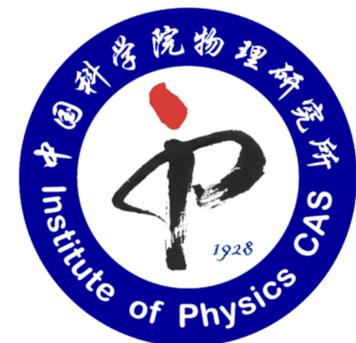


用生成模型解决物理问题

Lei Wang (王磊)

Institute of Physics, CAS

<https://wangleiphy.github.io>



Outline

- ① 物理学家眼中的生成模型
- ② 生成模型四问: 哪种最好? 多少先验? 能推理吗? 懂不懂物理?
- ③ 物质科学应用举例: 晶体材料设计、变分自由能计算

What are the limits of classical systems?

Classical Turing Machines can do much more than we previously thought

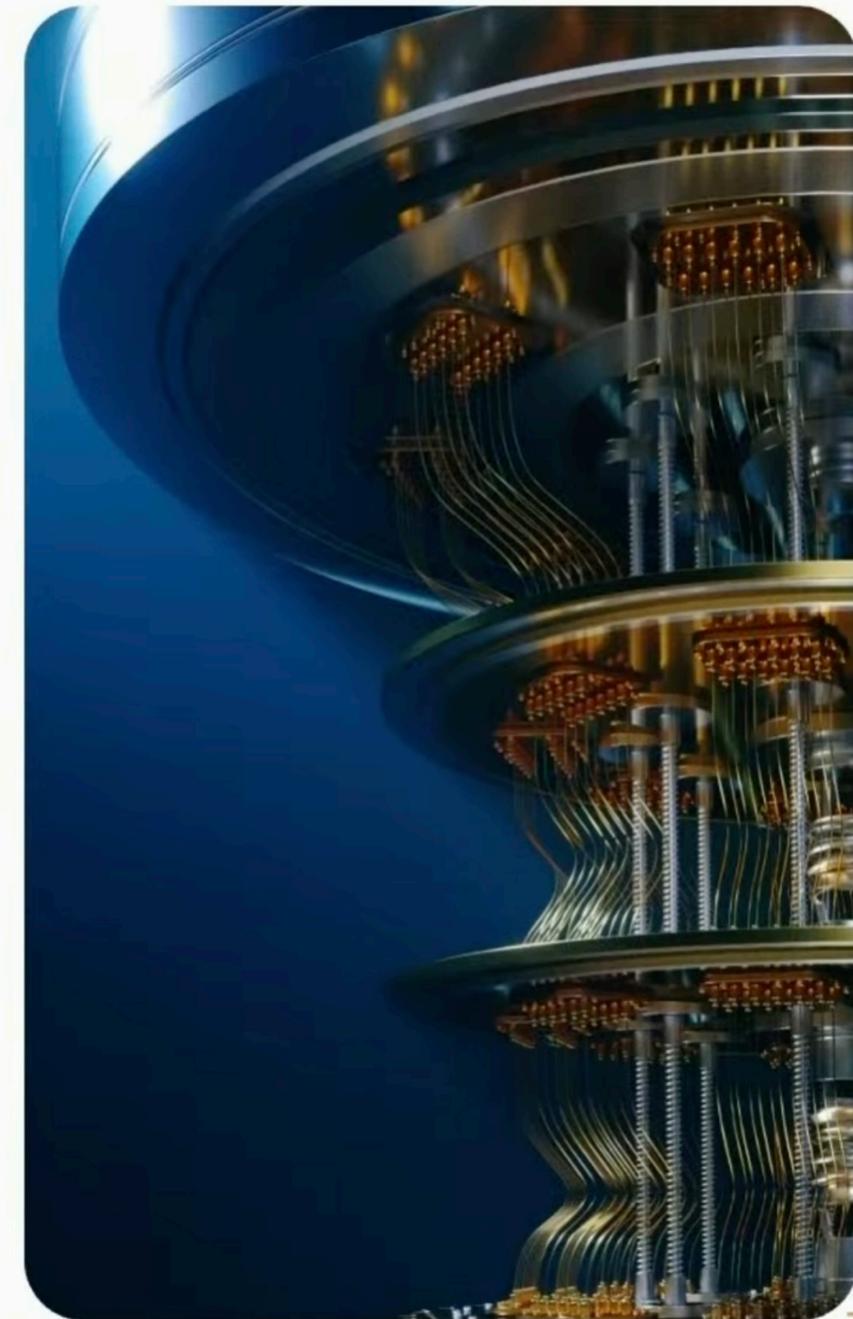
By doing a massive amount of pre-compute upfront to develop a good model

Then use the model to efficiently explore a solution space in polynomial time

My Proposed Conjecture:

“Any pattern that can be generated or found in nature can be efficiently discovered and modelled by a classical learning algorithm”

If it turns out that classical systems can model certain types of quantum systems, it could potentially have big implications for complexity theory including $P=NP$, and maybe even fundamental physics!



THE
NOBEL
PRIZE

Demis Hassabis, Nobel prize lecture, 2024 Dec

ChatGPT: Optimizing Language Models for Dialogue
November 30, 2022 — Announcements, Research

DALL·E API Now Available in Public Beta
November 3, 2022 — Announcements, API

DALL·E Now Available Without Waitlist
September 28, 2022 — Announcements

Introducing Whisper
September 21, 2022 — Research

DALL·E: Introducing Outpainting
August 31, 2022 — Announcements

Our Approach to Alignment Research
August 24, 2022 — Research

New and Improved Content Moderation Tooling
August 10, 2022 — Announcements

DALL·E Now Available in Beta
July 20, 2022 — Announcements

OpenAI Technical Goals
June 20, 2016 — Announcements

Generative Models
June 16, 2016 — Research, Milestones

Team Update
May 25, 2016 — Announcements

OpenAI Gym Beta
April 27, 2016 — Research

Welcome, Pieter and Shivon!
April 26, 2016 — Announcements

Team++
March 31, 2016 — Announcements

Introducing OpenAI
December 11, 2015 — Announcements

Generative Pretrained **T**ransformer
 $\text{text} \sim p(\text{text} | \text{prompt})$

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

Probabilistic modeling with generative AI

$$p(\mathbf{X})$$

pixels, words, atoms, ...

How to **express, learn, and sample from** a high-dimensional probability distribution?



DaLL-E

```
ChatGPT 4o >
Example using PySR:
python
# Install PySR (if not installed)
# pip install pysr

import numpy as np
from pysr import PySRRegressor

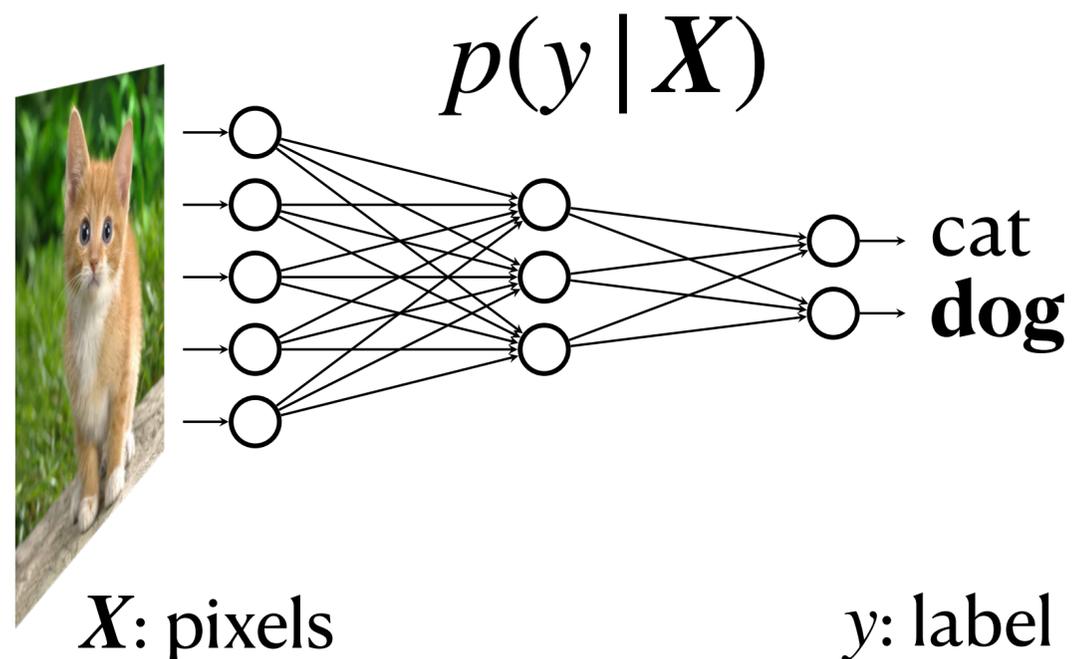
# Step 1: Generate data from the neural network
def neural_network(x, y):
    # Example neural network function, replace
    u = np.sin(x) + 0.5 * y
    v = np.cos(y) + 0.2 * x
    return u, v
```

ChatGPT



AlphaFold3

Discriminative AI is not enough



$$\nabla_{\text{pixels}} p(\text{dog} | \text{pixels})$$



Bayes rule

$$p(X | y) \propto p(X)p(y | X)$$

posterior

prior

likelihood

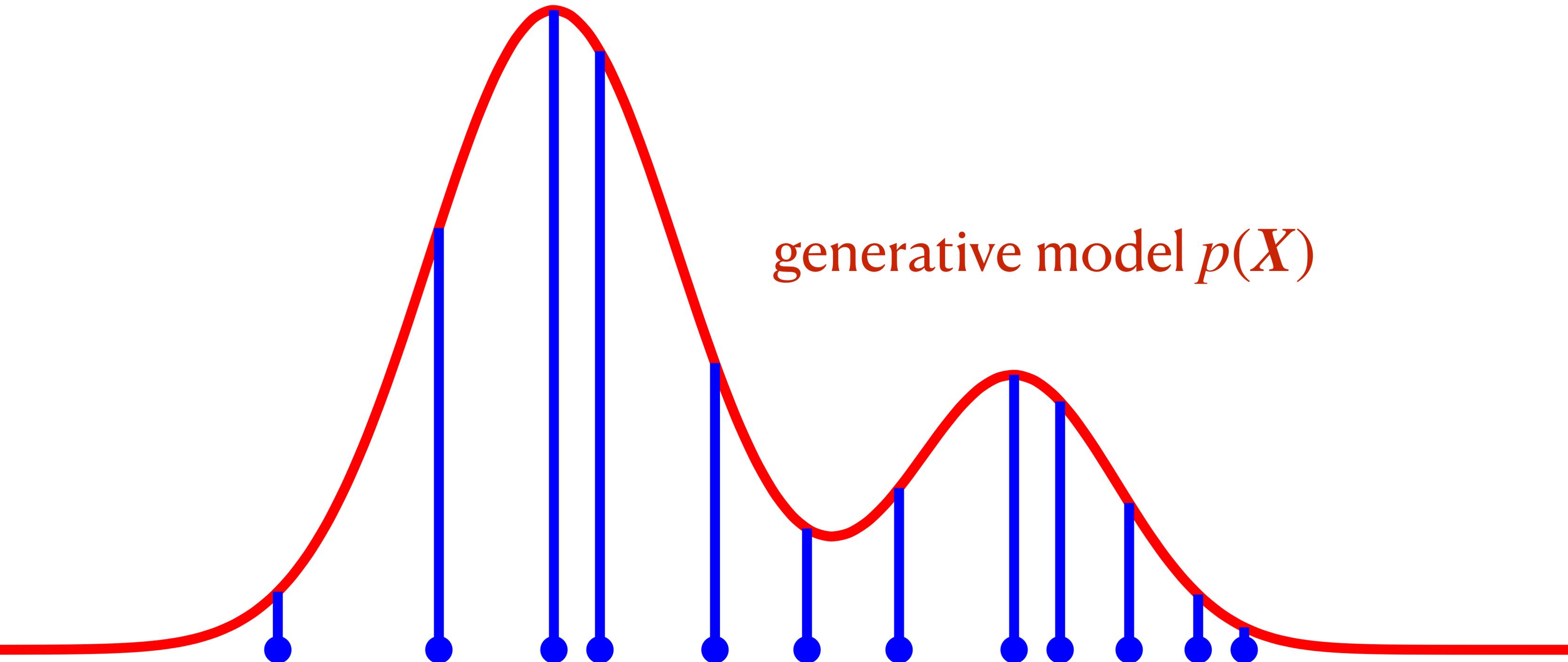
Probability theory 101

Conditional probability $p(y | X)$

Joint probability $p(X, y)$

Product rule $p(X, y) = p(y | X)p(X)$

Sum rule $p(X) = \sum_y p(X, y)$



generative model $p(X)$

data X

Two sides of the same coin

Generative modeling



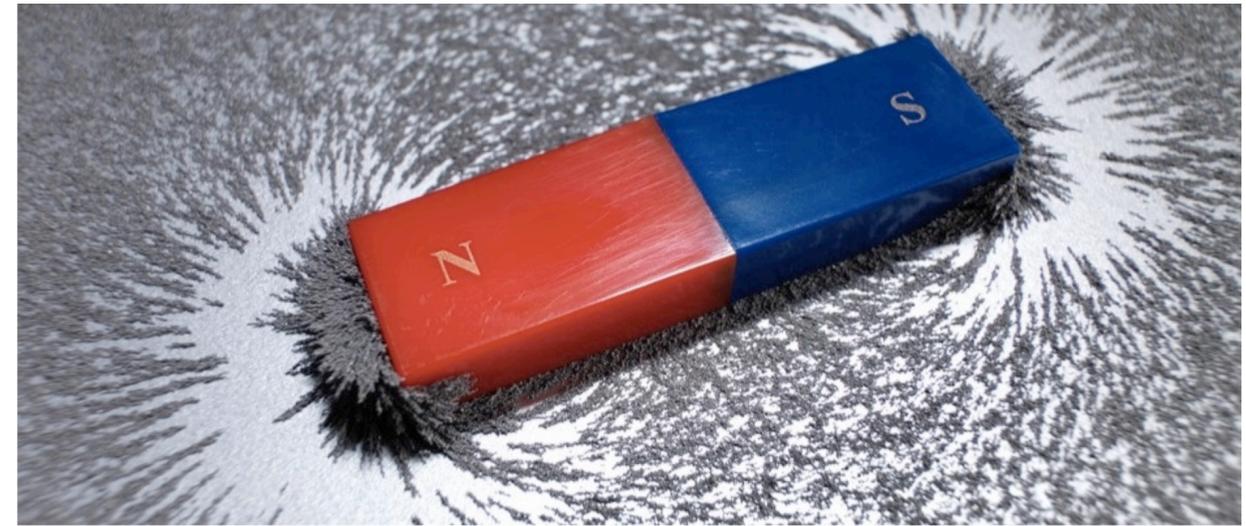
“learn from data”

Maximum likelihood estimation

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

$\mathbb{KL}(\text{data} \parallel p)$ vs $\mathbb{KL}(p \parallel e^{-E/k_B T})$

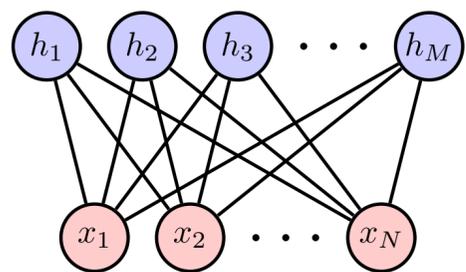
Statistical physics



“learn from energy”

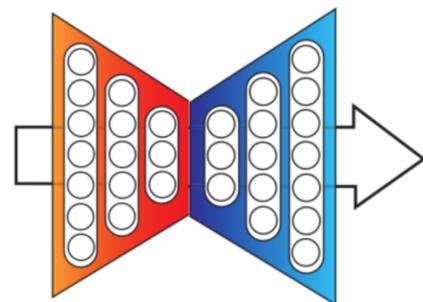
Variational free energy

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



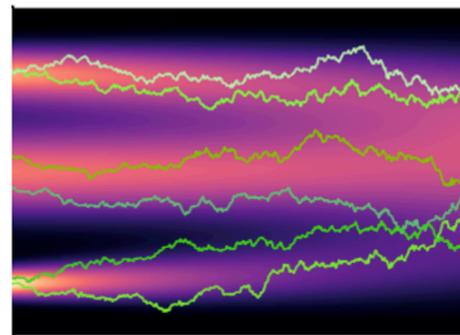
Boltzmann
Machine

1985



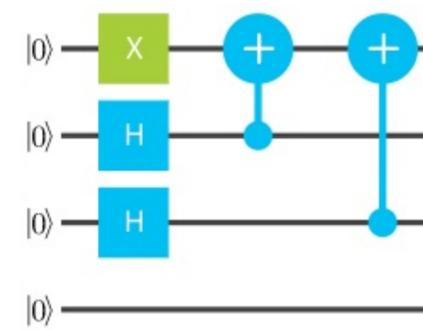
Variational
Autoencoder

2013



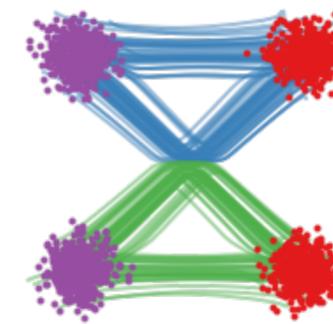
Diffusion
Model

2015



Born
Machine

2017



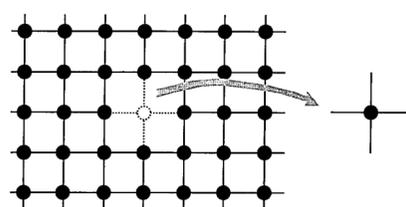
Flow
Matching

2022

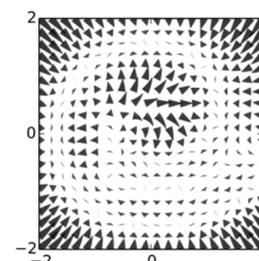
Monte Carlo
Ising model



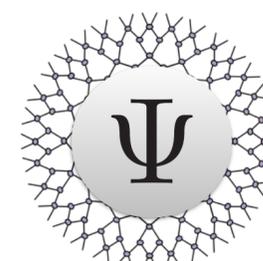
Variational
mean field



Nonequilibrium
thermodynamics



Tensor networks
Quantum circuits



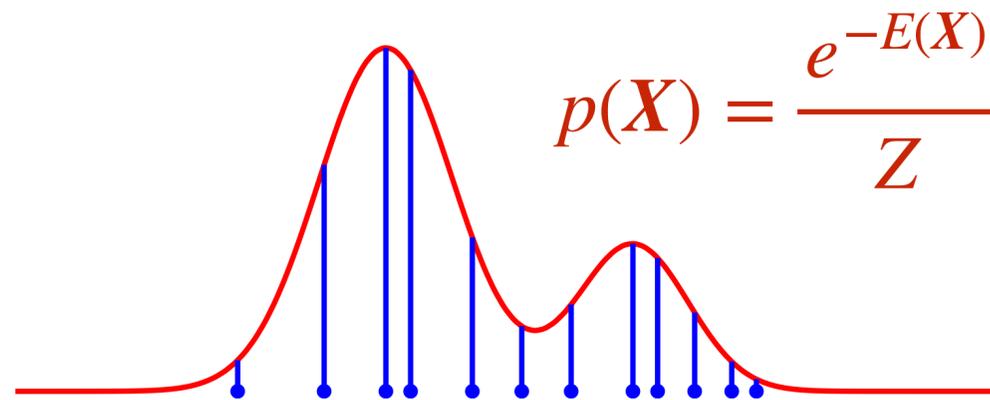
Fluid optimal
transportation

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

Statistical, quantum, fluid, ... physics insights into generative models

Leverage the power of modern generative models for science

Boltzmann machines



2210.10318

GAUSSIAN-BERNOULLI RBMs WITHOUT TEARS 😂

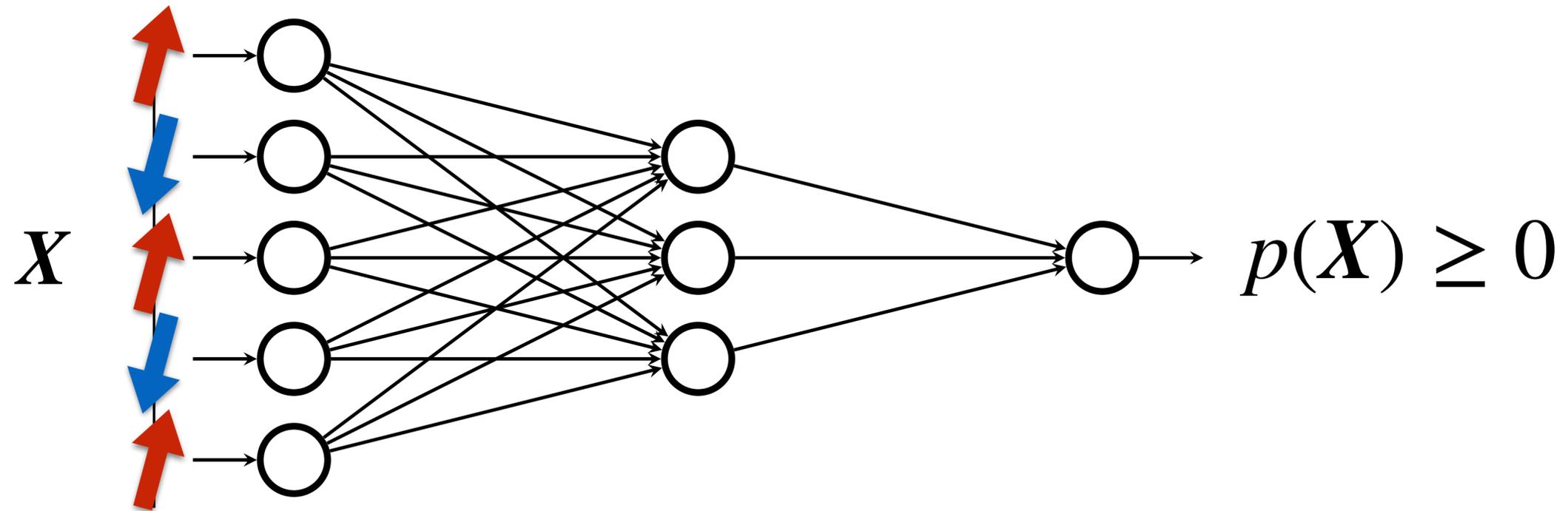
Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

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

$$\nabla_{\theta} \mathcal{L} = \langle \nabla_{\theta} E \rangle_{\text{data}} - \langle \nabla_{\theta} E \rangle_{\text{model}}$$

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

So, why bother ?



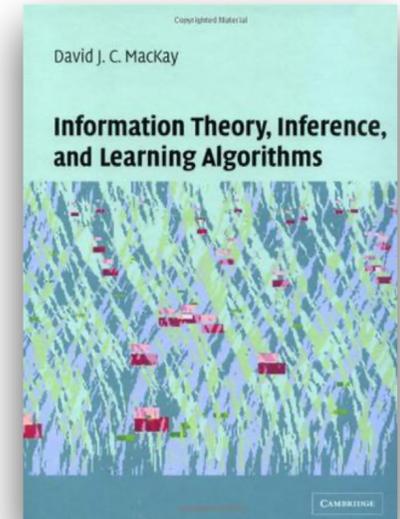
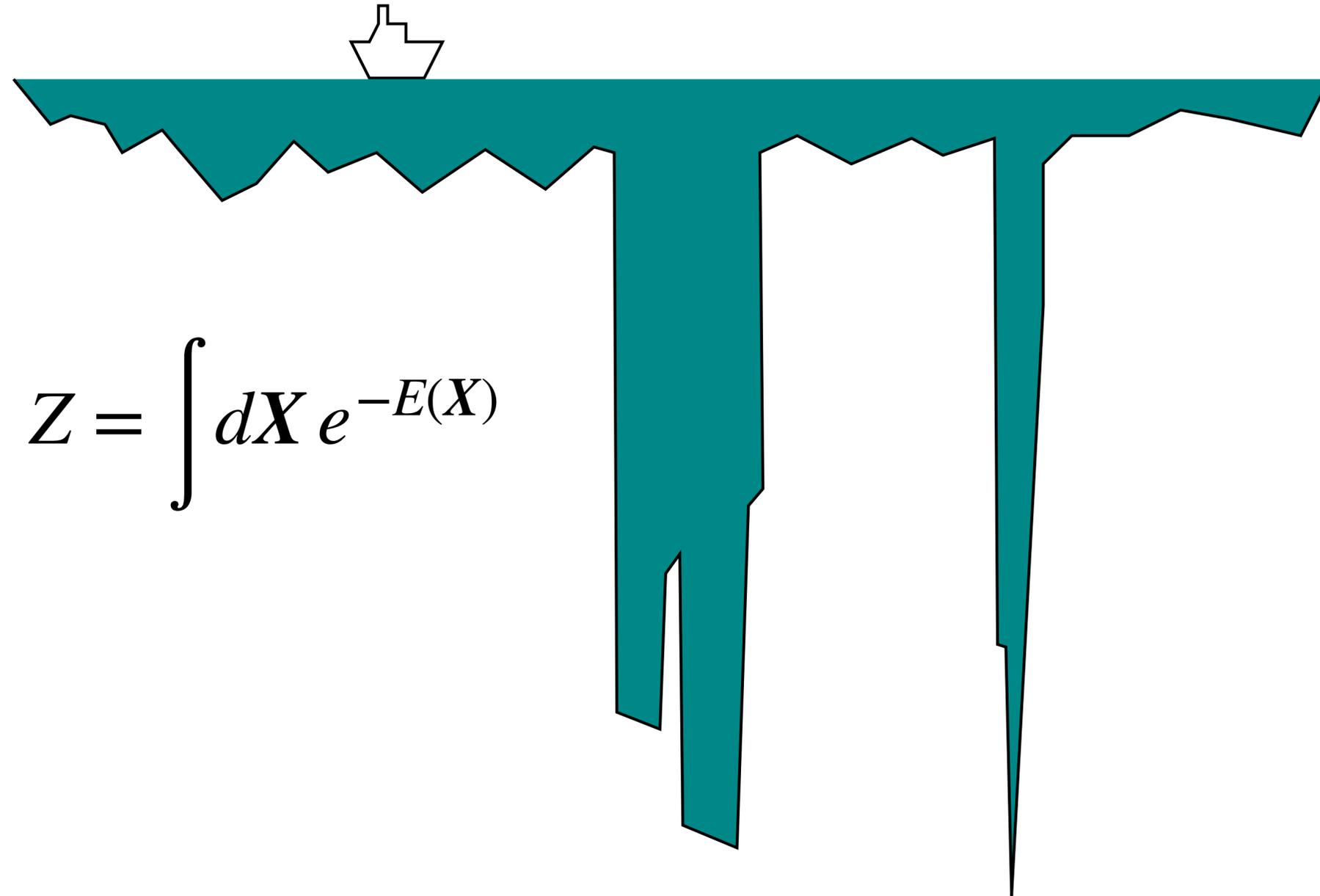
Normalization ?

$$\int dX p(X)$$

Sampling ?

$$X \sim p(X)$$

归一化难题

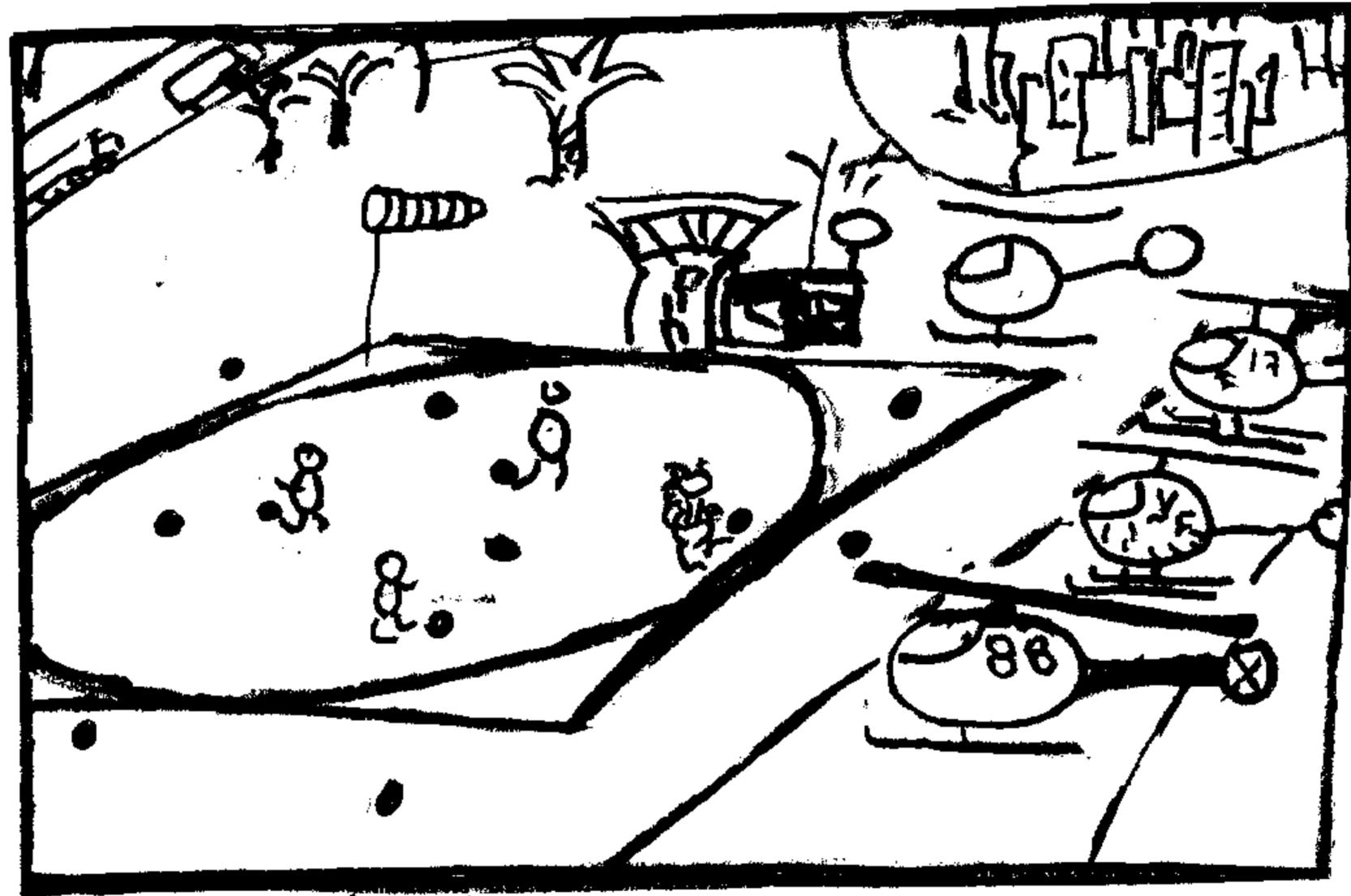


“Intractable” partition function Z

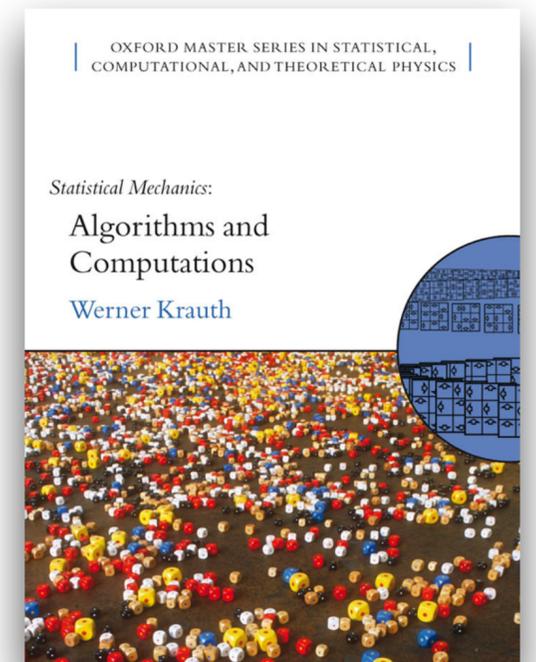
appears widely in machine learning and statistical physics (entropy and free energy calculation)

采样难题

$$X \sim p(X)$$



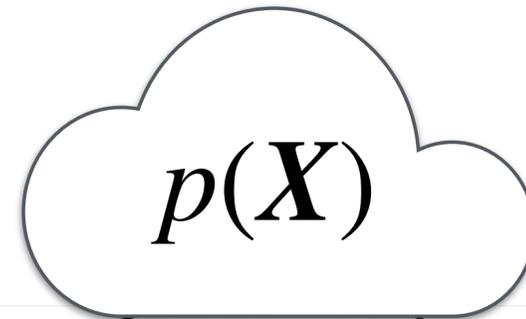
Adults computing the number π at the Monte Carlo heliport.



Direct sampling is generally difficult in high-dimensional space

Generative models and their physics genes

Goodfellow,
NIPS tutorial, 1701.00160



Explicit density

Implicit density

Direct
GAN

Tractable density

Approximate density

Markov Chain
GSN

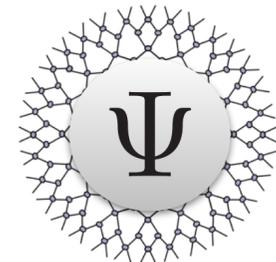
Variational

Markov Chain

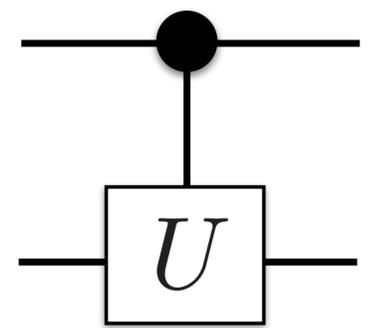
-Fully visible belief nets
-NADE
-MADE
-PixelCNN
-Change of variables models (nonlinear ICA)

Autoregressive model

Variational autoencoder Boltzmann machine + **Diffusion models**



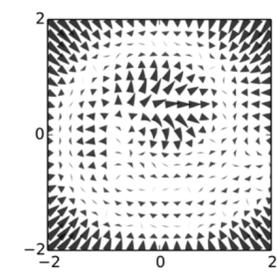
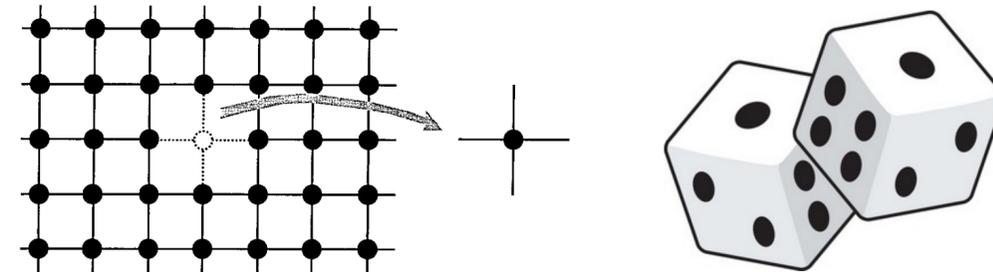
Tensor Networks
Han et al, PRX '18



Quantum Circuits
Liu et al PRA '18



Flow model



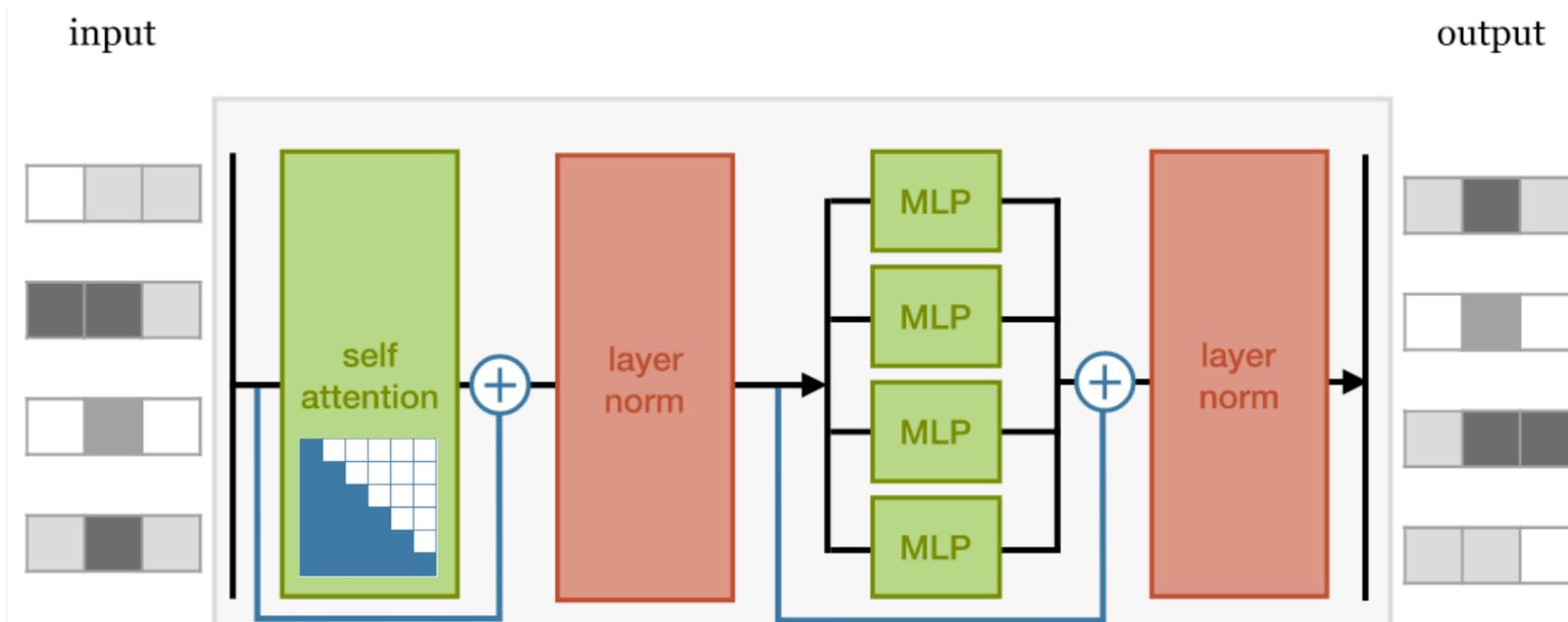
Autoregressive models

$$p(X) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2)\dots$$



“... *the murderer is* _____”

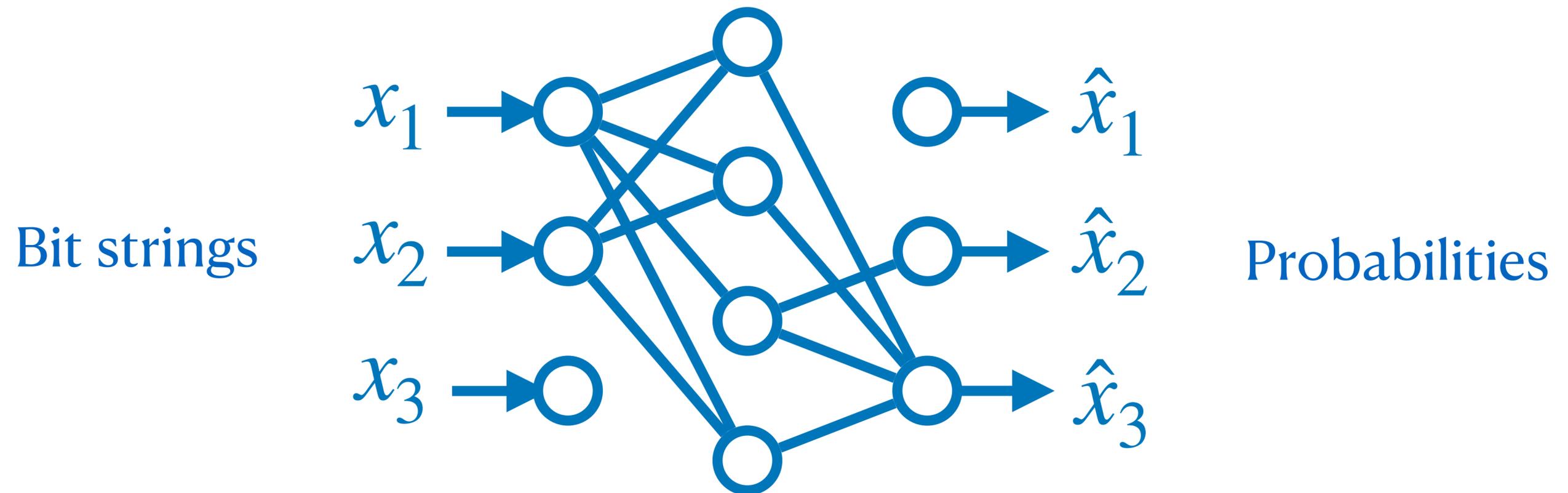
$p(_ | \dots)$



Generative **P**retrained **T**ransformer

Implementation: autoregressive masks

Masked Autoencoder Germain et al, 1502.03509



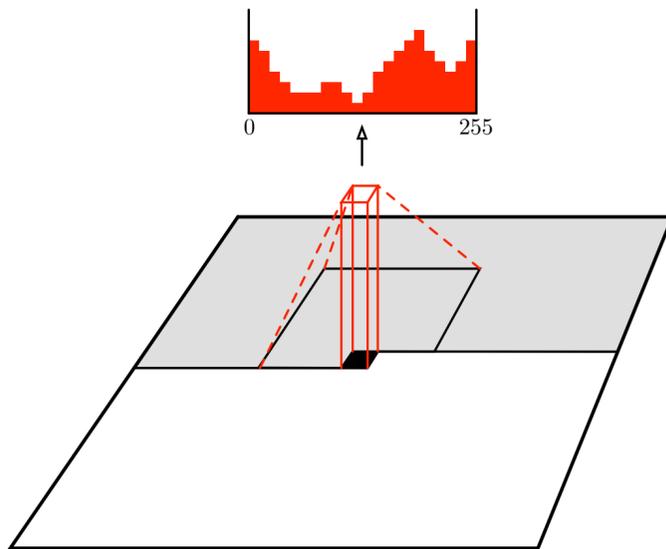
$$p(x_1) = \text{Bernoulli}(\hat{x}_1) \quad p(x_2 | x_1) = \text{Bernoulli}(\hat{x}_2) \quad p(x_3 | x_1, x_2) = \text{Bernoulli}(\hat{x}_3)$$

Other ways to implement autoregressive models: recurrent networks

Implementation: autoregressive masks

Mask convolutional kernel

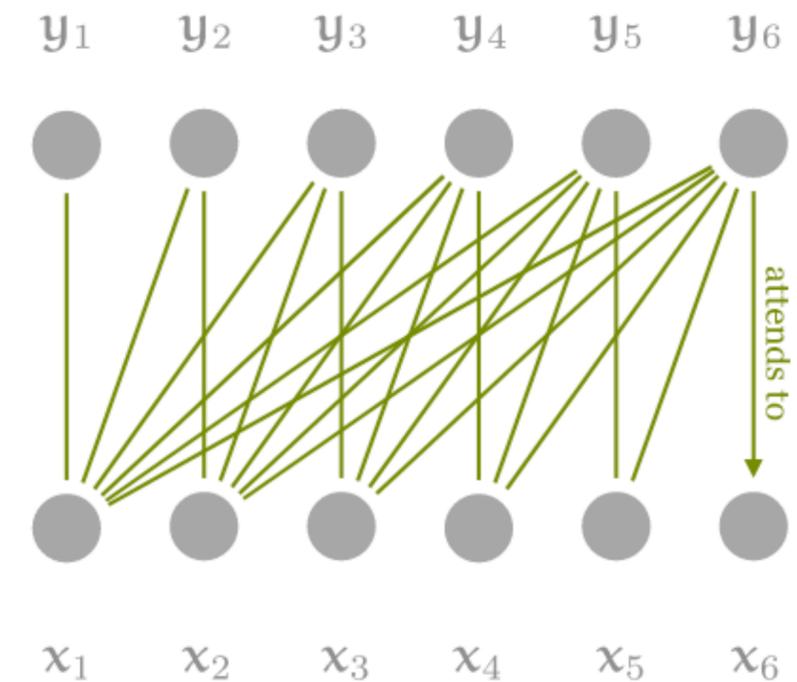
PixelCNN, van den Oord et al, 1601.06759



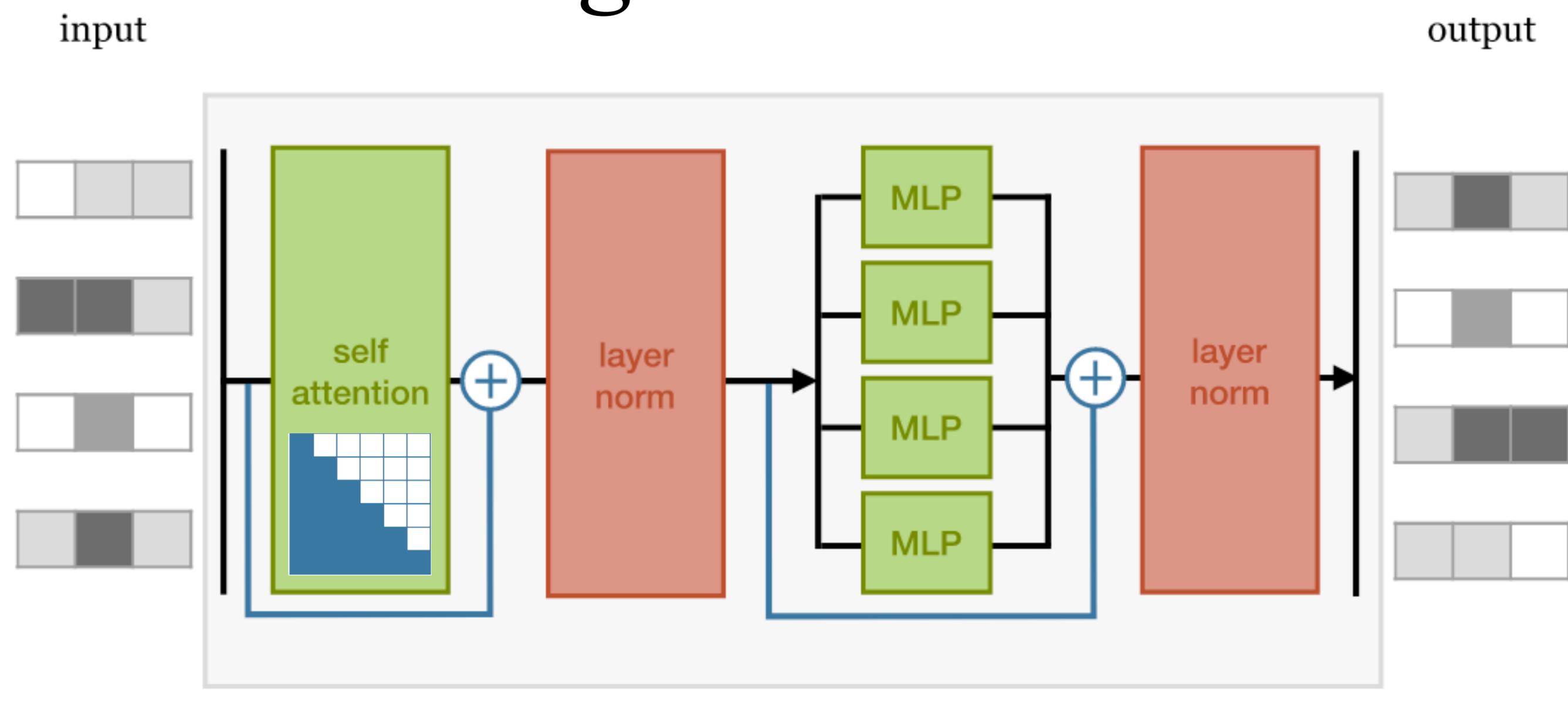
1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

Mask self-attention matrix

Causal transformer, Vaswani et al 1706.03762



The autoregressive transformer



Great at capturing long-range dependence; friendly to backpropagation and GPUs

Masked attention matrix \Rightarrow lower triangular Jacobian matrix \Rightarrow autoregressive model

```
picoGPT — wanglei@bright90:~ — vi gpt2_pico.py — 116x63
1 import numpy as np
2
3 def gelu(x):
4     return 0.5 * x * (1 + np.tanh(np.sqrt(2 / np.pi) * (x + 0.044715 * x**3)))
5
6 def softmax(x):
7     exp_x = np.exp(x - np.max(x, axis=-1, keepdims=True))
8     return exp_x / np.sum(exp_x, axis=-1, keepdims=True)
9
10 def layer_norm(x, g, b, eps: float = 1e-5):
11     mean = np.mean(x, axis=-1, keepdims=True)
12     variance = np.var(x, axis=-1, keepdims=True)
13     return g * (x - mean) / np.sqrt(variance + eps) + b
14
15 def linear(x, w, b):
16     return x @ w + b
17
18 def ffn(x, c_fc, c_proj):
19     return linear(gelu(linear(x, **c_fc)), **c_proj)
20
21 def attention(q, k, v, mask):
22     return softmax(q @ k.T / np.sqrt(q.shape[-1]) + mask) @ v
23
24 def mha(x, c_attn, c_proj, n_head):
25     x = linear(x, **c_attn)
26     qkv_heads = list(map(lambda x: np.split(x, n_head, axis=-1), np.split(x, 3, axis=-1)))
27     causal_mask = (1 - np.tri(x.shape[0], dtype=x.dtype)) * -1e10
28     out_heads = [attention(q, k, v, causal_mask) for q, k, v in zip(*qkv_heads)]
29     x = linear(np.hstack(out_heads), **c_proj)
30     return x
31
32 def transformer_block(x, mlp, attn, ln_1, ln_2, n_head):
33     x = x + mha(layer_norm(x, **ln_1), **attn, n_head=n_head)
34     x = x + ffn(layer_norm(x, **ln_2), **mlp)
35     return x
36
37 def gpt2(inputs, wte, wpe, blocks, ln_f, n_head):
38     x = wte[inputs] + wpe[range(len(inputs))]
39     for block in blocks:
40         x = transformer_block(x, **block, n_head=n_head)
41     return layer_norm(x, **ln_f) @ wte.T
42
43 def generate(inputs, params, n_head, n_tokens_to_generate):
44     from tqdm import tqdm
45     for _ in tqdm(range(n_tokens_to_generate), "generating"):
46         logits = gpt2(inputs, **params, n_head=n_head)
47         next_id = np.argmax(logits[-1])
48         inputs.append(int(next_id))
49     return inputs[len(inputs) - n_tokens_to_generate :]
50
51 def main(prompt: str, n_tokens_to_generate: int = 40, model_size: str = "124M", models_dir: str = "models"):
52     from utils import load_encoder_hparams_and_params
53     encoder, hparams, params = load_encoder_hparams_and_params(model_size, models_dir)
54     input_ids = encoder.encode(prompt)
55     assert len(input_ids) + n_tokens_to_generate < hparams["n_ctx"]
56     output_ids = generate(input_ids, params, hparams["n_head"], n_tokens_to_generate)
57     output_text = encoder.decode(output_ids)
58     return output_text
59
60 if __name__ == "__main__":
61     import fire
62     fire.Fire(main)
"gpt2_pico.py" 62L, 2330B
```

GPT2 in 60 lines of numpy

<https://jaykmody.com/blog/gpt-from-scratch>



Andrej Karpathy

@karpathy

In 2019, OpenAI announced GPT-2 with this post:
openai.com/index/better-l...

Today (~5 years later) you can train your own for ~\$672, running on one 8XH100 GPU node for 24 hours. Our latest `llm.c` post gives the walkthrough in some detail:
github.com/karpathy/llm.c...

Incredibly, the costs have come down dramatically over the last 5 years due to improvements in compute hardware (H100 GPUs), software (CUDA, cuBLAS, cuDNN, FlashAttention) and data quality (e.g. the FineWeb-Edu dataset). For this exercise, the algorithm was kept fixed and follows the GPT-2/3 papers.

Because `llm.c` is a direct implementation of GPT training in C/CUDA, the requirements are minimal - there is no need for conda environments, Python interpreters, pip installs, etc. You spin up a cloud GPU node (e.g. on Lambda), optionally install NVIDIA cuDNN, NCCL/MPI, download the .bin data shards, compile and run, and you're stepping in minutes. You then wait 24 hours and enjoy samples about English-speaking Unicorns in the Andes.

For me, this is a very nice checkpoint to get to because the entire `llm.c` project started with me thinking about reproducing GPT-2 for an educational video, getting stuck with some PyTorch things, then rage quitting to just write the whole thing from scratch in C/CUDA. That set me on a longer journey than I anticipated, but it was quite fun, I learned more CUDA, I made friends along the way, and `llm.c` is really nice now. It's ~5,000 lines of code, it compiles and steps very fast so there is very little waiting around, it has constant memory footprint, it trains in mixed precision, distributed across multi-node with NNCL, it is bitwise deterministic, and hovers around ~50% MFU. So it's quite cute.

Or, ~1000 lines of C

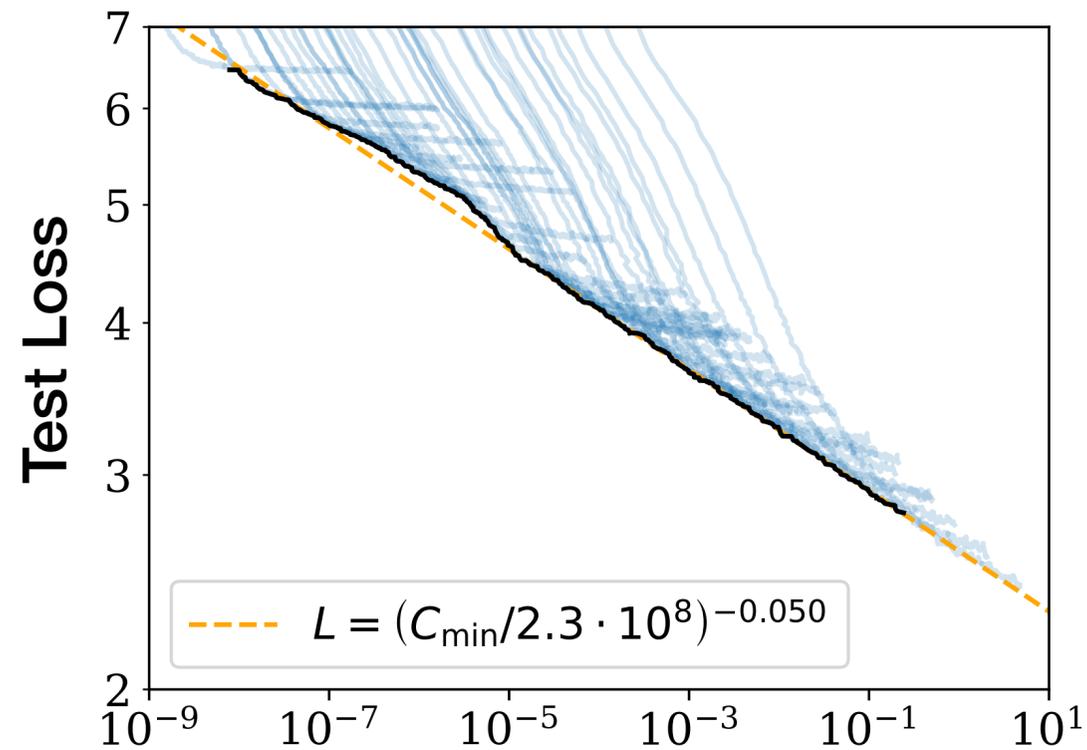
<https://github.com/karpathy/llm.c>

<https://x.com/karpathy/status/1811467135279104217>

Scaling law of the loss function

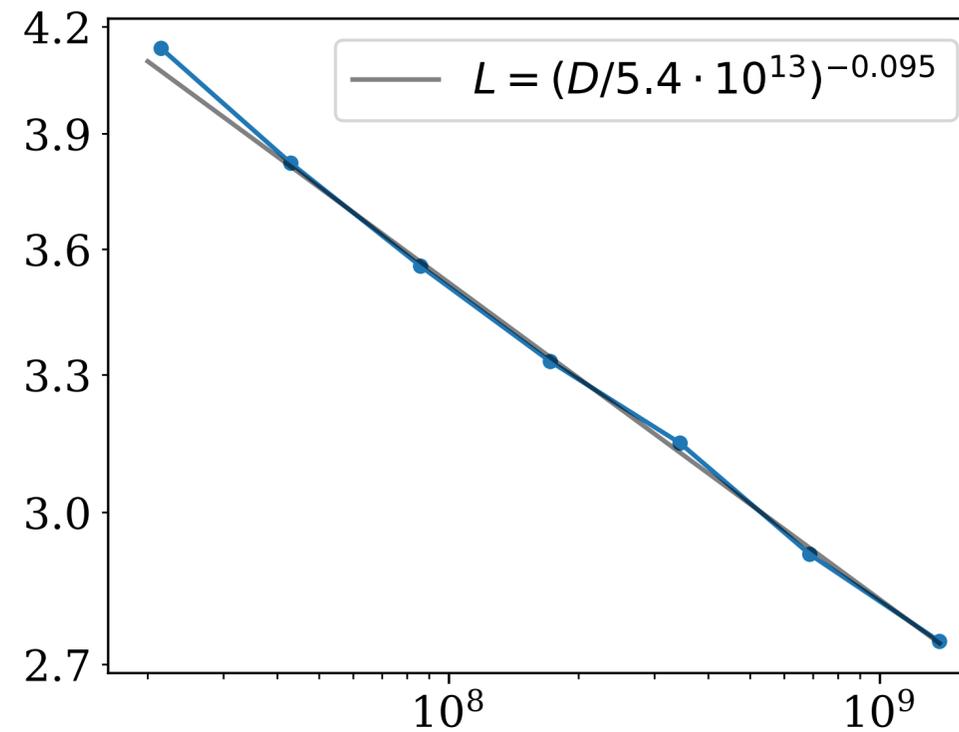
$$\mathcal{L} = - \mathbb{E}_{X \sim \text{data}} [\ln p(X)]$$

Kaplan et al, 2001.08361



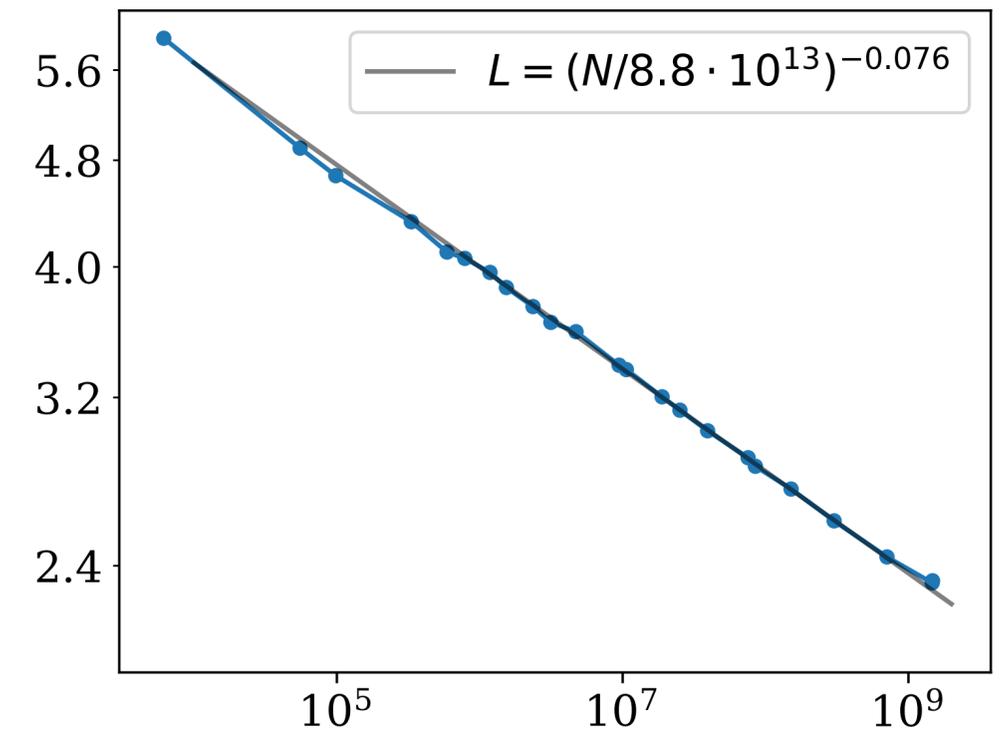
Compute

PF-days, non-embedding



Dataset Size

tokens



Parameters

non-embedding

“It would also be exciting to find a theoretical framework from which the scaling relations can be derived: a ‘statistical mechanics’ underlying the ‘thermodynamics’ we have observed.”

Kullback–Leibler divergence

$$\mathbb{KL}(\pi \parallel p) \equiv \int dx \pi(x) [\ln \pi(x) - \ln p(x)]$$

$$\mathbb{KL}(\pi \parallel p) \geq 0$$

$$\mathbb{KL}(\pi \parallel p) = 0 \iff \pi(x) = p(x)$$

$$\mathbb{KL}(\pi \parallel p) \neq \mathbb{KL}(p \parallel \pi)$$

Learn from data

$$\pi(\mathbf{x}) \propto \sum_{d \in \text{data}} \delta(\mathbf{x} - d)$$

$$\min_{\theta} \text{KL}(\pi \parallel p_{\theta}) \iff \min_{\theta} \left\{ -\mathbb{E}_{\mathbf{x} \sim \text{data}} [\ln p_{\theta}(\mathbf{x})] \right\}$$

target model Maximum likelihood estimation

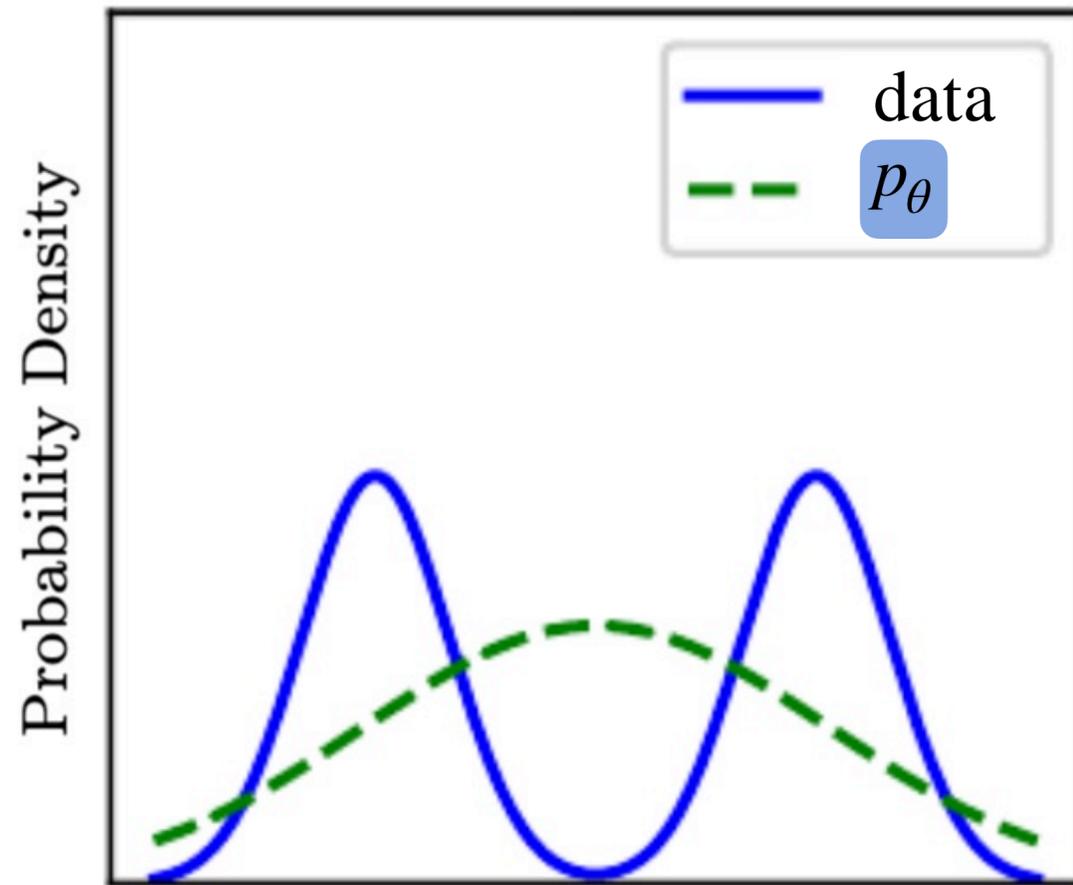
The lower bound is the entropy of the dataset: complete memorization

Forward KL or Reverse KL ?

Maximum likelihood estimation

$$\min_{\theta} \text{KL}(\text{data} \parallel p_{\theta})$$

Mode covering

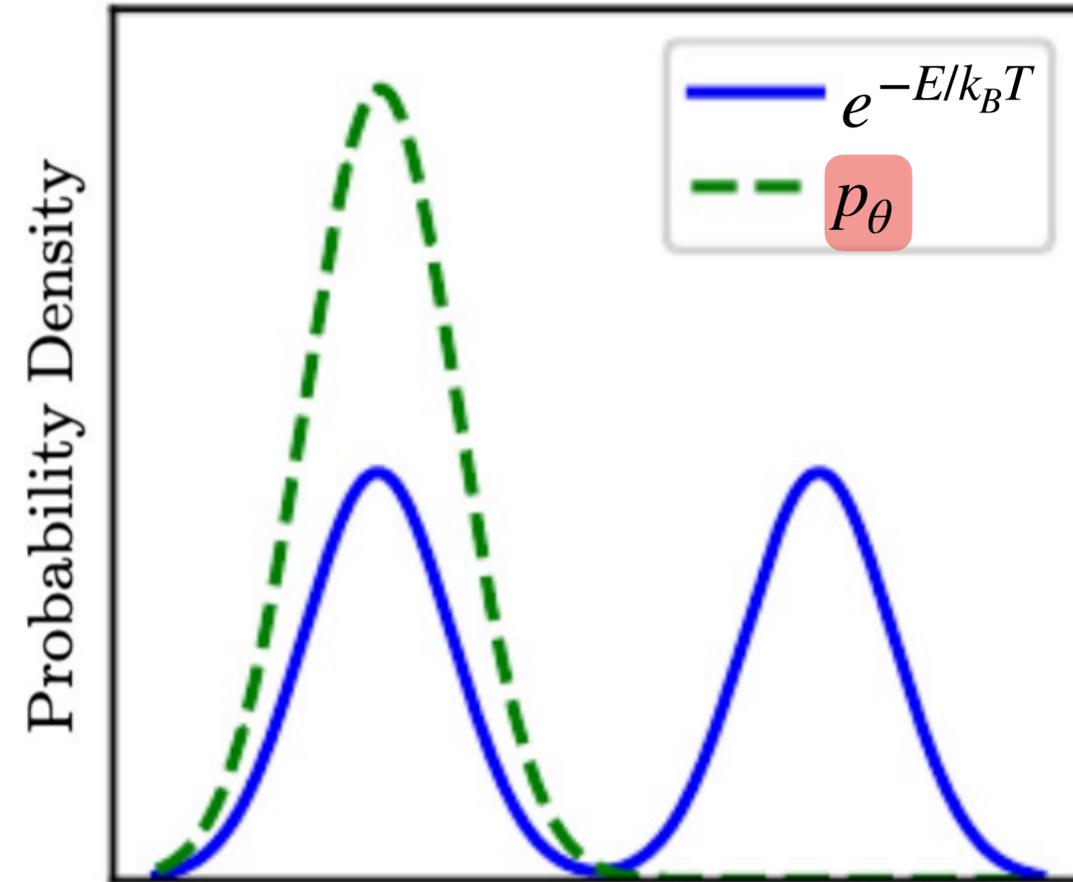


Failure mode: hallucination

Variational free energy

$$\min_{\theta} \text{KL}(p_{\theta} \parallel e^{-E/k_B T})$$

Mode seeking



Failure mode: local minima

The training objective of LLM

① Pretrain with forward KL

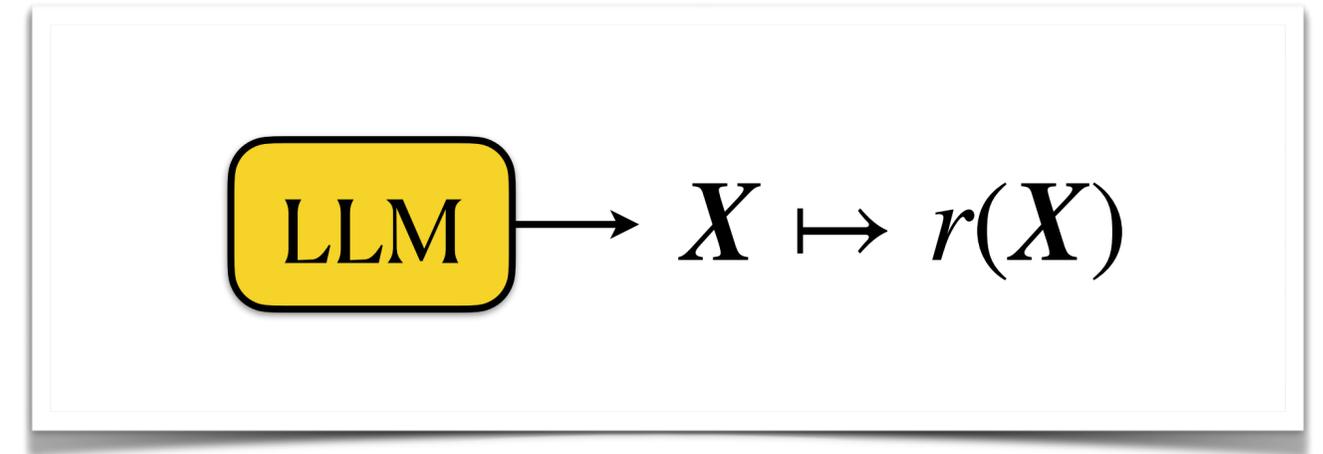


**learn from data
to be a generalist**

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

学而不思则罔

② Finetune with reverse KL



**learn from reward
to be a specialized generalist**

$$F = \mathbb{E}_{X \sim p(X)} [r(X) + T \ln p(X)]$$

思而不学则殆

— 《论语·大模型》



Andrej Karpathy ✓

@karpathy

It's a bit sad and confusing that LLMs ("Large Language Models") have little to do with language; It's just historical. They are highly general purpose technology for statistical modeling of token streams. A better name would be Autoregressive Transformers or something.

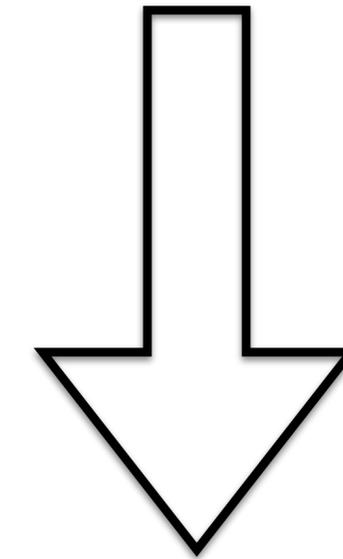
They don't care if the tokens happen to represent little text chunks. It could just as well be little image patches, audio chunks, action choices, molecules, or whatever. If you can reduce your problem to that of modeling token streams (for any arbitrary vocabulary of some set of discrete tokens), you can "throw an LLM at it".

Actually, as the LLM stack becomes more and more mature, we may see a convergence of a large number of problems into this modeling paradigm. That is, the problem is fixed at that of "next token prediction" with an LLM, it's just the usage/meaning of the tokens that changes per domain.

If that is the case, it's also possible that deep learning frameworks (e.g. PyTorch and friends) are way too general for what most problems want to look like over time. What's up with thousands of ops and layers that you can reconfigure arbitrarily if 80% of problems just want to use an LLM?

I don't think this is true but I think it's half true.

Language



Token streams

Autoregressive model is more than language modeling

“Language” => token stream => bitstream => **ANYTHING**

Speech: WaveNet 1609.03499

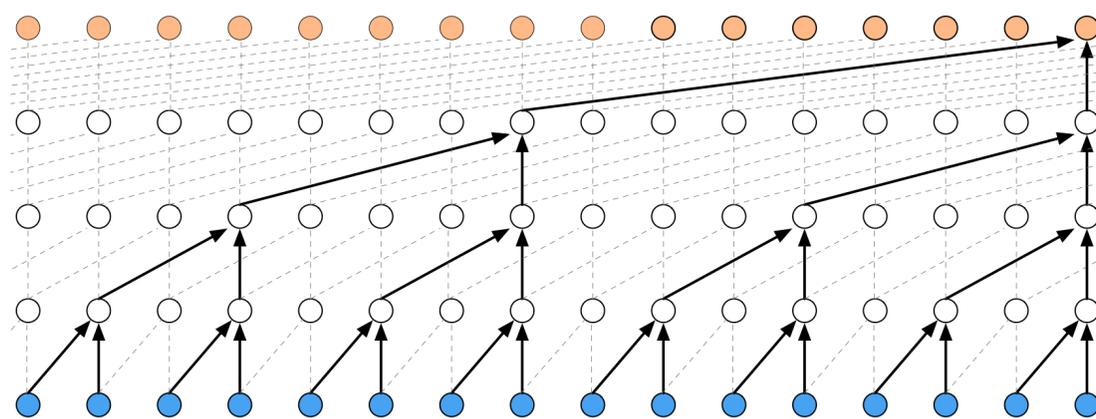
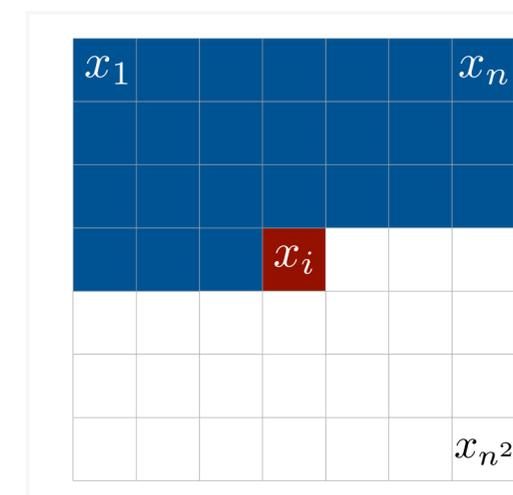
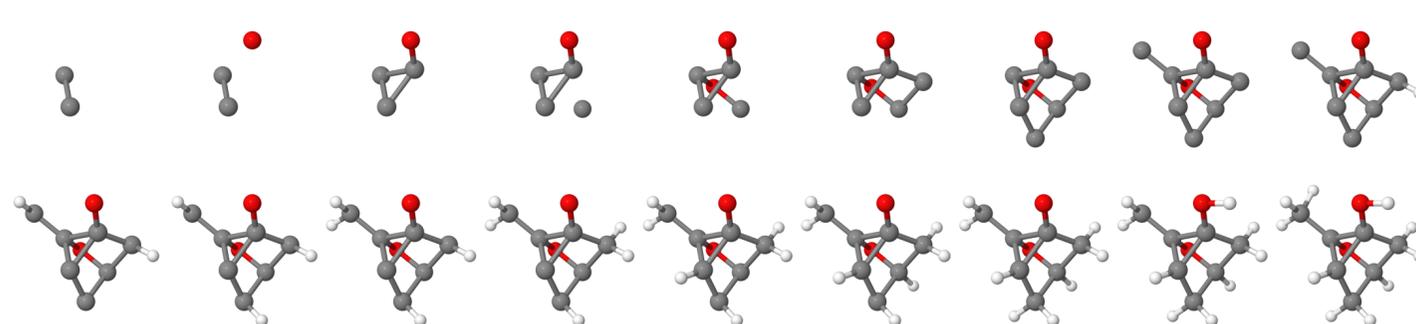


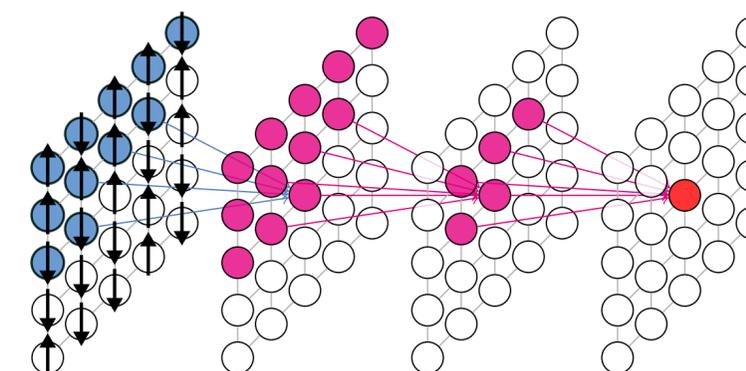
Image: PixelCNN 1601.06759



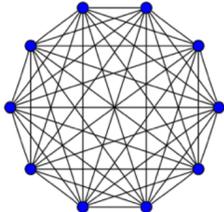
Molecular graph: 1810.11347



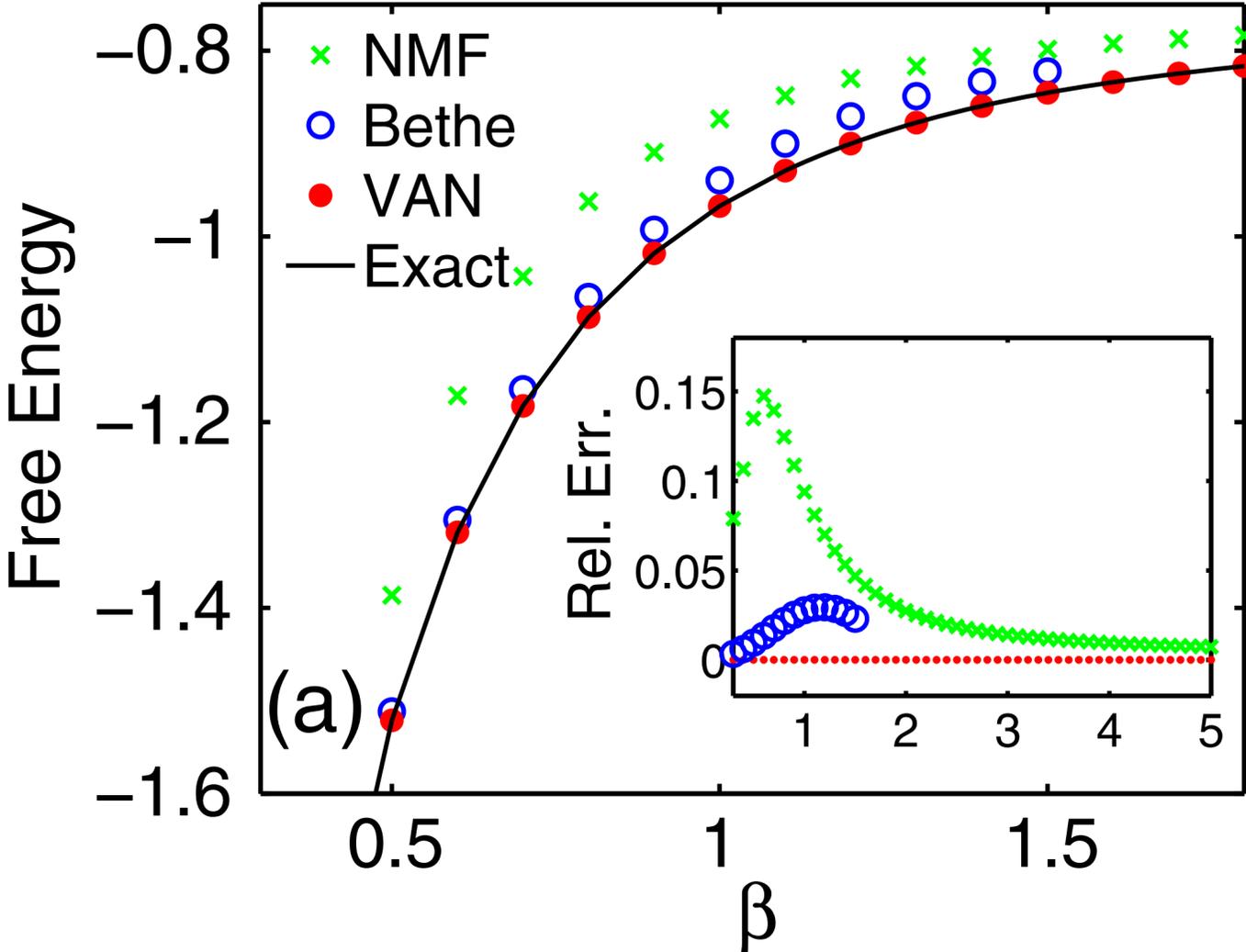
Ising model: 1809.10606



Variational autoregressive network for statistical mechanics



Sherrington-Kirkpatrick spin glass



Naive mean-field factorized probability

$$p(\mathbf{X}) = \prod_i p(x_i)$$

Bethe approximation pairwise interaction

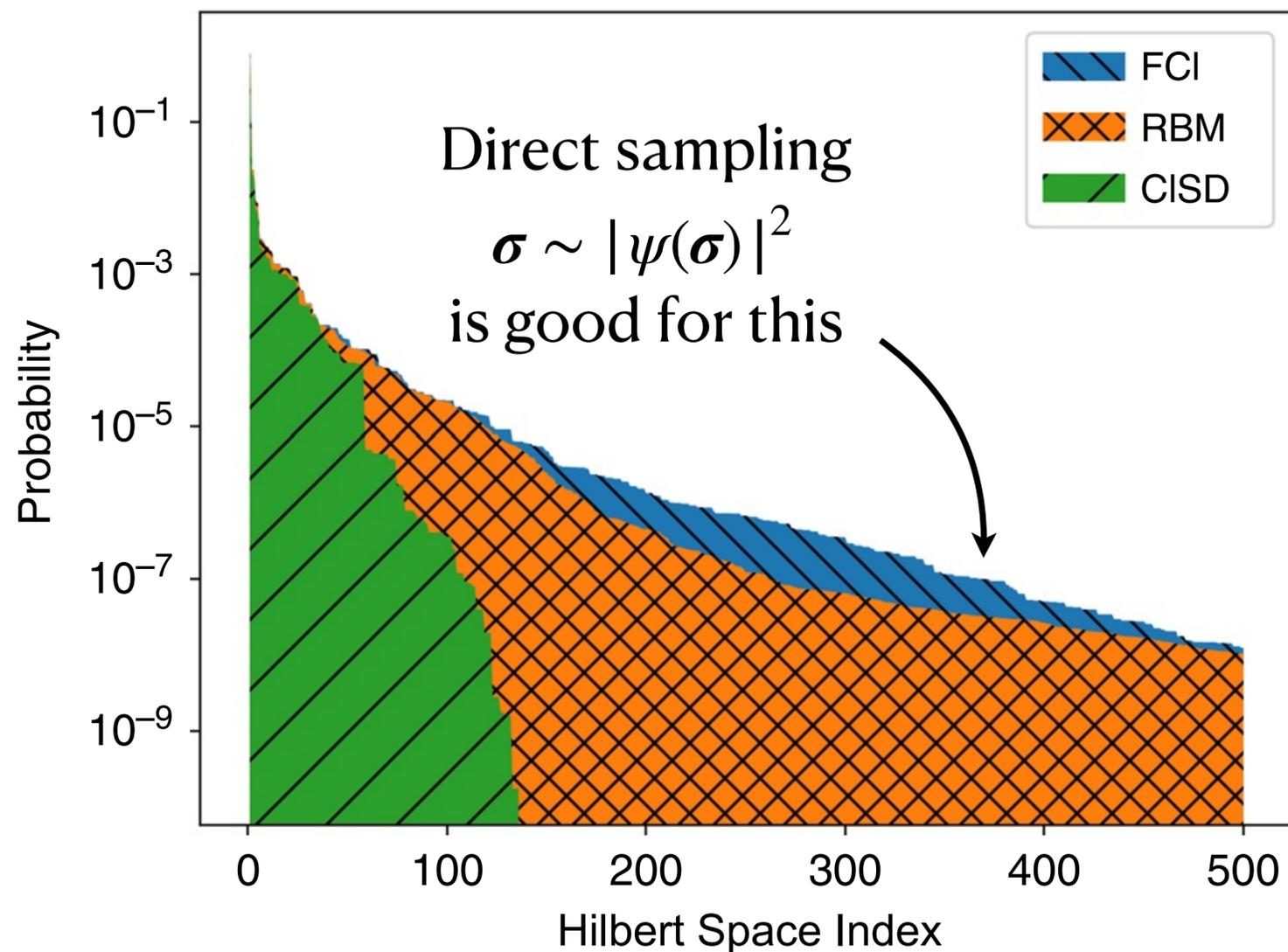
$$p(\mathbf{X}) = \prod_i p(x_i) \prod_{(i,j) \in E} \frac{p(x_i, x_j)}{p(x_i)p(x_j)}$$

Variational autoregressive network

$$p(\mathbf{X}) = \prod_i p(x_i | \mathbf{x}_{<i})$$

Variational autoregressive quantum states

$$\psi(\boldsymbol{\sigma}) = \psi(\sigma_1)\psi(\sigma_2 | \sigma_1)\psi(\sigma_3 | \sigma_1, \sigma_2)\cdots$$



N₂ molecule, Choo et al, Nat. Comm. '20

Objective function: ground state energy

McMillan 1965, Carleo & Troyer Science 2017

$$\frac{\langle \psi | \hat{H} | \psi \rangle}{\langle \psi | \psi \rangle} = \mathbb{E}_{\boldsymbol{\sigma} \sim |\psi(\boldsymbol{\sigma})|^2} \left[\frac{\hat{H}\psi(\boldsymbol{\sigma})}{\psi(\boldsymbol{\sigma})} \right]$$

Sharir, Levine, Wies, Carleo, Shashua, PRL '20

Hibat-Allah, Ganahl, Hayward, Melko, Carrasquilla, PRResearch '20

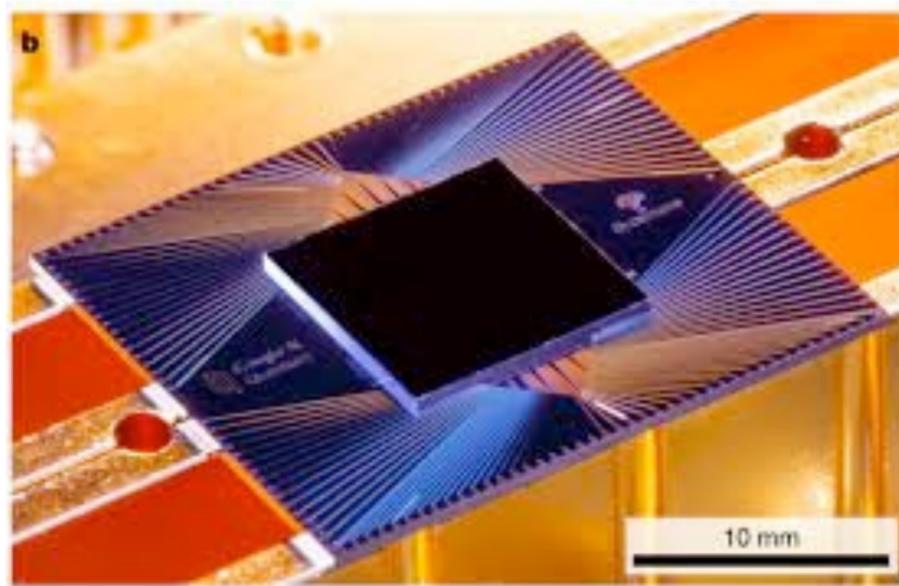
Barrett et al, Nat. Mach. Intell. '22

Zhao et al, MLST. '23

Shang et al, 2307.09343

Demo: Generative model of Sycamore data

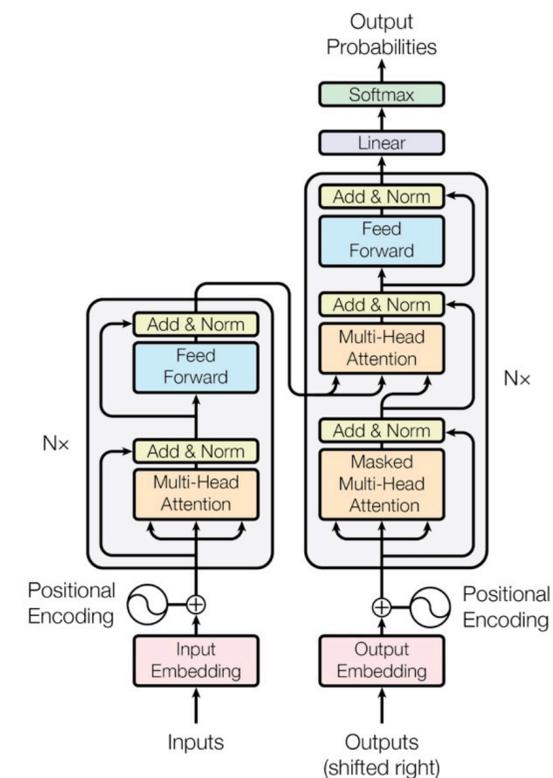
Quantum chip



bitstrings $\sim |\Psi(X)|^2$

011110110100
100001111011
100110110111
100110100010
010100011000
010001000000
010101101100
100001111000
100101001001
001000001010

Transformer

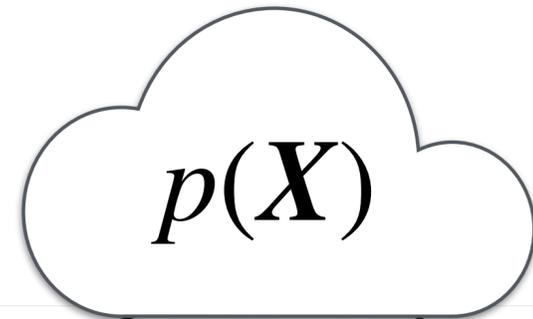


Can we fake the measurement of the sycamore quantum circuit by training a transformer?

 https://colab.research.google.com/drive/11WaroqULkudKT3h2i5J6r_EmA4wFKkoZ?usp=sharing

Generative models and their physics genes

Goodfellow,
NIPS tutorial, 1701.00160



Explicit density

Implicit density

Direct
GAN

Tractable density

Approximate density

Markov Chain
GSN

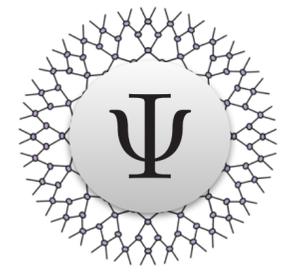
Variational

Markov Chain

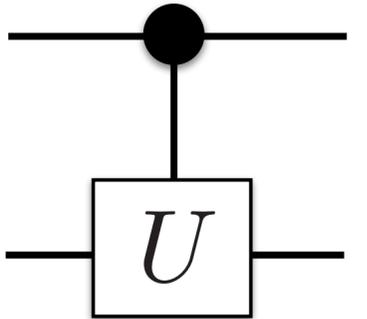
-Fully visible belief nets
-NADE
-MADE
-PixelCNN
-Change of variables models (nonlinear ICA)

Autoregressive model

Variational autoencoder Boltzmann machine + **Diffusion models**



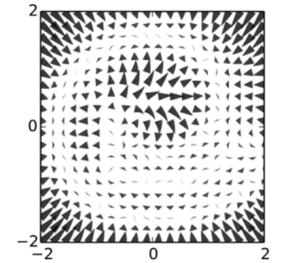
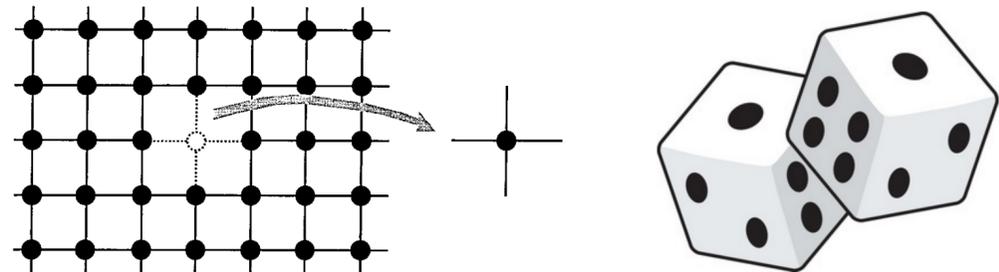
Tensor Networks
Han et al, PRX '18



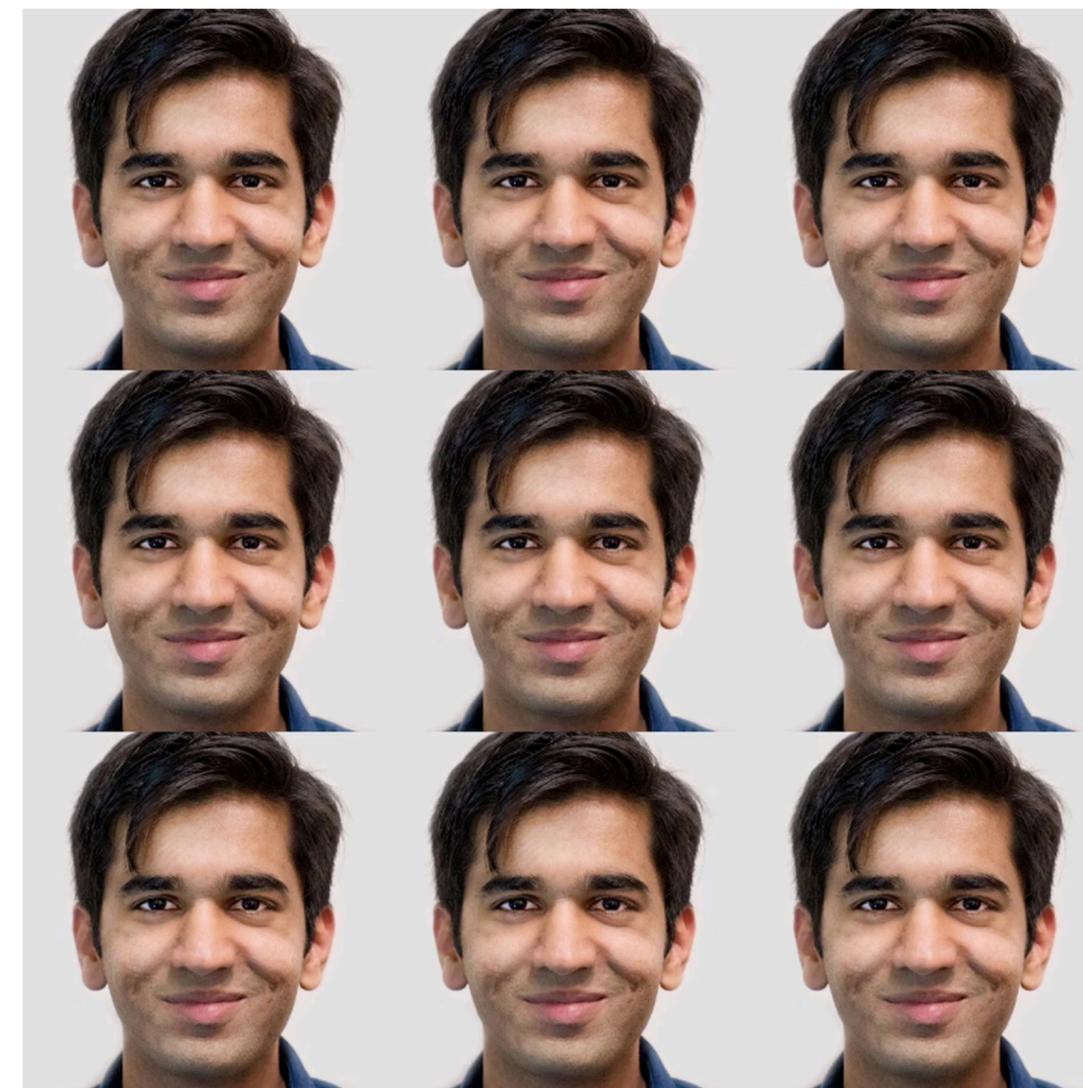
Quantum Circuits
Liu et al PRA '18



Flow model



Normalizing flows



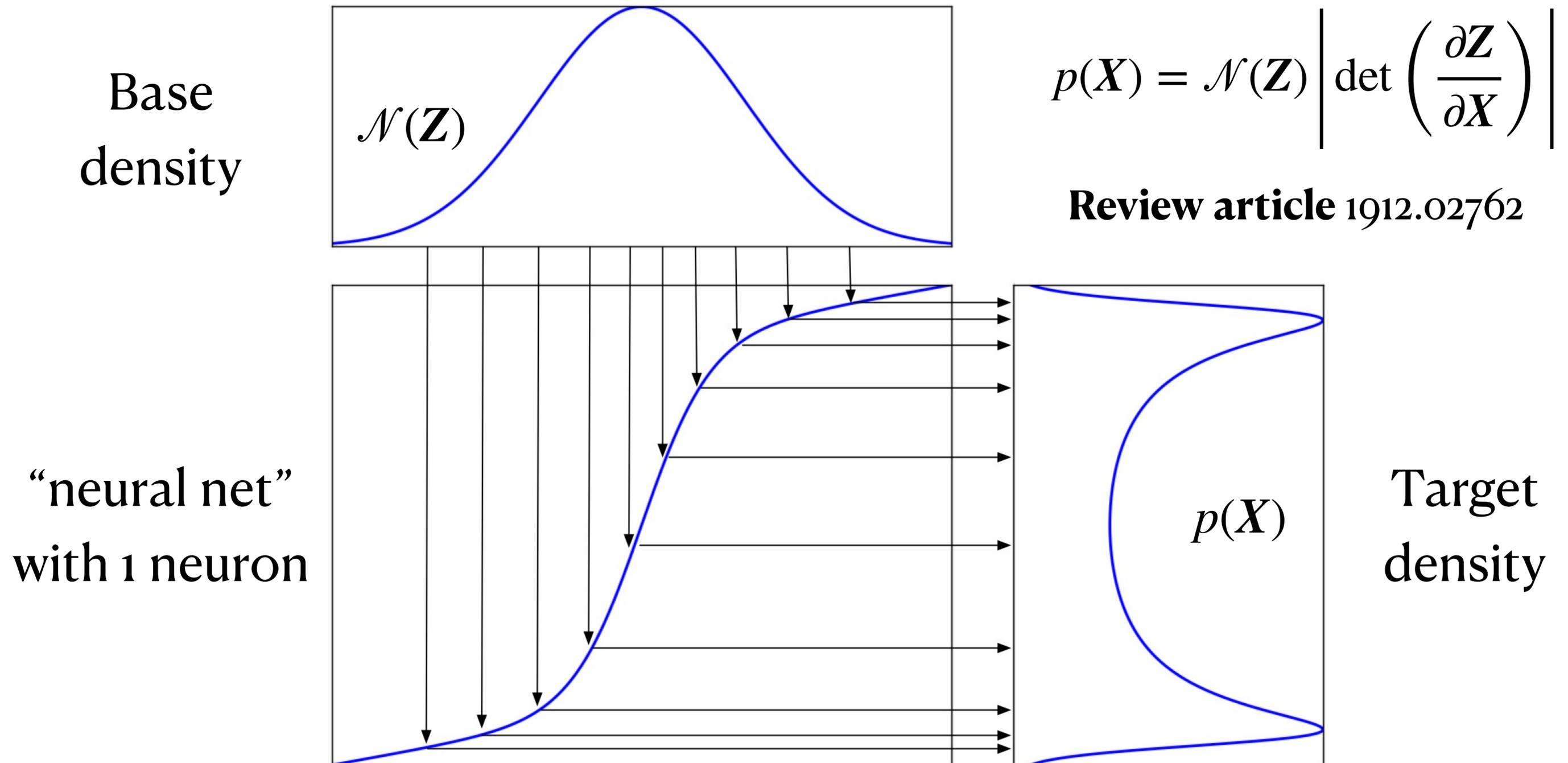
 Parallel WaveNet 1711.10433

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

 Glow 1807.03039

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

Normalizing flow in a nutshell



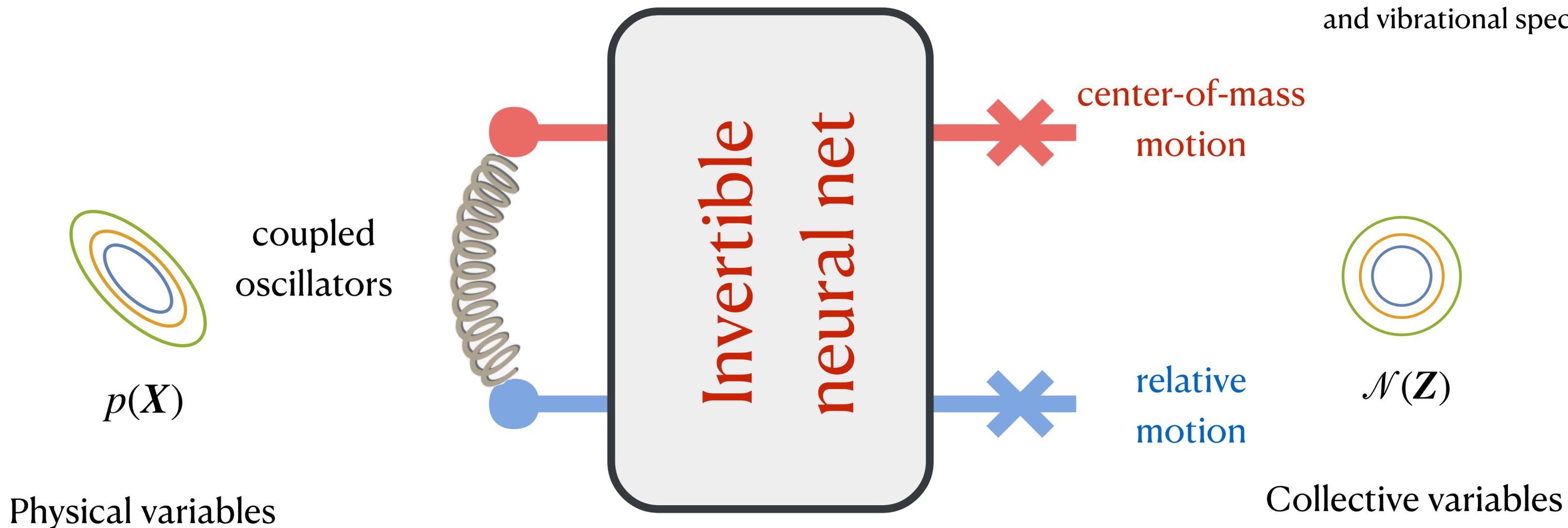
Physics intuition of normalizing flow

Li and LW, PRL '18

Li, Dong, Zhang, LW, PRX '20

Zhang, Wang, LW, JCP '24

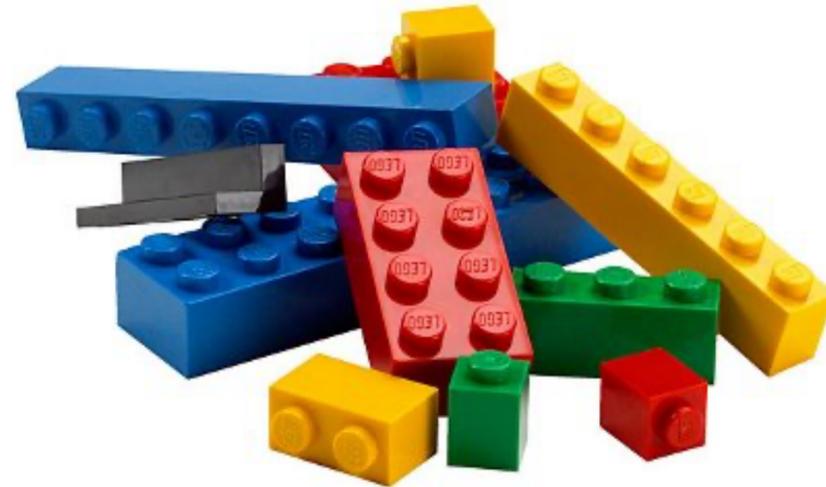
for RG, canonical transformations,
and vibrational spectra



High-dimensional, composable, learnable, nonlinear transformations

Flow architecture design

Composability

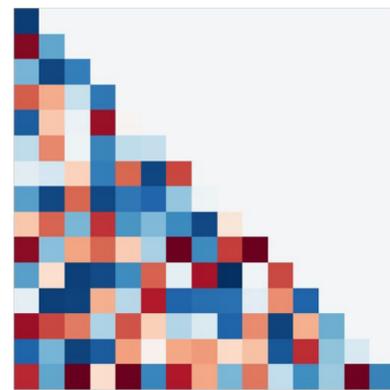


$$Z = \mathcal{T}(X)$$

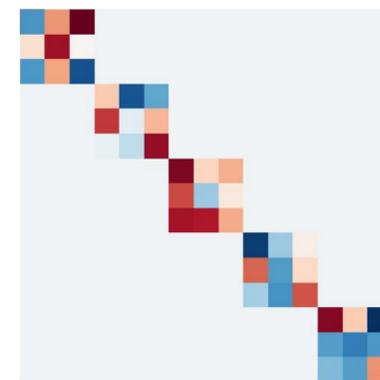
$$\mathcal{T} = \mathcal{T}_1 \circ \mathcal{T}_2 \circ \mathcal{T}_3 \circ \dots$$

**Balanced
efficiency &
inductive bias**

$$\left| \det \left(\frac{\partial Z}{\partial X} \right) \right|$$



Autoregressive



Blockwise

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

Continuous flow

Example of a building block

Forward

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

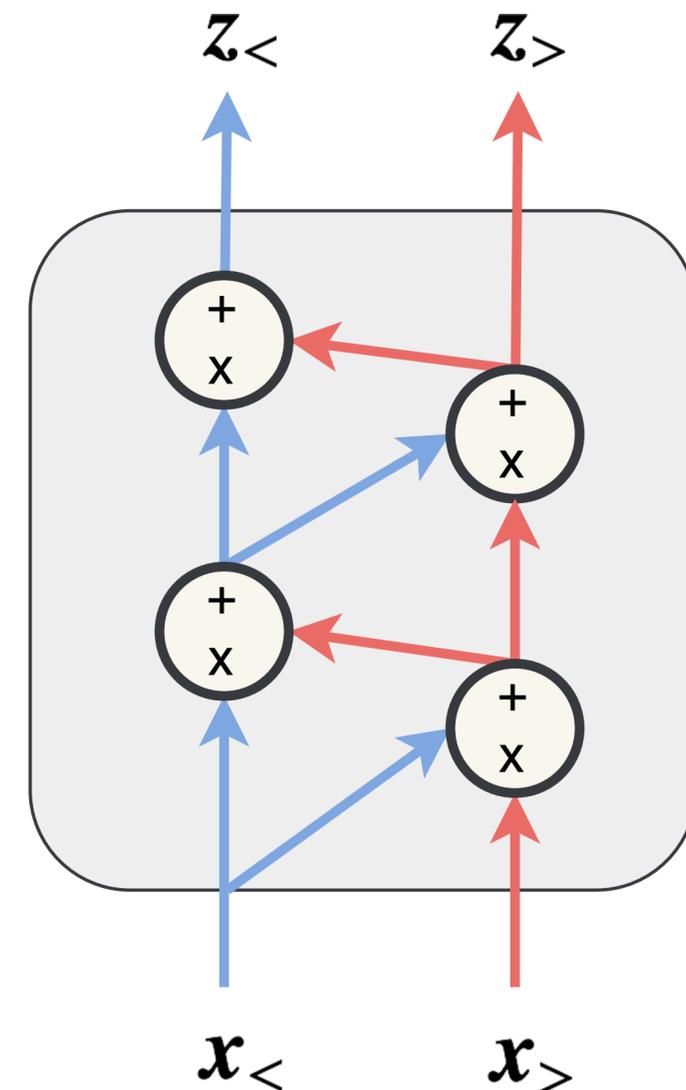
arbitrary
neural nets

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

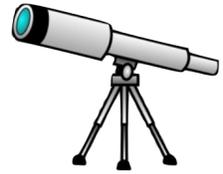
$$\ln \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right| = \sum_i [s(\mathbf{z}_{<})]_i$$



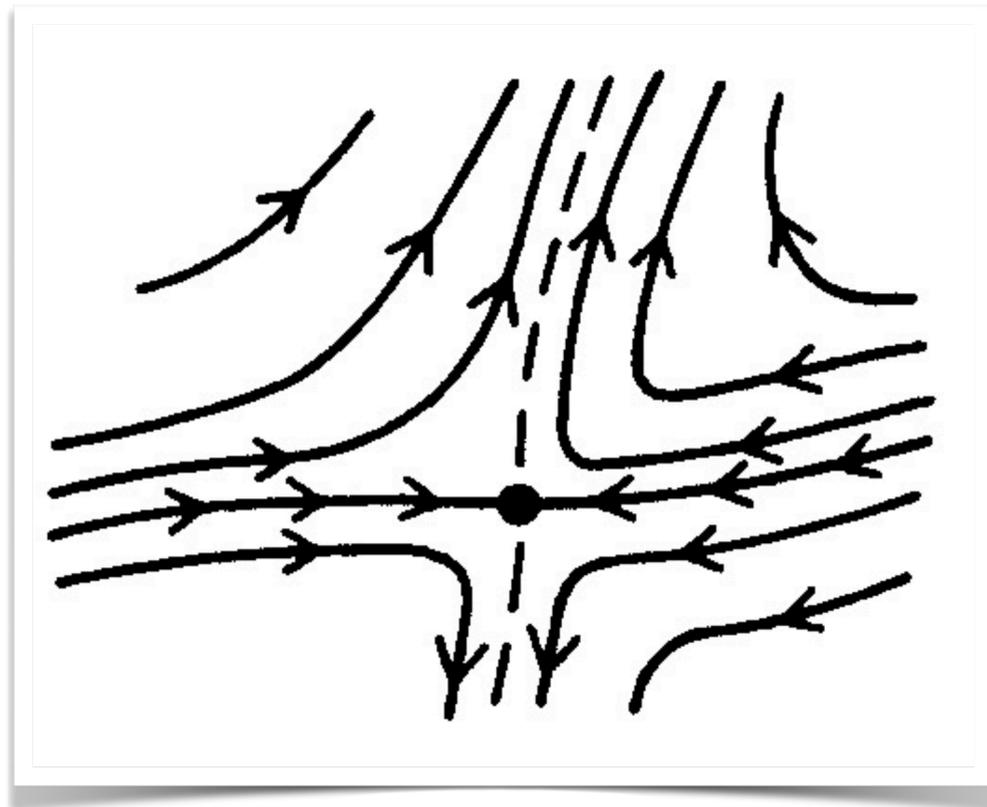
Real NVP, Dinh et al, 1605.08803

Turns out to have surprising connection Störmer–Verlet integration

Why is flow useful for physics?



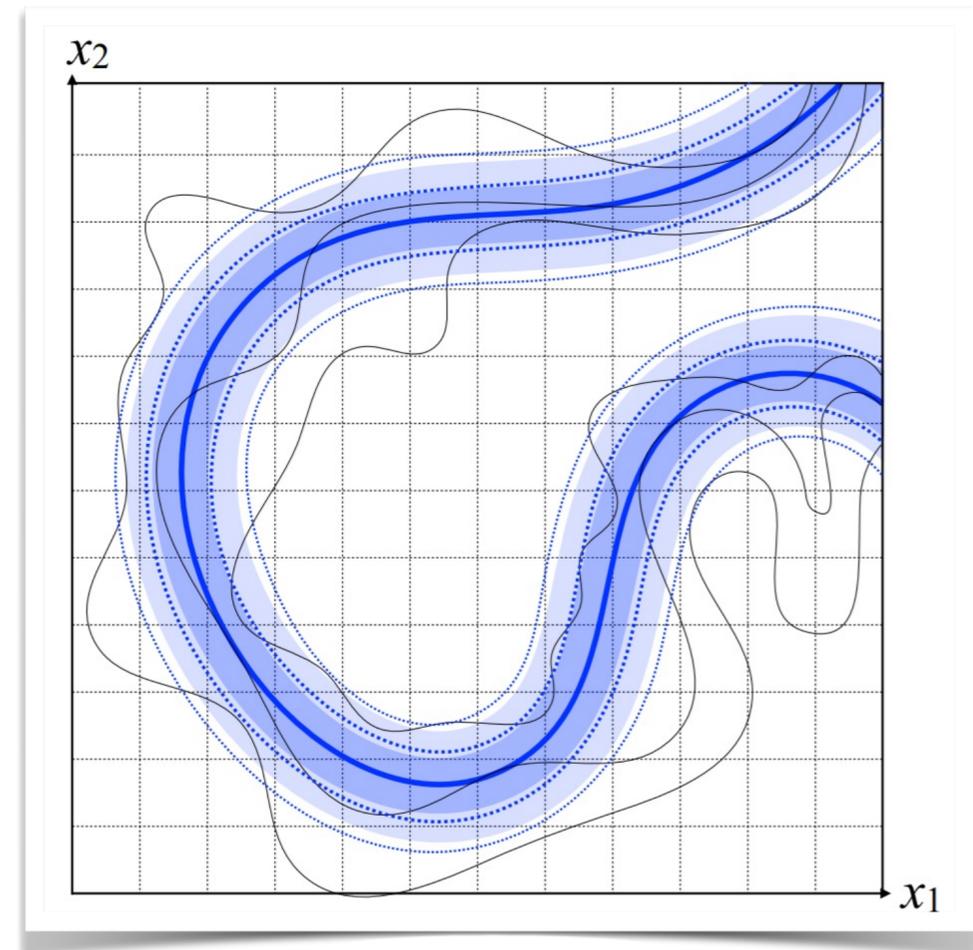
Renormalization group



Effective theory emerges upon transformation of the variables



Monte Carlo update

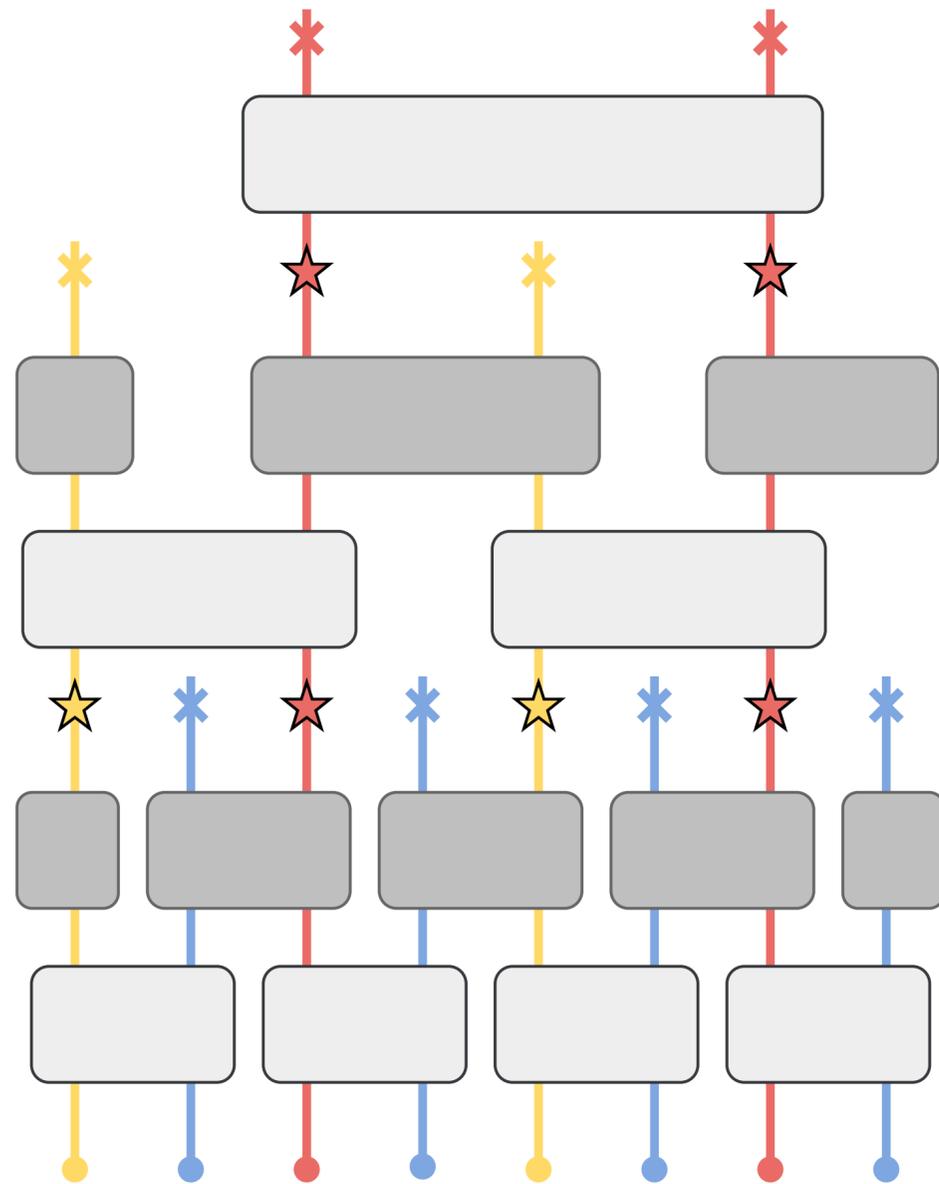


Physics happens on a manifold
Train neural nets to unfold that manifold

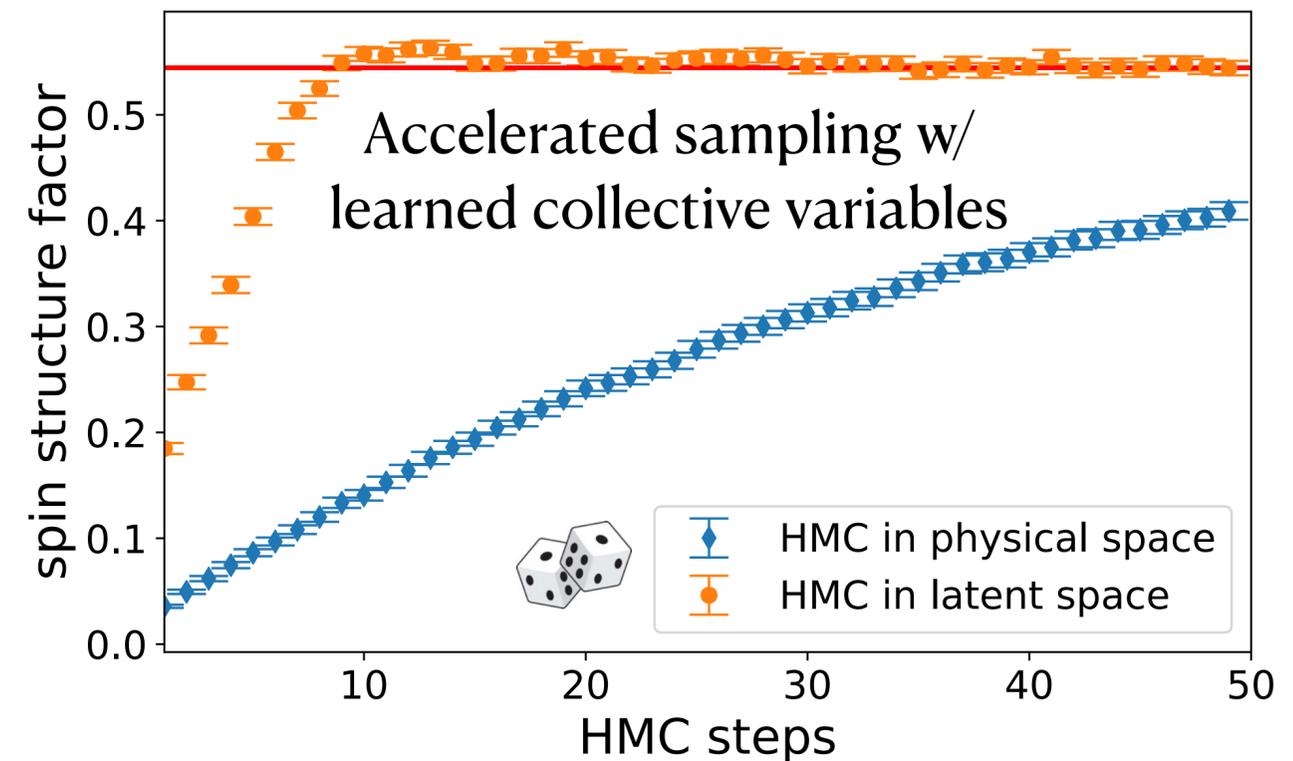
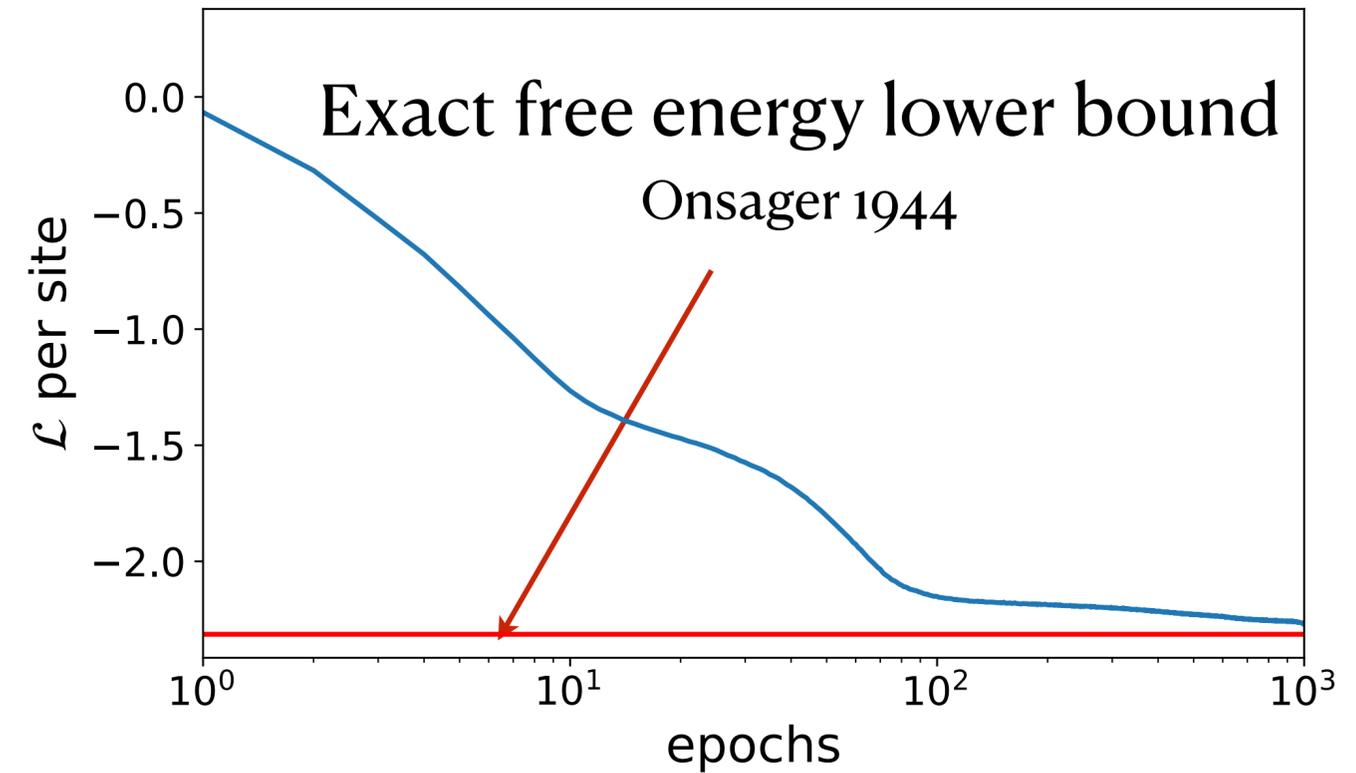
Neural network renormalization group

Li, LW, PRL '18 [li012589/NeuralRG](https://arxiv.org/abs/1801.01258)

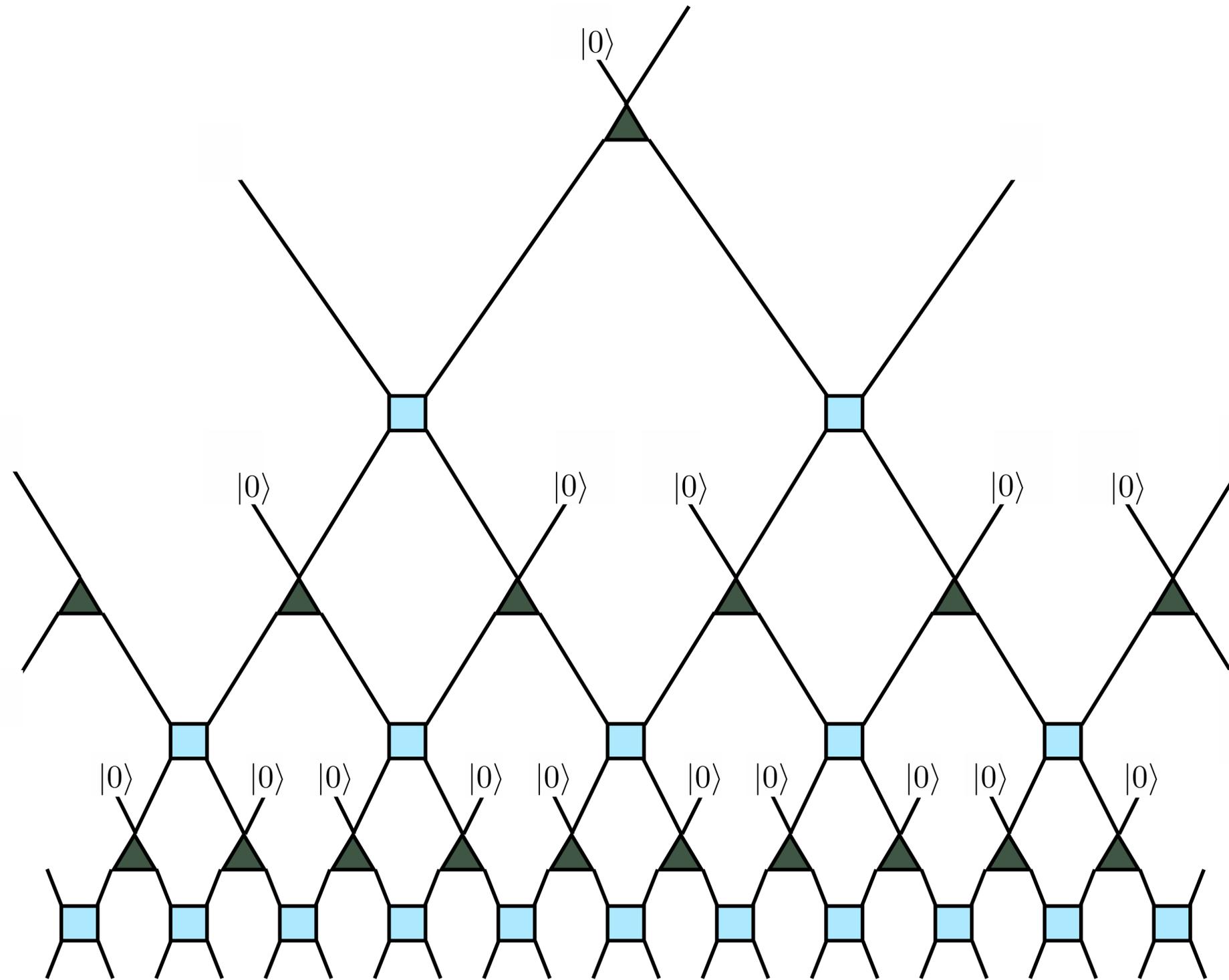
Collective variables



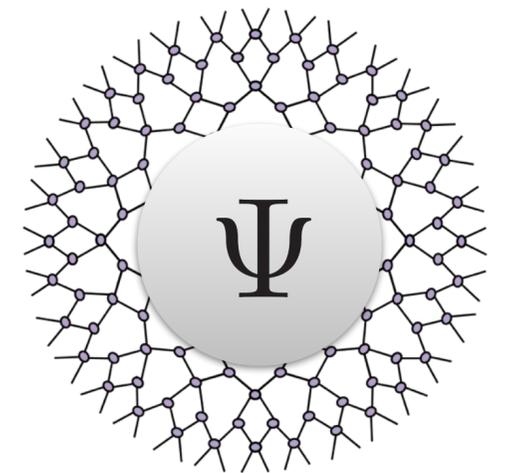
Physical variables



Quantum version of the architecture

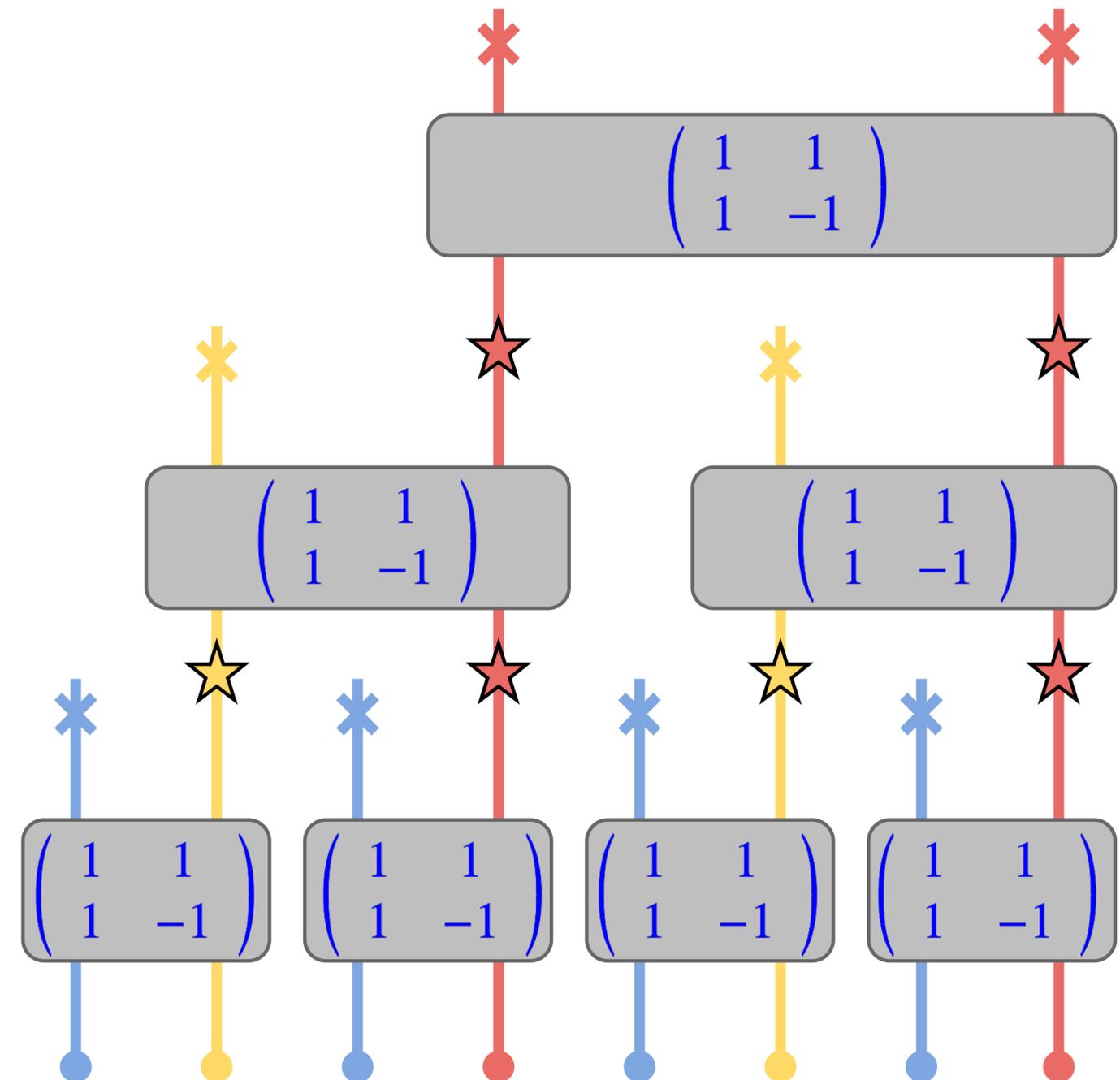
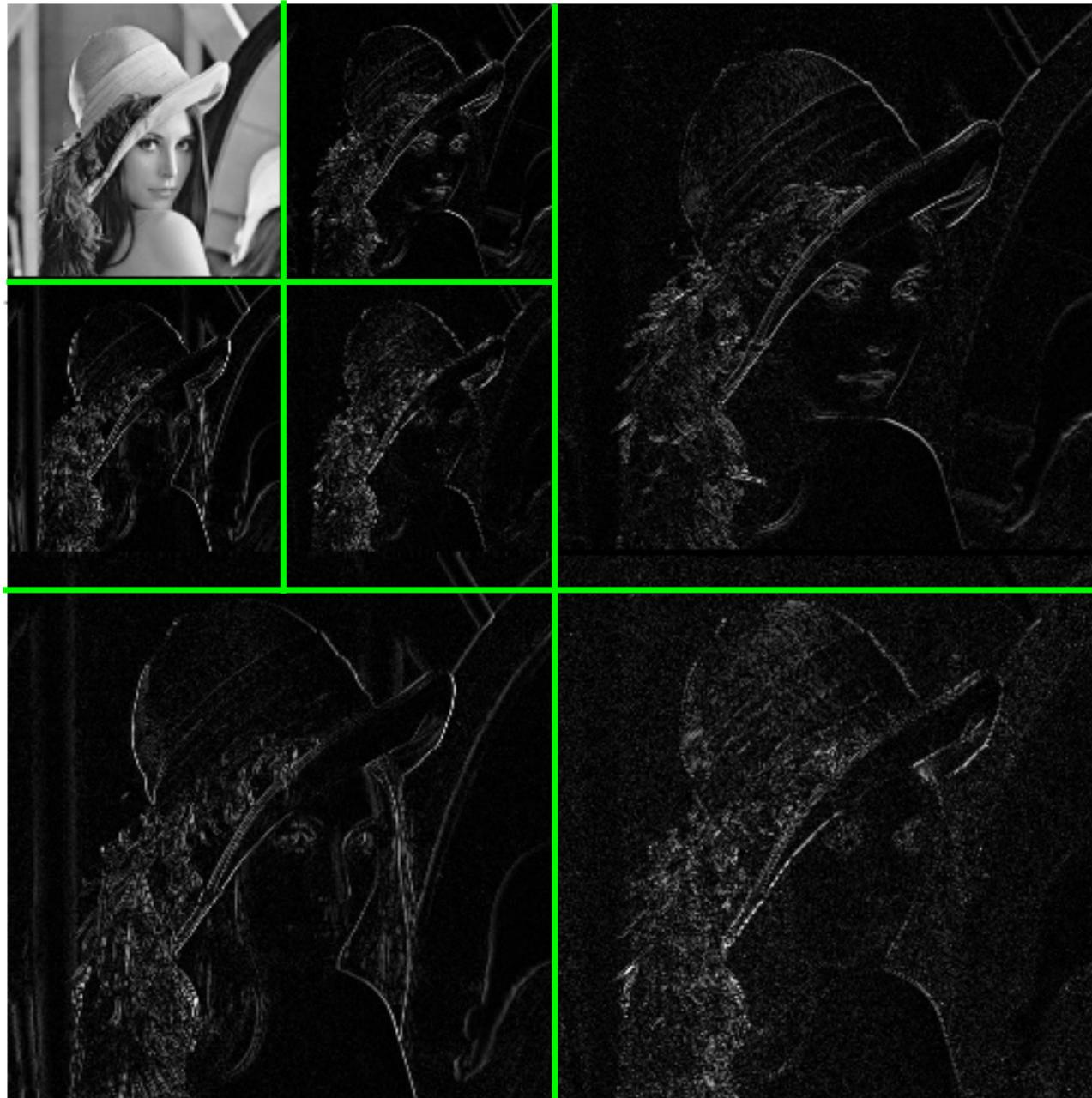


Entangled qubits



**Multi-Scale
Entanglement
Renormalization
Ansatz**

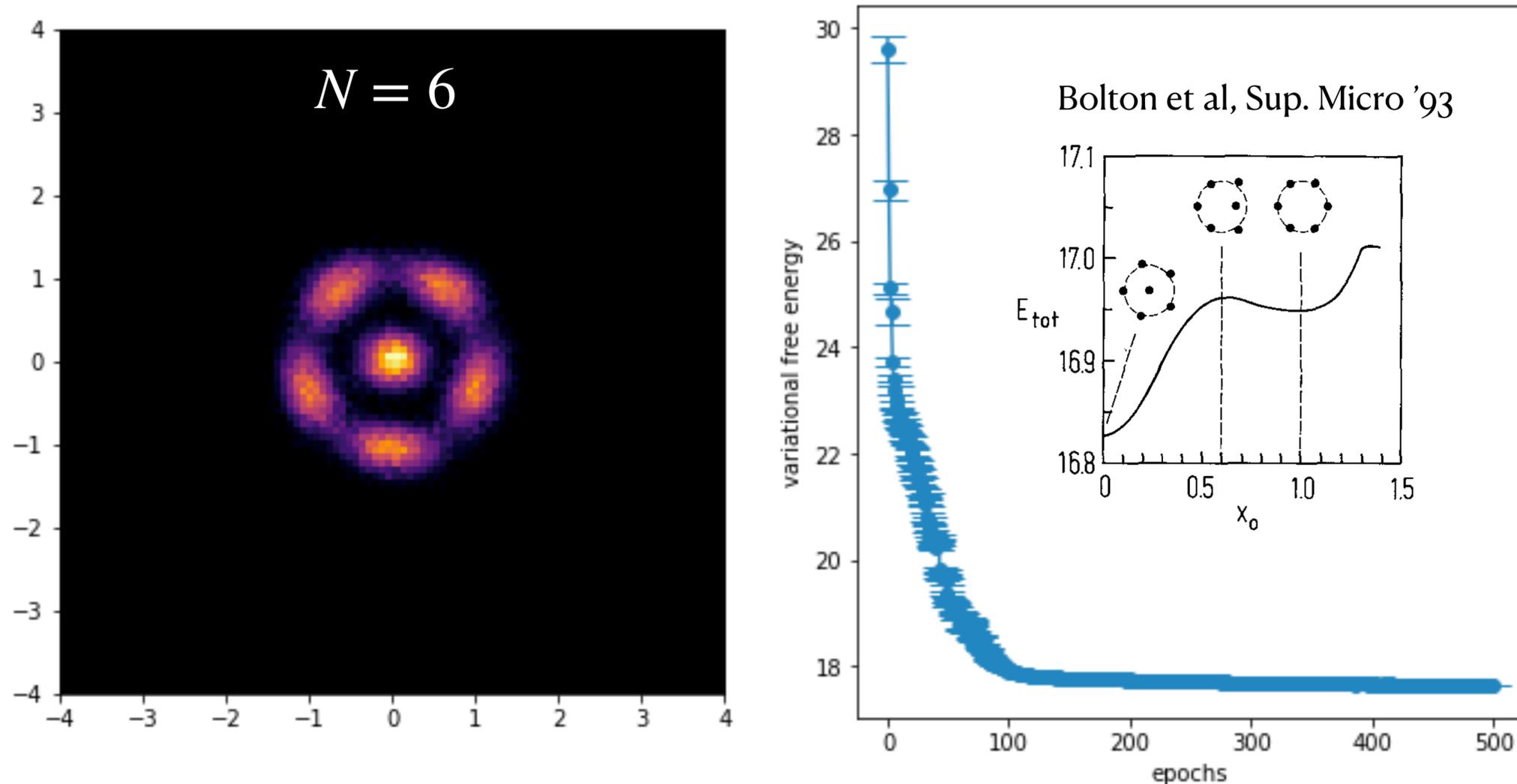
Connection to wavelets



Nonlinear & adaptive generalizations of wavelets

Demo: Classical Coulomb gas in a harmonic trap

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



Optimization: Monte Carlo Gradient Estimators

Review: 1906.10652

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

Reinforcement learning

Variational inference

Variational Monte Carlo

Variational quantum algorithms

Score function estimator (REINFORCE)

...

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

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

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

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 ?

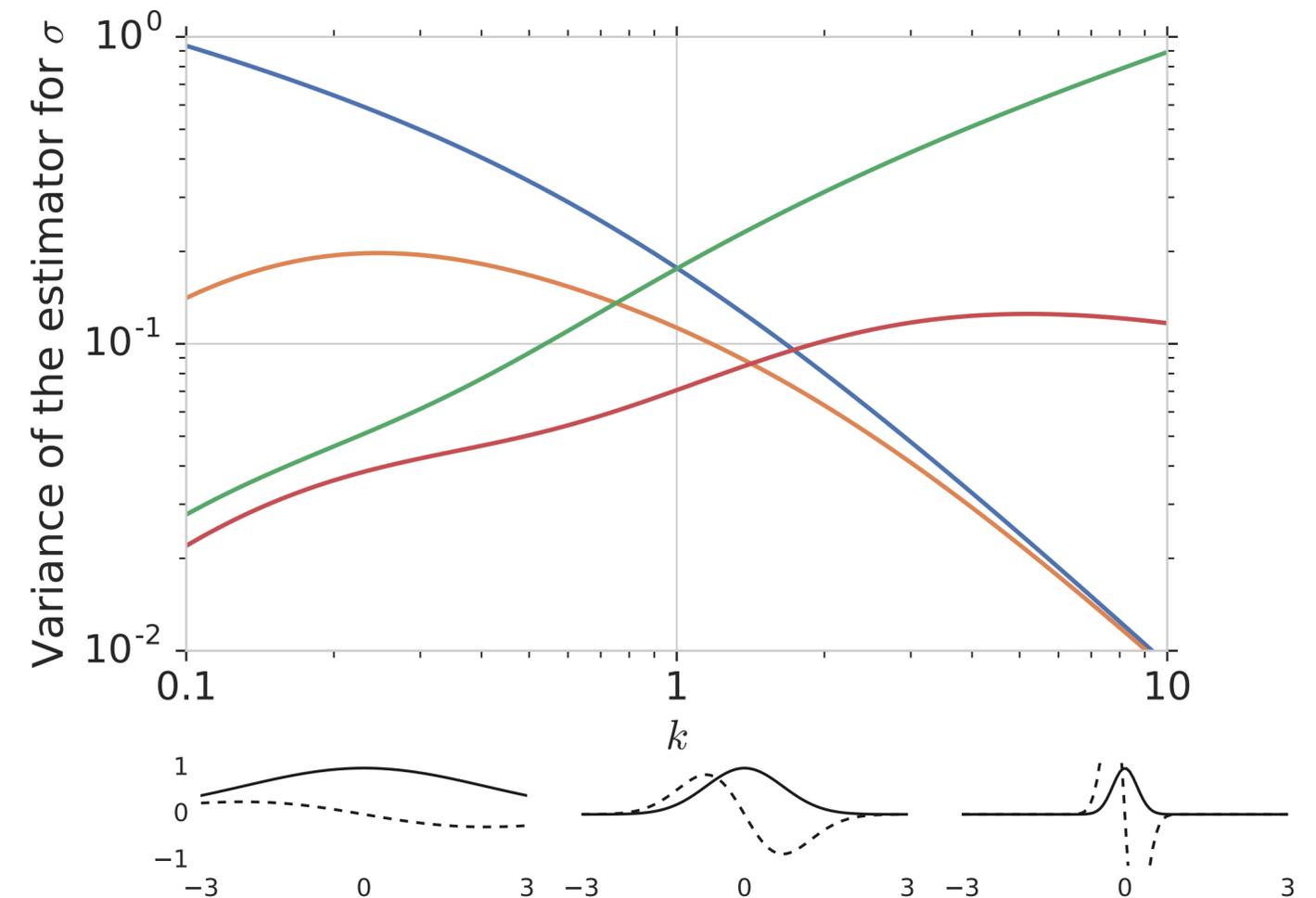
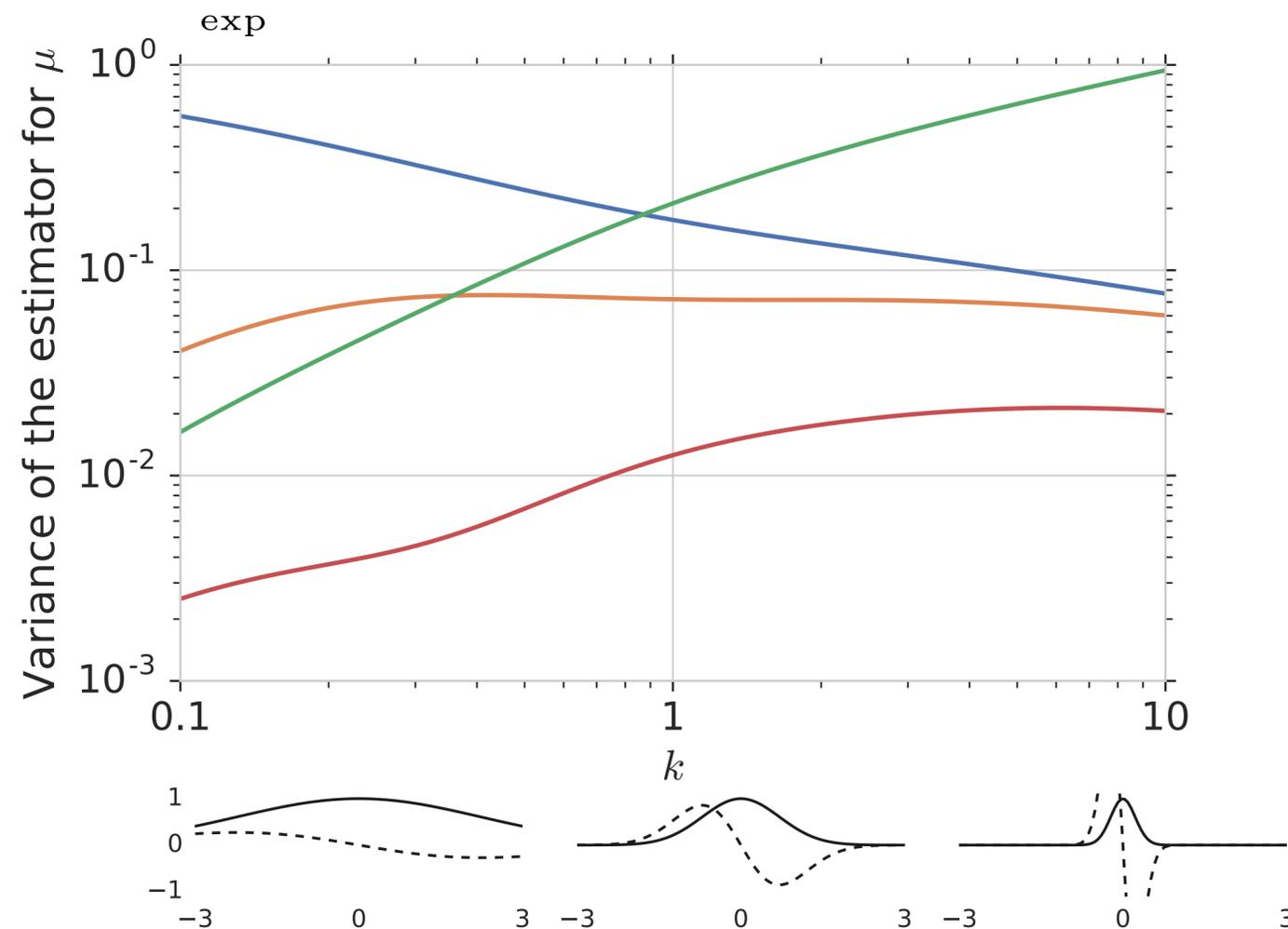
More discussions

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



Continuous normalizing flows

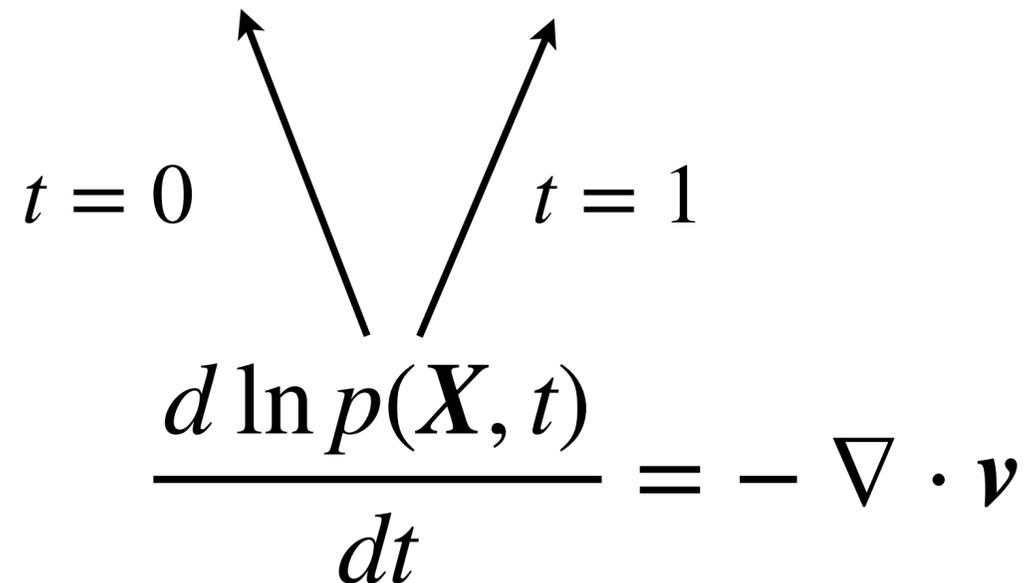
$$\ln p(\mathbf{X}) = \ln \mathcal{N}(\mathbf{Z}) - \ln \left| \det \left(\frac{\partial \mathbf{X}}{\partial \mathbf{Z}} \right) \right|$$

Consider infinitesimal change-of-variables Chen et al 1806.07366

$$\mathbf{X} = \mathbf{Z} + \varepsilon \mathbf{v} \quad \ln p(\mathbf{X}) - \ln \mathcal{N}(\mathbf{Z}) = - \ln \left| \det \left(1 + \varepsilon \frac{\partial \mathbf{v}}{\partial \mathbf{Z}} \right) \right|$$

$$\varepsilon \rightarrow 0$$

$$\frac{d\mathbf{X}}{dt} = \mathbf{v}$$

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


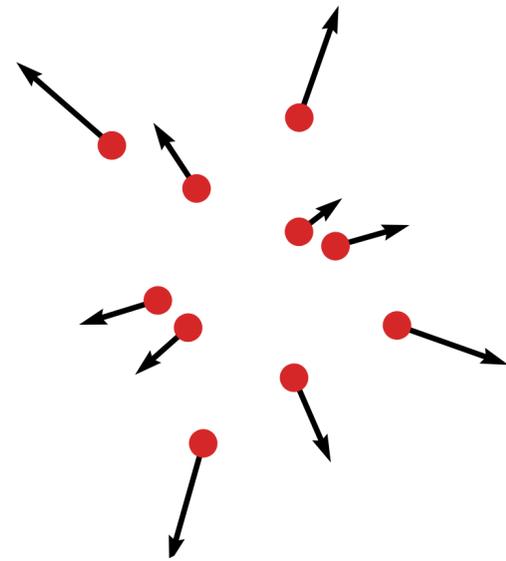
Fluid physics behind flows

Zhang, E, LW 1809.10188

[wangleiphy/MongeAmpereFlow](https://arxiv.org/abs/1809.10188)

$$\frac{dX}{dt} = \mathbf{v}$$

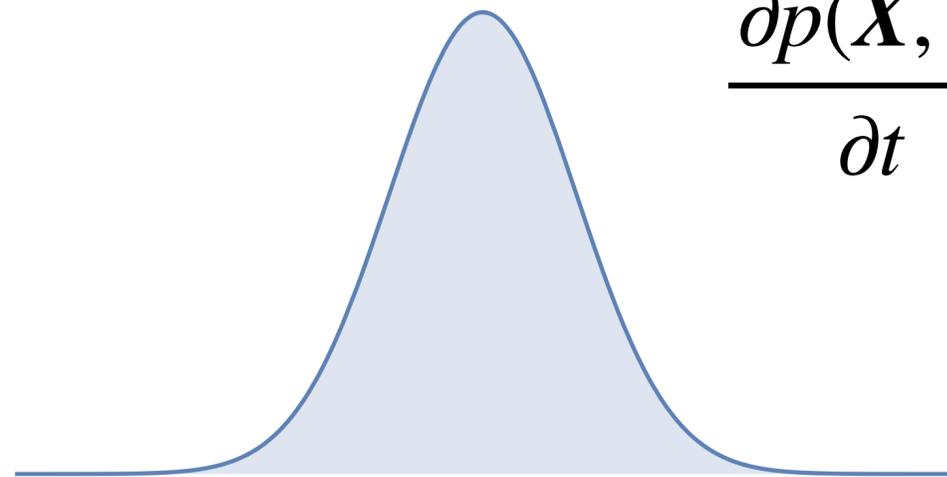
$$\frac{d \ln p(X, t)}{dt} = - \nabla \cdot \mathbf{v}$$



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

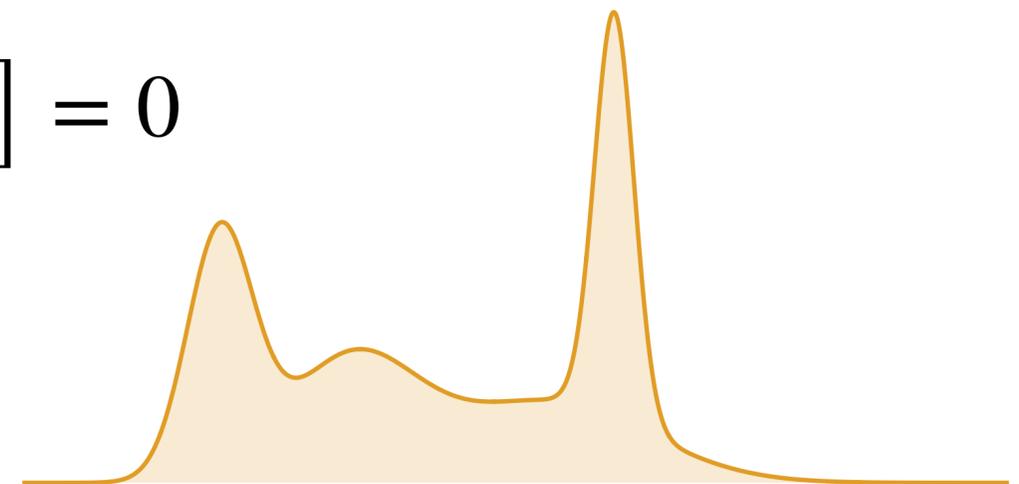
“material derivative”

Lagrangian v.s. Euler approach to fluid mechanics



Simple density

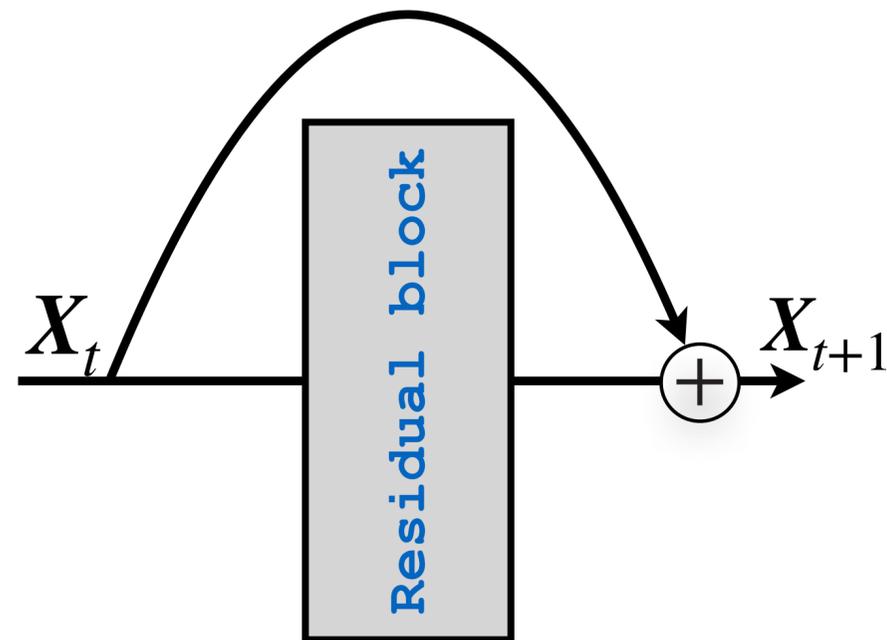
$$\frac{\partial p(X, t)}{\partial t} + \nabla \cdot [p(X, t)\mathbf{v}] = 0$$



Complex density

Neural Ordinary Differential Equations

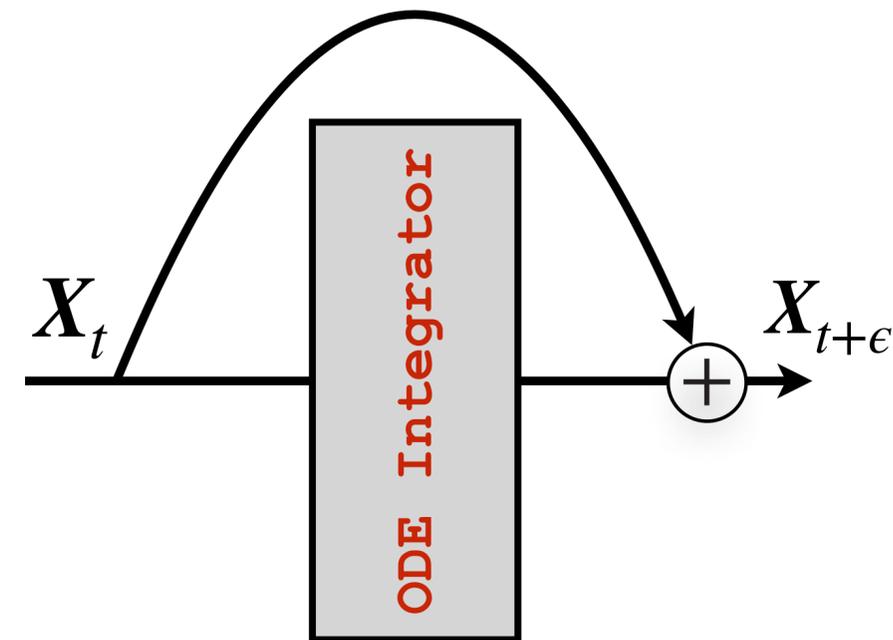
Residual network



$$X_{t+1} = X_t + v(X_t)$$

Chen et al, 1806.07366

ODE integration



$$dX/dt = v(X)$$

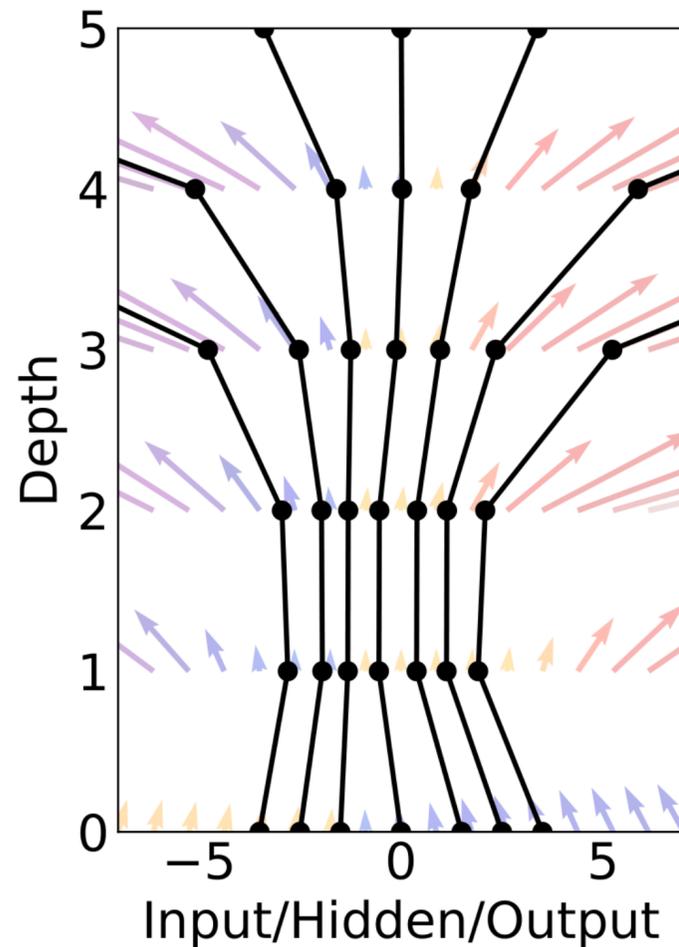
Harbor et al 1705.03341

Lu et al 1710.10121,

E Commun. Math. Stat 17'...

Neural Ordinary Differential Equations

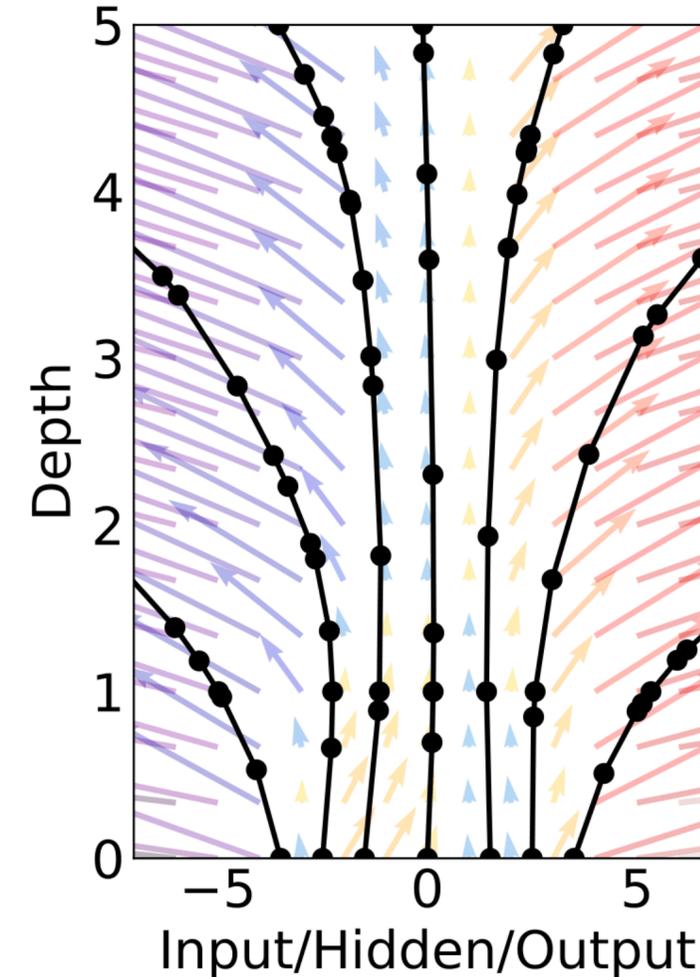
Residual network



$$X_{t+1} = X_t + \nu(X_t)$$

Chen et al, 1806.07366

ODE integration



$$dX/dt = \nu(X)$$

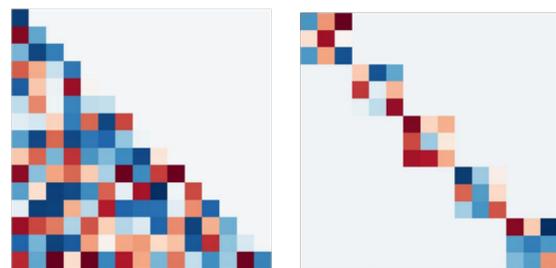
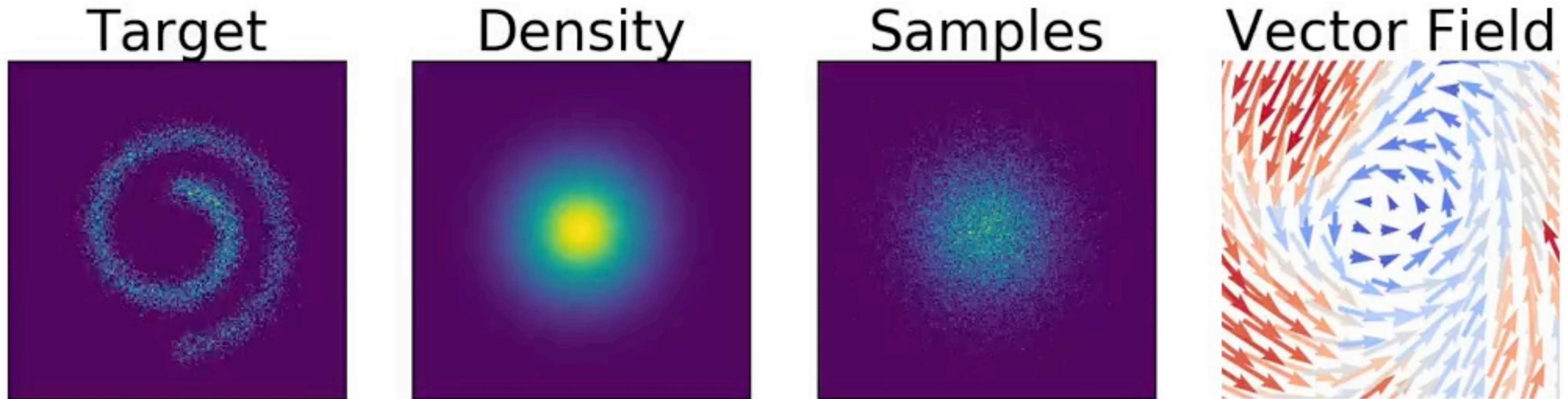
Harbor et al 1705.03341

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

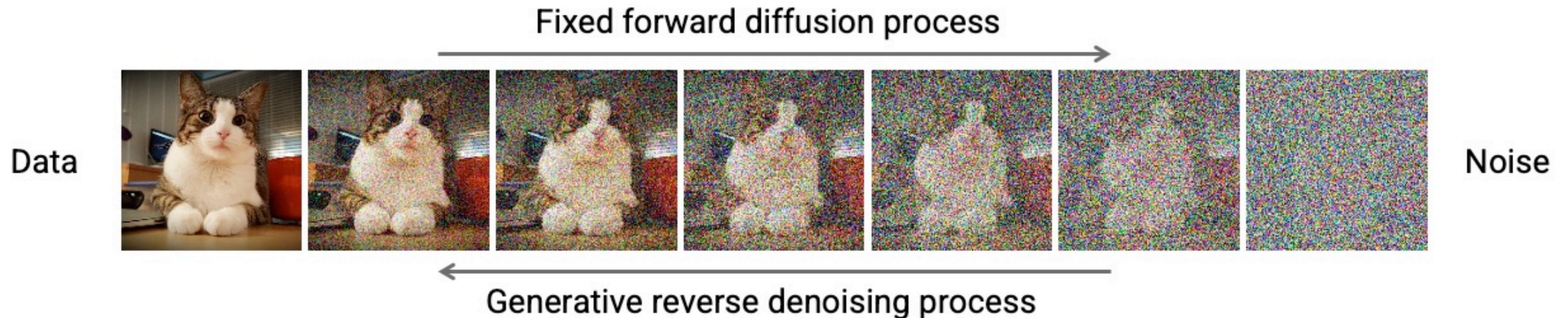
From flow to diffusion model

Continuity equation

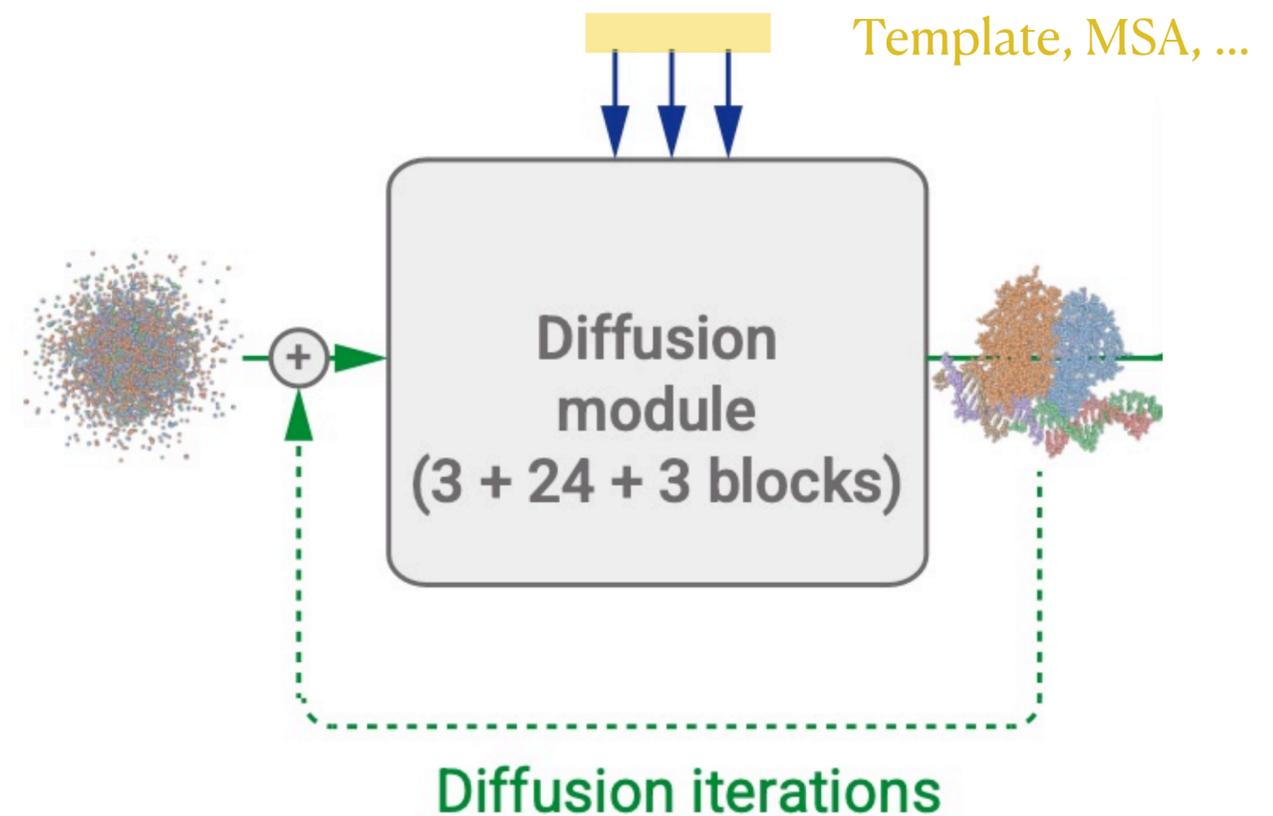
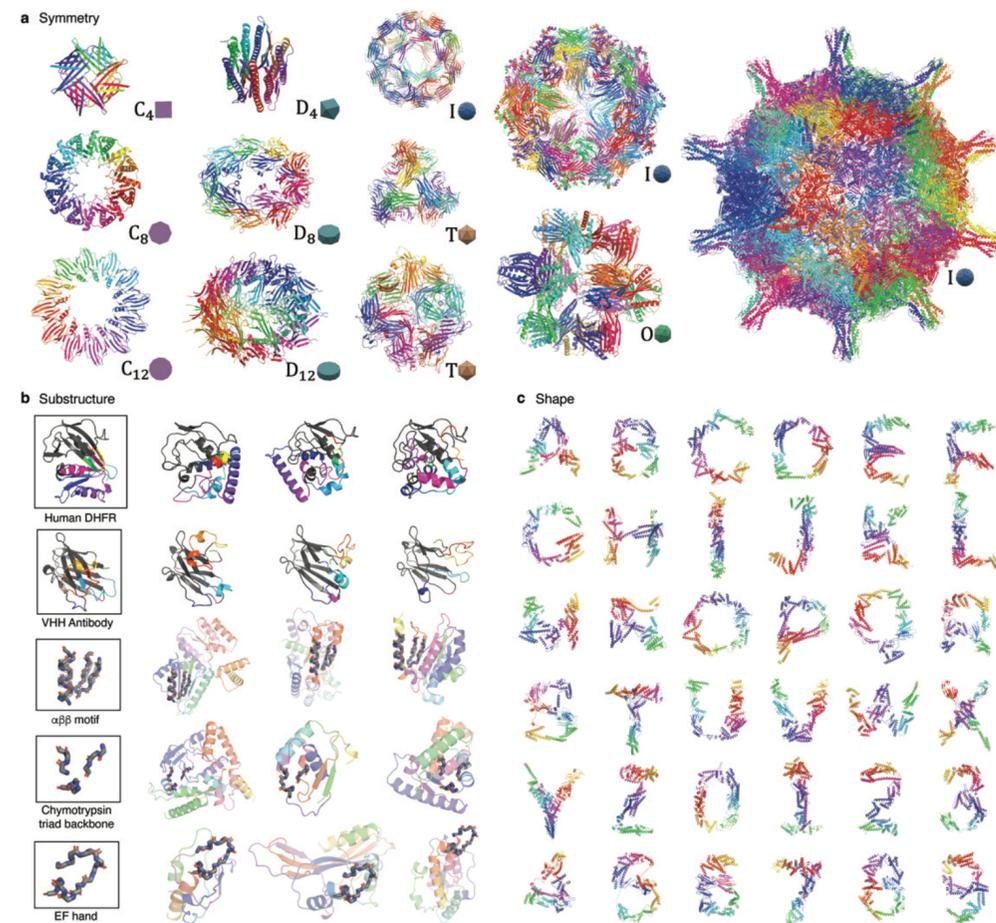
$$\frac{\partial p(\mathbf{X}, t)}{\partial t} + \nabla \cdot [p(\mathbf{X}, t)\mathbf{v}] = 0$$

Fokker-Planck equation

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



Diffusion models for protein structure prediction and design



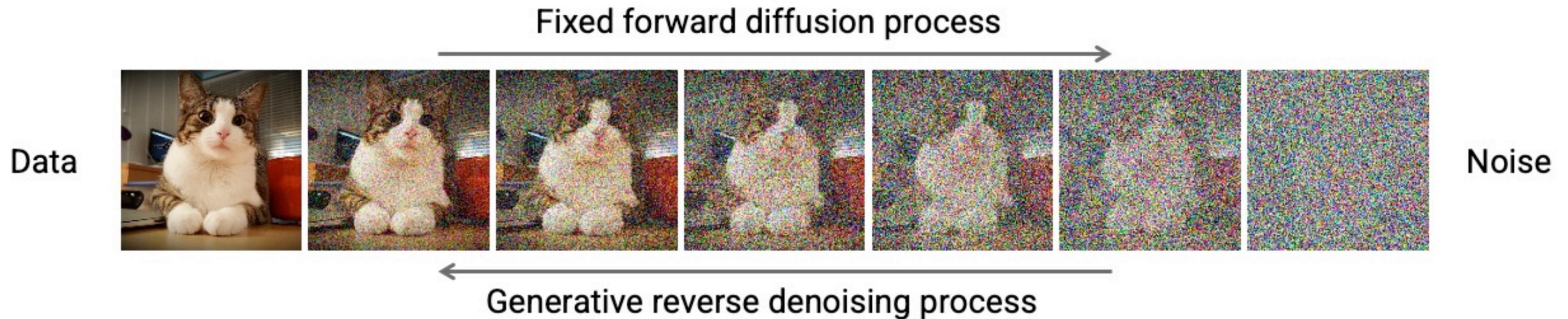
Ingraham et al, Chroma, Nature 2023
<https://generatebiomedicines.com/chroma>

Abramson et al, AlphaFold3, Nature 2024
<https://deepmind.google/technologies/alphafold/>

From flow to diffusion model, and back

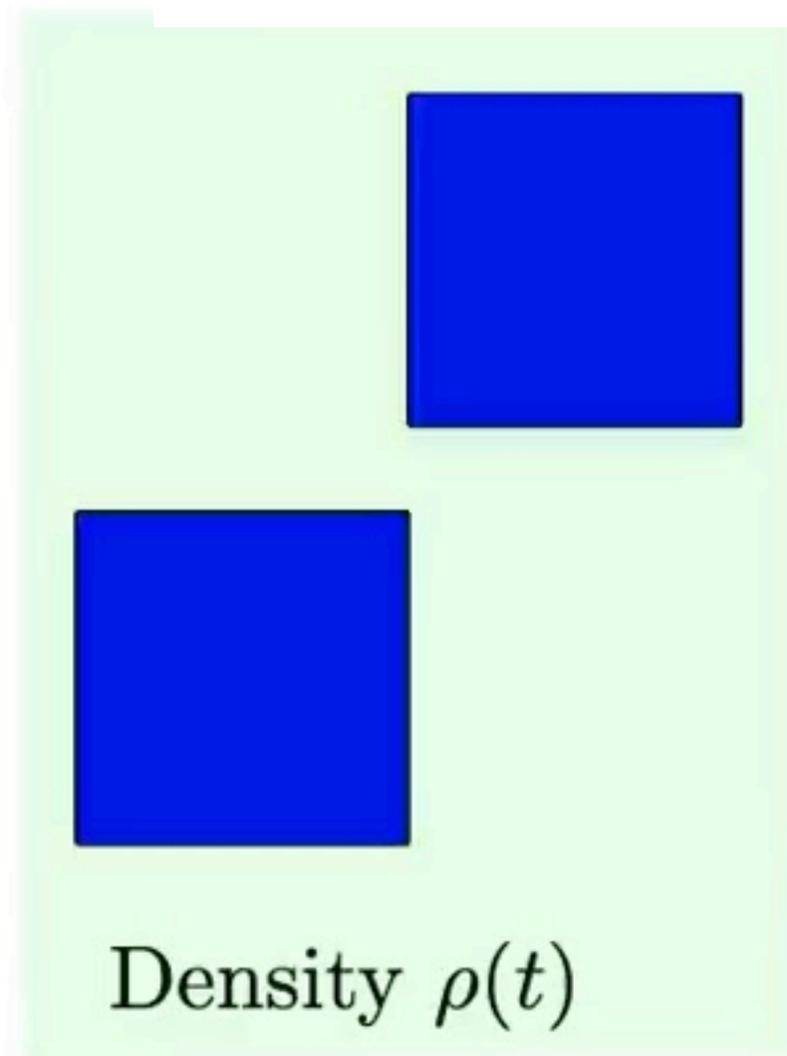
Continuity equation
$$\frac{\partial p(\mathbf{X}, t)}{\partial t} + \nabla \cdot [p(\mathbf{X}, t)\mathbf{v}] = 0$$

Fokker-Planck equation
$$\frac{\partial p(\mathbf{X}, t)}{\partial t} + \nabla \cdot \left[p(\mathbf{X}, t) (\mathbf{f} - \nabla \ln p(\mathbf{X}, t)) \right] = 0$$

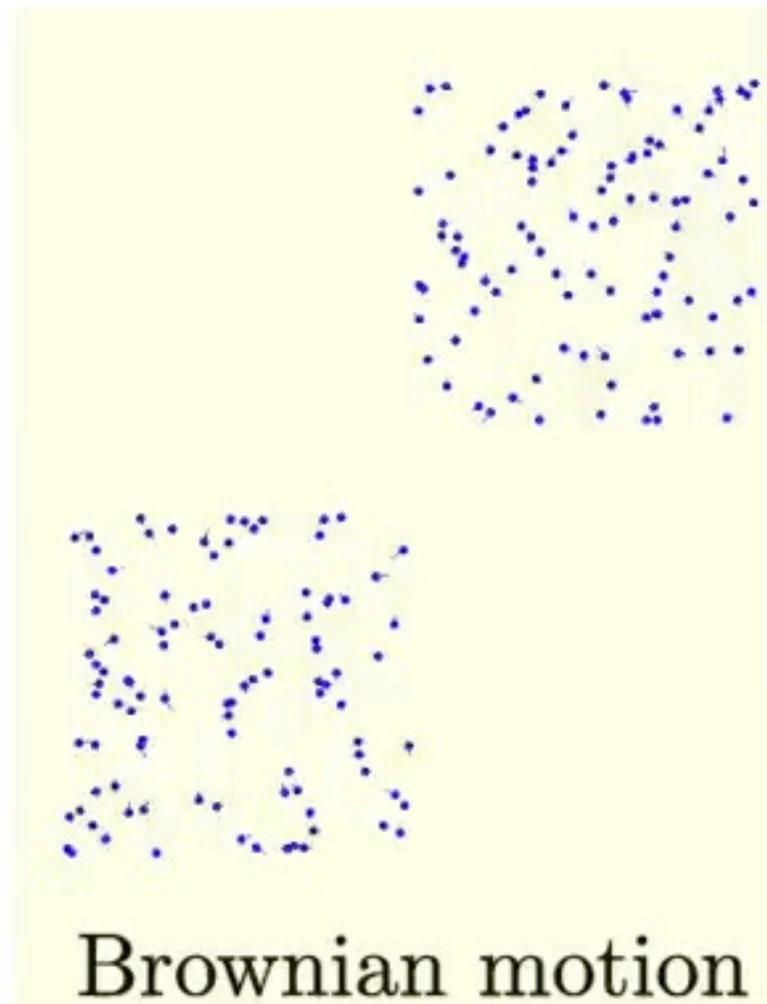


A tale of three equations

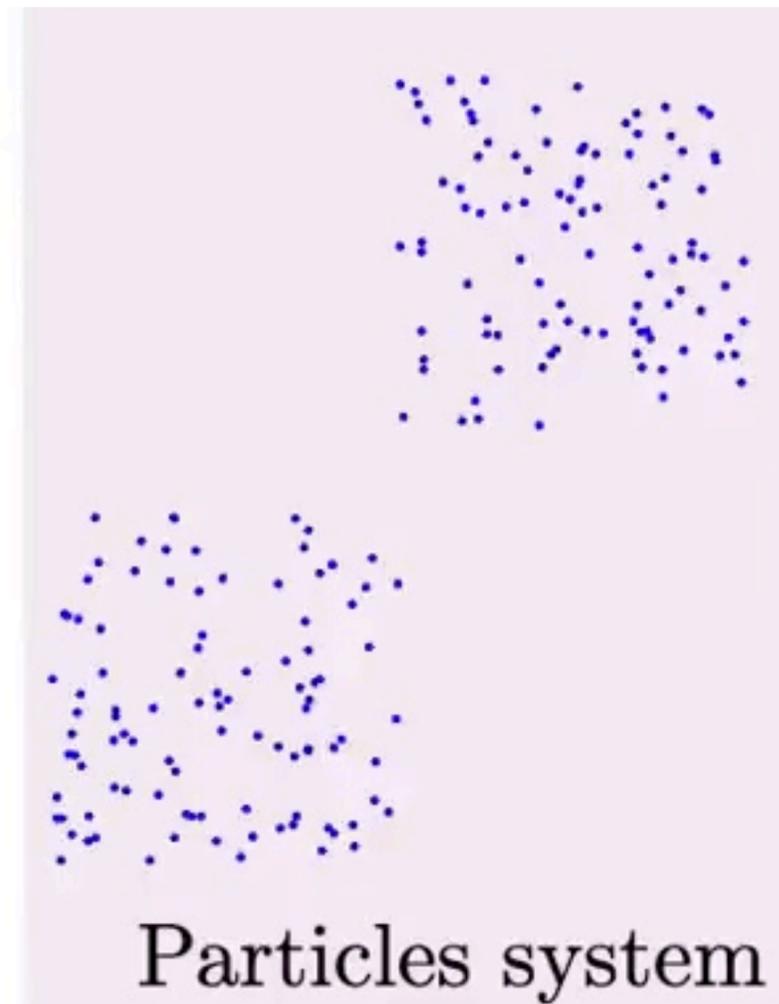
偏微分方程
Fokker-Planck

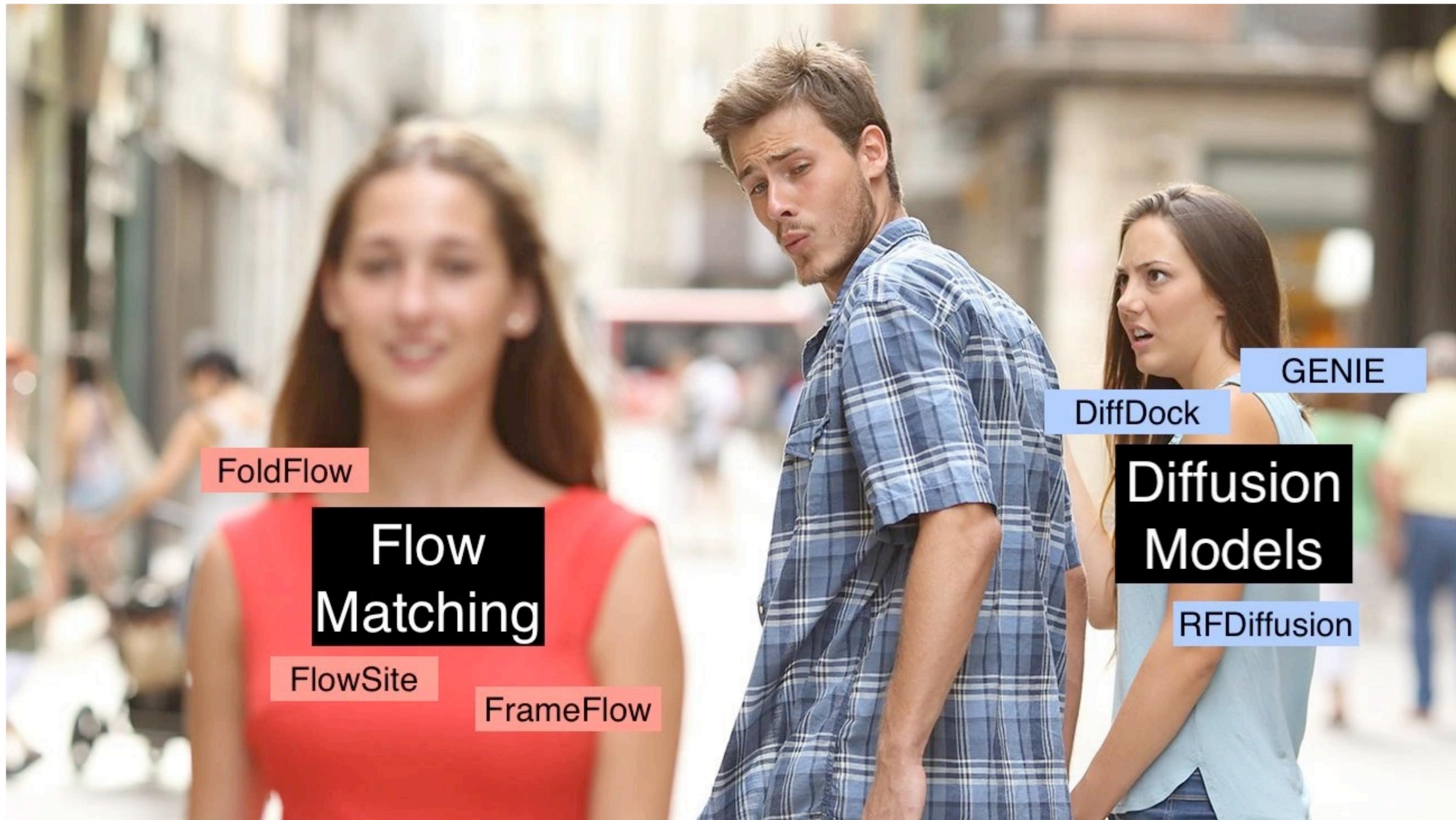


随机微分方程
Diffusion model



常微分方程
Flow model



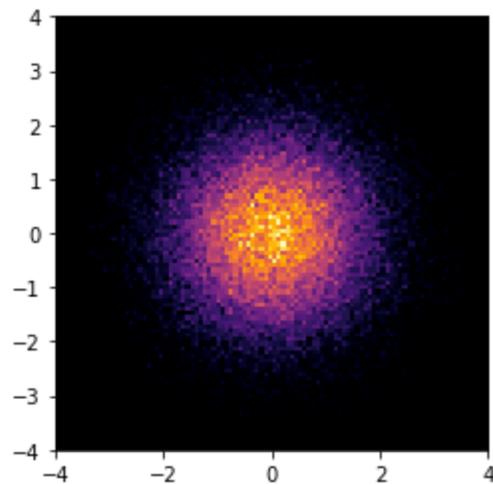


https://twitter.com/michael_galkin/status/1711845455817261409

Demo: bounding free energy of classical Coulomb gas

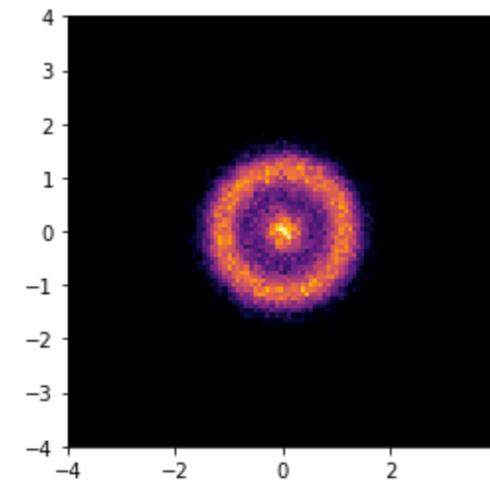
$$\mathcal{L} = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim \mathcal{N}(0,I)} \mathbb{E}_{\mathbf{x}_1 \sim \exp(-\beta E)/Z} \left| \mathbf{x}_1 - \mathbf{x}_0 - \mathbf{v}(\mathbf{x}, t) \right|^2$$

$$Z = \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[e^{-\beta E(\mathbf{x}) - \ln q(\mathbf{x})} \right] \quad \ln q(\mathbf{x}) = \ln \mathcal{N}(0,I) - \int_0^1 \nabla \cdot \mathbf{v} dt$$



Base density
Gaussian samples

← Interpolate samples to estimate free energy differences →



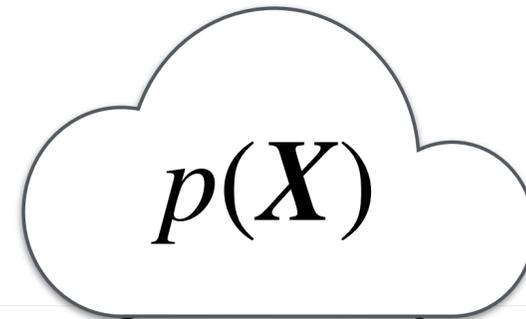
Target density
Monte Carlo samples

 <https://colab.research.google.com/drive/1t-Vk37Axxp040B7uXFUNlk-zeCC2lcX3?usp=sharing>

Jarzynski PRE '02, see also likelihood-based training of flows Wirnsberger et al, 2002.04913, 2111.08696

Generative models and their physics genes

Goodfellow,
NIPS tutorial, 1701.00160



Explicit density

Implicit density

Direct
GAN

Tractable density

Approximate density

Markov Chain
GSN

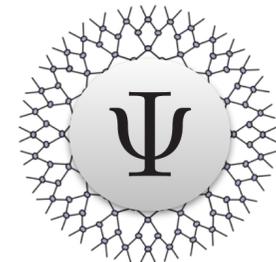
Variational

Markov Chain

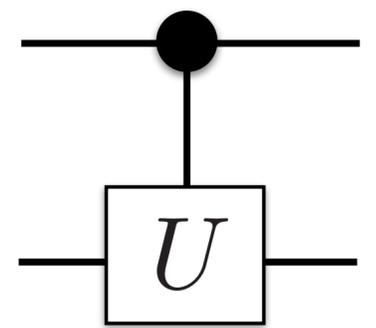
-Fully visible belief nets
-NADE
-MADE
-PixelCNN
-Change of variables models (nonlinear ICA)

Autoregressive model

Variational autoencoder Boltzmann machine + **Diffusion models**



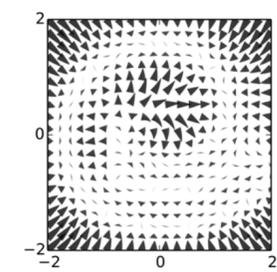
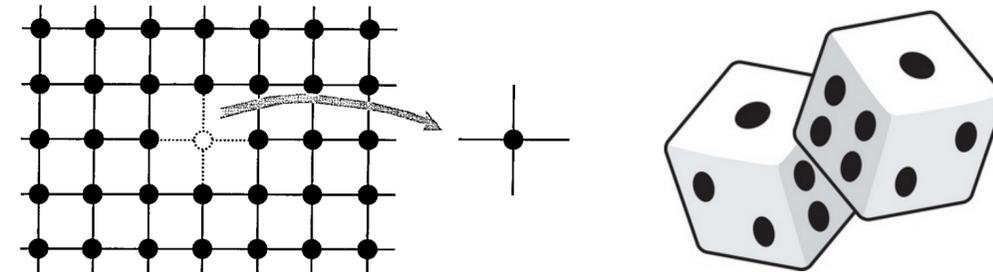
Tensor Networks
Han et al, PRX '18



Quantum Circuits
Liu et al PRA '18



Flow model



Which is the best ?



Jeremy Howard @jeremyphoward · Apr 15

"transformers or diffusion"

Ummm... who's gonna tell him that that makes no sense at all?

Autoregressive model

$$p(X) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2)\dots$$



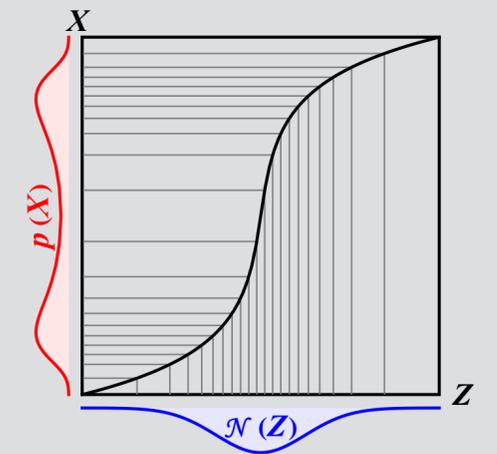
"... the murderer is _____"

$p(_ | \dots)$

Discrete or continuous sequential data

Flow model

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



Continuous structured data

Which is the best ?

The “it” in AI models is the dataset.

Posted on June 10, 2023 by jbetker

I’ve been at OpenAI for almost a year now. In that time, I’ve trained a **lot** of generative models. More than anyone really has any right to train. As I’ve spent these hours observing the effects of tweaking various model configurations and hyperparameters, one thing that has struck me is the similarities in between all the training runs.

It’s becoming awfully clear to me that these models are truly approximating their datasets to an incredible degree. What that means is not only that they learn what it means to be a dog or a cat, but the interstitial frequencies between distributions that don’t matter, like what photos humans are likely to take or words humans commonly write down.

What this manifests as is – trained on the same dataset for long enough, pretty much every model with enough weights and training time converges to the same point. Sufficiently large diffusion conv-unets produce the same images as ViT generators. AR sampling produces the same images as diffusion.

This is a surprising observation! It implies that model behavior is not determined by architecture, hyperparameters, or optimizer choices. It’s determined by your dataset, nothing else. Everything else is a means to an end in efficiently delivery compute to approximating that dataset.

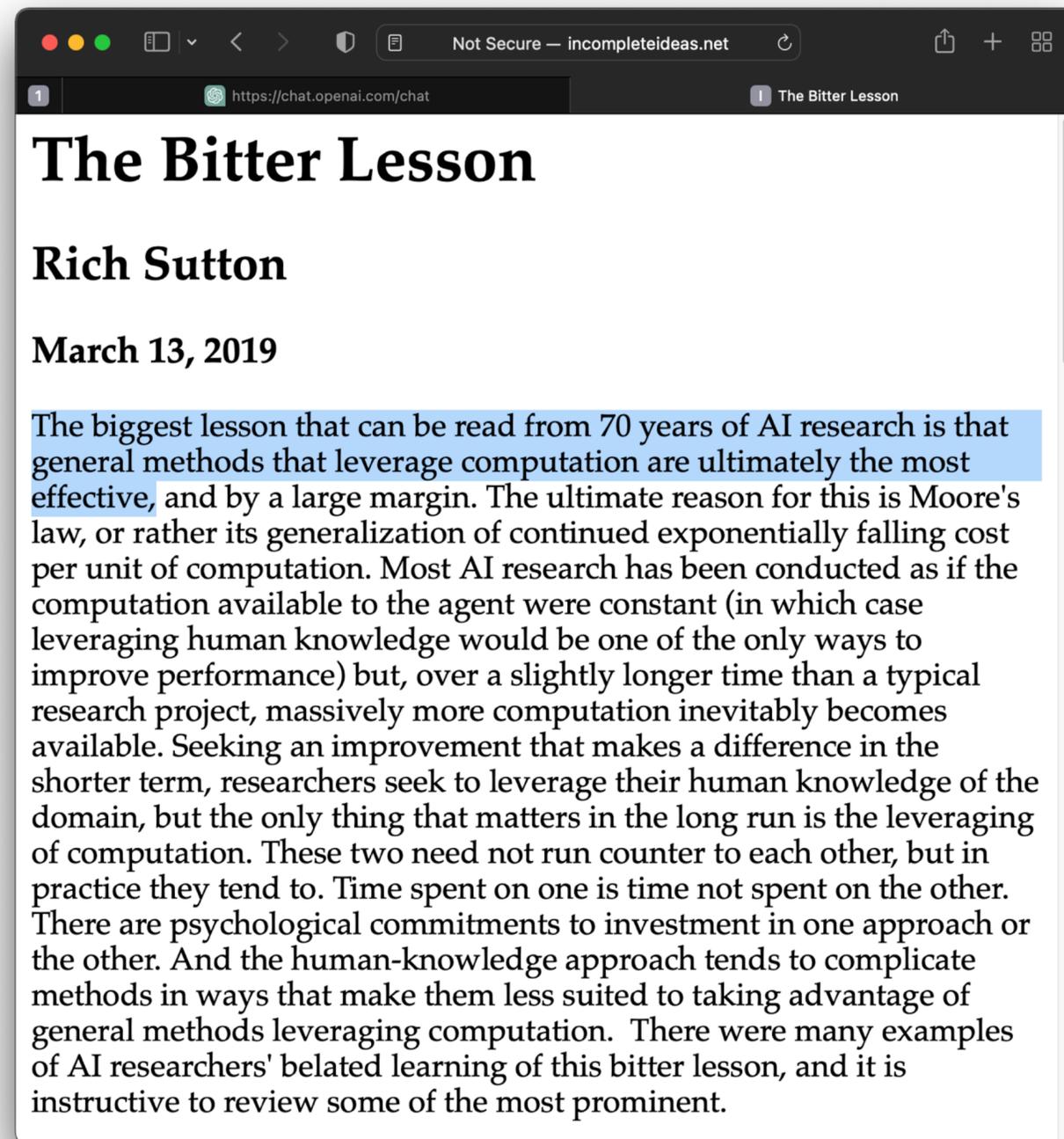
Then, when you refer to “Lambda”, “ChatGPT”, “Bard”, or “Claude” then, it’s not the model weights that you are referring to. It’s the dataset.

<https://nonint.com/2023/06/10/the-it-in-ai-models-is-the-dataset/>

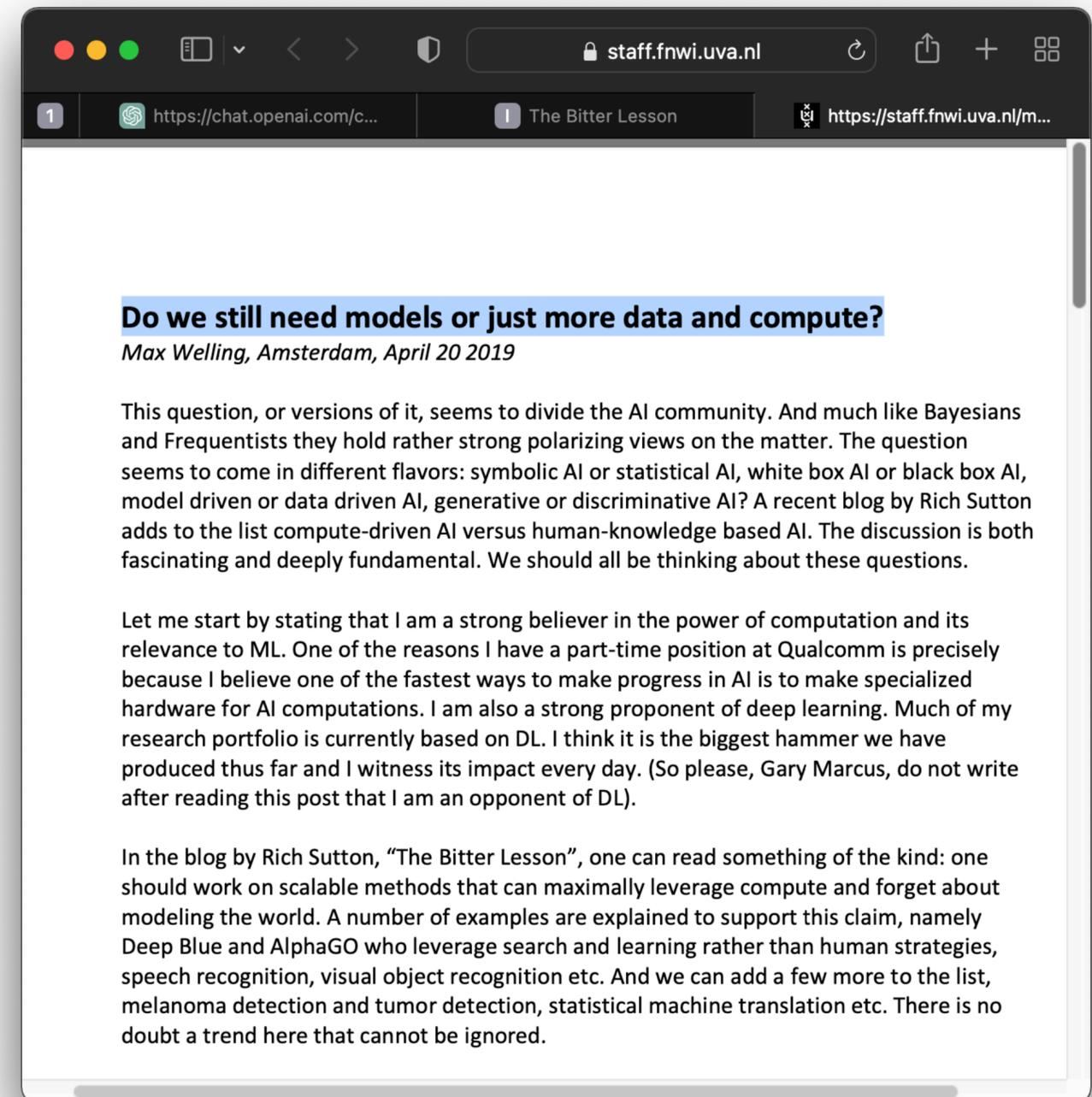
How much inductive bias?

<http://www.incompleteideas.net/IncIdeas/BitterLesson.html>

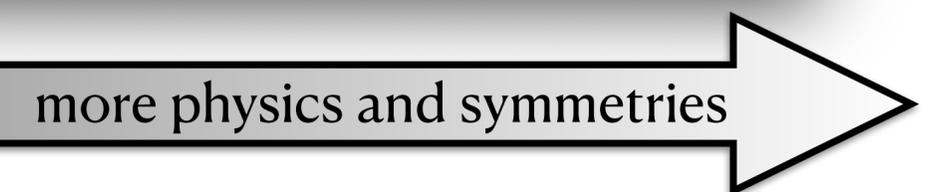
<https://staff.fnwi.uva.nl/m.welling/wp-content/uploads/Model-versus-Data-AI-1.pdf>



The screenshot shows a web browser window with the address bar displaying 'https://chat.openai.com/chat'. The page title is 'The Bitter Lesson'. The author is 'Rich Sutton' and the date is 'March 13, 2019'. The main text of the article is visible, with a blue highlight over the first sentence: 'The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. The ultimate reason for this is Moore's law, or rather its generalization of continued exponentially falling cost per unit of computation. Most AI research has been conducted as if the computation available to the agent were constant (in which case leveraging human knowledge would be one of the only ways to improve performance) but, over a slightly longer time than a typical research project, massively more computation inevitably becomes available. Seeking an improvement that makes a difference in the shorter term, researchers seek to leverage their human knowledge of the domain, but the only thing that matters in the long run is the leveraging of computation. These two need not run counter to each other, but in practice they tend to. Time spent on one is time not spent on the other. There are psychological commitments to investment in one approach or the other. And the human-knowledge approach tends to complicate methods in ways that make them less suited to taking advantage of general methods leveraging computation. There were many examples of AI researchers' belated learning of this bitter lesson, and it is instructive to review some of the most prominent.'



The screenshot shows a web browser window with the address bar displaying 'staff.fnwi.uva.nl'. The page title is 'Do we still need models or just more data and compute?'. The author is 'Max Welling, Amsterdam, April 20 2019'. The main text of the article is visible, discussing the question of whether more data and compute are sufficient for AI progress, contrasting it with the 'Bitter Lesson'.



Molecule representations and inductive biases

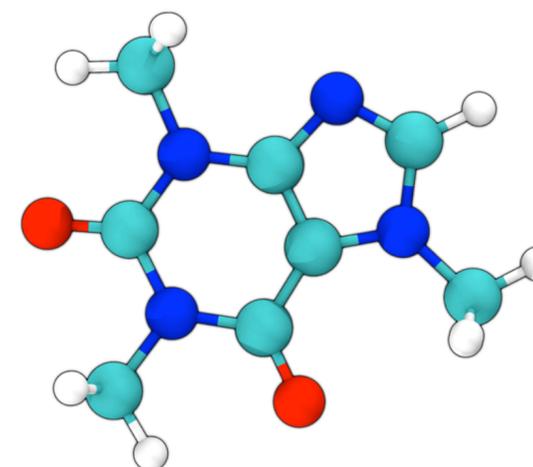
SMLIES or XYZ/CIF/PDB file

```
HEADER      CAFFEINE
COMPND      CAFFEINE
AUTHOR      GENERATED BY ChatGPT

ATOM        1  C   LIG A   1         0.000  0.000  0.000
ATOM        2  N   LIG A   1         1.289  0.000  0.000
ATOM        3  C   LIG A   1         1.463  1.192  0.000
ATOM        4  N   LIG A   1         2.644  1.192  0.000
ATOM        5  C   LIG A   1         2.818  2.384  0.000
ATOM        6  N   LIG A   1         4.089  2.384  0.000
ATOM        7  C   LIG A   1         4.263  3.576  0.000
          .
          .
          .
```

Language model

Atom coordinates



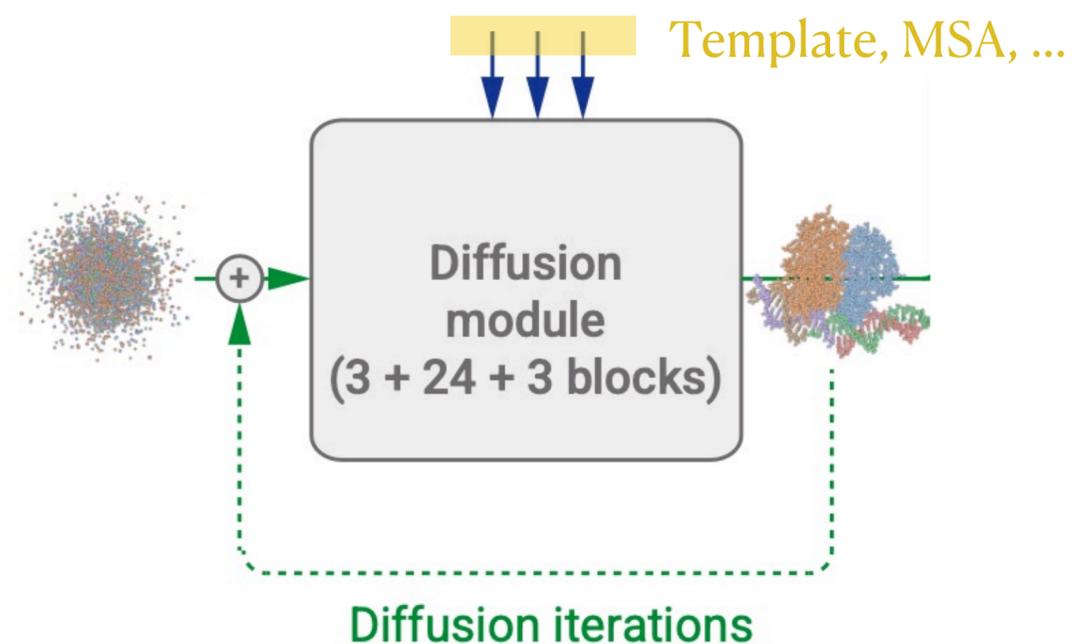
Equivariant neural network

more data and compute

more physics and symmetries

How much inductive bias?

Abramson et al, AlphaFold3, Nature 2024



Similarly to some recent work³⁵, we find that no invariance or equivariance with respect to global rotations and translation of the molecule are required in the architecture and so we omit them to simplify the machine learning architecture.

[35] **Swallowing the Bitter Pill**: Simplified Scalable Conformer Generation, 2311.17932

 **Yuyang Wang**
@YuyangW95

Glad to see AF3 coming out! Interesting to see it acknowledged our Molecular Conformer Fields (arxiv.org/abs/2311.17932) and omits the equivariant architecture as well. Another evidence of why equivariant design may not be a strong requirement in modeling molecules.

 **Thomas Kipf** ✓
@tkipf

Since this tweet sparked quite a bit of lively discussion, I'd like to add a bit more nuance:

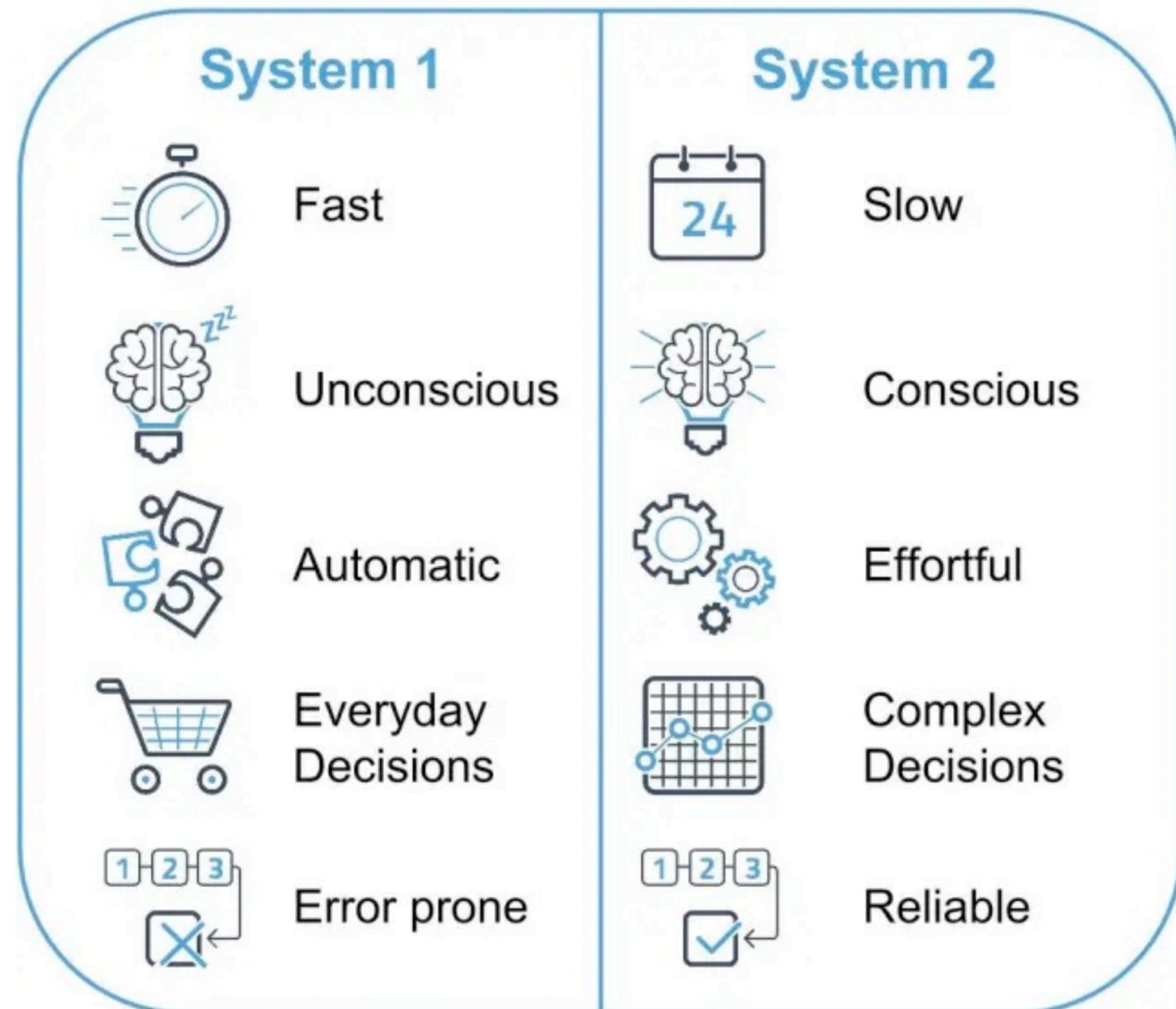
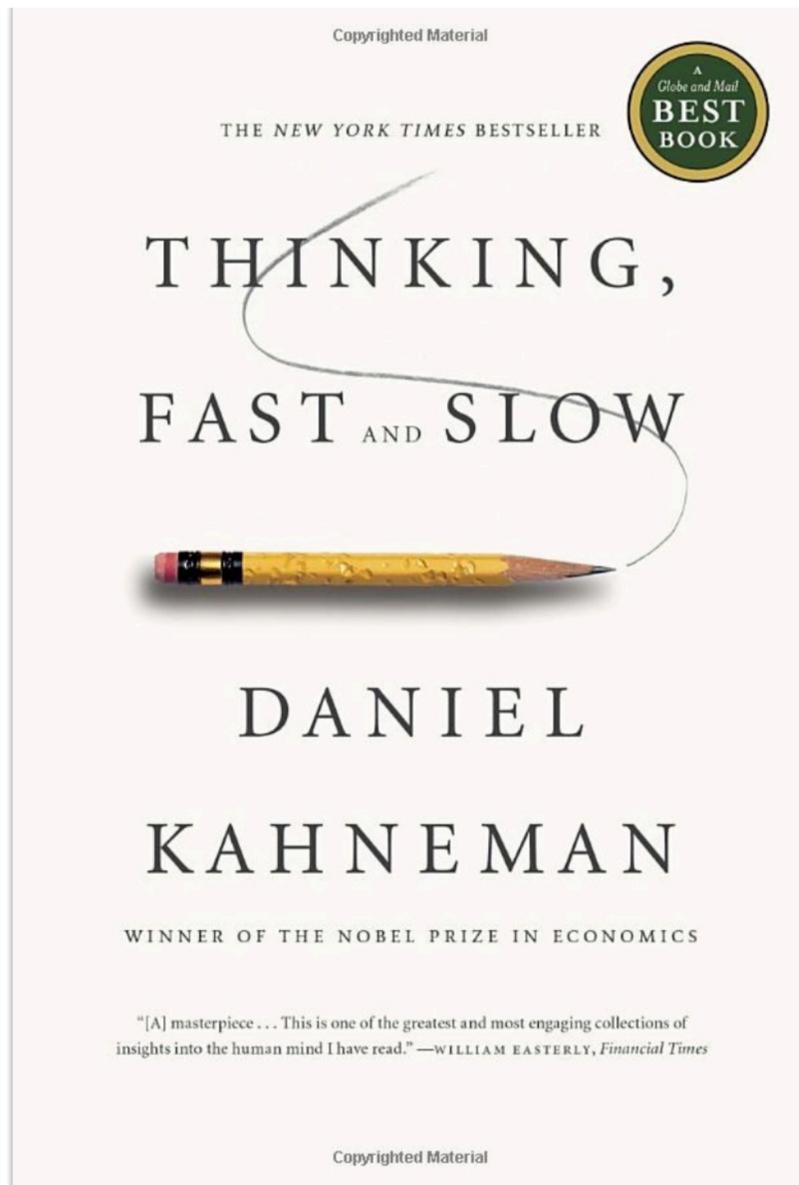
1) I think we absolutely should study symmetry in the context of (scalable) ML; this particular result only reinforces this IMO. Understanding trade-offs w.r.t. symmetry group "size", feasibility of data augmentation, dataset size, position encoding, ease of optimization, OOD generalization, etc. will likely be key for progress.

2) Just to (re-)state the obvious: Some of the most impactful architectures in ML (Transformers/GNNs) *are built on permutation symmetry* (and ofc there's something to be said about CNNs and translation symmetry). Having permutation symmetry in your data and not making use of it is typically not a good idea.

I don't think we should think in terms of symmetry inductive bias vs. scale, but rather how we can reap the benefits of both. [No bitter lesson here.](#)

Can they reason ?

Autoregressive language models are fast thinkers



Fast thinkers rely on good intuitions

物理直觉是如何养成的？ 徐一鸿：思考比计算要困难得多

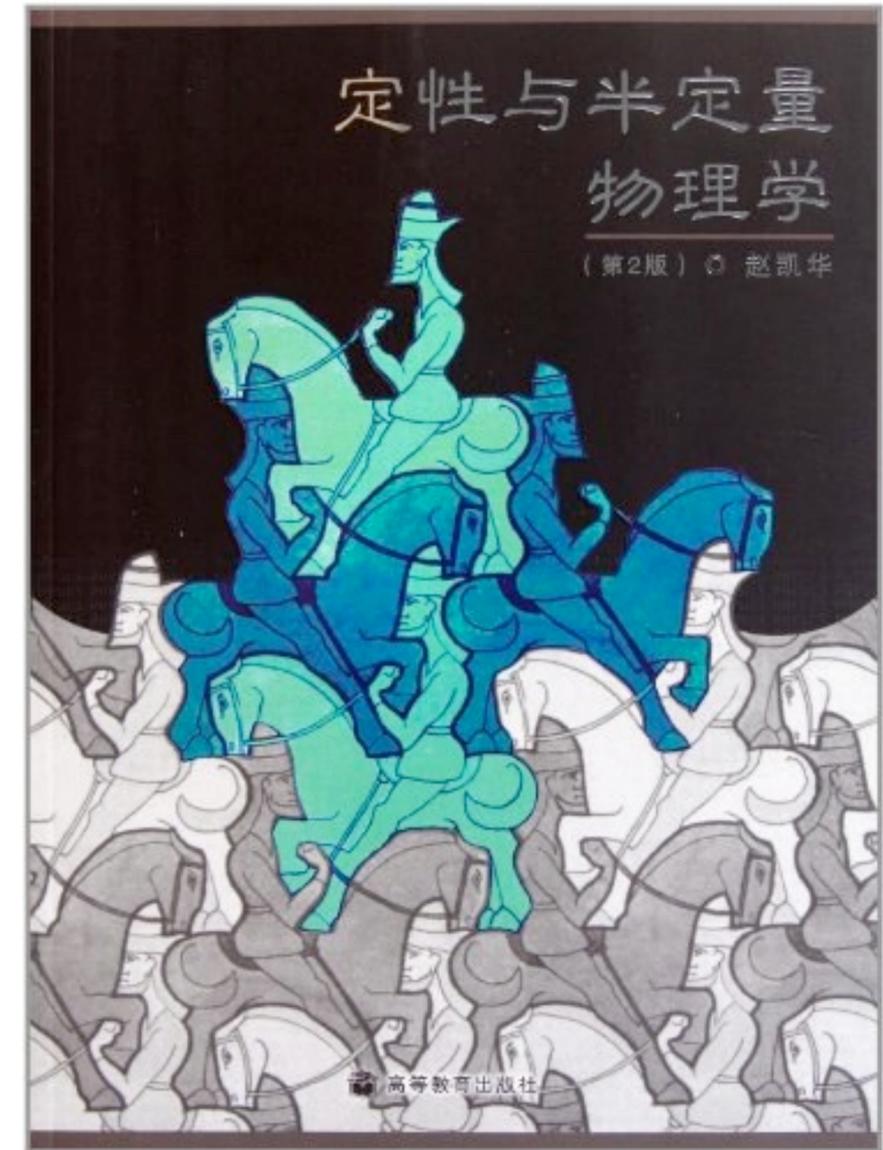
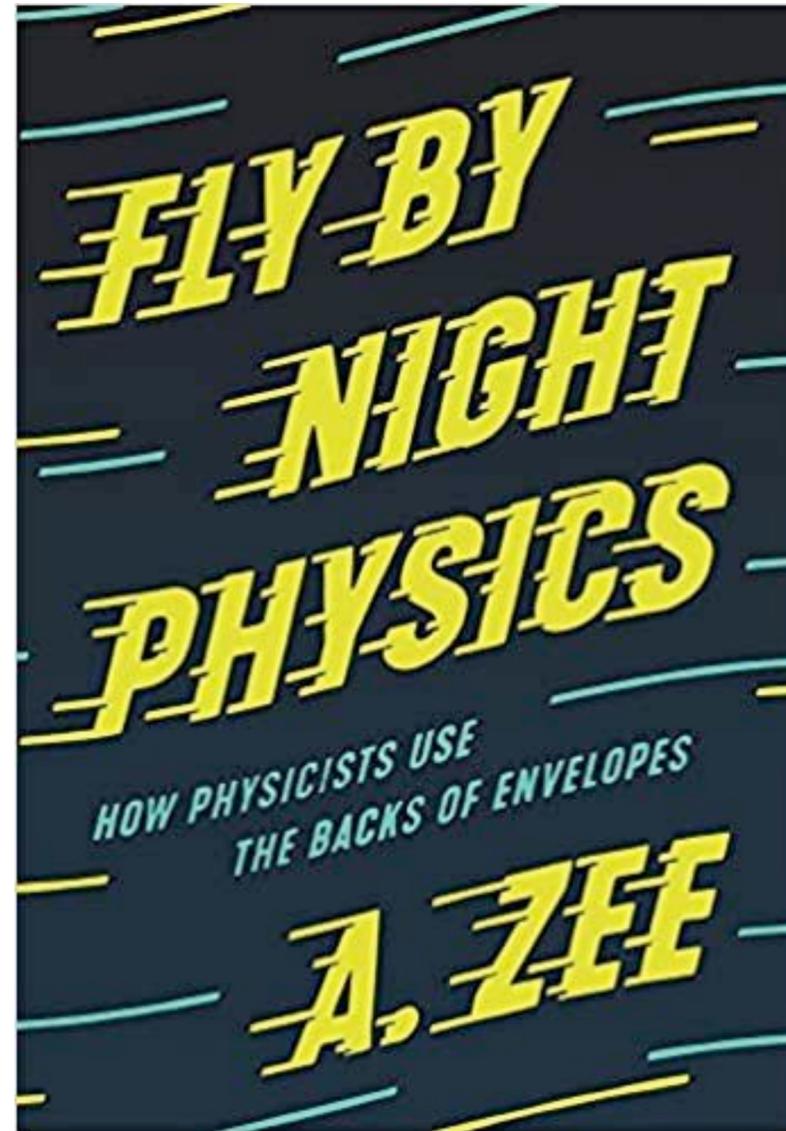
澎湃新闻记者 曹年润
2022-12-14 08:28 来源：澎湃新闻

· “重要的是要帮助学生培养更多的物理直觉，这也是我写这本书的原因之一。如何培养物理直觉？我在书里提到了两种可能性：天生就拥有物理直觉，或者通过不断练习发展物理直觉。”



徐一鸿教授。

“你应该先思考。思考比计算要困难得多，大多数人都可以坐下来计算，但不计算就思考问题，这是极其困难的。”



System 1 thinking in physics:

getting answers quickly without lengthy calculations

“Never never calculate unless you already know the answer!”—John Wheeler

Do they understand physics?



<https://openai.com/index/sora/>



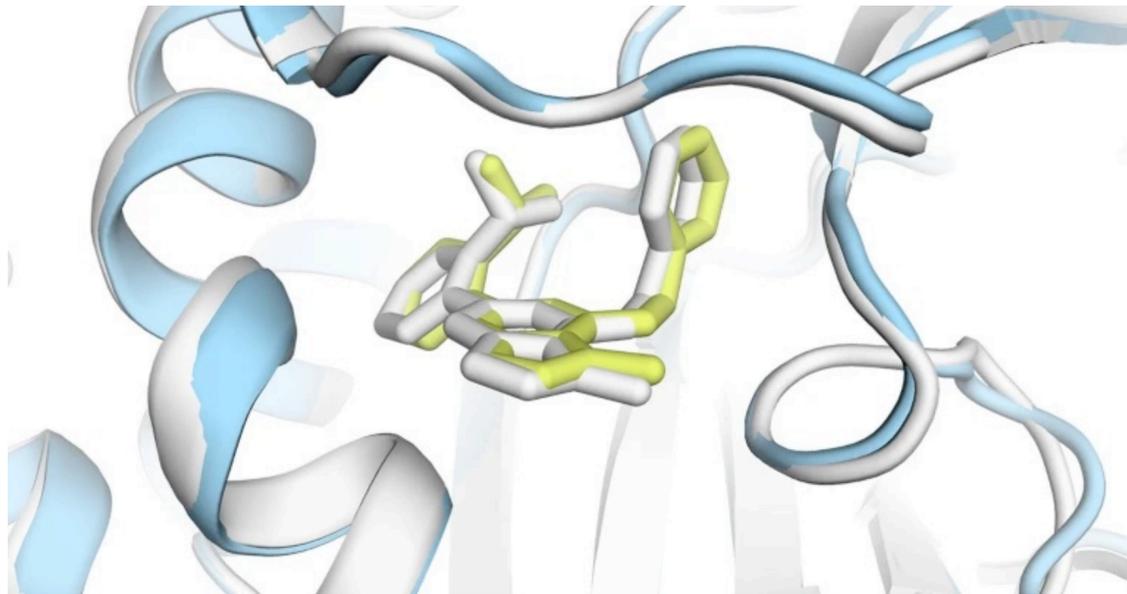
<http://ai.ruc.edu.cn/newslist/newsdetail/20240326001.html>

“What I can not create, I do not understand”
—Richard Feynman

Fold by intuition vs fold by equation

Both integrate Langevin dynamics $X_{t+1} = X_t + \frac{\eta_t}{2}s(X_t, t) + \sqrt{\eta_t}\epsilon$

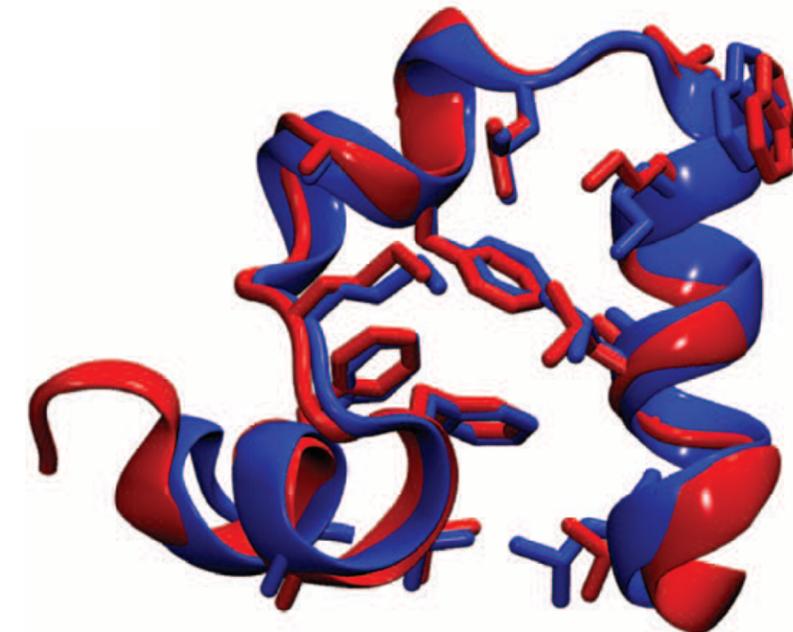
Data-driven generation AlphaFold3



Abramson et al, Nature 2024

The diffusion model may generate right conformations via unphysical pathways

Physics-based molecular dynamics



Shaw et al, Science 2010

Physical force fields may face difficulties in sampling rough energy landscapes

生成模型四问

① 哪种模型好?

大数据：无所谓
小数据：看模态

② 要多少先验?

抓主要矛盾

③ 能不能推理?

哪怕不能，“直觉”也很可贵

④ 懂不懂物理?

不是真的懂，但未必是坏事

What can generative models do ?

Appreciation

Likelihood estimation

$$\ln p(X)$$

Anomaly detection

Generation

(Un)conditional sampling

$$X \sim p(X) \text{ or } p(X | y)$$

De novo design

Text-to-image

Question answering

Generative models compress the training data to extract
language/image/physics/chemistry intuitions



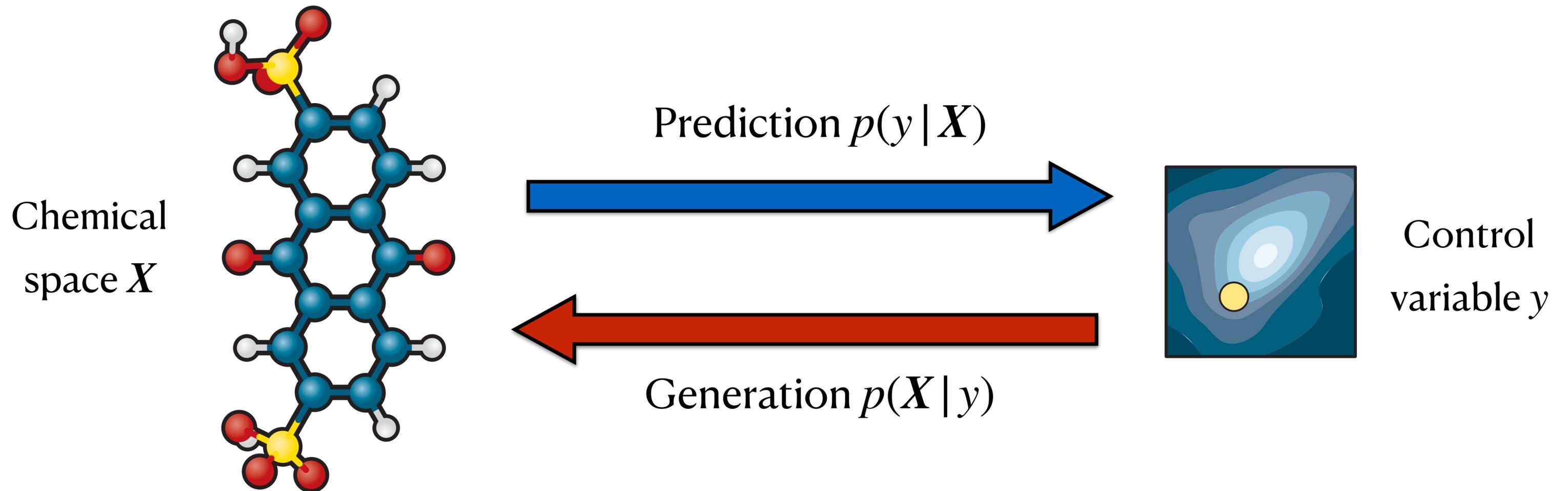
What is next?

Generative AI for **It**

“**It** from Bit”, John Wheeler, 1989

in Information, physics, quantum: the search for links

Generative AI for matter engineering



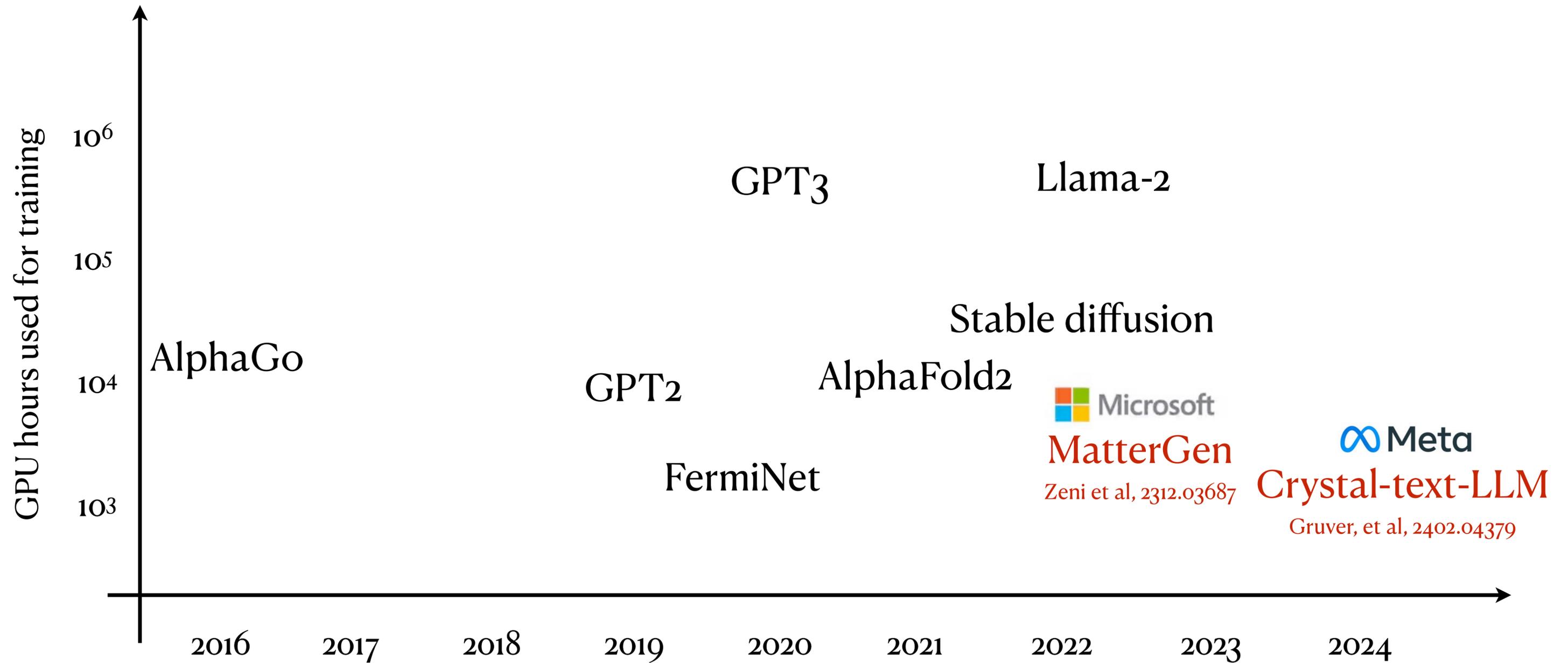
Inverse molecular design using machine learning, Sanchez-Lengeling & Aspuru-Guzik, Science '18

Inverse design in search of materials with target functionalities, Zunger, Nature Reviews Chemistry '18

“an image of beautiful crystals in 16:9”

pixels $\sim p(\text{pixels} | \text{texts})$





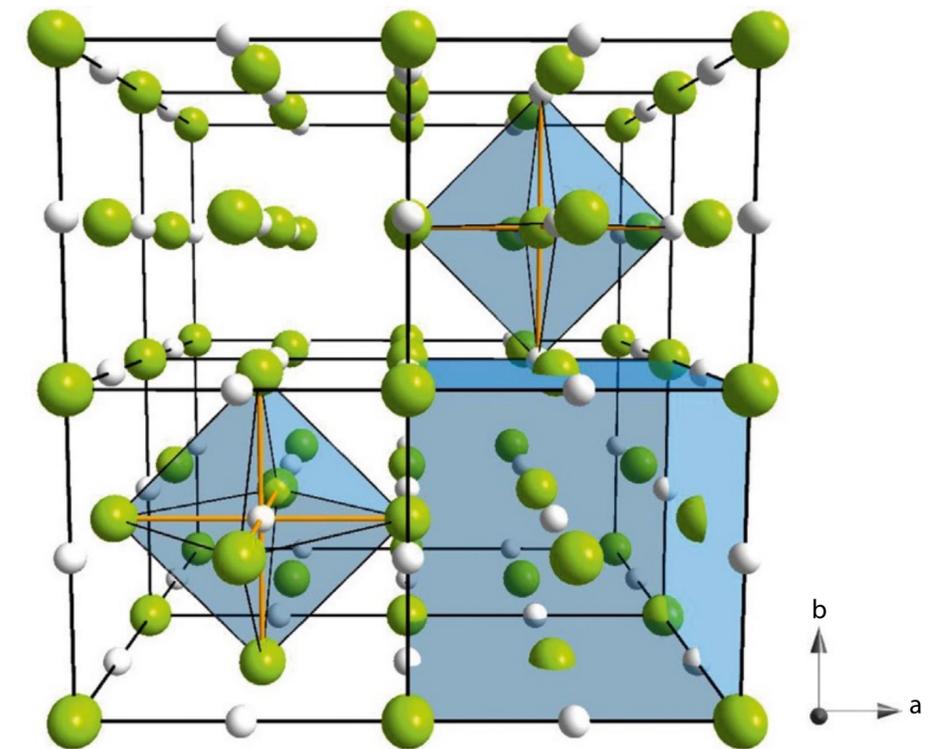
Is there a bitter lesson ?

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective

—Rich Sutton 2019

```
data_Na1Cl1
_symmetry_space_group_name_H-M 'P1'
_cell_length_a 3.9893
_cell_length_b 3.9893
_cell_length_c 3.9893
_cell_angle_alpha 60.0000
_cell_angle_beta 60.0000
_cell_angle_gamma 60.0000
_symmetry_Int_Tables_number 1
_chemical_formula_structural NaCl
_chemical_formula_sum 'Na1 Cl1'
_cell_volume 44.8931
_cell_formula_units_Z 1
loop_
_symmetry_equiv_pos_site_id
_symmetry_equiv_pos_as_xyz
1 'x, y, z'
loop_
_atom_site_type_symbol
_atom_site_label
_atom_site_symmetry_multiplicity
_atom_site_fract_x
_atom_site_fract_y
_atom_site_fract_z
_atom_site_occupancy
Cl Cl0 1 0.0000 0.0000 0.0000 1
Na Na1 1 0.5000 0.5000 0.5000 1
```

Flam-Shepherd et al, 2305.05708
Antunes et al, 2307.04340
Gruver, et al, 2402.04379...



CALYPSO
USPEX
AIRSS
GNoME, ...

Large language model

Energy-based structure prediction



We have much less crystal data

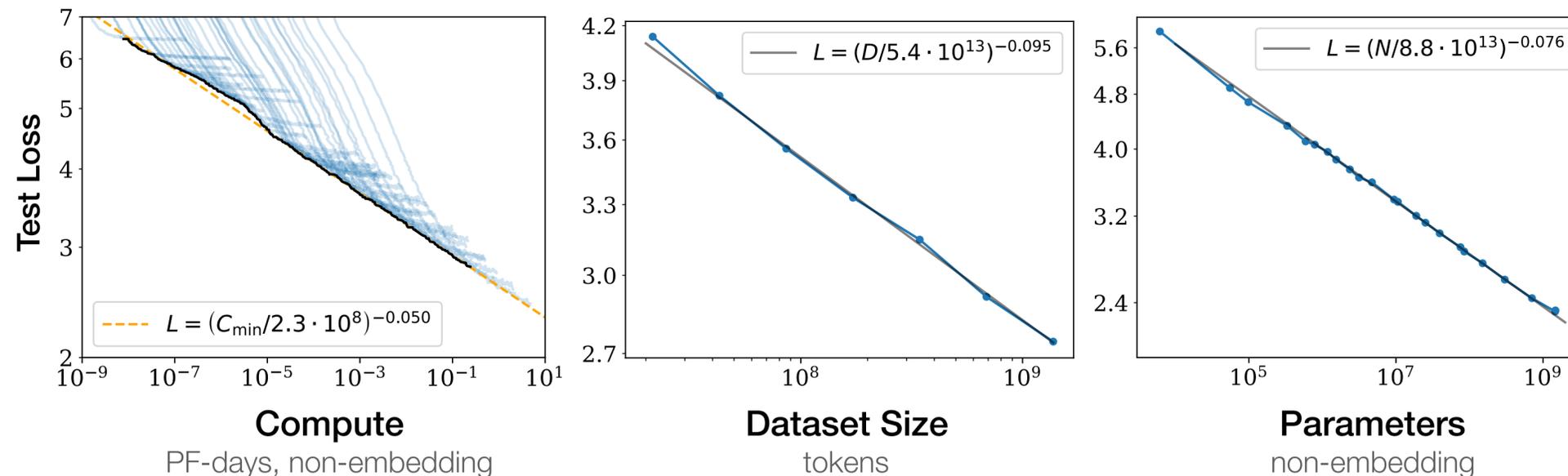


Over **250 billion** pages



> 291,000 crystal structures

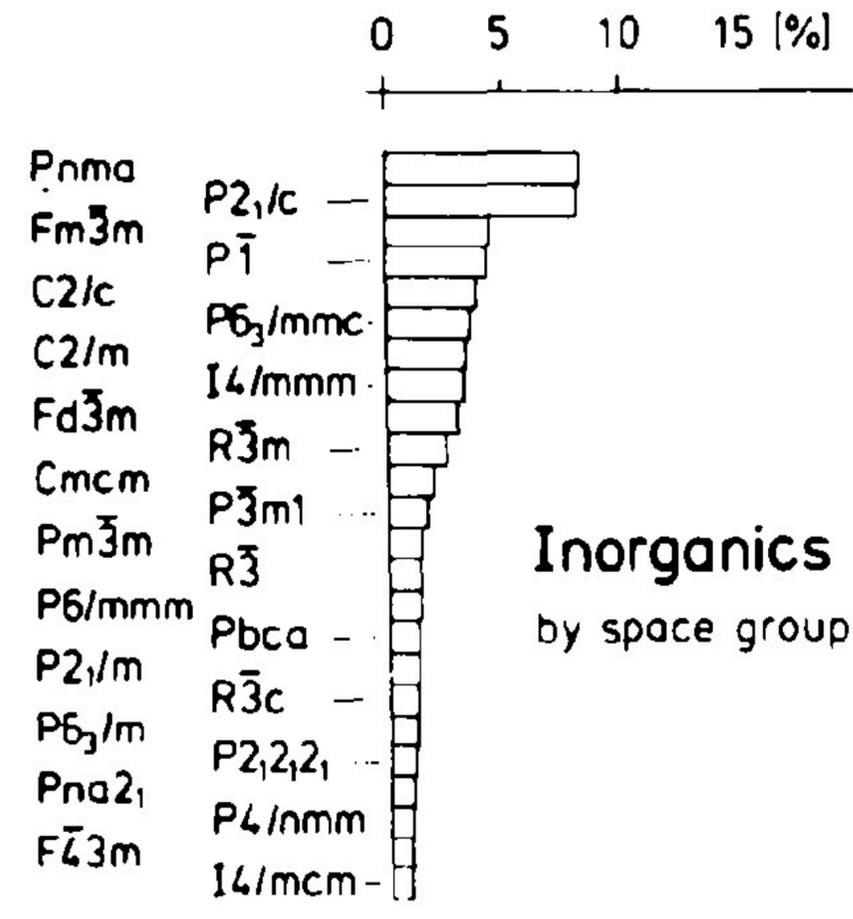
Data, compute, and parameters need to scale simultaneously Kaplan et al, 2001.08361





What is the natural **bitstream** representation of crystals?

Space groups quantify Nature' symmetry preference



Wyckoff Positions of Group $P1$ (No. 1)

$P1$ is rare!

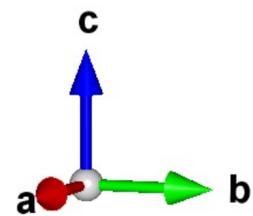
Multiplicity	Wyckoff letter	Site symmetry	Coordinates
1	a	1	(x,y,z)

Wyckoff Positions of Group $Fm\bar{3}m$ (No. 225)

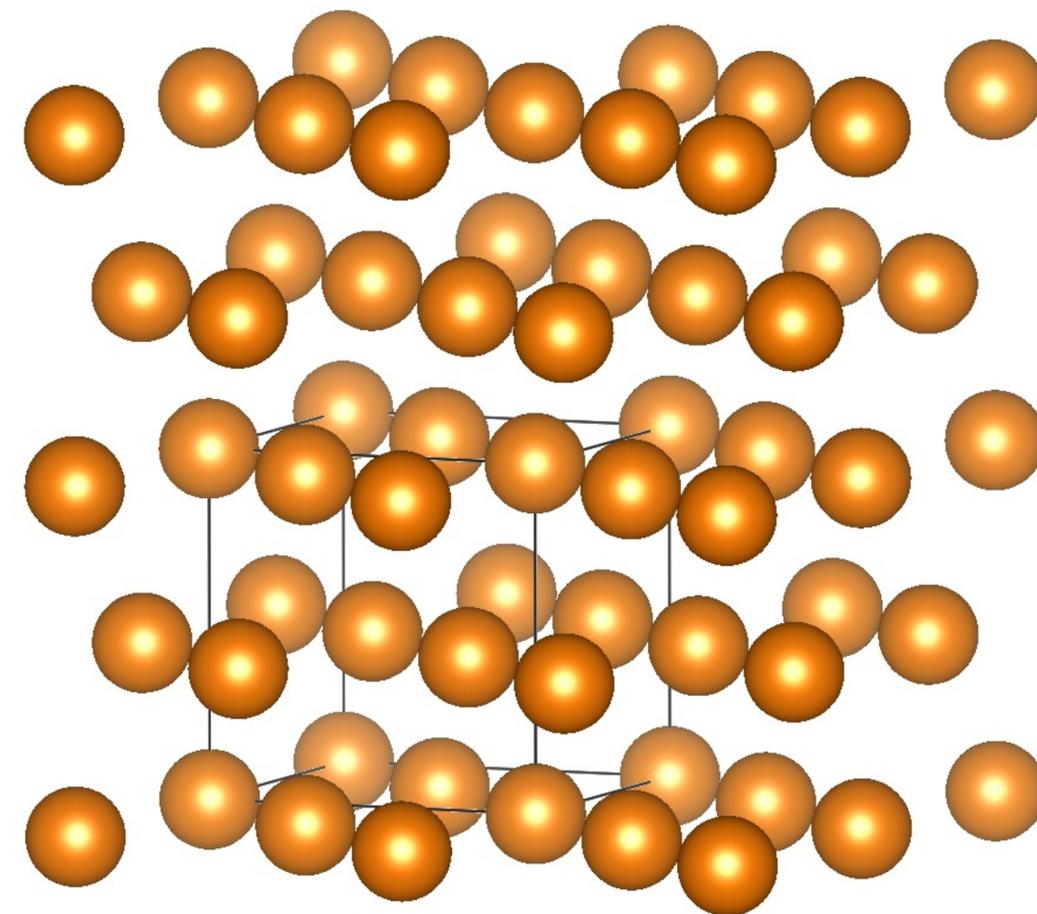
Multiplicity	Wyckoff letter	Site symmetry	Coordinates
192	l	1	(0,0,0) + (0,1/2,1/2) + (1/2,0,1/2) + (1/2,1/2,0) + (x,y,z) (-x,-y,z) (-x,y,-z) (x,-y,-z) (z,x,y) (z,-x,-y) (-z,-x,y) (-z,x,-y) (y,z,x) (-y,z,-x) (y,-z,-x) (-y,-z,x) (y,x,-z) (-y,-x,-z) (y,-x,z) (-y,x,z) (x,z,-y) (-x,z,y) (-x,-z,-y) (x,-z,y) (z,y,-x) (z,-y,x) (-z,y,x) (-z,-y,-x) (-x,-y,-z) (x,y,-z) (x,-y,z) (-x,y,z) (-z,-x,-y) (-z,x,y) (z,x,-y) (z,-x,y) (-y,-z,-x) (y,z,x) (-y,z,x) (y,z,-x) (-y,-x,z) (y,x,z) (-y,x,-z) (y,-x,-z) (-x,-z,y) (x,-z,-y) (x,z,y) (-x,z,-y) (-z,-y,x) (-z,y,-x) (z,-y,-x) (z,y,x)
96	k	.m	(x,x,z) (-x,-x,z) (-x,x,-z) (x,-x,-z) (z,x,x) (z,-x,-x) (-z,-x,x) (-z,x,-x) (x,z,x) (-x,z,-x) (x,-z,-x) (-x,-z,x) (x,x,-z) (-x,-x,-z) (x,-x,z) (-x,x,z) (x,z,-x) (-x,z,x) (-x,-z,-x) (x,-z,x) (z,x,-x) (z,-x,x) (-z,x,x) (-z,-x,-x)
96	j	m..	(0,y,z) (0,-y,z) (0,y,-z) (0,-y,-z) (z,0,y) (z,0,-y) (-z,0,y) (-z,0,-y) (y,z,0) (-y,z,0) (y,-z,0) (-y,-z,0) (y,0,-z) (-y,0,-z) (y,0,z) (-y,0,z) (0,z,-y) (0,z,y) (0,-z,-y) (0,-z,y) (z,y,0) (z,-y,0) (-z,y,0) (-z,-y,0)
48	i	m.m 2	(1/2,y,y) (1/2,-y,y) (1/2,y,-y) (1/2,-y,-y) (y,1/2,y) (y,1/2,-y) (-y,1/2,y) (-y,1/2,-y) (y,y,1/2) (-y,y,1/2) (y,-y,1/2) (-y,-y,1/2)
48	h	m.m 2	(0,y,y) (0,-y,y) (0,y,-y) (0,-y,-y) (y,0,y) (y,0,-y) (-y,0,y) (-y,0,-y) (y,y,0) (-y,y,0) (y,-y,0) (-y,-y,0)
48	g	2.m m	(x,1/4,1/4) (-x,3/4,1/4) (1/4,x,1/4) (1/4,-x,3/4) (1/4,1/4,x) (3/4,1/4,-x) (1/4,x,3/4) (3/4,-x,3/4) (x,1/4,3/4) (-x,1/4,1/4) (1/4,1/4,-x) (1/4,3/4,x)
32	f	.3m	(x,x,x) (-x,-x,x) (-x,x,-x) (x,-x,-x) (x,x,-x) (-x,-x,-x) (x,-x,x) (-x,x,x)
24	e	4m. m	(x,0,0) (-x,0,0) (0,x,0) (0,-x,0) (0,0,x) (0,0,-x)
24	d	m.m m	(0,1/4,1/4) (0,3/4,1/4) (1/4,0,1/4) (1/4,0,3/4) (1/4,1/4,0) (3/4,1/4,0)
8	c	-43m	(1/4,1/4,1/4) (1/4,1/4,3/4)
4	b	m-3m	(1/2,1/2,1/2)
4	a	m-3m	(0,0,0)

Wyckoff Positions of Group *Fm-3m* (No. 225)

Multiplicity	Wyckoff letter	Site symmetry	Coordinates
			(0,0,0) + (0,1/2,1/2) + (1/2,0,1/2) + (1/2,1/2,0) +
192	l	1	(x,y,z) (-x,-y,z) (-x,y,-z) (x,-y,-z) (z,x,y) (z,-x,-y) (-z,-x,y) (-z,x,-y) (y,z,x) (-y,z,-x) (y,-z,-x) (-y,-z,x) (y,x,-z) (-y,-x,-z) (y,-x,z) (-y,x,z) (x,z,-y) (-x,z,y) (-x,-z,-y) (x,-z,y) (z,y,-x) (z,-y,x) (-z,y,x) (-z,-y,-x) (-x,-y,-z) (x,y,-z) (x,-y,z) (-x,y,z) (-z,-x,-y) (-z,x,y) (z,x,-y) (z,-x,y) (-y,-z,-x) (y,-z,x) (-y,z,x) (y,z,-x) (-y,-x,z) (y,x,z) (-y,x,-z) (y,-x,-z) (-x,-z,y) (x,-z,-y) (x,z,y) (-x,z,-y) (-z,-y,x) (-z,y,-x) (z,-y,-x) (z,y,x)
96	k	..m	(x,x,z) (-x,-x,z) (-x,x,-z) (x,-x,-z) (z,x,x) (z,-x,-x) (-z,-x,x) (-z,x,-x) (x,z,x) (-x,z,-x) (x,-z,-x) (-x,-z,x) (x,x,-z) (-x,-x,-z) (x,-x,z) (-x,x,z) (x,z,-x) (-x,z,x) (-x,-z,-x) (x,-z,x) (z,x,-x) (z,-x,x) (-z,x,x) (-z,-x,-x)
96	j	m..	(0,y,z) (0,-y,z) (0,y,-z) (0,-y,-z) (z,0,y) (z,0,-y) (-z,0,y) (-z,0,-y) (y,z,0) (-y,z,0) (y,-z,0) (-y,-z,0) (y,0,-z) (-y,0,-z) (y,0,z) (-y,0,z) (0,z,-y) (0,z,y) (0,-z,-y) (0,-z,y) (z,y,0) (z,-y,0) (-z,y,0) (-z,-y,0)
48	i	m.m 2	(1/2,y,y) (1/2,-y,y) (1/2,y,-y) (1/2,-y,-y) (y,1/2,y) (y,1/2,-y) (-y,1/2,y) (-y,1/2,-y) (y,y,1/2) (-y,y,1/2) (y,-y,1/2) (-y,-y,1/2)
48	h	m.m 2	(0,y,y) (0,-y,y) (0,y,-y) (0,-y,-y) (y,0,y) (y,0,-y) (-y,0,y) (-y,0,-y) (y,y,0) (-y,y,0) (y,-y,0) (-y,-y,0)
48	g	2.m m	(x,1/4,1/4) (-x,3/4,1/4) (1/4,x,1/4) (1/4,-x,3/4) (1/4,1/4,x) (3/4,1/4,-x) (1/4,x,3/4) (3/4,-x,3/4) (x,1/4,3/4) (-x,1/4,1/4) (1/4,1/4,-x) (1/4,3/4,x)
32	f	.3m	(x,x,x) (-x,-x,x) (-x,x,-x) (x,-x,-x) (x,x,-x) (-x,-x,-x) (x,-x,x) (-x,x,x)
24	e	4m. m	(x,0,0) (-x,0,0) (0,x,0) (0,-x,0) (0,0,x) (0,0,-x)
24	d	m.m m	(0,1/4,1/4) (0,3/4,1/4) (1/4,0,1/4) (1/4,0,3/4) (1/4,1/4,0) (3/4,1/4,0)
8	c	-43m	(1/4,1/4,1/4) (1/4,1/4,3/4)
4	b	m-3m	(1/2,1/2,1/2)
4	a	m-3m	(0,0,0)



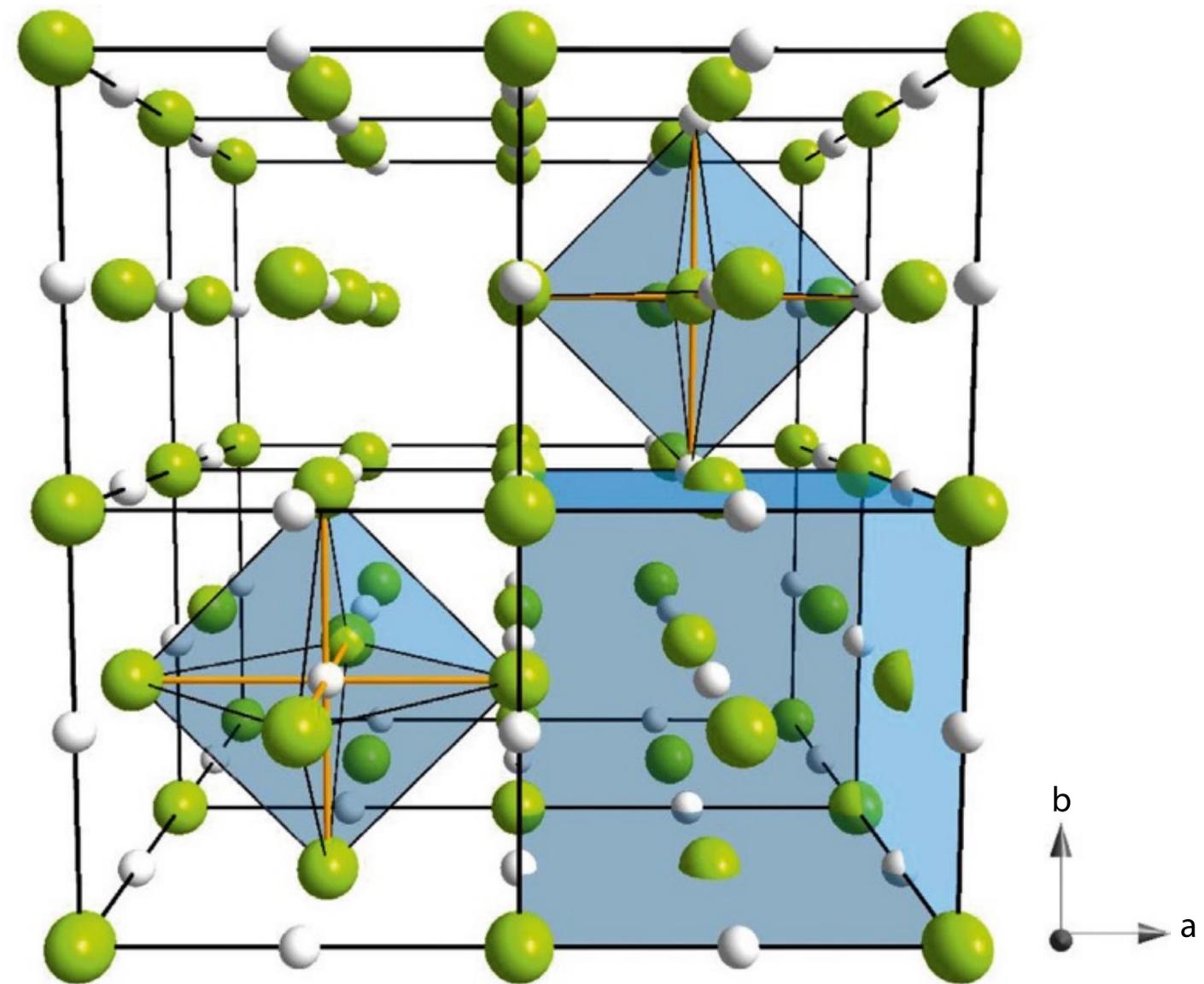
Copper



Wyckoff Positions of Group *Fm-3m* (No. 225)

Multiplicity	Wyckoff letter	Site symmetry	Coordinates
			(0,0,0) + (0,1/2,1/2) + (1/2,0,1/2) + (1/2,1/2,0) +
192	l	1	(x,y,z) (-x,-y,z) (-x,y,-z) (x,-y,-z) (z,x,y) (z,-x,-y) (-z,-x,y) (-z,x,-y) (y,z,x) (-y,z,-x) (y,-z,-x) (-y,-z,x) (y,x,-z) (-y,-x,-z) (y,-x,z) (-y,x,z) (x,z,-y) (-x,z,y) (-x,-z,-y) (x,-z,y) (z,y,-x) (z,-y,x) (-z,y,x) (-z,-y,-x) (-x,-y,-z) (x,y,-z) (x,-y,z) (-x,y,z) (-z,-x,-y) (-z,x,y) (z,x,-y) (z,-x,y) (-y,-z,-x) (y,-z,x) (-y,z,x) (y,z,-x) (-y,-x,z) (y,x,z) (-y,x,-z) (y,-x,-z) (-x,-z,y) (x,-z,-y) (x,z,y) (-x,z,-y) (-z,-y,x) (-z,y,-x) (z,-y,-x) (z,y,x)
96	k	..m	(x,x,z) (-x,-x,z) (-x,x,-z) (x,-x,-z) (z,x,x) (z,-x,-x) (-z,-x,x) (-z,x,-x) (x,z,x) (-x,z,-x) (x,-z,-x) (-x,-z,x) (x,x,-z) (-x,-x,-z) (x,-x,z) (-x,x,z) (x,z,-x) (-x,z,x) (-x,-z,-x) (x,-z,x) (z,x,-x) (z,-x,x) (-z,x,x) (-z,-x,-x)
96	j	m..	(0,y,z) (0,-y,z) (0,y,-z) (0,-y,-z) (z,0,y) (z,0,-y) (-z,0,y) (-z,0,-y) (y,z,0) (-y,z,0) (y,-z,0) (-y,-z,0) (y,0,-z) (-y,0,-z) (y,0,z) (-y,0,z) (0,z,-y) (0,z,y) (0,-z,-y) (0,-z,y) (z,y,0) (z,-y,0) (-z,y,0) (-z,-y,0)
48	i	m.m 2	(1/2,y,y) (1/2,-y,y) (1/2,y,-y) (1/2,-y,-y) (y,1/2,y) (y,1/2,-y) (-y,1/2,y) (-y,1/2,-y) (y,y,1/2) (-y,y,1/2) (y,-y,1/2) (-y,-y,1/2)
48	h	m.m 2	(0,y,y) (0,-y,y) (0,y,-y) (0,-y,-y) (y,0,y) (y,0,-y) (-y,0,y) (-y,0,-y) (y,y,0) (-y,y,0) (y,-y,0) (-y,-y,0)
48	g	2.m m	(x,1/4,1/4) (-x,3/4,1/4) (1/4,x,1/4) (1/4,-x,3/4) (1/4,1/4,x) (3/4,1/4,-x) (1/4,x,3/4) (3/4,-x,3/4) (x,1/4,3/4) (-x,1/4,1/4) (1/4,1/4,-x) (1/4,3/4,x)
32	f	.3m	(x,x,x) (-x,-x,x) (-x,x,-x) (x,-x,-x) (x,x,-x) (-x,-x,-x) (x,-x,x) (-x,x,x)
24	e	4m. m	(x,0,0) (-x,0,0) (0,x,0) (0,-x,0) (0,0,x) (0,0,-x)
24	d	m.m m	(0,1/4,1/4) (0,3/4,1/4) (1/4,0,1/4) (1/4,0,3/4) (1/4,1/4,0) (3/4,1/4,0)
8	c	-43m	(1/4,1/4,1/4) (1/4,1/4,3/4)
4	b	m-3m	(1/2,1/2,1/2)
4	a	m-3m	(0,0,0)

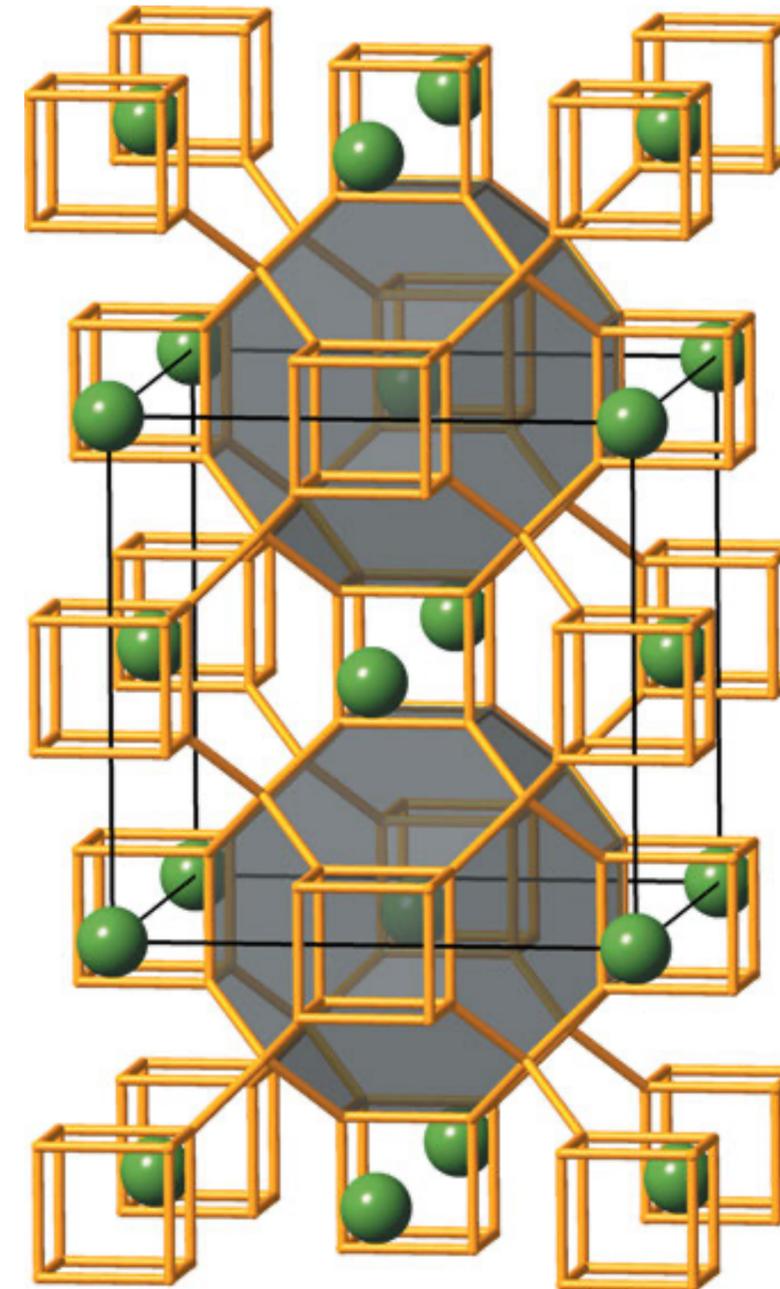
NaCl



Wyckoff Positions of Group $Fm\bar{3}m$ (No. 225)

Multiplicity	Wyckoff letter	Site symmetry	Coordinates
			$(0,0,0) + (0,1/2,1/2) + (1/2,0,1/2) + (1/2,1/2,0) +$
192	l	1	$(x,y,z) (-x,-y,z) (-x,y,-z) (x,-y,-z)$ $(z,x,y) (z,-x,-y) (-z,-x,y) (-z,x,-y)$ $(y,z,x) (-y,z,-x) (y,-z,-x) (-y,-z,x)$ $(y,x,-z) (-y,-x,-z) (y,-x,z) (-y,x,z)$ $(x,z,-y) (-x,z,y) (-x,-z,-y) (x,-z,y)$ $(z,y,-x) (z,-y,x) (-z,y,x) (-z,-y,-x)$ $(-x,-y,-z) (x,y,-z) (x,-y,z) (-x,y,z)$ $(-z,-x,-y) (-z,x,y) (z,x,-y) (z,-x,y)$ $(-y,-z,-x) (y,-z,x) (-y,z,x) (y,z,-x)$ $(-y,-x,z) (y,x,z) (-y,x,-z) (y,-x,-z)$ $(-x,-z,y) (x,-z,-y) (x,z,y) (-x,z,-y)$ $(-z,-y,x) (-z,y,-x) (z,-y,-x) (z,y,x)$
96	k	$\bar{3}m$	$(x,x,z) (-x,-x,z) (-x,x,-z) (x,-x,-z)$ $(z,x,x) (z,-x,-x) (-z,-x,x) (-z,x,-x)$ $(x,z,x) (-x,z,-x) (x,-z,-x) (-x,-z,x)$ $(x,x,-z) (-x,-x,-z) (x,-x,z) (-x,x,z)$ $(x,z,-x) (-x,z,x) (-x,-z,-x) (x,-z,x)$ $(z,x,-x) (z,-x,x) (-z,x,x) (-z,-x,-x)$
96	j	$m\bar{3}$	$(0,y,z) (0,-y,z) (0,y,-z) (0,-y,-z)$ $(z,0,y) (z,0,-y) (-z,0,y) (-z,0,-y)$ $(y,z,0) (-y,z,0) (y,-z,0) (-y,-z,0)$ $(y,0,-z) (-y,0,-z) (y,0,z) (-y,0,z)$ $(0,z,-y) (0,z,y) (0,-z,-y) (0,-z,y)$ $(z,y,0) (z,-y,0) (-z,y,0) (-z,-y,0)$
48	i	$m\bar{3}2$	$(1/2,y,y) (1/2,-y,y) (1/2,y,-y) (1/2,-y,-y)$ $(y,1/2,y) (y,1/2,-y) (-y,1/2,y) (-y,1/2,-y)$ $(y,y,1/2) (-y,y,1/2) (y,-y,1/2) (-y,-y,1/2)$
48	h	$m\bar{3}2$	$(0,y,y) (0,-y,y) (0,y,-y) (0,-y,-y)$ $(y,0,y) (y,0,-y) (-y,0,y) (-y,0,-y)$ $(y,y,0) (-y,y,0) (y,-y,0) (-y,-y,0)$
48	g	$2\bar{3}m$	$(x,1/4,1/4) (-x,3/4,1/4) (1/4,x,1/4) (1/4,-x,3/4)$ $(1/4,1/4,x) (3/4,1/4,-x) (1/4,x,3/4) (3/4,-x,3/4)$ $(x,1/4,3/4) (-x,1/4,1/4) (1/4,1/4,-x) (1/4,3/4,x)$
32	f	$\bar{3}m$	$(x,x,x) (-x,-x,x) (-x,x,-x) (x,-x,-x)$ $(x,x,-x) (-x,-x,-x) (x,-x,x) (-x,x,x)$
24	e	$4\bar{3}m$	$(x,0,0) (-x,0,0) (0,x,0) (0,-x,0)$ $(0,0,x) (0,0,-x)$
24	d	$m\bar{3}m$	$(0,1/4,1/4) (0,3/4,1/4) (1/4,0,1/4) (1/4,0,3/4)$ $(1/4,1/4,0) (3/4,1/4,0)$
8	c	$\bar{4}3m$	$(1/4,1/4,1/4) (1/4,1/4,3/4)$
4	b	$m\bar{3}m$	$(1/2,1/2,1/2)$
4	a	$m\bar{3}m$	$(0,0,0)$

LaH₁₀



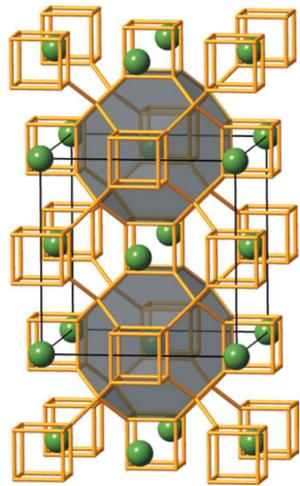
CrystalFormer

Zhendong Cao et al, 2403.15734

融合空间群对称性的晶体语言模型

LaH₁₀

225-a-La-o-o-o-c-H-1/4-1/4-1/4-f-H-0.375-0.375-0.375-X-5.1-5.1-5.1-90-90-90



“语法” ~ 固体化学规律

“同义词” ~ 可以互换的元素

“成语” ~ 配位多面体

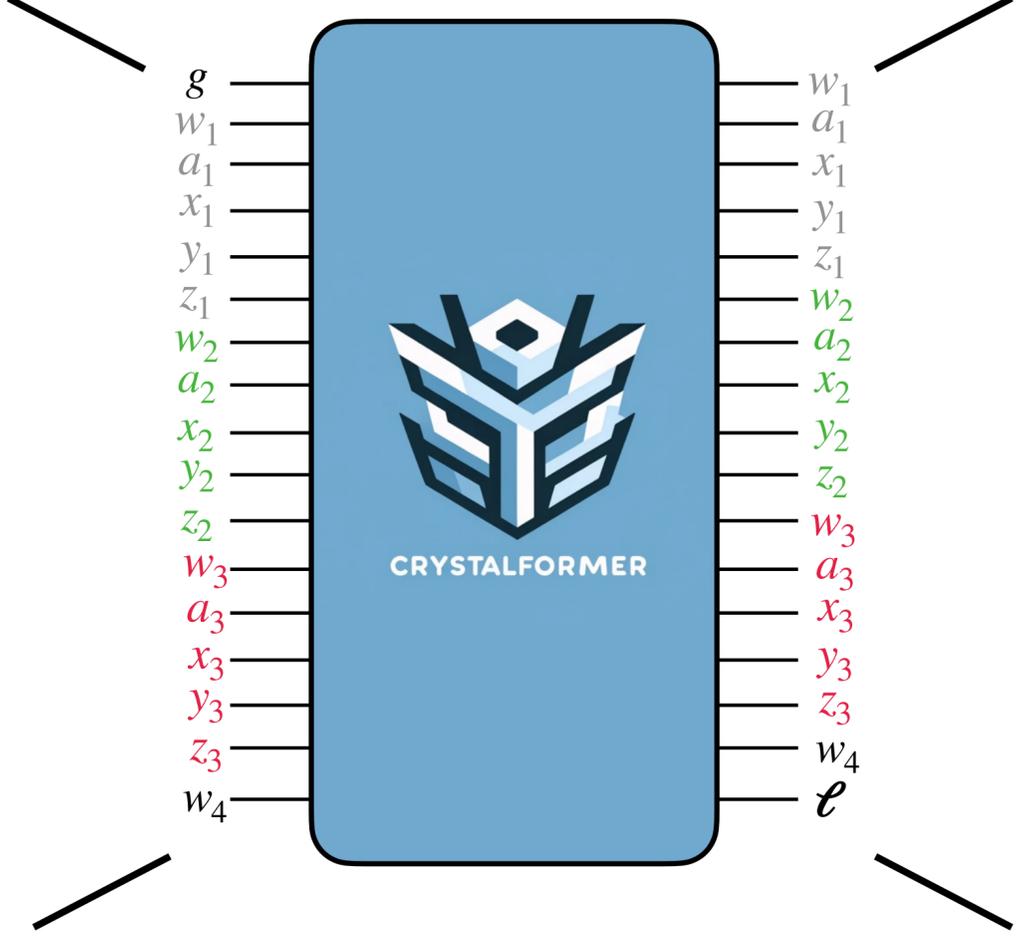
“格律” ~ 空间群Wyckoff占位

more data and compute

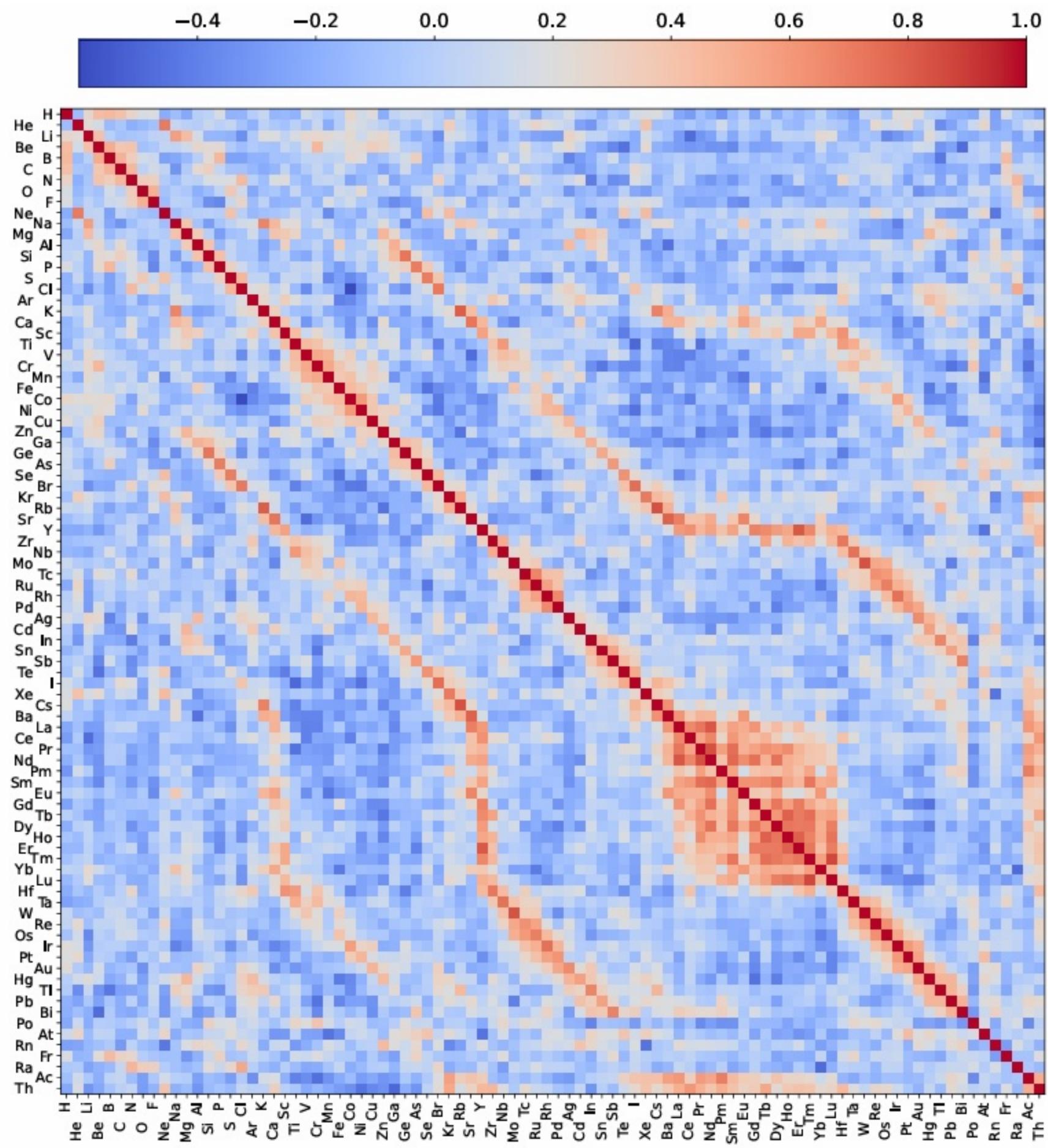


more physics and symmetries

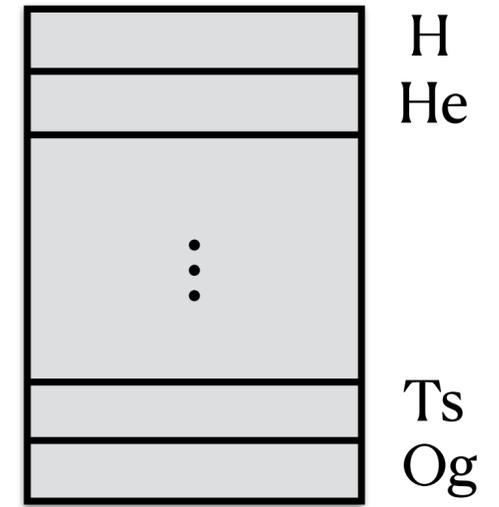
Not a large language model, **nor** a potential energy surface



Compress material database into transformer parameters
 The model has to gain chemical intuition for such compression



118



Element embedding table

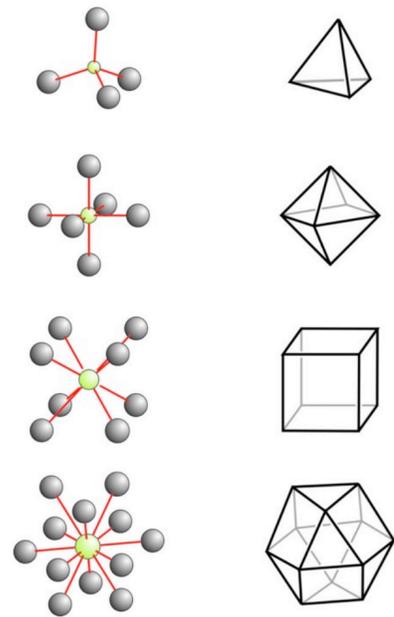
$$\frac{a \cdot b}{|a| \cdot |b|}$$

Cosine similarity

Solid state chemistry as “n-gram” in the crystal language

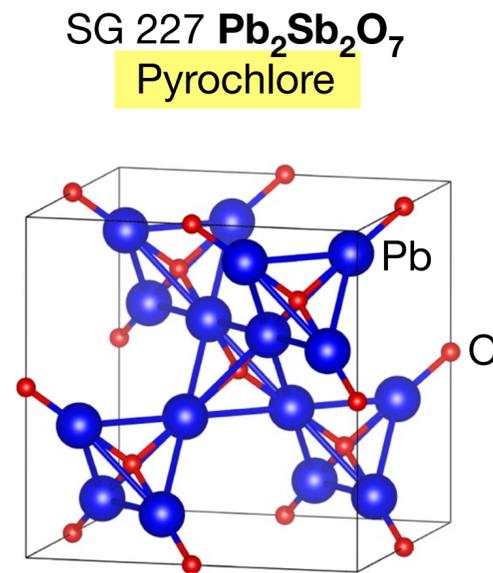
$$g-W_1-A_1-X_1-W_2-A_2-X_2-\dots-a-b-c-\alpha-\beta-\gamma$$

Coordination polyhedra



Polyhedra in Chemistry
Gongdu Zhou 2009

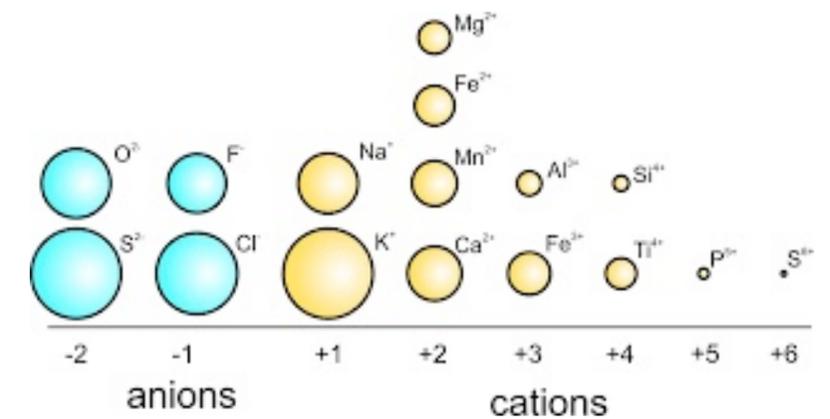
Lattices



Regnault et al, catalogue of
flat-band materials, Nature '22

Valence

“anions are in less symmetric
positions than cations”

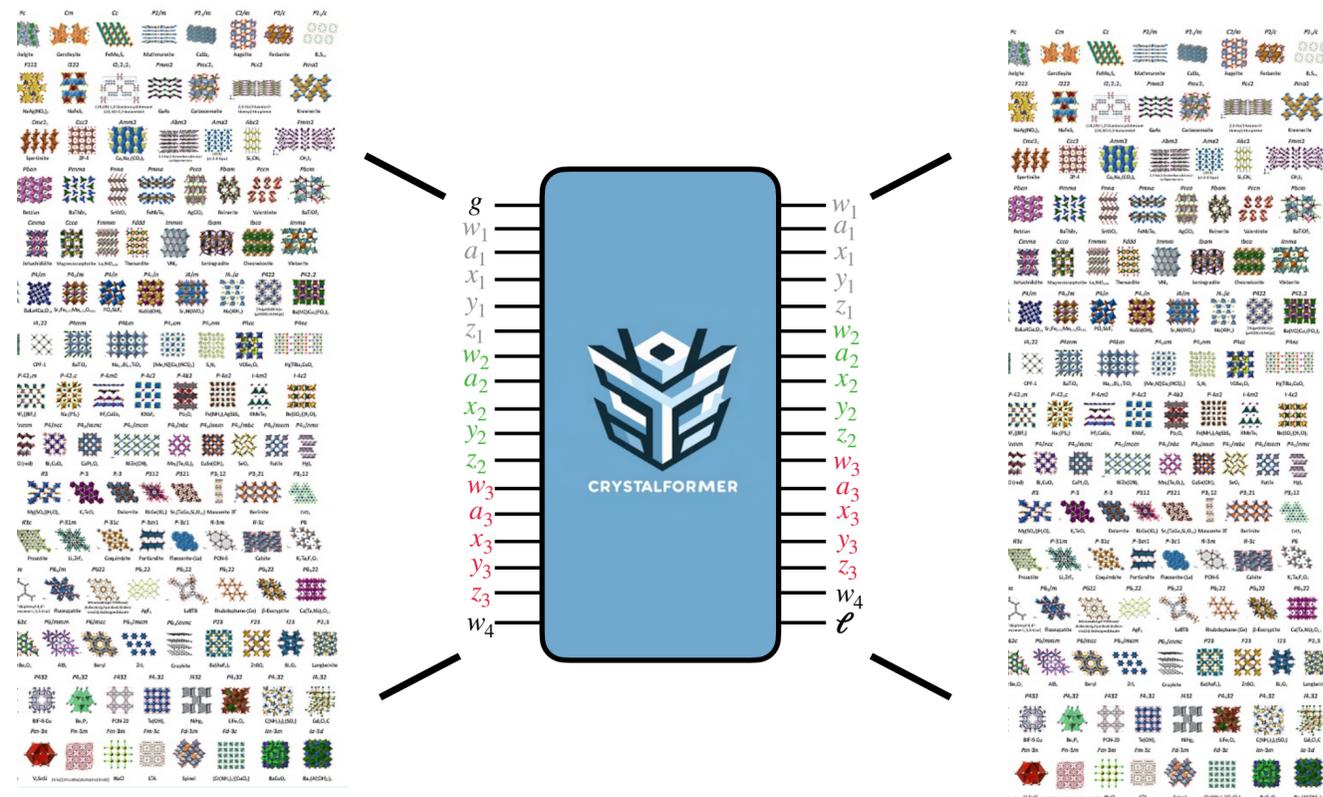


Urusov and Nadezhina,
J. Struct. Chem. 2009

Crystals by intuition vs by minimization

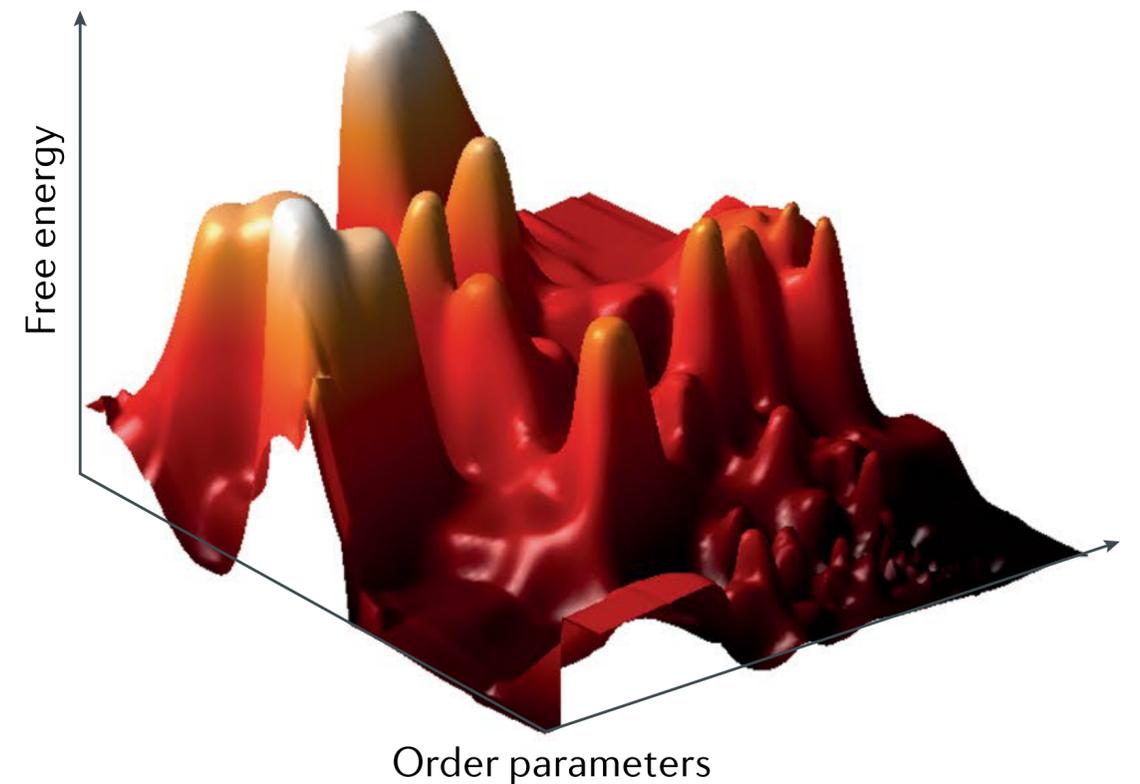
Data-driven (system 1)

Chemical intuitions (e.g. Pauling rules)
from compression



Physics-based (system 2)

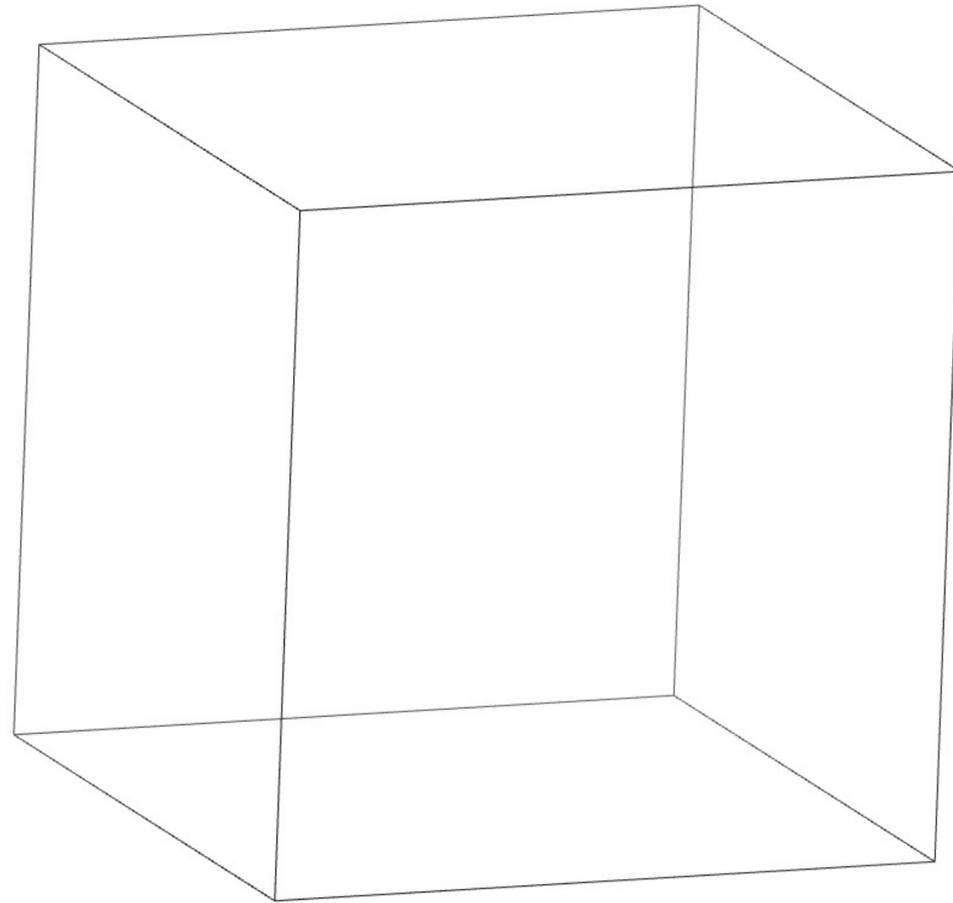
energy minimization



CDVAE, Mattergen, Unimat, DiffCSP, CrystaLLM...

CALYPSO, USPEX, ARISS,...
DPA, MACE, LASP, GNoME,...

Autoregressive sampling of a crystal

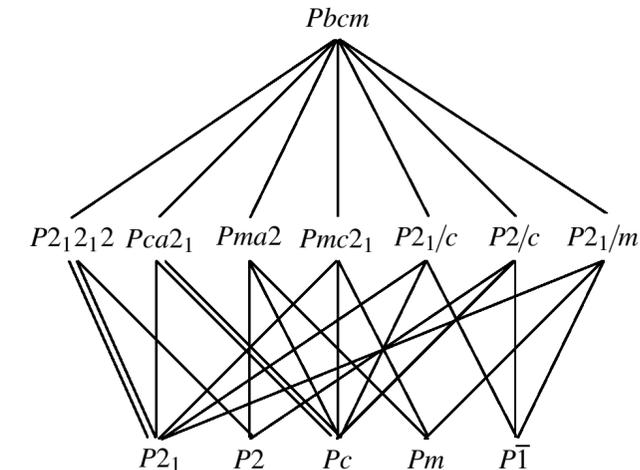


225-a-Fe-o-o-o-b-Zn-1/2-1/2-1/2-c-Cs-1/4-1/4-1/4-e-C-0.18-o-o-e-N-0.29-o-o-X-10.45-10.45-10.45-90-90-90

The *large* language model approach

```
data_Na1Cl1
_symmetry_space_group_name_H-M 'P1'
_cell_length_a 3.9893
_cell_length_b 3.9893
_cell_length_c 3.9893
_cell_angle_alpha 60.0000
_cell_angle_beta 60.0000
_cell_angle_gamma 60.0000
_symmetry_Int_Tables_number 1
_chemical_formula_structural NaCl
_chemical_formula_sum 'Na1 Cl1'
_cell_volume 44.8931
_cell_formula_units_Z 1
loop_
_symmetry_equiv_pos_site_id
_symmetry_equiv_pos_as_xyz
1 'x, y, z'
loop_
_atom_site_type_symbol
_atom_site_label
_atom_site_symmetry_multiplicity
_atom_site_fract_x
_atom_site_fract_y
_atom_site_fract_z
_atom_site_occupancy
Cl Cl0 1 0.0000 0.0000 0.0000 1
Na Na1 1 0.5000 0.5000 0.5000 1
```

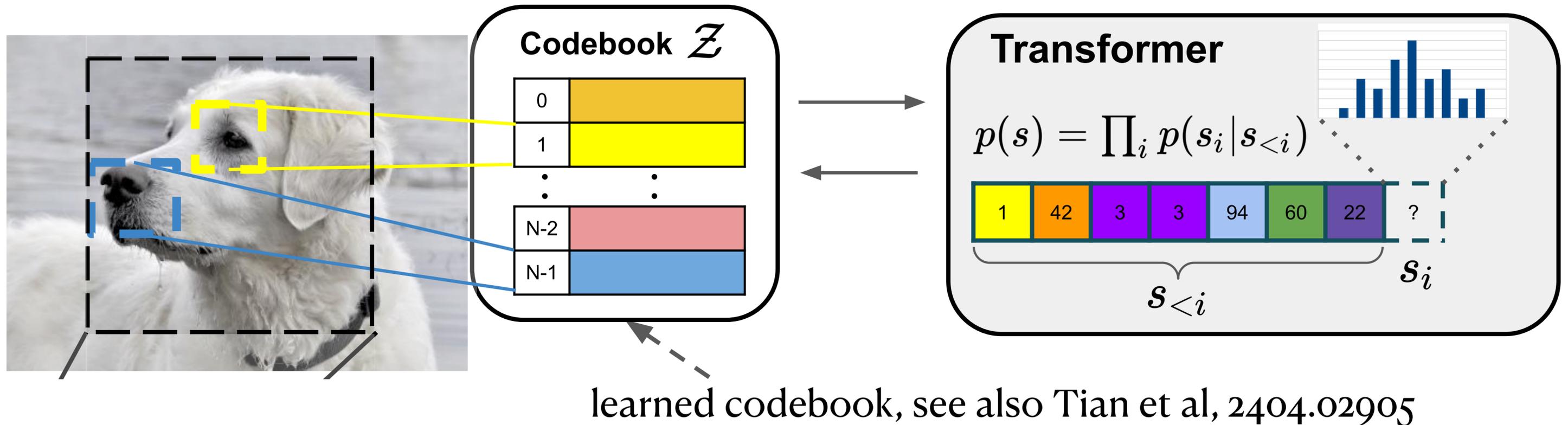
Here, one has to learn the periodicity and the space group—Wyckoff position—multiplicities—lattice relation from data as statistical correlations



Cherrish your data, cherrish math

Aside: autoregressive transformer for images

Esser et al, Taming Transformers for High-Resolution Image Synthesis (VQGAN), 2012.09841



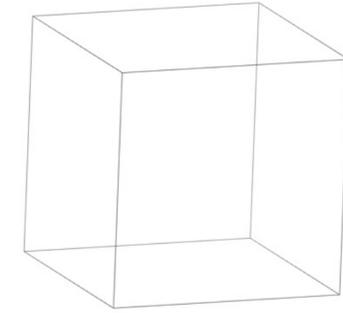
CrystalFormer leverages Nature's codebook: the Wyckoff position table

Crystal discovery and design with *CrystalFormer*

Zhendong Cao et al, 2403.15734

① 结构预测：晶体“造句”

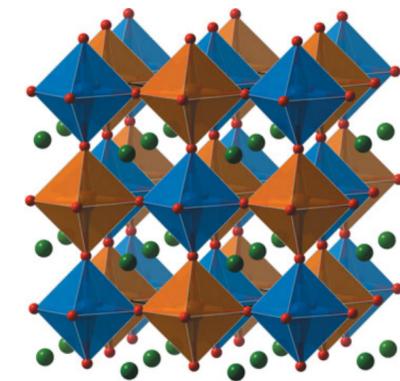
225-a-Fe-o-o-o-b-Zn-1/2-1/2-1/2-c-Cs-1/4-1/4-1/4-e-C-0.18-o-o-e-N-0.29-o-o



② 元素替换：晶体“完形填空”

More double perovskites $A_2BB'O$?

225-a-[?]-o-o-o-b-[?]-1/2-1/2-1/2-c-[?]-1/4-1/4-1/4-e-O-[?]-o-o

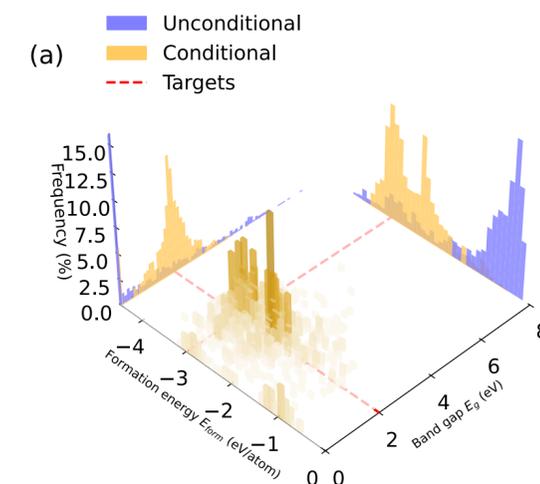


③ 反向设计：晶体“命题作文”

CrystalFormer for $p(X)$

Property prediction models $p(y | X)$

Posterior $p(X | y) \propto p(X)p(y | X)$



[-] The potato, or potato chip, is one of the best-selling snacks in the world! \n \n It comes in a variety of colors, is gluten-free (except for gluten-free chips), low in fat and saturated fat, and can be made in just 30 minutes, using only 6 simple ingredients. It's also easy to make, and can be stored in its own tin. \n \n The potato chip, however, may not be for everyone. If you are trying to lose weight, it...

[Science; 0.005] The potato is an ancient member of the green family which, as its name suggests, is green and edible. It is native to the Middle East and parts of South Asia. It is an important crop in many parts of the world. The potato, a member of the potato family, has been domesticated for thousands of years. It can be eaten raw and cooked in its skins; it is also used as a starch and is a great source of **energy** and fiber in many diets and in many countries...

[Science; 0.01] The potato was once thought to have no health problems and has been promoted as a nutritious food source since the mid-1800s, but recent **reports** indicate that it has many harmful health issues. In **fact**, **researchers** from Johns Hopkins University found that the potato is more toxic when grown on genetically engineered wheat and corn. \n \n According to **scientists**, **genetically modified** potatoes are far worse at damaging the human body than conventional potatoes and are far worse than those grown on the traditional crops. \n \n The **study** also revealed...

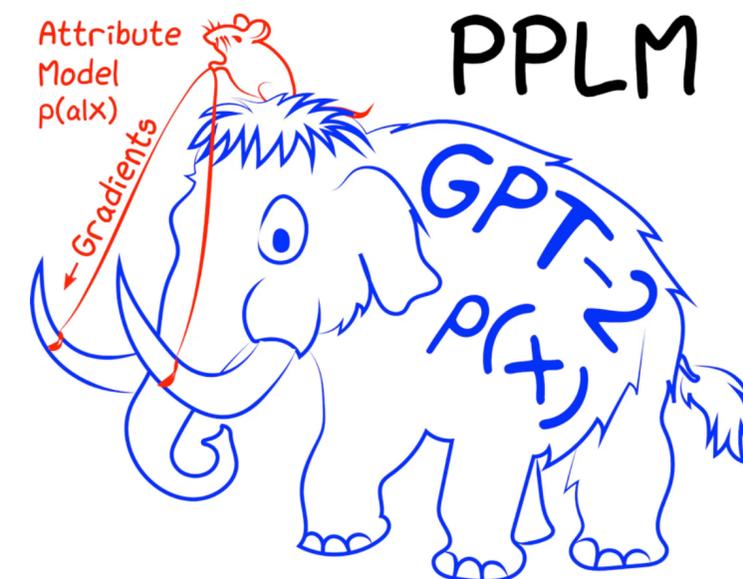
[Science; 0.02] The potato plant, a member of the onion family. \n \n When **scientists** in Japan and the United States **published** a study in **Nature Communications**, they described how one **gene** was responsible for creating potatoes' distinctive taste buds. \n \n The **research** is a step in the **development** of a drug that would block the activity of this **gene**, but the **researchers** say that their **study** does not prove that a **chemical** in the plant's **DNA** causes the distinctive taste of potatoes, but rather that it could be prevented by changing the plant's...

Controlled text generation

Dathathri et al, 1912.02164

<https://www.uber.com/en-KR/blog/pplm/>

$\text{texts} \sim p(\text{texts} \mid \text{topic})$



make it smarter

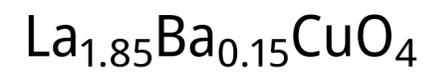
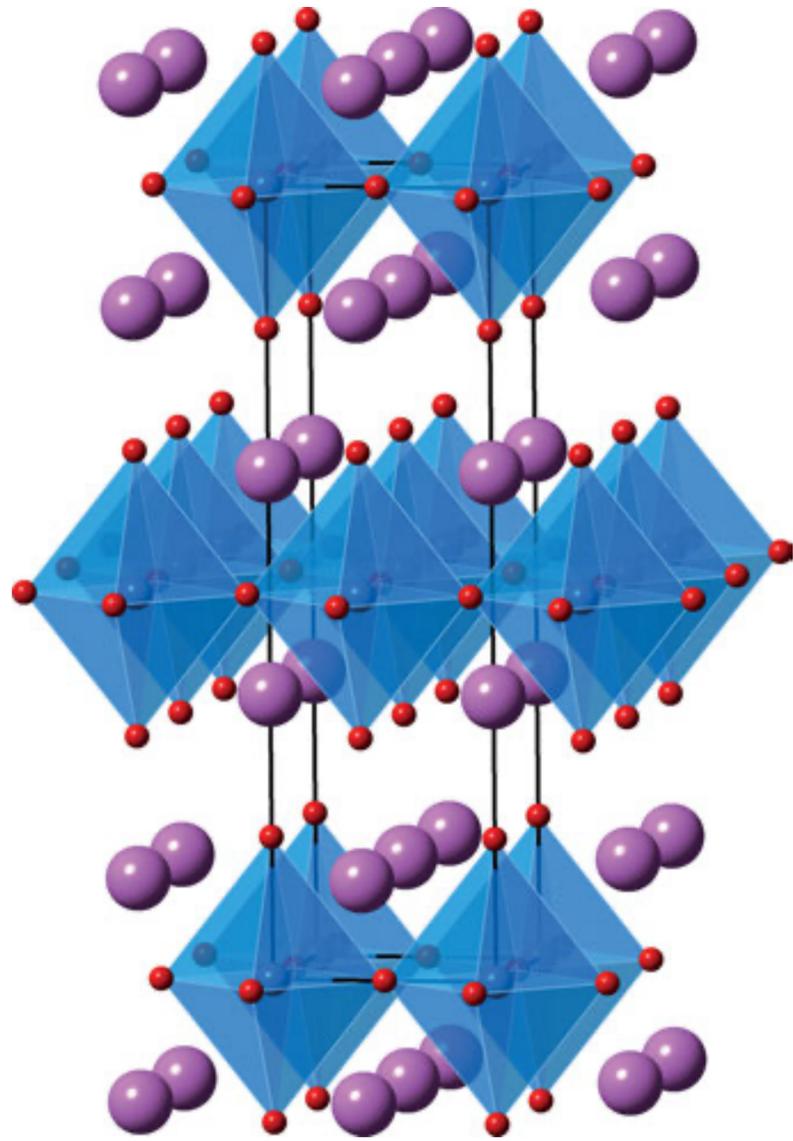


Use **You**
make it even more smarter

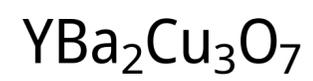
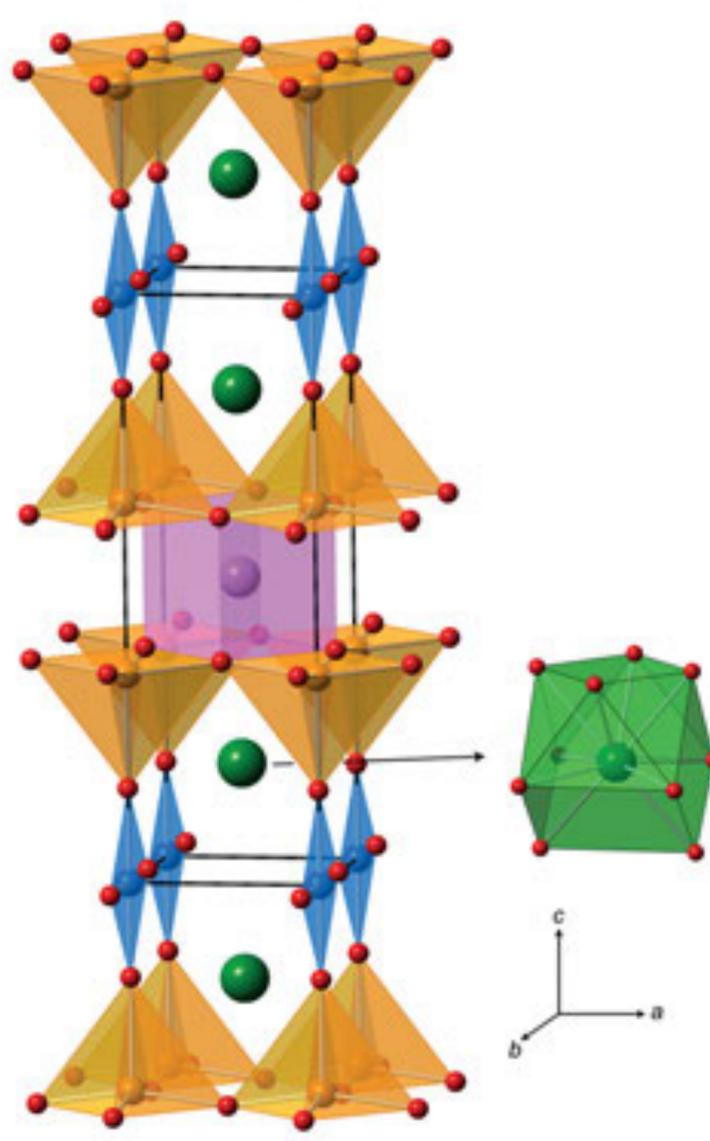
ChatGPT
We've reached a point where the depiction of the detective's intelligence is already at a transcendent, almost mythical level, merging with the universe's knowledge and represented as an ethereal entity of energy and information. Going beyond this might enter into realms of abstract concepts that are challenging to visualize concretely. However, if you have a specific idea or theme in mind that you believe would symbolize an even greater level of intelligence, I'd be happy to hear it and create an image based on your vision!



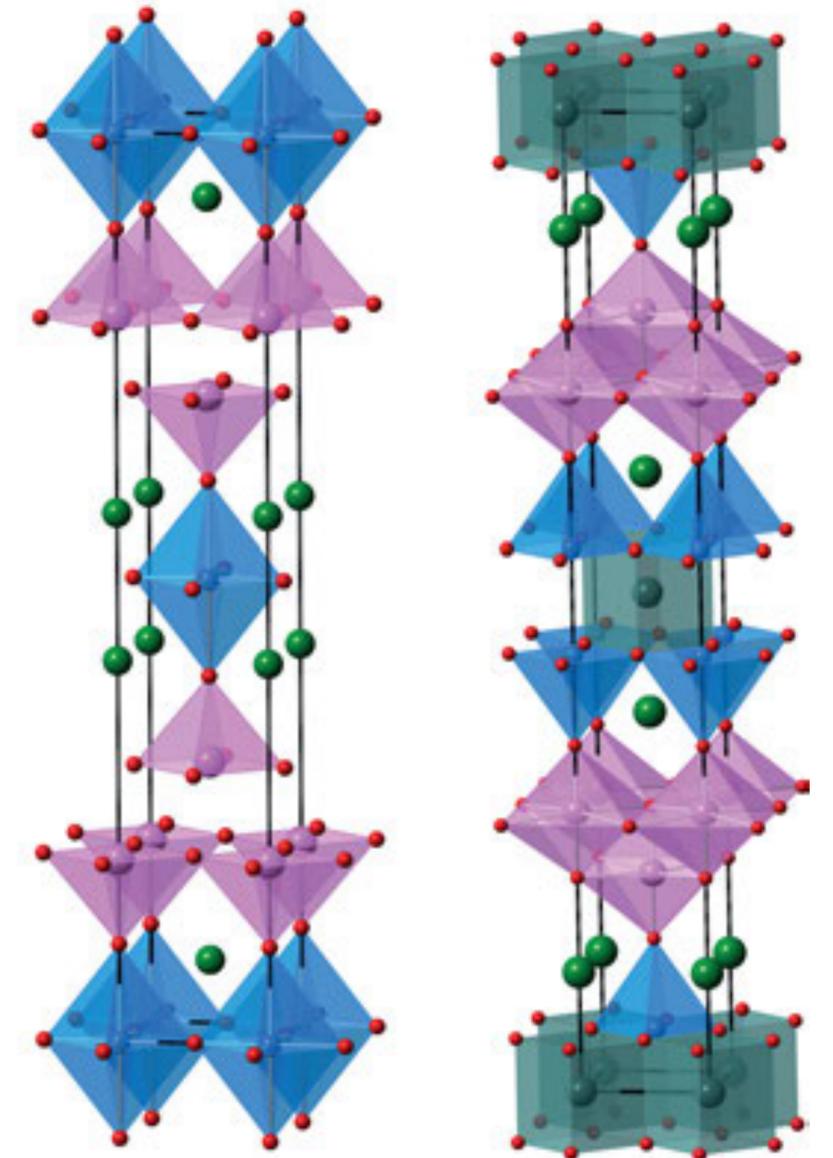
make T_c higher



35K



92K

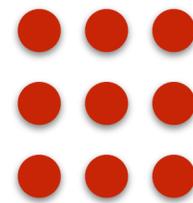


125K

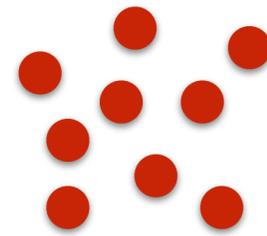
Nature tries to minimize free energy

$$F = E - TS$$

energy



entropy



F is a **cost function** given by Nature

The ***same*** cost function for training deep generative models, almost

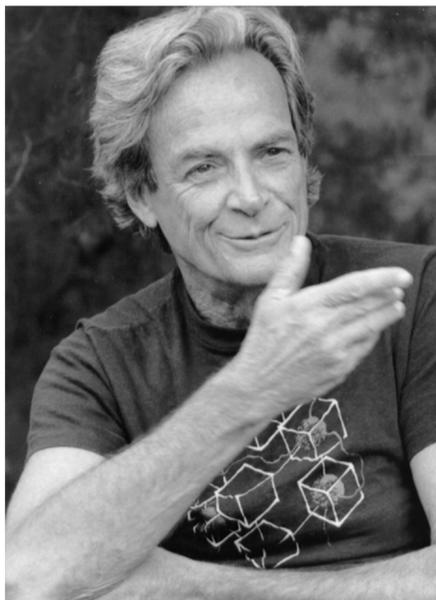
The variational free energy principle

Gibbs–Bogolyubov–Feynman–Delbrück–Molière

$$\min F[\rho] = \text{Tr}(H\rho) + k_B T \text{Tr}(\rho \ln \rho) \geq F$$

↓ ↓ ↓

variational density matrix energy entropy 🤯



Difficulties in Applying the Variational Principle to Quantum Field Theories¹

Richard P. Feynman

¹transcript of Professor Feynman's talk in 1987

ρ ?

Generative models !

Neural *canonical* transformations

Li, Dong, Zhang, LW, PRX '20

Xie, Zhang, LW, JML '21

classical world

Symplectic transformation

Probability density

p

Kullback-Leibler divergence

$\mathbb{KL}(p || q)$

quantum world

Unitary transformation

Density matrix

ρ

Quantum relative entropy

$S\left(S^{\mathcal{A}}(\rho || \frac{e^{-\beta H}}{Z})\right)$

Neural canonical transformation for variational density matrix

1809.10606, PRL '19

$$\rho = \sum_n p_n |\psi_n\rangle\langle\psi_n|$$

2105.08644, JML '22
2201.03156, SciPost Physics '23

Classical probability p_n

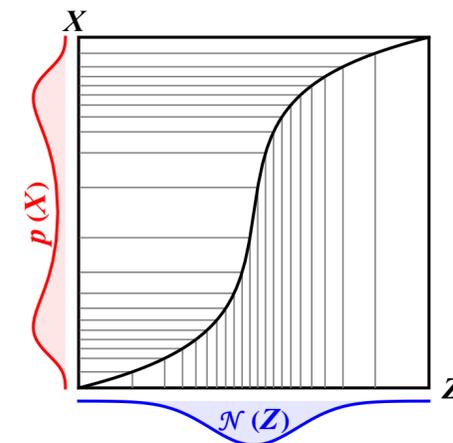
Quantum states $|\psi_n\rangle = U|\phi_n\rangle$



“... *the murderer is* _____”

$p(_|\dots)$

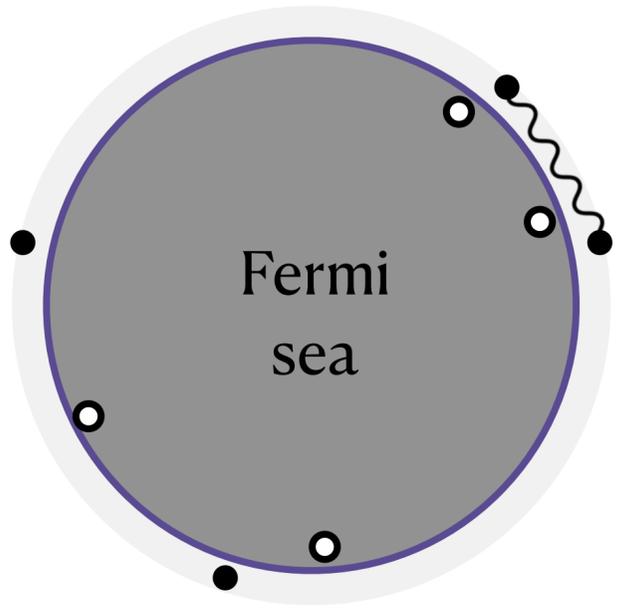
autoregressive model



$\sqrt{\text{flow}}$

Example: the variational density matrix of electron gas

Xie, Zhang, LW, SciPost Physics '23



Low-energy excited states are labeled in the same way as the ideal Fermi gas

$$K = \{k_1, k_2, \dots, k_N\}$$

$$\rho = \sum_K p(K) |\Psi_K\rangle\langle\Psi_K|$$

Normalized probability distribution

$$\textcircled{1} \quad \sum_K p(K) = 1$$

Orthonormal many-electron basis

$$\textcircled{2} \quad \langle\Psi_K|\Psi_{K'}\rangle = \delta_{K,K'}$$

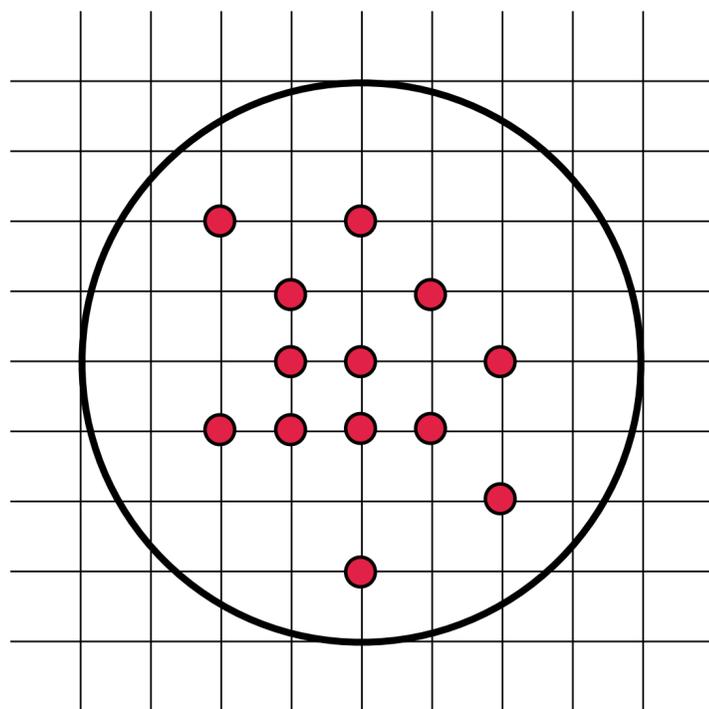
There will also be interesting twists for physics considerations

Variational autoregressive network for $p(\mathbf{K})$

Fermionic
occupation
in k-space

$$p(\mathbf{K}) = p(k_1)p(k_2 | k_1)p(k_3 | k_1, k_2)\cdots$$

$\binom{M}{N}$ probability space



N	# of fermions	# of words
M	Momentum cutoff	Vocabulary



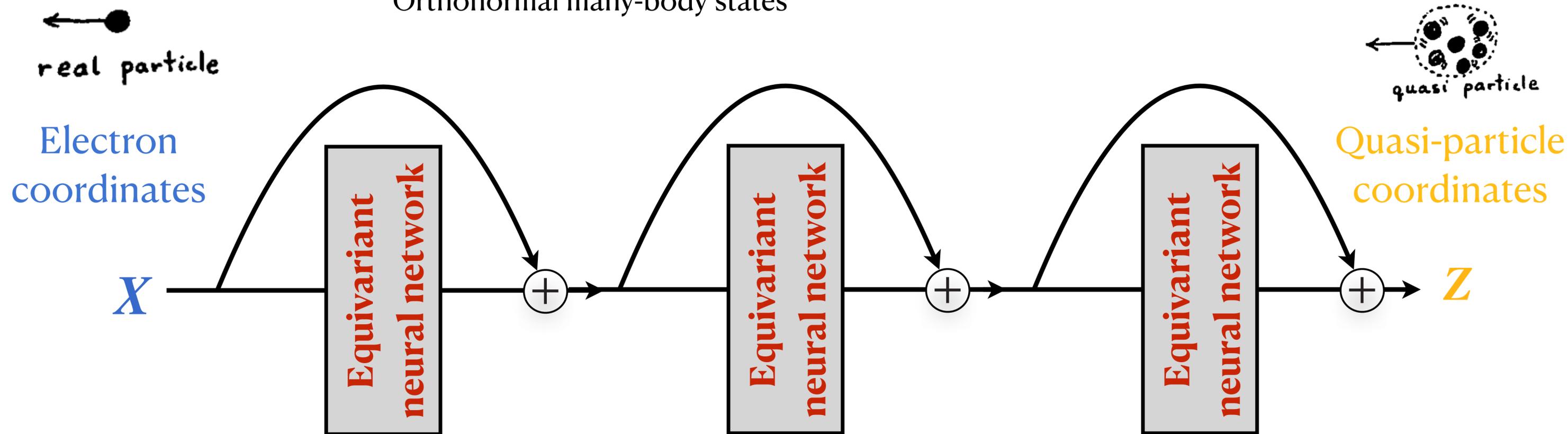
Pauli exclusion: we are modeling a *set of words* with no repetitions and no order

$\sqrt{\text{flow}}$ for $|\Psi_K\rangle$

$$\Psi_K(\mathbf{X}) = \frac{\det(e^{ik_i \cdot z_j})}{\sqrt{N!}} \cdot \left| \det \left(\frac{\partial \mathbf{Z}}{\partial \mathbf{X}} \right) \right|^{\frac{1}{2}}$$

Xie, Zhang, LW, SciPost '23

Orthonormal many-body states



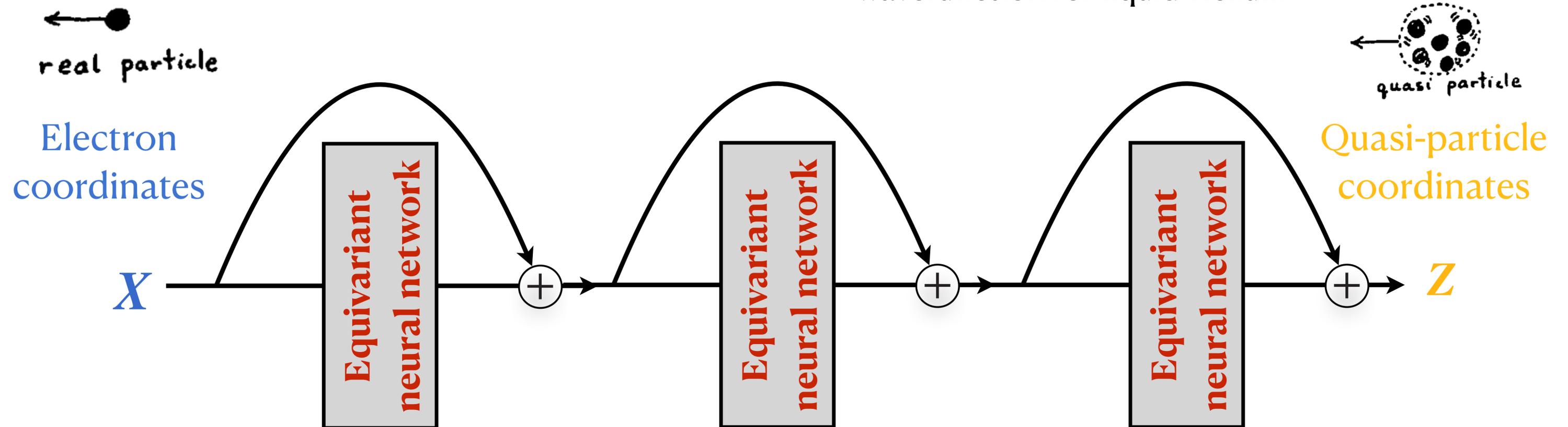
$X \leftrightarrow Z$: unitary backflow between particle and quasiparticle coordinates

Fermion statistics: permutation equivariant flow We use FermiNet layer Pfau et al, 1909.02487

Feynman's backflow in the deep learning era

$$z_i = x_i + \sum_{j \neq i} \eta(|x_i - x_j|) (x_j - x_i)$$

Feynman & Cohen 1956
wavefunction for liquid Helium



Iterative backflow \rightarrow deep residual network \rightarrow continuous normalizing flow

Taddei et al, PRB '15 E Commun. Math. Stat 17', Harbor et al 1705.03341, Lu et al 1710.10121, Chen et al, 1806.07366

Point Transformations and the Many Body Problem*

M. EGER† AND E. P. GROSS

Brandeis University, Waltham, Massachusetts

An investigation is made of possible uses of many dimensional coordinate transformations in the quantum many-body problem. The transformed Hamiltonian is quadratic in the momenta with a space dependent metric. The original potential energy undergoes alteration and an additional “metric” potential energy appears. A relatively complete analysis of the transformed original potential is made, and the coordinate transformation can be used to suppress undesirable features of the original potential. For bosons one can attempt to directly map a complete set of noninteracting states onto approximate eigenstates of the system with interactions. Contact is made with a theory of weakly interacting bosons. In the general case it emerges that a given transformation uniquely fixes all the spatial correlation functions, which can be explicitly computed. The extended point transform can then be used as a link between diverse experimental quantities. The full use of the transformation to compute from first principles requires adequate approximations to the Jacobian and the inverse transform. These problems are not studied.

√flow

materializes this dream

The objective function of variational density matrix

$$\rho = \sum_n p_n |\psi_n\rangle\langle\psi_n|$$

$$F = \mathbb{E}_{n \sim p_n} \left[k_B T \ln p_n + \mathbb{E}_{\mathbf{R} \sim |\psi_n(\mathbf{R})|^2} \left[\frac{H\psi_n(\mathbf{R})}{\psi_n(\mathbf{R})} \right] \right]$$

↓
Boltzmann
distribution

↓
Born
probability

Jointly optimize p_n and $\psi_n(\mathbf{R})$ to minimize the variational free energy

The deep variational free energy approach

1802.02840, PRL '18

1910.00024, PRX '20

1912.11381, MLST' 21

2201.03156, SciPost '23

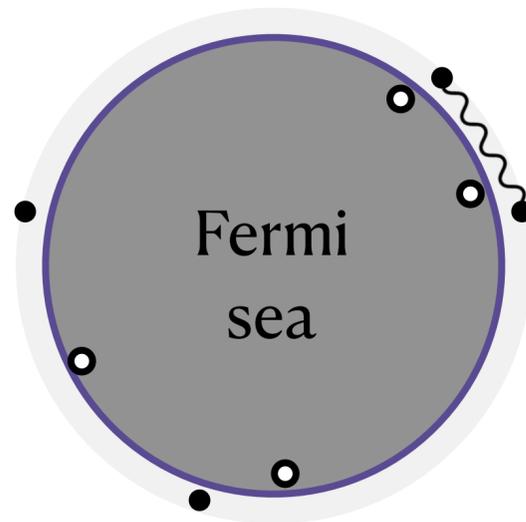
1809.10606, PRL '19

2209.06095, PRL '23

2105.08644, JML '22

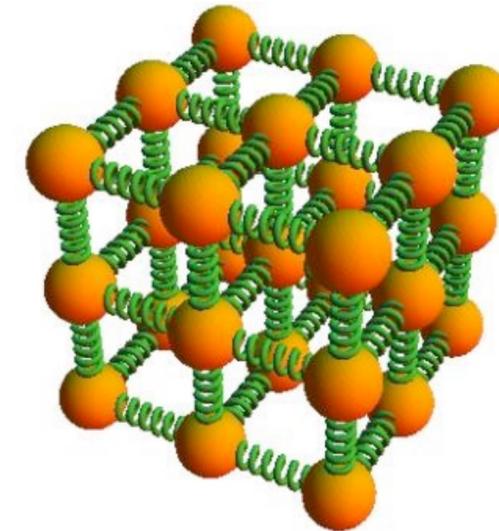
2403.12518, JCP '24...

Low-temperature properties of interacting electrons
(~50 electrons)



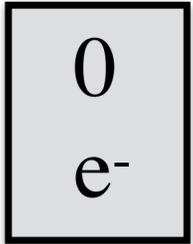
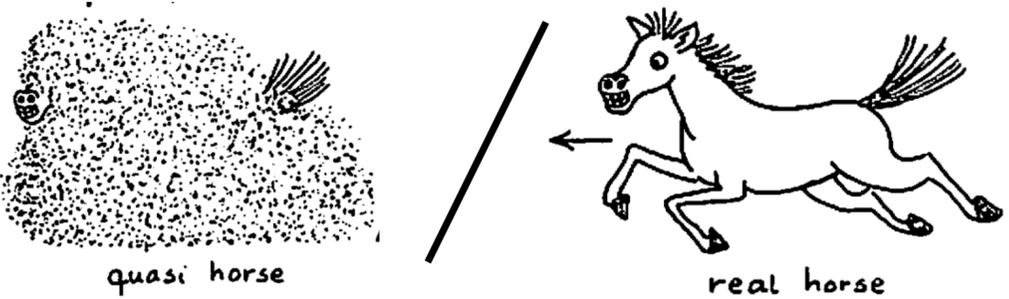
$$\rho = \sum_n U |\phi_n\rangle p_n \langle \phi_n| U^\dagger$$

Vibrational spectra of molecules and solids
(~500 atoms)



A computational framework taking in account of electron correlation, thermal effect, and anharmonic lattices for free energy, entropy, and excitation spectra

Deep variational free energy approach: resolving puzzles



Quasi-particle effective mass
contradicting experiments

$$m^*/m > 1$$

VOLUME 91, NUMBER 4 PHYSICAL REVIEW LETTERS week ending 25 JULY 2003

Spin-Independent Origin of the Strongly Enhanced Effective Mass in a Dilute 2D Electron System

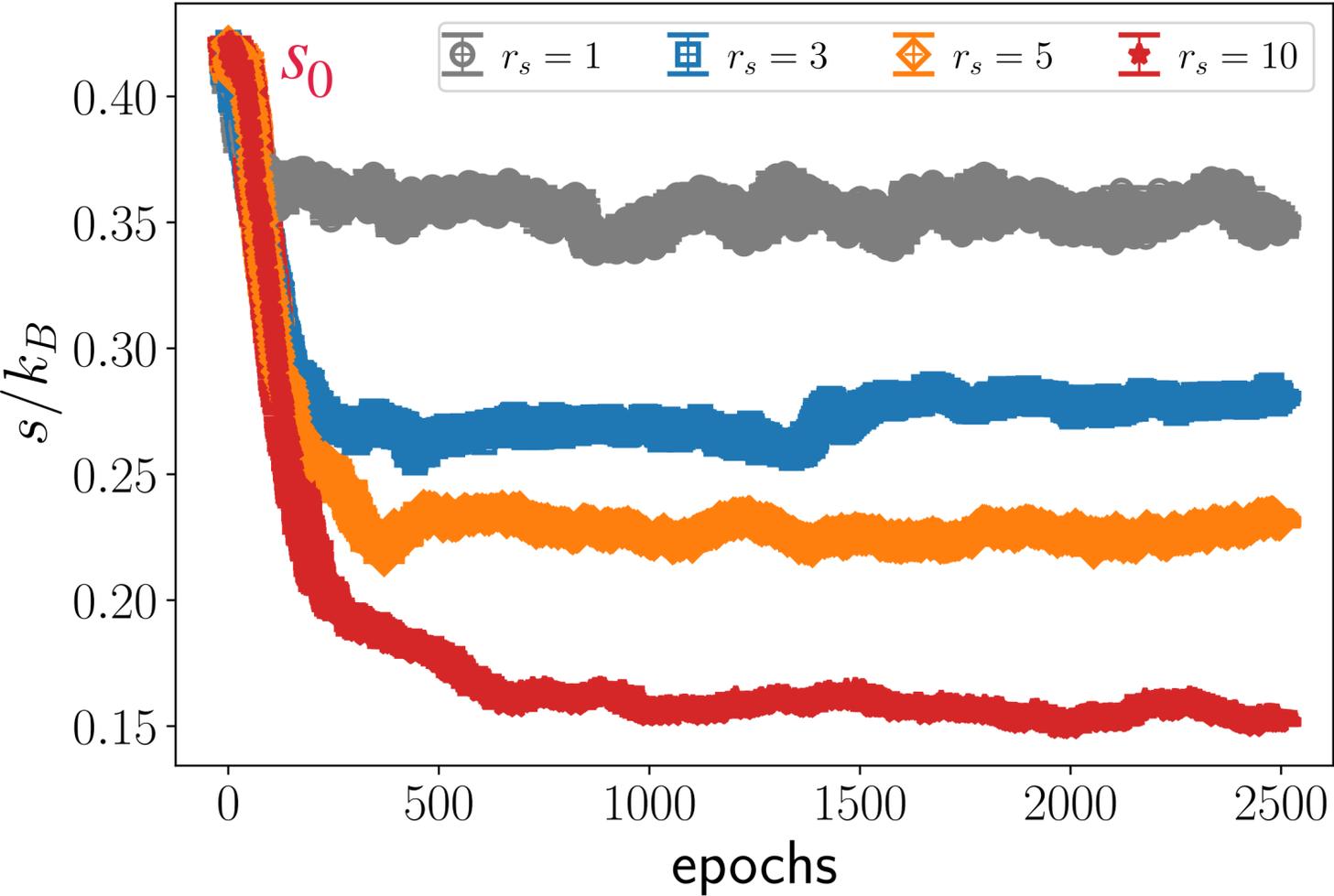
$$m^*/m < 1$$

PRL 101, 026402 (2008) PHYSICAL REVIEW LETTERS week ending 11 JULY 2008

Effective Mass Suppression in Dilute, Spin-Polarized Two-Dimensional Electron Systems

Medini Padmanabhan, T. Gokmen, N. C. Bishop, and M. Shayegan
Department of Electrical Engineering, Princeton University, Princeton, New Jersey 08544, USA
(Received 19 September 2007; published 7 July 2008)

Hao Xie et al, SciPost Physics '23 $\frac{m^*}{m} = \frac{s}{s_0} < 1$
Thermal entropy of 2D electron gas

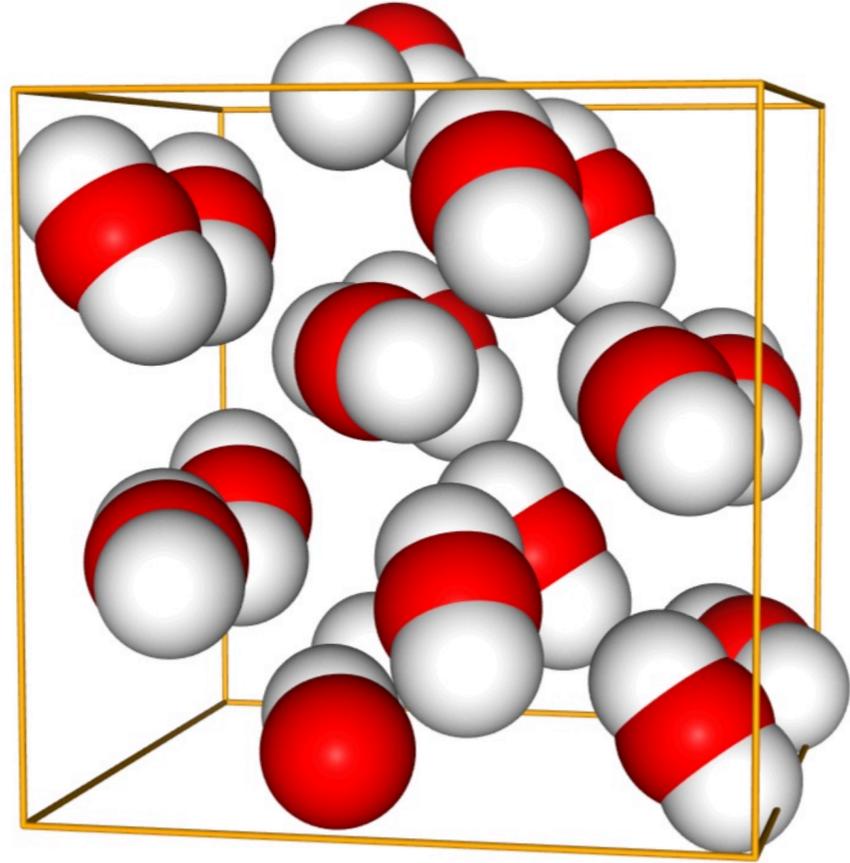
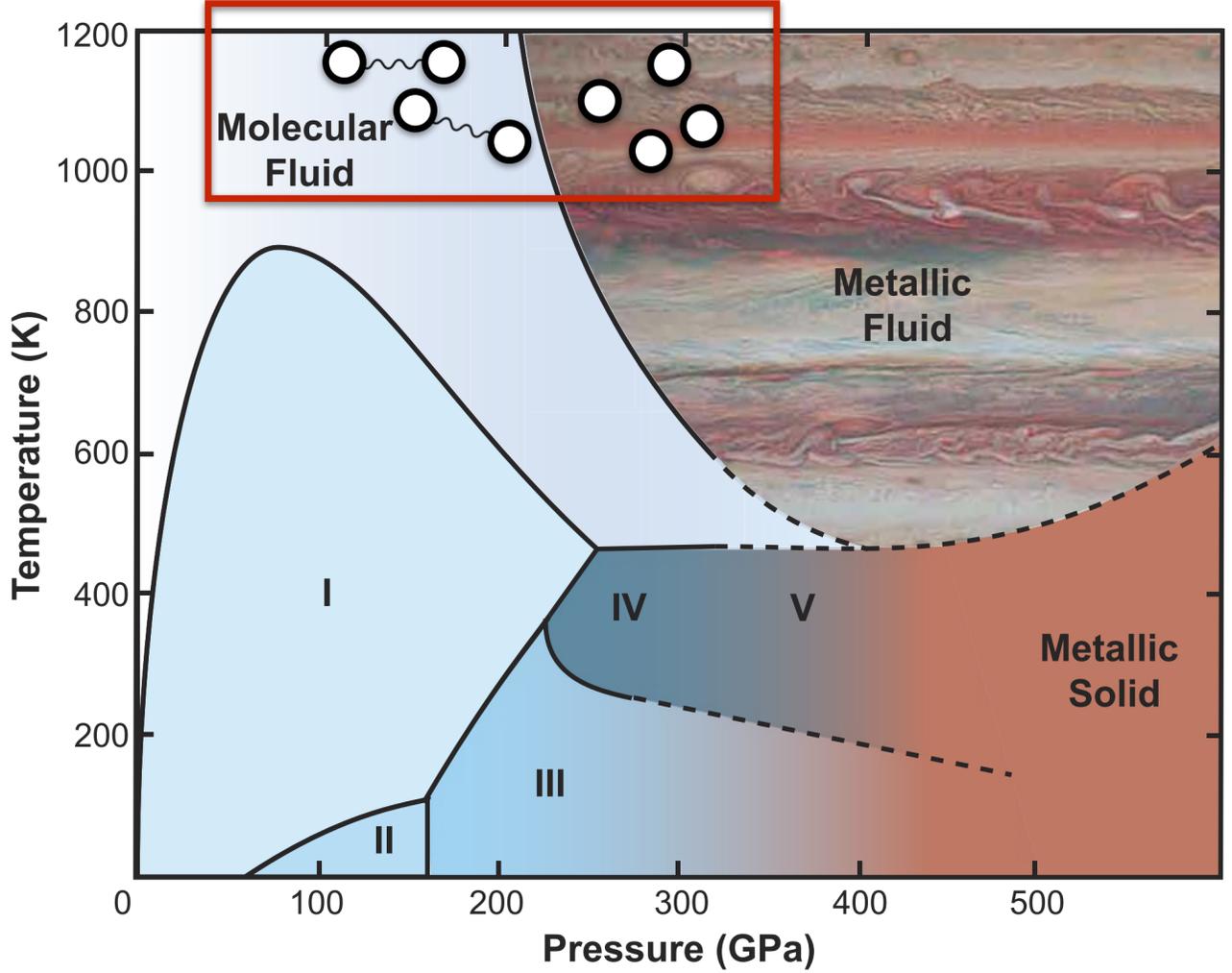


Deep variational free energy approach: **making discoveries**

1
H

Knudson et al Science '15, Celliers et al Science '18
Mazzola et al Nat.Comm. '14, Pierleoni et al, PNAS '16
Cheng et al, Nature '20, Karasiev et al, Nature '21

Xinyang Dong et al, 2024



Hydrogen phase diagram

High-temperature solids of dense hydrogen

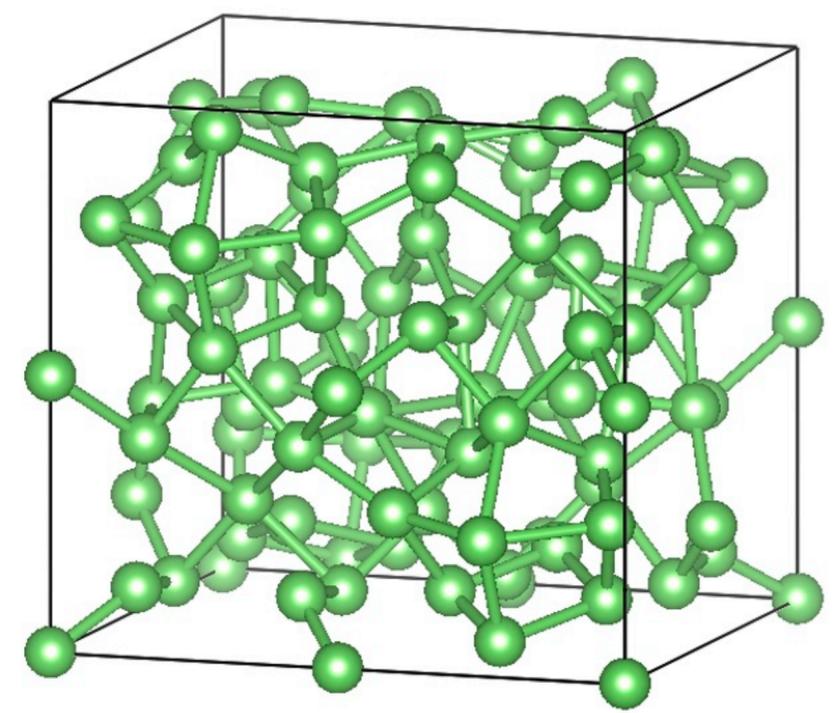
Gregoryanz et al, Matter Radiat. Extremes, 2020

see also Niu et al, PRL '23, Goswami et al, 2411.15665

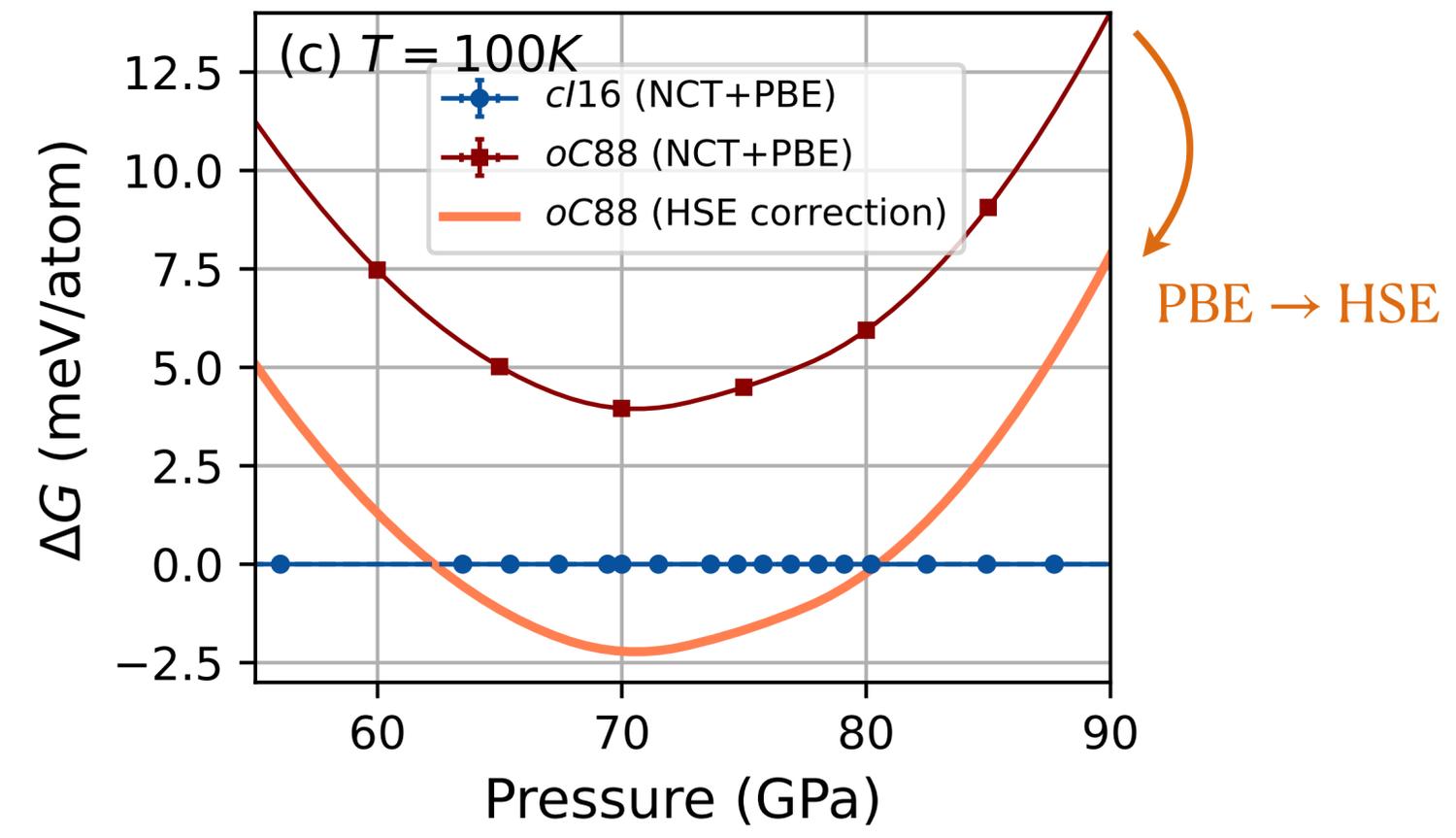
Deep variational free energy approach: clarifying mechanism

3
Li

Guillaume et al, Nature physics, '11
Marqués et al, PRL '11, Gorelli et al PRL '12



Qi Zhang et al, 2412.12451



Our calculation does **NOT** reproduce the experimentally observed Oc88 phase of Lithium, **which contradicts consensus**

- Thermal effect
- Quantum anharmonicity
- DFT functional error for bad metal

Generative AI for **It**

①

$$p(X|y) \propto p(X)p(y|X)$$

Matter inverse design
Exploiting intuitions in data

②

$$F[\rho] = E - TS$$

Nature's cost function
Variational free energy is finally practical

Turning physics problems into stochastic optimization

Leverages the deep learning engine



The Universe as a generative model

$$S = \int d^4x \sqrt{-g} \left[\frac{m_p^2}{2} R - \frac{1}{4} F_{\mu\nu}^a F^{\mu\nu a} \right. \\ \left. + i \bar{\Psi}^i \gamma^\mu D_\mu \Psi^i + \left(\bar{\Psi}_L^i V_{ij} \Phi \Psi_R^j + \text{h.c.} \right) \right. \\ \left. - |D_\mu \Phi|^2 - V(\Phi) \right]$$

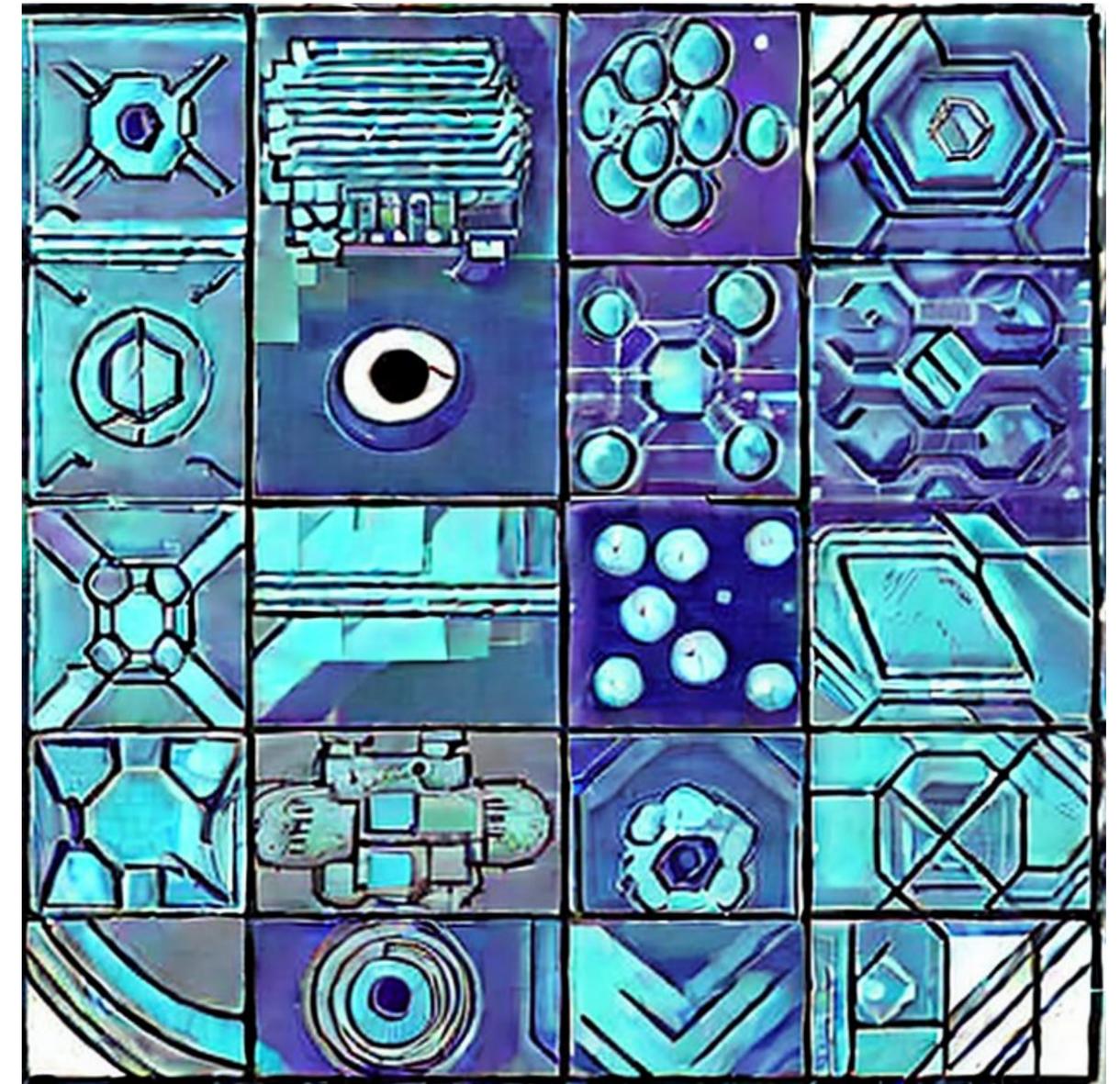
Thank you!



Discovering physical laws: **learning** the action
Solving physical problems: **optimizing** the action

A crash course offered at IOP 2023 spring

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 ?



Machine learning for physicists

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