauge fields, ...

Gemici et al 1611.02304, Rezende et al, 2002.02428, Boyda et al, 2008.05456

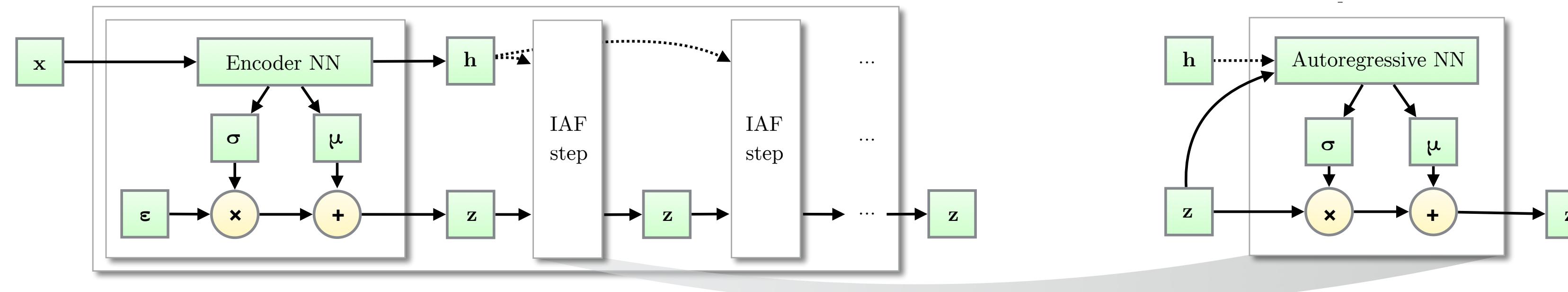
Neural ODE on manifolds, Falorsi et al, 2006.06663, Lou et al, 2006.10254, Mathieu et al, 2006.10605

Obstructions

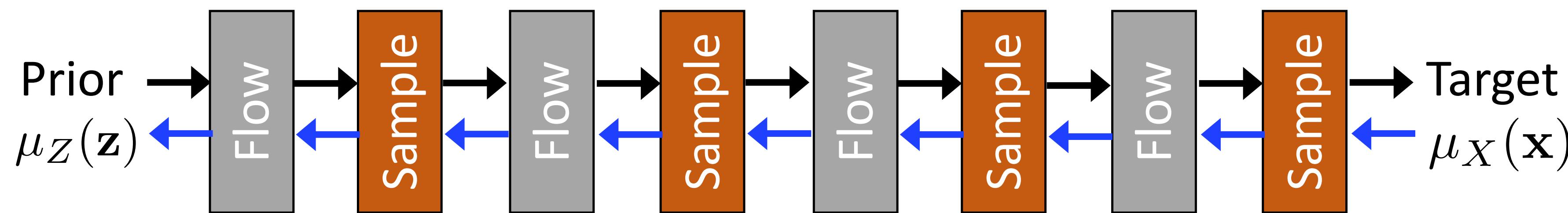


Dupont et al 1904.01681, Cornish et al, 1909.13833, Zhang et al, 1907.12998, Zhong et al, 2006.00392...

Mix with other approaches



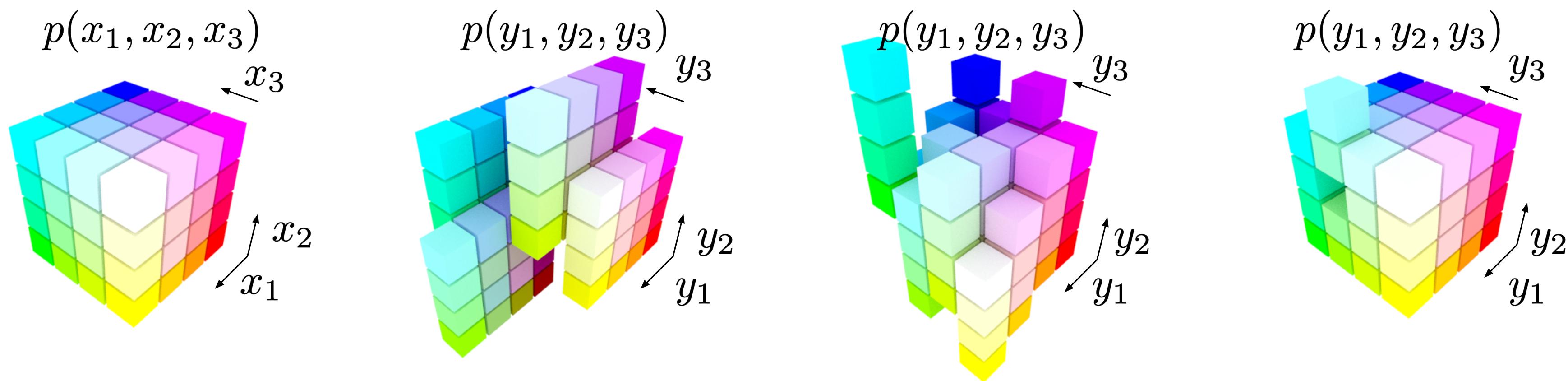
Kingma et al, 1606.04934, ...



Levy et al, 1711.09268, Wu et al 2002.06707, ...

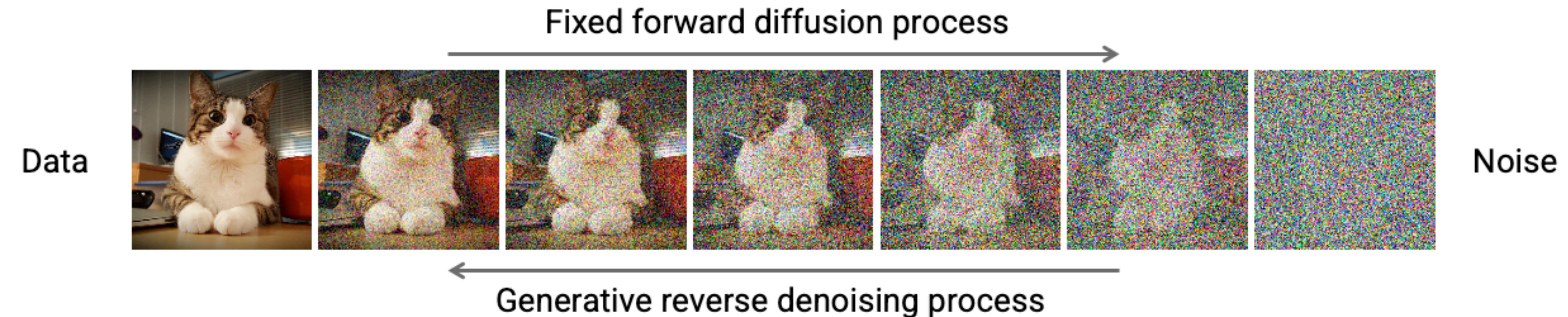
Discrete flows

$$p(\mathbf{x}) = p(\mathbf{y} = \mathcal{T}(\mathbf{x}))$$



Tran et al, 1905.10347, Hoogeboom et al, 1905.07376, van den Berg 2006.12459

Diffusion models



Score matching

Minimizing Fisher divergence by-passes intractable partition function

$$\mathbb{F}(\pi \parallel p) = \int d\mathbf{x} \, \pi(\mathbf{x}) \left| \nabla \ln \pi(\mathbf{x}) - \nabla \ln p(\mathbf{x}) \right|^2$$

Draw samples with Langevin dynamics

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{\epsilon}{2} \nabla \ln p(\mathbf{x}_{t-1}) + \sqrt{\epsilon} \mathcal{N}(0, I)$$

$$\mathbf{x}_T \sim p(\mathbf{x}_T) \quad \text{for } \epsilon \rightarrow 0, \quad T \rightarrow \infty$$

From score matching to diffusion model

Denoising score matching Vincent 2011

$$\frac{1}{2} \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}}|\mathbf{x}) p_{\text{data}}(\mathbf{x})} [\|\mathbf{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})\|_2^2].$$

Diffusion generative model

Sohl-Dickstein et al, 1503.03585

$$\mathbb{E} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_q \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

Denoising diffusion probabilistic model

Song et al, 1907.05600, Ho et al, 2006.11239

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t|\mathbf{x}_0)} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t|\mathbf{x}_0)\|_2^2$$

diffusion time t data sample \mathbf{x}_0 diffused data sample \mathbf{x}_t neural network score of diffused data sample

A tale of three equations

Langevin equation (SDE)

$$q(x_{t+dt} | x_t) = \mathcal{N}(x_t + f dt, 2 T dt I) \quad \text{or} \quad x_{t+dt} - x_t = f dt + \sqrt{2 T dt} \mathcal{N}(0, 1)$$

Fokker-Planck equation (PDE)

$$\frac{\partial p(x, t)}{\partial t} + \nabla \cdot [p(x, t) f] - T \nabla^2 p(x, t) = 0$$

“Particle method” (ODE)

$$\frac{dx}{dt} = f - T \nabla \ln p(x, t)$$

(Another way to reverse the diffusion is via the reverse-time SDE Anderson 1982)

Maoutsou et al, 2006.00702
Song et al, 2011.13456

$$\mathcal{P}(\vec{x}, t) = \int d^3\vec{x}' \left(\frac{1}{4\pi D\epsilon} \right)^{3/2} \exp \left[-\frac{(\vec{x} - \vec{x}' - \epsilon \vec{v}(\vec{x}'))^2}{4D\epsilon} \right] \mathcal{P}(\vec{x}', t - \epsilon), \quad (9.18)$$

and simplified by the change of variables,

$$\begin{aligned} \vec{y} &= \vec{x}' + \epsilon \vec{v}(\vec{x}') - \vec{x} \implies \\ d^3\vec{y} &= d^3\vec{x}' (1 + \epsilon \nabla \cdot \vec{v}(\vec{x}')) = d^3\vec{x}' (1 + \epsilon \nabla \cdot \vec{v}(\vec{x}) + \mathcal{O}(\epsilon^2)). \end{aligned} \quad (9.19)$$

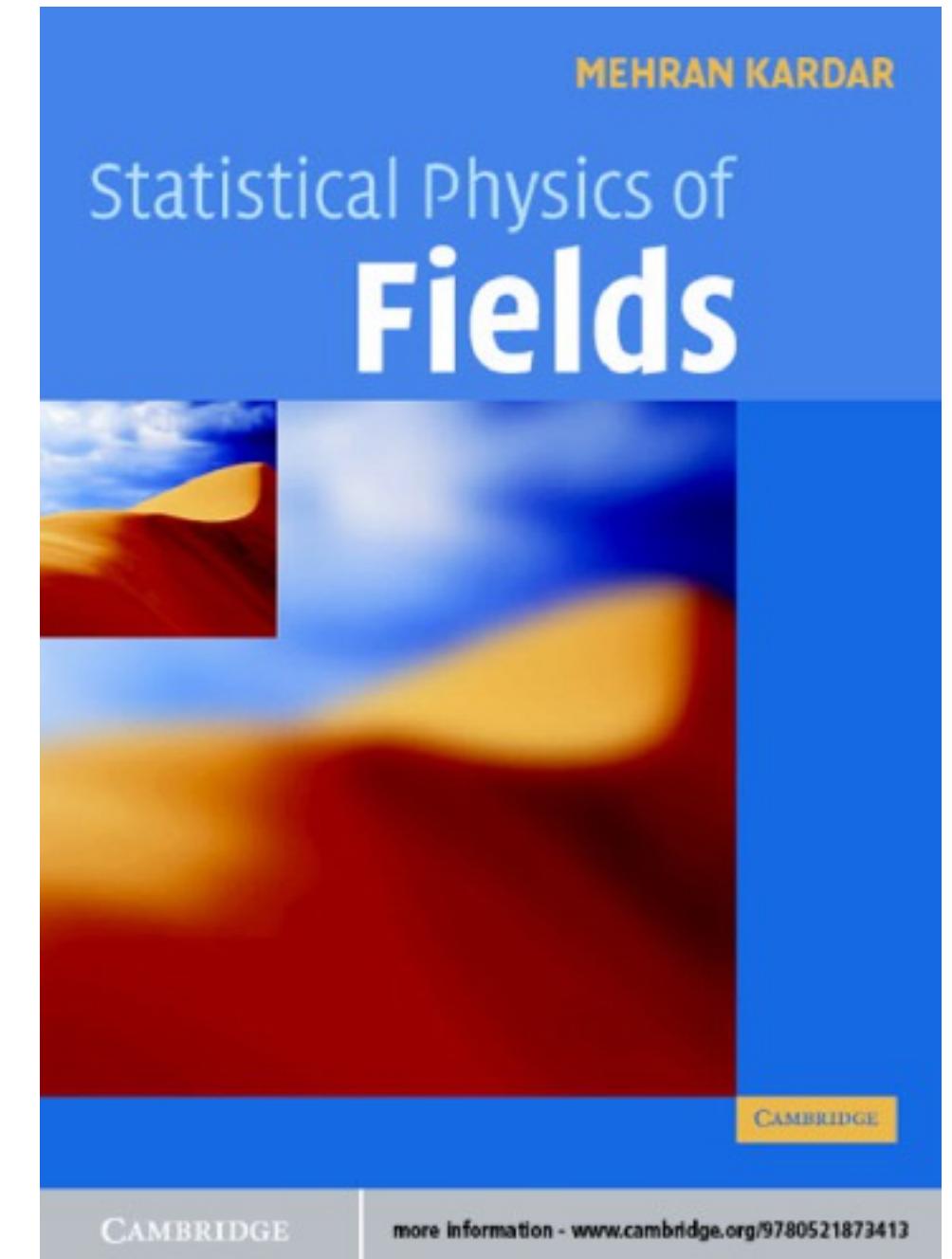
Keeping only terms at order of ϵ , we obtain

$$\begin{aligned} \mathcal{P}(\vec{x}, t) &= [1 - \epsilon \nabla \cdot \vec{v}(\vec{x})] \int d^3\vec{y} \left(\frac{1}{4\pi D\epsilon} \right)^{3/2} e^{-\frac{y^2}{4D\epsilon}} \mathcal{P}(\vec{x} + \vec{y} - \epsilon \vec{v}(\vec{x}), t - \epsilon) \\ &= [1 - \epsilon \nabla \cdot \vec{v}(\vec{x})] \int d^3\vec{y} \left(\frac{1}{4\pi D\epsilon} \right)^{3/2} e^{-\frac{y^2}{4D\epsilon}} \\ &\times \left[\mathcal{P}(\vec{x}, t) + (\vec{y} - \epsilon \vec{v}(\vec{x})) \cdot \nabla \mathcal{P} + \frac{y_i y_j - 2\epsilon y_i v_j + \epsilon^2 v_i v_j}{2} \nabla_i \nabla_j \mathcal{P} - \epsilon \frac{\partial \mathcal{P}}{\partial t} + \mathcal{O}(\epsilon^2) \right] \\ &= [1 - \epsilon \nabla \cdot \vec{v}(\vec{x})] \left[\mathcal{P} - \epsilon \vec{v} \cdot \nabla + \epsilon D \nabla^2 \mathcal{P} - \epsilon \frac{\partial \mathcal{P}}{\partial t} + \mathcal{O}(\epsilon^2) \right]. \end{aligned} \quad (9.20)$$

Equating terms at order of ϵ leads to the *Fokker–Planck equation*,

$$\frac{\partial \mathcal{P}}{\partial t} + \nabla \cdot \vec{J} = 0, \quad \text{with} \quad \vec{J} = \vec{v} \mathcal{P} - D \nabla \mathcal{P}. \quad (9.21)$$

from Langevin
to Fokker-Planck



Lessons from diffusion models

Continuous normalizing flow has great potential: diffusion model is an “existence proof”

Going beyond maximum likelihood estimation training (even if we can)

Break the loss into small pieces, sample them (kind of layer-wise training)

[https://blog.alexalemi.com/
diffusion.html](https://blog.alexalemi.com/diffusion.html)

The conditional trick (originated from denoising score matching Vincent 2011)

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t | \mathbf{x}_0)} \| \mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x}_0) \|_2^2$$

diffusion time t data sample \mathbf{x}_0 diffused data sample \mathbf{x}_t neural network score of diffused data sample

Claim:

$$\mathbb{E}_{x \sim q(x)} |s_\theta(x) - \nabla_x \ln q(x)|^2 = \mathbb{E}_{x_0 \sim q_0(x_0)} \mathbb{E}_{x \sim q(x|x_0)} |s_\theta(x) - \nabla_x \ln q(x|x_0)|^2 + \text{const.}$$



$$q(x) = \int q(x|x_0) q_0(x_0) dx_0$$

Independent
of θ

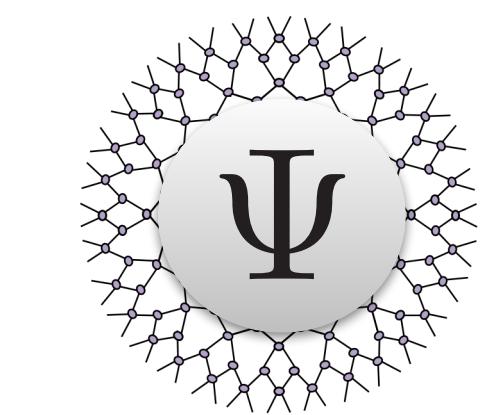
Proof:

$$\mathbb{E}_{x_0 \sim q_0(x_0)} \mathbb{E}_{x \sim q(x|x_0)} |s|^2 = \int dx_0 \int dx q_0(x_0) q(x|x_0) |s|^2 = \int dx q(x) |s|^2 = \mathbb{E}_{x \sim q(x)} |s|^2$$

$$\begin{aligned} \mathbb{E}_{x_0 \sim q_0(x_0)} \mathbb{E}_{x \sim q(x|x_0)} [s \cdot \nabla \ln q(x|x_0)] &= \int dx_0 \int dx q_0(x_0) q(x|x_0) \frac{s \cdot \nabla q(x|x_0)}{q(x|x_0)} \\ &= \int dx_0 \int dx q_0(x_0) s \cdot \nabla q(x|x_0) \\ &= \int dx s \cdot \nabla q(x) = \mathbb{E}_{x \sim q(x)} [s \cdot \nabla \ln q(x)] \end{aligned}$$

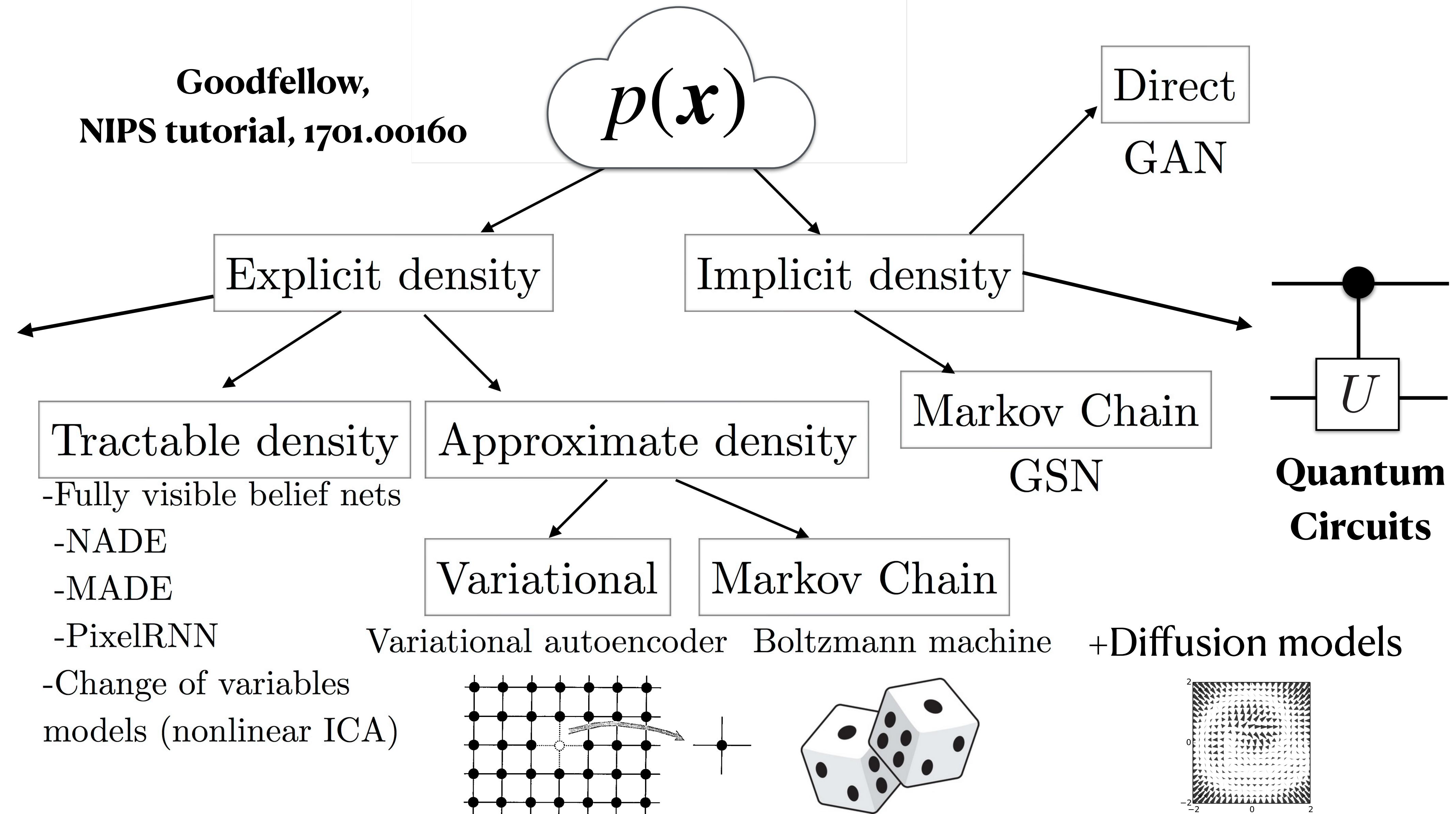
要让 $s_\theta(\mathbf{x}_t, t)$ 等于 $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)$ 的加权平均 【即 $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ 】 只需要最小化 $\|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)\|^2$ 的加权平均 <https://spaces.ac.cn/archives/9209>

Generative models and their physics genes



**Tensor
Networks**

**Goodfellow,
NIPS tutorial, 1701.00160**



GAN

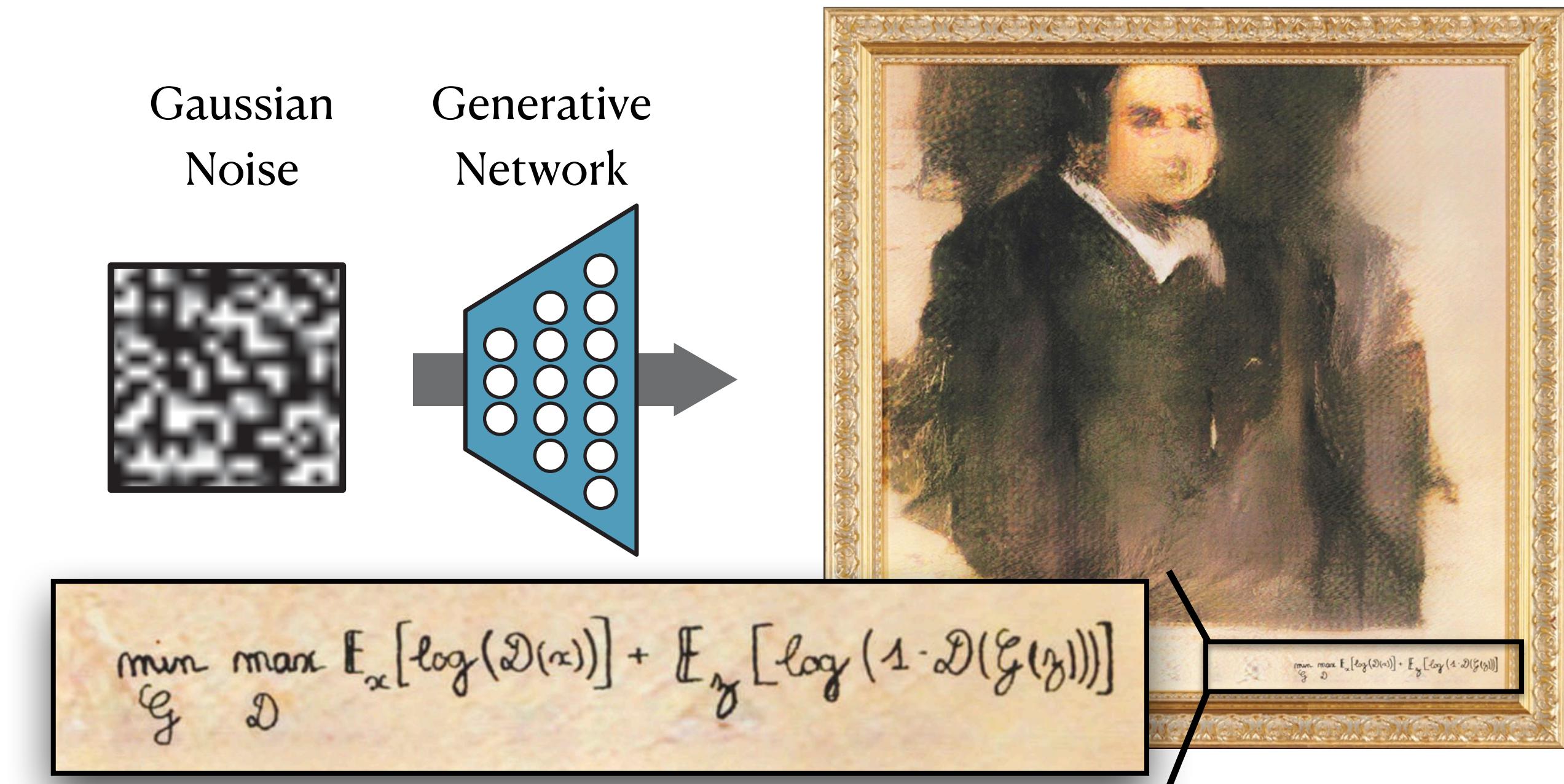
Likelihood free simulator

Prone to mode collapse

More tricky to train than others

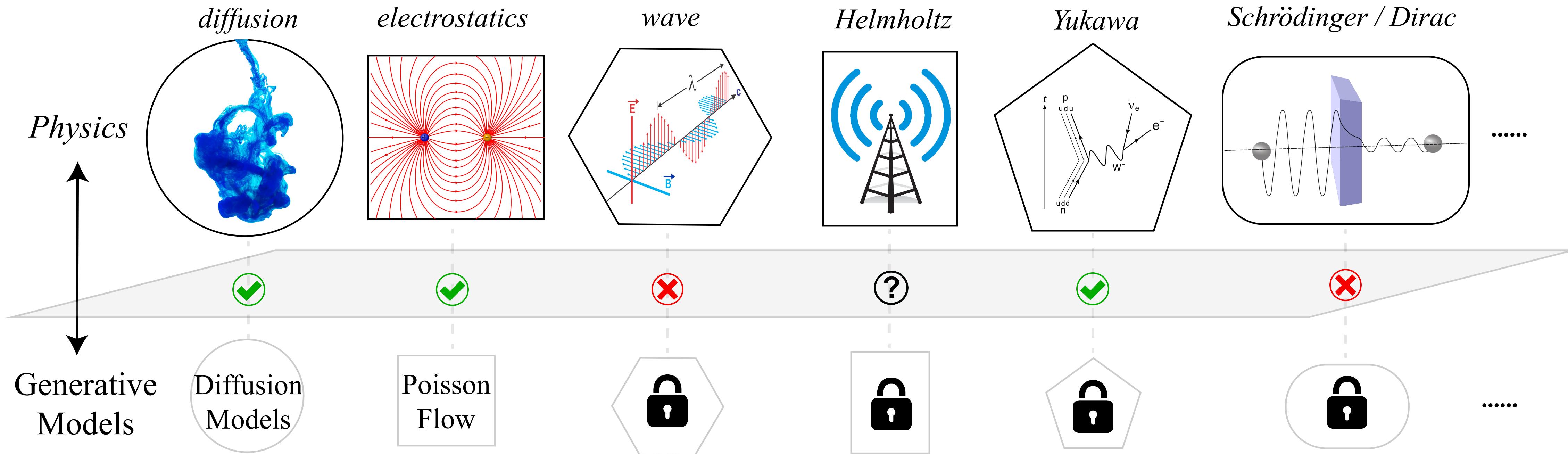
Performance have been surpassed by diffusion models

I found GAN to be less useful for serious scientific applications



<https://www.christies.com/Features/A-collaboration-between-two-artists-one-human-one-a-machine-9332-1.aspx>

Physical processes for generative modeling



Liu et al, 2304.02637