# Autoregressive model: alphabets, actions, and atoms

## A unified perspective to LLM, RL, and atomistic modeling

Lei Wang (王磊)

Institute of Physics, CAS

https://wangleiphy.github.io

# Plan

① Generative AI

② Autoregressive models

-principle

-architecture

-training

-inference

③ Applications

# Probabilistic modeling with generative AI

$$p(X)$$

pixels, words, atoms, ...

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



DaLL-E
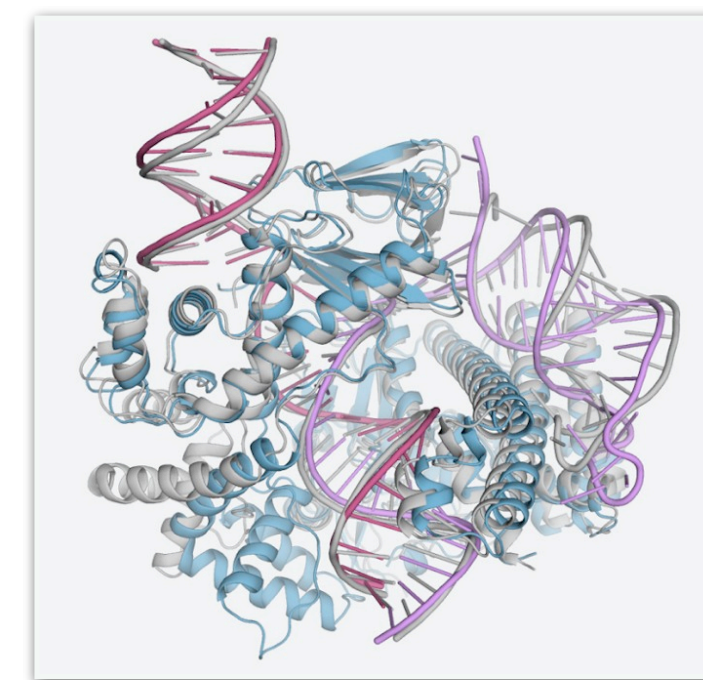
ChatGPT

AlphaFold3

# Discriminative AI is not enough

$$p(y \mid X)$$



cat
**dog**

$X$: pixels　　　　　$y$: label

$$\nabla_{\text{pixels}}\, p(\text{dog} \mid \text{pixels})$$



Deepdream: http://googleresearch.blogspot.ch/2015/06/inceptionism-going-deeper-into-neural.html

# Bayes rule

posterior  prior  likelihood

$$p(\boldsymbol{X} \,|\, y) \propto p(\boldsymbol{X}) p(y \,|\, \boldsymbol{X})$$

Inverse design          Forward prediction

# Probability theory 101

Conditional probability $\qquad\qquad p(y \mid X)$

Joint probability $\qquad\qquad\qquad p(X, y)$

Product rule $\qquad p(X, y) = p(y \mid X)p(X)$

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

generative model $p(X)$

data $X$

# Two sides of the same coin

**Generative modeling**



"learn from data"

**Mamixmum likelihood estimation**

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

**Statistical physics**



"learn from energy"

**Variational free energy**

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

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

# Kullback–Leibler divergence

$$\mathbb{KL}(\pi \parallel p) \equiv \sum_X \pi(X) \big[\ln \pi(X) - \ln p(X)\big]$$

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

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

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

# Learn from data

$$\pi(X) \propto \sum_{d \in \text{dataset}} \delta(X - d)$$

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

target     model        Maximum likelihood estimation

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

# Learn from Energy

$$\pi(X) \propto e^{-E/k_B T}$$

$$\min_\theta \mathbb{KL}(p_\theta \parallel \pi) \iff \min_\theta \left\{ \mathop{\mathbb{E}}_{X \sim p_\theta(X)} \left[ E(X) + k_B T \ln p_\theta(X) \right] \right\}$$

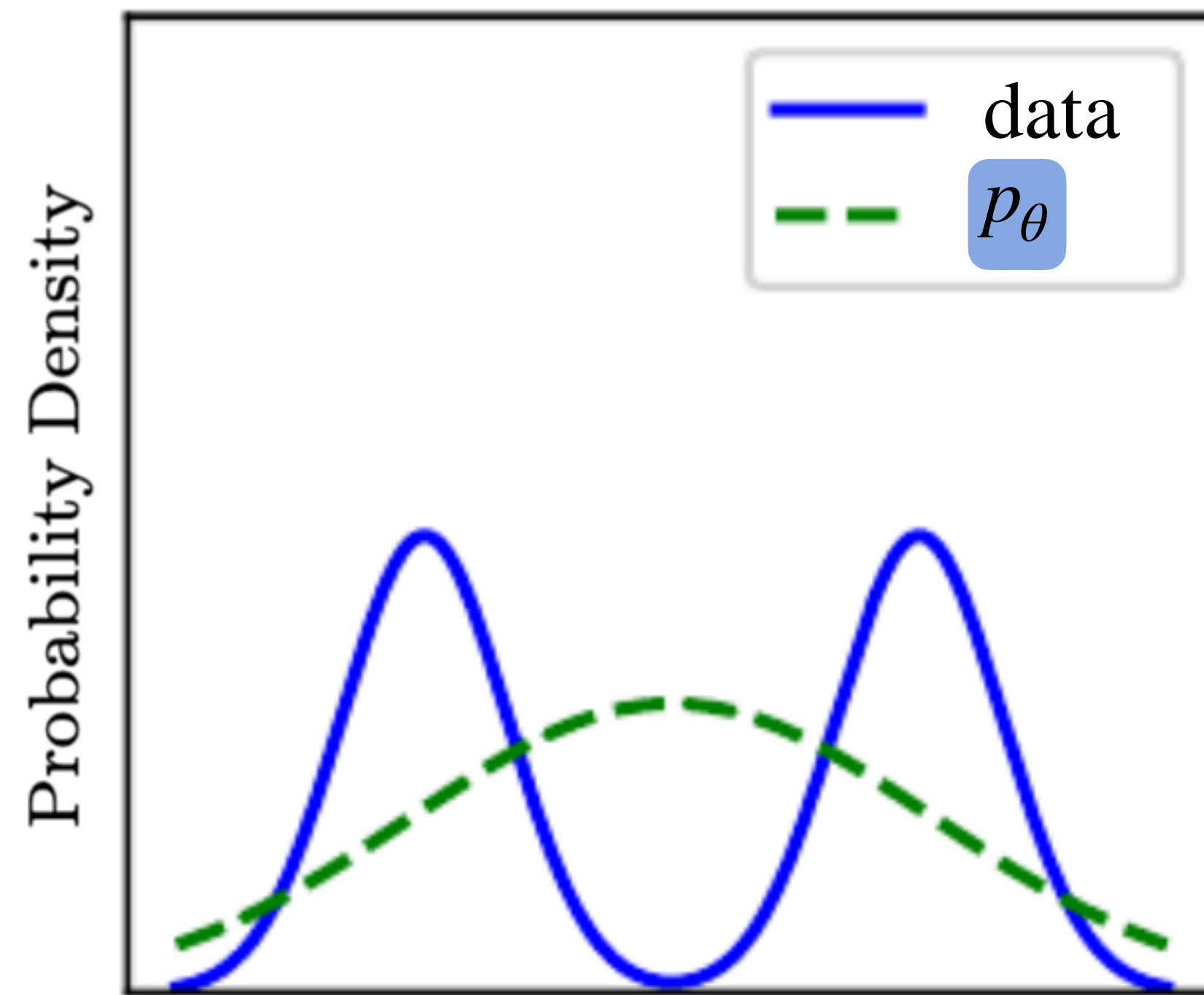model      target

Variational free energy

The lower bound is the true free energy: exact solution

# Forward KL or Reverse KL ?

## Maximum likelihood estimation

$$\min_{\theta} \mathbb{KL}(\text{data} \parallel \boxed{p_\theta})$$
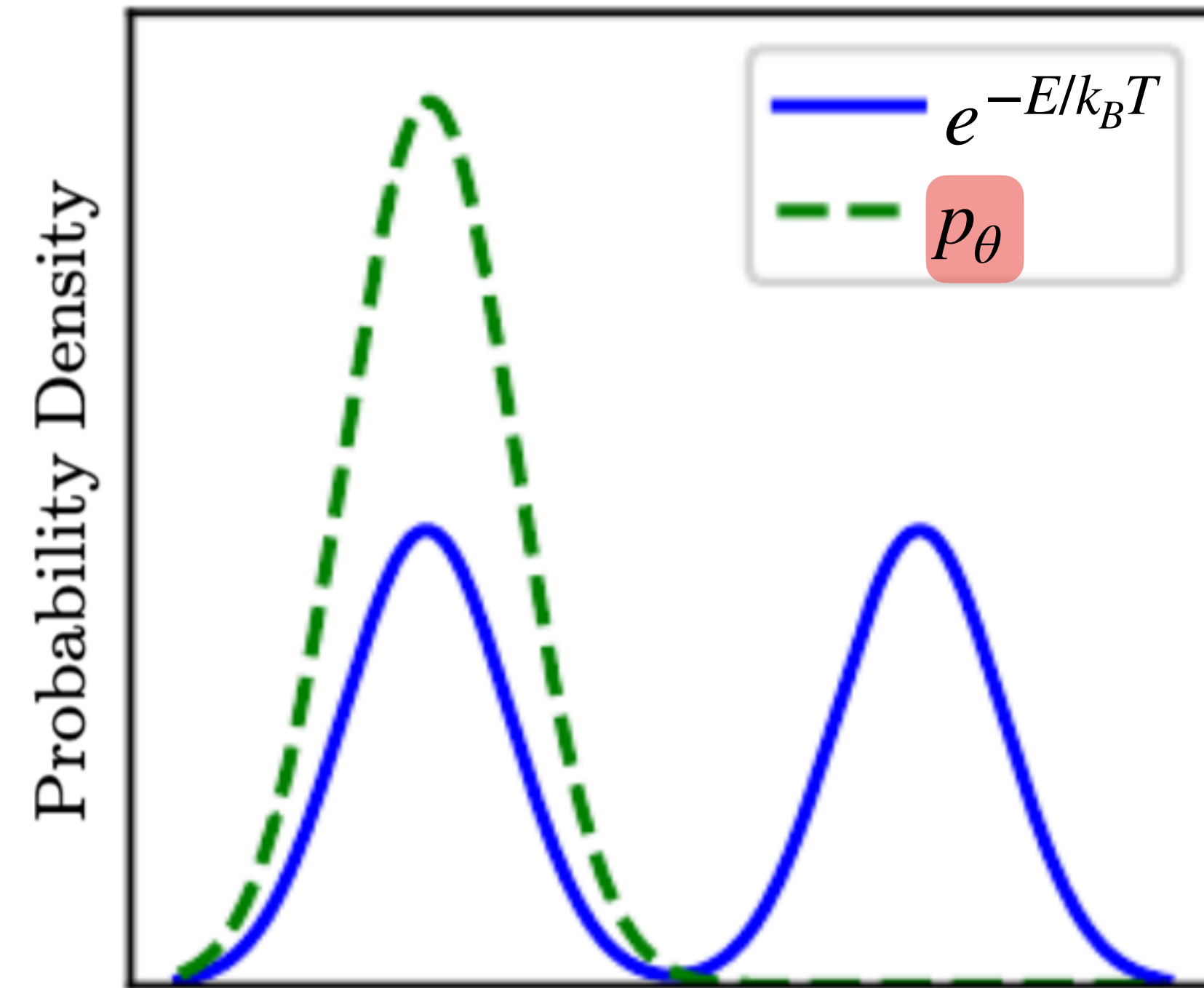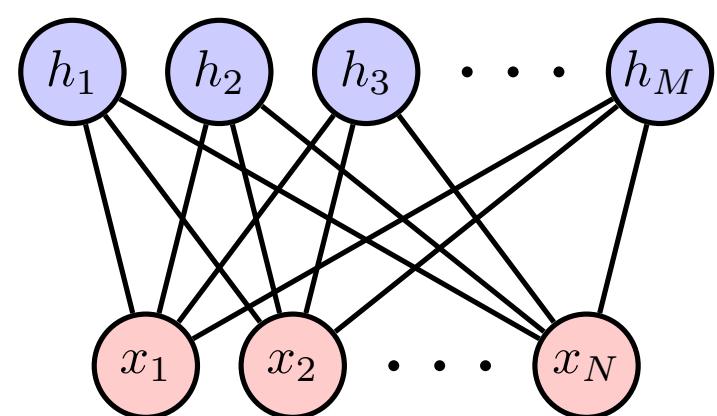
Mode covering



Failure mode: hallucination

## Variational free energy

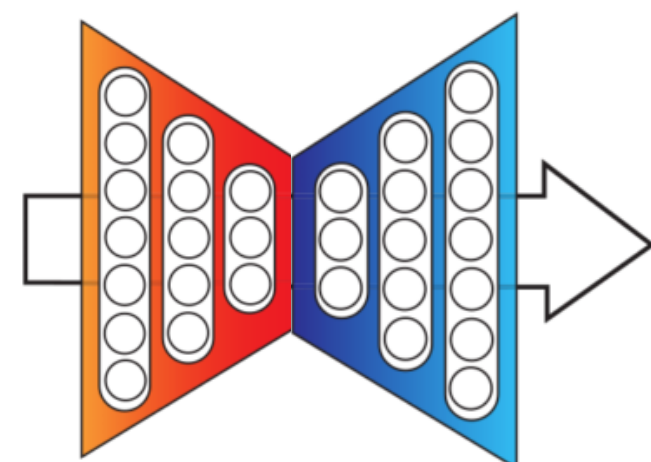$$\min_{\theta} \mathbb{KL}(\boxed{p_\theta} \parallel e^{-E/k_B T})$$
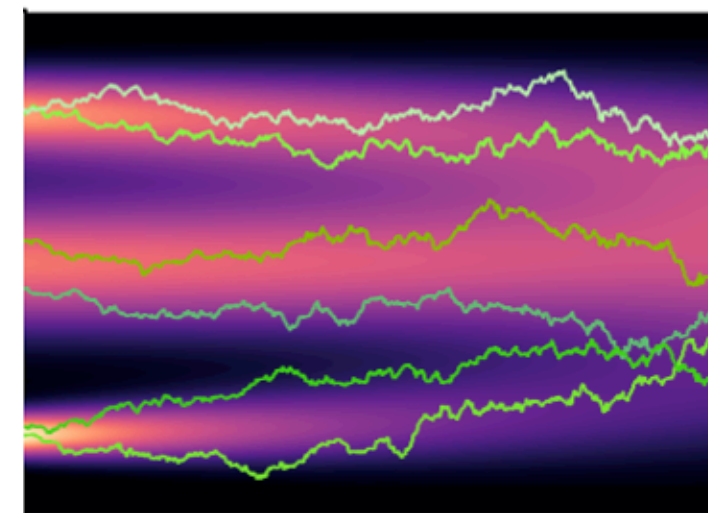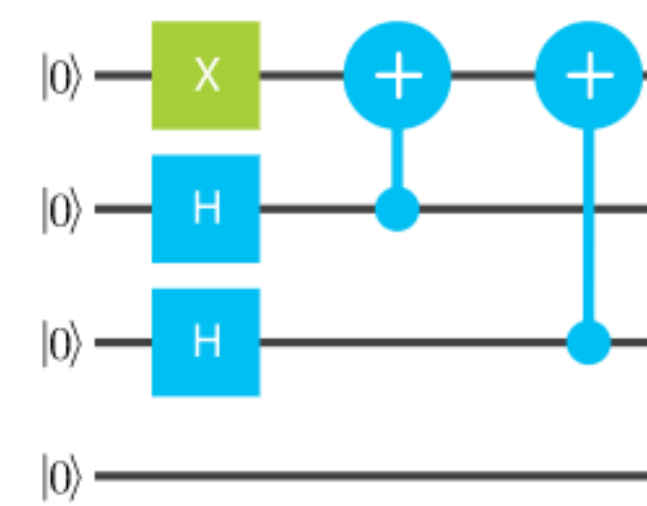
Mode seeking



Failure mode: local minima

Goodfellow et al, Deep Learning

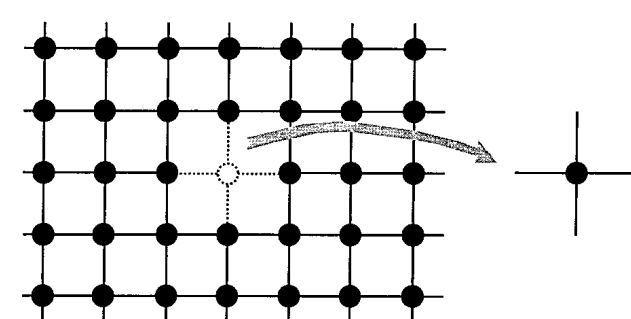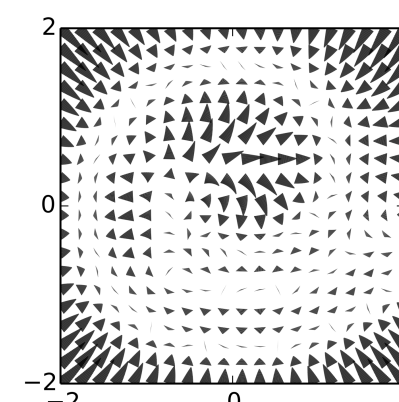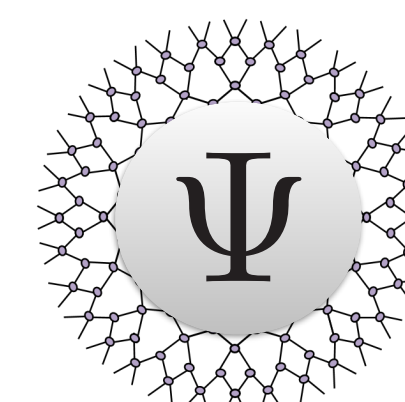| Boltzmann Machine | Variational Autoencoder | Diffusion Model | Born Machine | Flow Matching |
|---|---|---|---|---|
| 1985 | 2013 | 2015 | 2017 | 2022 |
| Monte Carlo Ising model | Variational mean field | Nonequilibrium thermodynamics | Tensor networks Quantum circuits | Fluid optimal transportation |

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

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

**Leverage the power of modern generative models for science**

# Boltzmann machines

$$p(X) = \frac{e^{-E(X)}}{Z}$$

Ackley, Hinton, Sejnowski, 1985

**Learn**

$h_1$  $h_2$  $h_3$  $\cdots$  $h_M$

$x_1$  $x_2$  $\cdots$  $x_N$

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

$$\nabla_\theta \mathscr{L} = \mathop{\mathbb{E}}_{X\sim\text{dataset}}\left[\nabla_\theta E\right] - \mathop{\mathbb{E}}_{X\sim p(X)}\left[\nabla_\theta E\right]$$

# Boltzmann machines

$$p(X) = \frac{e^{-E(X)}}{Z}$$

Ackley, Hinton,
Sejnowski, 1985

**Learn**

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

$$h_1 \quad h_2 \quad h_3 \quad \cdots \quad h_M$$

$$x_1 \quad x_2 \quad \cdots \quad x_N$$

$$\nabla_\theta \mathcal{L} = \mathop{\mathbb{E}}_{X\sim\text{dataset}}\left[\nabla_\theta E\right] - \mathop{\mathbb{E}}_{X\sim p(X)}\left[\nabla_\theta E\right]$$

# Boltzmann machines

$$p(X) = \frac{e^{-E(X)}}{Z}$$

Ackley, Hinton,
Sejnowski, 1985

**Learn**

**Generate**

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

$$X \sim p(X)$$

$h_1$ $h_2$ $h_3$ $\cdots$ $h_M$

$x_1$ $x_2$ $\cdots$ $x_N$

$$\nabla_\theta \mathscr{L} = \mathbb{E}_{X \sim \text{dataset}} \left[ \nabla_\theta E \right] - \mathbb{E}_{X \sim p(X)} \left[ \nabla_\theta E \right]$$
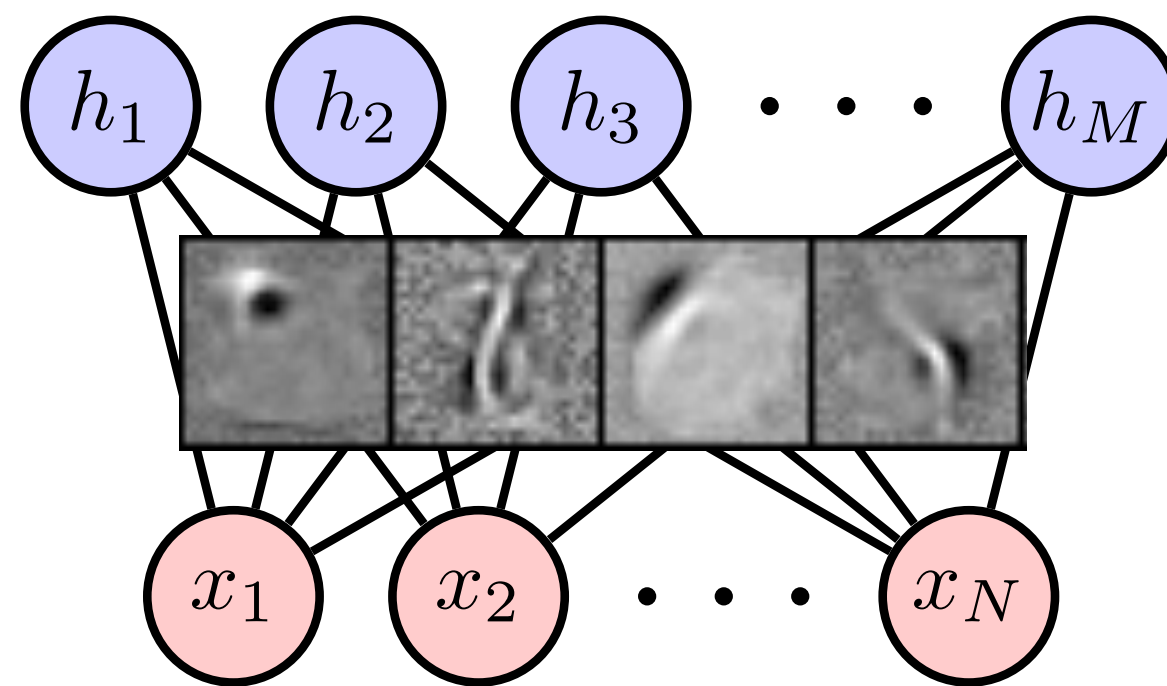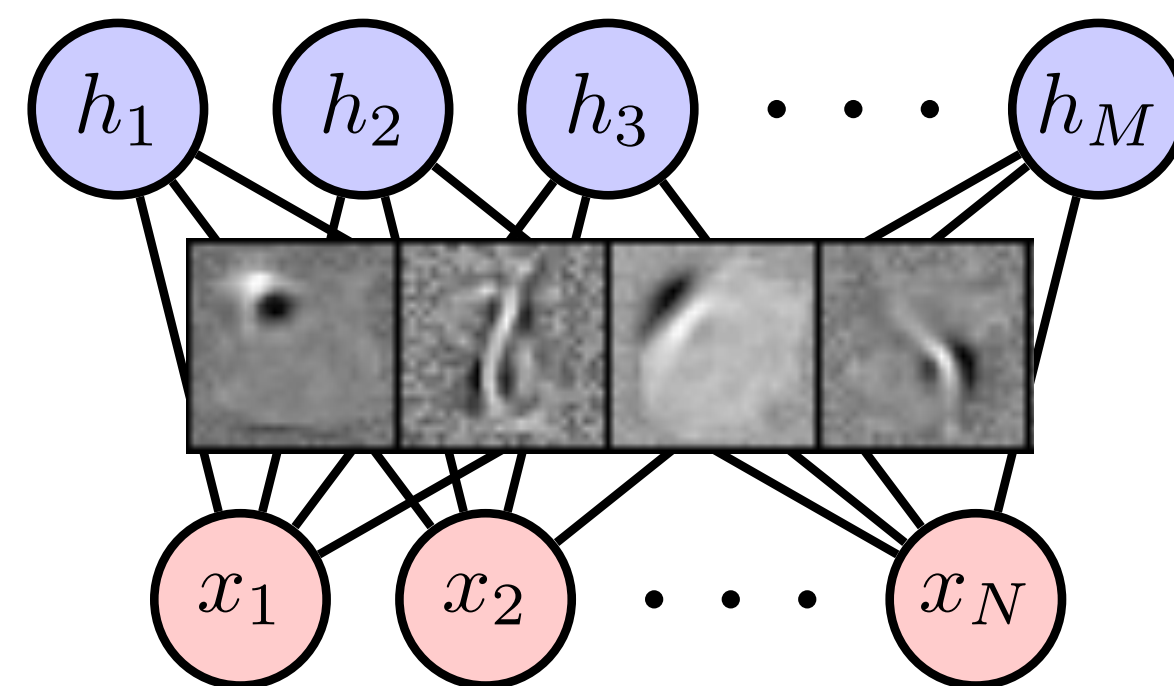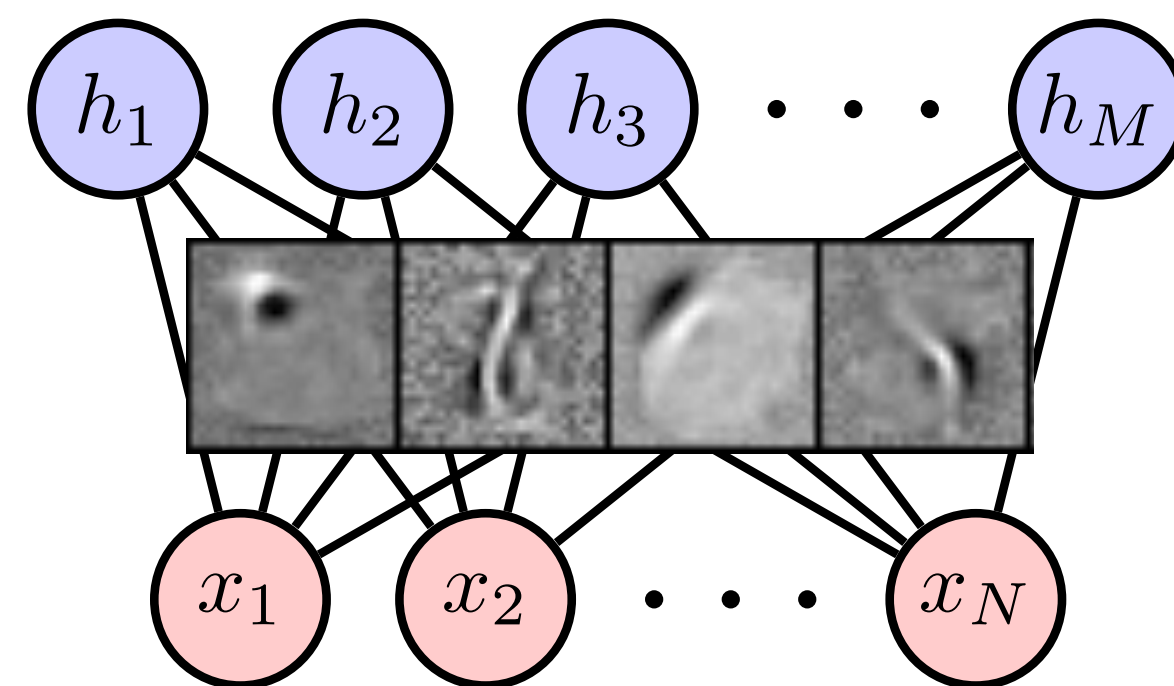
# Boltzmann machines

$$p(X) = \frac{e^{-E(X)}}{Z}$$

Ackley, Hinton,
Sejnowski, 1985

**Learn**

**Generate**

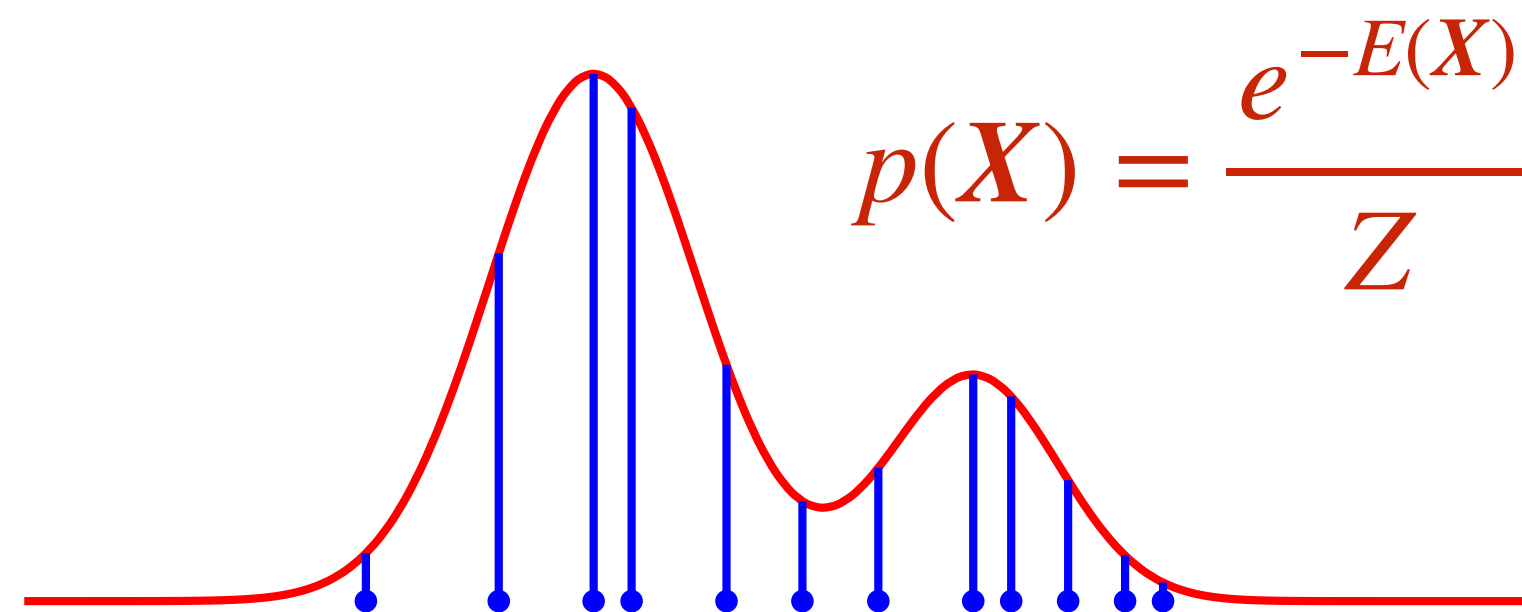$$\mathscr{L} = \mathbb{E}_{X \sim \text{data}} \left[ -\ln p(X) \right]$$

$$X \sim p(X)$$

$$\nabla_\theta \mathscr{L} = \mathop{\mathbb{E}}_{X \sim \text{dataset}} \left[ \nabla_\theta E \right] - \mathop{\mathbb{E}}_{X \sim p(X)} \left[ \nabla_\theta E \right]$$

# Boltzmann machines 🥇

$$p(X) = \frac{e^{-E(X)}}{Z}$$

Ackley, Hinton, Sejnowski, 1985

## GAUSSIAN-BERNOULLI RBMs WITHOUT TEARS 😂

**Renjie Liao**[*,1]**, Simon Kornblith**[2]**, Mengye Ren**[3]**, David J. Fleet**[2,4,5]**, Geoffrey Hinton**[2,4,5]

$$\nabla_\theta \mathcal{L} = \mathop{\mathbb{E}}_{X \sim \text{dataset}} \left[\nabla_\theta E\right] - \mathop{\mathbb{E}}_{X \sim p(X)} \left[\nabla_\theta E\right]$$

# So, why bother ?

$$p(X) \geq 0$$

**Normalization ?**   **Sampling ?**

$$\sum_X p(X)$$   $$X \sim p(X)$$

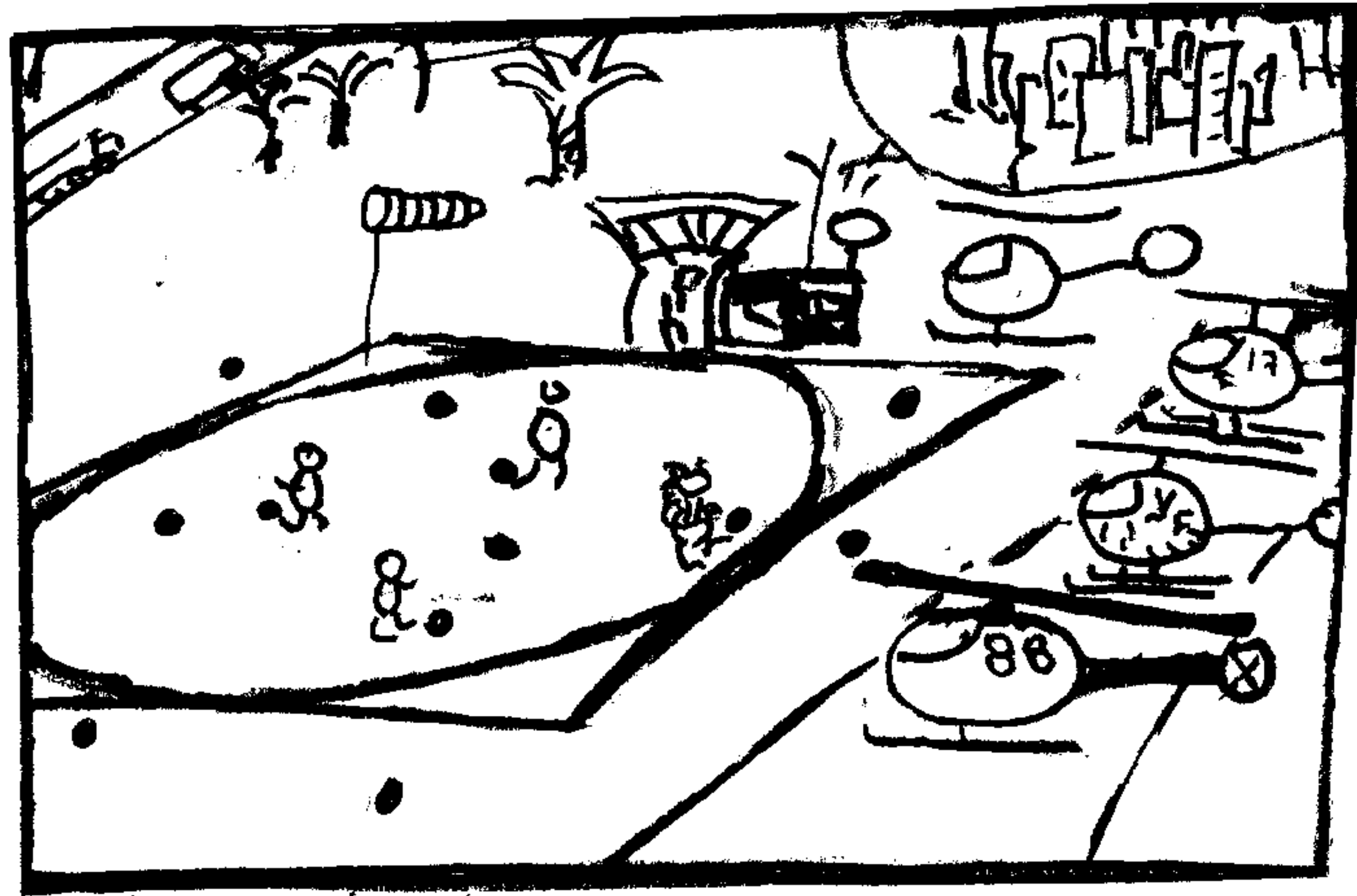# The difficulty of normalization



$$Z = \sum_X e^{-E(X)}$$

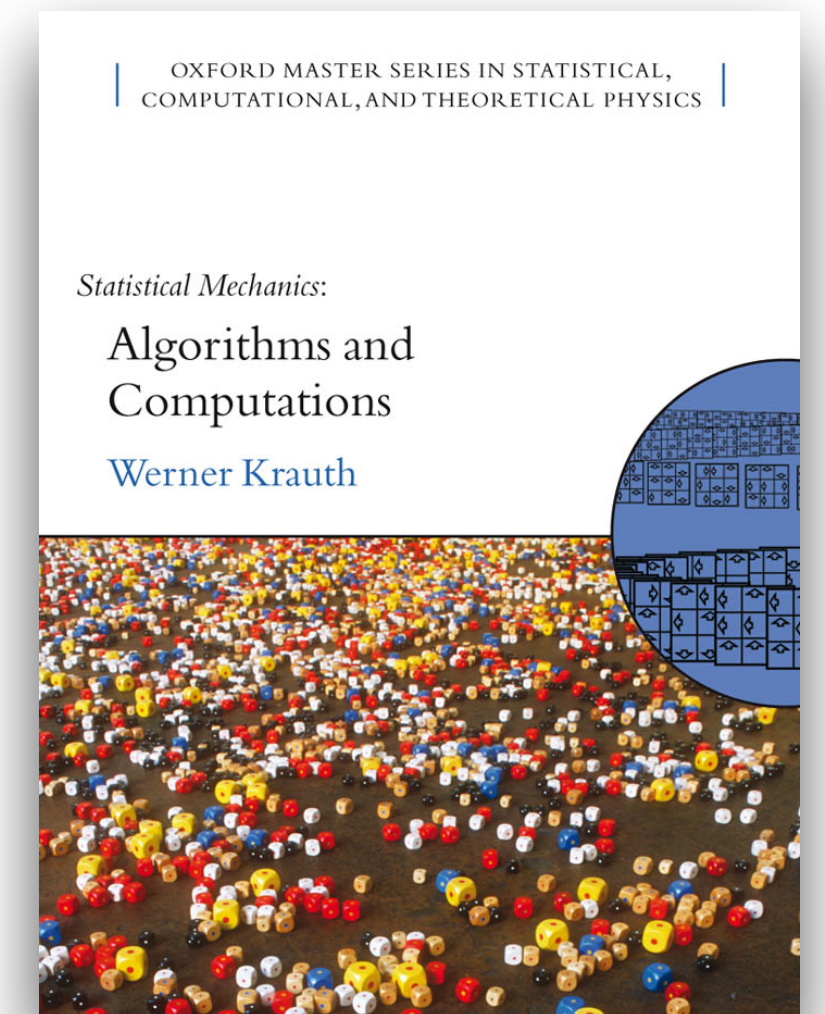"Intractable" partition function $Z$
appears widely in machine learning and statistical physics (entropy and free energy calculation)

# The difficulty of sampling

$$X \sim p(X)$$



Adults computing the number $\pi$ 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**

$p(X)$

Direct
GAN

Explicit density

Implicit density

**Tensor
Networks**

Han et al, PRX '18

Tractable density

-Fully visible belief nets
-NADE
**Autoregressive**
-MADE
**model**
-Pixel

-Change of variables
models (nonlinear ICA)

**Flow model**

Approximate density

Variational

Markov Chain

Variational autoencoder    Boltzmann machine

Markov Chain
GSN

**Quantum
Circuits**

Liu et al PRA '18

$U$

+**Diffusion models**

# Autoregressive model

$$p(X) = p(x_1)p(x_2 \mid x_1)p(x_3 \mid x_1, x_2)\cdots$$

*"... the murderer is ___"*

$$p(\_ \mid \dots)$$

**Normalization**

$$\sum_{x_1} p(x_1) \sum_{x_2} p(x_2 \mid x_1) \sum_{x_3} p(x_3 \mid x_1, x_2)\cdots$$

**Sampling**

$$x_1 \sim p(x_1)$$

$$x_2 \sim p(x_2 \mid x_1)$$

$$\vdots$$

# Implementation: autoregressive masks

Masked Autoencoder Germain et al, 1502.03509



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

https://github.com/karpathy/pytorch-made

# Implementation: autoregressive masks

Mask convoluti                                                                trix

PixelCNN, van den Oord et al, 1601.06759                    Causal transformer, Vaswani et al 1706.03762

# The autoregressive transformer



Masked attention matrix => lower triangular Jacobian matrix => autoregressive model
Great at capturing long-range dependence; friendly to backpropagation and GPUs

Transformers from scratch, https://peterbloem.nl/blog/transformers

$x_i$  Sequence  Embedding  Embedding  Sequence  $y_i$

① Self-attention

$$y_i = \sum_j \alpha(x_i, x_j) x_j$$

attention weights

Mixing signal along
sequence direction

$x_i$  Sequence  Embedding  Embedding  Sequence  $y_i$

② Multi-layer perceptrons

$$y_i = \sigma \circ \cdots \sigma(x_i W + b)$$

Nonlinear activation

Transforming
signals locally

# Feynman's backflow as an attention layer

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

Quasi-particle
coordinates

Feynman & Cohen 1956
wavefunction for liquid Helium

Electron
coordinates

Each particle attends to its surrounding
and dresses up as a quasi-particle

c.f. Lu et al, 1906.02762  transformer as convection-diffusion multi-particle dynamics

```python
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

# Params count in GPT3

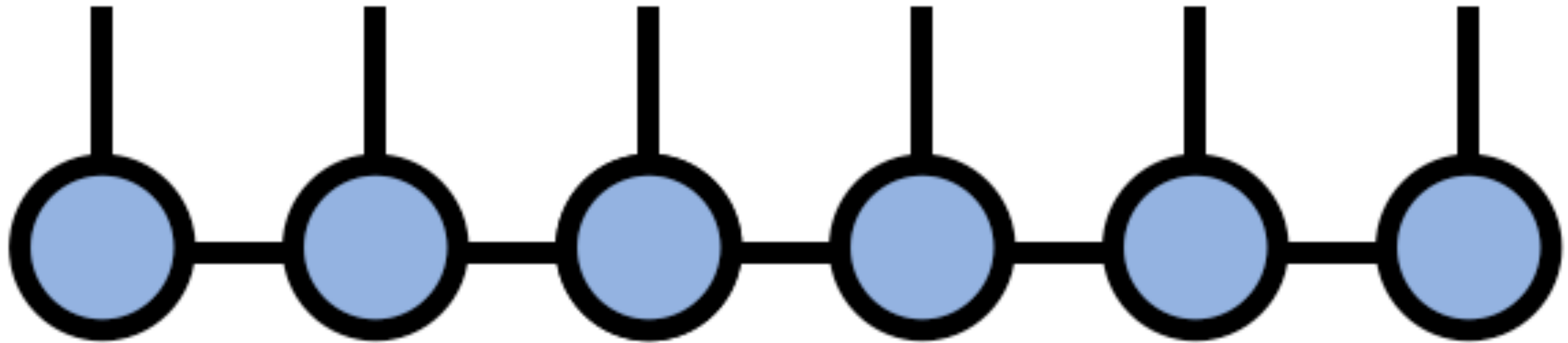3Blue1Brown, https://youtu.be/9-Jl0dxWQs8?t=940



175B in total

Independent of the context length (4096) and vocabulary size (50257, almost)

# An MPS analog

Vocabulary

$d = 50257$



Model size $\chi = 12288$

Contex length $L = 4096$

# Aside: SVD attack



**What is $\chi$?** 🔒

If you only have access to $d$-dim logits
via the LLM API

$$\text{svd} \begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix}_{\vdots \times d}$$

$\chi < d$

Carlini et al, Stealing Part of a Production Language Model, 2403.06634
Finlayson et al, Logits of API-Protected LLMs Leak Proprietary Information, 2403.09539

# Scaling law of the loss function

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

Kaplan et al, 2001.08361



**Compute**
PF-days, non-embedding

$L = (C_{\min}/2.3 \cdot 10^8)^{-0.050}$

**Dataset Size**
tokens

$L = (D/5.4 \cdot 10^{13})^{-0.095}$

**Parameters**
non-embedding

$L = (N/8.8 \cdot 10^{13})^{-0.076}$

"Predict resouces needed to sovle increasingly difficult tasks" — Sam McCandlish, Aspen talk '19

https://sites.google.com/view/phys4ml/home

# A trillion $ plot

# A visionary discussion section

Scaling Laws for Neural Language Models, Kaplan et al, 2001.08361
Scaling Laws for Autoregressive Generative Modeling, Henighan et al, 2010.14701

It is natural to conjecture that the scaling relations will apply to other generative modeling tasks with a maximum likelihood loss, and perhaps in other settings as well. To this purpose, it will be interesting to test these relations on other domains, such as images, audio, and video models, and perhaps also for random network distillation. At this point we do not kno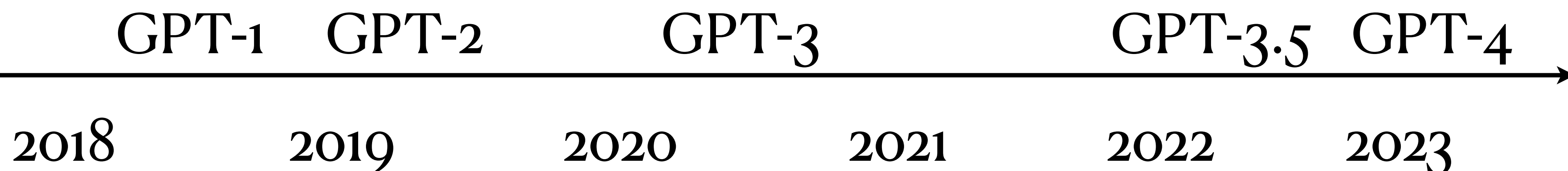w which of our results depend on the structure of natural language data, and which are universal. It would also be exciting to find a theoretical framework from which the scaling relations can be derived: a 'statistical mechanics' underlying the 'thermodynamics' we have observed. Such a theory might make it possible to derive other more precise predictions, and provide a systematic understanding of the limitations of the scaling laws.

In the domain of natural language, it will be important to investigate whether continued improvement on the loss translates into improvement on relevant language tasks. Smooth quantitative change can mask major qualitative improvements: "more is different". For example, the smooth aggregate growth of the economy provides no indication of the specific technological developments that underwrite it. Similarly, the smooth improvements in language model loss may hide seemingly qualitative changes in capability.

GPT-1   GPT-2          GPT-3                    GPT-3.5   GPT-4

2018          2019          2020          2021          2022          2023

# Emergent abilities: more is different

Wei et al, 2206.07682

https://www.jasonwei.net/
blog/emergence

Legend: LaMDA • GPT-3 ■ Gopher ◆ Chinchilla ▲ PaLM ⬠ Random ---



(A) Mod. arithmetic
(B) IPA transliterate
(C) Word unscramble
(D) Persian QA
(E) TruthfulQA
(F) Grounded mappings
(G) Multi-task NLU
(H) Word in context

Model scale (training FLOPs)

Avogadro
constant
number
of FLOPs

# Autoregressive model is more than language modeling

"Language" => token sequence => bitstream => **ANYTHING**
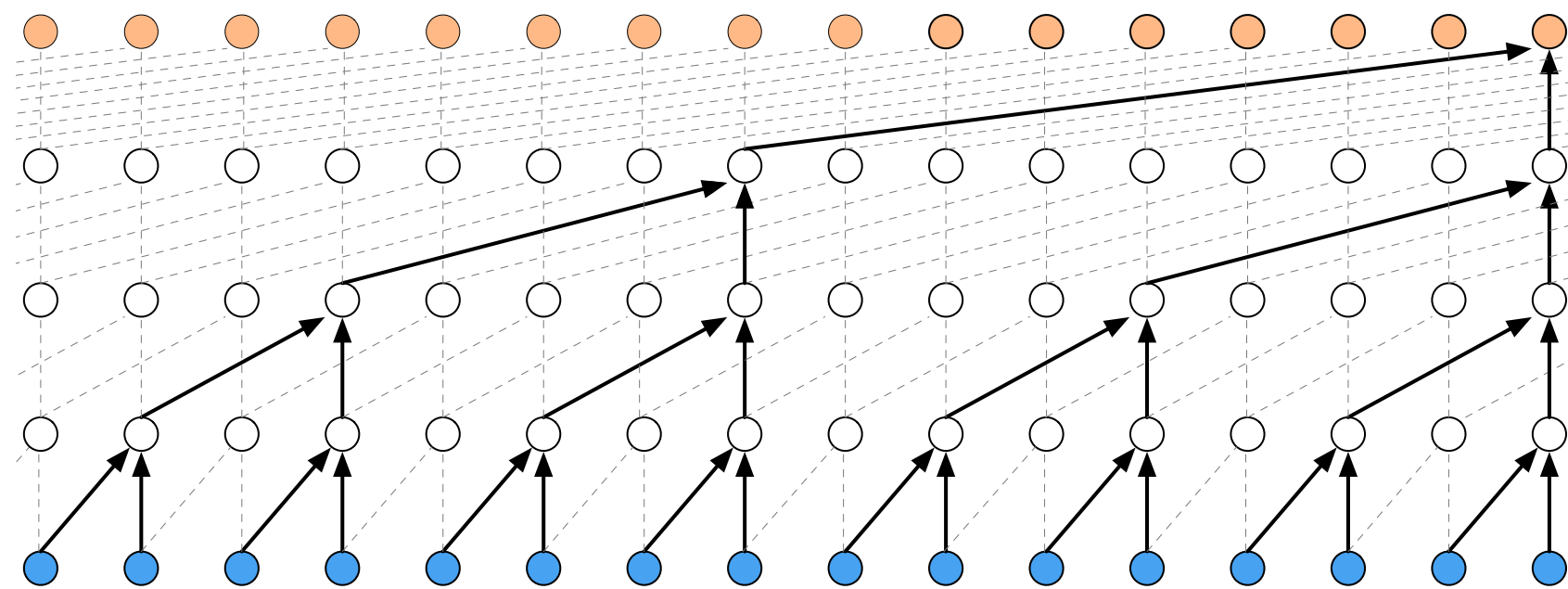
**Speech**: WaveNet 1609.03499



**Molecules**: 1810.11347



Multi-scale c

# Autoregressive models for images



Chen et al, PMLR '20, Esser et al, 2012.09841

Reed et al, 1703.03664
Tian et al, 2404.02905, Li et al, 2502.17437

Next pixel (patch) prediction

Next **scale** prediction

What is the suitable 1D ordering of 2D images ?

# Autoregressive model for images

"Language" => token sequence => bitstream => **ANYTHING**



Compress

Further compress

JPEG-LM

westlake.jpeg
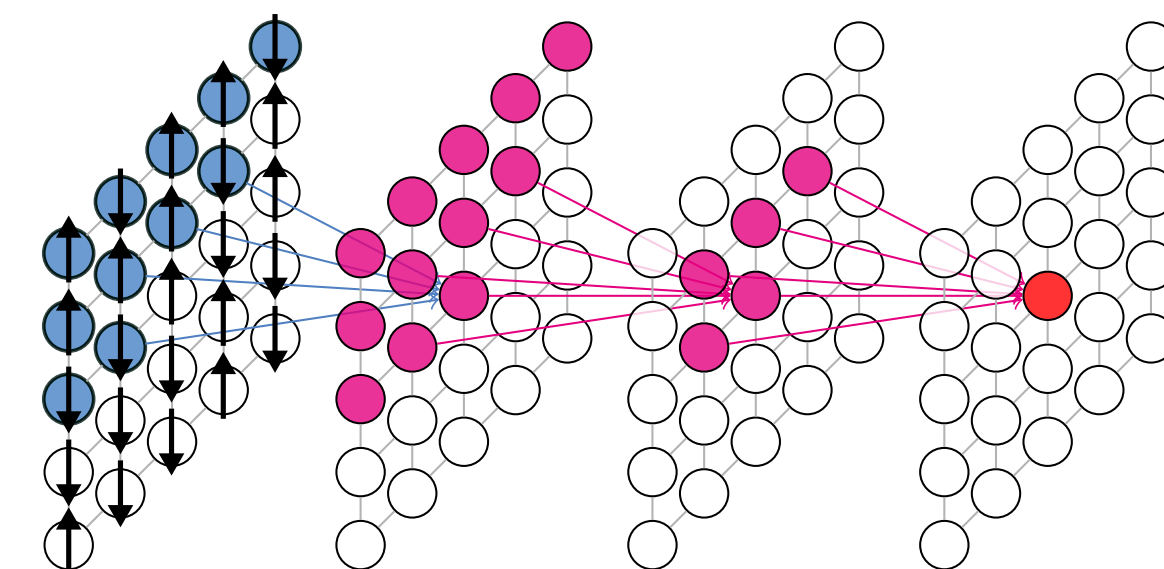
jpeg is a common lossy compression format for digital images
1) compute weights on predefined high-and-low frequency patches
2) throw away high-frequency weights; lossless compress low-frequency weights

# Demo: Generative model of Sycamore data

Quantum chip

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

Transformer



0111110110100
100001111011
100110110111
100110100010
010100011000
010001000000
010101101100
100001111000
100101001001
001000001010

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

https://colab.research.google.com/drive/11War0qULkudKT3h2i5J6r_EmA4wFKkoZ?usp=sharing

# Rydberg GPT

Fitzek et al, 2405.21052

## Prompts
$x$: Hamiltonian parameters

## Answers
$\sigma$: Projective measurements

$$p(\boldsymbol{\sigma}|\boldsymbol{x}) = p(\sigma_1|\boldsymbol{x})p(\sigma_2|\sigma_1,\boldsymbol{x})\ldots$$

"an image of beautiful crystals in 16:9"

pixels ~ $p(\text{pixels} \mid \text{texts})$

GPU hours used for training

$10^8$ — Grok-3

$10^7$ — Llama-3 / DeepeekV3

$10^6$ — GPT3 / Llama-2

$10^5$

$10^4$ — AlphaGo / GPT2 / AlphaFold2 / Stable diffusion

Microsoft MatterGen (Zeni et al, 2312.03687)

Meta Crystal-text-LLM (Gruver, et al, 2402.04379)

$10^3$ — FermiNet

2016  2017  2018  2019  2020  2021  2022  2023  2024  2025

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

more data and compute

more physics and symmetries

# We have much less crystal data

**COMMON CRAWL**

Over 250 billion pages

**ICSD**

> 291,000 crystal structures

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

# Symmetry Preference
# Space groups quantify Nature's preference over symmetry



Inorganics by space group

**Wyckoff Positions of Group _Ia_-3_d_ (No. 230)**

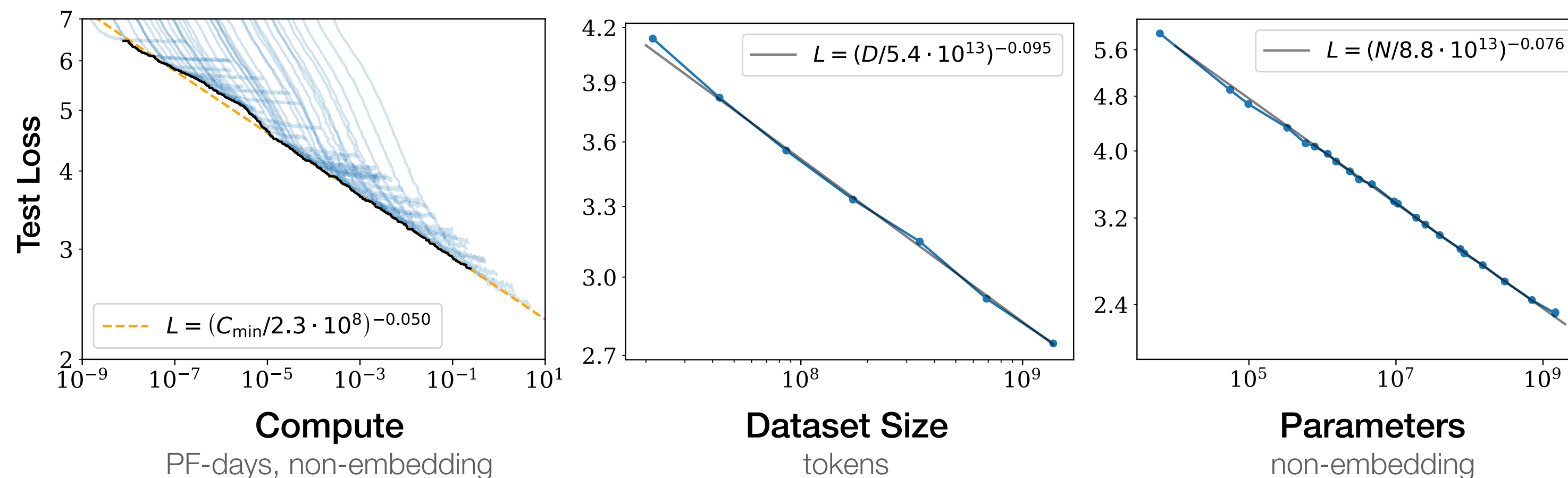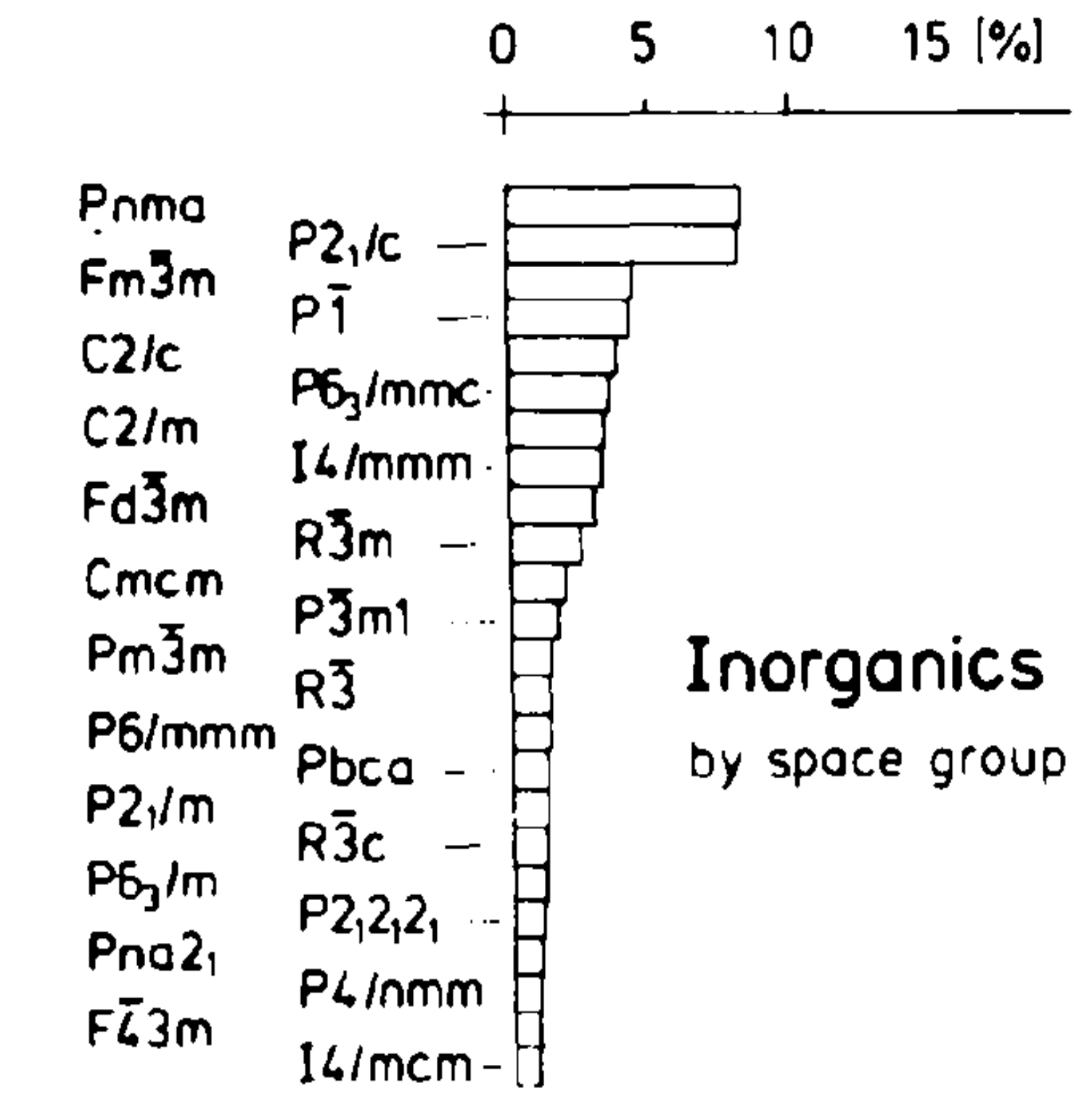| Multiplicity | Wyckoff letter | Site symmetry | Coordinates (0,0,0) + (1/2,1/2,1/2) + | | | |
|---|---|---|---|---|---|---|
| 96 | h | 1 | (x,y,z) | (-x+1/2,-y,z+1/2) | (-x,y+1/2,-z+1/2) | (x+1/2,-y+1/2,-z) |
| | | | (z,x,y) | (z+1/2,-x+1/2,-y) | (-z+1/2,-x,y+1/2) | (-z,x+1/2,-y+1/2) |
| | | | (y,z,x) | (-y,z+1/2,-x+1/2) | (y+1/2,-z+1/2,-x) | (-y+1/2,-z,x+1/2) |
| | | | (y+3/4,x+1/4,-z+1/4) | (-y+3/4,-x+3/4,-z+3/4) | (y+1/4,-x+1/4,z+3/4) | (-y+1/4,x+3/4,z+1/4) |
| | | | (x+3/4,z+1/4,-y+1/4) | (-x+1/4,z+3/4,y+1/4) | (-x+3/4,-z+3/4,-y+3/4) | (x+1/4,-z+1/4,y+3/4) |
| | | | (z+3/4,y+1/4,-x+1/4) | (z+1/4,-y+1/4,x+3/4) | (-z+1/4,y+3/4,x+1/4) | (-z+3/4,-y+3/4,-x+3/4) |
| | | | (-x,-y,-z) | (x+1/2,y,-z+1/2) | (x,-y+1/2,z+1/2) | (-x+1/2,y+1/2,z) |
| | | | (-z,-x,-y) | (-z+1/2,x+1/2,y) | (z+1/2,x,-y+1/2) | (z,-x+1/2,y+1/2) |
| | | | (-y,-z,-x) | (y,-z+1/2,x+1/2) | (-y+1/2,z+1/2,x) | (y+1/2,z,-x+1/2) |
| | | | (-y+1/4,-x+3/4,z+3/4) | (y+1/4,x+1/4,z+1/4) | (-y+3/4,x+3/4,-z+1/4) | (y+3/4,-x+1/4,-z+3/4) |
| | | | (-x+1/4,-z+3/4,y+3/4) | (x+3/4,-z+1/4,-y+3/4) | (x+1/4,z+1/4,y+1/4) | (-x+3/4,z+3/4,-y+1/4) |
| | | | (-z+1/4,-y+3/4,x+3/4) | (-z+3/4,y+3/4,-x+1/4) | (z+3/4,-y+1/4,-x+3/4) | (z+1/4,y+1/4,x+1/4) |
| 48 | g | ..2 | (1/8,y,-y+1/4) (3/8,-y,-y+3/4) (7/8,y+1/2,y+1/4) (5/8,-y+1/2,y+3/4) | | | |
| | | | (-y+1/4,1/8,y) (-y+3/4,3/8,-y) (y+1/4,7/8,y+1/2) (y+3/4,5/8,-y+1/2) | | | |
| | | | (y,-y+1/4,1/8) (-y,-y+3/4,3/8) (y+1/2,y+1/4,7/8) (-y+1/2,y+3/4,5/8) | | | |
| | | | (7/8,-y,y+3/4) (5/8,y,y+1/4) (1/8,-y+1/2,-y+3/4) (3/8,y+1/2,-y+1/4) | | | |
| | | | (y+3/4,7/8,-y) (y+1/4,5/8,y) (-y+3/4,1/8,-y+1/2) (-y+1/4,3/8,y+1/2) | | | |
| | | | (-y,y+3/4,7/8) (y,y+1/4,5/8) (-y+1/2,-y+3/4,1/8) (y+1/2,-y+1/4,3/8) | | | |
| 48 | f | 2.. | (x,0,1/4) (-x+1/2,0,3/4) (1/4,x,0) (3/4,-x+1/2,0) | | | |
| | | | (0,1/4,x) (0,3/4,-x+1/2) (3/4,x+1/4,0) (3/4,-x+3/4,1/2) | | | |
| | | | (x+3/4,1/2,1/4) (-x+1/4,0,1/4) (0,1/4,-x+1/4) (1/2,1/4,x+3/4) | | | |
| | | | (-x,0,3/4) (x+1/2,0,1/4) (3/4,-x,0) (1/4,x+1/2,0) | | | |
| | | | (0,3/4,-x) (0,1/4,x+1/2) (1/4,-x+3/4,0) (1/4,x+1/4,1/2) | | | |
| | | | (-x+1/4,1/2,3/4) (x+3/4,0,3/4) (0,3/4,x+3/4) (1/2,3/4,-x+1/4) | | | |
| 32 | e | .3. | (x,x,x) (-x+1/2,-x,x+1/2) (-x,x+1/2,-x+1/2) (x+1/2,-x+1/2,-x) | | | |
| | | | (x+3/4,x+1/4,-x+1/4) (-x+3/4,-x+3/4,-x+3/4) (x+1/4,-x+1/4,x+3/4) (-x+1/4,x+3/4,x+1/4) | | | |
| | | | (-x,-x,-x) (x+1/2,x,-x+1/2) (x,-x+1/2,x+1/2) (-x+1/2,x+1/2,x) | | | |
| | | | (-x+1/4,-x+3/4,x+3/4) (x+1/4,x+1/4,x+1/4) (-x+3/4,x+3/4,-x+1/4) (x+3/4,-x+1/4,-x+3/4) | | | |
| 24 | d | -4.. | (3/8,0,1/4) (1/8,0,3/4) (1/4,3/8,0) (3/4,1/8,0) | | | |
| | | | (0,1/4,3/8) (0,3/4,1/8) (3/4,5/8,0) (3/4,3/8,1/2) | | | |
| | | | (1/8,1/2,1/4) (7/8,0,1/4) (0,1/4,7/8) (1/2,1/4,1/8) | | | |
| 24 | c | 2.2 2 | (1/8,0,1/4) (3/8,0,3/4) (1/4,1/8,0) (3/4,3/8,0) | | | |
| | | | (0,1/4,1/8) (0,3/4,3/8) (7/8,0,3/4) (5/8,0,1/4) | | | |
| | | | (3/4,7/8,0) (1/4,5/8,0) (0,3/4,7/8) (0,1/4,5/8) | | | |
| 16 | b | .32 | (1/8,1/8,1/8) (3/8,7/8,5/8) (7/8,5/8,3/8) (5/8,3/8,7/8) | | | |
| | | | (7/8,7/8,7/8) (5/8,1/8,3/8) (1/8,3/8,5/8) (3/8,5/8,1/8) | | | |
| 16 | a | .-3. | (0,0,0) (1/2,0,1/2) (0,1/2,1/2) (1/2,1/2,0) | | | |
| | | | (3/4,1/4,1/4) (3/4,3/4,3/4) (1/4,1/4,3/4) (1/4,3/4,1/4) | | | |

**Wyckoff Positions of Group _P_1 (No. 1)**

P1 is rare!

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

...

https://www.cryst.ehu.es/cryst/get_wp.html

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

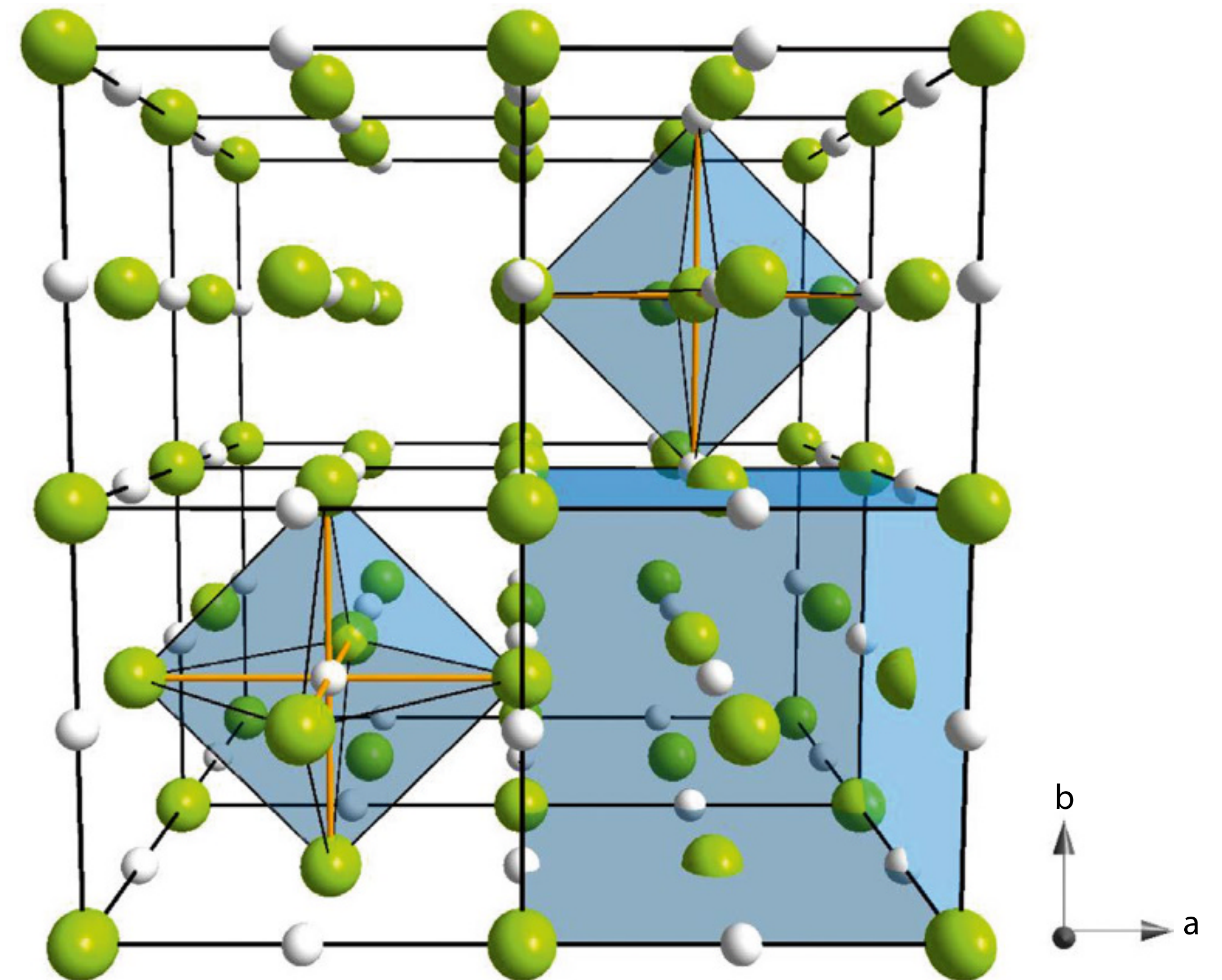| Multiplicity | Wyckoff letter | Site symmetry | Coordinates (0,0,0) + (0,1/2,1/2) + (1/2,0,1/2) + (1/2,1/2,0) + |
|---|---|---|---|
| 192 | l | 1 | (x,y,z) (-x,-y,z) (-x,y,-z) (x,-y,-z)<br>(z,x,y) (z,-x,-y) (-z,-x,y) (-z,x,-y)<br>(y,z,x) (-y,z,-x) (y,-z,-x) (-y,-z,x)<br>(y,x,-z) (-y,-x,-z) (y,-x,z) (-y,x,z)<br>(x,z,-y) (-x,z,y) (-x,-z,-y) (x,-z,y)<br>(z,y,-x) (z,-y,x) (-z,y,x) (-z,-y,-x)<br>(-x,-y,-z) (x,y,-z) (x,-y,z) (-x,y,z)<br>(-z,-x,-y) (-z,x,y) (z,x,-y) (z,-x,y)<br>(-y,-z,-x) (y,-z,x) (-y,z,x) (y,z,-x)<br>(-y,-x,z) (y,x,z) (-y,x,-z) (y,-x,-z)<br>(-x,-z,y) (x,-z,-y) (x,z,y) (-x,z,-y)<br>(-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)<br>(z,x,x) (z,-x,-x) (-z,-x,x) (-z,x,-x)<br>(x,z,x) (-x,z,-x) (x,-z,-x) (-x,-z,x)<br>(x,x,-z) (-x,-x,-z) (x,-x,z) (-x,x,z)<br>(x,z,-x) (-x,z,x) (-x,-z,-x) (x,-z,x)<br>(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)<br>(z,0,y) (z,0,-y) (-z,0,y) (-z,0,-y)<br>(y,z,0) (-y,z,0) (y,-z,0) (-y,-z,0)<br>(y,0,-z) (-y,0,-z) (y,0,z) (-y,0,z)<br>(0,z,-y) (0,z,y) (0,-z,-y) (0,-z,y)<br>(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)<br>(y,1/2,y) (y,1/2,-y) (-y,1/2,y) (-y,1/2,-y)<br>(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)<br>(y,0,y) (y,0,-y) (-y,0,y) (-y,0,-y)<br>(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)<br>(1/4,1/4,x) (3/4,1/4,-x) (1/4,x,3/4) (3/4,-x,3/4)<br>(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)<br>(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)<br>(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)<br>(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*-3*m* (No. 225)

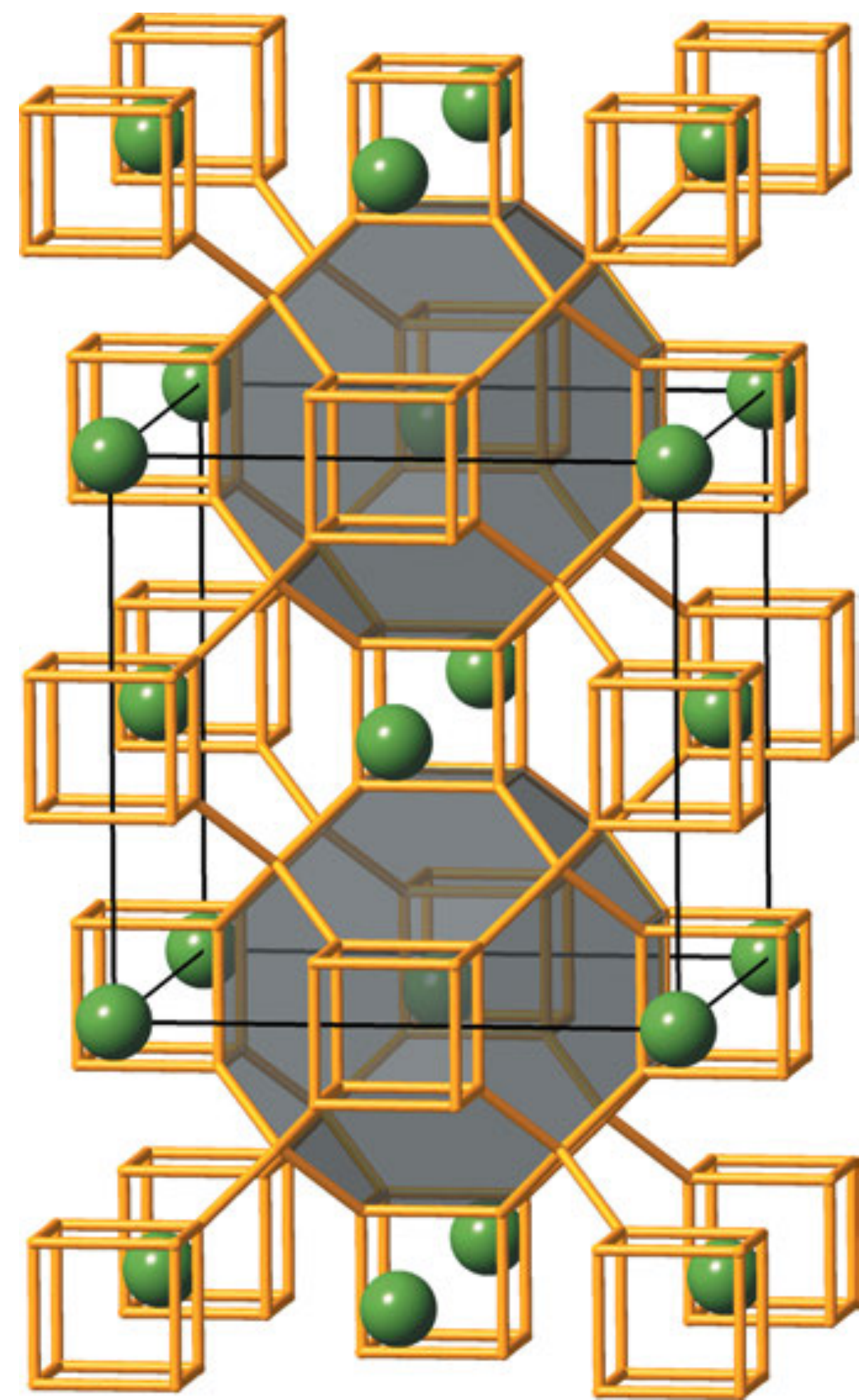| Multiplicity | Wyckoff letter | Site symmetry | Coordinates (0,0,0) + (0,1/2,1/2) + (1/2,0,1/2) + (1/2,1/2,0) + |
|---|---|---|---|
| 192 | l | 1 | (x,y,z) (-x,-y,z) (-x,y,-z) (x,-y,-z)<br>(z,x,y) (z,-x,-y) (-z,-x,y) (-z,x,-y)<br>(y,z,x) (-y,z,-x) (y,-z,-x) (-y,-z,x)<br>(y,x,-z) (-y,-x,-z) (y,-x,z) (-y,x,z)<br>(x,z,-y) (-x,z,y) (-x,-z,-y) (x,-z,y)<br>(z,y,-x) (z,-y,x) (-z,y,x) (-z,-y,-x)<br>(-x,-y,-z) (x,y,-z) (x,-y,z) (-x,y,z)<br>(-z,-x,-y) (-z,x,y) (z,x,-y) (z,-x,y)<br>(-y,-z,-x) (y,-z,x) (-y,z,x) (y,z,-x)<br>(-y,-x,z) (y,x,z) (-y,x,-z) (y,-x,-z)<br>(-x,-z,y) (x,-z,-y) (x,z,y) (-x,z,-y)<br>(-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)<br>(z,x,x) (z,-x,-x) (-z,-x,x) (-z,x,-x)<br>(x,z,x) (-x,z,-x) (x,-z,-x) (-x,-z,x)<br>(x,x,-z) (-x,-x,-z) (x,-x,z) (-x,x,z)<br>(x,z,-x) (-x,z,x) (-x,-z,-x) (x,-z,x)<br>(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)<br>(z,0,y) (z,0,-y) (-z,0,y) (-z,0,-y)<br>(y,z,0) (-y,z,0) (y,-z,0) (-y,-z,0)<br>(y,0,-z) (-y,0,-z) (y,0,z) (-y,0,z)<br>(0,z,-y) (0,z,y) (0,-z,-y) (0,-z,y)<br>(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)<br>(y,1/2,y) (y,1/2,-y) (-y,1/2,y) (-y,1/2,-y)<br>(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)<br>(y,0,y) (y,0,-y) (-y,0,y) (-y,0,-y)<br>(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)<br>(1/4,1/4,x) (3/4,1/4,-x) (1/4,x,3/4) (3/4,-x,3/4)<br>(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)<br>(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)<br>(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)<br>(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) |

3  )

NaCl

# 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)<br>(z,x,y) (z,-x,-y) (-z,-x,y) (-z,x,-y)<br>(y,z,x) (-y,z,-x) (y,-z,-x) (-y,-z,x)<br>(y,x,-z) (-y,-x,-z) (y,-x,z) (-y,x,z)<br>(x,z,-y) (-x,z,y) (-x,-z,-y) (x,-z,y)<br>(z,y,-x) (z,-y,x) (-z,y,x) (-z,-y,-x)<br>(-x,-y,-z) (x,y,-z) (x,-y,z) (-x,y,z)<br>(-z,-x,-y) (-z,x,y) (z,x,-y) (z,-x,y)<br>(-y,-z,-x) (y,-z,x) (-y,z,x) (y,z,-x)<br>(-y,-x,z) (y,x,z) (-y,x,-z) (y,-x,-z)<br>(-x,-z,y) (x,-z,-y) (x,z,y) (-x,z,-y)<br>(-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)<br>(z,x,x) (z,-x,-x) (-z,-x,x) (-z,x,-x)<br>(x,z,x) (-x,z,-x) (x,-z,-x) (-x,-z,x)<br>(x,x,-z) (-x,-x,-z) (x,-x,z) (-x,x,z)<br>(x,z,-x) (-x,z,x) (-x,-z,-x) (x,-z,x)<br>(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)<br>(z,0,y) (z,0,-y) (-z,0,y) (-z,0,-y)<br>(y,z,0) (-y,z,0) (y,-z,0) (-y,-z,0)<br>(y,0,-z) (-y,0,-z) (y,0,z) (-y,0,z)<br>(0,z,-y) (0,z,y) (0,-z,-y) (0,-z,y)<br>(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)<br>(y,1/2,y) (y,1/2,-y) (-y,1/2,y) (-y,1/2,-y)<br>(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)<br>(y,0,y) (y,0,-y) (-y,0,y) (-y,0,-y)<br>(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)<br>(1/4,1/4,x) (3/4,1/4,-x) (1/4,x,3/4) (3/4,-x,3/4)<br>(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)<br>(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)<br>(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)<br>(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) |



LaH$_{10}$

# CrystalFormer

Zhendong Cao, Xiaoshan Luo, Jian Lv, and LW, 2403.15734

deepmodeling/CrystalFormer

## Space Group Informed Transformer for Crystals

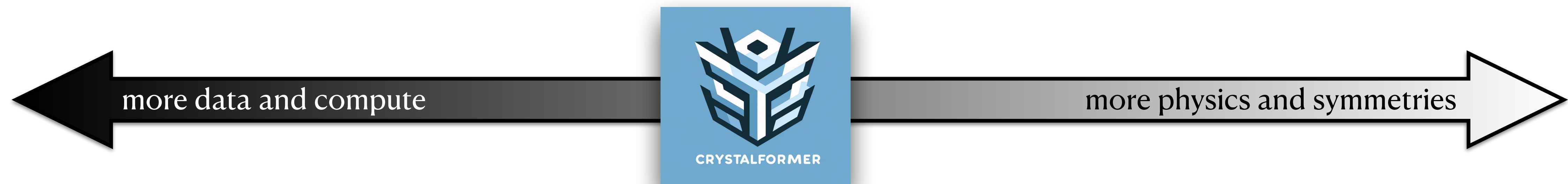225-a-La-0-0-0-c-H-1/4-1/4-1/4-f-H-0.375-0.375-0.375-X-5.1-5.1-5.1-90-90-90

"Grammar" ~ Solid state chemistry

"Synonyms" ~ Elememt substitution
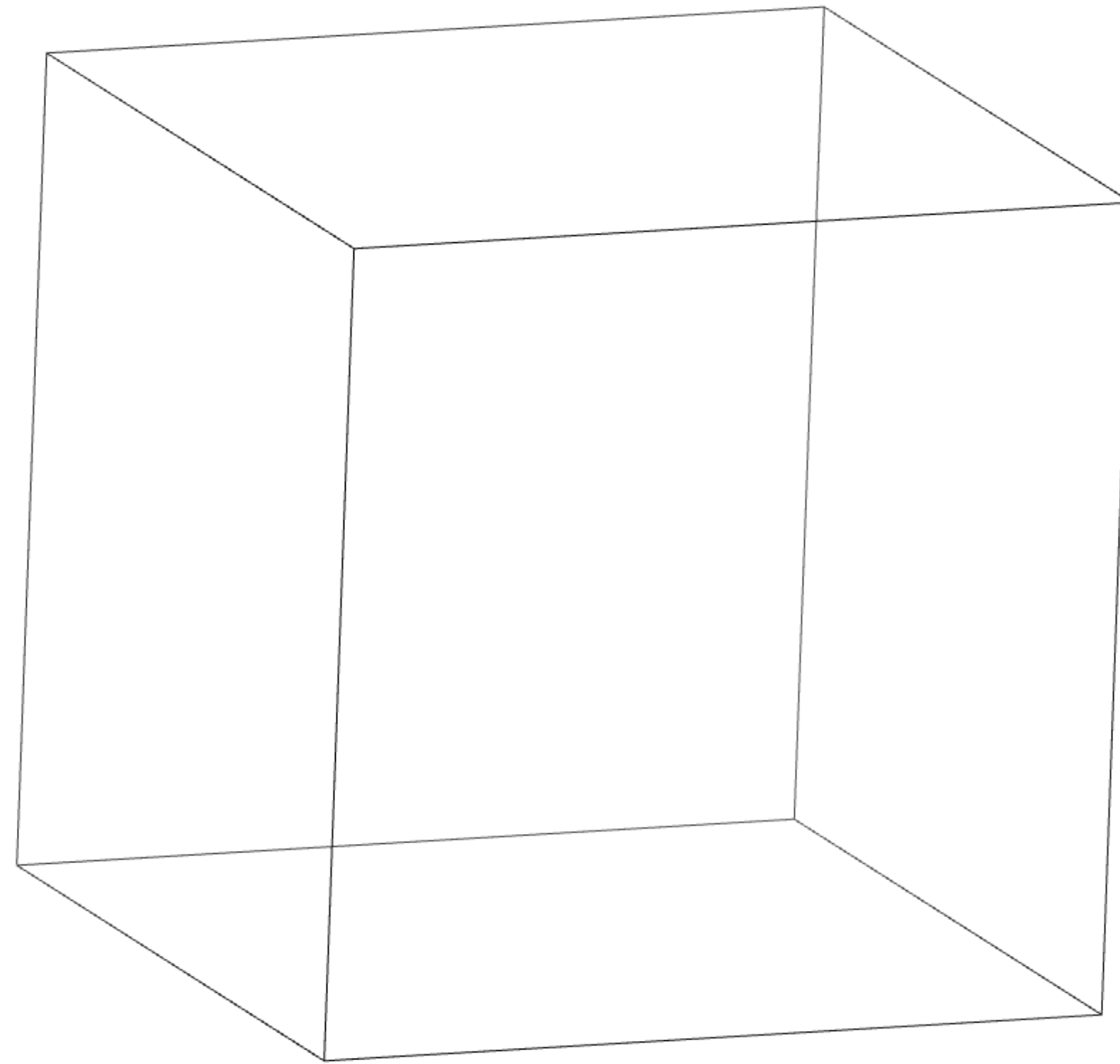
"Idioms" ~ Coordination polyhedra

"Rhythm" ~ Wyckoff positions

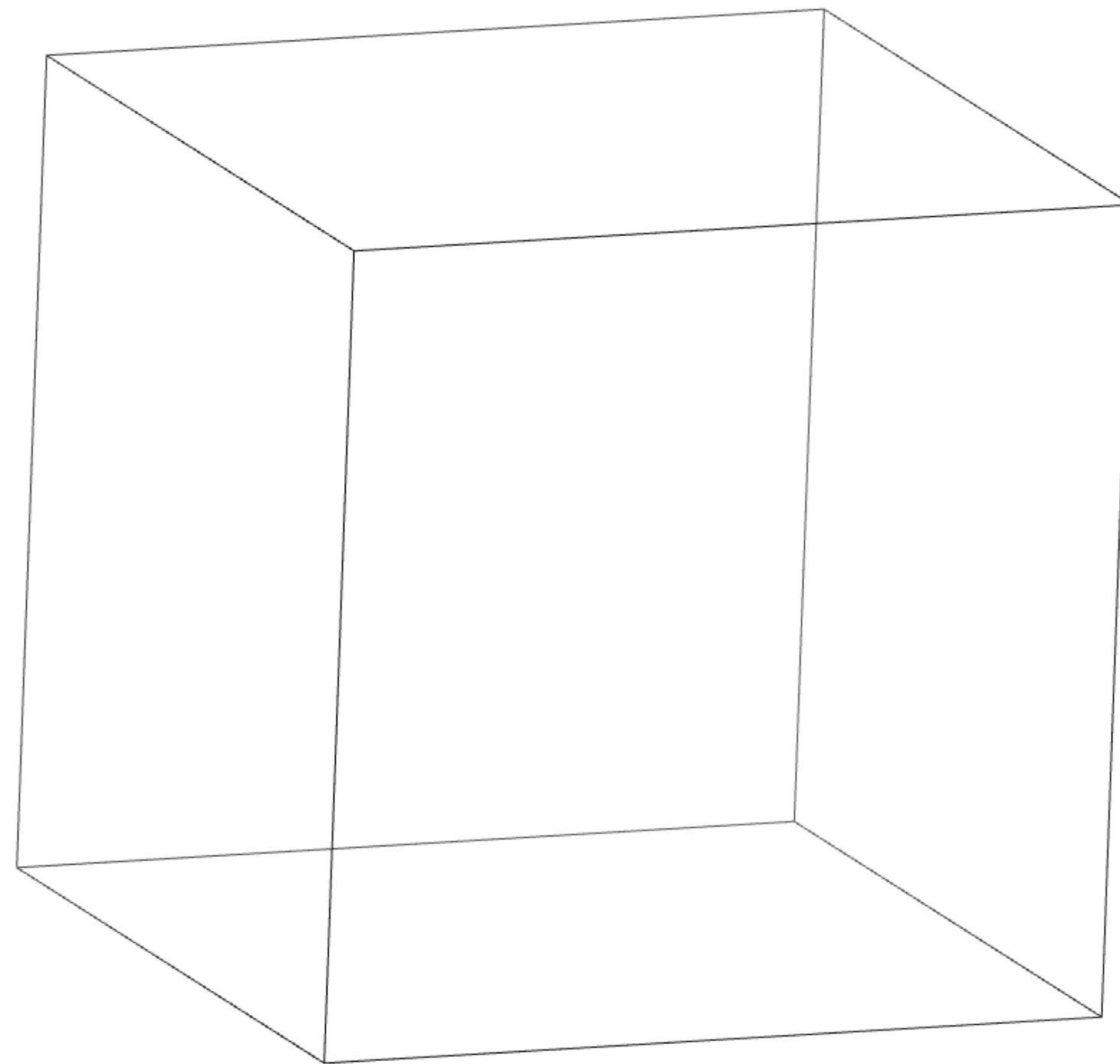**Nature's codebook for tokenization discrete, pre-compression**

more data and compute ⟷ more physics and symmetries

CRYSTALFORMER

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

# Autoregressive sampling of a crystal

$Cs_2ZnFe(CN)_6$

# Autoregressive sampling of a crystal



$Cs_2ZnFe(CN)_6$

225-a-Fe-0-0-0-b-Zn-1/2-1/2-1/2-c-Cs-1/4-1/4-1/4-e-C-0.18-0-0-e-N-0.29-0-0-X-10.45-10.45-10.45-90-90-90

# Aside: autoregressive transformer for images

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



**Codebook** $\mathcal{Z}$

**Transformer**

$$p(s) = \prod_i p(s_i | s_{<i})$$

$$s_{<i} \qquad s_i$$

learned codebook, see also Tian et al, 2404.02905

CrystalFormer leverages Nature's codebook: the Wyckoff position table

# Bayes rule for materials inverse design



$p(X)$

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

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

How to sample from $p(X \,|\, y)$? Two approaches originated in physics

**Markov chain Monte Carlo**
Metropolis et al, 1953, Hastings 1970

$\nabla F$

**Variational inference**
Gibbs, Feynman, Bogoliubov,..., Jordan et al 1999

# MCMC sampling from the posterior

Generate more double perovskites $A_2BB'O_6$



225-a-[?]-0-0-0-b-[?]-1/2-1/2-1/2-c-[?]-1/4-1/4-1/4-e-O-[?]-0-0

Solve crystal cloze test via MCMC sweep through the "crystal string"

$$A(X \to X') = \min \left[ 1, \frac{p(X')}{p(X)} \right]$$

opstring[p]

4
0
9
12
6
0
0
4
12
0
9
14

MCMC sweep through the "operator string" in Sandvik's SSE algorithm

# Aside: Constrained sentence generation in languate modeling

$$\pi(x) \propto P_{\text{LM}}(x) \cdot \text{Constraint}(x)$$

① Paris is located in France.       : Deletion
② Paris is located in France.
③ Paris located in France.
④ Is Paris located in France?

"traverses the probabilistic space of high-quality sentences more effectively"

Miao et al, 1811.10996, Zhang et al, 2011.12334



The sequence length for inorganic crystals is ~100 with vocabulary size ~100

So, even naive Metropolis-Hastings with annealing works fine

# Variational inference the posterior

$$\mathbb{KL}\left(q(X) \parallel p(X|y)\right) = \underset{X \sim q(X)}{\mathbb{E}}\left[-\ln p(y|X)\right] + \mathbb{KL}\left(q(X)\|p(X)\right)$$

↑ Variational probability

↑ Likelihood function

↑ Prior

$q(X)$ is easy to sample, e.g. another autoregressive model

Variational inference turns a sampling problem into a stochastic optimization problem

# Also known as: reinforcement fine-tuning

$$\mathbb{KL}\left(q(X) \parallel p(X|y)\right) = \underset{X \sim q(X)}{\mathbb{E}}\left[-r(X)\right] + \mathbb{KL}\left(q(X)\|p(X)\right)$$

↑        ↑        ↑

Fine-tuned    Reward       Remain close to
model        function     the pretrained model

"RL with KL penalties is better viewed as Bayesian inference" Korbar et al, 2205.11275

# Two sides of the same coin

### ① Pre-training



**learn from data
to be a generalist**

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

### ② Post-training



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

$$\mathscr{L} = \mathop{\mathbb{E}}_{X\sim q(X)} \big[ -r(X) \big] + \mathbb{KL}\big( q(X) \| p(X) \big)$$

$$\mathbb{KL}(\text{data} \| p) \quad \text{vs} \quad \mathbb{KL}(q \| p e^r)$$

# Reinforcement fine-tuning fo ... sign



(a)

Crys

$\mathbb{E}$

$X \sim q(X$

Reward = Band gap × d

(Two usually anti-

$\varepsilon_{total} = 21.32$

$E_g = 5.38$eV

$E_{hull} = 0.04$eV/atom

225-a-Sr-b-Ba-c-Cs-e-F

$Cs_2BaSrF_6$

225-a-Pb-c-Cl-d-Li-e-Cl

$Li_6PbCl_8$

$\varepsilon_{total} = 25.19$

$E_g = 4.64$ eV

$E_{hull} = 0.08$eV/atom

# 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

# Variational autoregressive network for statistical mechanics



Sherrington-Kirkpatrick spin glass

Objective function: variational free-energy

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

Naive mean-field factorized probability

$$p(X) = \prod_i p(x_i)$$

Bethe approximation pairwise interaction

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

Variational autoregressive network

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

Wu, LW, Zhang, PRL '19

# VAN for triangular Ising

Wu, LW, Zhang, PRL '19

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



(b)

Residual entropy

S/N = 0.323 Wannier 1950

Legend: L = 4, L = 6, L = 8, L = 10, L = 12, L = 14, L = 16

$\beta = 1/k_B T$

Hot configuration

Cold configuration
MacKay, 2006

# VAN (aka RL) for 8-queens problem

$$\mathscr{L} = \mathop{\mathbb{E}}_{X \sim q(X)} \left[ -r(X) + \ln q(X) \right]$$

Energy
exploitation

Entropy
exploration

Reward $r(X) = \begin{cases} 1 & \text{if no attack} \\ 0 & \text{otherwise} \end{cases}$

Policy network
$q(X) = q(x_1)q(x_2 | x_1)\ldots$

$X$: a sequence of actions
a1—b7—c4—d6—e8—f2—g5—h3

# VAN (aka RL) for 8-queens problem

$$\mathscr{L} = \mathbb{E}_{X \sim q(X)} \left[ -r(X) + \ln q(X) \right]$$

Energy exploitation

Entropy exploration

| Board size | Solutions |
|:---:|:---:|
| 8 | 92 |
| 12 | 14,200 |
| 16 | 14,772,512 |
| 20 | 39,029,188,884 |
| 24 | 227,514,171,973,736 |
| 28 | ??? |

Can you solve it ?

# Variational autoregressive quantum states

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

## Objective function: ground state energy

McMillan 1965, Carleo & Troyer Science 2017

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

### Heisenberg and Hubbard models

Sharir et al, PRL '20, Hibat-Allah et al, PRResarch '20

Humeniuk et al, SciPost '23

Ibarra-García-Padilla et al, 2411.07144   Moss et al, 2502.17144

### Quantum chemistry problems

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

Zhao et al, MLST. '23   Shang et al, 2307.09343

Malyshev et al, 2310.04166   Malyshev et al, 2408.07625

# Deep learning for variational calculations

Turning physics problems into stochastic optimization

Leverages the deep learning engine



RBM, FermiNet,...

**Representation**

Carleo and Troyer, Science '17

Pfau et al, PR Research '20

Automatic differentiation,
Wasserstein gradient, KFAC, ...

**Optimization**

Liao et al, PRX '19

Neklyudov, et al, 2307.07050

Chen et al, Nat.Phys. '24

**Sampling**

Wu et al, PRL '19

Humeniuk et al, SciPost '23

Malyshev et al, 2408.07625

Autoregressive sampling,

Gumbel-top-k sampling,...

# Two kinds of variational Monte Carlo

**Variational ground state energy $T = 0$**

McMillan 1965, Carleo & Troyer Science 2017, Pfau et al, FermiNet, ...

$$E[\Psi] = \mathop{\mathbb{E}}_{X \sim |\psi(X)|^2} \left[ \frac{\hat{H}\Psi(X)}{\Psi(X)} \right]$$

$\Psi$: ANY neural network that respects physical symmetries

**Variational free energy $T > 0$**

Gibbs–Bogolyubov-Feynman, Li and LW, PRL '18, Wu, LW, Zhang, PRL '19, ...

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

$p$: probabilistic models with tractable normalization

# Three kinds of variational Monte Carlo

| Quantum Ground state | Quantum Stat-Mech | Classical Stat-Mech |

McMillan 1965
Carleo & Troyer Science 2017, Pfau et al, FermiNet, ...

Gibbs–Bogolyubov-Feynman
Li and LW, PRL '18, Wu, LW, Zhang, PRL '19, ...

$$E[\Psi] = \mathop{\mathbb{E}}_{X \sim |\psi(X)|^2} \left[ \frac{\hat{H}\Psi(X)}{\Psi(X)} \right]$$

$$\rho$$

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

$\Psi$: ANY neural network that respects physical symmetries

$p$: probabilistic models with tractable normalization

# The variational free energy principle

Gibbs–Bogolyubov-Feynman-Delbrück–Molière

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

↓         ↓              ↓

variational density matrix    energy         entropy 😱

**Difficulties in Applying the Variational Principle to Quantum Field Theories**[1]

Richard P. Feynman

$\rho$ ?

Generative models !

# Example: the variational density matrix of electron gas

Fermi sea

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

$$\boldsymbol{K} = \{\boldsymbol{k}_1, \boldsymbol{k}_2, \ldots, \boldsymbol{k}_N\}$$

$$\rho = \sum_{\boldsymbol{K}} p(\boldsymbol{K}) \left| \Psi_{\boldsymbol{K}} \right\rangle \left\langle \Psi_{\boldsymbol{K}} \right|$$

Normalized probability
distribution

Orthonormal
many-electron basis

❶ $\sum_{\boldsymbol{K}} p(\boldsymbol{K}) = 1$

❷ $\langle \Psi_{\boldsymbol{K}} | \Psi_{\boldsymbol{K}'} \rangle = \delta_{\boldsymbol{K},\boldsymbol{K}'}$

autoregressive model

$\sqrt{\text{flow}}$

There will also be interesting twists for physics considerations

# ① Variational autoregressive network for $p(\boldsymbol{K})$

Fermionic
occupation
in k-space

$$p(\boldsymbol{K}) = p(\boldsymbol{k}_1)p(\boldsymbol{k}_2\,|\,\boldsymbol{k}_1)p(\boldsymbol{k}_3\,|\,\boldsymbol{k}_1,\boldsymbol{k}_2)\cdots$$

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



|   |   |   |
|---|---|---|
| N | # of fermions | # of words |
| M | Momentum cutoff | Vocabulary |

*quick*

*brown*   *fox*

*jumps*

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

We use masked casual self-attention Vaswani et al 1706.03762; Alternative solution: Hibat-Allah et al, 2002.02793, Barrett et al, 2109.12606

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

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

Xie, Zhang, LW, SciPost '23

Orthonormal many-body states



real particle

quasi particle

Electron coordinates

Quasi-particle coordinates

$X$    Equivariant neural network   +   Equivariant neural network   +   Equivariant neural network   +   $Z$

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

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

# The objective function of variational density matrix

$$\rho = \sum_{K} p(K) \left| \Psi_K \right\rangle \left\langle \Psi_K \right|$$

$$F = \mathop{\mathbb{E}}_{K \sim p(K)} \left[ k_B T \ln p(K) + \mathop{\mathbb{E}}_{X \sim \left| \Psi_K(X) \right|^2} \left[ \frac{H \Psi_K(X)}{\Psi_K(X)} \right] \right]$$

Boltzmann distribution

Born probability

Jointly optimize $p(K)$ and $\Psi_K(X)$ to minimize the variational free energy

# Benchmarks on uniform electron gas

Xie, Zhang, LW, SciPost Physics '23

$r_s = 10$, $T/T_F = 0.0625$, $N = 33$



metals: $2 < r_s < 6$

| $r_s$ | $\Theta$ | $\langle sign \rangle$ | $E_{tot}^{exact}$ | $E_{tot}$ |
|---|---|---|---|---|
| 4.0 | 0.0625 | $-0.00055(62)$ | $-0.5(1)$ | $-0.1023(7)$ |
| 10.0 | 0.0625 | $-0.002(1)$ | $-0.16(2)$ | $-0.1010(1)$ |

Brown et al, PRL '13 Restricted PIMC

see also Schoof et al PRL '15, Malone et al PRL '16

# Application to 2DEG effective mass


quasi horse / real horse

Quansi-particle effective mass
contradicting experiments

Hao Xie et al, SciPost Physics '23
Thermal entropy of 2D electron gas

$$\frac{m^*}{m} = \frac{s}{s_0} < 1$$

$$m*/m > 1$$

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

$$m*/m < 1$$

**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)

# Deep variational free-energy for electrons and atoms



**Ideal Fermi gas**

$\updownarrow$

**Fermi liquid**

Low-temperature properties
of Coulomb gas
JML '22 and SciPost Physics '23
**(~50 electrons)**

**Hartree-Fock states**

$\updownarrow$

**Interacting electrons**

Equation of states of
dense hydrogen
PRL '23 and ongoing
**(~50 e-p pairs)**

Poster by Zihang Li

**Harmonic oscillators**

$\updownarrow$

**Anharmonic crystal**

Vibrational spectra of
molecules and quantum solids
JCP '24 and 2412.12451
**(~500 atoms)**

Poster by Qi Zhang

# 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

# Autoregressive modeling

① Ordering     ② Tokenization     ③ Objective function     ④ Inference



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

# ⤶ Comparision/Connection to MPS
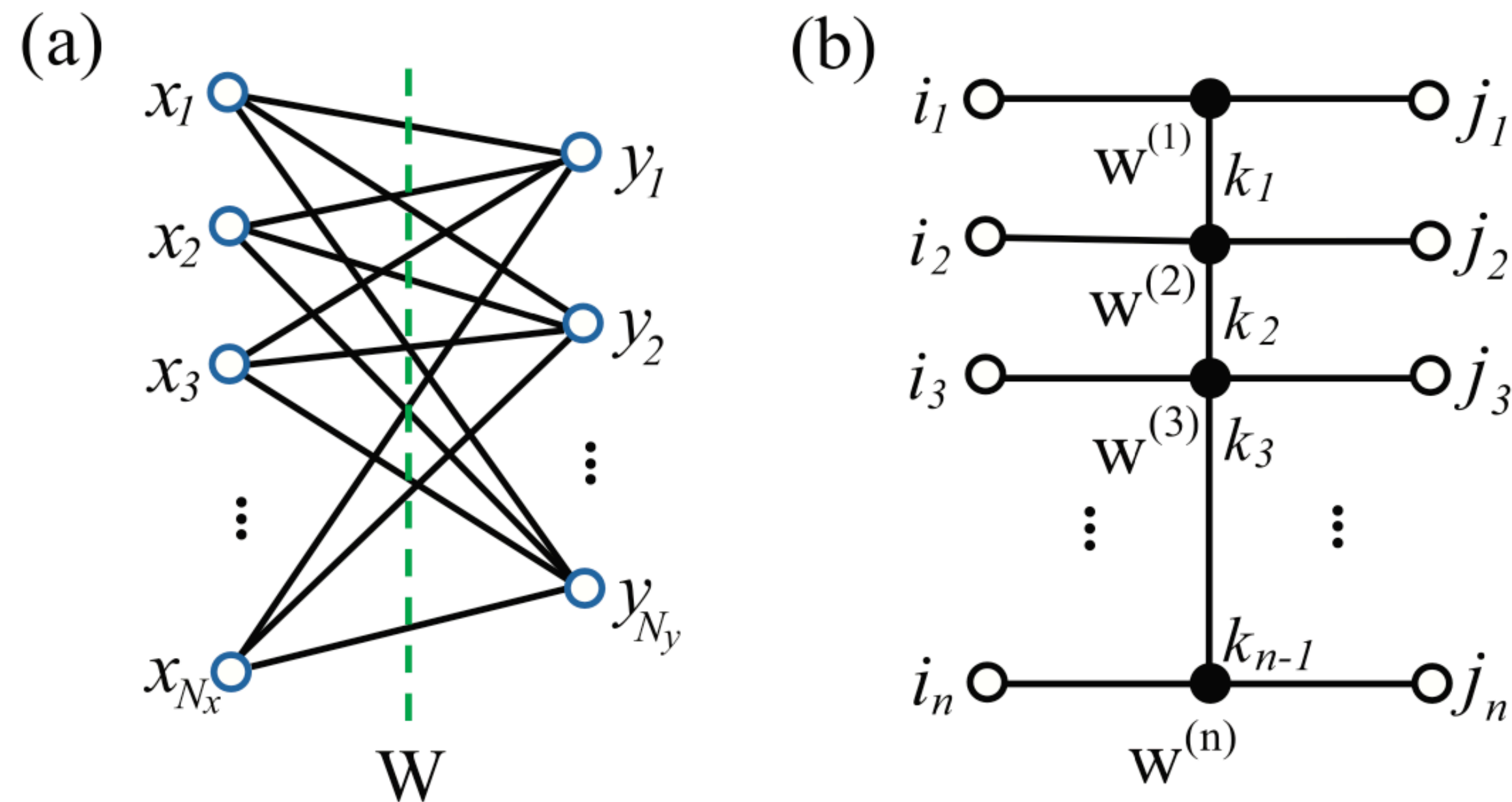


Han et al, PRX '18

Both require 1D ordering

Both support direct sampling
and tractable normaliztion

MPS is bidirectional

MPS mediates long range
correlation via virtual bonds

Similar to recurrent neural net

Can we tensorize GPT ? What does it good for ?
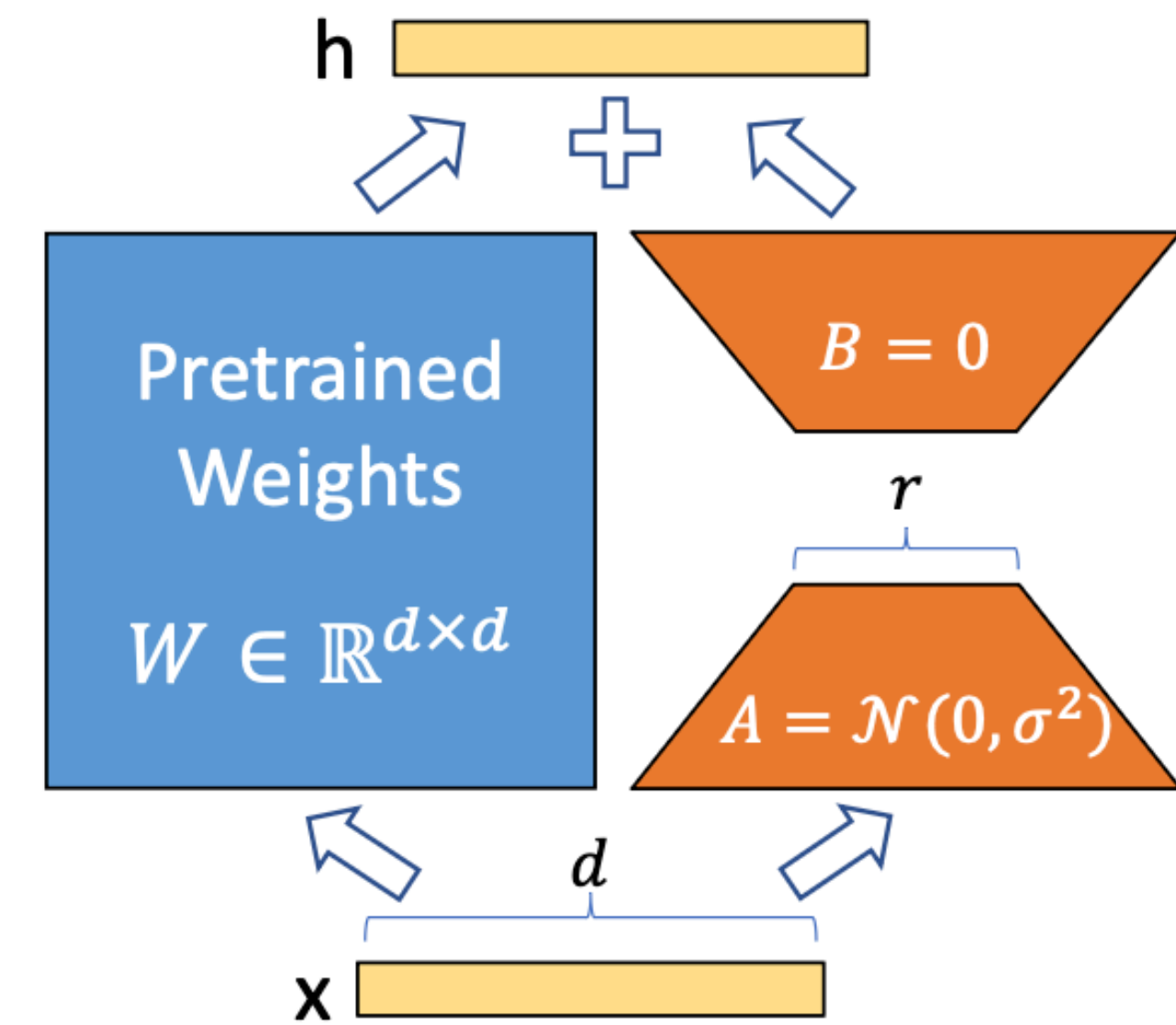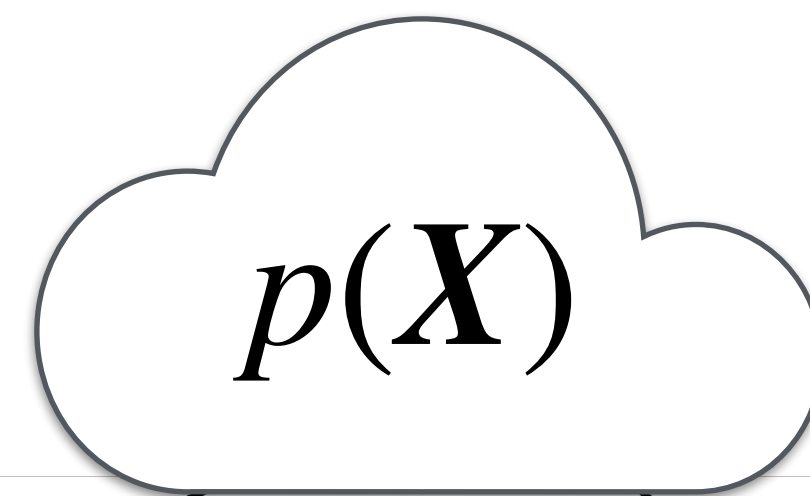
# Tensor network-based compression and finetuning



(a)

$x_1$, $x_2$, $x_3$, ..., $x_{N_x}$ — $y_1$, $y_2$, ..., $y_{N_y}$

W

(b)

$i_1$ — $W^{(1)}$ $k_1$ — $j_1$
$i_2$ — $W^{(2)}$ $k_2$ — $j_2$
$i_3$ — $W^{(3)}$ $k_3$ — $j_3$
...
$i_n$ — $k_{n-1}$ — $j_n$
$W^{(n)}$

Gao et al, PRResearch '20



h

Pretrained Weights
$W \in \mathbb{R}^{d \times d}$

$B = 0$
$r$
$A = \mathcal{N}(0, \sigma^2)$

$d$

x

Figure 1: Our reparametrization. We only train $A$ and $B$.

Low-Rank Adaptation (LoRA)

Hu et al, 2106.09685

# Generative models and their physics genes

**Goodfellow,**
**NIPS tutorial, 1701.00160**

$p(X)$

Direct
GAN

Explicit density

Implicit density

$\Psi$

**Tensor Networks**

Han et al, PRX '18

Tractable density

-Fully visible belief nets
-NADE
**Autoregressive model**
-MADE
-Pixel
-Change of variables models (nonlinear ICA)

**Flow model**

Approximate density

Markov Chain

GSN

$U$

**Quantum Circuits**

Liu et al PRA '18

Variational

Markov Chain

Variational autoencoder    Boltzmann machine

$+$**Diffusion models**

A crash course offered at IOP 2023 spring

# Machine learning for physicists

https://github.com/wangleiphy/ml4p