# Introduction to Deep Learning

## Michael Kagan

### SLAC

CERN Openlab Summer Student Lectures
July 5, 2024

People are now building a **new kind of software** by assembling networks of **parameterized functional blocks** and by **training them from examples using some form of gradient-based optimization**.
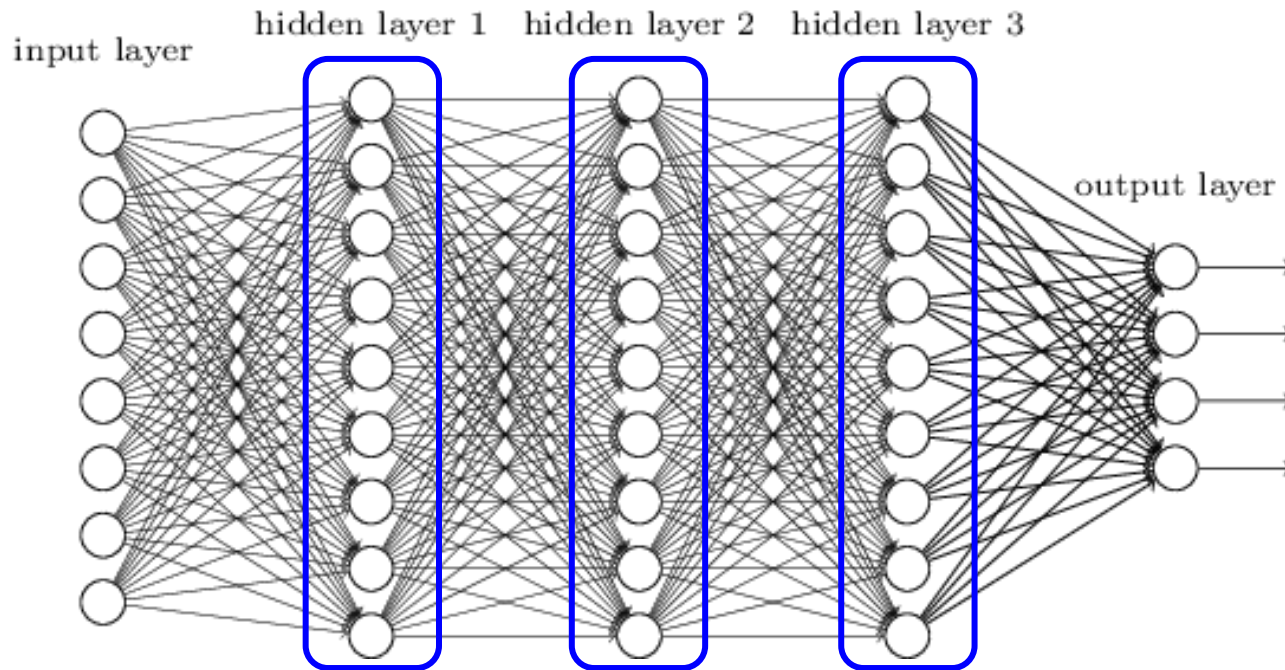
- Yann LeCun, 2018

# Modern Neural Networks

People are now building a **new kind of software** by assembling networks of **parameterized functional blocks** and by **training them from examples using some form of gradient-based optimization**.
- Yann LeCun, 2018

- Non-linear operations of data with parameters

- Layers (set of operations) designed to perform specific mathematical operations

- Chain together layers to perform desired computation

- Train system (with examples) for desired computation using gradient descent
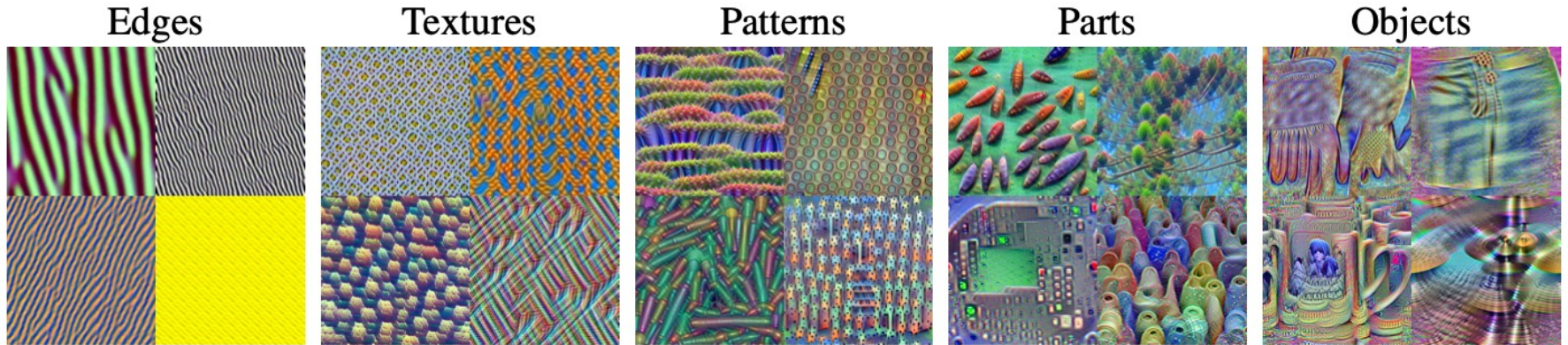
# Deep Neural Networks

- As data complexity grows, need exponentially large number of neurons in a single-hidden-layer network to capture all structure in data

- Deep networks *factorize learning* of structure in data across layers

- Large datasets, fast computing (GPU / TPU) and new training procedures / network structures made training possible
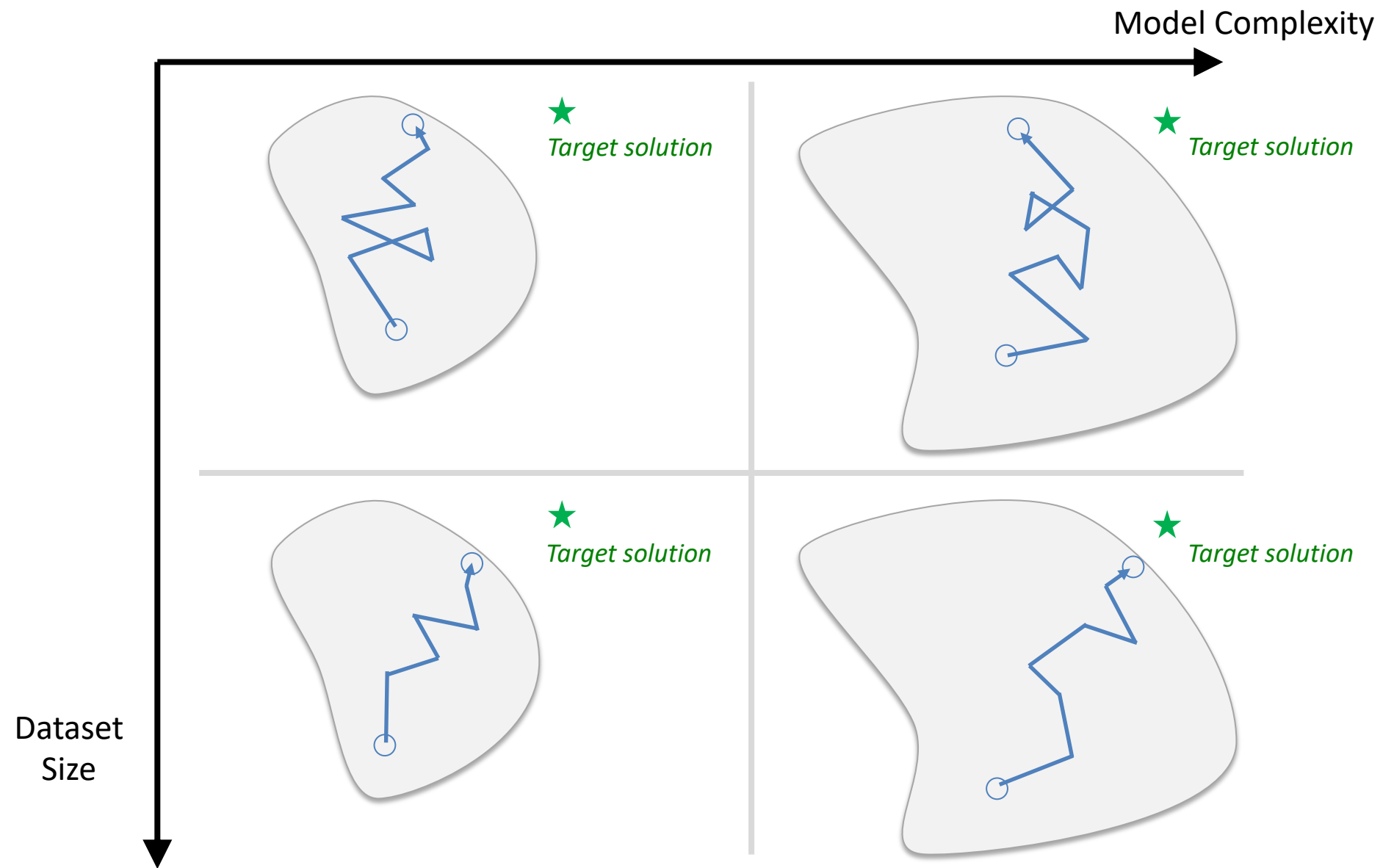
# Hierarchical Learning of Features

Edges     Textures     Patterns     Parts     Objects

*Depth*

2212.06727

Model Complexity
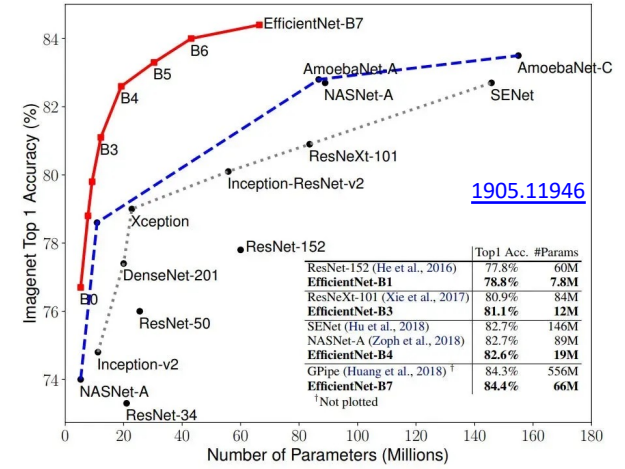
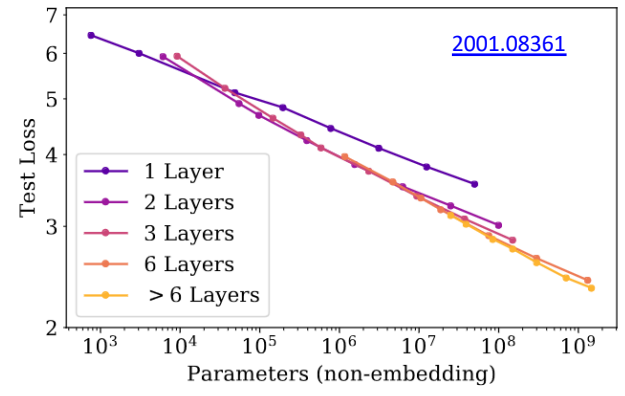Target solution

Target solution

Target solution

Target solution

Dataset Size

Image credit: D. McCandless, T. Evans, P. Barton

MODE CONNECTIVITY
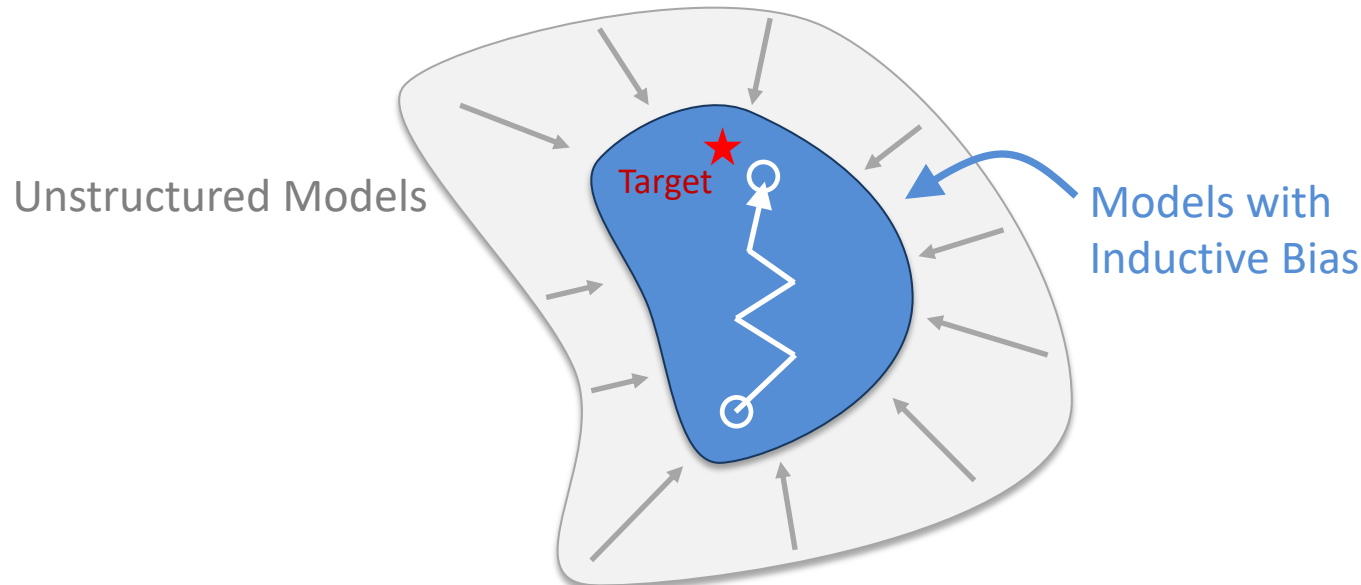
OPTIMA OF COMPLEX LOSS FUNCTIONS CONNECTED BY SIMPLE CURVES OVER WHICH TRAINING AND TEST ACCURACY ARE NEARLY CONSTANT

BASED ON THE PAPER BY TIMUR GARIPOV, PAVEL IZMAILOV, DMITRII PODOPRIKHIN, DMITRY VETROV, ANDREW GORDON WILSON
VISUALIZATION & ANALYSIS IS A COLLABORATION BETWEEN TIMUR GARIPOV, PAVEL IZMAILOV AND JAVIER IDEAMI @LOSSLANDSCAPE.COM

NeurIPS 2018, ARXIV:1802.10026 | LOSSLANDSCAPE.COM

LOSS (TRAIN MODE)

REAL DATA, RESNET-20 NO-SKIP,
CIFAR10, SGD-MOM, BS=128
WD=3e-4, MOM=0.9
BN, TRAIN MOD, 90K PTS
LOG SCALED (ORIG LOSS NUMS)
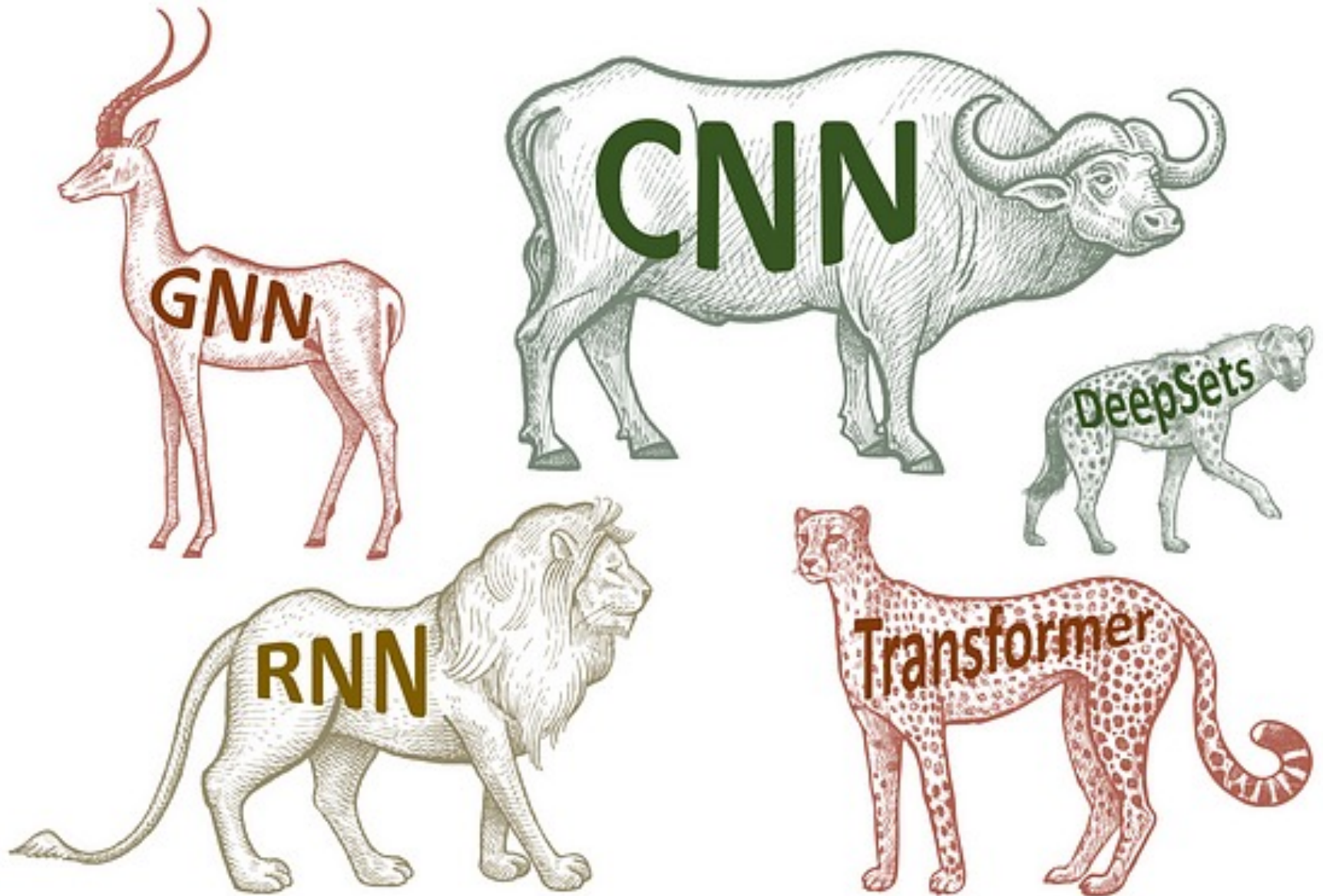
https://arxiv.org/abs/1802.10026

# Choosing the right function...

- We know a lot about our data
  - What transformations shouldn't affect predictions
  - Symmetries, structures, geometry, …

- **Inductive Bias:** we can match models to this knowledge
  - Throw out irrelevant functions we know aren't the solution
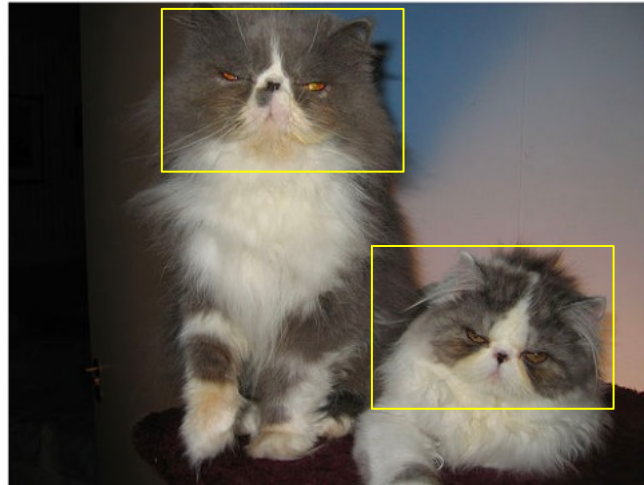  - Bias the learning process towards good solutions
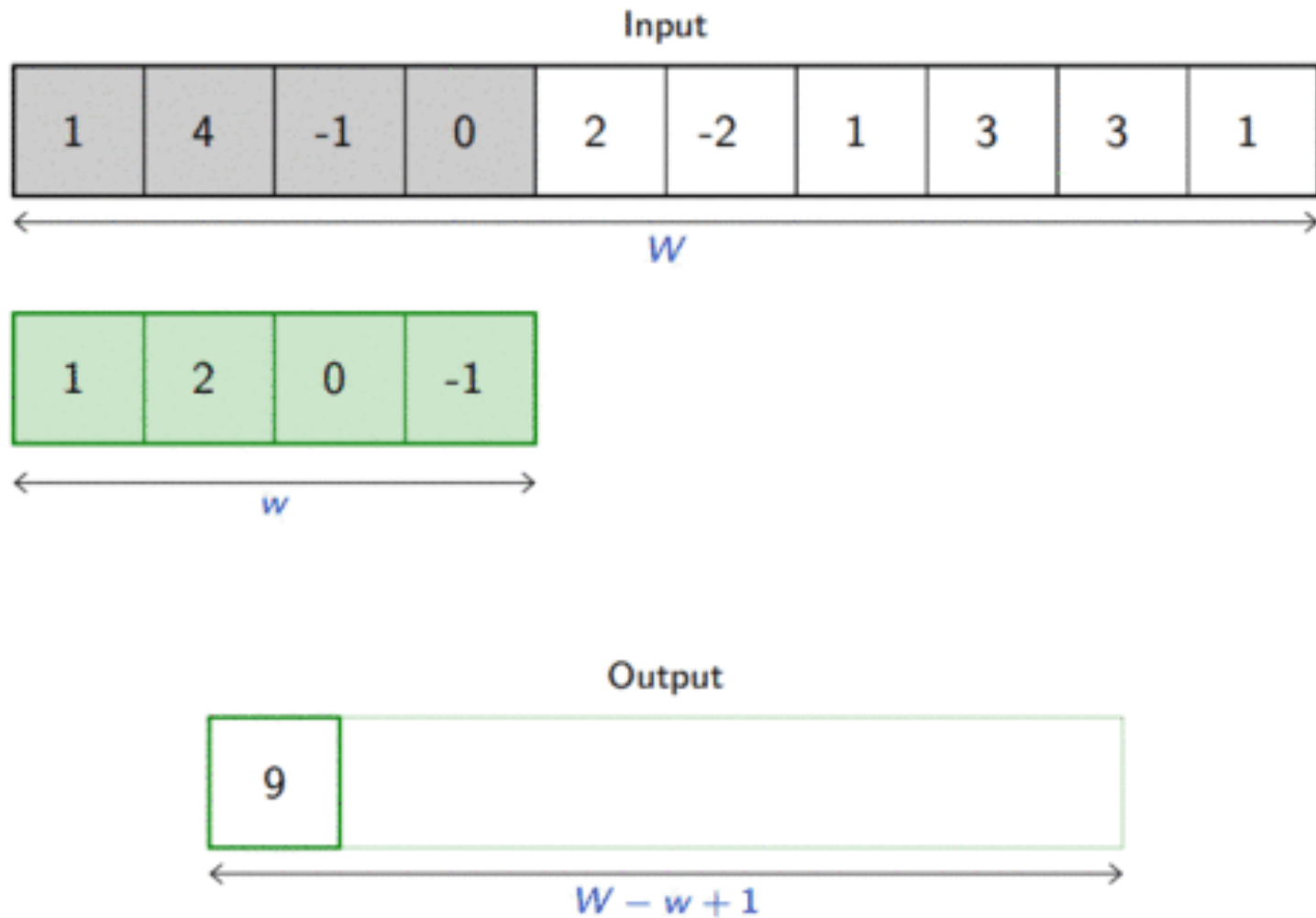
Image credit: Michael Bronstein

- When the structure of data includes "invariance to translation", a representation meaningful at a certain location can / should be used everywhere
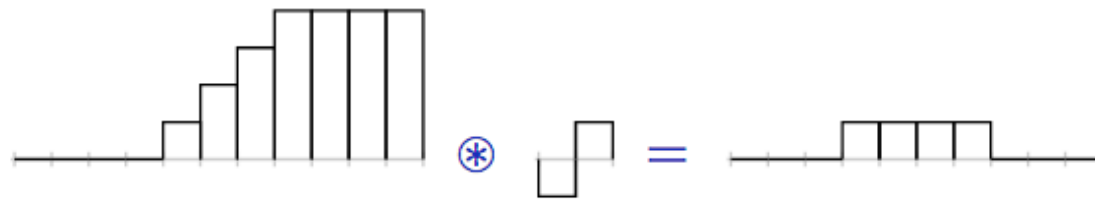


- Convolutional layers build on this idea, that the same "local" transformation is applied everywhere and preserves the signal structure

Fleuret, Deep Learning Course

# 1D Convolutional Layer Example

Input

| 1 | 4 | -1 | 0 | 2 | -2 | 1 | 3 | 3 | 1 |
|---|---|----|---|---|----|---|---|---|---|

$W$

| 1 | 2 | 0 | -1 |
|---|---|---|----|

$w$

Output

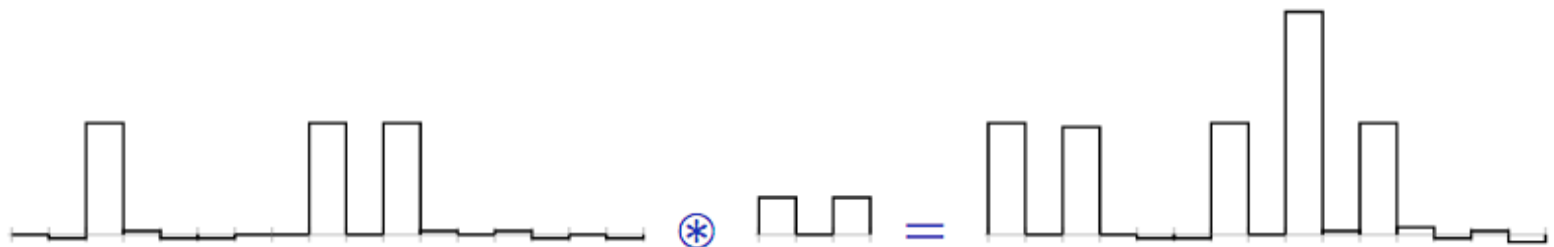| 9 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|

$W - w + 1$

Fleuret, Deep Learning Course

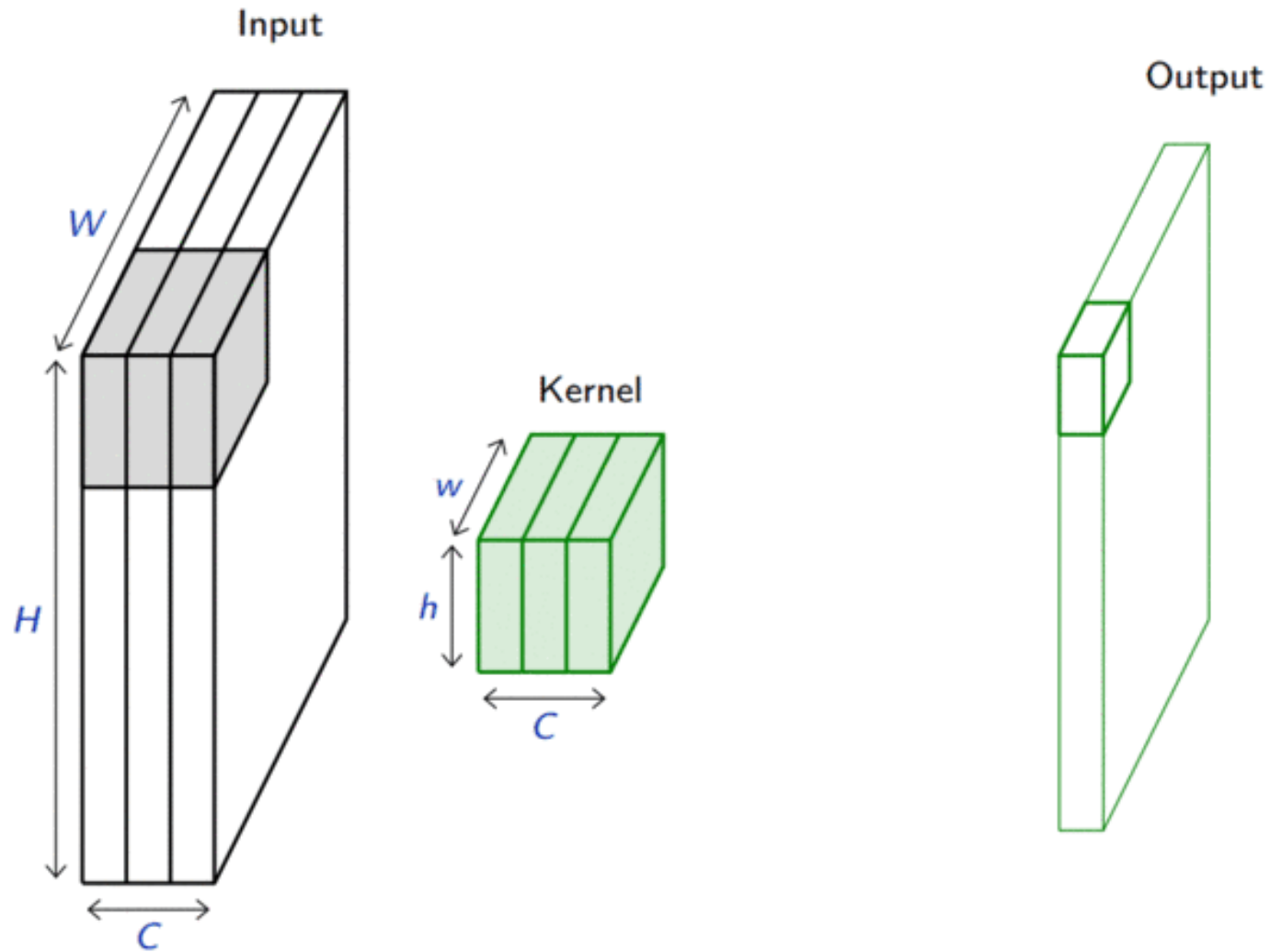Convolution can implement in particular differential operators, *e.g.*

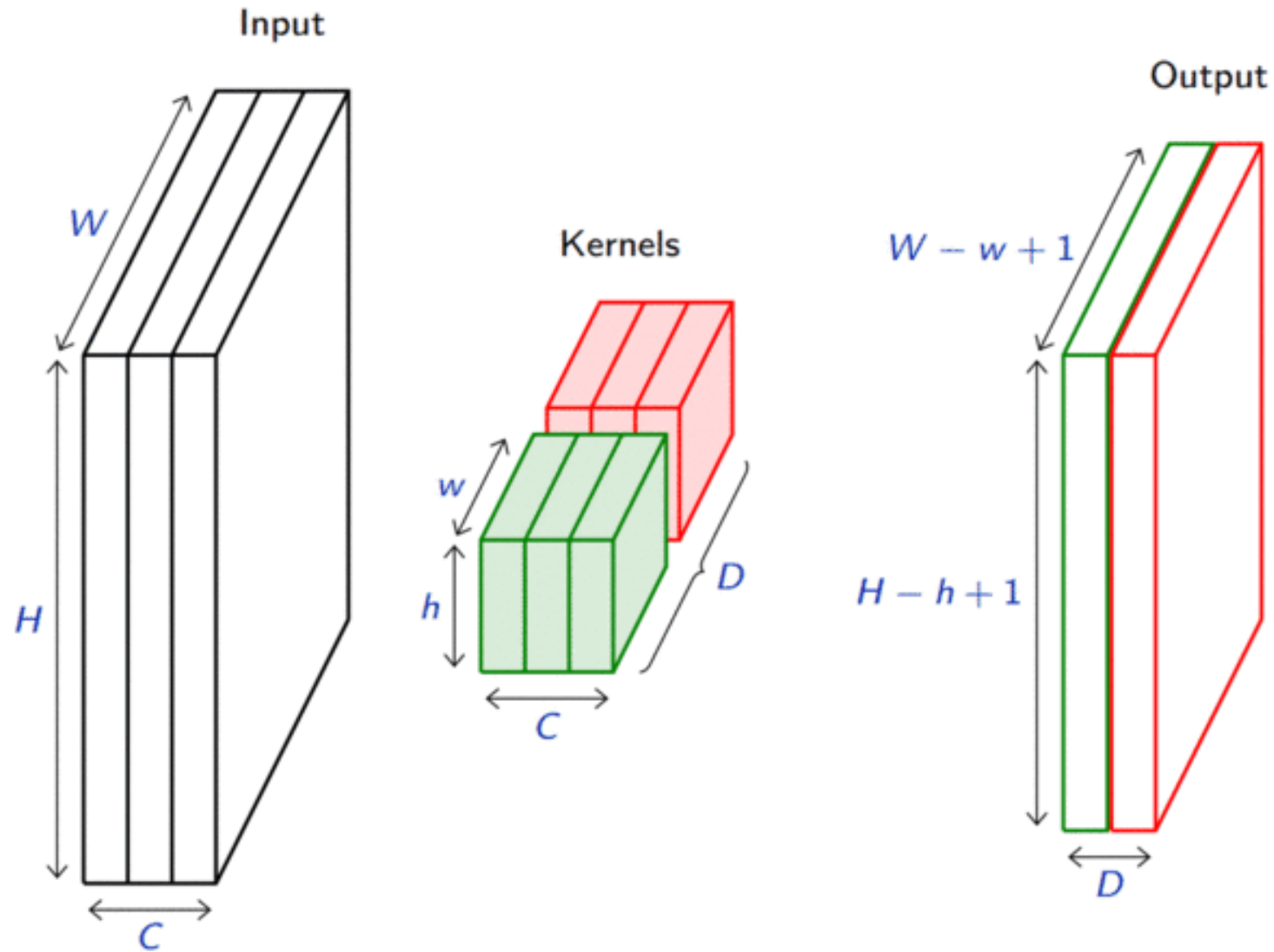$$(0, 0, 0, 0, 1, 2, 3, 4, 4, 4, 4) \circledast (-1, 1) = (0, 0, 0, 1, 1, 1, 1, 0, 0, 0).$$



or crude "template matcher", *e.g.*

Input

Kernel

Output

# 2D Convolution Over Multiple Channels
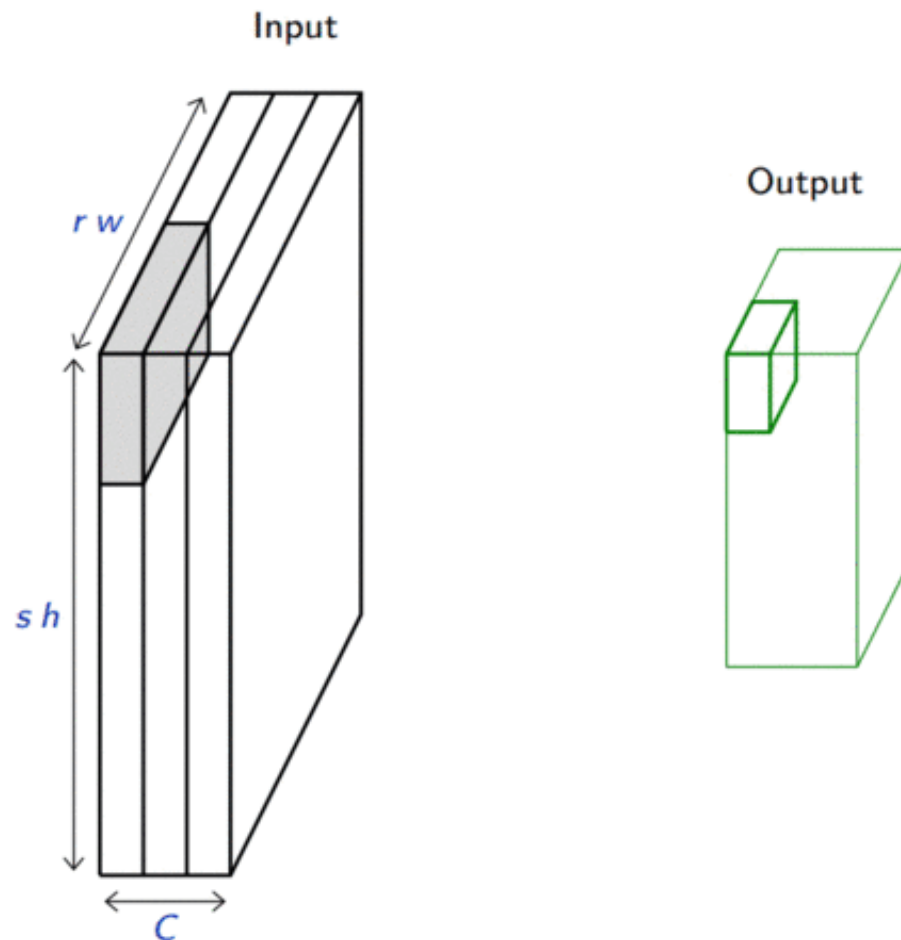
Fleuret, [Deep Learning Course](#)

- Parameters are *shared* by each neuron producing an output in the activation map

- Dramatically reduces number of weights needed to produce an activation map
  - Data: 256×256×3 RGB image
  - Kernel: 3×3×3 → 27  weights
  - Fully connected layer:
    - 256×256×3 inputs → 256×256×3 outputs → $O(10^{10})$ weights

Y. LeCun et. al. 1998

- Parameters are *shared* by each neuron producing an output in the activation map

- Dramatically reduces number of weights needed to produce an activation map

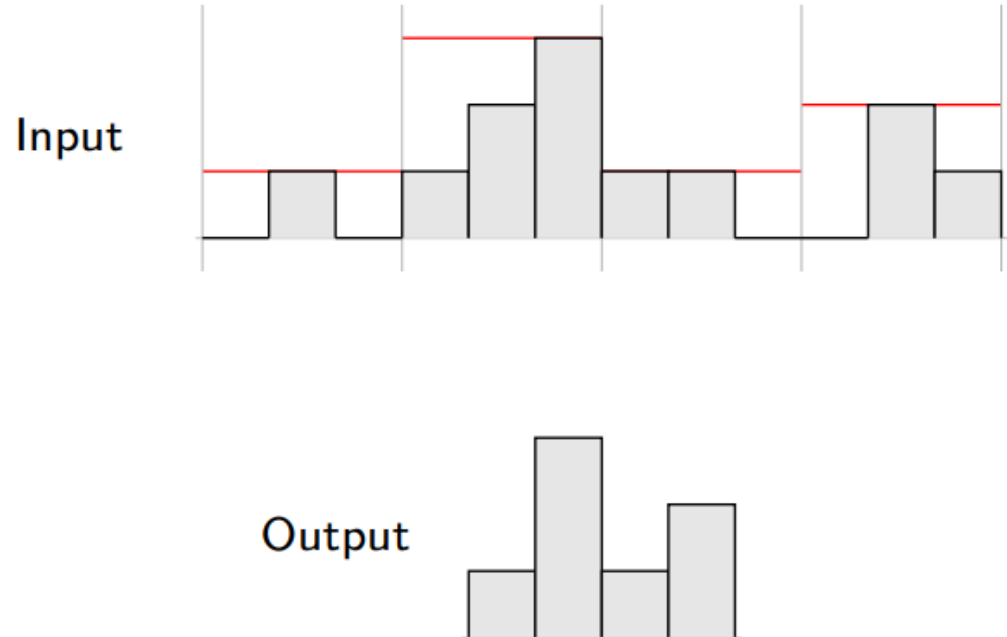- Convolutional layer does pattern matching at any location → Equivariant to translation



Y. LeCun et. al. 1998

- In each channel, find *max* or *average* value of pixels in a pooling area of size $h \times w$
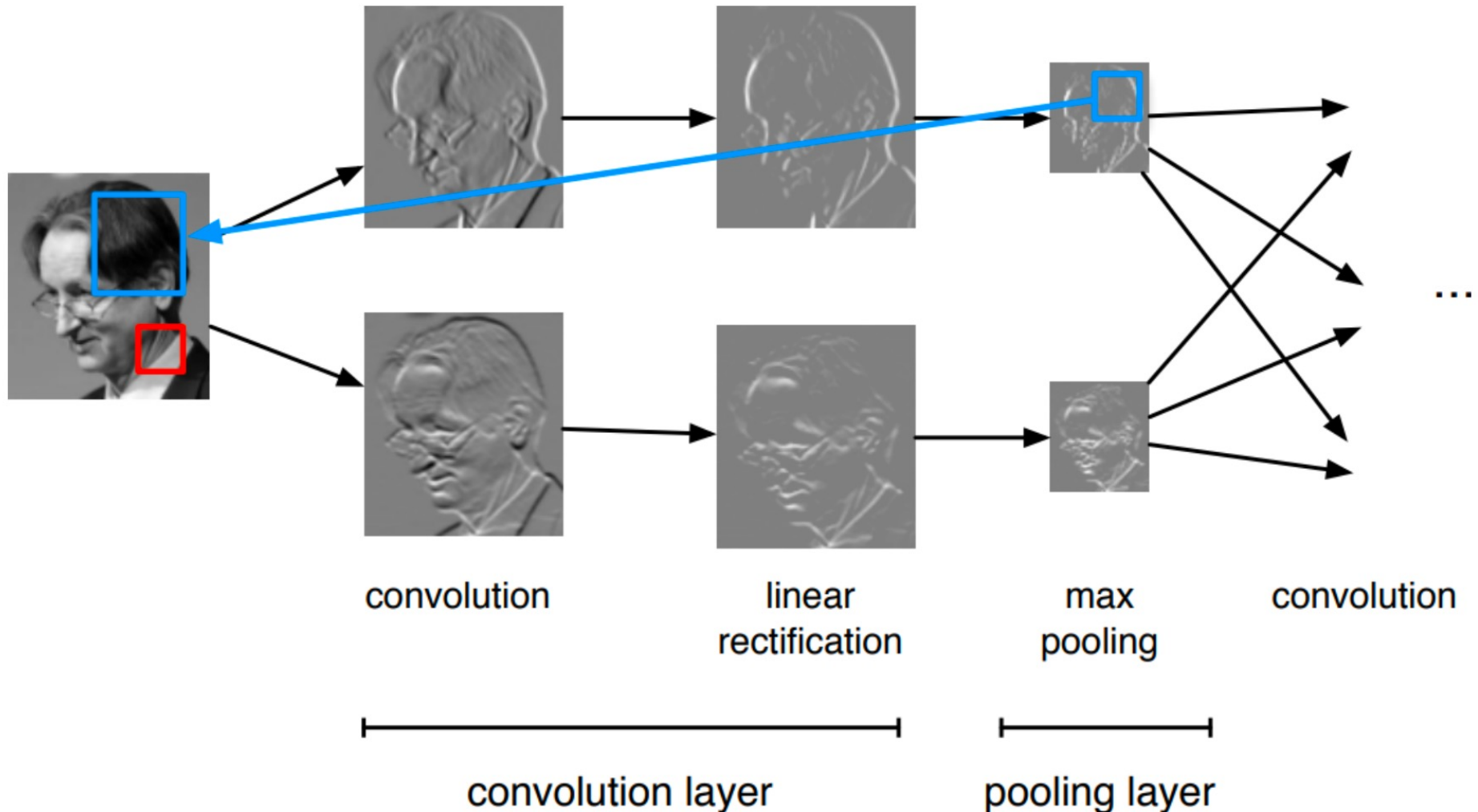
# Pooling

- In each channel, find *max* or *average* value of pixels in a pooling area of size $h \times w$

- Invariance to permutation within pooling area

Input
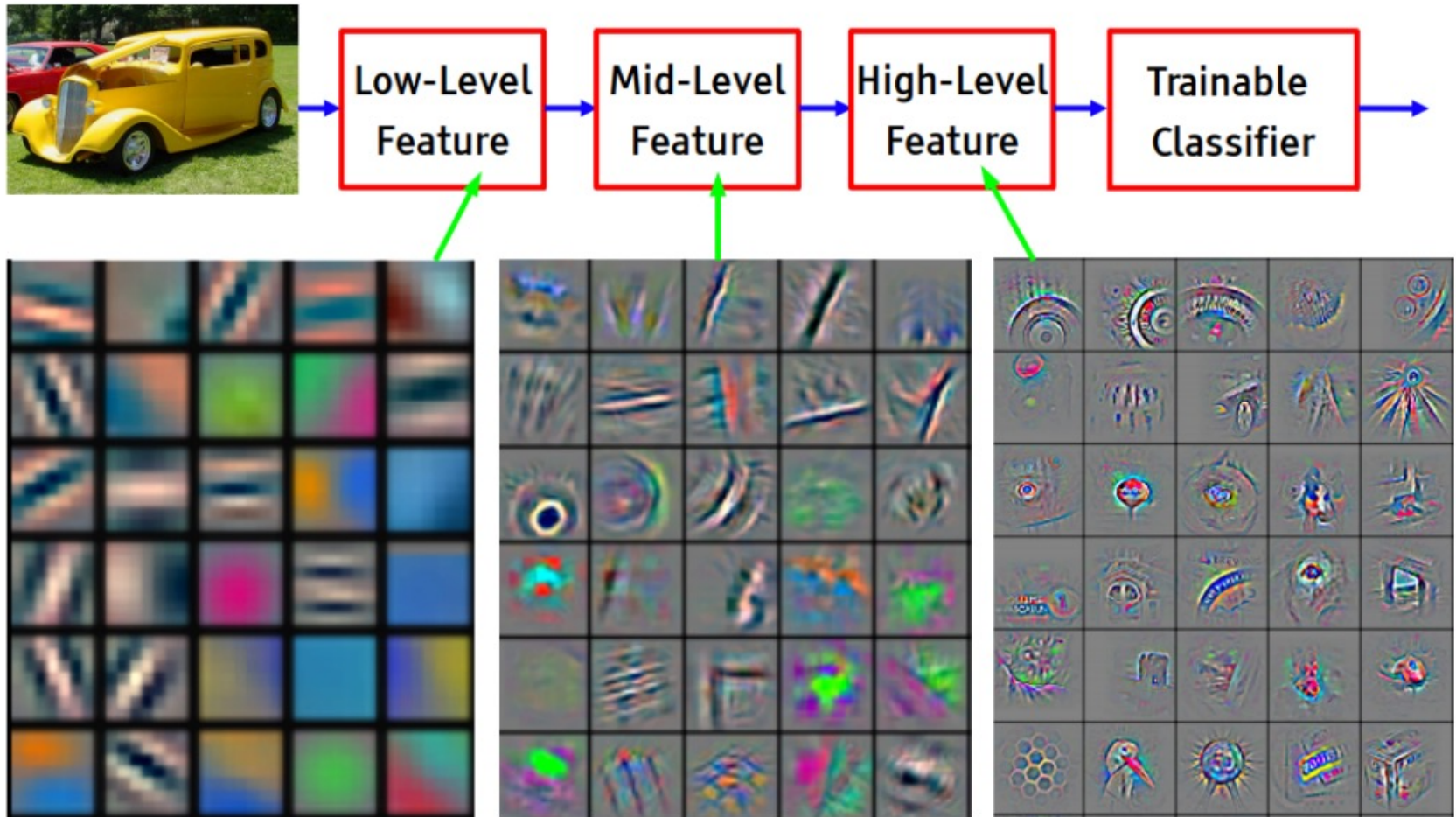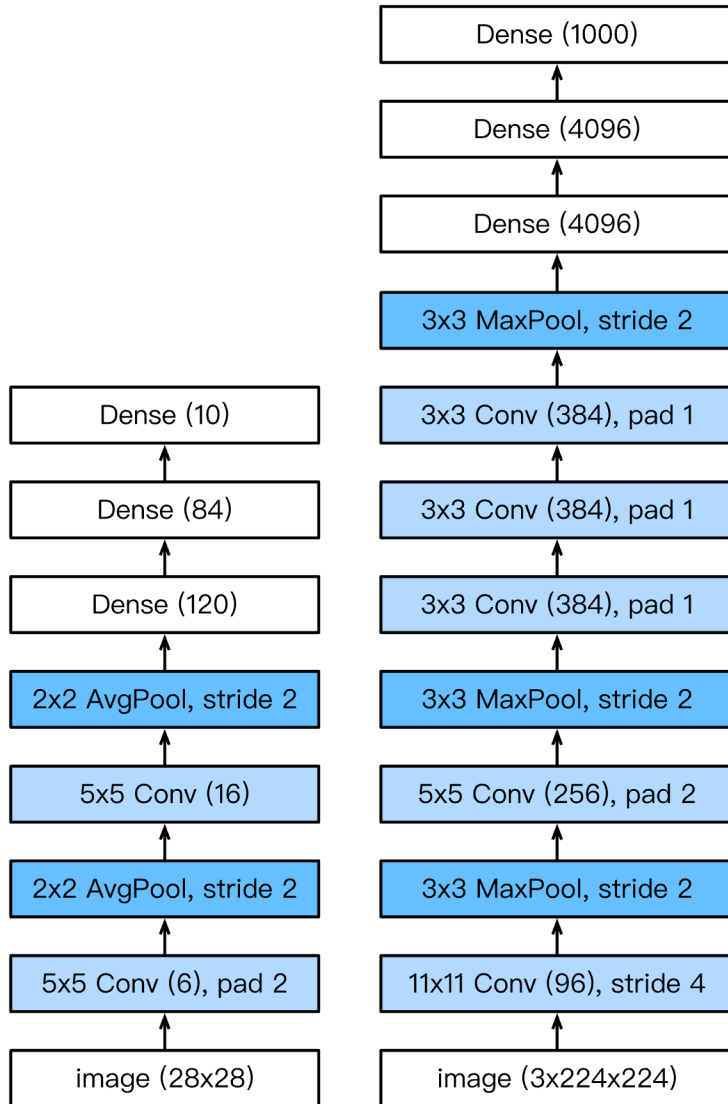
- Invariance to local perturbations

Output

- A combination of convolution, pooling, ReLU, and fully connected layers



convolution      linear rectification      max pooling      convolution

convolution layer      pooling layer

Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

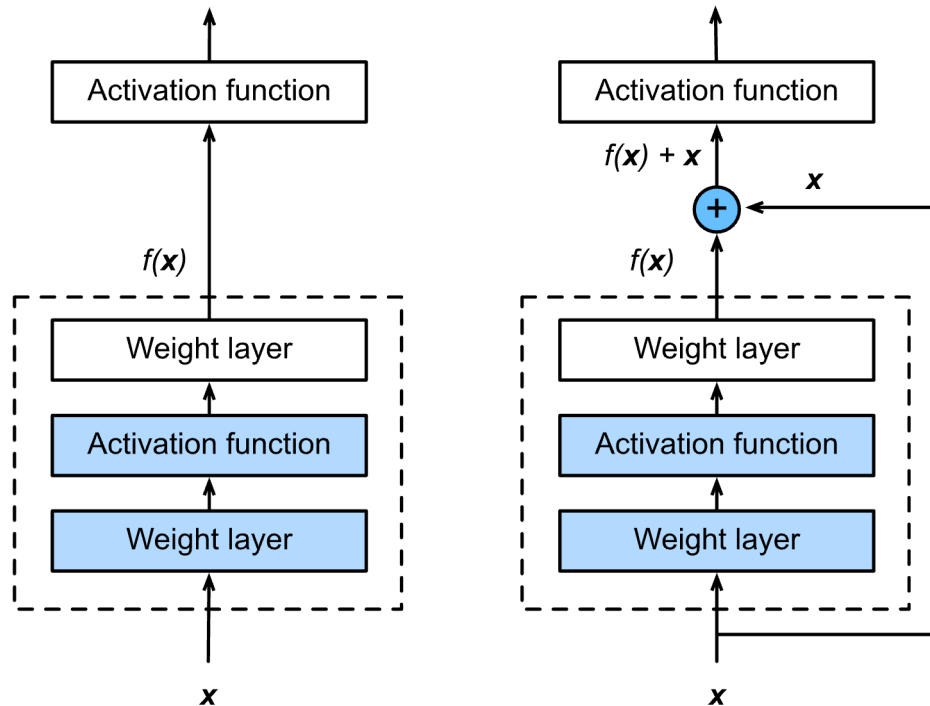# Convolutional Networks
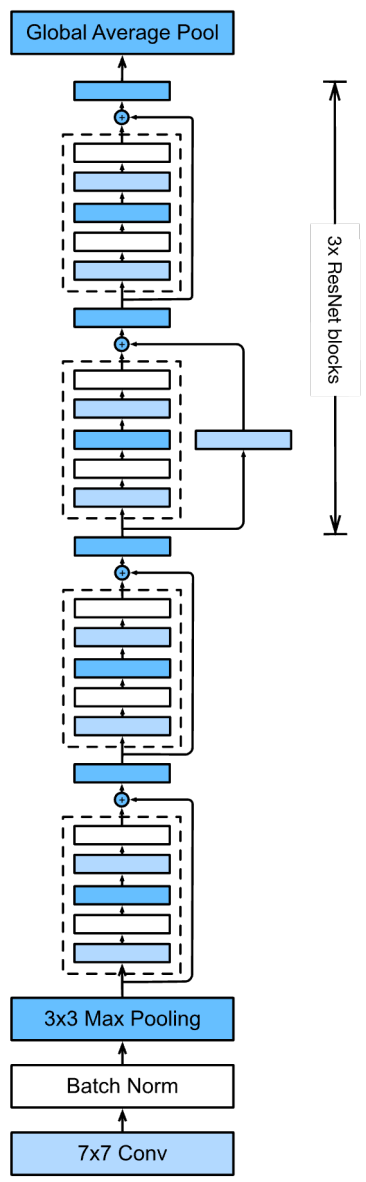
ImageNet Classification

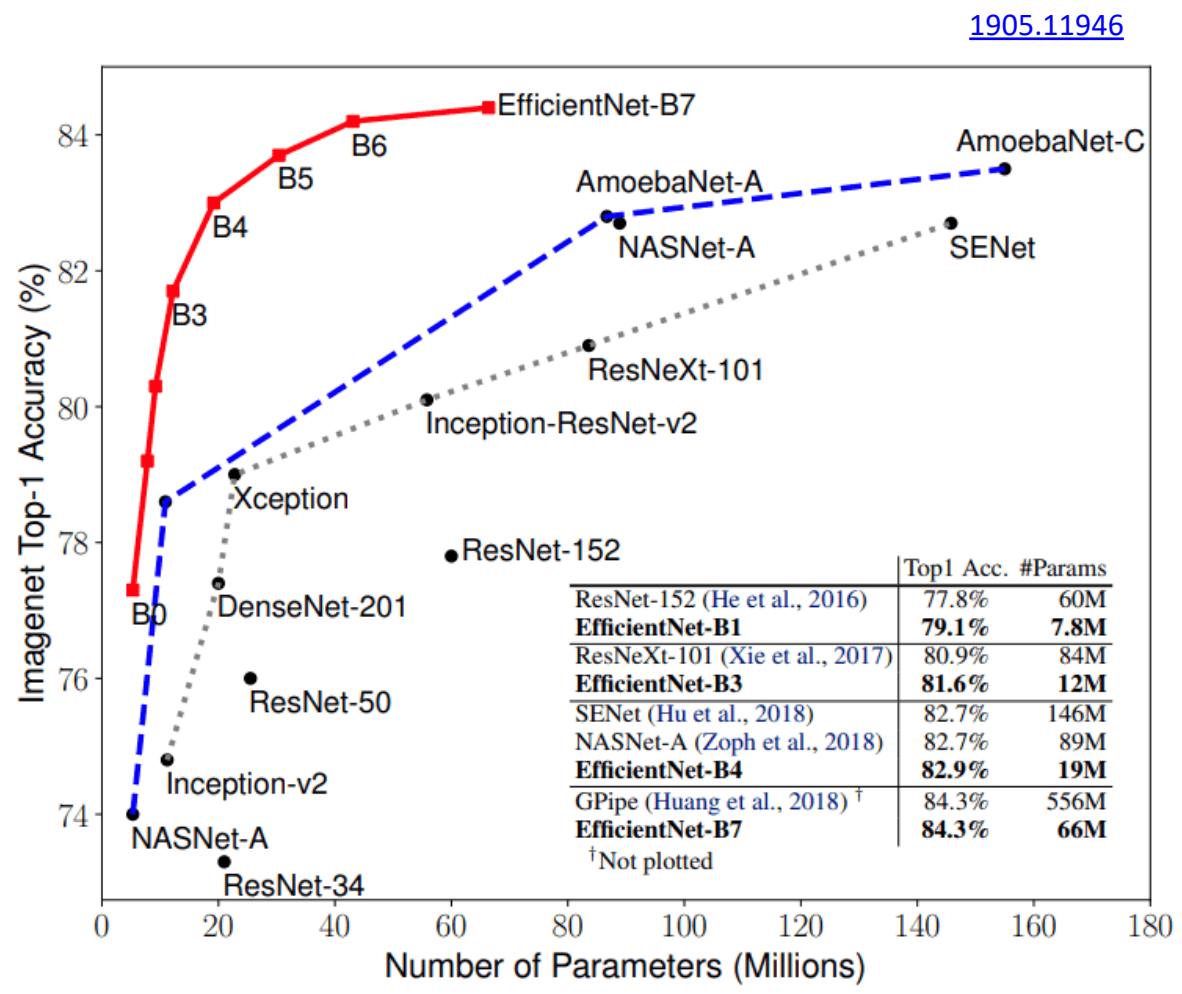**LeNet**
(LeCun et al, 1998)

**AlexNet**
(Krizhevsky et al, 2012)

- Training very deep networks is made possible because of the **skip connections** in the residual blocks. Gradients can shortcut the layers and pass through without vanishing.

# Deep CNNs

ResNet
(He et al, 2015)

[1905.11946](1905.11946)

| | Top1 Acc. | #Params |
|---|---|---|
| ResNet-152 (He et al., 2016) | 77.8% | 60M |
| **EfficientNet-B1** | **79.1%** | **7.8M** |
| ResNeXt-101 (Xie et al., 2017) | 80.9% | 84M |
| **EfficientNet-B3** | **81.6%** | **12M** |
| SENet (Hu et al., 2018) | 82.7% | 146M |
| NASNet-A (Zoph et al., 2018) | 82.7% | 89M |
| **EfficientNet-B4** | **82.9%** | **19M** |
| GPipe (Huang et al., 2018) [†] | 84.3% | 556M |
| **EfficientNet-B7** | **84.3%** | **66M** |

[†] Not plotted

# Sequential Data

- Many types of data are not fixed in size

- Many types of data have a temporal or sequence-like structure
  - Text
  - Video
  - Speech
  - DNA
  - ...

- MLP expects fixed size data

- How to deal with sequences?

# Sequential Data

- Given a set $\mathcal{X}$, let $S(\mathcal{X})$ be the set of sequences, where each element of the sequence $x_i \in \mathcal{X}$
  - $\mathcal{X}$ could reals $\mathbb{R}^M$, integers $\mathbb{Z}^M$, etc.
  - Sample sequence $x = \{x_1, x_2, \dots, x_T\}$

- Tasks related to sequences:
  - Classification $\qquad\qquad f\colon S(\mathcal{X}) \to \{\boldsymbol{p} \mid \sum_{c=1}^{N} p_i = 1\}$
  - Generation $\qquad\qquad\quad f\colon \mathbb{R}^d \to S(\mathcal{X})$
  - Seq.-to-seq. translation $\quad f\colon S(\mathcal{X}) \to S(\mathcal{Y})$

- Input sequence $x \in S(\mathbb{R}^m)$ of *variable* length $T(x)$

- Recurrent model maintain a **recurrent state $\boldsymbol{h}_t \in \mathbb{R}^q$** updated at each time step $t$. For $t = 1, \dots, T(x)$:

$$\boldsymbol{h}_{t+1} = \phi(\boldsymbol{x}_t, \boldsymbol{h}_t; \theta)$$

  – Simplest model:

$$\phi(\boldsymbol{x}_t, \boldsymbol{h}_t; W, U) = \sigma(W\boldsymbol{x}_t + U\boldsymbol{h}_t)$$

- Predictions can be made at any time $t$ from the recurrent state

$$\boldsymbol{y}_t = \psi(\boldsymbol{h}_t; \theta)$$

Credit: F. Fleuret

*Recurrent Model*

$$h_{t+1} = \phi(x_t, h_t; \theta)$$

*Recurrent Model*

$$h_{t+1} = \phi(x_t, h_t; \theta)$$

# Recurrent Neural Networks

*Prediction*

$$\boldsymbol{y}_t = \psi(\boldsymbol{h}_t; \theta)$$



Credit: F. Fleuret

# Recurrent Neural Networks



[0.98] → Positive Sentiment

Sentiment Analysis

The          movie          was          great

Credit: F. Fleuret

Prediction per sequence element



Although the number of steps $T(x)$ depends on $x$, this is a standard computational graph and automatic differentiation can deal with it as usual. This is known as "backpropagation through time" (Werbos, 1988)

Credit: F. Fleuret

- Gating:
  - network can grow very deep, in time → vanishing gradients.
  - *Critical component*: add pass-through (additive paths) so recurrent state does not go repeatedly through squashing non-linearity.
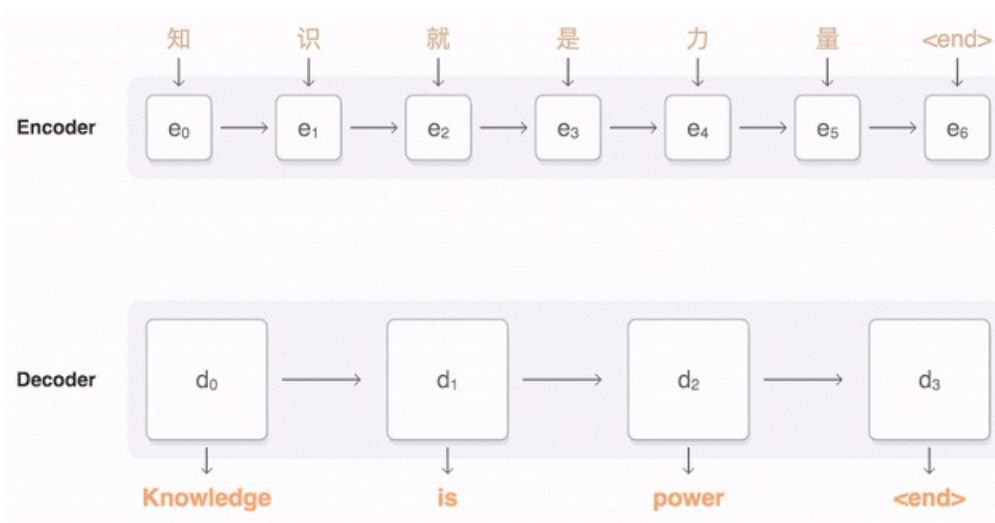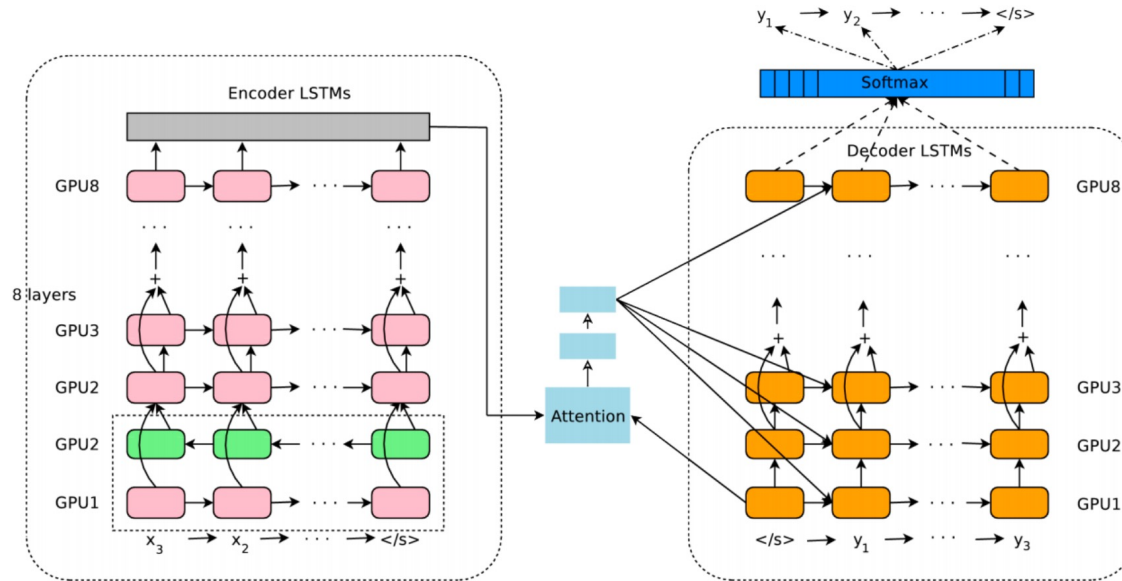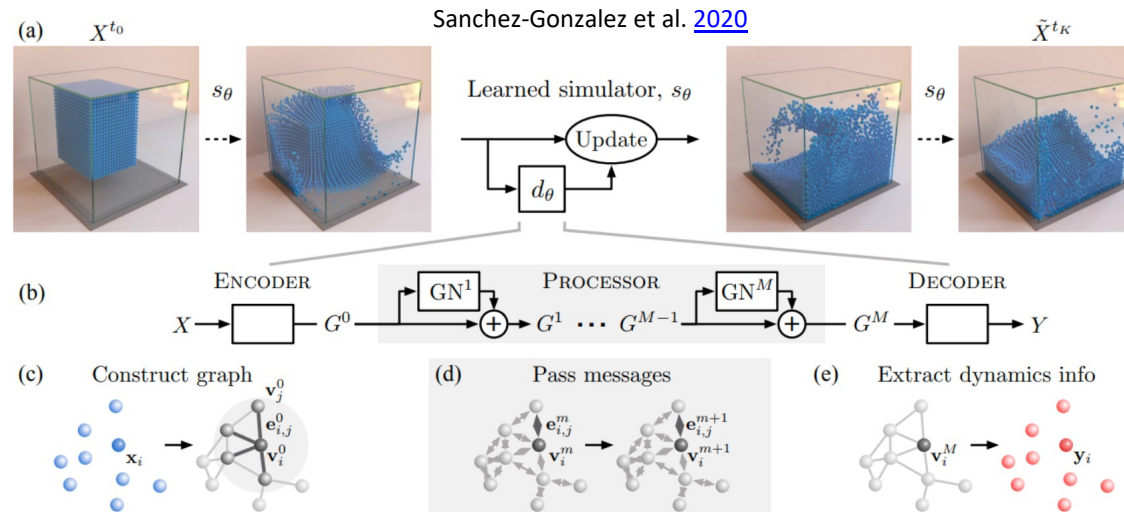
- Gating:

  – network can grow very deep, in time → vanishing gradients.

  – *Critical component*: add pass-through (additive paths) so recurrent state does not go repeatedly through squashing non-linearity.
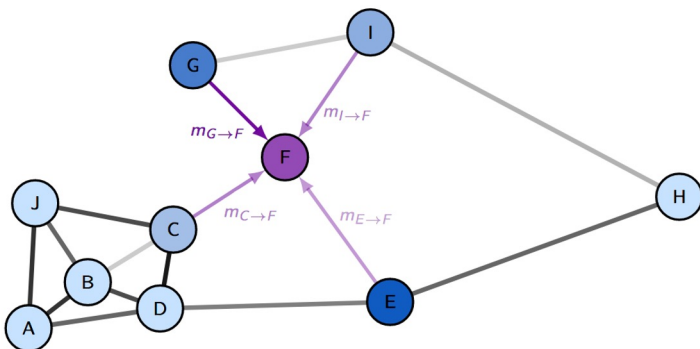
- LSTM:

  – Add internal state separate from output state

  – Add input, output, and forget gating

# Examples

**Neural machine translation**



Y. Wu et al, 2016

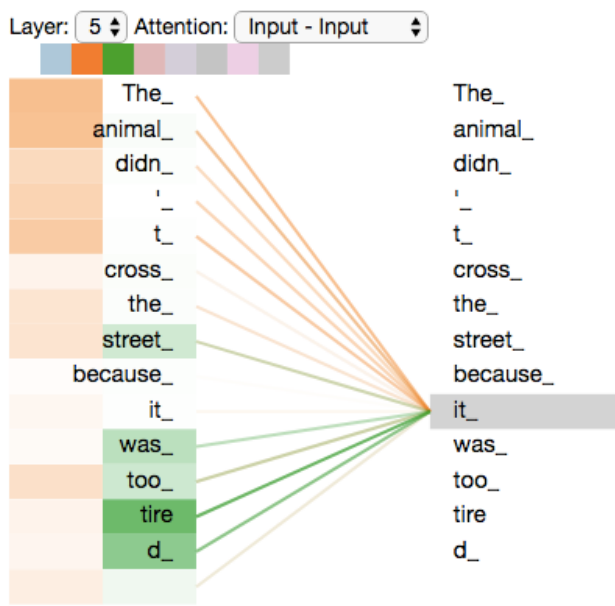- ## Permutation invariant data with geometric relationships
  - Features can be local on graph, but meaningful anywhere on graph

- ## Graph layers can encode these relationships on nodes & edges
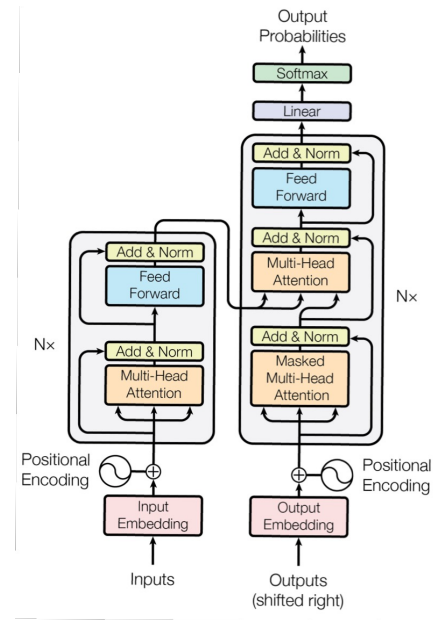


Sanchez-Gonzalez et al. 2020

- **Deep Sets** and **Transformers** can process permutation invariant sets of data

- *Transformers are very adaptable*: Built using layers of *attention*, they can also process sequences, images, and other data



**Attention Is All You Need**

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com
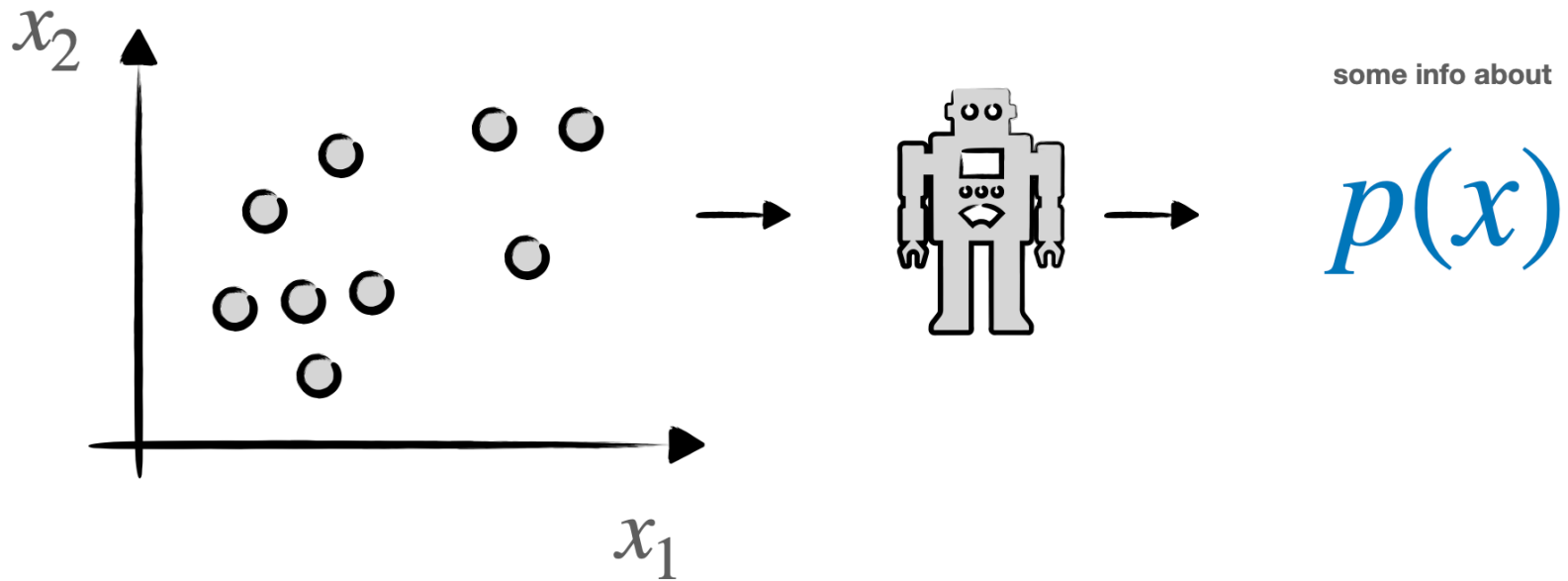
# Beyond Regression and Classification

# Beyond Regression and Classification

- Not all tasks are predicting a label from features, as in classification and regression

- May want to model a high-dim. signal
  - Data synthesis / simulation
  - Density estimation
  - Anomaly detection
  - Denoising, super resolution
  - Data compression
  - …

- Often don't have labels → Unsupervised Learning

- Our goal is to study the data density $p(x)$

- Even w/o labels, aim to characterize the distribution

# Probability Models

**A process**

$$\text{🎲} \rightarrow \mathbb{R}^2$$

**Generating new samples
from randomness**

**A formula**

$$\mathbb{R}^2 \rightarrow \mathbb{R}$$

$$p_{\mu,\Sigma}(x) = \frac{1}{\sqrt{|2\pi\Sigma|}} exp\left(-\frac{1}{2}(x-\mu)^T\Sigma(x-\mu)\right)$$
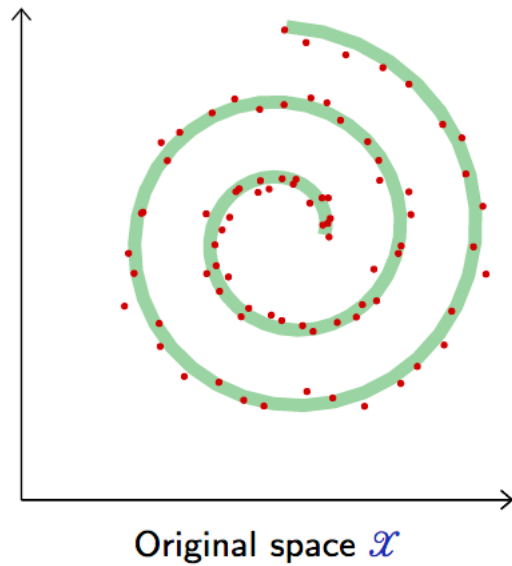
**Evaluating the Probability
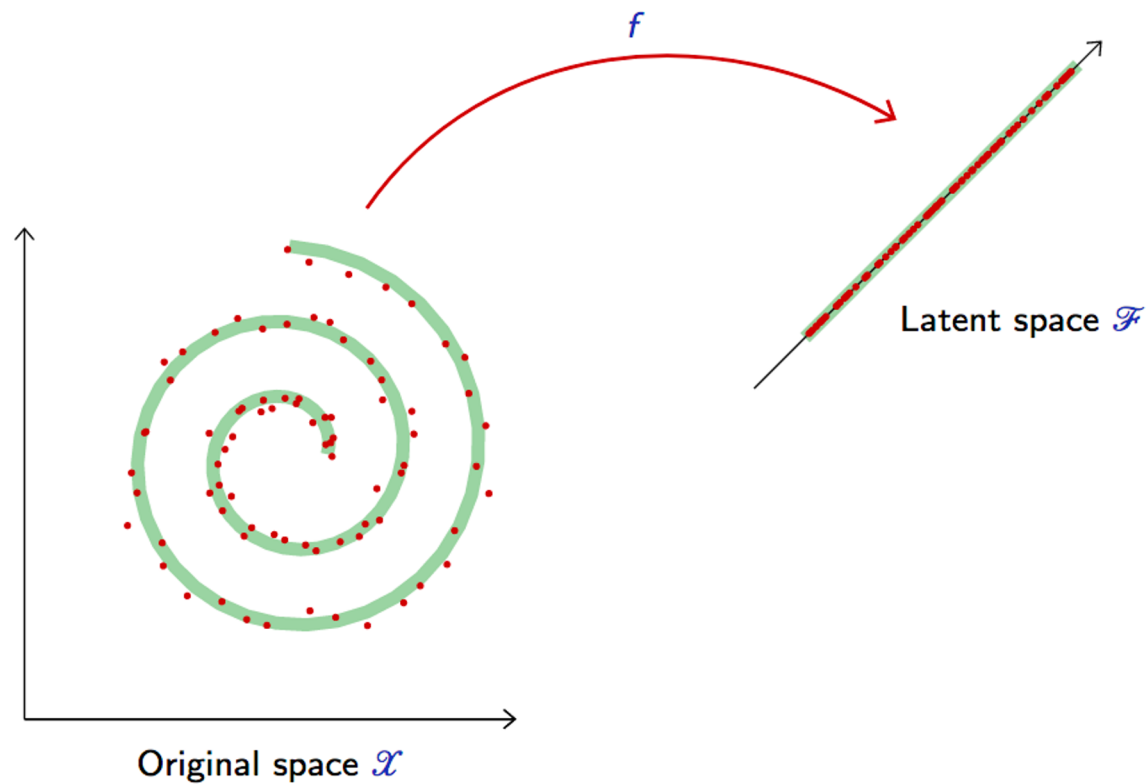for a given sample**

"Understanding $p(x)$" – ability to do either or both of these

Image credit: L. Heinrich

- In many cases, we don't have a theory of the underlying process → *Can still learn to sample*

- Deep learning can be very good at this!

$$\text{face} \sim p(\text{face})$$
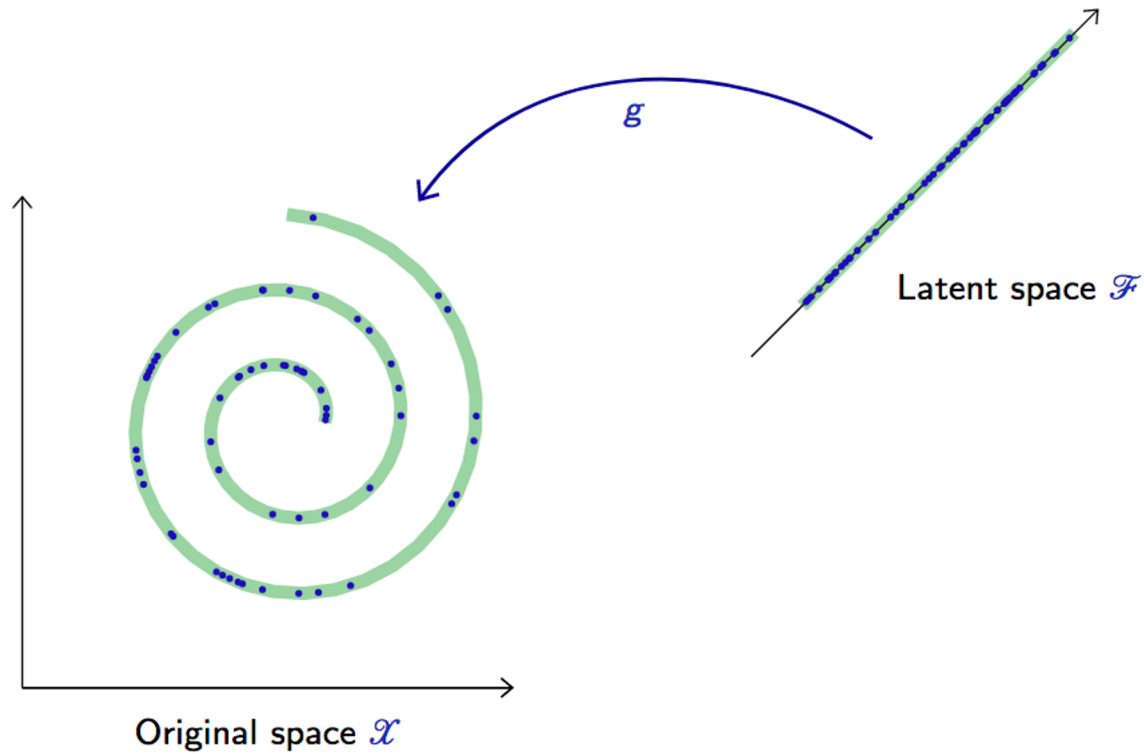


https://thispersondoesnotexist.com/

- Unsupervised learning is more heterogeneous than supervised learning

- Many architectures, losses, learning strategies

- Often constructed so model converges to $p(x)$
  - Variational inference, Adversarial learning, Self-supervision, …

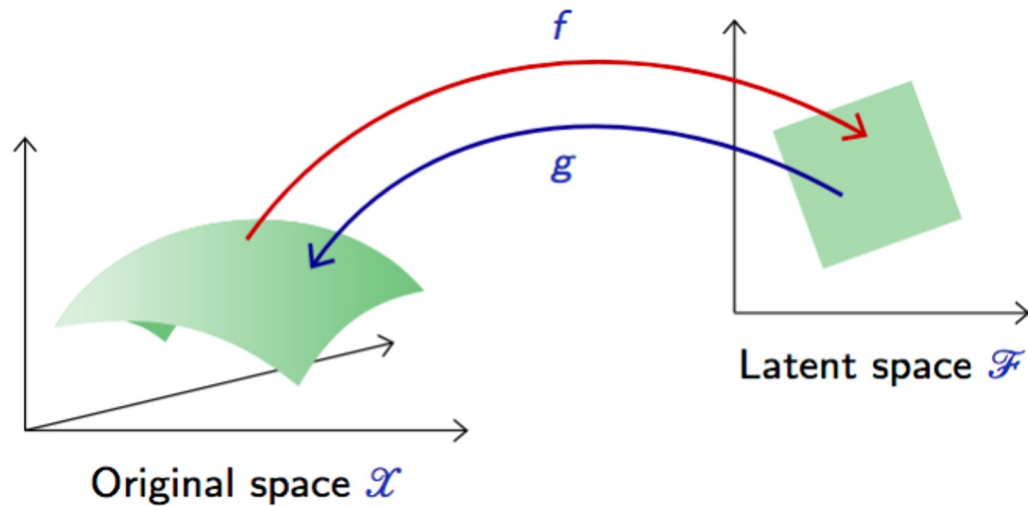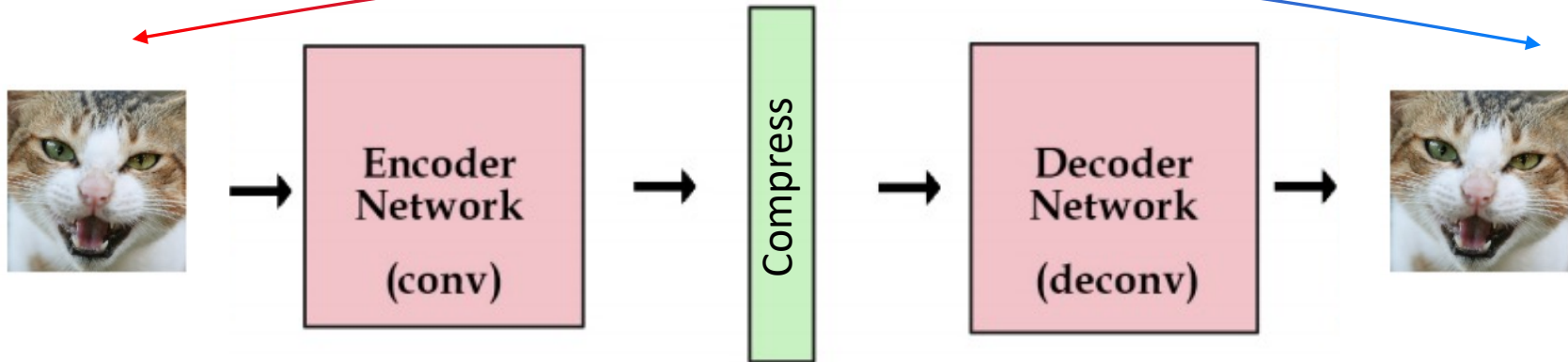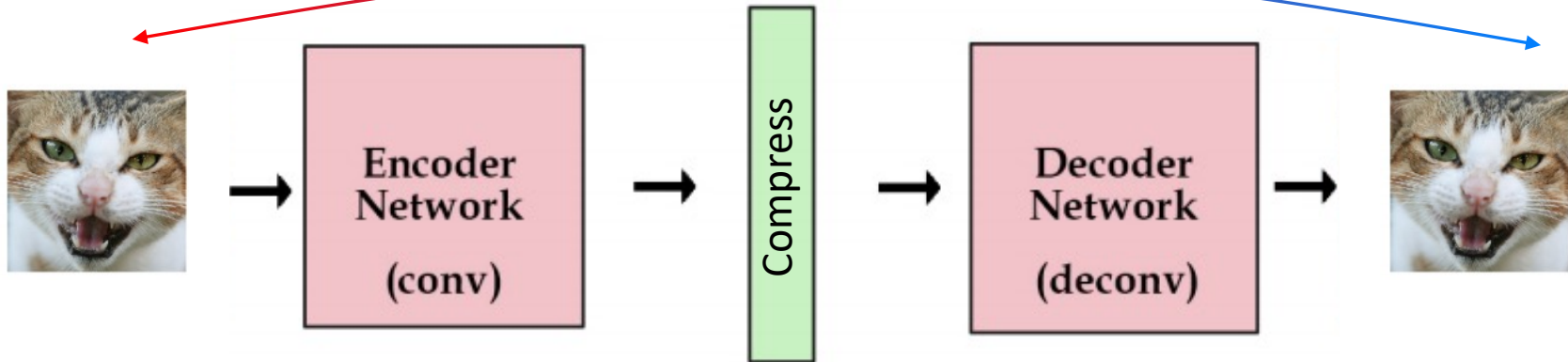- Often framed as **modeling the lower dimensional "meaningful degrees of freedom"** that describe the data

Original space $\mathcal{X}$

$f$

Latent space $\mathscr{F}$

Original space $\mathscr{X}$

Original space $\mathcal{X}$

Latent space $\mathcal{F}$

Fleuret, [Deep Learning Course](#)

# AutoEncoders

$$L(\theta, \psi) = \frac{1}{N} \sum_{n} \left\| x_n - g_\psi\big(f_\theta(x_n)\big) \right\|^2$$

$X$ (original samples)      $g \circ f(X)$ (CNN, $d = 16$)

# Can We Generate Data with Decoder?

- Can we sample in latent space and decode to generate data?



Original space $\mathcal{X}$

Latent space $\mathcal{F}$

- What distribution to sample from in latent space?
  - Try Gaussian with mean and variance from data



Autoencoder sampling ($d = 16$)

- Don't know the right latent space density

# Generative Models Goal

A <span style="color:red">generative model</span> is a probabilistic model $q$ that can be <span style="color:blue">used as a simulator of the data</span>.

**Goal**: generate synthetic, realistic high-dimensional data

$$x \sim q(x; \theta)$$

that is <span style="color:green">as close as possible to the unknown data distribution $p(x)$</span> for which we have empirical samples.

i.e. want to recreate the raw data distribution (such as the distribution of natural images).

**Prompt:**

*street style photo of a woman selling pho at a Vietnamese street market, sunset, shot on fujifilm*

Design of new molecules with desired chemical properties.
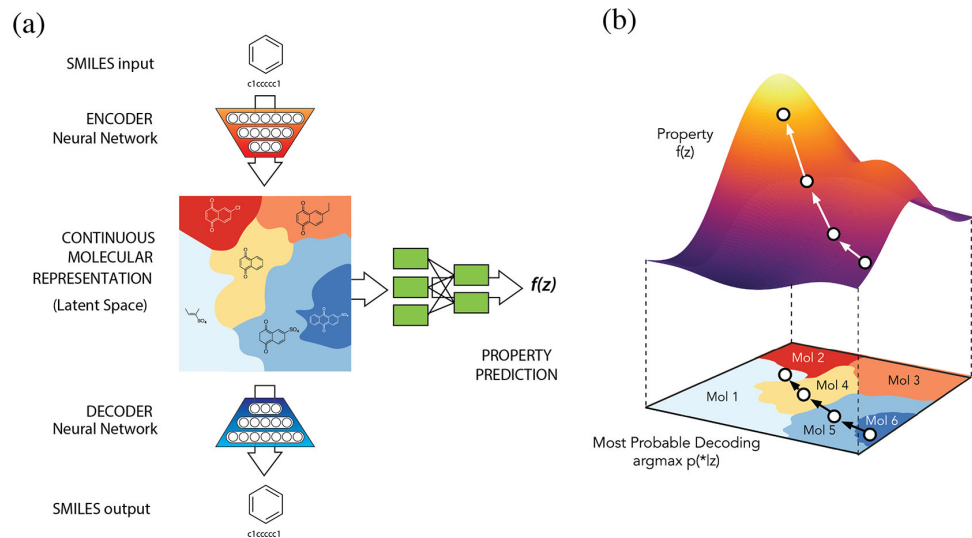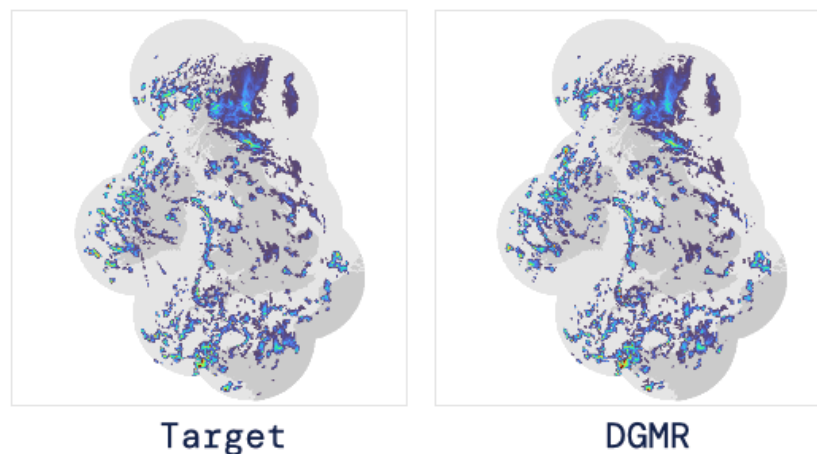(Gomez-Bombarelli et al, 2016)



**Figure 1**: Samples from N-body simulation and from GAN for the box size of 500 Mpc. Note that the transformation in Equation 3.1 with $a = 20$ was applied to the images shown above for better clarity.
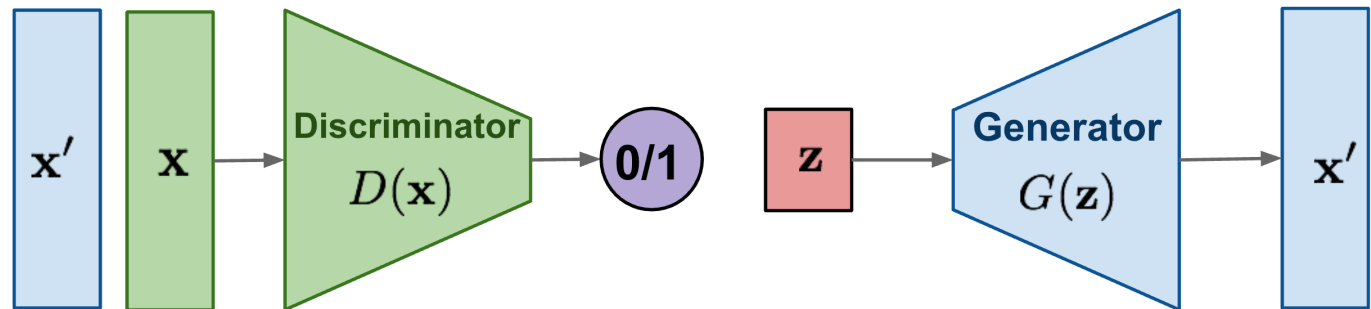
Learning cosmological models
(Rodriguez et al, 2018)

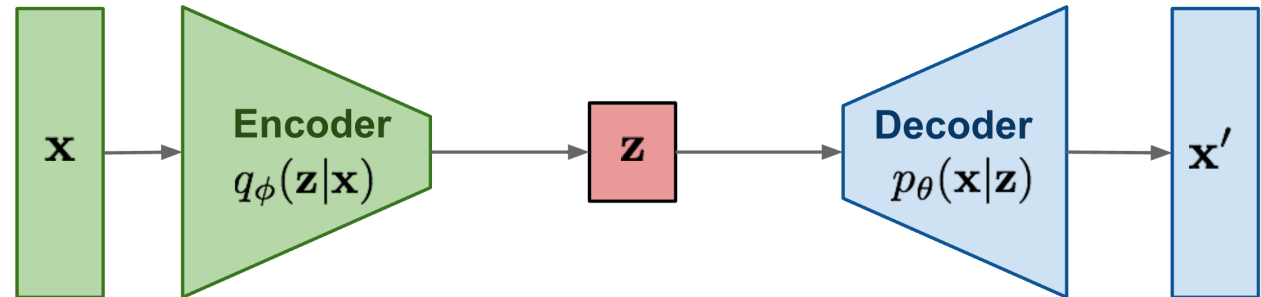

Deep Generative Model of Rainfall (Ravuri et. al. 2021)

# What Deep Generative Models Are There?

**GAN:** Adversarial training
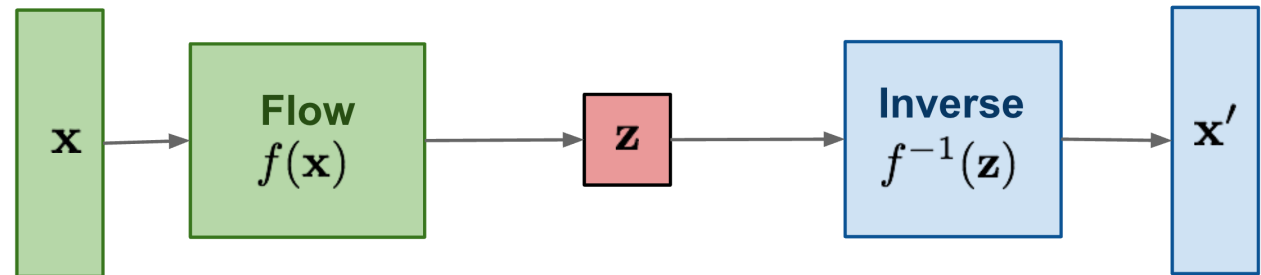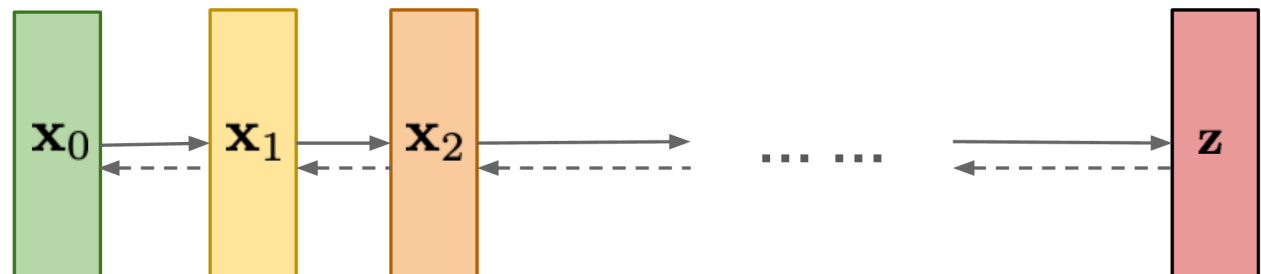
$\mathbf{x}'$ $\mathbf{x}$ → **Discriminator** $D(\mathbf{x})$ → 0/1 $\mathbf{z}$ → **Generator** $G(\mathbf{z})$ → $\mathbf{x}'$

**VAE:** maximize variational lower bound

$\mathbf{x}$ → **Encoder** $q_\phi(\mathbf{z}|\mathbf{x})$ → $\mathbf{z}$ → **Decoder** $p_\theta(\mathbf{x}|\mathbf{z})$ → $\mathbf{x}'$

**Flow-based models:** Invertible transform of distributions

$\mathbf{x}$ → **Flow** $f(\mathbf{x})$ → $\mathbf{z}$ → **Inverse** $f^{-1}(\mathbf{z})$ → $\mathbf{x}'$

**Diffusion models:** Gradually add Gaussian noise and then reverse

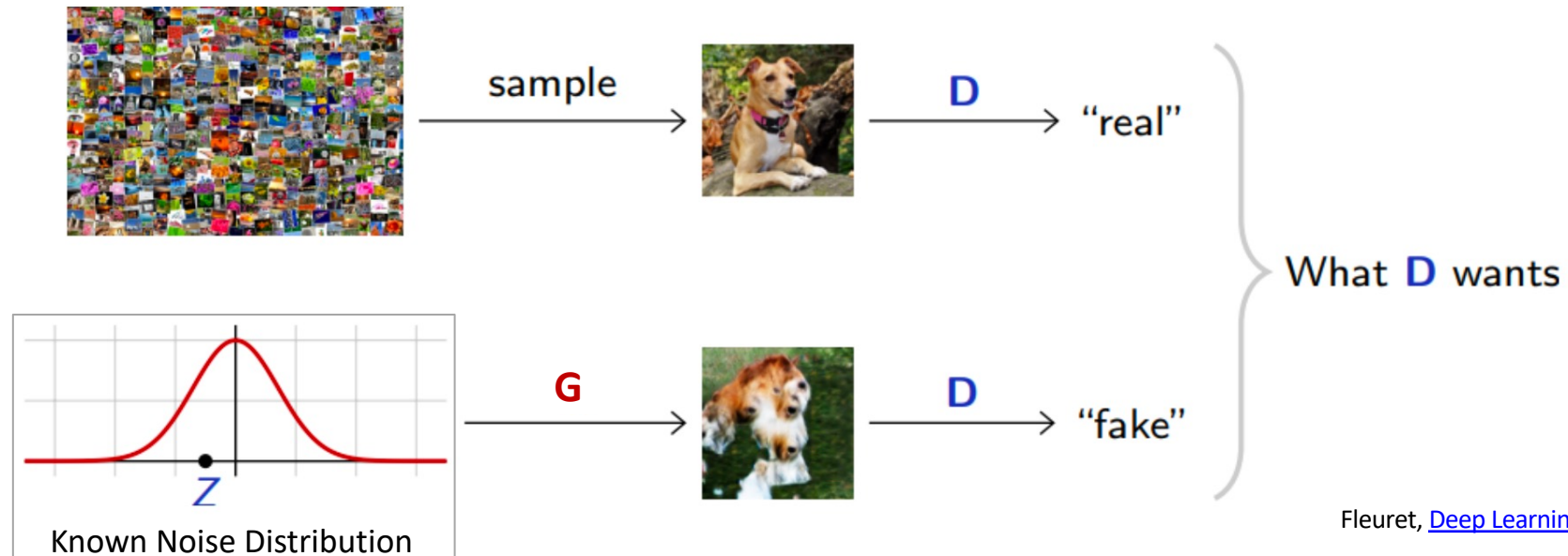$\mathbf{x}_0$ ⇄ $\mathbf{x}_1$ ⇄ $\mathbf{x}_2$ ⋯ ⋯ ⇄ $\mathbf{z}$

Image from Lilian Weng

# Variational AutoEncoders

Choose known distribution for latent space and learn map to data space

# Generative Adversarial Networks (GAN)

Known Noise Distribution

Fleuret, Deep Learning Course

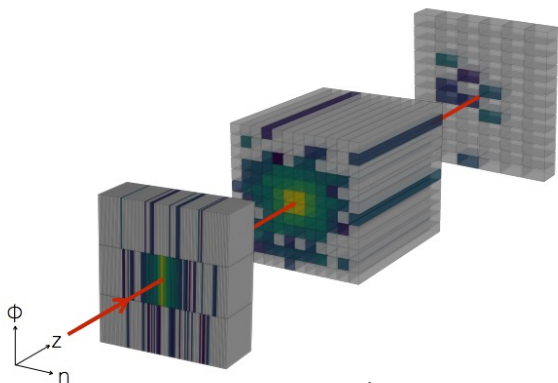- **Generator** creates data from noise, trained to trick **Discriminator** that classifies data as real or fake



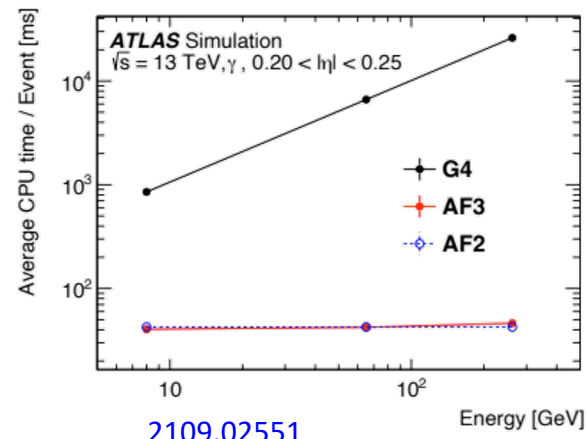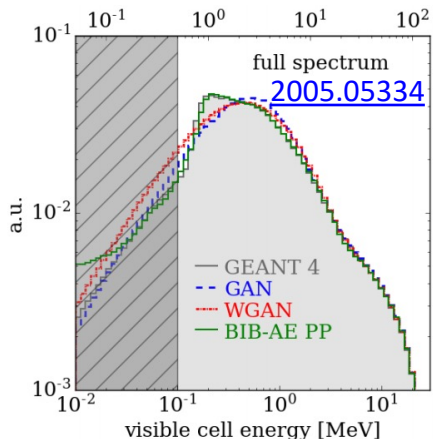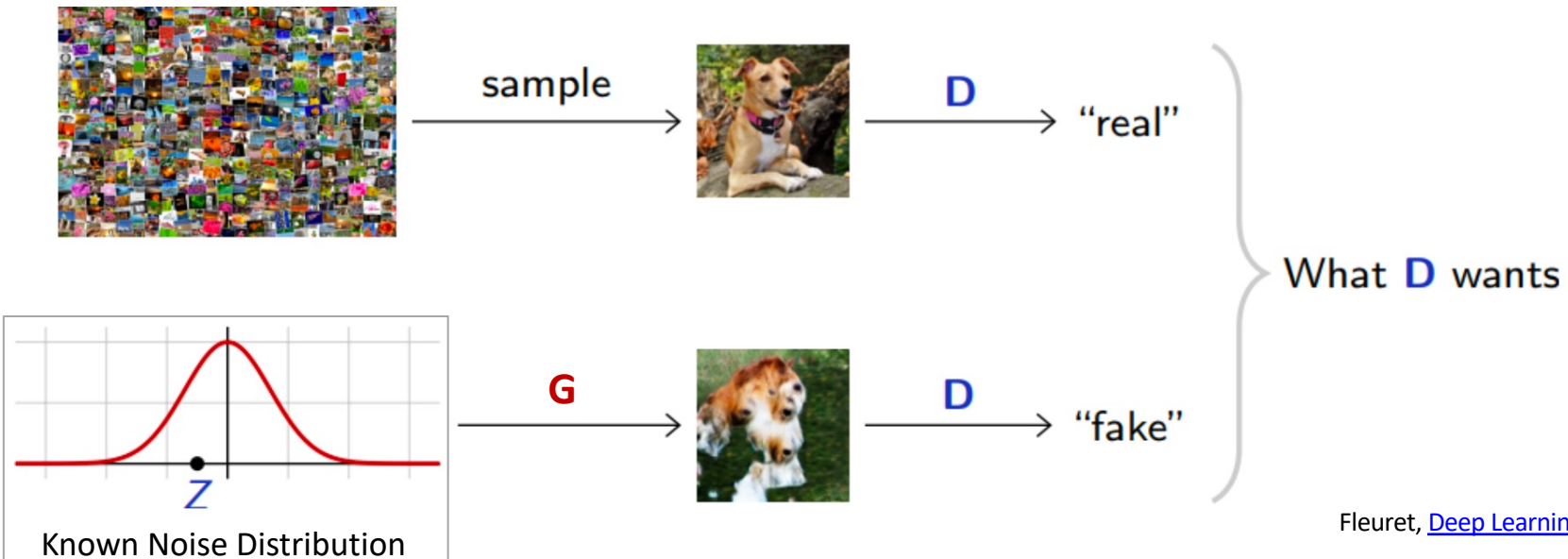Image credit: 1705.02355



2005.05334



2109.02551

# Generative Adversarial Networks (GAN)

sample → D → "real"

What **D** wants

**G** → D → "fake"

*Z*

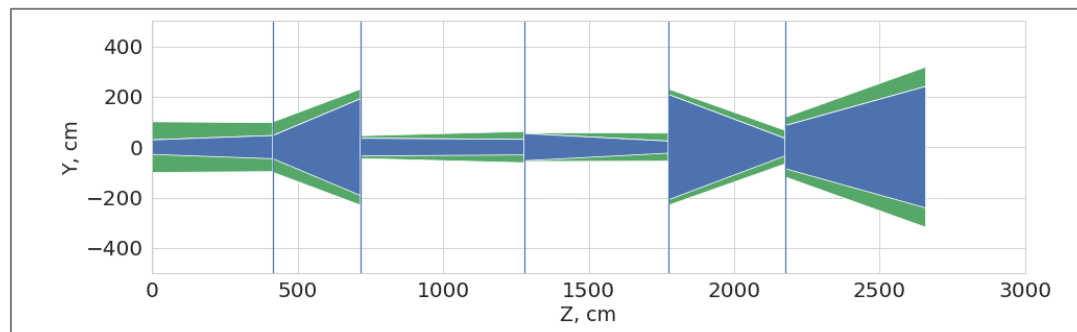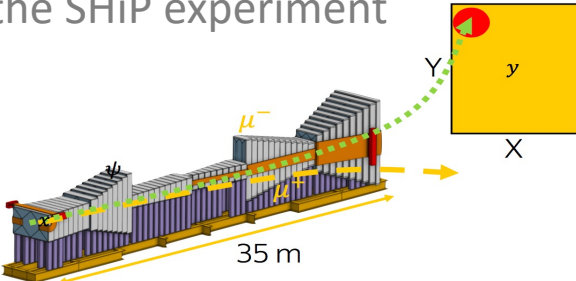Known Noise Distribution

Fleuret, Deep Learning Course

- **Generator** creates data from noise, trained to trick **Discriminator** that classifies data as real or fake

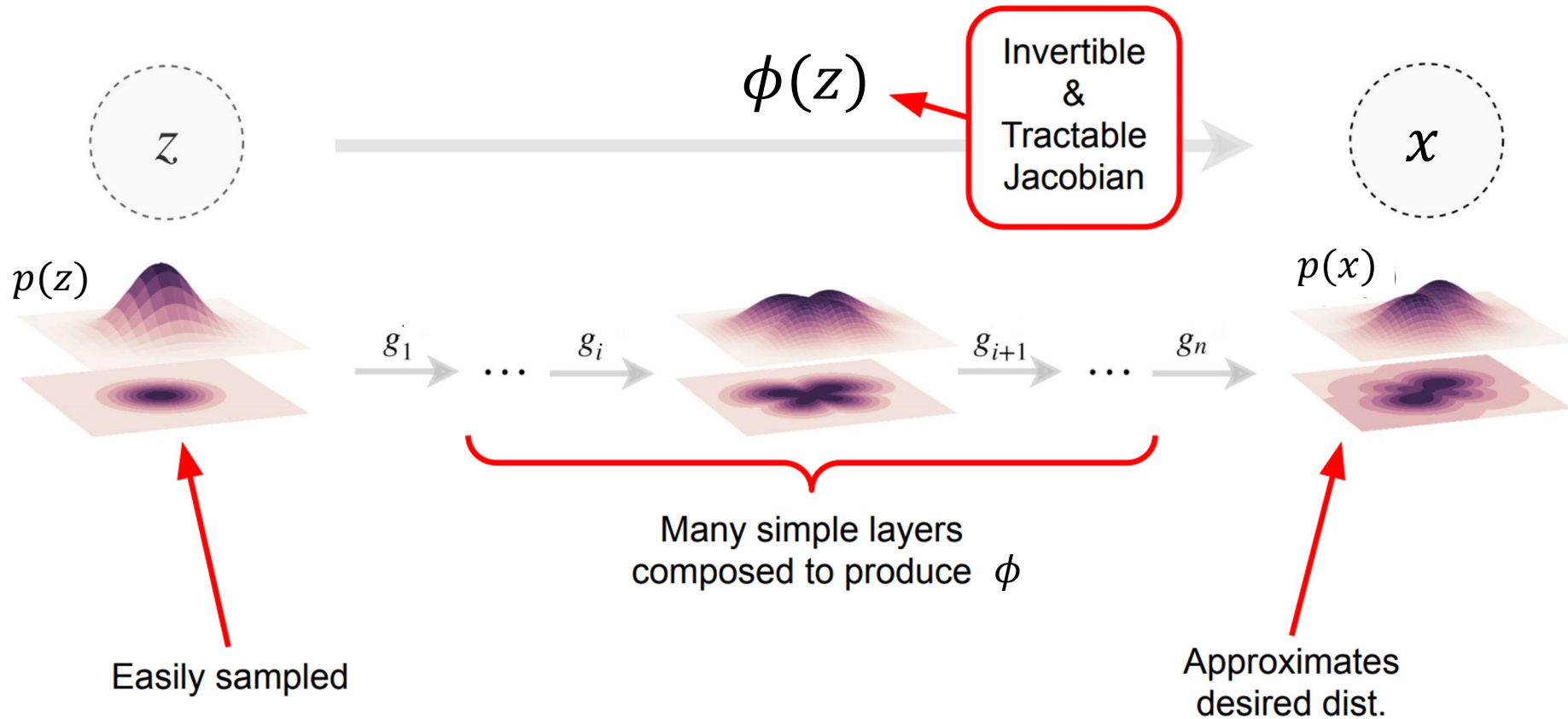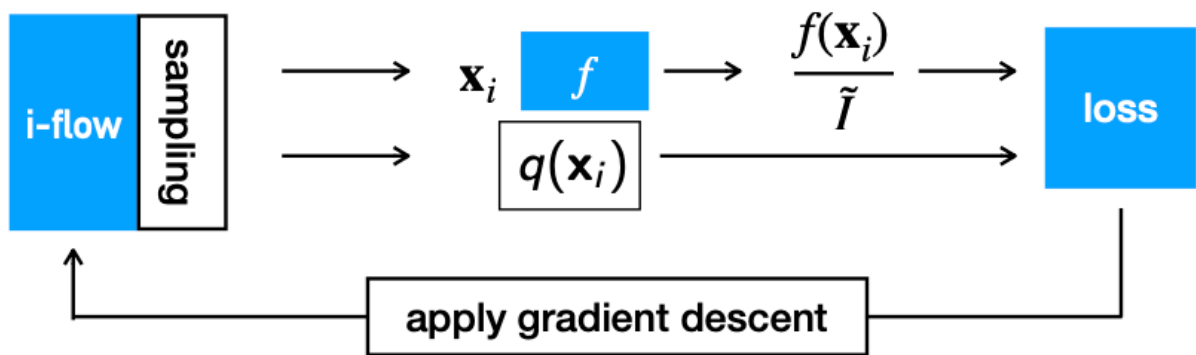Optimization of the magnet system
For the SHiP experiment



Shirbokov, **MK**, et al., NeurIPS **33**, 14650-14662 (2020)

*Explicit density estimation*
We can evaluate density $p(x)$

$$p_x(\boldsymbol{x}) = p_z(\boldsymbol{z}) \left| \det\left( \frac{\partial \phi(\boldsymbol{z})}{d\boldsymbol{z}} \right)^{-1} \right|$$



$\phi(z)$

Invertible
&
Tractable
Jacobian

$z$

$x$

$p(z)$

$p(x)$

$g_1$ $\cdots$ $g_i$ $g_{i+1}$ $\cdots$ $g_n$

Many simple layers
composed to produce $\phi$

Easily sampled

Approximates
desired dist.

arXiv: 2001.05486, ML:ST
arXiv: 2001.10028, PRD
Slide credit: C. Krause

**Example: Learning $e^+e^- \to 3j$**



$\leftarrow$ g color
$\leftarrow$ q color
$\leftarrow$ g color spectator
$\leftarrow \cos\vartheta$ of decaying fermion with beam
$\leftarrow \varphi$ of decaying fermion with beam
$\leftarrow \cos\vartheta$ of decay
$\leftarrow \varphi$ of decay
$\leftarrow$ propagator of decaying fermion
$\leftarrow$ multichannel

Target distribution
with learning color



$\leftarrow$ g color
$\leftarrow$ q color
$\leftarrow$ g color spectator
$\leftarrow \cos\vartheta$ of decaying fermion with beam
$\leftarrow \varphi$ of decaying fermion with beam
$\leftarrow \cos\vartheta$ of decay
$\leftarrow \varphi$ of decay
$\leftarrow$ propagator of decaying fermion
$\leftarrow$ multichannel

Learned distribution
with learning color

Use variational lower bound



$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is unknown

Image credit: Lilian Weng

- Iteratively add noise to data,
  Train model to learn how to denoise step by step



noise ⟶ data

# Wrapping Up

1989

~10 years

~25 years

# Still Early for Deep Learning, Where Will We be in 25 Years?

65



**~10 years**

**~25 years**

**?**

Prompt: `Several giant wooly mammoths approach treading through a snowy meadow [...]`
OpenAI Sora

2012
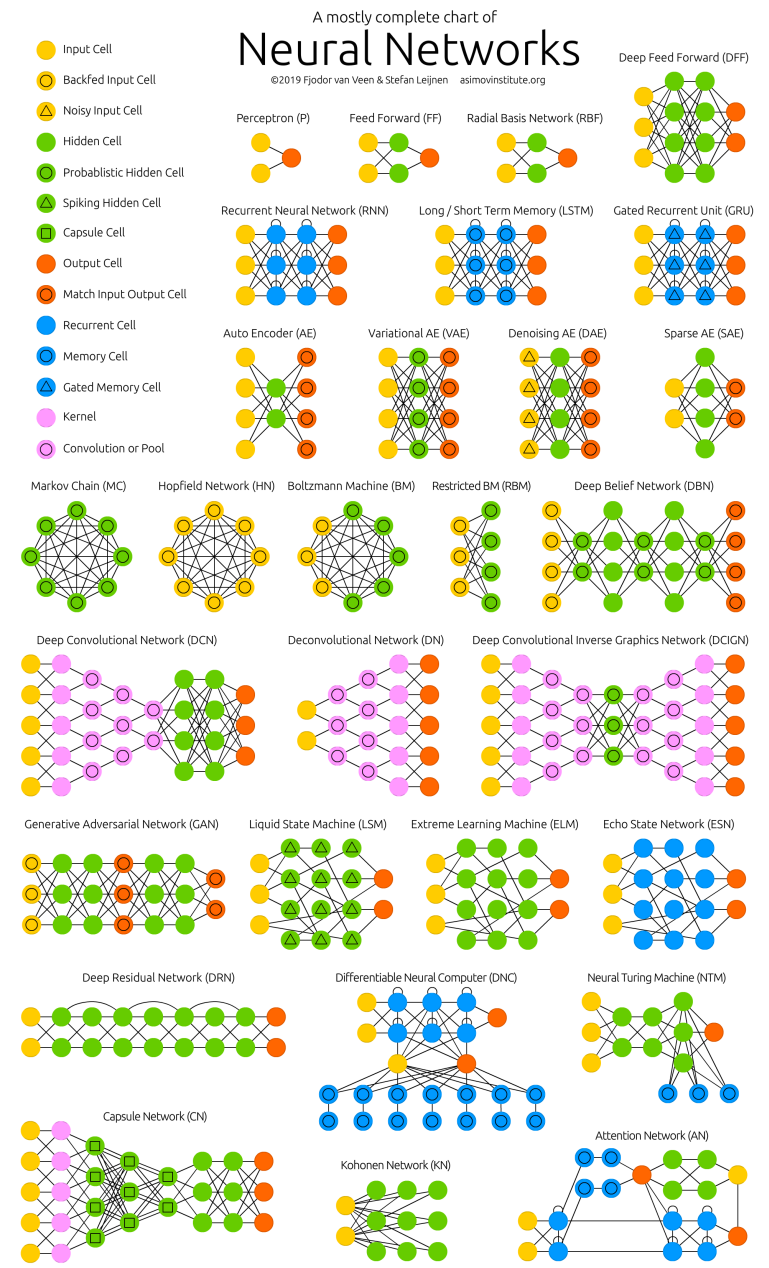
0:06

Credit: <u>Jim Fan</u> + Sora

# Summary

- Deep neural networks allow us to learn complex function by hierarchically structuring the feature learning

- We can express our inductive bias about a system in terms of model design, and can be adapted to a many types of data

- Many neural networks structures are available for training models on a wide array of data types.

- Beyond classification and regression, deep neural networks allow for powerful generative models to enable us to model and generate data

# Backup

# Deep Neural Networks

- Structure of the networks, and the node connectivity can be adapted for problem at hand

- Moving inductive bias from feature engineering to model design

  - *Inductive bias*:
    Knowledge about the problem

  - *Feature engineering*:
    Hand crafted variables

  - *Model design*:
    The data representation and the structure of the machine learning model / network



Image credit: neural-network-zoo

- A single layer network may need a width exponential in D to approximate a depth-D network's output
    - Simplified version of Telgarsky ([2015](#), [2016](#))

- A single layer network may need a width exponential in D to approximate a depth-D network's output
  - Simplified version of Telgarsky ([2015](#), [2016](#))

- Over-parametrizing a deep model often improves test performance, contrary to bias-variance tradeoff prediction
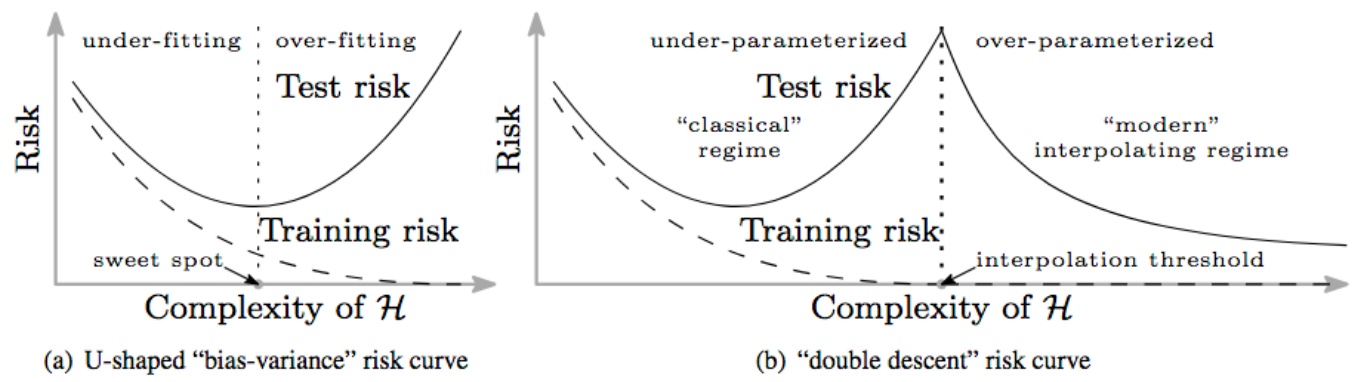
Belkin et. al. 2018



Figure 1: Curves for training risk (dashed line) and test risk (solid line). (a) The classical *U-shaped risk curve* arising from the bias-variance trade-off. (b) The *double descent risk curve*, which incorporates the U-shaped risk curve (i.e., the "classical" regime) together with the observed behavior from using high complexity function classes (i.e., the "modern" interpolating regime), separated by the interpolation threshold. The predictors to the right of the interpolation threshold have zero training risk.

- A single layer network may need a width exponential in D to approximate a depth-D network's output
  - Simplified version of Telgarsky ([2015](#), [2016](#))

- Over-parametrizing a deep model often improves test performance, contrary to bias-variance tradeoff prediction

  - But we must control that:
    - Gradients don't vanish
    - Gradient amplitude is homogeneous across network
    - Gradients are under control when weights change

- A single layer network may need a width exponential in D to approximate a depth-D network's output
  - Simplified version of Telgarsky ([2015](#), [2016](#))

- Over-parametrizing a deep model often improves test performance, contrary to bias-variance tradeoff prediction

- Major part of deep learning is choosing the right function

  - Need to make gradient descent work, even if substantial engineering required

# Convolutional Neural Networks

# 1D Convolutional Layers

- **Data:** $x \in \mathbb{R}^M$

- **Convolutional kernel of width k:** $u \in \mathbb{R}^k$

- Convolution $x \circledast u$ is vector of size M-k+1

$$(x \circledast \text{u})_i = \sum_{b=0}^{k-1} x_{i+b} u_b$$
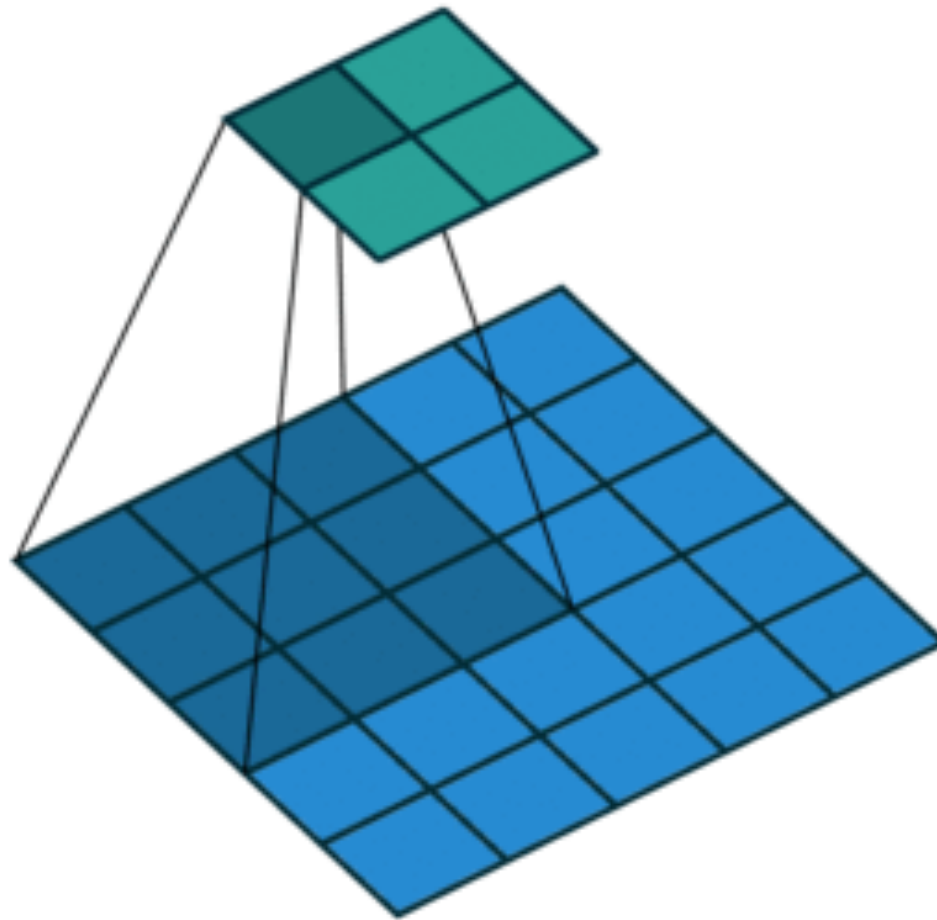
- Scan across data and multiply by kernel elements

# 2D Convolutional Layer

- Input data (tensor) **x** of size $C \times H \times W$
  - C channels (e.g. RGB in images)

- Learnable Kernel **u** of size $C \times h \times w$
  - The size $h \times w$ is the *receptive field*

$$(\boldsymbol{x} \circledast \boldsymbol{u})_{i,j} = \sum_{c=0}^{C-1} (\boldsymbol{x}_c \circledast \boldsymbol{u}_c)_{i,j} = \sum_{c=0}^{C-1} \sum_{n=0}^{h-1} \sum_{m=0}^{w-1} \boldsymbol{x}_{c,n+i,m+j} \boldsymbol{u}_{c,n,m}$$

- Output size $(H - h + 1) \times (W - w + 1)$ for each kernel
  - Often called *Activation Map* or *Output Feature Map*
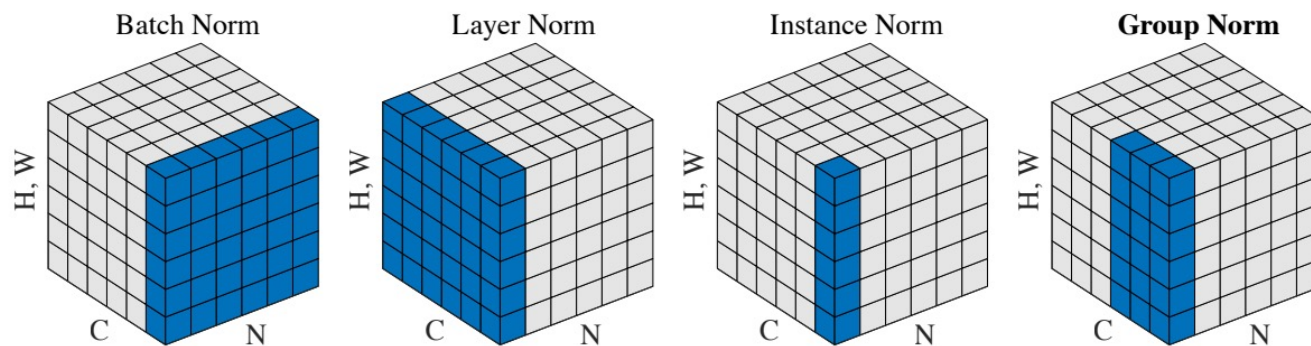
# Normalization

- Maintaining proper statistics of the activations and derivatives is a critical issue to allow the training of deep architectures

"Training Deep Neural Networks is complicated by the fact that **the distribution of each layer's inputs changes during training, as the parameters of the previous layers change**. This slows down the training by requiring lower learning rates and careful parameter initialization …"

Ioffe, Szegedy,
*Batch Normalization*, ICML 2015



Wu, He, *Group Normalization*, CoRR 2018

- During training batch normalization shifts and rescales according to the mean and variance estimated on the batch.
  - During test, use empirical moments estimated during training

- Per-component mean and variance on the batch

$$m_{batch} = \frac{1}{B} \sum_{b=1}^{B} x_b$$

$$v_{batch} = \frac{1}{B} \sum_{1}^{B} (x_b - m_{batch})^2$$

- Normalize and compute output $\forall b = 1 \dots B$

$$z_b = \frac{x_b - m_{batch}}{\sqrt{v_{batch} + \epsilon}}$$

$$y_b = \gamma \odot z_b + \beta$$



Figure 2: *Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.*

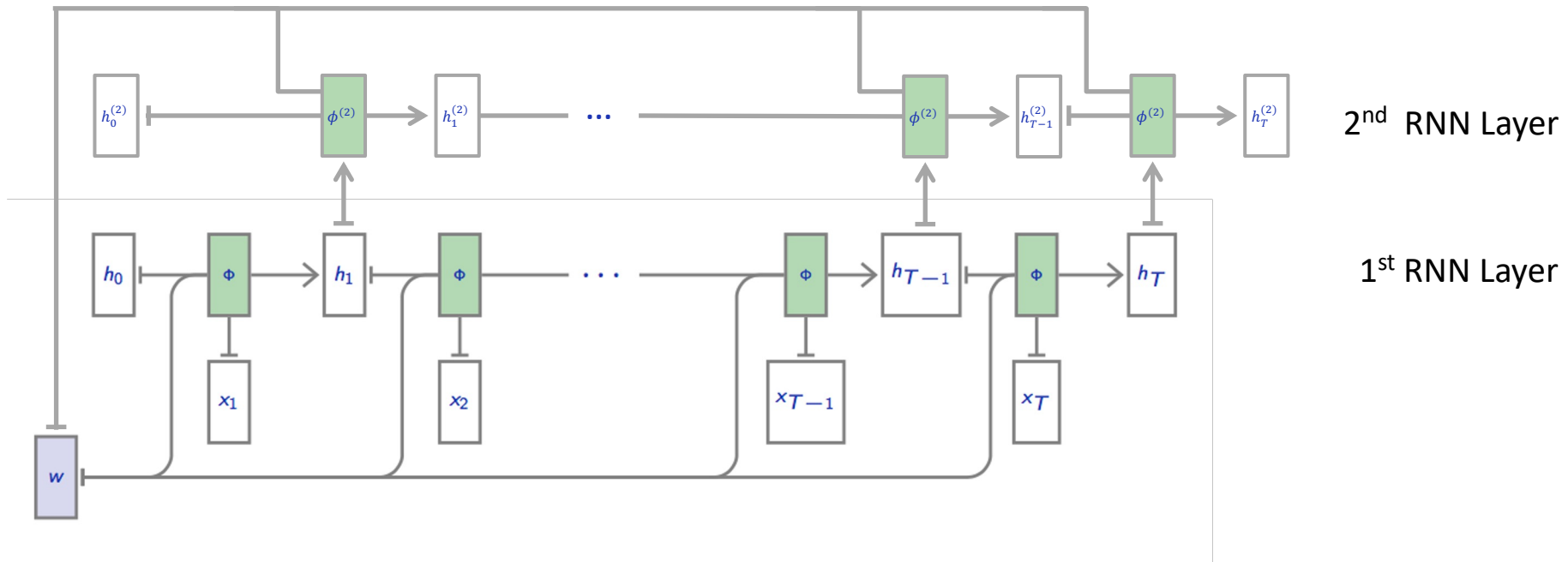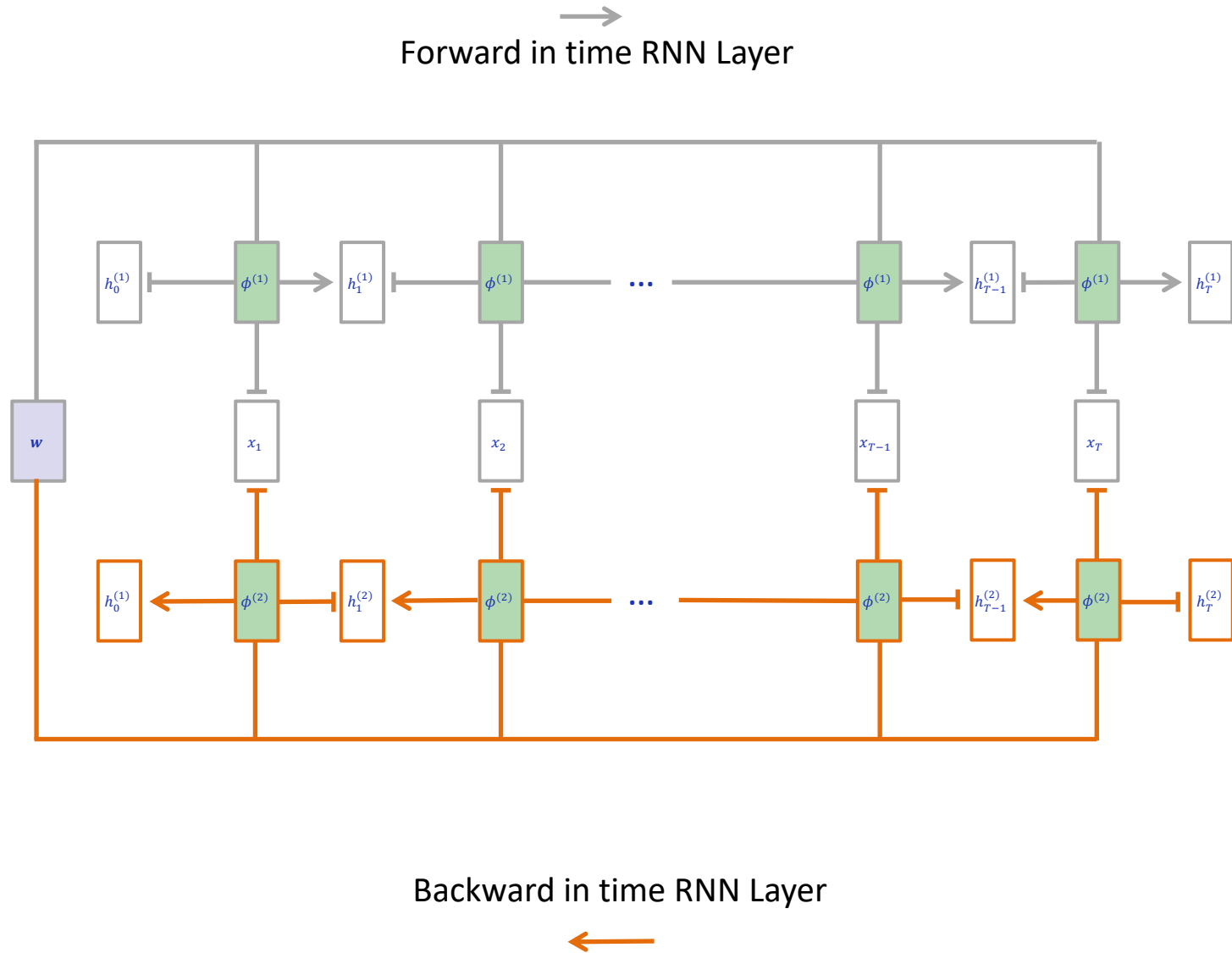  - $\gamma$ and $\beta$ are parameters to optimize

# RNN

$x_{1:T}$   RNN   $h_{1:T}$   RNN   $h_{1:T}^{(2)}$   $\cdots$   RNN   $h_{1:T}^{(N)}$

$w$

Two Stacked LSTM Layers

Forward in time RNN Layer

Backward in time RNN Layer

# Comparison on Toy Problem

Learn to recognize palindrome
Sequence size between 1 to 10

| $\mathbf{x}$ | $y$ |
|:---:|:---:|
| $(1, 2, 3, 2, 1)$ | 1 |
| $(2, 1, 2)$ | 1 |
| $(3, 4, 1, 2)$ | 0 |
| $(0)$ | 1 |
| $(1, 4)$ | 0 |

Self-driving Mario Kart with RNN: YouTube video

## Text-to-speech synthesis



Shen et al., 2017

# Deep Sets

# What if our data has no time structure?

- Data may be variable in length but have no temporal structure → *Data are sets of values*

- *One option*: If we know about the data domain, could try to impose an ordering, then use RNN

- *Better option*: use system that can operate on variable length sets in permutation invariant way

  – Why permutation invariant → so order doesn't matter

# Deep Sets

$$\tilde{h}_{1:T}$$

$\Sigma$

Permutation invariant operation: Sum, Max, …

$h_1$ $\quad$ $h_2$ $\quad$ ... $\quad$ $h_T$

$\phi$ $\quad$ $\phi$ $\quad$ $\phi$

$x_1$ $\quad$ $x_2$ $\quad$ $x_T$

w

## Outlier detection



black hair &
brown hair

M. Zaheer et. al 2017

## Medical Imaging

With more complex architecture



(a)  (b)

(c)  (d)

*Figure 5.* (a) H&E stained histology image. (b) $27 \times 27$ patches centered around all marked nuclei. (c) Ground truth: Patches that belong to the class epithelial. (d) Heatmap: Every patch from (b) multiplied by its corresponding attention weight, we rescaled the attention weights using $a'_k = (a_k - \min(\mathbf{a}))/(\max(\mathbf{a}) - \min(\mathbf{a}))$.

M. Ilse et al., 2018

# Transformers

- Gradients may not explode or vanish, but managing a meaningful context over a long sequence is challenging.

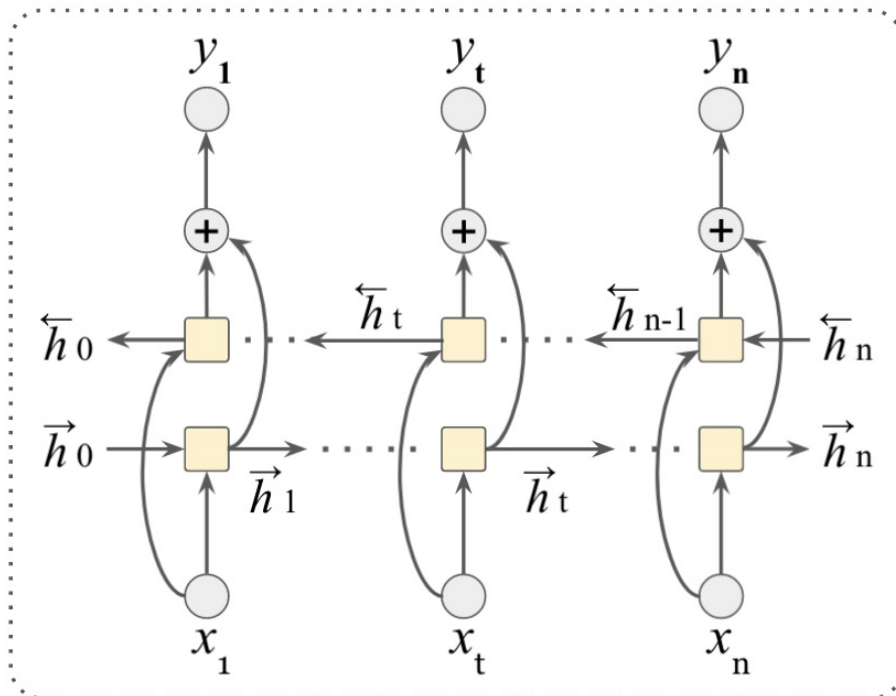- Bottleneck: fixed length array in model with long input



Bi-Directional

RNN Encoder-Decoder

- Idea: allow RNN to look at all the hidden state sequence when producing an output. Output is generated from context $c$

$$c_i = \sum_{j=1}^{T} \alpha_{ij} h_j \quad \text{where} \quad \alpha_{ij} = softmax(\beta_{ij})_{over\ j}$$

$$\text{and} \quad \beta_{ij} = U \tanh(W s_{i-1} + \widetilde{W} h_j + b_i)$$



[1409.0473](1409.0473)

- Idea: Get rid of the RNN and only use attention

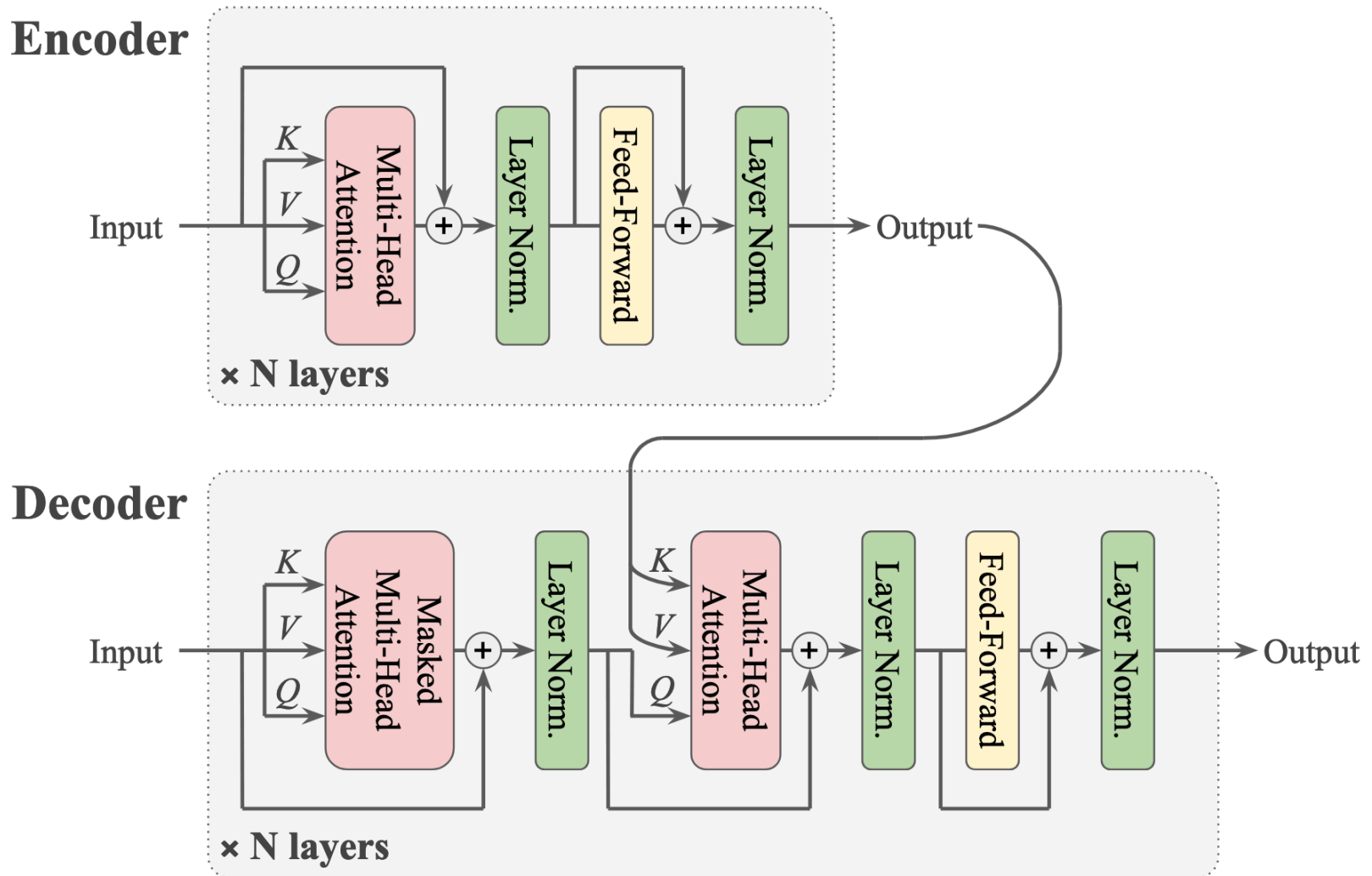$$\text{Attention}\,(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) V \qquad \text{where} \qquad \begin{aligned} Q &\in \mathbb{R}^{n \times d} \\ K &\in \mathbb{R}^{m \times d} \\ V &\in \mathbb{R}^{m \times d_v} \end{aligned}$$

Query

$$Q = \begin{pmatrix} \vec{q}_1 \\ \vdots \\ \vec{q}_n \end{pmatrix}_{n \times d}$$

Key

$$K = \begin{pmatrix} \vec{k}_1 \\ \vdots \\ \vec{k}_m \end{pmatrix}_{m \times d}$$

Value

$$V = \begin{pmatrix} \vec{v}_1 \\ \vdots \\ \vec{v}_m \end{pmatrix}_{m \times d_v}$$

- Project the input "query" onto a "key" to compute the weights for the corresponding "value"

- Return the weighted value

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) V \qquad \text{where} \qquad \begin{aligned} Q &\in \mathbb{R}^{n \times d} \\ K &\in \mathbb{R}^{m \times d} \\ V &\in \mathbb{R}^{m \times d_v} \end{aligned}$$

- **Self-Attention**: using input $X$ to define Q,K,V

$$Q = XW_Q \qquad\qquad K = XW_K \qquad\qquad V = XW_V$$



(A) attention

Legend
$\underline{x}$ = word vector for "that"
A = 3 neural networks in parallel
α = vector of soft weights for "that".
  = softmax( $xW_q$ $xW_k^T$ / sqrt(100) )
X = word vectors stacked together as
  sentence matrix

- Lets look at a single query

$$\frac{qK^T}{\sqrt{d}} = \left( \frac{\vec{q}_1 \cdot \vec{k}_1}{\sqrt{d}}, \frac{\vec{q}_1 \cdot \vec{k}_2}{\sqrt{d}}, \cdots, \frac{\vec{q}_1 \cdot \vec{k}_m}{\sqrt{d}} \right)_{1 \times m}$$

$$\text{softmax} \left( \frac{qK^T}{\sqrt{d}} \right) = (p_1, p_2, ..., p_m)_{1 \times m} = \vec{p} \quad \text{where} \quad p_i = \frac{\exp \frac{\vec{q}_1 \cdot \vec{k}_i}{\sqrt{d}}}{\sum_{j=1}^{m} \exp \frac{\vec{q}_1 \cdot \vec{k}_j}{\sqrt{d}}}$$

$$\text{Attention}(q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V = \vec{p}V = \sum_{i=1}^{n} p_i \vec{v}_i$$
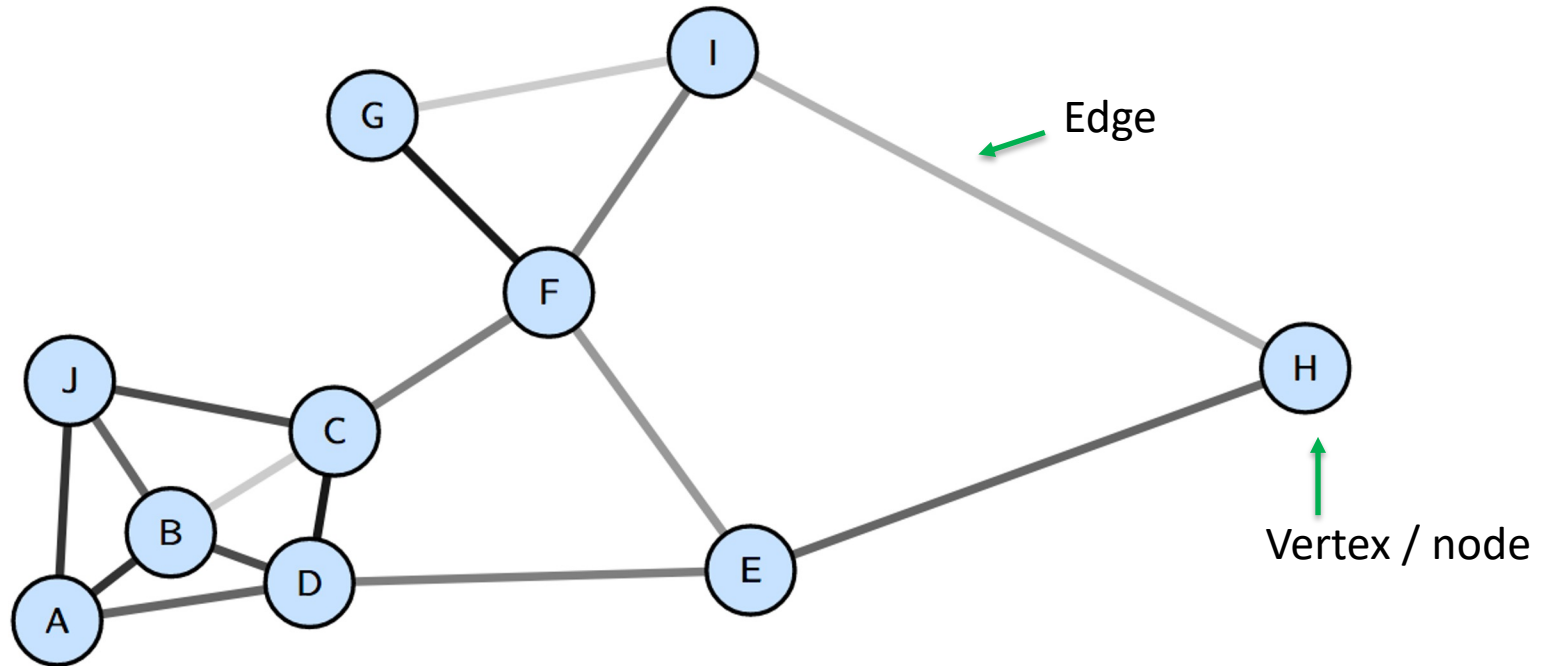
- Generalize input to length $n$

$$\text{Attention}(Q, K, T) = \begin{pmatrix} p_{11}\vec{v}_1 + p_{12}\vec{v}_2 + \cdots + p_{1m}\vec{v}_m \\ p_{21}\vec{v}_1 + p_{22}\vec{v}_2 + \cdots + p_{2m}\vec{v}_m \\ \vdots \\ p_{n1}\vec{v}_1 + p_{n2}\vec{v}_2 + \cdots + p_{nm}\vec{v}_m \end{pmatrix} = \begin{pmatrix} \sum_{i}^{m} p_{1i}\vec{v}_i \\ \sum_{i}^{m} p_{2i}\vec{v}_i \\ \vdots \\ \sum_{i}^{m} p_{ni}\vec{v}_i \end{pmatrix}_{n \times d_v}$$
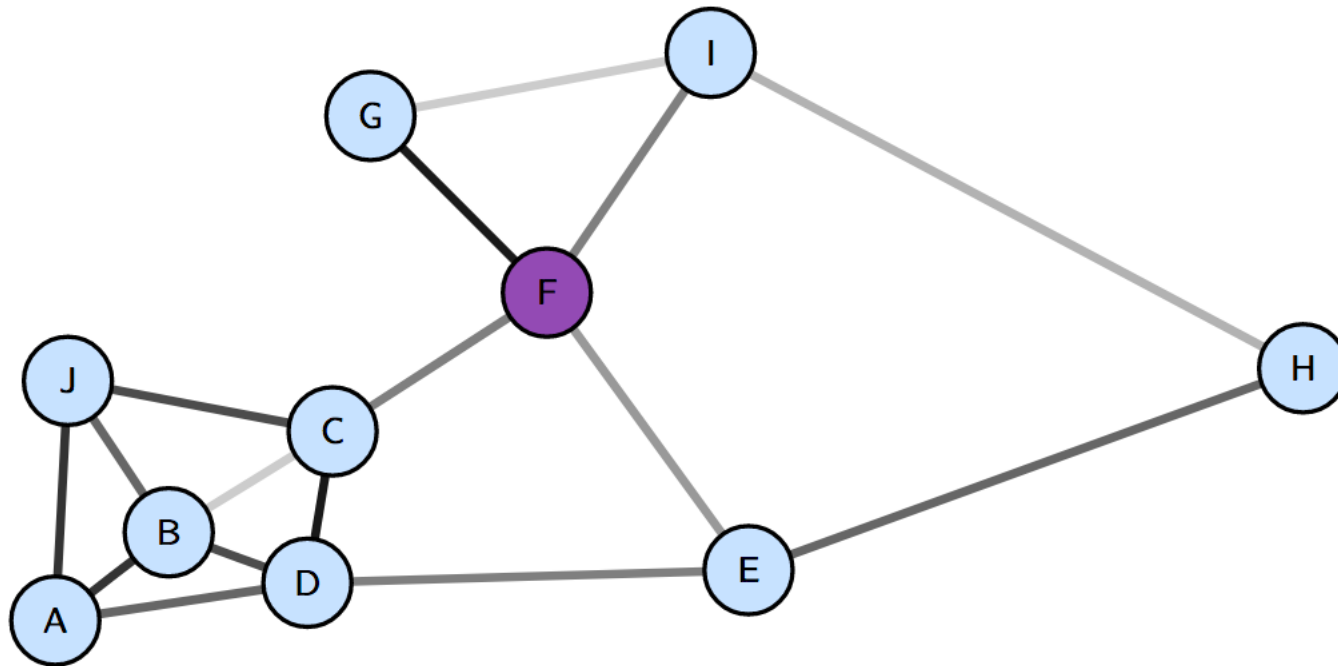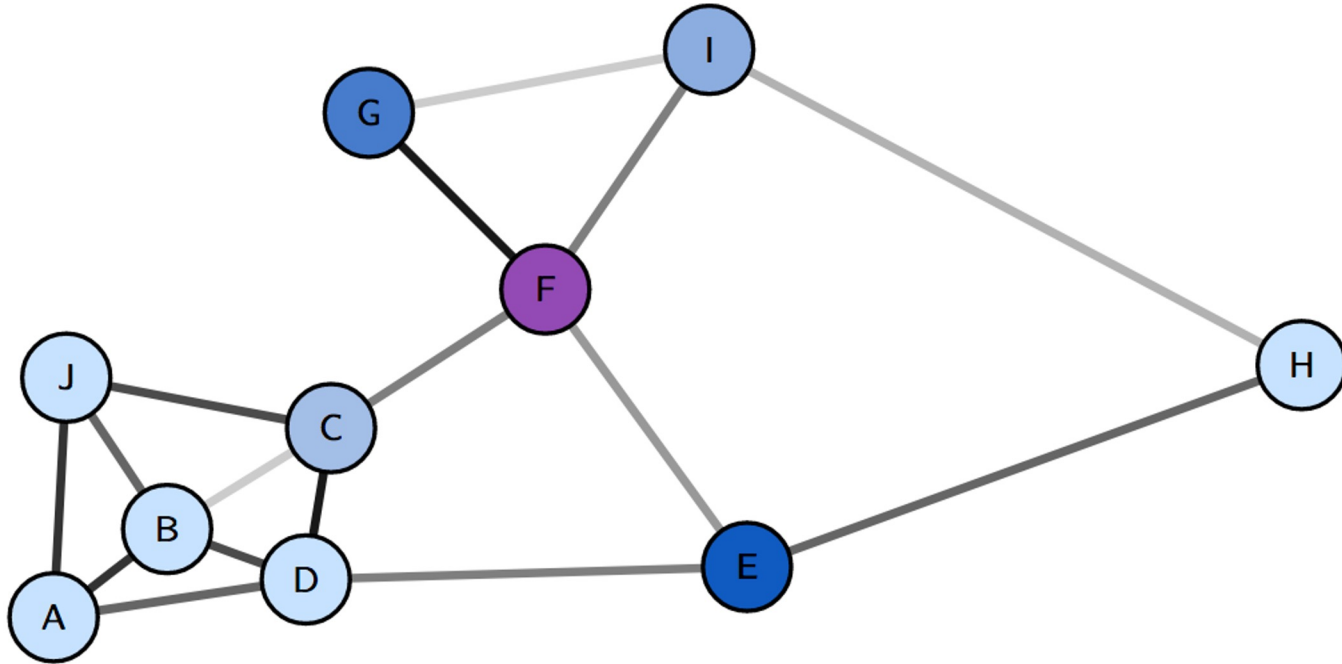
# Graph Neural Networks

# Graph Data

- Sequential data has single (directed) connections from data at current time to data at next time

- What about data with more complex dependencies

Image credit: N. Wang et al., 2018

# Graphs

Edge

Vertex / node

- Adjacency matrix: $A_{ij} = \delta(edge\ between\ vertex\ i\ and\ j)$

- Each node can have features

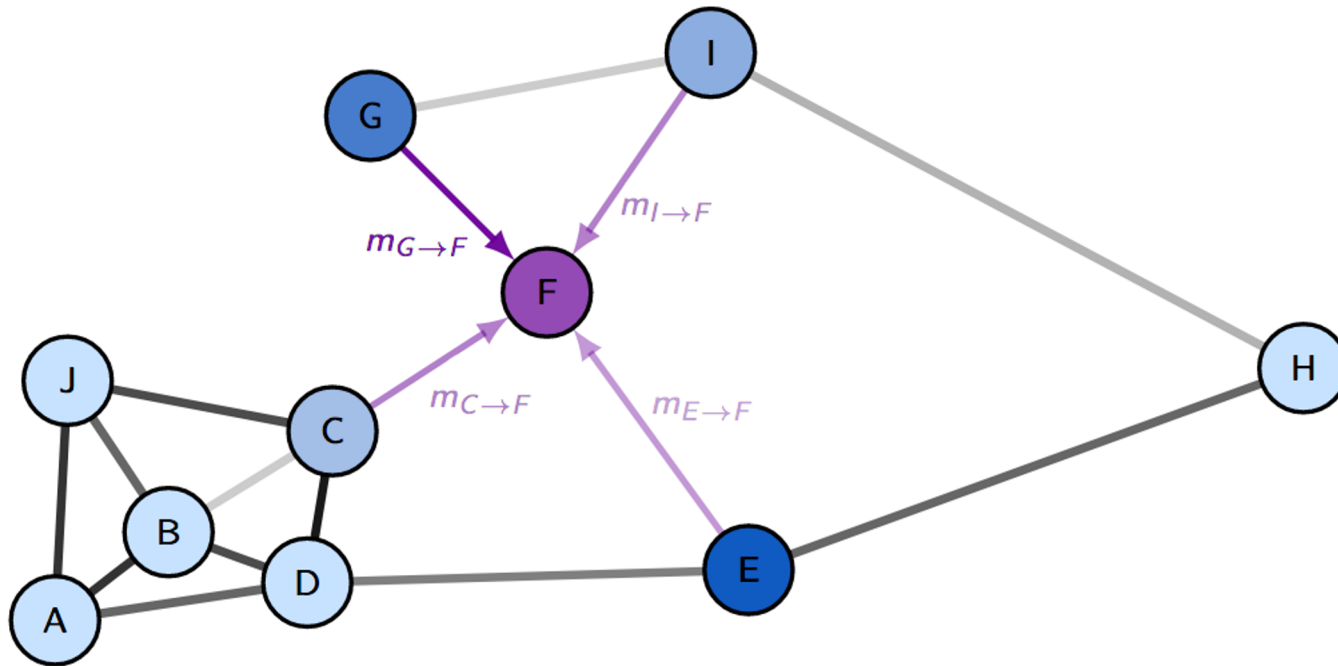- Each edge can have features, e.g. distance between nodes
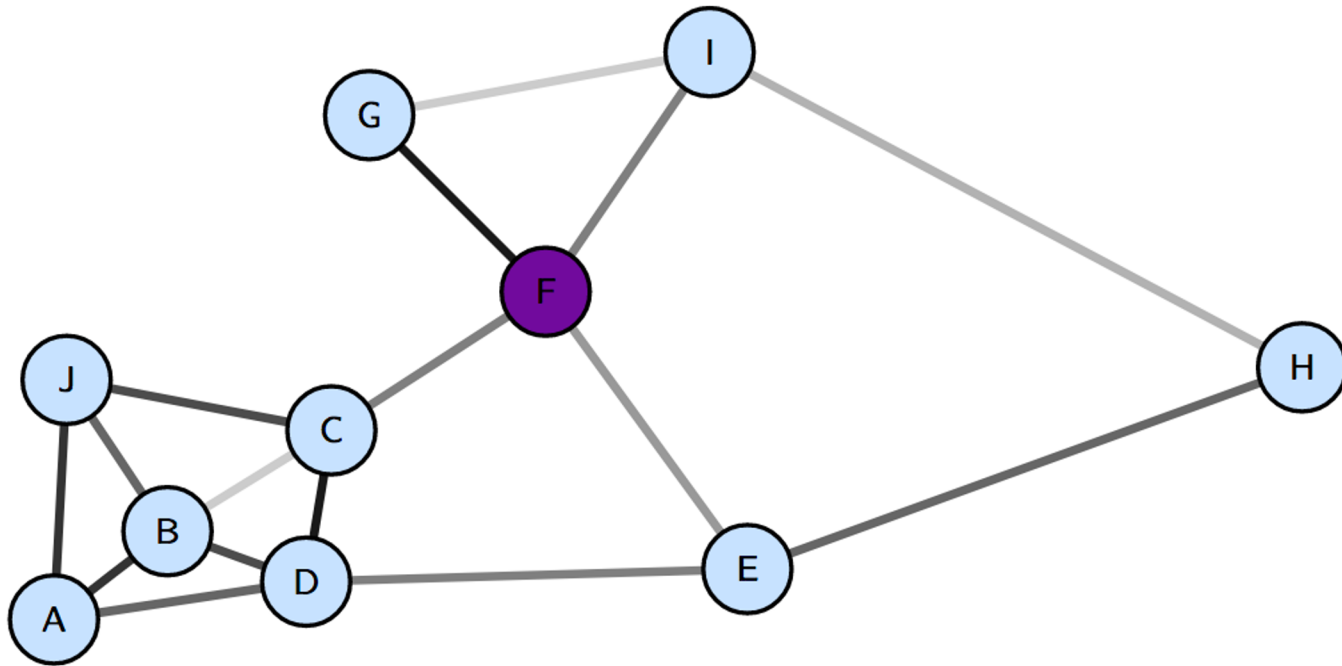
Image Credit: I. Henrion

# Neural Message Passing

$$\tilde{m}_j^t = f(h_j^{t-1})$$

Image Credit: I. Henrion

$$\tilde{m}_j^t = f(h_j^{t-1})$$

$$m_{j \to i}^t = \sigma(A_{ij} \tilde{m}_j^t)$$

$$\tilde{m}_j^t = f(h_j^{t-1})$$

$$m_{j \to i}^t = \sigma(A_{ij} \tilde{m}_j^t)$$

$$h_i^t = \text{GRU}(h_i^{t-1}, \Sigma_j m_{j \to i}^t)$$

Image Credit: I. Henrion

# Neural Message Passing

---

**Algorithm 1** Message passing neural network

---

**Require:** $N \times D$ nodes $\mathbf{x}$, adjacency matrix $A$

  $\mathbf{h} \leftarrow \text{Embed}(\mathbf{x})$

  **for** $t = 1, \ldots, T$ **do**

    $\mathbf{m} \leftarrow \text{Message}(A, \mathbf{h})$

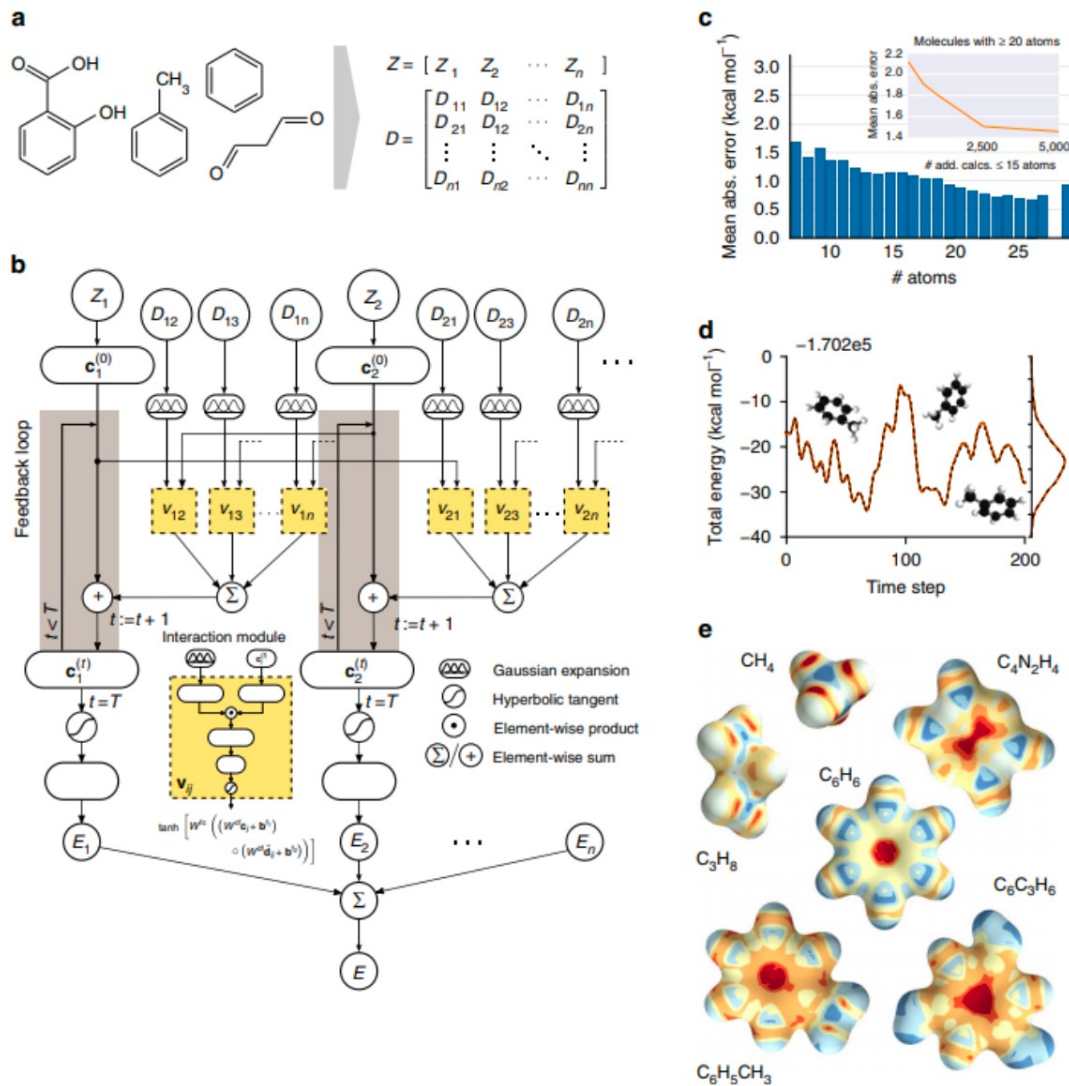    $\mathbf{h} \leftarrow \text{VertexUpdate}(\mathbf{h}, \mathbf{m})$

  **end for**

  $\mathbf{r} = \text{Readout}(\mathbf{h})$

  **return** $\text{Classify}(\mathbf{r})$

---

Image Credit: I. Henrion

## Quantum chemistry with graph networks



Schutt et al. 2017
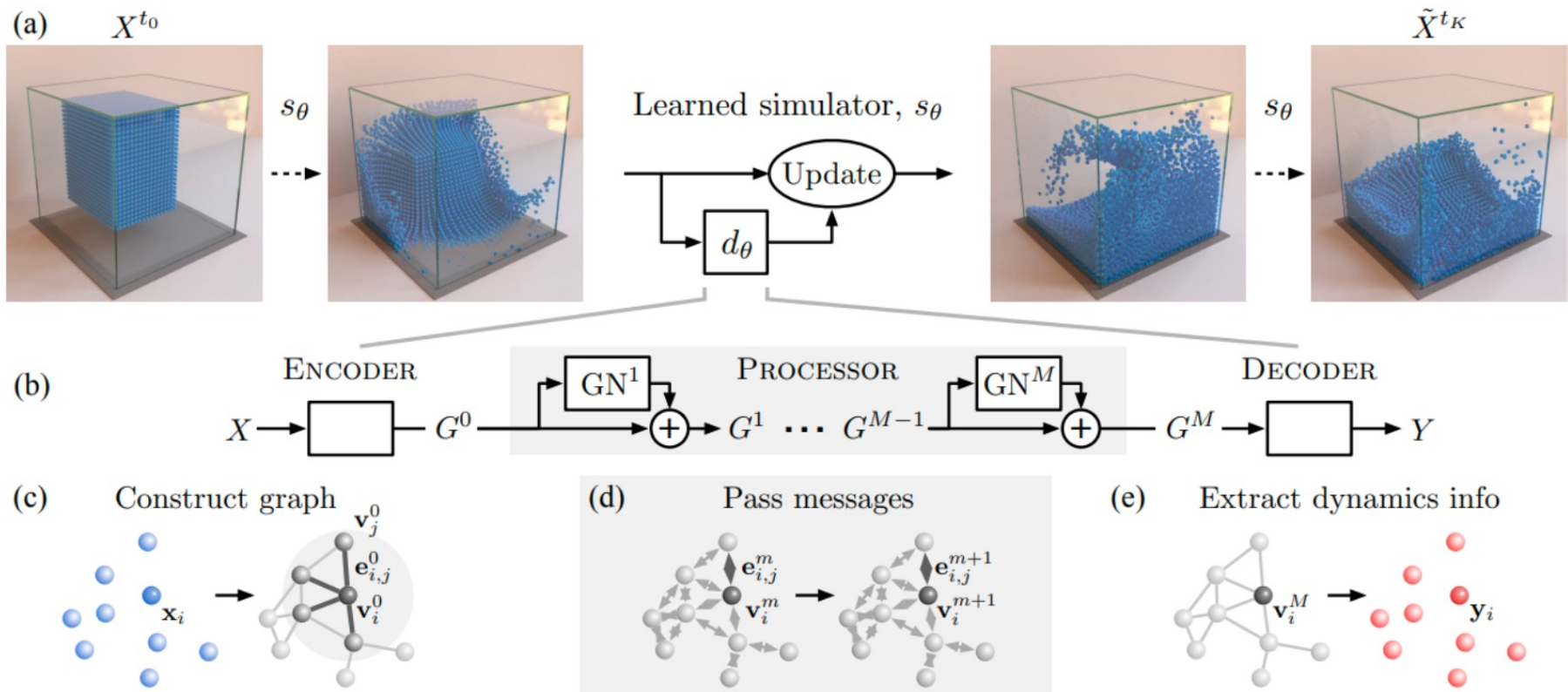
## Learning to simulate physics with graph networks



Figure 2. (a) Our GNS predicts future states represented as particles using its learned dynamics model, $d_\theta$, and a fixed update procedure. (b) The $d_\theta$ uses an "encode-process-decode" scheme, which computes dynamics information, $Y$, from input state, $X$. (c) The ENCODER constructs latent graph, $G^0$, from the input state, $X$. (d) The PROCESSOR performs $M$ rounds of learned message-passing over the latent graphs, $G^0, \ldots, G^M$. (e) The DECODER extracts dynamics information, $Y$, from the final latent graph, $G^M$.

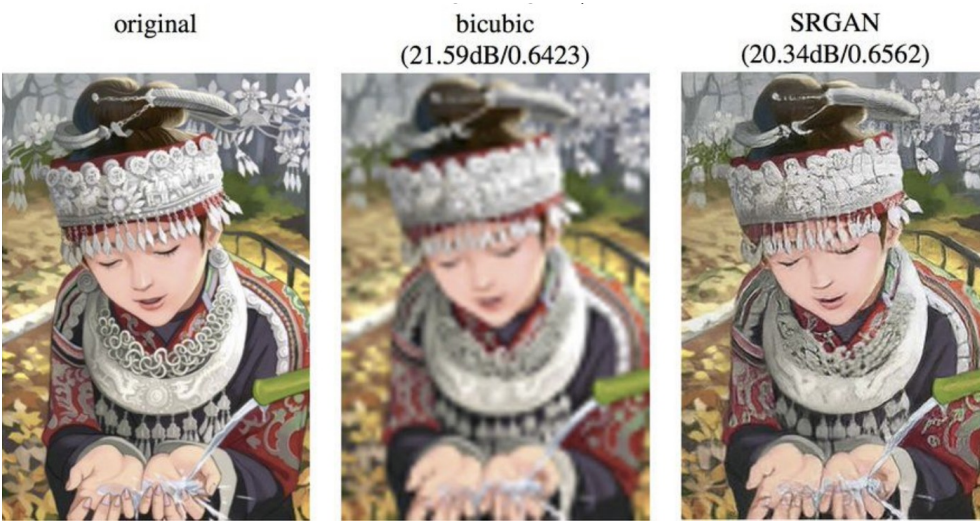# Deep Generative Model Examples

## BigGan



(Brock et al, 2018)



## Image-to-Image Translation with CycleGAN
Zhu et. al. 2017



Simulate future trajectories of environments based on actions for planning. (Finn et al, 2016)
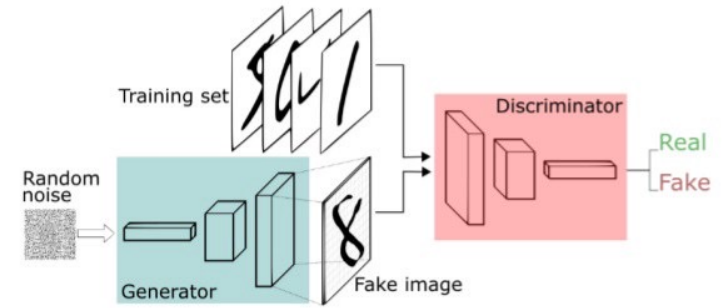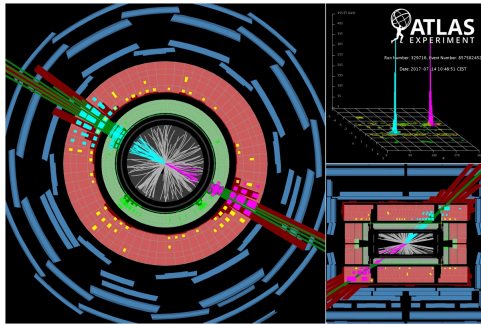
## Single Image Super-Resolution (Ledig et al, 2016)



Text-to-Image Synthesis with StackGAN (Zhang et. al. 2017)

# Generative Models

# Generative Models approximate and simulate the data generation process
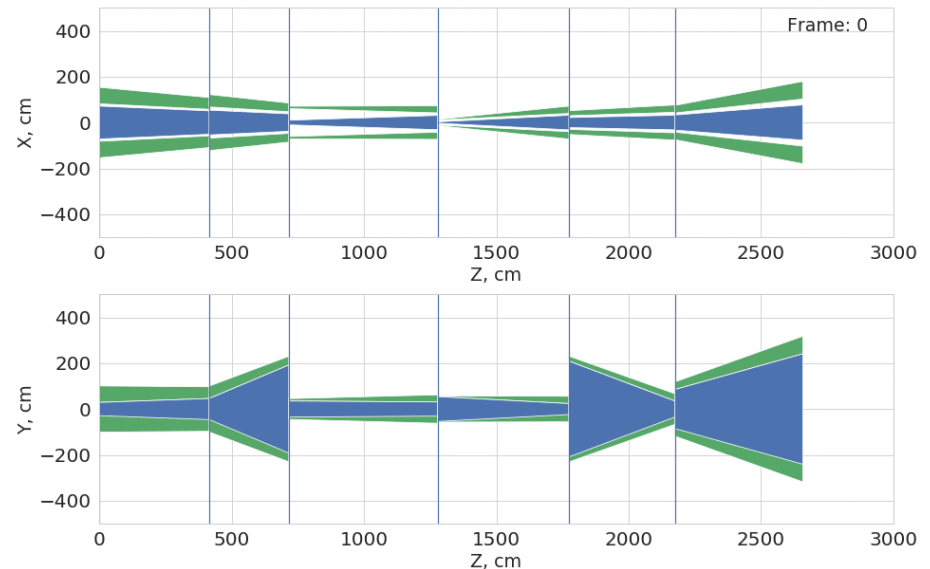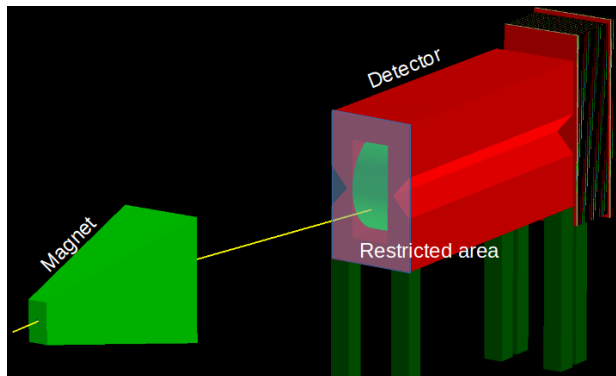


## Scientific Simulators

- Built from science knowledge

- Relatively few parameters, often interpretable

- High Fidelity, often computationally costly

## Machine Learning

- Fit to data, inductive bias in model design / optimization

- Can have $>10^6$ parameters, often not interpretable

- Often slow to train, fast to evaluate

# GANs for Detector Design

- GAN to emulate detector simulation $\tilde{x} = g(z|\psi)$
  given detector parameters $\psi$ (e.g. magnet shape below)

- Design objective $C$ to minimize: $\min\limits_{\psi} \mathbb{E}_{\tilde{x}}[C(\tilde{x} = g(z|\psi))]$

- GAN is differentiable → Minimize with gradient descent



Magnet Optimization

NeurIPS **33**, 14650-14662 (2020)

# Variational Autoencoders

- Learn a mapping from corrupted data space $\widetilde{\mathcal{X}}$ back to original data space

  - Mapping $\phi_w(\widetilde{\mathcal{X}}) = \mathcal{X}$
  - $\phi_w$ will be a neural network with parameters $w$

- Loss:

$$L = \frac{1}{N} \sum_n \|x_n - \phi_w(x_n + \epsilon_n)\|$$
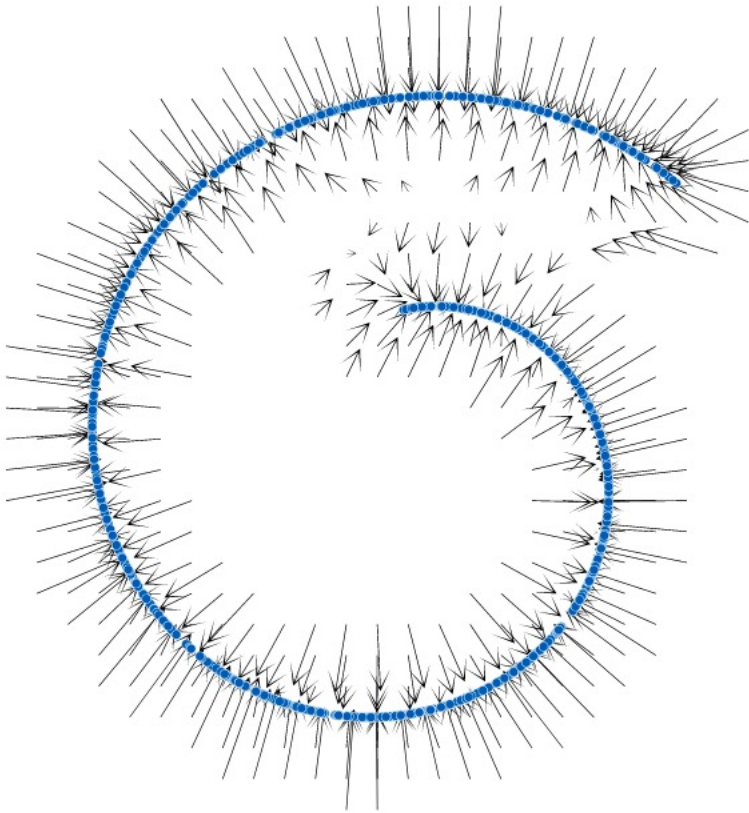
Perturbation, e.g. Gaussian noise

- **Autoencoder learns the average behavior**

- What if we care about these variations?

- Can we add a notion of variation in the autoencoder?

# Autoencoder

$f$

$g$

**x**

Latent space $\mathscr{F}$

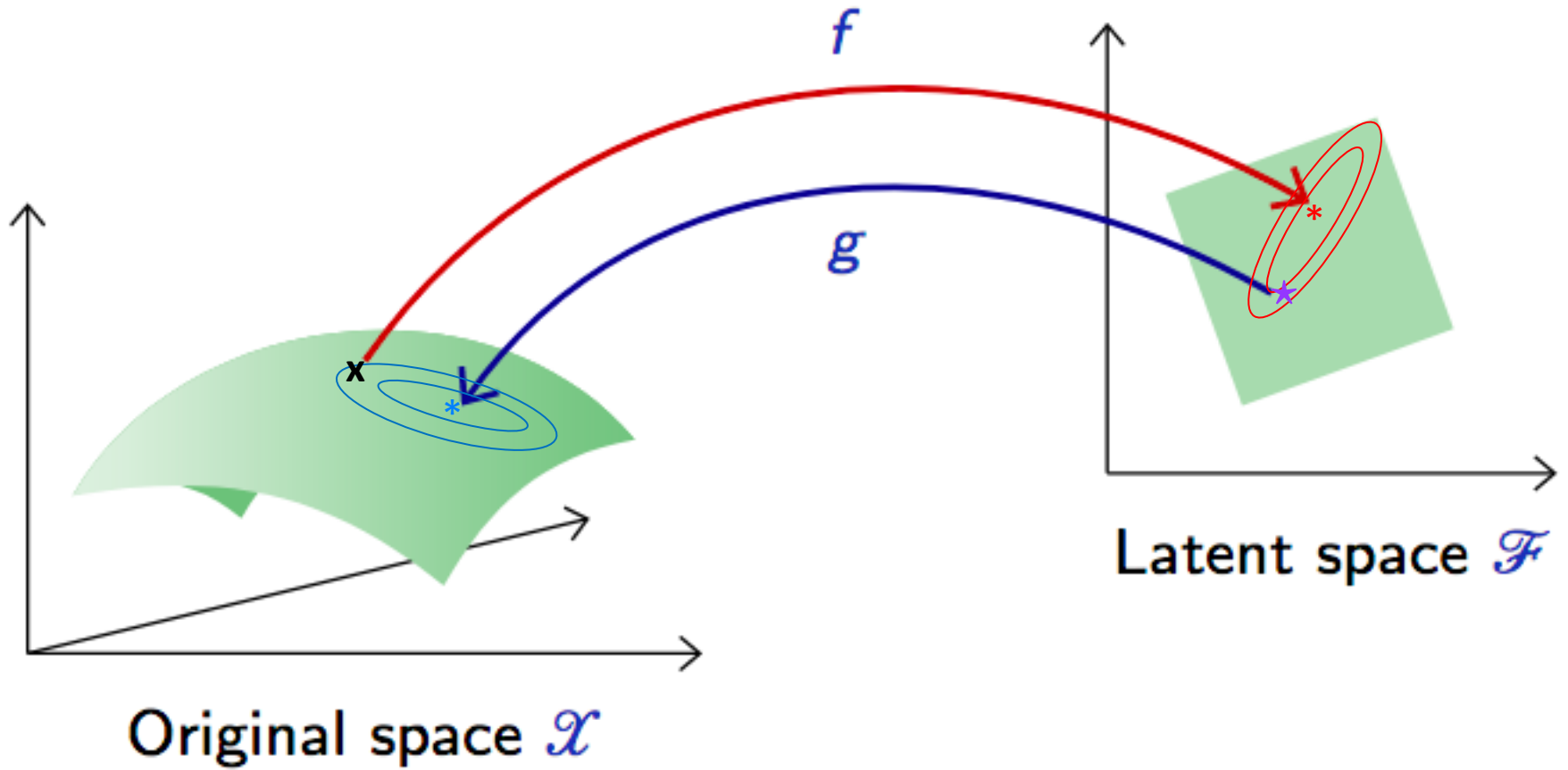Original space $\mathscr{X}$

# Variational Autoencoder

$f$

$g$

Draw sample

Latent space $\mathscr{F}$

Original space $\mathscr{X}$

- Observed random variable $x$ depends on unobserved latent random variable $z$

- Joint probability: $p(x, z) = p(x|z)p(z)$

- $p(x|z)$ is stochastic generation process from $z \to x$
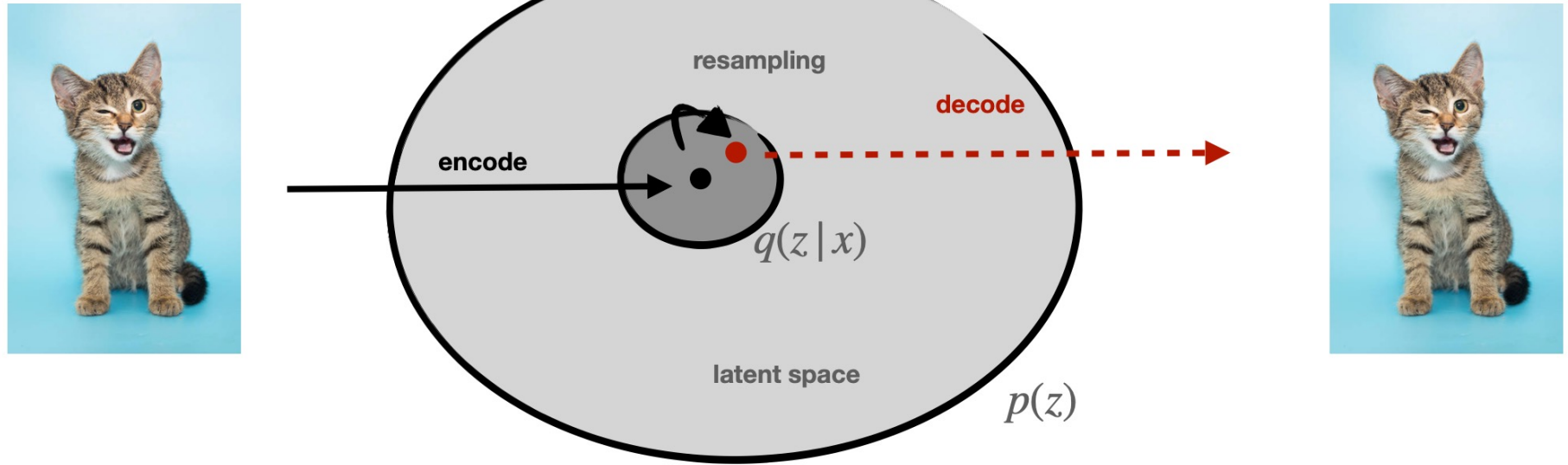
- Probabilistic relationship between data and latents

$$x, z \sim p(x, z) = p(x|z)p(z)$$

- Autoencoding

$$x \rightarrow q(z|x) \underset{sample}{\Longrightarrow} z \rightarrow p(x|z)$$

– **Encoder:** Learn what latents can produced data: $q(z|x)$

– **Decoder:** Learn what data is produced by latent: $p(x|z)$

# Variational Autoencoder
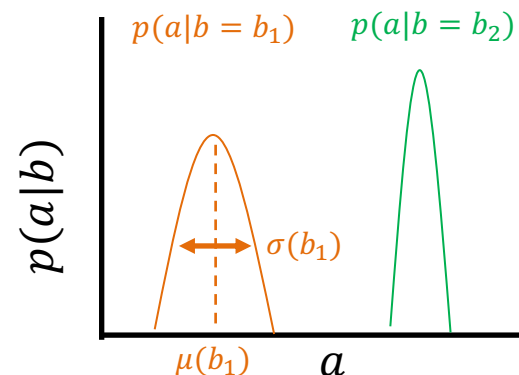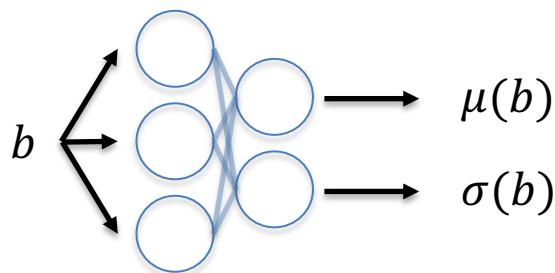
- Close-by points must decode to similar images

- Classification / regression models make single predictions…

  How to model a conditional density $p(a|b)$ ?

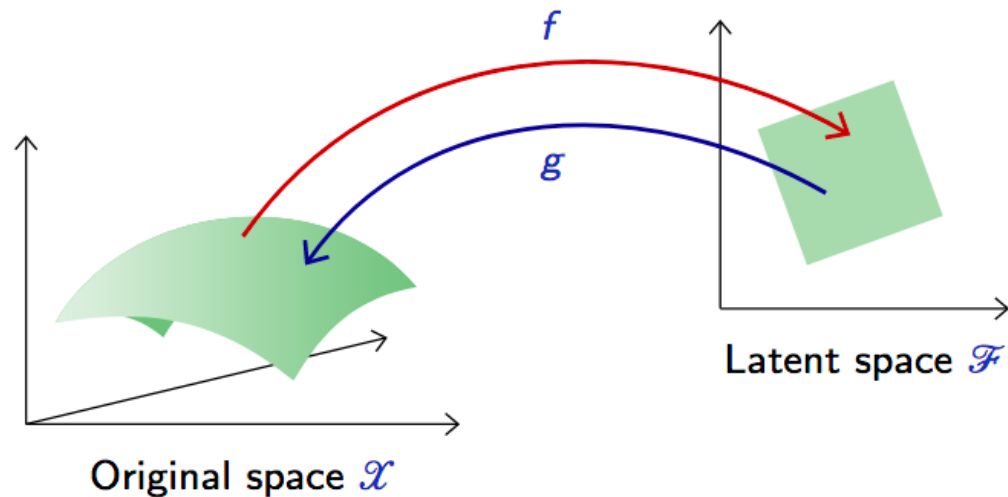- Assume a known form of density, e.g. normal

$$p(a|b) = \mathcal{N}\big(a; \mu(b), \sigma(b)\big)$$

  - Parameters of density depend on conditioned variable

- Use neural network to model density parameters

- Typical encoder maps input $x$ to "average" point in latent space

$$f(x) = \mu(x)$$



Original space $\mathcal{X}$

Latent space $\mathcal{F}$

$f$

$g$

- A VAE Encoder has two outputs: mean & variance function

$$f_\psi(x) = \{\mu_\psi(x), \sigma^2_\psi(x)\}$$

$\psi$ are parameters of the NN



*f*

Density

*g*

Draw sample

x

Latent space $\mathscr{F}$

Original space $\mathscr{X}$
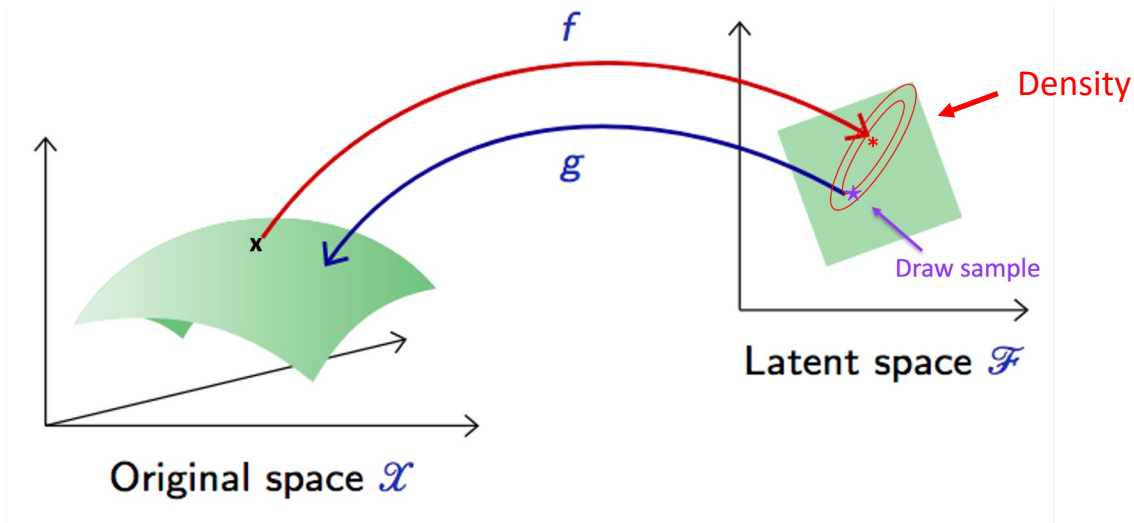
- A VAE Encoder has two outputs: mean & variance function

$$f_\psi(x) = \{\mu_\psi(x), \sigma_\psi^2(x)\}$$

$\psi$ are parameters of the NN

- What is the probability of a point in latent space?

$$p_\psi(z|x) = N(z \mid \mu_\psi(x), \sigma_\psi^2(x))$$

Could choose different density
Gaussian is easiest

- Given $x \sim p(x|\theta)$

- Sometimes, we can rewrite $x$ as a function of the parameters and a simpler distribution without parameter dependence

$$x = g(\epsilon, \theta) \qquad \epsilon \sim p(\epsilon)$$

- Example:

$$x \sim N(x|\mu, \sigma) \rightarrow x = \sigma * \epsilon + \mu \text{ with } \epsilon \sim N(0,1)$$

# Encoding

- A VAE Encoder has two outputs: mean & variance function

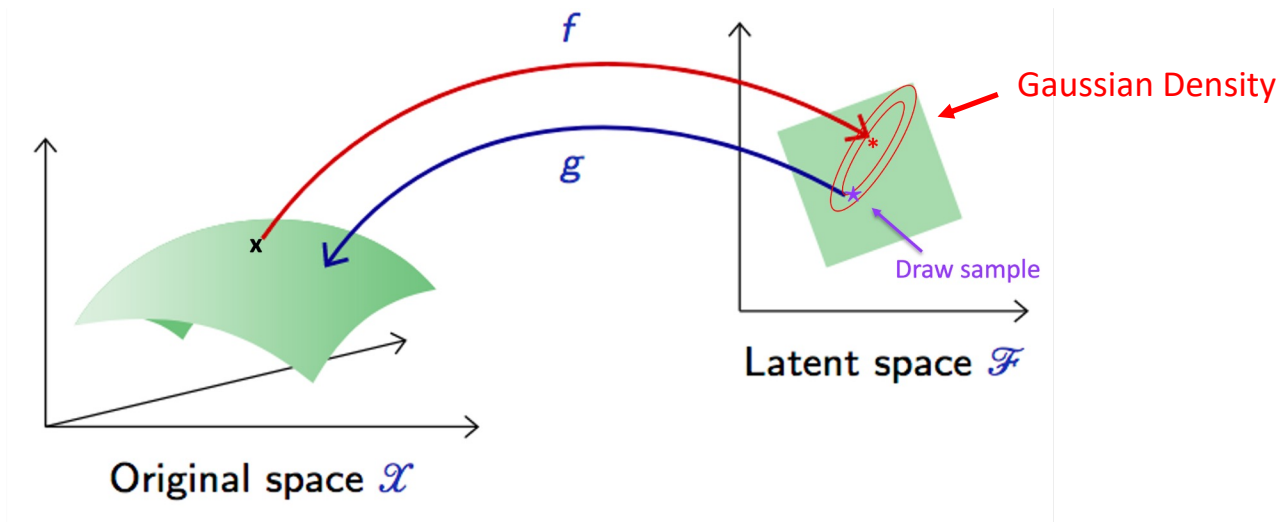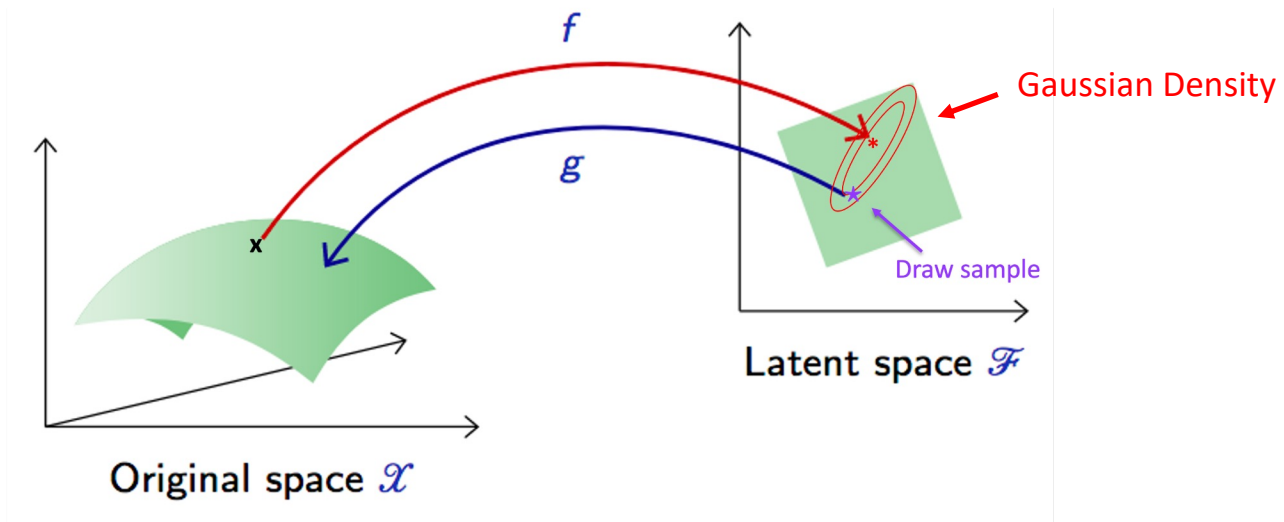$$f_\psi(x) = \{\mu_\psi(x), \sigma_\psi^2(x)\}$$

$\psi$ are parameters of the NN
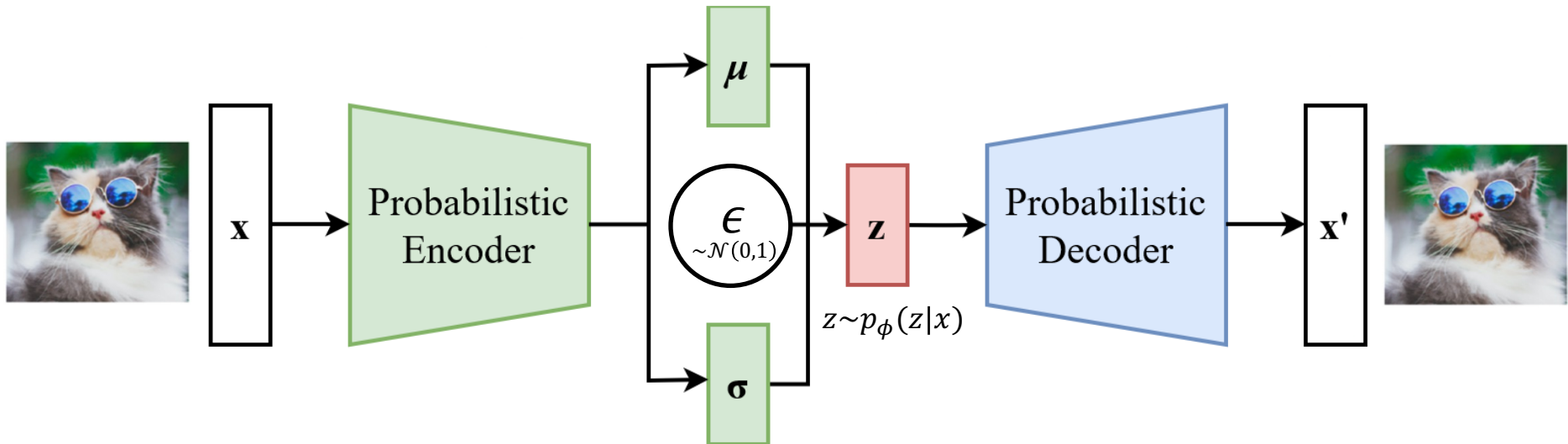
- What is the probability of a point in latent space?

$$p_\psi(z|x) = N(z \mid \mu_\psi(x), \sigma_\psi^2(x))$$

Could choose different density
Gaussian is easiest

- How do we draw a sample in latent space?

$$z = \sigma_\psi(x) * \epsilon + \mu_\psi(x) \qquad \epsilon \sim N(0, I)$$

Re-parameterization trick

Kingma, Welling, 1312.6114
Rezende, Mohamed, Wierstra, 1401.4082

# Decoding

- Same as autoencoder

$$g_\theta(z) \equiv \mu_\theta(z)$$

$\theta$ are parameters of the NN

- Likelihood of an observation $x$

$$p_\theta(x|z) = N(x \mid \mu_\theta(z), I)$$

# Decoding
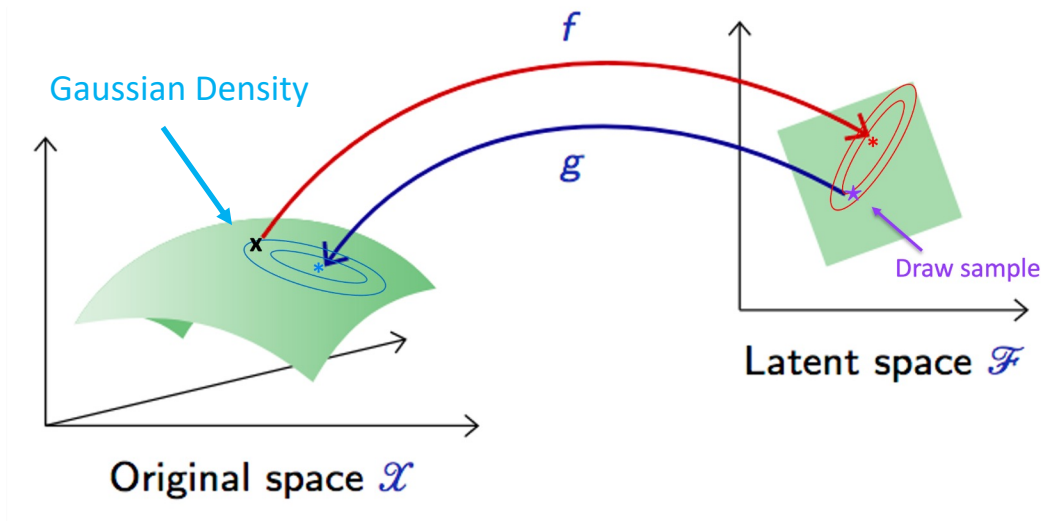
- Same as autoencoder

$$g_\theta(z) \equiv \mu_\theta(z)$$

$\theta$ are parameters of the NN

- Likelihood of an observation $x$

$$p_\theta(x|z) = N(x \mid \mu_\theta(z), I)$$

- "**Reconstruction Loss**": Maximum likelihood

$$L_{reco} = \mathbb{E}_{z \sim q(z|x)}[\log p(x|z)]$$

# Decoding

- Same as autoencoder
$$g_\theta(z) \equiv \mu_\theta(z)$$
$\theta$ are parameters of the NN

- Likelihood of an observation $x$
$$p_\theta(x|z) = N(x \mid \mu_\theta(z), I)$$

- "**Reconstruction Loss**": Maximum likelihood

$$L_{reco} = \mathbb{E}_{z \sim q(z|x)}[\log p(x|z)] \approx \frac{1}{N} \sum_{z_i \sim q(z|x)} \log N(x \mid g_\theta(z_i), I)$$

# Decoding

- Same as autoencoder

$$g_\theta(z) \equiv \mu_\theta(z)$$

$\theta$ are parameters of the NN

- Likelihood of an observation $x$

$$p_\theta(x|z) = N(x \mid \mu_\theta(z), I)$$

- "**Reconstruction Loss**": Maximum likelihood

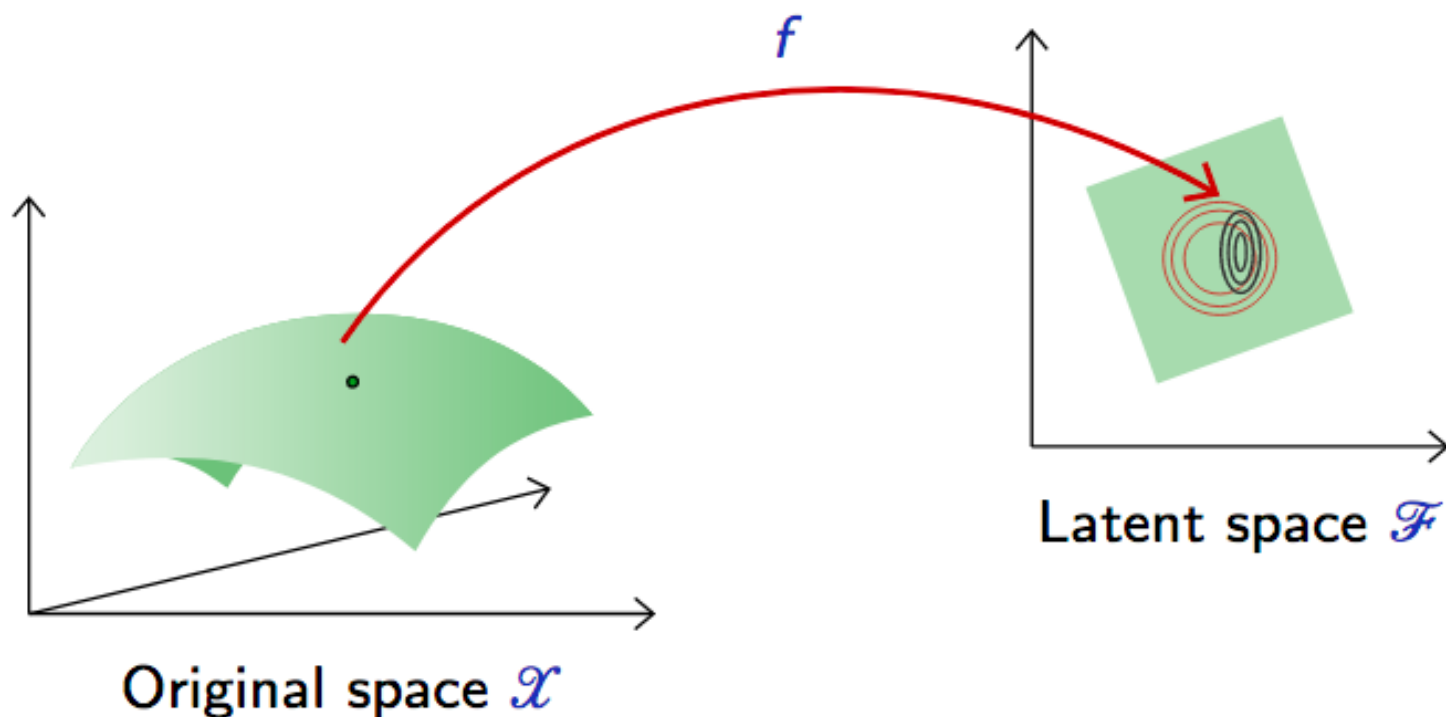$$L_{reco} = \mathbb{E}_{z \sim q(z|x)}[\log p(x|z)] \approx -\frac{1}{N} \sum_{z_i \sim q(z|x)} \left(x - g_\theta(z_i)\right)^2$$

Same as the autoencoder loss

- How do we make sure system doesn't collapse to an autoencoder (i.e. VAE encoder only predicts mean)?

- How do we make sure system doesn't collapse to an autoencoder (i.e. VAE encoder only predicts mean)?

- Use prior $p(z)$ for the latent space distribution, **need to ensure the encoder is consistent with prior**



Original space $\mathcal{X}$

Latent space $\mathcal{F}$

$f$

- Constrain difference between distributions with **Kullback–Leibler divergence**

$$D_{KL}[q(z|x)|p(z)] = \mathbb{E}_{q(z|x)}\left[\log\frac{q(z|x)}{p(z)}\right] = \int q(z|x)\log\frac{q(z|x)}{p(z)} \, dz$$
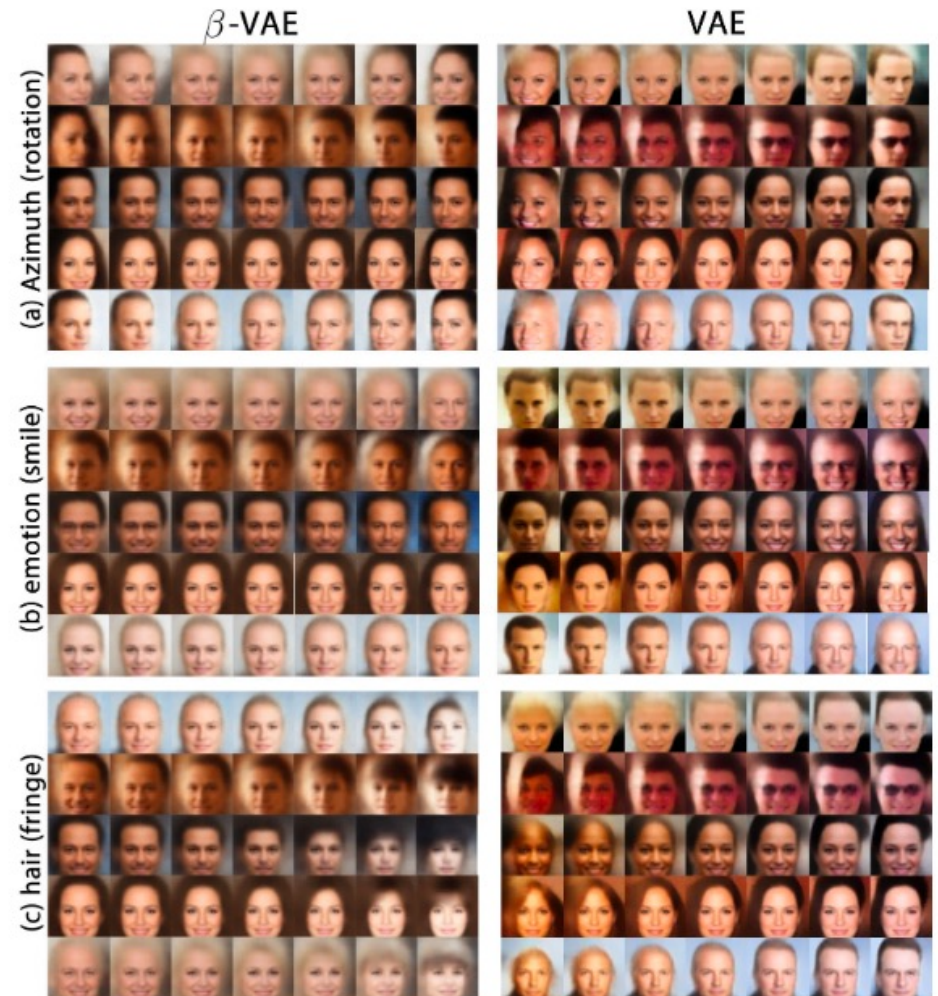
- $D_{KL}[q|p] \geq 0$ and is only 0 when $q = p$

- Constrain difference between distributions with **Kullback–Leibler divergence**

$$D_{KL}[q(z|x)|p(z)] = \mathbb{E}_{q(z|x)}\left[\log\frac{q(z|x)}{p(z)}\right] = \int q(z|x)\log\frac{q(z|x)}{p(z)}\,dz$$
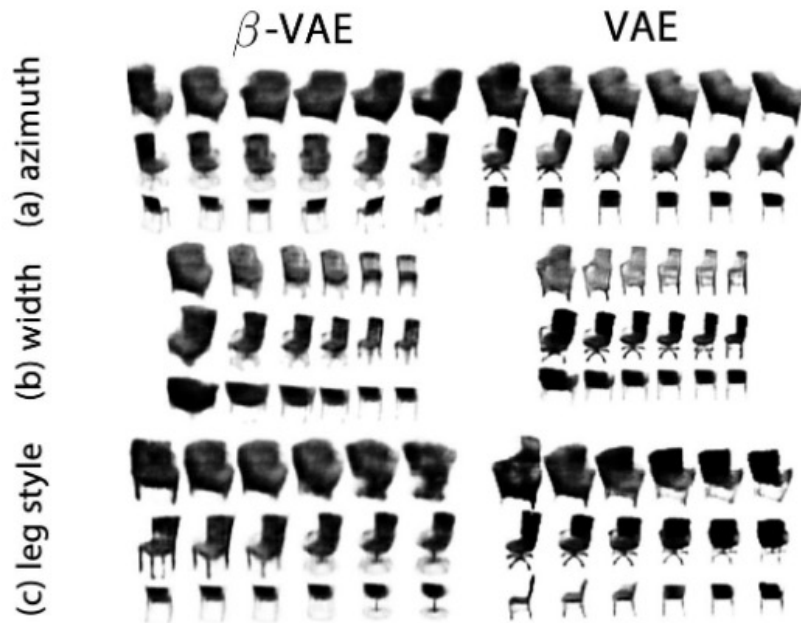
- VAE full objective

**Reconstruction Loss**　　　**Regularization of Encoder**

$$\max_{\theta,\psi} L(\theta,\psi) = \max_{\theta,\psi}\left[\mathbb{E}_{q_\psi(z|x)}[\log p_\theta(x|z)] - D_{KL}[q_\psi(z|x)|p(z)]\right]$$

Higgins et al., 2017
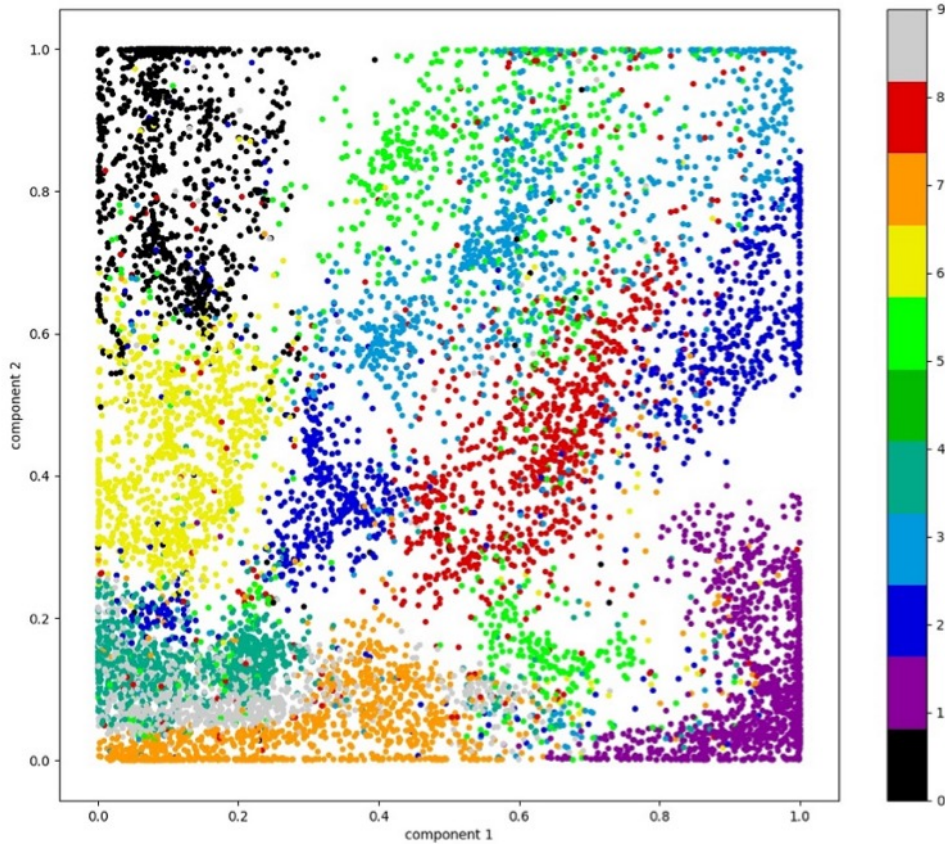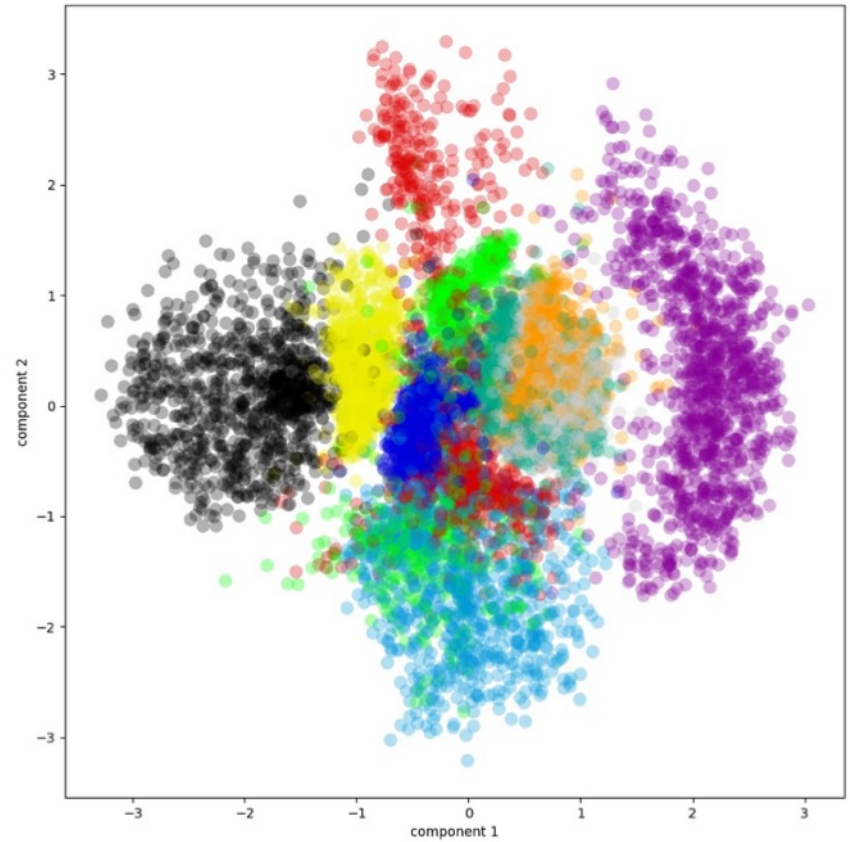
# Comparing Latent Spaces

## Autoencoder

## Variational Autoencoder



Data: MNIST data set of hand-written digits