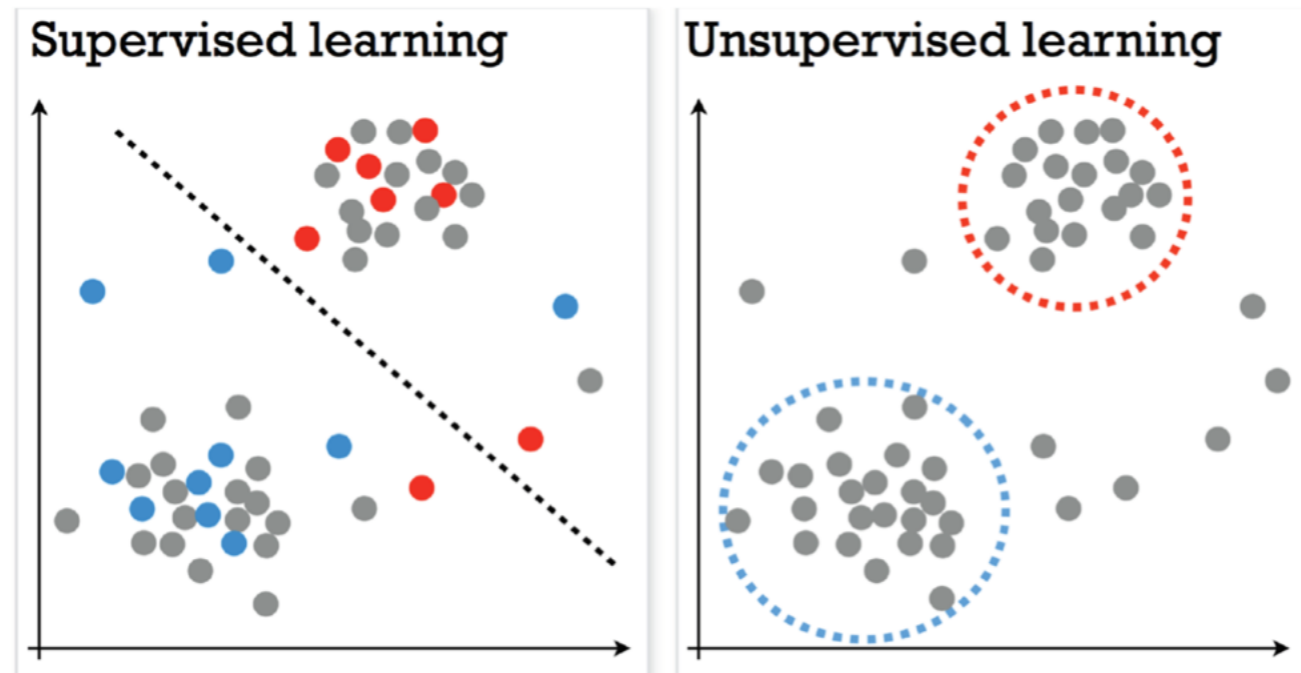


Quantum-enhanced unsupervised (generative) learning
(with near-term devices)

Machine Learning: the **WHAT**



or



?

- generative models*
- clustering (discriminative)*
- feature extraction*

Learning $P(\text{labels}|\text{data})$ given samples from $P(\text{data}, \text{labels})$
(also regression)

Learning structure in $P(\text{data})$
give samples from $P(\text{data})$

Generative models

“What I cannot create, I do not understand”

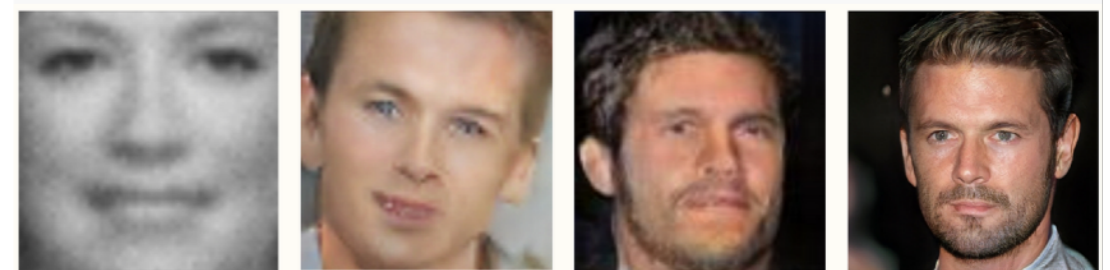
R. Feynman

$$\mathbf{x} \in S \subseteq \mathbb{R}^n$$

$$P(\mathbf{x})$$

$$D = \{\mathbf{x}_i\} \sim P^{\times |D|}$$

(algorithmically) generate new samples \mathbf{x}
(approximately) distributed according to P



2014

2017

[arXiv:1710.10196](https://arxiv.org/abs/1710.10196)

[arXiv:1802.07228](https://arxiv.org/abs/1802.07228)

Name one thing in this photo



Why care about generative models

- generating new data; e.g. medicinal new drugs
- completing missing data; image recovery

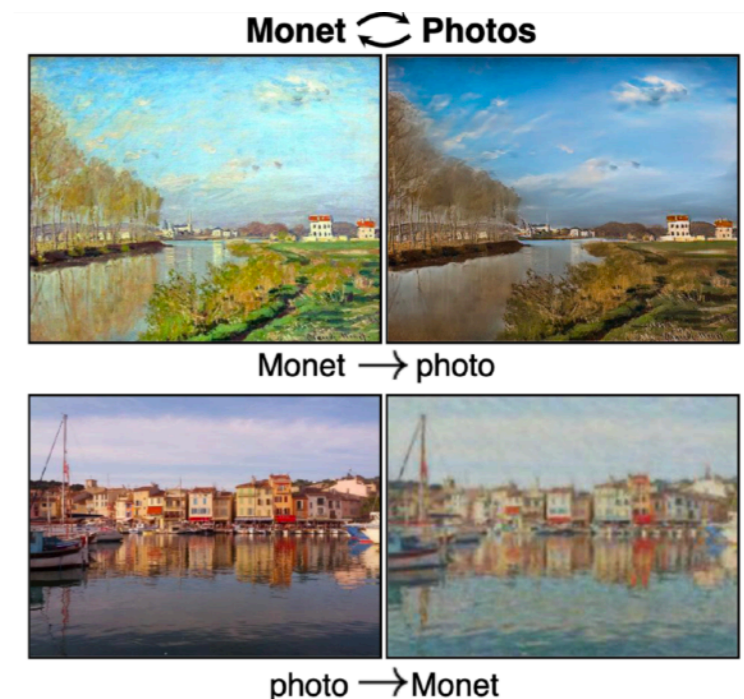
$$\mathcal{P}_{\theta}(x_1 \dots x_M) \approx \mathcal{P}(\vec{x}) \quad x = (x_1 \dots x_M)$$

↖ EQ. PIXEL VALUES

CORRUPT : $(\bar{x}_1 \dots \bar{x}_k, ??? \dots ?)$

RECOVERY : $\mathcal{P}_{\theta}(x_{k+1} \dots x_M \mid \bar{x}_1, \bar{x}_2 \dots \bar{x}_k)$

- modality translation
- ...modelling reality...predict the future?



Full generative problems strictly more general than supervised learning

Let $Z = X \times Y$

SL: Given $D \sim P(X, Y)$, approximate $P(Y|X)$

UL: Given $D \sim P(X, Y)$, approximate $P(X, Y)$

Given access to $P(X, Y)$, the conditional can be derived....

(Also: distributions generalize functions)

$f: X \rightarrow Y \longrightarrow P(Y|X)$ (is FUNCTIONAL if $P(Y|X)$ is DIRAC-DELTA $\forall x$)

$$P(Y|X) \xrightarrow{+ P(X)} P(x, y) = P(Y|X) \cdot P(X)$$

$$\hookrightarrow P(X) = \sum_y P(x, y)$$

Full generative problems strictly more general than supervised learning

Let $Z = X \times Y$

SL: Given $D \sim P(X, Y)$, approximate $P(Y | X)$

UL: Given $D \sim P(X, Y)$, approximate $P(X, Y)$

Given access to $P(X, Y)$, the conditional can be derived....

In generative models we find parameters that *explain* (“cause”) *all of the data*

In discriminative models we find parameters that explain *only the (for SL) relevant aspects* of data

How do we do it?

$$\mathbf{x} \in S \subseteq \mathbb{R}^n$$

$$P(\mathbf{x})$$

(algorithmically) generate new samples \mathbf{x}
(approximately) distributed according to P

$$D = \{\mathbf{x}_i\} \sim P^{\times |D|}$$

How (1): parametrized distribution family $\{d^\theta(\mathbf{x})\}$

Find θ such that $d^\theta(\mathbf{x}) \approx P_{data}(\mathbf{x})$

Hard because we need to efficiently:

- represent
- sample from
- learn

complex and high dimensional probability distribution

How (2): the training and the metrics

- have parametrized distribution family $\{d^\theta(\mathbf{x})\}$
- How do we find θ such that $d^\theta(\mathbf{x}) \approx P_{data}(\mathbf{x})$ $D = \{x_i\}$

Rather: *what does “ \approx ” mean* given that we have samples only!

1) Maximum likelihood estimation (MLE):

$$\text{Likelihood } L(\theta, D) = \prod_{x \in D} d^\theta(x)$$

$$\text{Find } \theta_{opt} = \operatorname{argmin}_\theta \{-\log L(\theta, D)\}$$

*optimizing negative log-likelihood is equivalent to optimizing the so called Kullback–Leibler divergence (a.k.a. relative entropy, a measure of distinction between distributions) between the source and target distribution

How (2): the training and the metrics

- have parametrized distribution family $\{d^\theta(\mathbf{x})\}$
 - How do we find θ such that $d^\theta(\mathbf{x}) \approx P_{data}(\mathbf{x})$ $D = \{x_i\}$

Rather: what does “ \approx ” mean given that we have samples only!

1) Maximum likelihood estimation (MLE):

$$\text{Likelihood } L(\theta, D) = \prod_{x \in D} d^\theta(x)$$

$$\text{Find } \theta_{opt} = \operatorname{argmin}_\theta \{-\log L(\theta, D)\}$$

Updates can go via (approximate) gradient descent...
or derivative-free methods

How (2): the training and the metrics

- have parametrized distribution family $\{d^\theta(\mathbf{x})\}$
- How do we find θ such that $d^\theta(\mathbf{x}) \approx P_{data}(\mathbf{x})$

$$D = \{x_i\}$$

Other options (2)

2) Maximum a-posteriori (MAP):

$$P(\theta | D) = \frac{P(D | \theta)P(\theta)}{P(D)}$$

$$\theta_{MAP} = \operatorname{argmax}_\theta P(D | \theta)P(\theta)$$

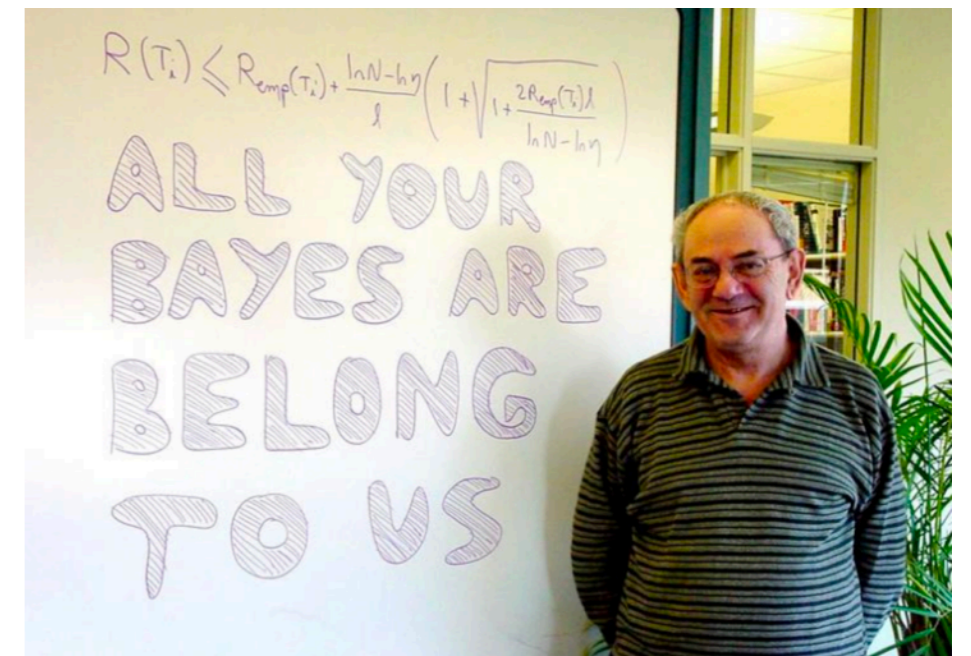
$$= \operatorname{argmax}_\theta \log P(D | \theta)P(\theta)$$

$$= \operatorname{argmax}_\theta \underbrace{\log P(D | \theta)}_{MLE} + \underbrace{\log P(\theta)}$$

MLE

=> MAP IS MLE

UNDER UNIFORM PRIOR



Vladimir Vapnik

How (2): the training and the metrics

- have parametrized distribution family $\{d^\theta(\mathbf{x})\}$
- How do we find θ such that $d^\theta(\mathbf{x}) \approx P_{data}(\mathbf{x})$ $D = \{x_i\}$

A number of options

3) Full Bayes:

$$P(\theta | D) = \frac{P(D | \theta)P(\theta)}{P(D)}$$

4) Adversarial methods (can be MLE)

All in general intractable, exactly

...some require values we may not be able to compute ($d^\theta(x)$)

more on this in a second....

But which parametrized distribution family $\{d^\theta(\mathbf{x})\}$?

How to specify distributions in general...

by characterizing the probabilities:

$$f : [N] \rightarrow \mathbb{R}^+; d(x) = \frac{f(x)}{\sum_x f(x)} \quad \text{c.f. sqashing functions in ML}$$

by characterizing a generating process

e.g. the stationary distribution of a *Markov chain*

$$P = \begin{bmatrix} P_{1,1} & P_{1,2} & \dots & P_{1,j} & \dots & P_{1,S} \\ P_{2,1} & P_{2,2} & \dots & P_{2,j} & \dots & P_{2,S} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ P_{i,1} & P_{i,2} & \dots & P_{i,j} & \dots & P_{i,S} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ P_{S,1} & P_{S,2} & \dots & P_{S,j} & \dots & P_{S,S} \end{bmatrix}.$$

But which parametrized distribution family $\{d^\theta(\mathbf{x})\}$?

Boltzmann machines

- energy based models
- stochastic recurrent NNs

→ GENERATE DISTRIBUTION OVER $\{-1, 1\}^n$

$$E(x_1, \dots, x_n, x_{n+1}, \dots, x_m) = \sum_{i,j} w_{ij} x_i \cdot x_j$$

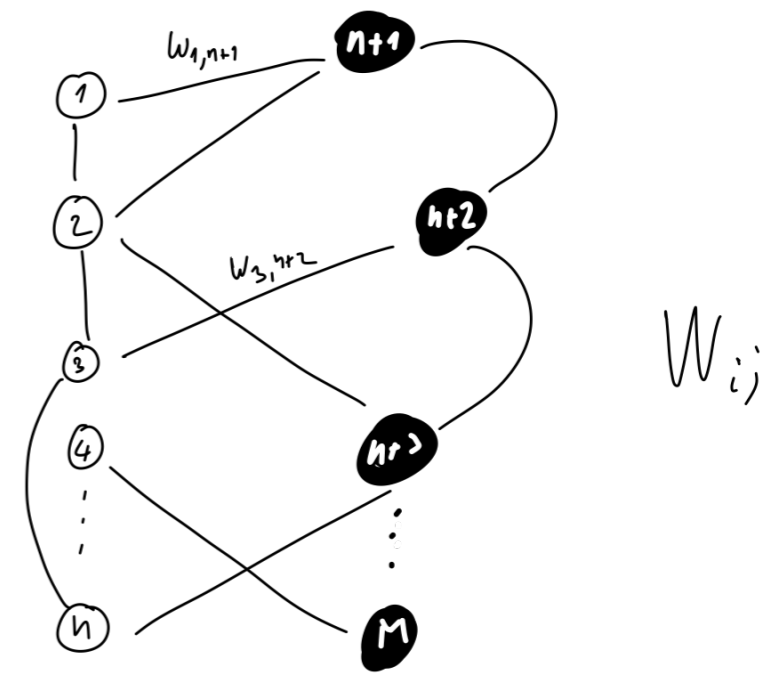
⇒ ISING MODEL!

THERMAL STATE: $P(\bar{x}) = \frac{\exp[\beta E(\bar{x})]}{Z}$ $Z = \sum_x \exp(\beta E(x))$

"Gibbs state"

"Boltzmann distribution"

$$\beta = \frac{1}{T} = \text{"inverse temperature"}$$



But which parametrized distribution family $\{d^\theta(\mathbf{x})\}$?

Boltzmann machines

→ GENERATE DISTRIBUTION OVER $\{-1, 1\}^n$

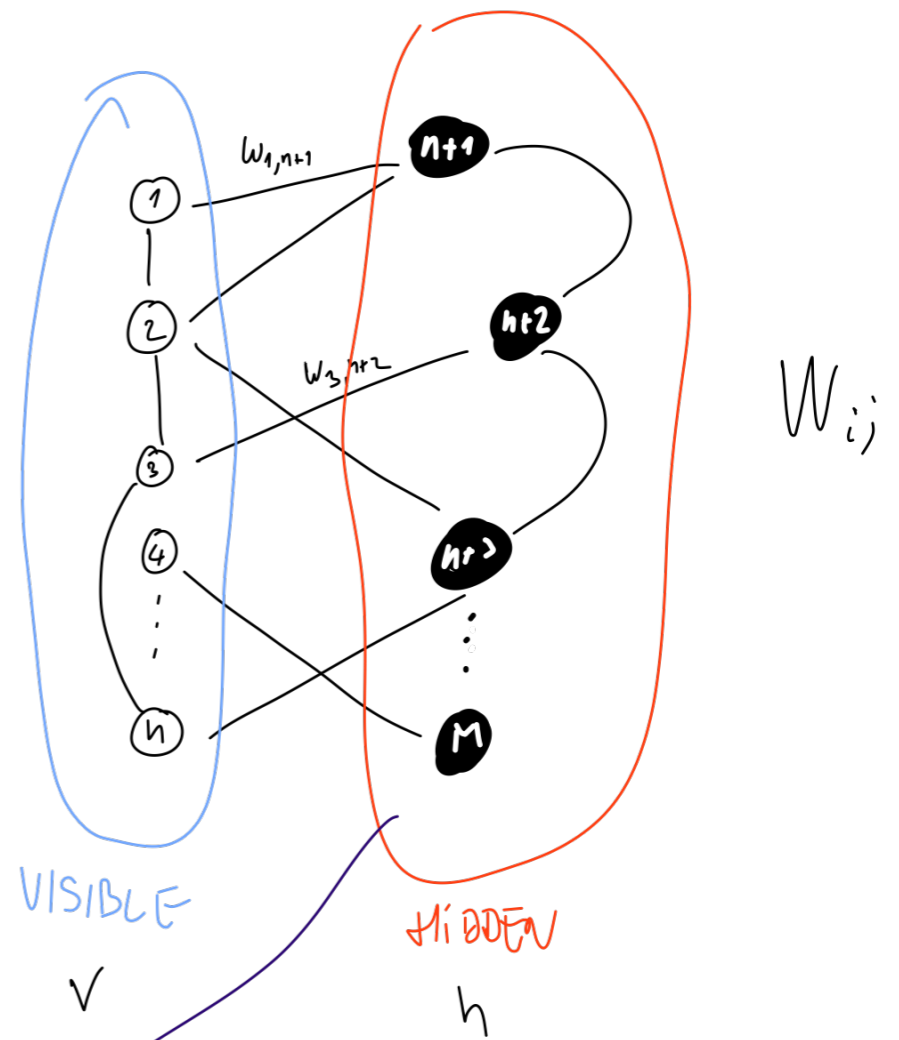
$$E(x_1 \dots x_n, x_{n+1} \dots x_m) = \sum_{i,j} w_{ij} x_i \cdot x_j$$

⇒ ISING MODEL!

THERMAL STATE:
$$P(\bar{x}) = \frac{\exp[\beta E(\bar{x})]}{Z}$$

$$P(v, h) \dots$$

$$BM := P_w(v) = \sum_h P_w(v, h)$$



POINT: CORRELATIONS!

But which parametrized distribution family $\{d^\theta(\mathbf{x})\}$?

Boltzmann machines

→ GENERATE DISTRIBUTION OVER $\{-1, 1\}^n$

$$E(x_1 \dots x_n, x_{n+1} \dots x_m) = \sum_{i,j} w_{ij} x_i \cdot x_j$$

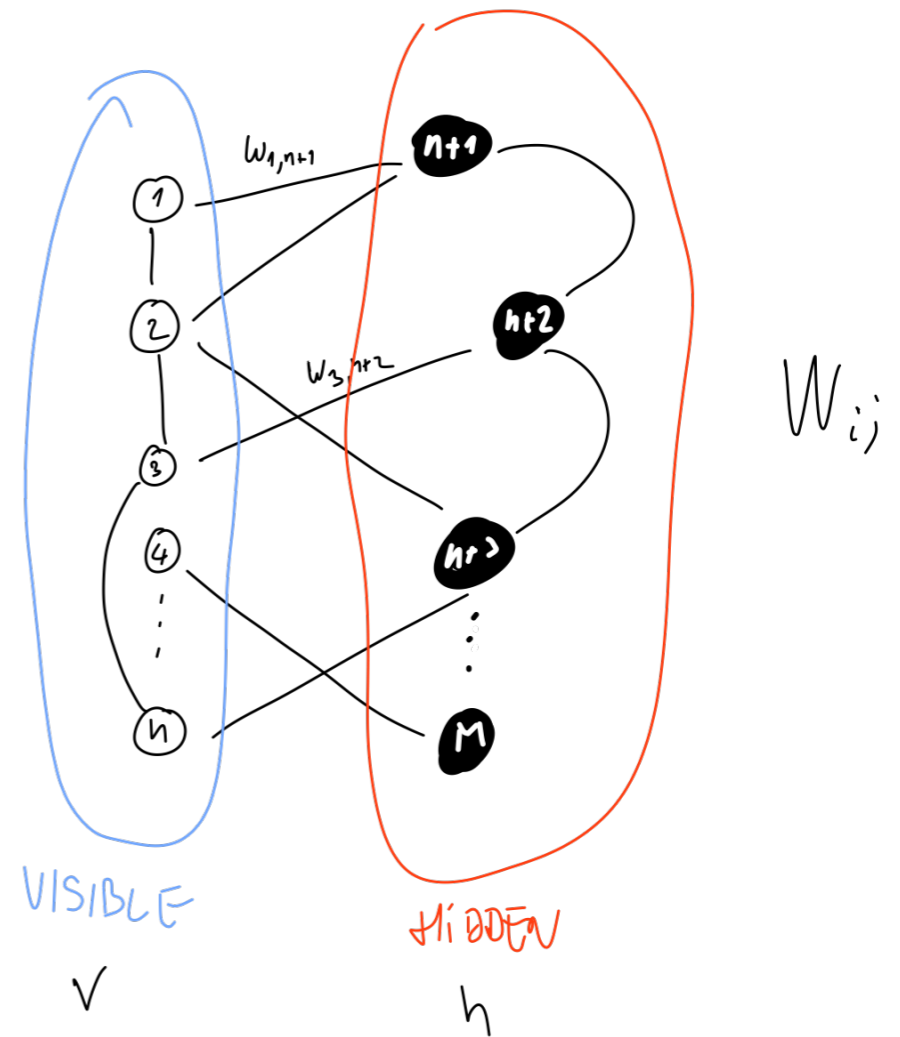
⇒ ISING MODEL!

THERMAL STATE: $P(\bar{x}) = \frac{\exp[\beta E(\bar{x})]}{Z}$

C.S.

"SOFTMAX"

"SQUASHING FUNCTION"



But which parametrized distribution family $\{d^\theta(\mathbf{x})\}$?

Boltzmann Machines (BM)

How do we use this? Run? Train?

Powerful but heavily intractable

Cannot run:

CS: Low-temperature sampling NP-hard

Phys: Need to compute the partition function

Cannot train:

Training requires conditional Gibbs sampling: intractable

Even approximating log-likelihood too hard

"ISING MODEL"

"

"QUADRATIC UNCONSTRAINED BINARY OPTIMIZATION"

IS "NP-HARD"

- Q.S. IS FNP-HARD

- LOW-TEMP. SAMPLING SOLVES Q.S.

But which parametrized distribution family $\{d^\theta(\mathbf{x})\}$?

Enter Restricted Boltzmann Machines (RBM) More tractable!

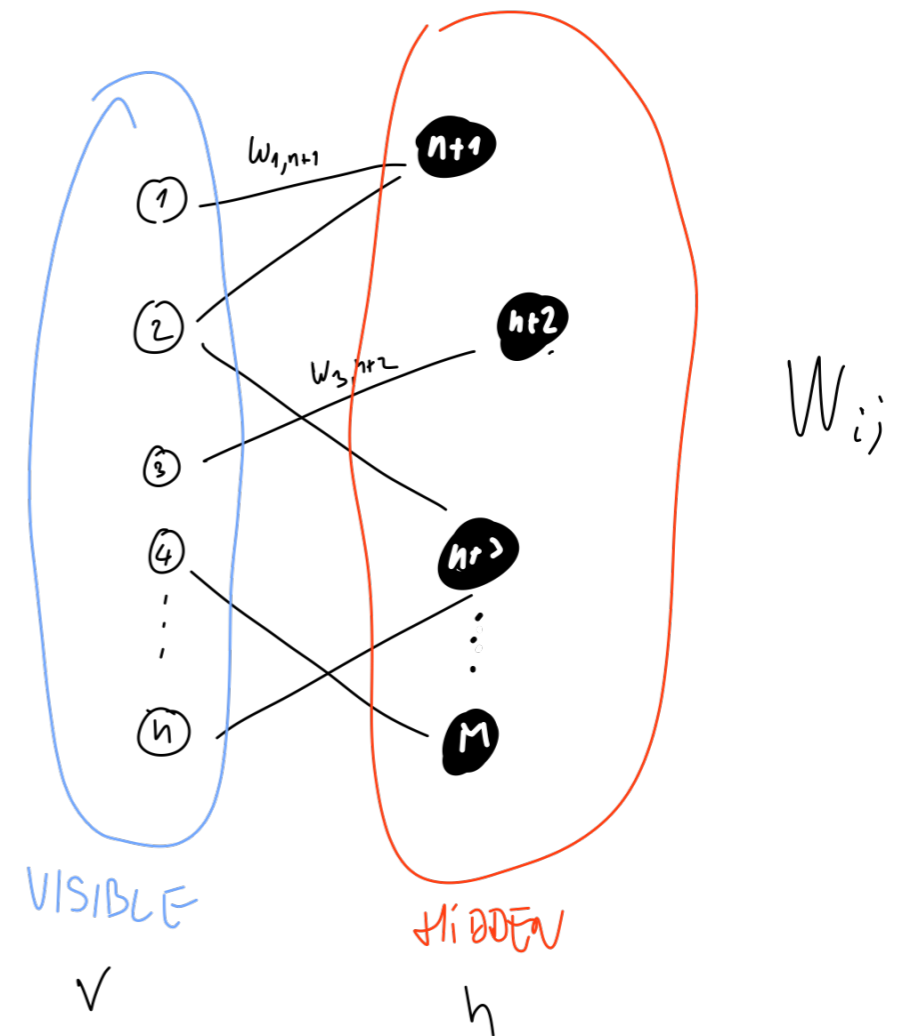
Bipartite nature makes things conditionally factorize, e.g.

$$P(\mathbf{v}) = \frac{1}{Z} \prod_{j=1}^m e^{b_j v_j} \prod_{i=1}^n \left(1 + e^{c_i + \sum_{j=1}^m w_{ij} v_j} \right)$$

This allows for more tractable training algorithms

$$\frac{\partial \log p(\mathbf{v})}{\partial w_{ij}} = \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model}$$

Still “model “ is hard. NB: Z is still HARD!
See “contrastive divergence”... approximation...

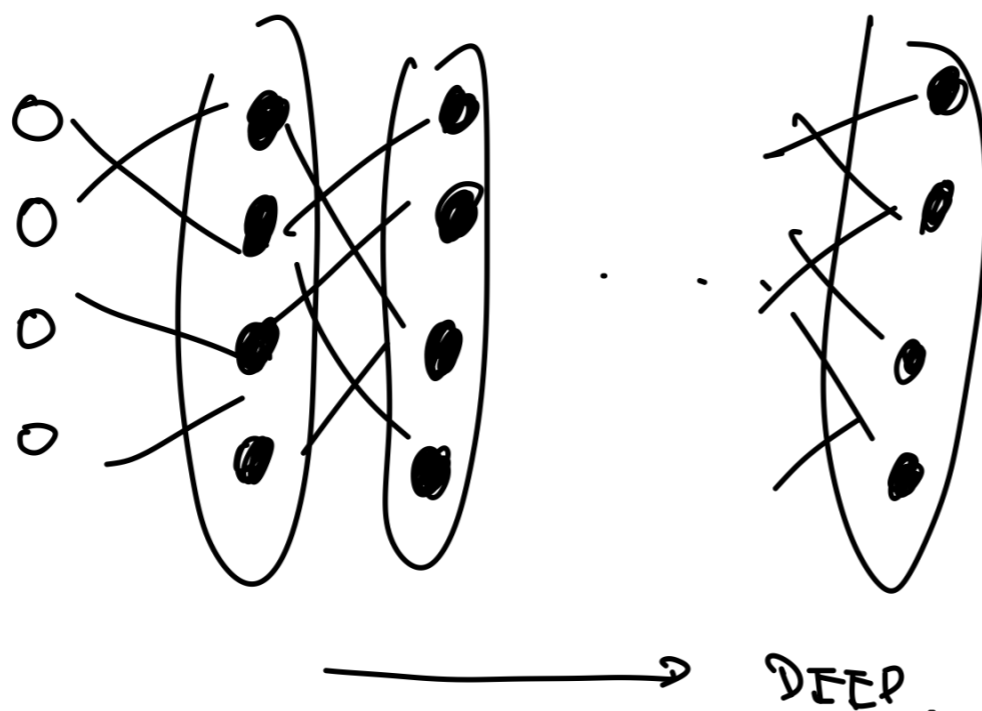


=> BIPARTITE GRAPH!

RBM is universal (Freund, Haussler, '94) (analog of Cybenko theorem)

RBM vs BM is analogous to NN vs deep NN

What we mean by deep BM (DBM):



FIXUP

STILL RBM IS LESS POWERFUL

DBM IS ESSENTIALLY EQUIVALENT TO BM.

RBM vs (D)BM (summary)

BM and RBM both energy models: distribution specified via an (bipartite) Ising model over “visible” and “hidden” units, by marginalizing over the hidden units

Both universal, but BM significantly more expressive for same number of units (has to do with the independence of hidden variables)

Training much less costly in RBM case (but still very expensive)

In both cases, to even *run* the model, you need to use a sampler, e.g. Markov Chain Monte Carlo

QC/QM meets (R)BM

Quantum-applied BMs:

(R)BMs can be used to parametrize quantum states (wavefunctions) [Carleo, Troyer '16]
(weights allowed to be complex-valued)

Quantum-enhanced BMs:

Quantum computers **can help train (R)BMs**: key step in training is (approximate) sampling from Gibbs (Boltzmann, $P(v, h)$) distribution; QCs help get those: Quantum walks; Q. semidefinite programs; Annealing. Quantum Approximate Optimization algorithms

QM can generalize (R)BM: move away from classical Ising

$$H = - \sum_a \Gamma_a \sigma_a^x - \sum_a b_a \sigma_a^z - \sum_{a,b} w_{ab} \sigma_a^z \sigma_b^z$$

QC/QM meets (R)BM

Quantum-enhanced BMs:

Quantum computers **can help train (R)BM**s: key step in training is (approximate) sampling from Gibbs (Boltzmann, $P(v,h)$) distribution; QCs help get those: Quantum walks; Q. semidefinite programs; Annealing. Quantum Approximate Optimization algorithms

QM can generalize (R)BM: move away from classical Ising

$$H = - \sum_a \Gamma_a \sigma_a^x - \sum_a b_a \sigma_a^z - \sum_{a,b} w_{ab} \sigma_a^z \sigma_b^z$$

In all cases. The physical intuition is: encode distribution in density matrices: mixed states.

In the quantum world... already **pure states** encode (many) distributions!

Next: quantum circuit Born machine

Boltzmann machine \Rightarrow distribution is a (marginalized) Boltzmann distribution

(distr which maximizes the entropy, subject to mean energy condition)

statistical mechanics

Born machine \Rightarrow distribution governed by the **Born rule** of quantum mechanics

$$P_O(i | |\psi\rangle) = |\langle i | \psi \rangle|^2 = \text{Tr}[|i\rangle\langle i| |\psi\rangle\langle\psi|]$$

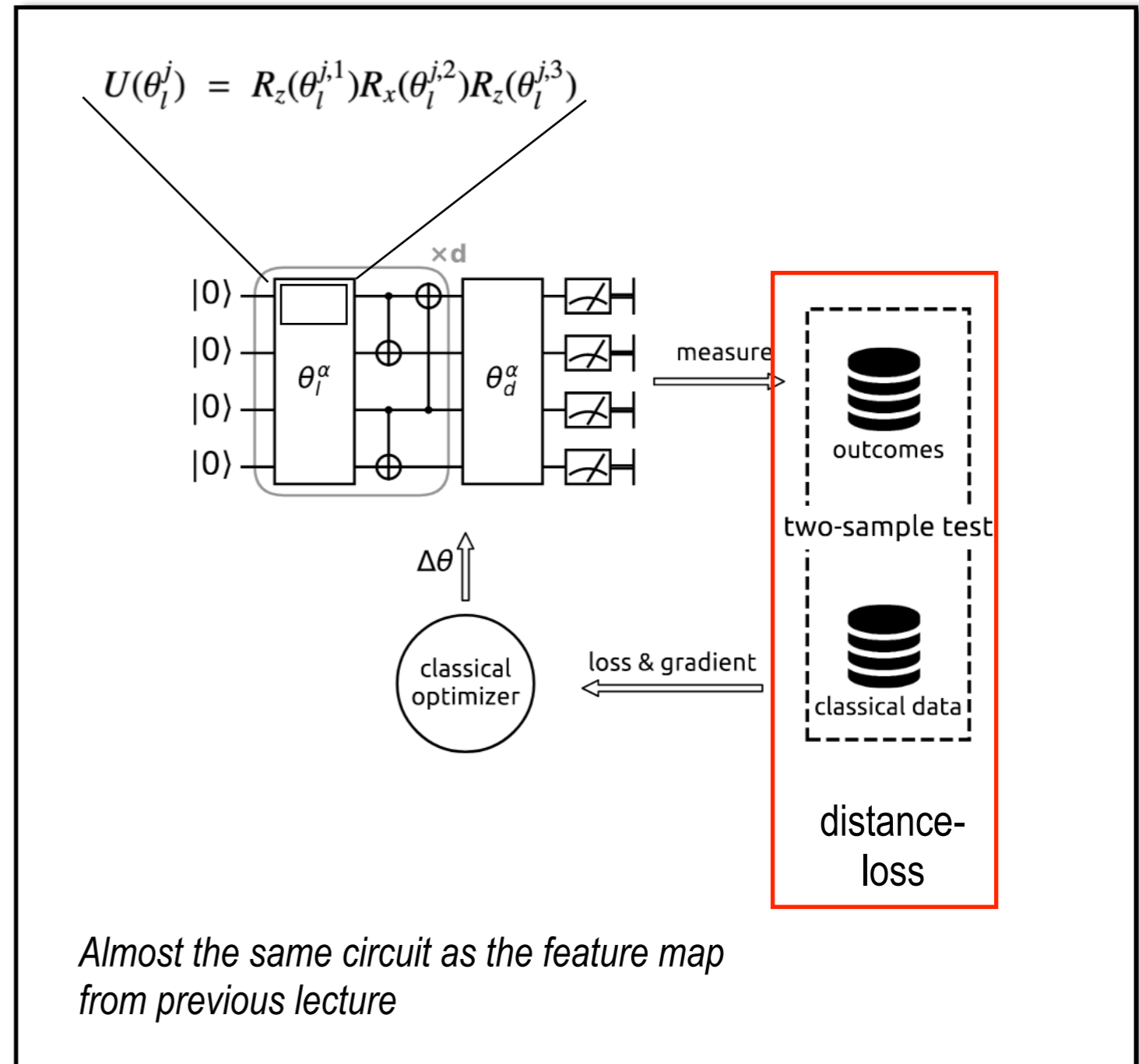
$$P_O(i | \rho) = \text{Tr}[P_i \rho]$$

Advantage: sampling (inference) is (quantum) easy (no nasty MCMC or similar)

Hypothesis family: parametrized circuit+ computational basis measurement

Plus: evaluation is easy, does not involve hard sampling procedures

“Minus”: “implicit” model. No direct access to likelihood, no way to compute $d_{\theta}(x)$, for a chosen x



How to measure distances between distributions based on sampling... deep waters.

MLE optimization/KL optimization?

$$\mathcal{L}_{\text{KL}} = - \sum_{\mathbf{z}} \pi(\mathbf{z}) \log(p_{\theta}(\mathbf{z})) = -\mathbb{E}_{\mathbf{z} \sim \pi}(\log(p_{\theta}(\mathbf{z})))$$

$$\frac{\partial \mathcal{L}_{\text{KL}}}{\partial \theta} \sim \sum_{\mathbf{z}} \frac{\pi(\mathbf{z})}{p_{\theta}(\mathbf{z})} (p_{\theta}^{-}(\mathbf{z}) - p_{\theta}^{+}(\mathbf{z}))$$

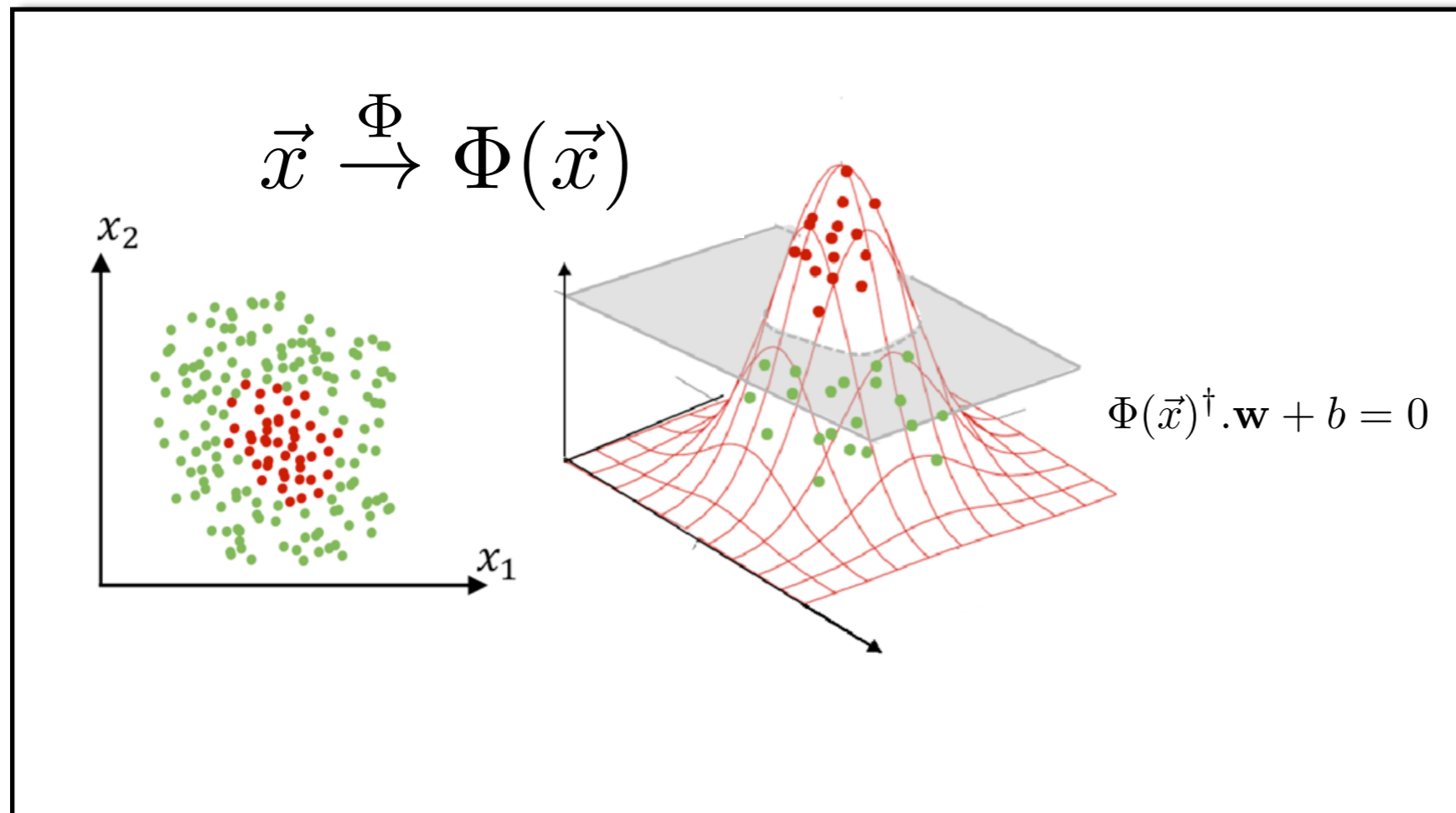
Intractable

Here: Squared maximum mean discrepancy (MMD)

$$\mathcal{L} = \left\| \sum_x p_{\theta}(x) \phi(x) - \sum_x \pi(x) \phi(x) \right\|^2 = \mathbb{E}_{x \sim p_{\theta}, y \sim p_{\theta}} [K(x, y)] - 2 \mathbb{E}_{x \sim p_{\theta}, y \sim \pi} [K(x, y)] + \mathbb{E}_{x \sim \pi, y \sim \pi} [K(x, y)].$$

Squared maximum mean discrepancy (MMD)

$$\text{MMD}(P, Q) = \|\mathbb{E}_{X \sim P}[\varphi(X)] - \mathbb{E}_{Y \sim Q}[\varphi(Y)]\|_{\mathcal{H}}.$$



Feature map

Key point: express differences of distributions JUST in terms of expected values (first moment)

Easier to measure. If Φ is powerful enough... zero iff same distribution.

$$\mathcal{L} = \left\| \sum_x p_\theta(x) \phi(x) - \sum_x \pi(x) \phi(x) \right\|^2 = \mathbb{E}_{x \sim p_\theta, y \sim p_\theta} [K(x, y)] - 2 \mathbb{E}_{x \sim p_\theta, y \sim \pi} [K(x, y)] + \mathbb{E}_{x \sim \pi, y \sim \pi} [K(x, y)].$$

$$K(x, y) = \phi(x)^T \phi(y)$$

Everything boils down to sampling and estimating the mean...

$$k(x, x') = \exp(-\|x - x'\|^2 / (2 \sigma^2))$$

Gradients also can be expressed in terms of samples + parameter shift rule.

$$\frac{\partial \mathcal{L}}{\partial \theta_l^\alpha} = \mathbb{E}_{x \sim p_{\theta^+}, y \sim p_\theta} [K(x, y)] - \mathbb{E}_{x \sim p_{\theta^-}, y \sim p_\theta} [K(x, y)] - \mathbb{E}_{x \sim p_{\theta^+}, y \sim \pi} [K(x, y)] + \mathbb{E}_{x \sim p_{\theta^-}, y \sim \pi} [K(x, y)].$$

$$\frac{\partial f_{\mathcal{P}}(\vec{\theta})}{\partial \theta_j} = \frac{f_{\mathcal{P}}(\vec{\theta} + \frac{\pi}{2} e_j) + f_{\mathcal{P}}(\vec{\theta} - \frac{\pi}{2} e_j)}{2}$$

For parameter shift rule, recall 7.3.1 of

<http://liacs.leidenuniv.nl/~dunjikov/aQa/aQa-Lecture-3-vqe1.pdf>

Some results (bottom line it works)

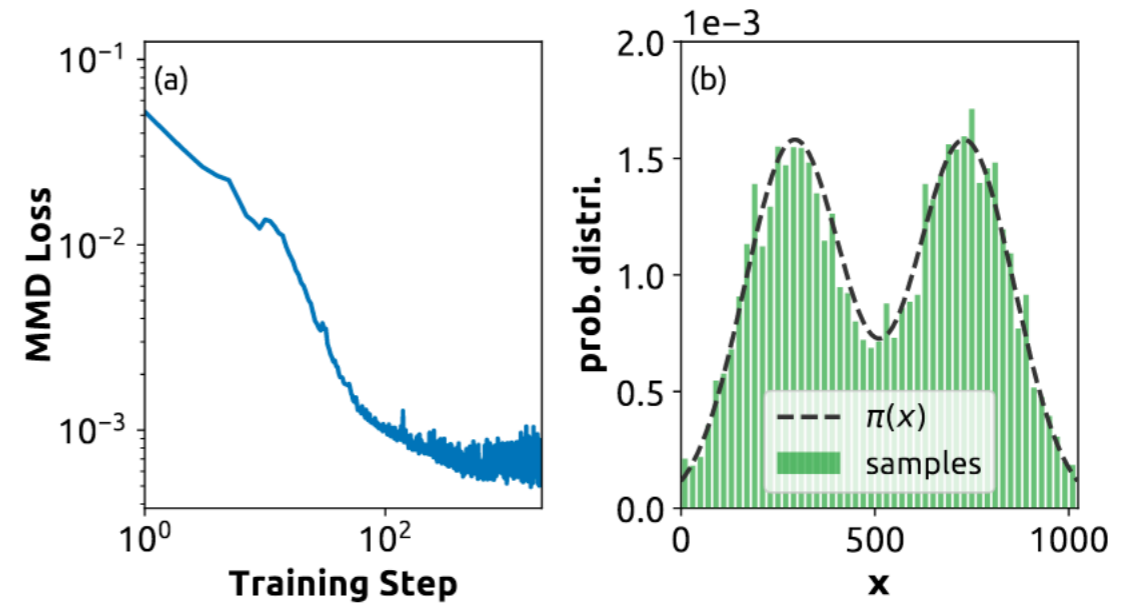
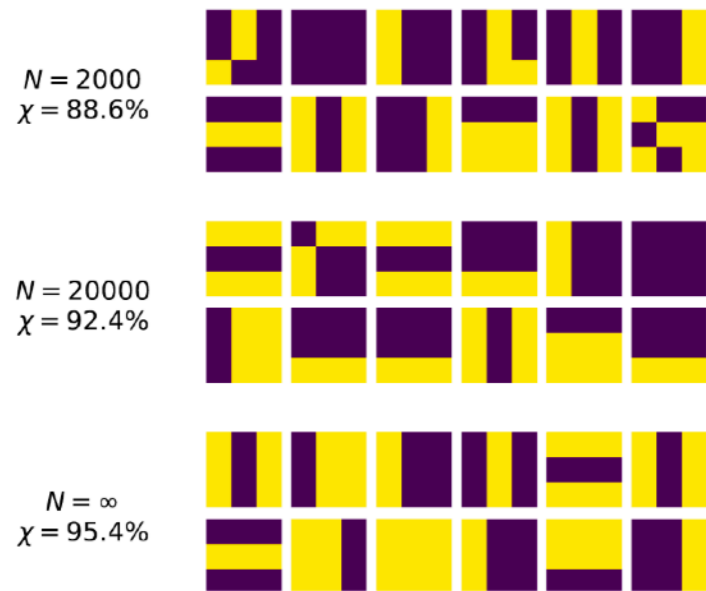


Figure 6. (a) The MMD loss as a function of Adam training step. (b) Histogram for samples generated by a trained QCBM with a bin width 20 (green bars), in comparison with the exact probability density function (black dashed line).

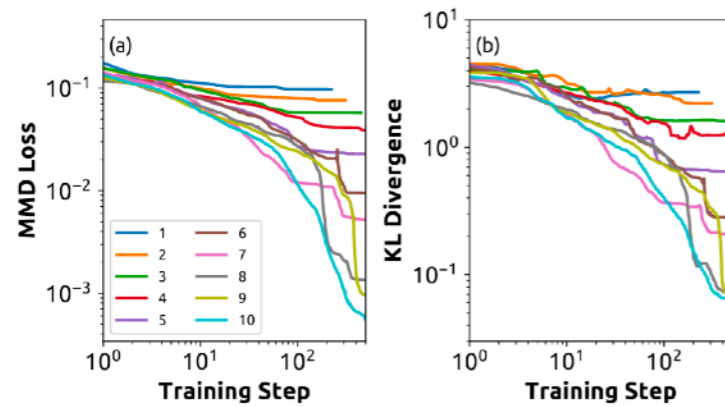


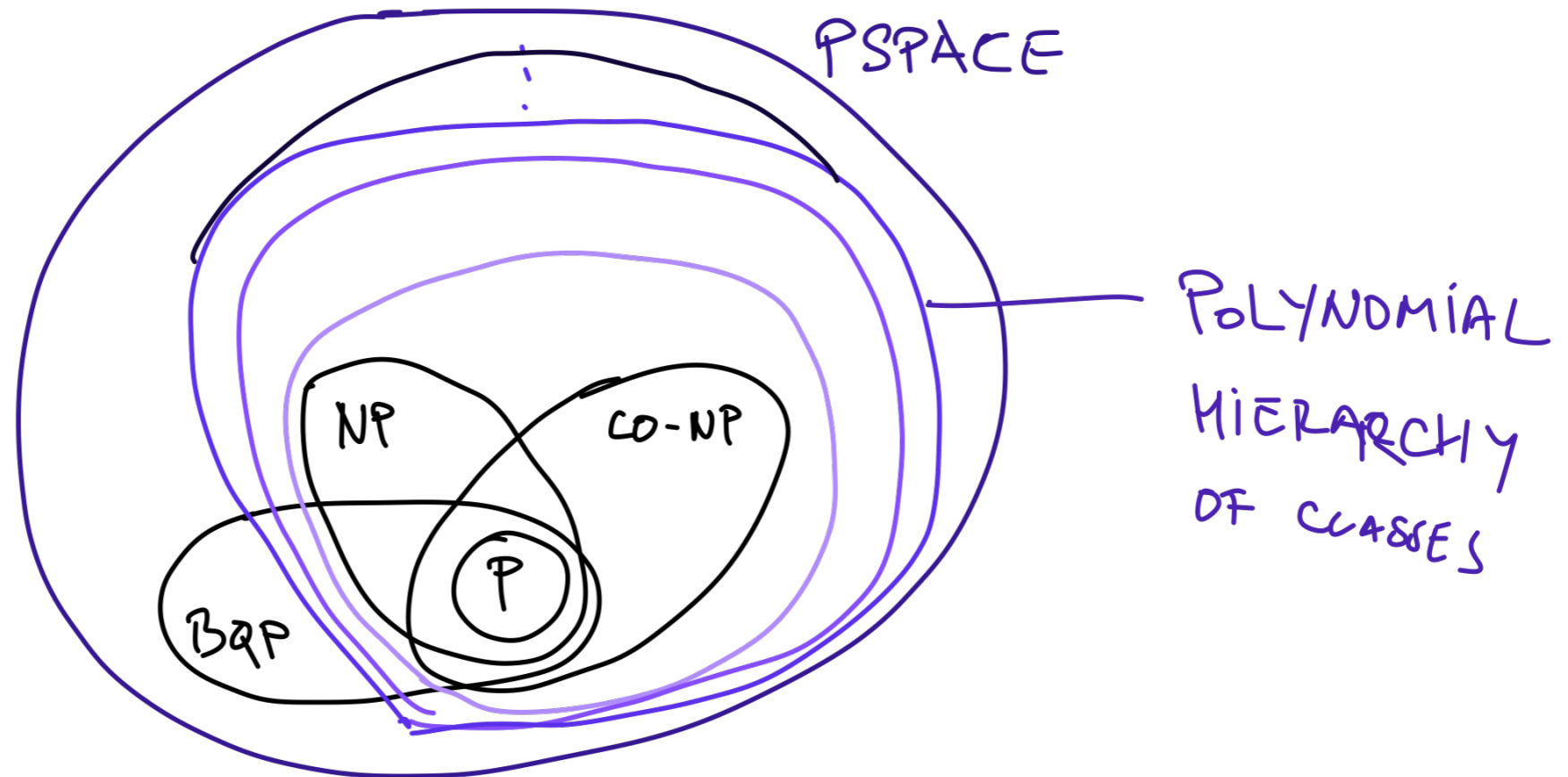
Figure 5. Losses as a function of training step for circuit depth $d = 1, \dots, 10$. (a) The MMD loss Eq. (1), and (b) the corresponding KL divergence. Here, we use L-BFGS-B optimizer with exact gradient.

Quantum supremacy of IQP sampling

Theorem 3.2 ([35]). *Assume it is $\#P$ -hard to approximate $|\mathcal{Z}|^2$ up to a relative error $1/4 + o(1)$ for $1/24$ fraction of instances over the choice of the weights and biases, J_{ij}, b_k . If it is possible to classically sample from the output probability distribution of any IQP circuit in polynomial time, up to an additive error of $1/384$ in total variation distance, then there is a BPP^{NP} algorithm to solve any problem in $P^{\#P}$, and hence the polynomial hierarchy collapses to the third level)*

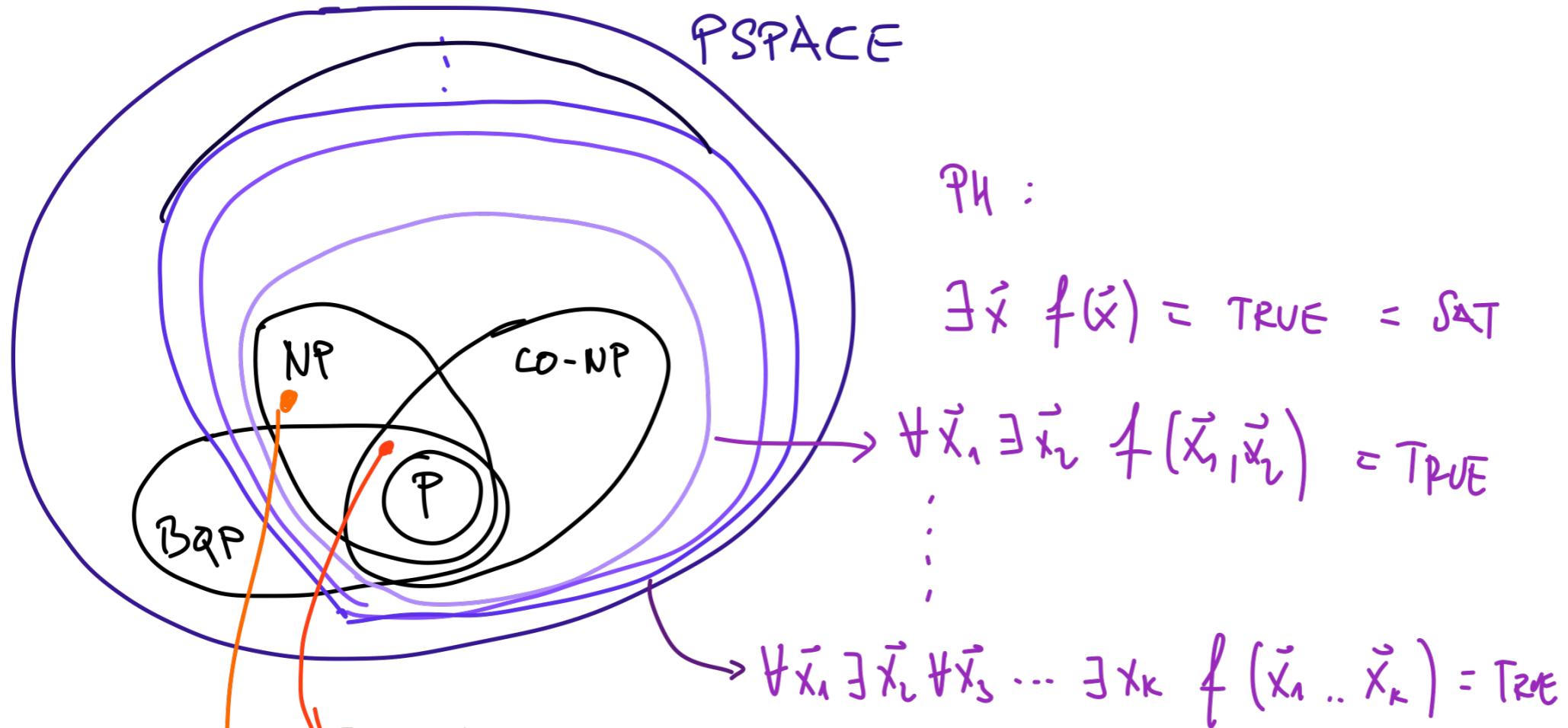
The "quantum supremacy" business..

Complexity theory (of decision problems)



The "quantum supremacy" business..

Complexity theory (of decision problems)



FACTORIZING

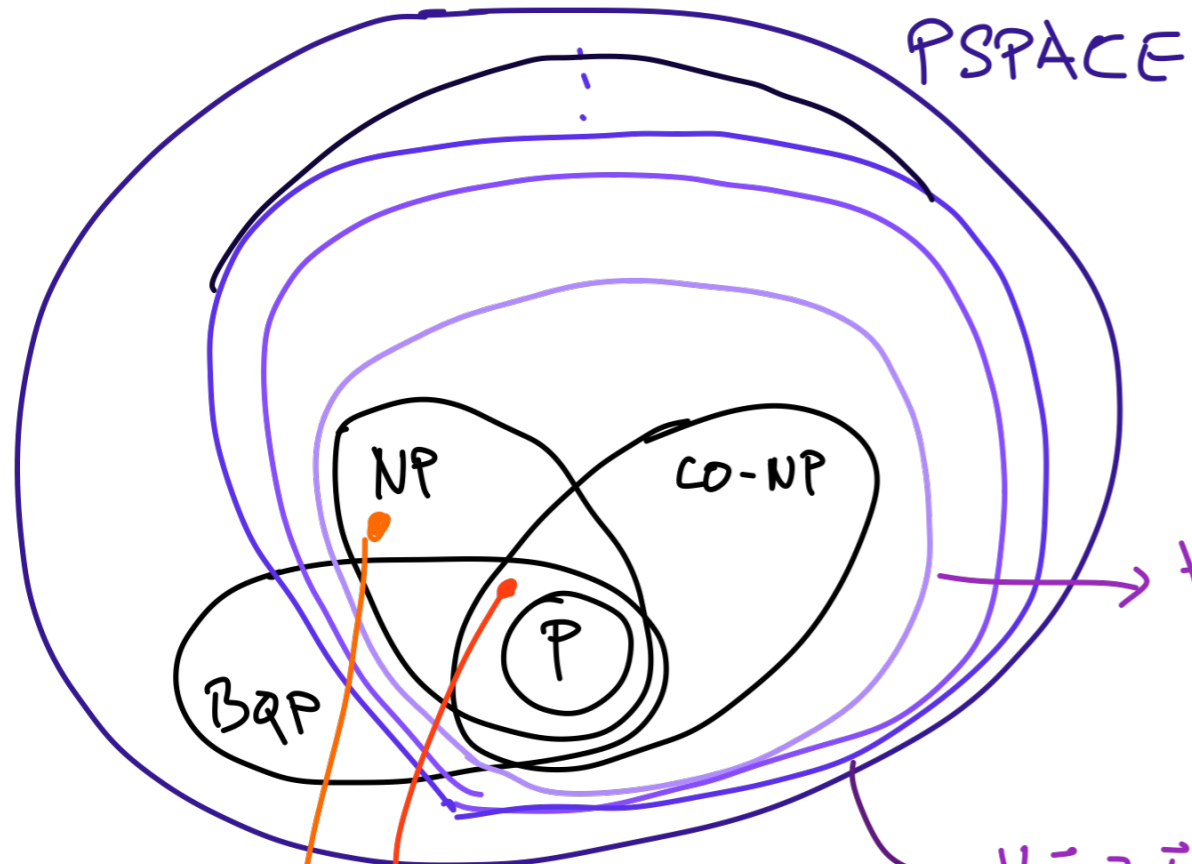
? $\exists \vec{x} = (x_1 \dots x_n)$

$f(\vec{x}) = (x_1 \vee \neg x_2 \vee x_3) \wedge \dots \wedge (x_{n-2} \vee \neg x_3 \dots) = \text{TRUE}$

SATISFIABILITY

The "quantum supremacy" business..

Complexity theory (of decision problems)



PH :

$$\exists \vec{x} f(\vec{x}) = \text{TRUE} = \text{SAT} \quad \text{PH} - 0$$

$$\forall \vec{x}_1 \exists \vec{x}_2 f(\vec{x}_1, \vec{x}_2) = \text{TRUE} \quad \text{PH} - 1$$

⋮

$$\forall \vec{x}_1 \exists \vec{x}_2 \forall \vec{x}_3 \dots \exists \vec{x}_k f(\vec{x}_1 \dots \vec{x}_k) = \text{TRUE} \quad \text{PH} - k$$

FACTORIZING

$$? \exists \vec{x} = (x_1 \dots x_n)$$

$$f(\vec{x}) = (x_1 \vee \neg x_2 \vee x_3) \wedge \dots \wedge (x_{n-2} \vee \neg x_{n-1}) = \text{TRUE}$$

SATISFIABILITY

"PH DOES NOT COLLAPSE"

NOTE ...

PH

HAS NOTHING TO DO WITH QC

PER-SE

• IT CAN BE : $PH-1 = PH-n \quad \forall n \quad \&$
 $BQP \neq P$

• IT CAN BE $BQP = P \quad \& \quad PH-k \subsetneq PH-k+1 \quad \forall k$

The "quantum supremacy" business..

Sampling problems v.s. decision problems

BQP contains Factoring. Assume BQP in P.

Factoring is then in P. Unlikely but no other consequences

The "quantum supremacy" business..

Sampling problems v.s. decision problems

*BQP contains Factoring. Assume BQP in P.
Factoring is then in P. Unlikely but no other consequences*

Surprisingly: If we can sample from IQP circuits, PH collapses to 3rd level

Even though Sampling is not decision stuff! and PH is about that...

The "quantum supremacy" business..

Sampling problems v.s. decision problems

BQP contains Factoring. Assume BQP in P.

Factoring is then in P. Unlikely but no other consequences

Surprisingly: If we can sample from IQP circuits, PH collapses to 3rd level

Even though Sampling is not decision stuff! and PH is about that...

Surprisingly: Even if the sampling is with multiplicative error and even additive error (worst case)

Open: Average case hardness for additive error; has unproven conjectures atop of PH collapse

Instantaneous Quantum Polynomial-time

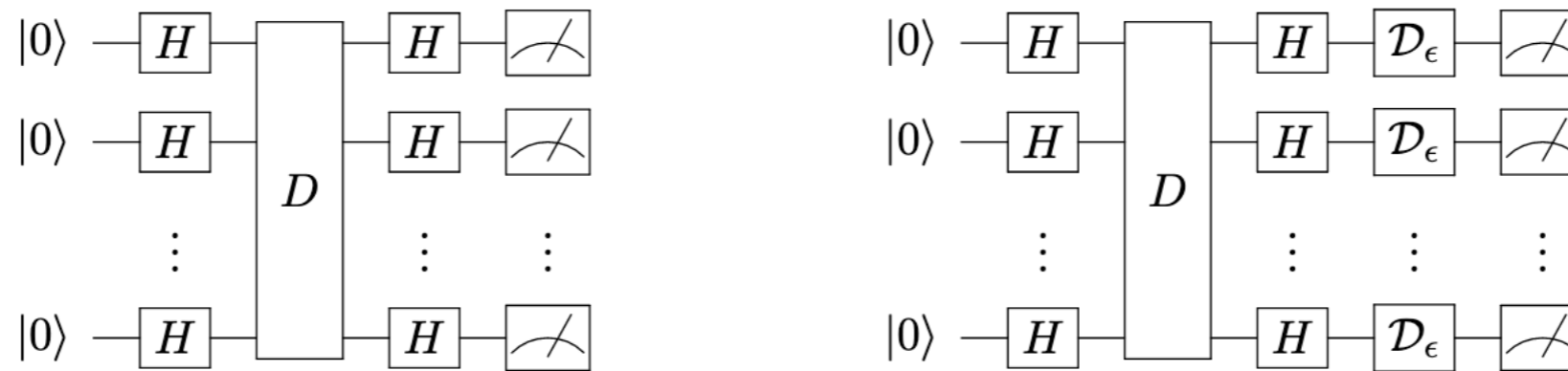


Figure 1: A standard IQP circuit, and an IQP circuit with depolarising noise. D is a circuit made up of $\text{poly}(n)$ diagonal gates.

These are our QML circuits...

The framework of Generative Adversarial Nets

Elements of Generative models:

- a) model
(something that can generate samples from distributions / specification of distribution)
- b) distance measure (metric, divergence)
(something to tell you how “far we are from real data”, to define the *loss*)
- c) computational method to minimize loss

The framework of Generative Adversarial Nets

Elements of Generative models:

- a) model
(something that can generate samples from distributions / specification of distribution)
- b) distance measure (metric, divergence)
(something to tell you how “far we are from real data”, to define the *loss*)
- c) computational method to minimize loss

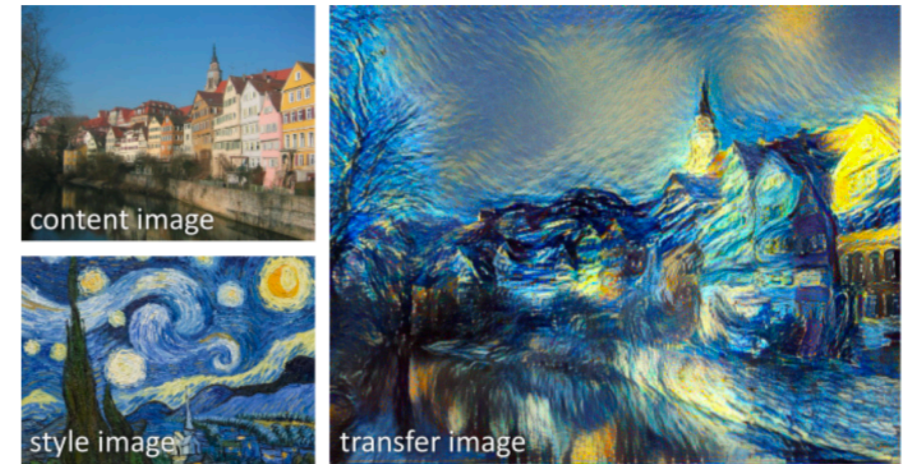
Both b) and c) can be very problematic (involving likelihoods).
How to choose metric...

The framework of Generative Adversarial Nets

Generative models about distribution $P(X)$...
when are we *faking it well*?

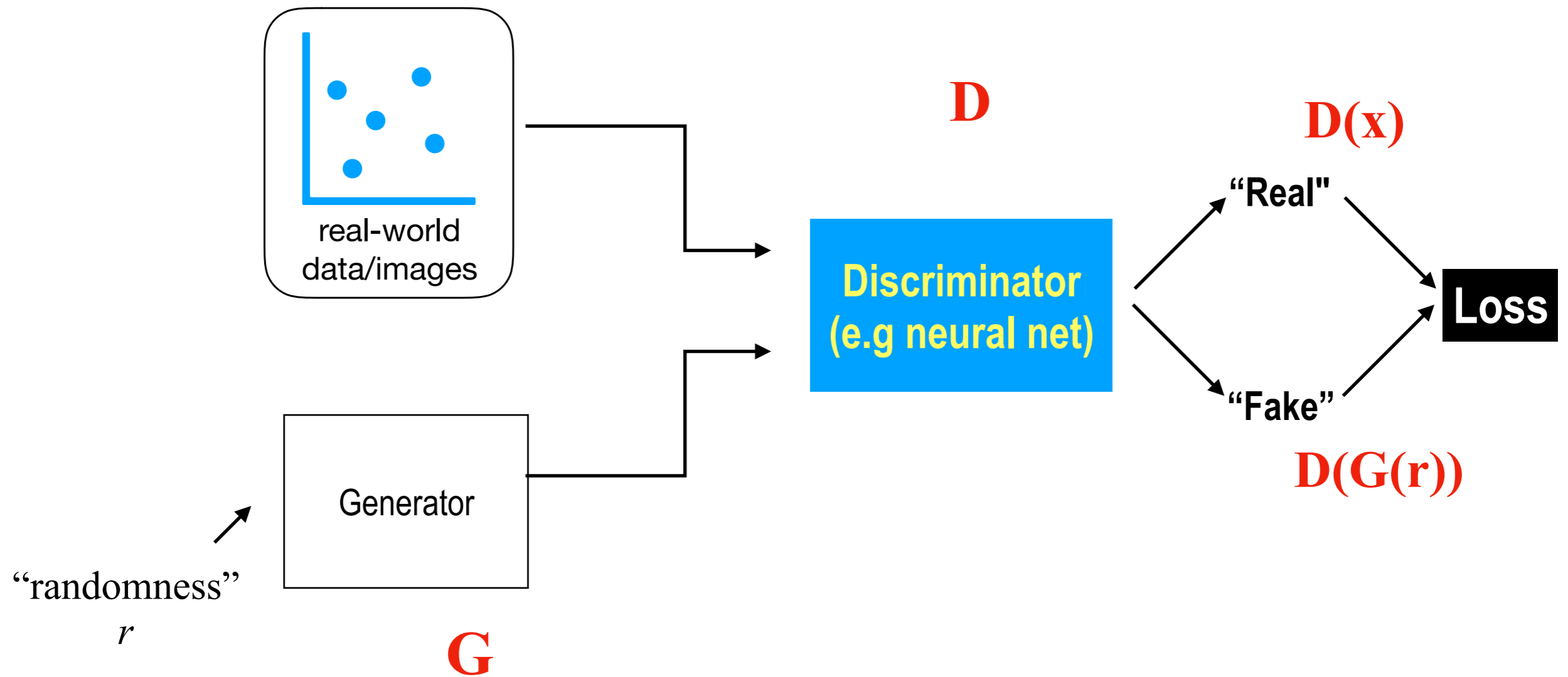
Mebbe: When an expert cannot tell real from fake...

Who is the relevant expert?
How about ML itself?

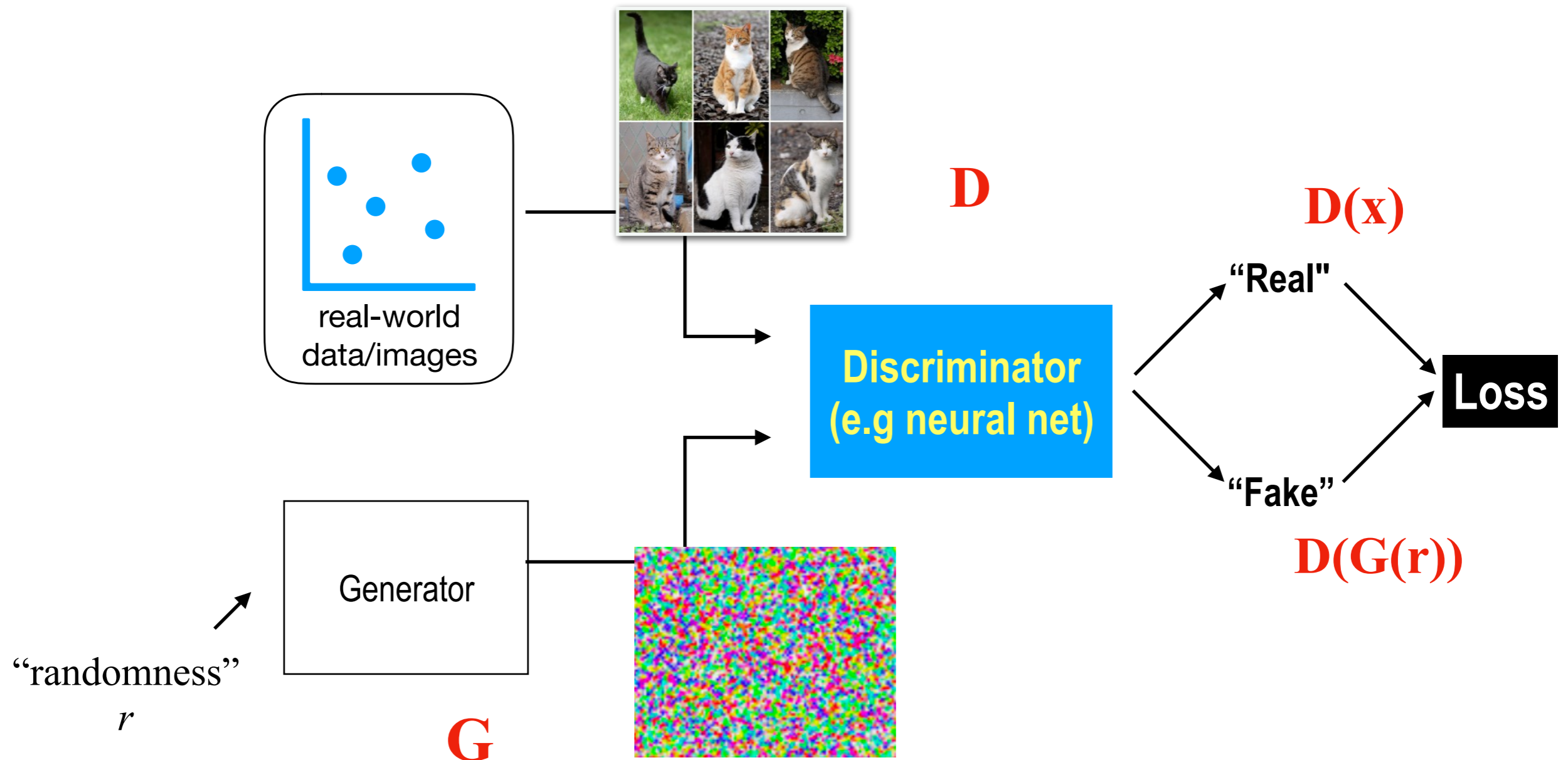


Lost Van Gogh?

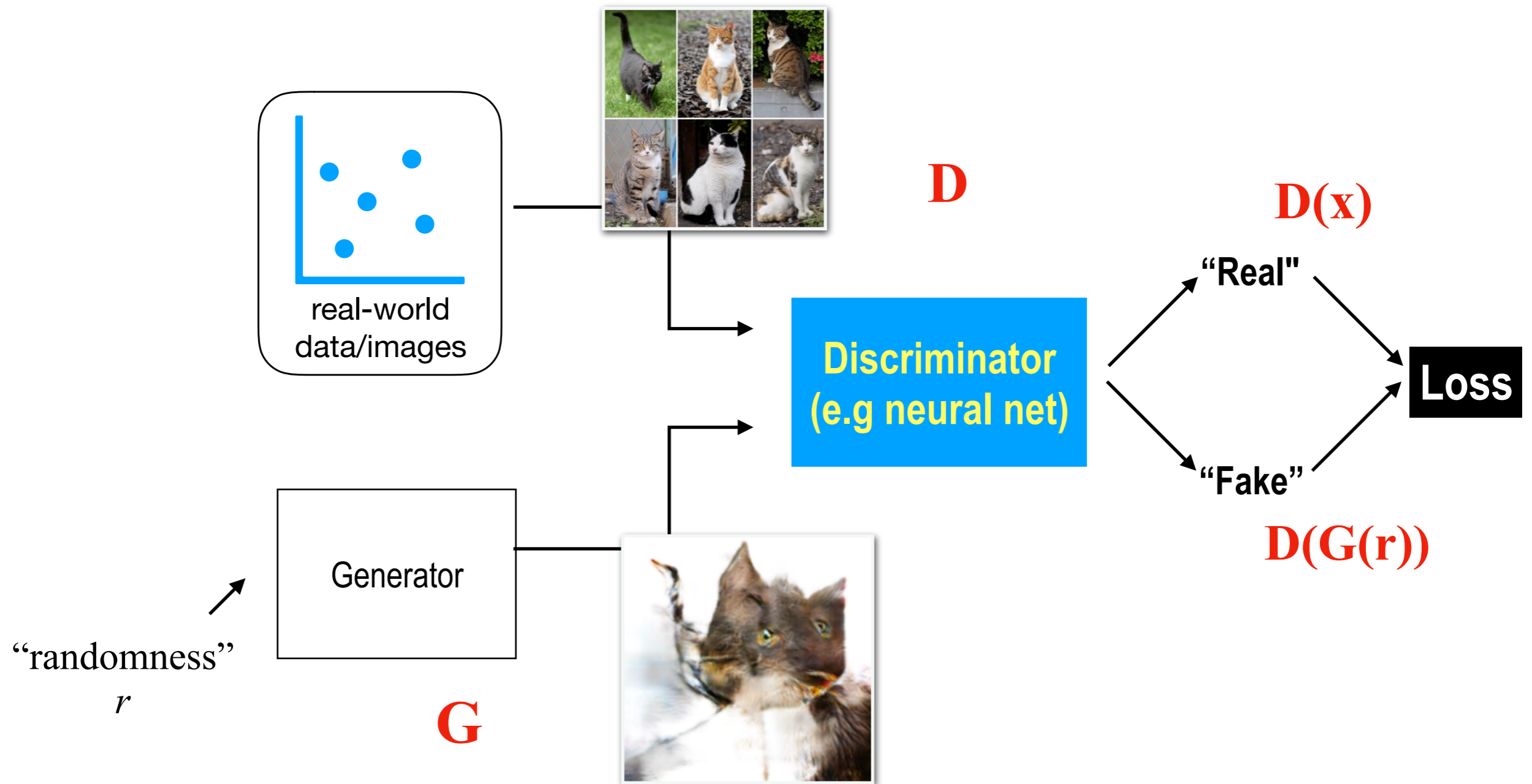
The framework of Generative Adversarial Nets



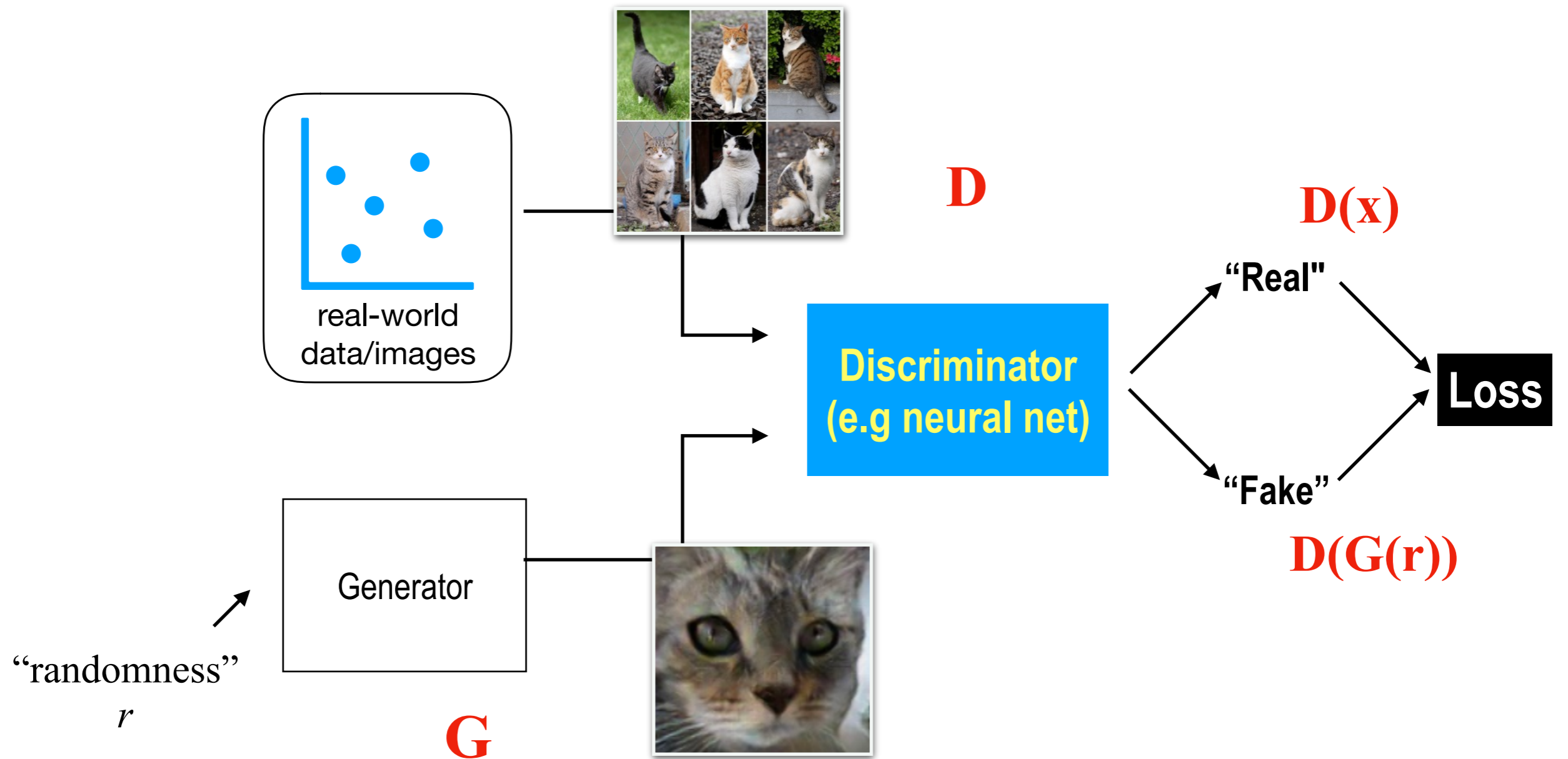
The framework of Generative Adversarial Nets



The framework of Generative Adversarial Nets



The framework of Generative Adversarial Nets

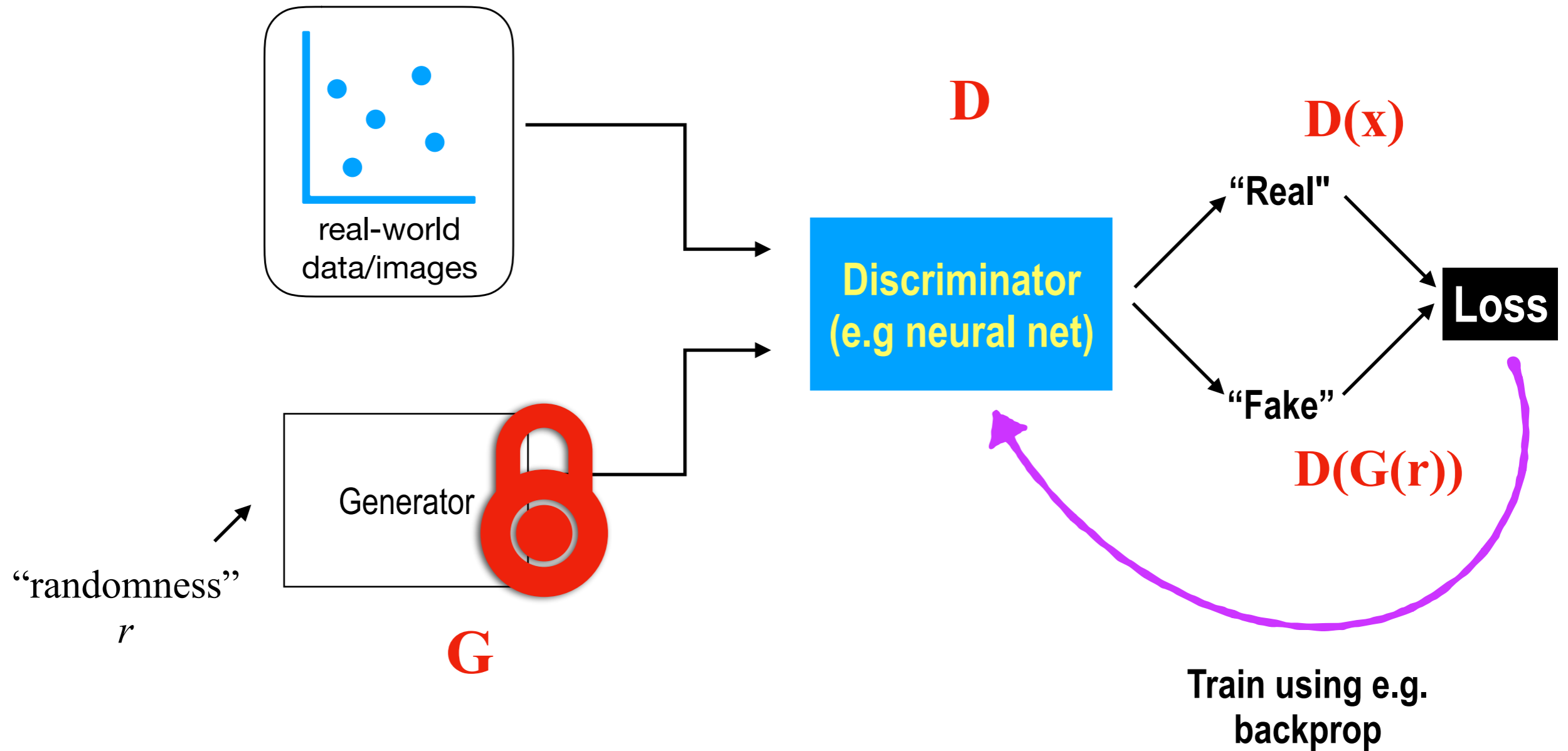


Google "Meow generator"...

The framework of Generative Adversarial Nets

Training

Phase: training discriminator

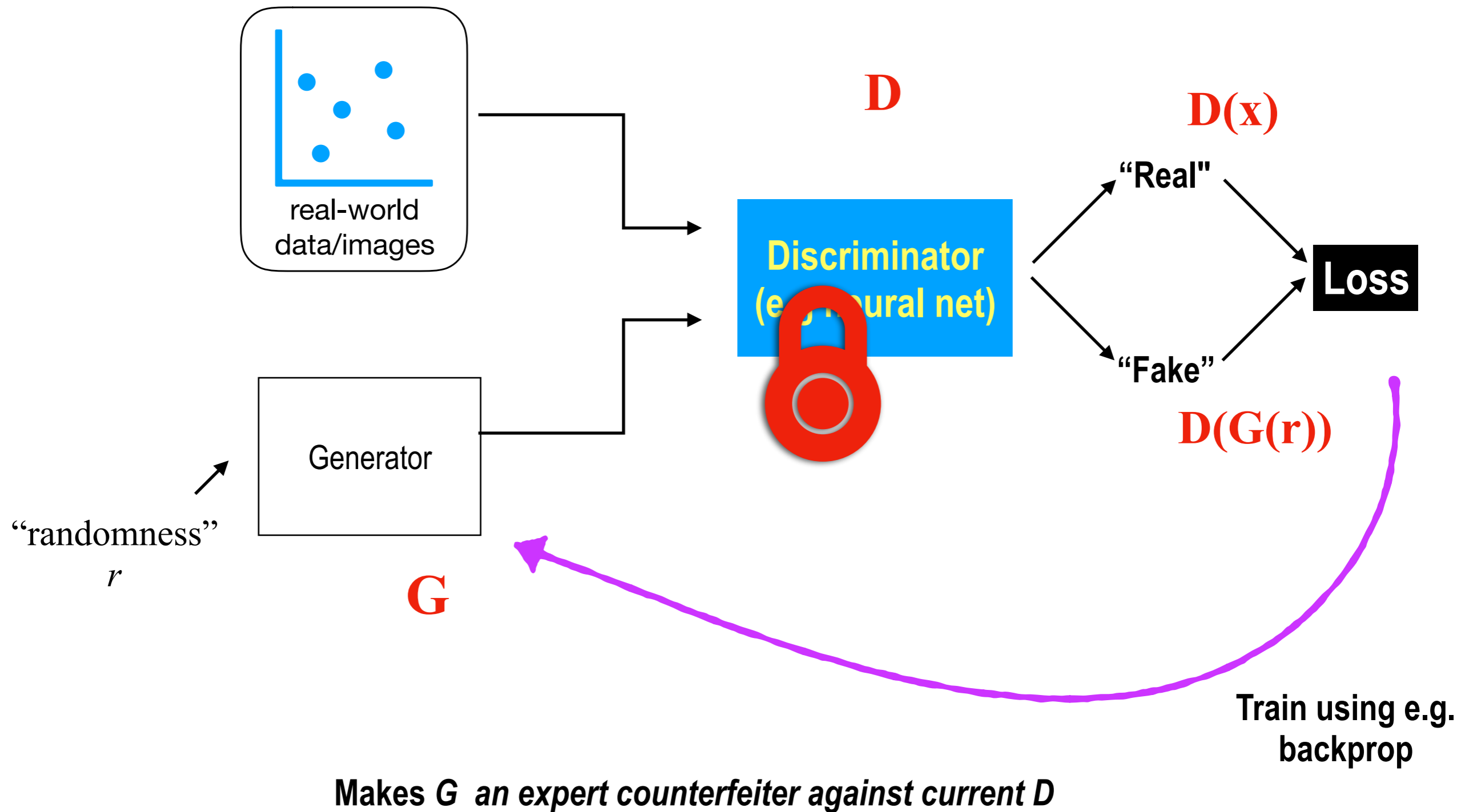


Makes D an expert counterfeit sensor against current G

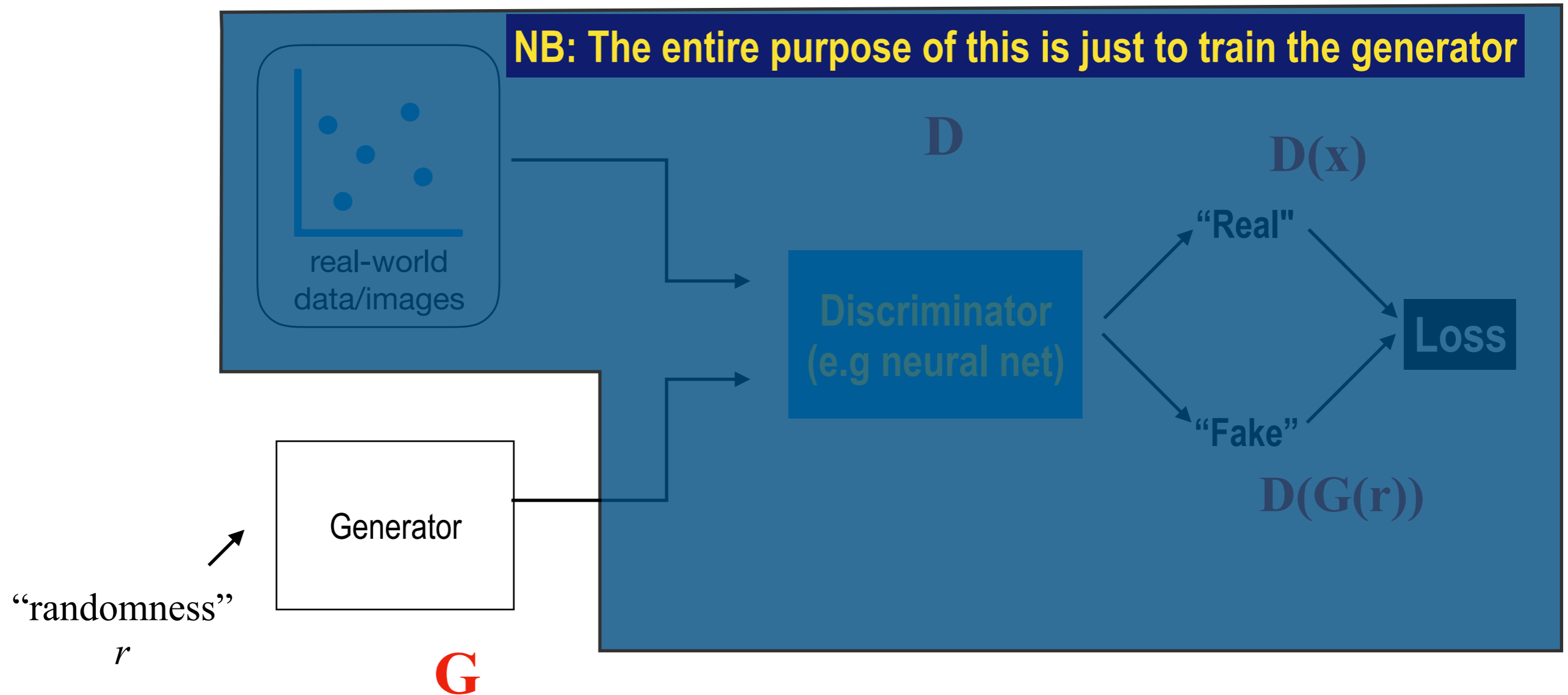
The framework of Generative Adversarial Nets

Training

Phase: training generator



The framework of Generative Adversarial Nets



The framework of Generative Adversarial Nets

Phases are iterated

Sequential solution to a *minimax game*; $V(D, G)$ is success of D

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] .$$

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(\mathbf{x}^{(i)}) + \log(1 - D(G(\mathbf{z}^{(i)})))] .$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(\mathbf{z}^{(i)}))) .$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Train D

Train G

The framework of Generative Adversarial Nets

Theory says it works:

The Nash equilibrium is achieved at:

- $p_{data}(x) = p_G(x)$
- $D(x) = 1/2$ (random guess)

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

$$C(G) = \max_D V(G, D)$$

Theorem 1. *The global minimum of the virtual training criterion $C(G)$ is achieved if and only if $p_g = p_{data}$. At that point, $C(G)$ achieves the value $-\log 4$.*

Proposition 2. *If G and D have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given G , and p_g is updated so as to improve the criterion*

$$\mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))]$$

then p_g converges to p_{data}

The framework of Generative Adversarial Nets

Practice says it sometimes works.

Advantages:

- Works with models where sampling is easy (to train feed-forward NN with noise input)
- No need to compute Maximum Likelihood estimation
- Could be robust to overfitting
- Empirical success

| | Deep directed graphical models | Deep undirected graphical models | Generative autoencoders | Adversarial models |
|-------------------|--|---|--|--|
| Training | Inference needed during training. | Inference needed during training. MCMC needed to approximate partition function gradient. | Enforced tradeoff between mixing and power of reconstruction generation | Synchronizing the discriminator with the generator. Helvetica. |
| Inference | Learned approximate inference | Variational inference | MCMC-based inference | Learned approximate inference |
| Sampling | No difficulties | Requires Markov chain | Requires Markov chain | No difficulties |
| Evaluating $p(x)$ | Intractable, may be approximated with AIS | Intractable, may be approximated with AIS | Not explicitly represented, may be approximated with Parzen density estimation | Not explicitly represented, may be approximated with Parzen density estimation |
| Model design | Models need to be designed to work with the desired inference scheme — some inference schemes support similar model families as GANs | Careful design needed to ensure multiple properties | Any differentiable function is theoretically permitted | Any differentiable function is theoretically permitted |

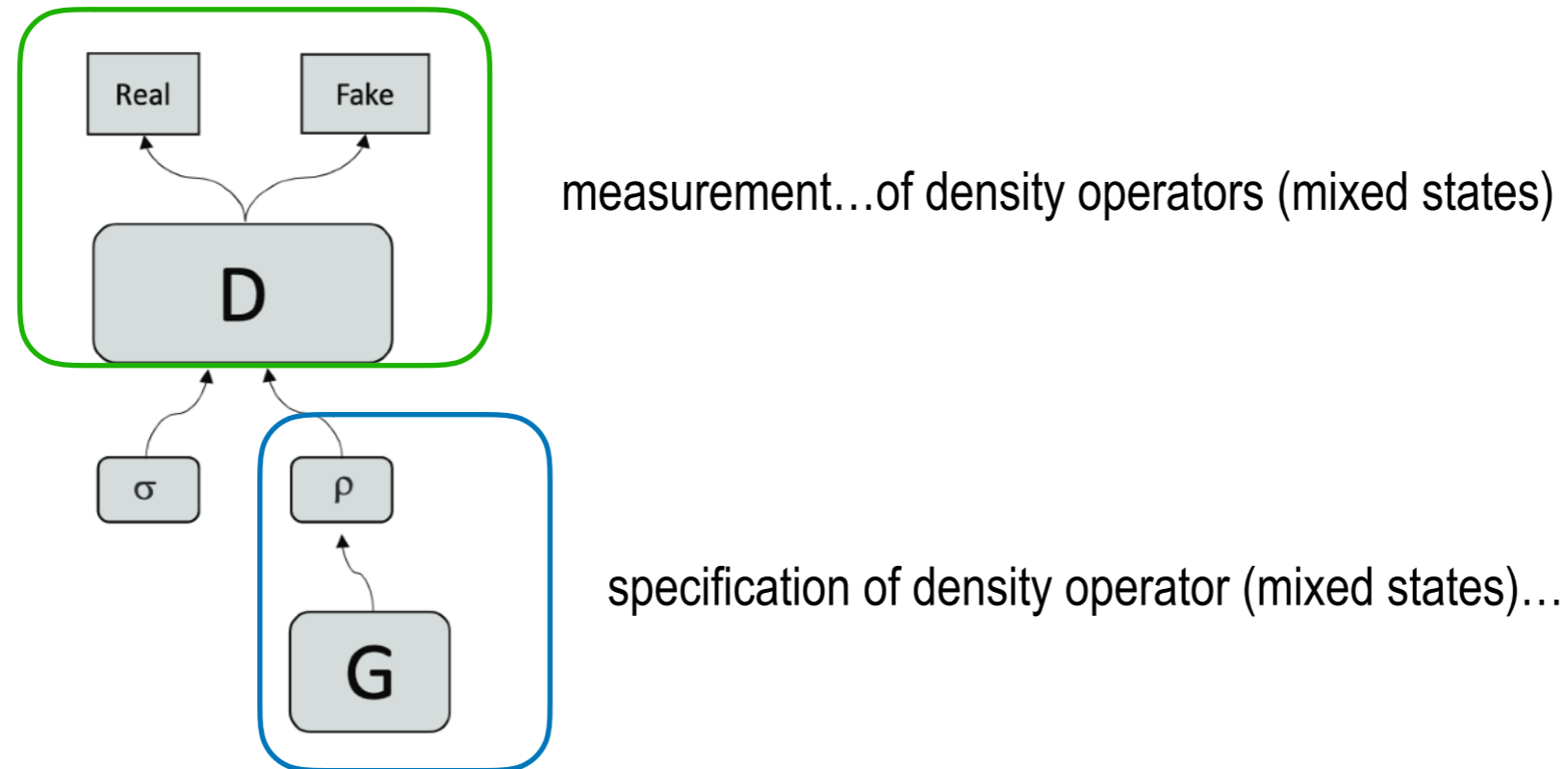
Table 2: Challenges in generative modeling: a summary of the difficulties encountered by different approaches to deep generative modeling for each of the major operations involving a model.

From GANs to Q. GANS

Note GANs is a framework which can in principle work with (almost) any pair of supervised learner and a generator.

Any part can be made quantum

The entire formalism generalizes to quantum systems



Fixed-point theorems can be proven in a vastly generalized setting, and they are actually simpler.
(theory of optimal measurements was developed by Helstrom in '76)

Quantum cases allow speed-ups, applications in many-body physics and chemistry

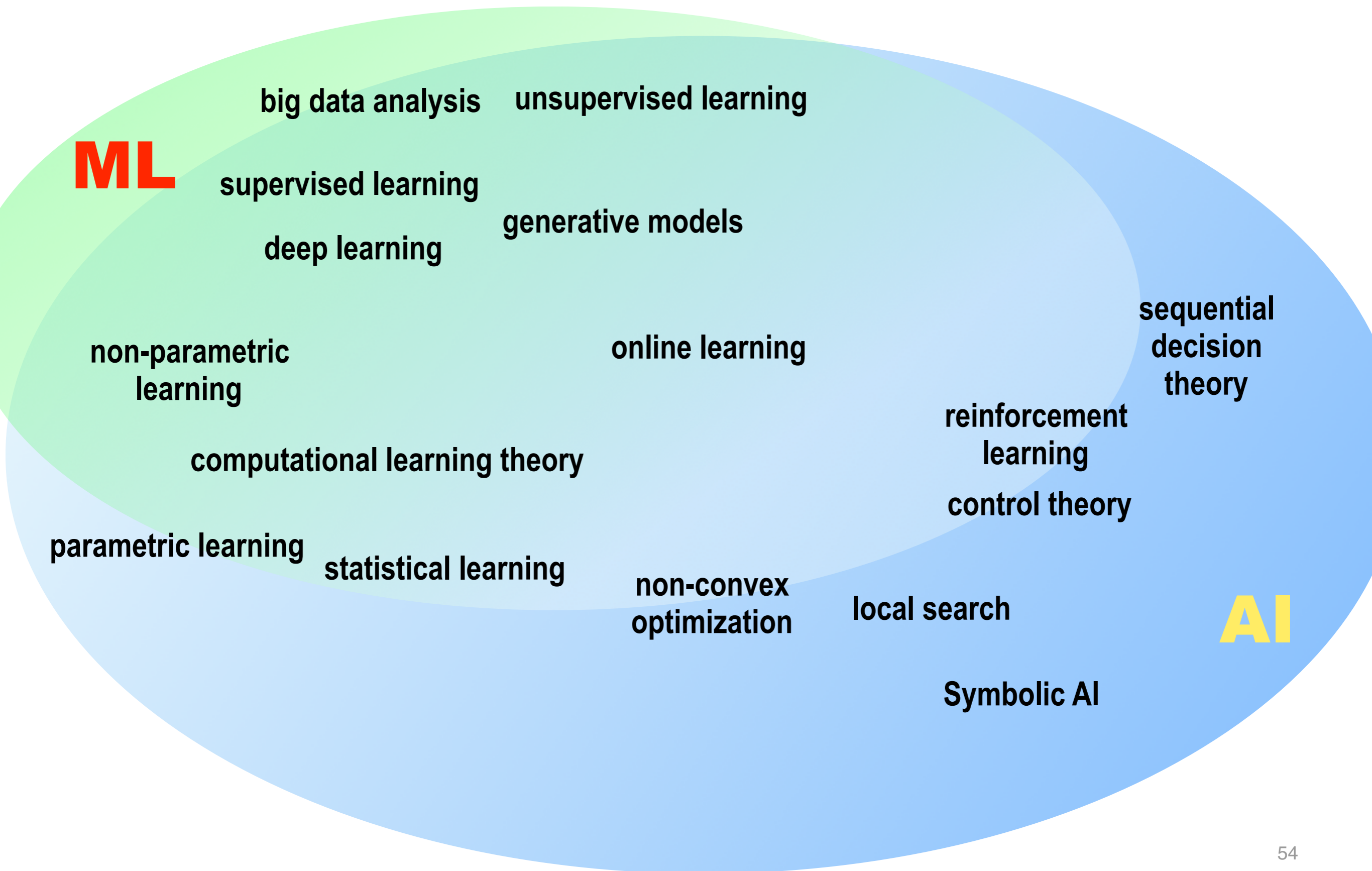
QML — the rest

Fun with variational circuits

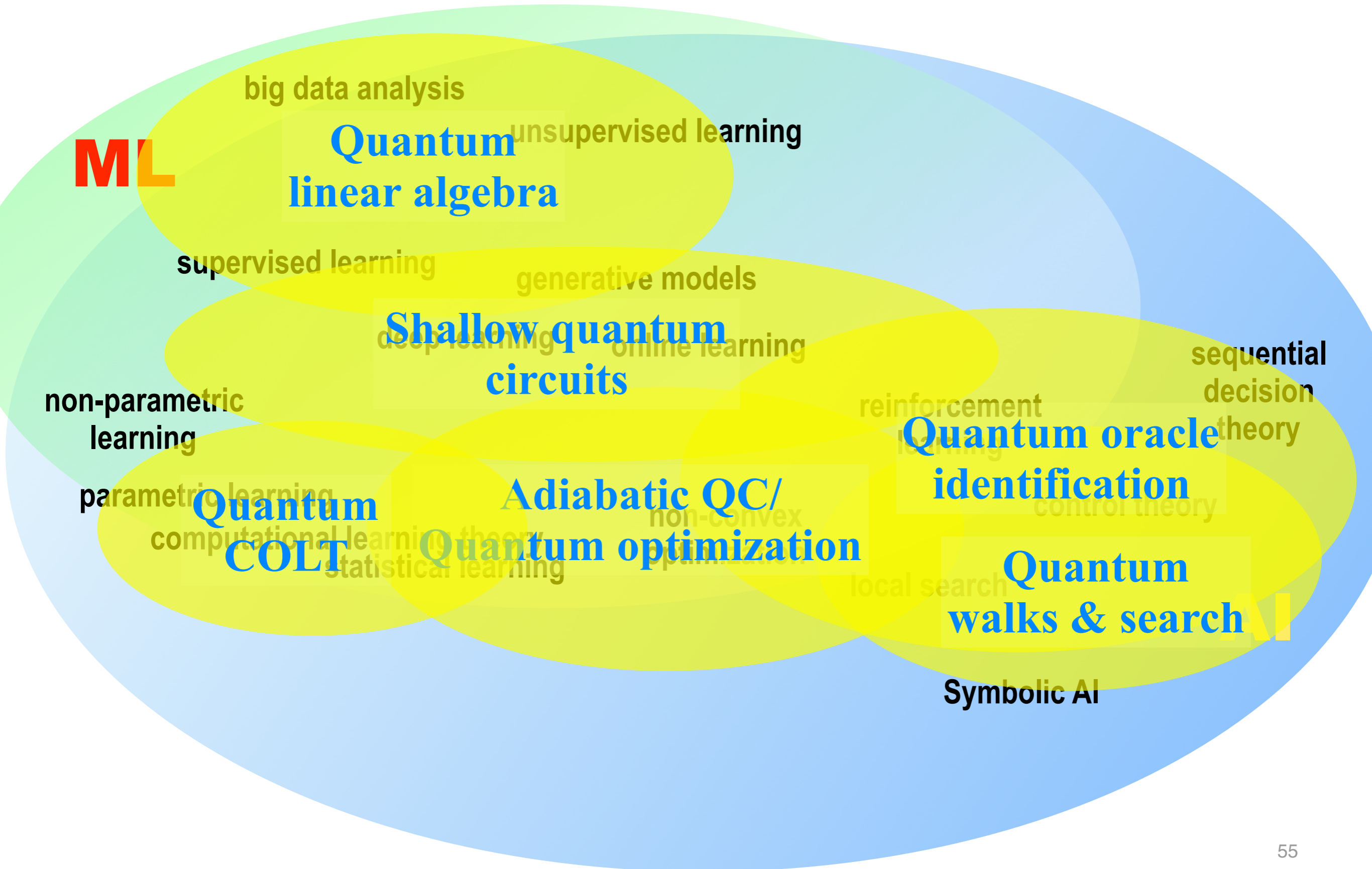
- quantum convolutional NNs
- data reuploading
- quantum transfer learning
- quantum reinforcement learning

Then there is all the rest...

Machine learning is not one thing.
AI is not even a few things.



Quantum-enhanced ML



And then there's Quantum-applied ML!

