

Deep Learning for Computer Vision: A Crash Course

Yu-Chiang Frank Wang

Professor, Dept. Electrical Engineering, National Taiwan University

Research Director for Deep Learning & Computer Vision, NVIDIA

About Myself

• Education



- **Ph.D./M.S. in Electrical & Computer Engineering** 2002 – 2009
Carnegie Mellon University, Pittsburgh, USA
- **B.S. in Electrical Engineering** 1997 – 2001
National Taiwan University, Taipei, Taiwan

• Experiences



- **Research Director, NVIDIA** 2022/08 ~
- **Professor, Dept. EE, National Taiwan University**
- **Principal AI Consultant, Inventec** 2021 – 2022
- **AI Advisory Consultant, ASUS Intel. Cloud Services (AICS)** 2019 – 2022
- **Associate Professor, National Taiwan University** 2017 – 2019
- **Deputy Director, CITI, Academia Sinica** 2015 – 2017
- **Associate/Assistant Research Fellow** 2009 – 2017
Research Center for IT Innovation (CITI), Academia Sinica

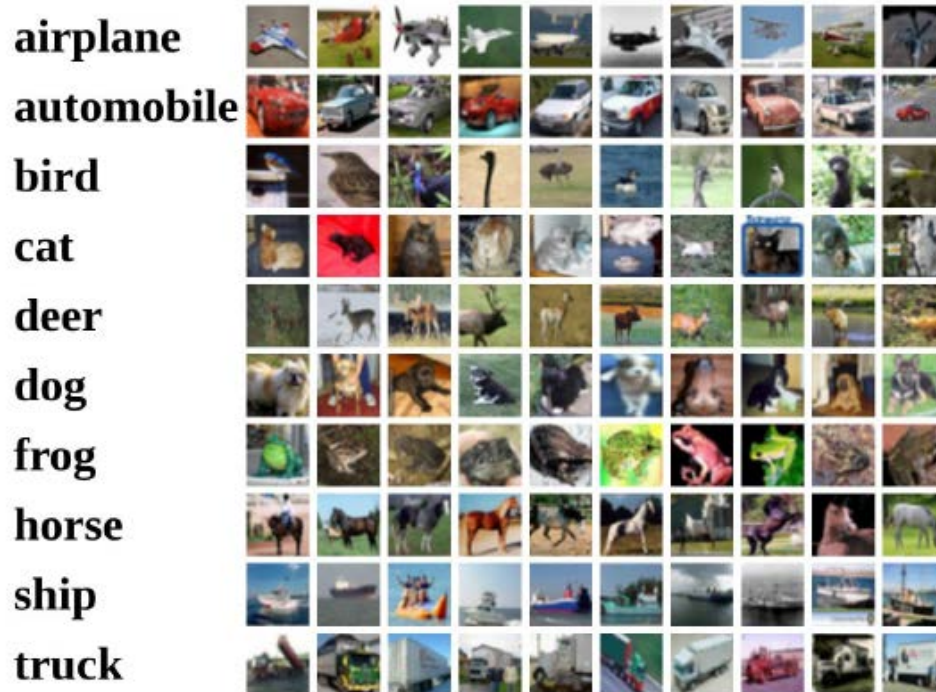


What'd Be Covered in This Crash Course...

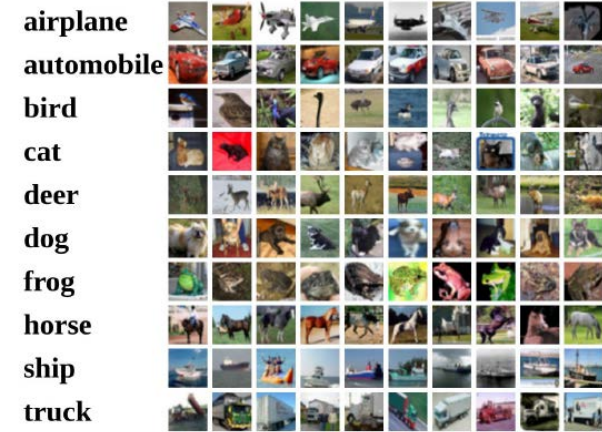
- **Learning-based Computer Vision**
 - From Linear to Non-Linear Classifiers
- **Start of Deep Learning for Computer Vision**
 - Convolutional Neural Networks
 - Self-Supervised Learning
 - Segmentation & Detection
- **Generative Models**
 - Autoencoder, Variational Autoencoder, Generative Adversarial Networks & Diffusion Models
- **Sequence-to-Sequence Learning**
 - Attention is All You Need: Transformer
- **Vision & Language Foundation Models**
 - Image-to-Text vs. Text-to-Image
 - Parameter-Efficient Fine-tuning

Linear Classification

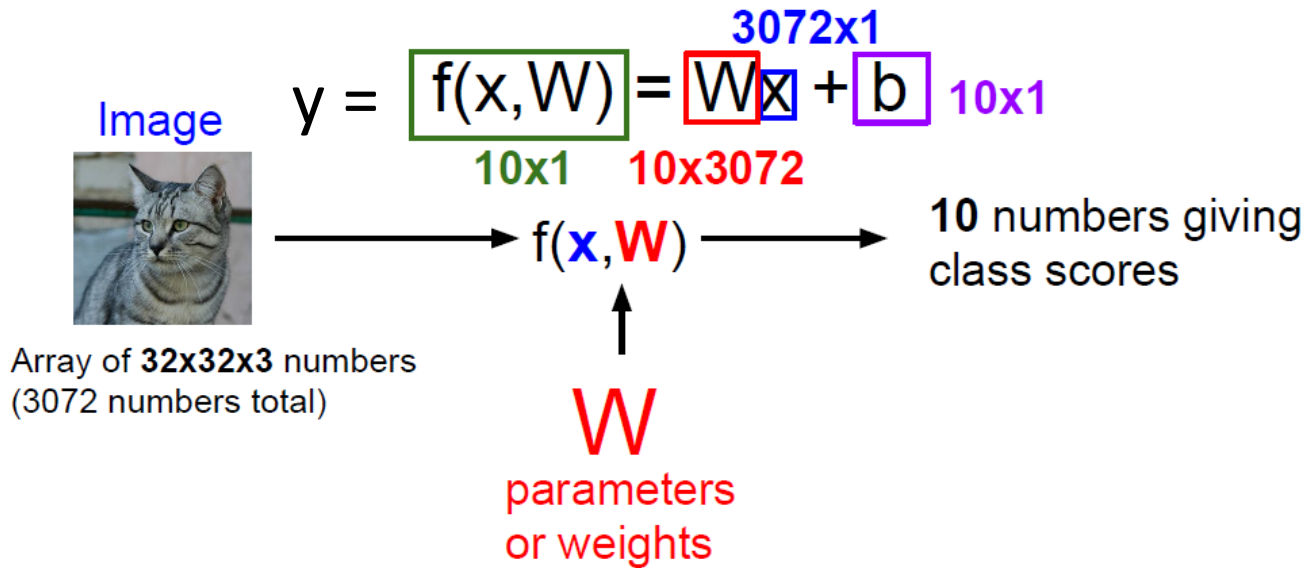
- Linear Classifier
 - Can be viewed as a **parametric or algebraic approach**.
 - Consider that we have 10 object categories of interest
 - E.g., CIFAR10 with 50K training & 10K test images of 10 categories.
Each image is of size 32 x 32 x 3 pixels.



Linear Classification (cont'd)

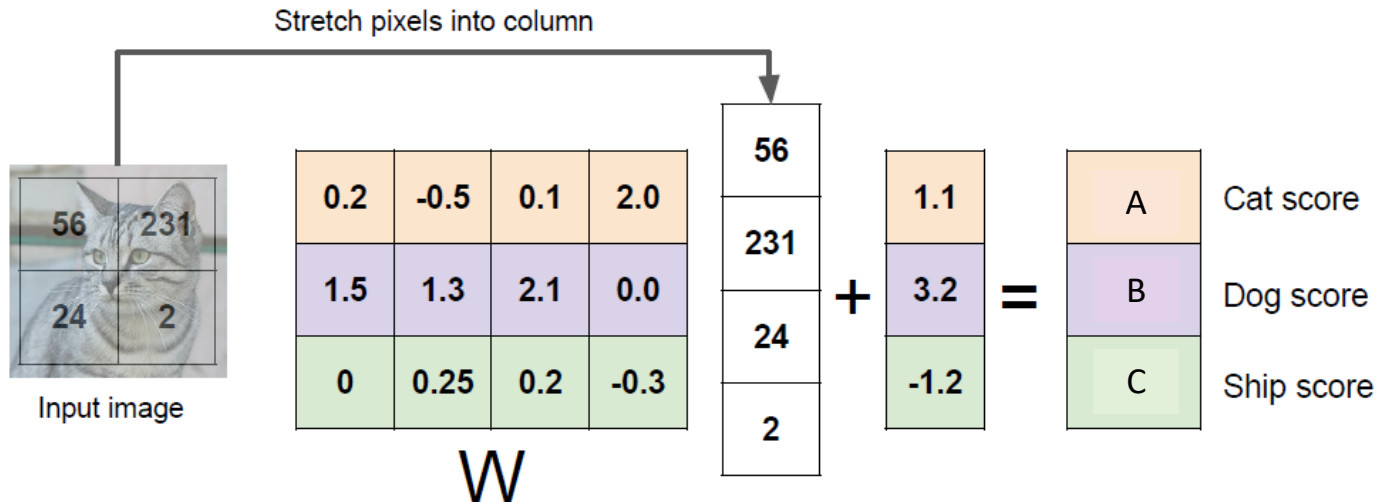


- Linear Classifier
 - Can be viewed as a parametric or algebraic approach. Why?
 - Consider that we have 10 object categories of interest
 - Let's take the input image as \mathbf{x} , and the linear classifier as \mathbf{W} . We need $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$ as a 10-dimensional output vector, indicating the score for each class.



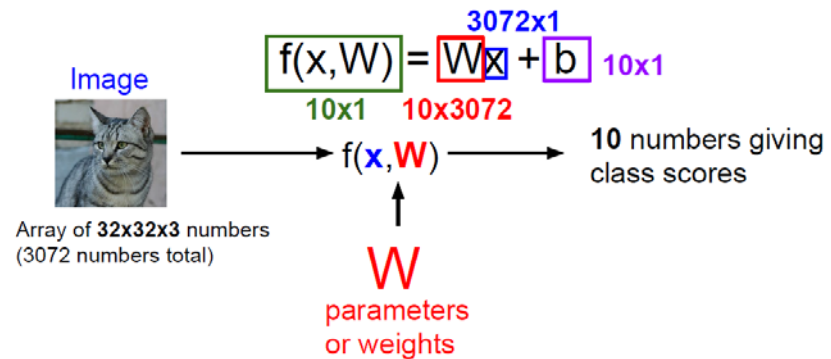
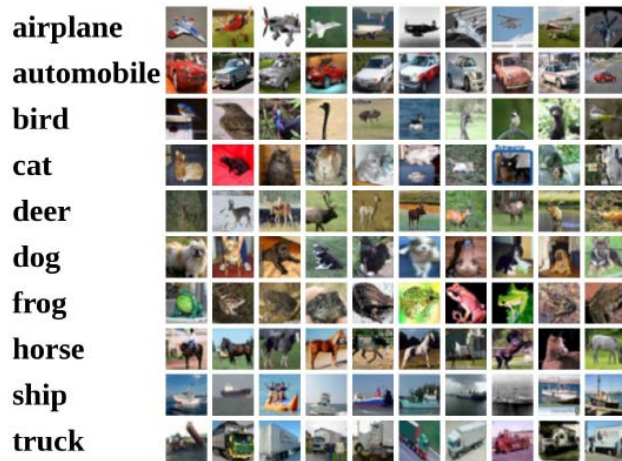
Linear Classification (cont'd)

- Linear Classifier
 - Can be viewed as a parametric or algebraic approach. Why?
 - Consider that we have 10 object categories of interest
 - Let's take the input image as \mathbf{x} , and the linear classifier as \mathbf{W} . We need $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$ as a 10-dimensional output vector, indicating the score for each class.
 - For example, an image with 2 x 2 pixels & 3 classes of interest we need to learn a linear classifier \mathbf{W} (plus a bias \mathbf{b}), so that desirable outputs $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$ can be expected.



Remarks

- Interpreting \mathbf{W} in $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$
 - Weights \mathbf{W} are learned by observing (training) data \mathbf{X} and their ground truth labels \mathbf{Y} .
 - Each row in \mathbf{W} can be viewed as an **exemplar** of the corresponding class.
 - Equivalently, we perform **inner product** between \mathbf{x} and each row of \mathbf{W} -> class similarity
 - How to determine a proper objective/loss function for deriving \mathbf{W} ?



Loss Function

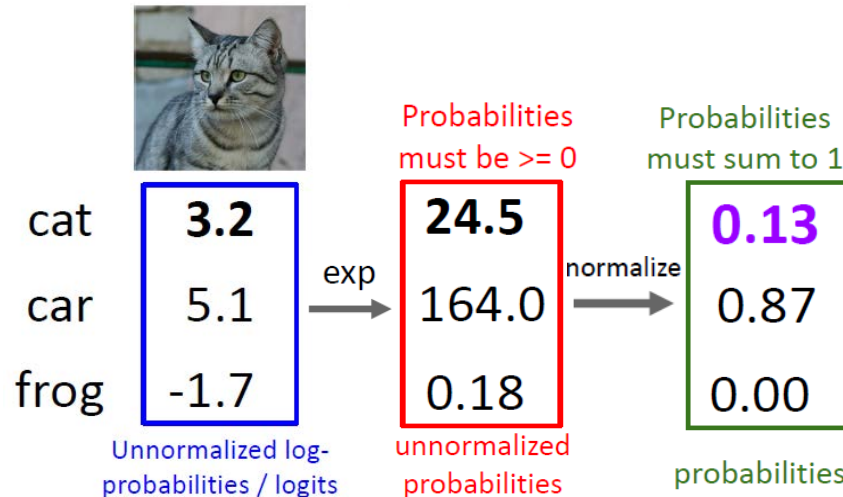
- **Cross-Entropy Loss (Multinomial Logistic Regression)**

- Interpret classifier scores as **probabilities**

- **Softmax function:**

$$P(Y = k | X = x_i) = \frac{\exp(s_k)}{\sum_j \exp(s_j)} \quad \text{with } s = f(x_i; W) \text{ as the classifier output for input } \mathbf{x}_i$$

- See example below



$$L_i = -\log P(Y = y_i | X = x_i)$$

$$L_{\text{cat}} = -\log(0.13) = 2.04$$

What about this L?

What are its possibly **min/max** value??

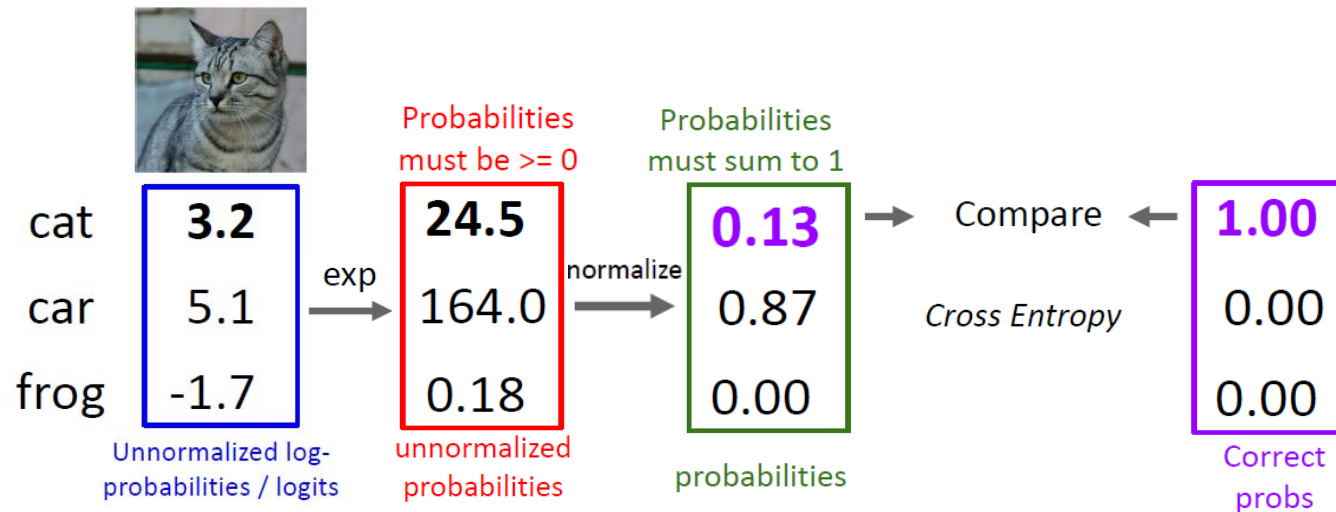
- **Cross-Entropy Loss (cont'd)**

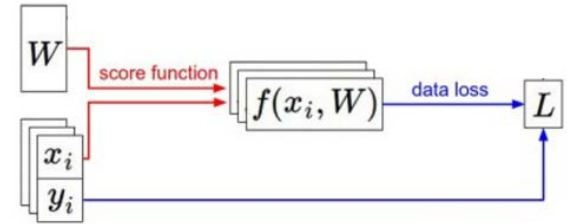
- **Softmax function:**

$$P(Y = k | X = x_i) = \frac{\exp(s_k)}{\sum_j \exp(s_j)} \quad \text{with } s = f(x_i; W) \text{ as the classifier output for input } \mathbf{x}_i$$

→ $L_i = -\log P(Y = y_i | X = x_i)$ or $L_i = -\log \left(\frac{\exp(s_{y_i})}{\sum_j \exp(s_j)} \right)$

- **(Binary) Cross Entropy Loss (or L_{BCE} ; see example below):**





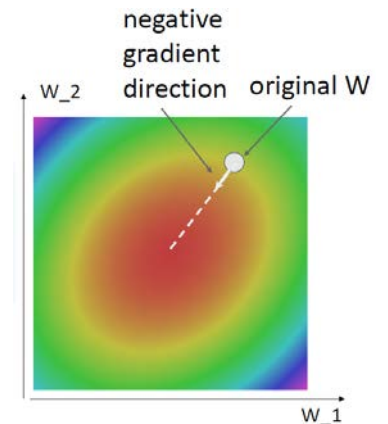
- **Searching for W from L_{BCE}**

- Computing **gradients**:
Following the slope to reach the (hopefully global) minimum for W .

- **Gradient Descent** via numeric or analytic gradients:

- Iteratively step in the direction of the **negative gradient** & search for W
- Hyperparameters: weight initialization, # of steps, learning rate, etc.

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```



- **Stochastic Gradient Descent**

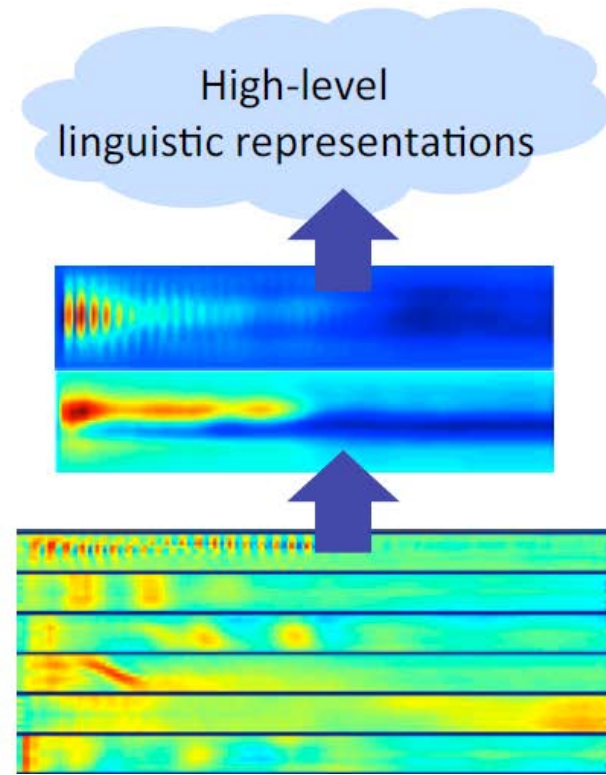
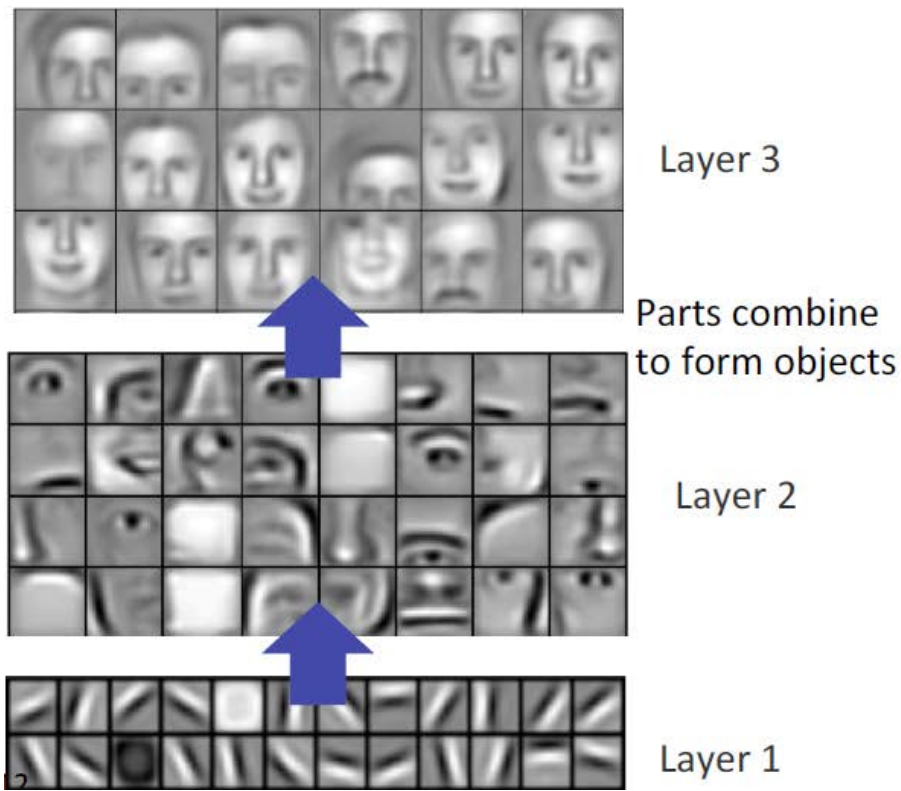
- Full sum in L is **expensive when large N**
- *Approximate* sum using a minibatch of instances (e.g., 32, 64, 128, etc.)
- Additional hyperparameters of batch size and data sampling

```
# Stochastic gradient descent
w = initialize_weights()
for t in range(num_steps):
    minibatch = sample_data(data, batch_size)
    dw = compute_gradient(loss_fn, minibatch, w)
    w -= learning_rate * dw
```



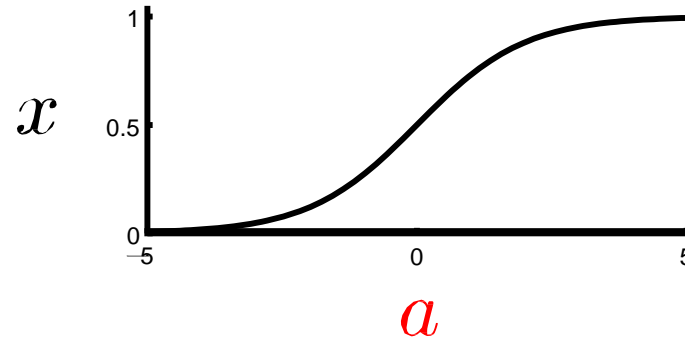
Hierarchical Representation Learning

- Successive model layers learn deeper intermediate representations.



Let's Take a Closer Look...

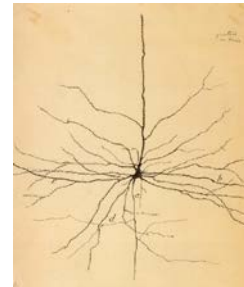
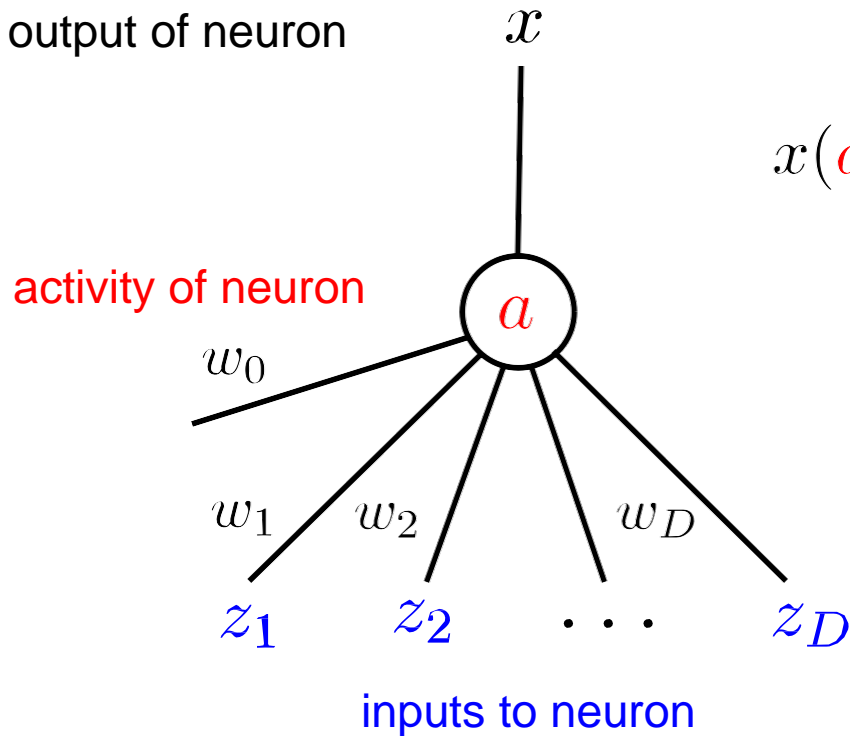
- A single neuron in a single layer



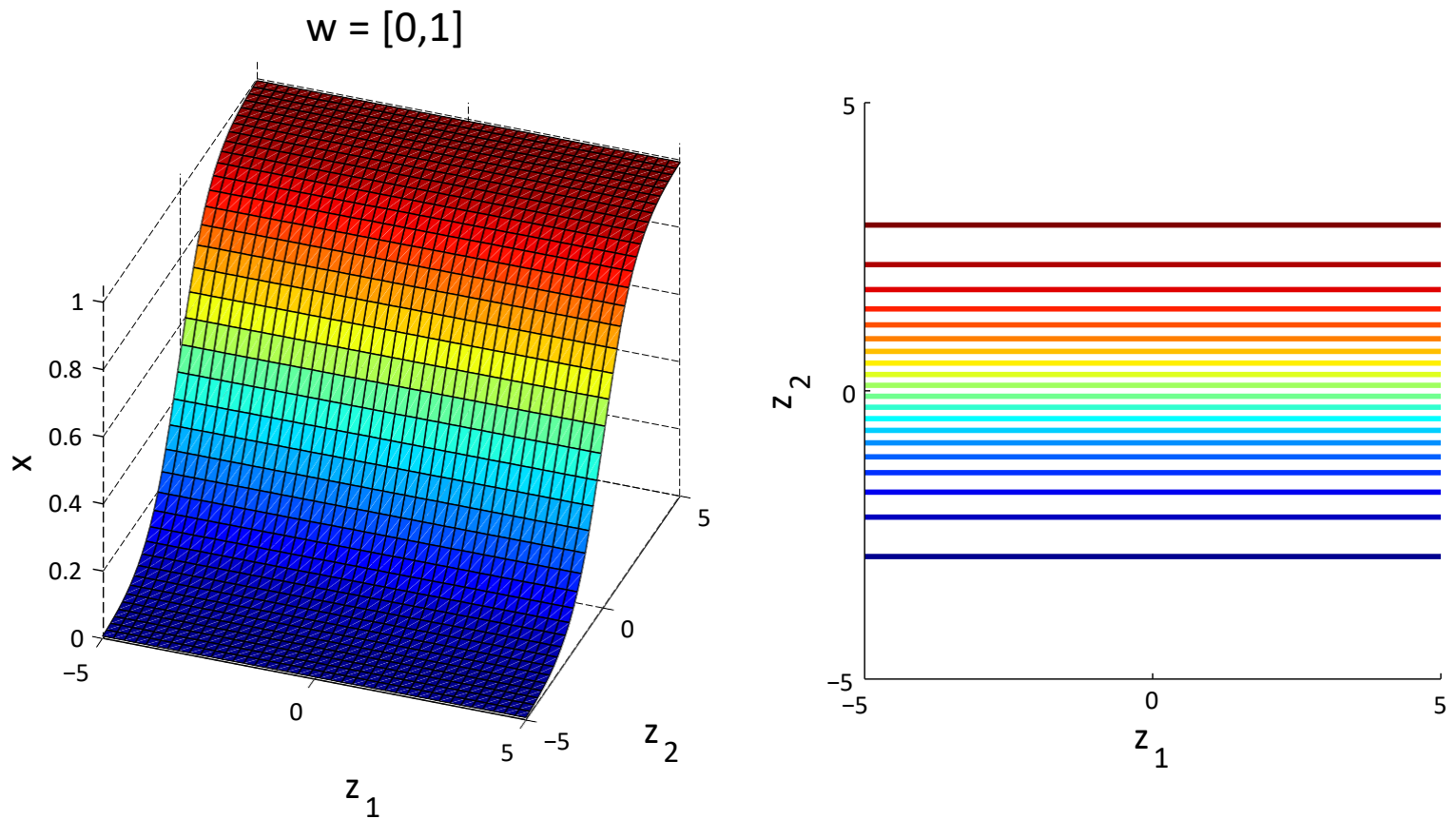
$$x(a) = \frac{1}{1 + \exp(-a)} \quad x \in (0, 1)$$

$$a = w_0 + \sum_{d=1}^D w_d z_d$$

$$= \sum_{d=0}^D w_d z_d$$

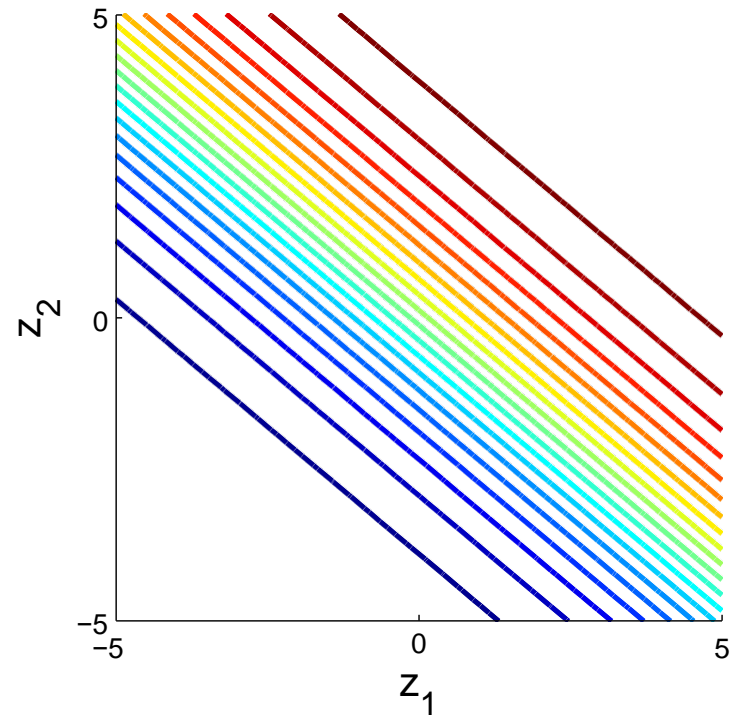
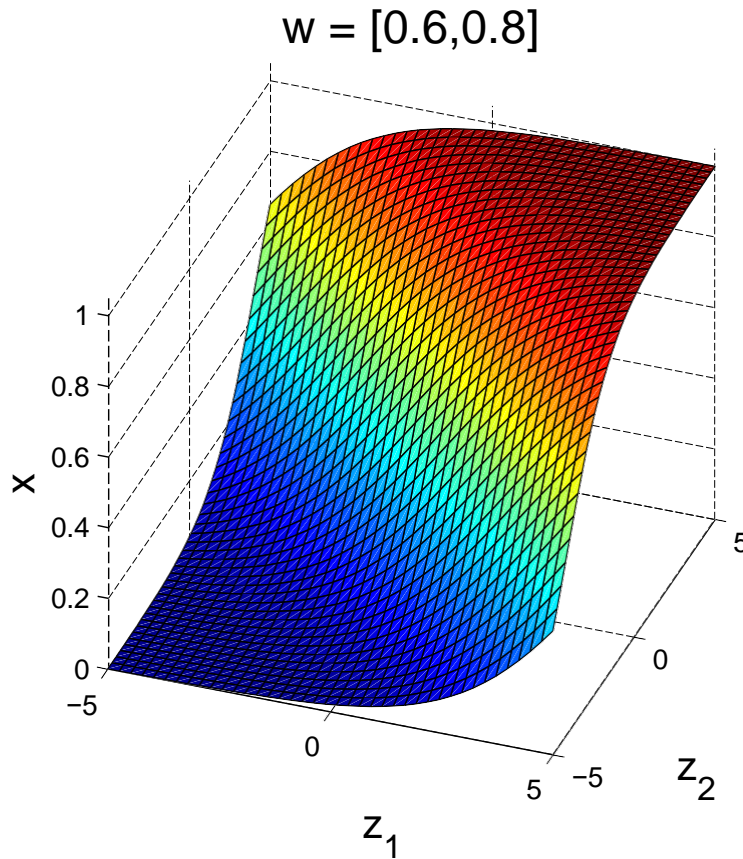


Input-Output Function of a Single Neuron



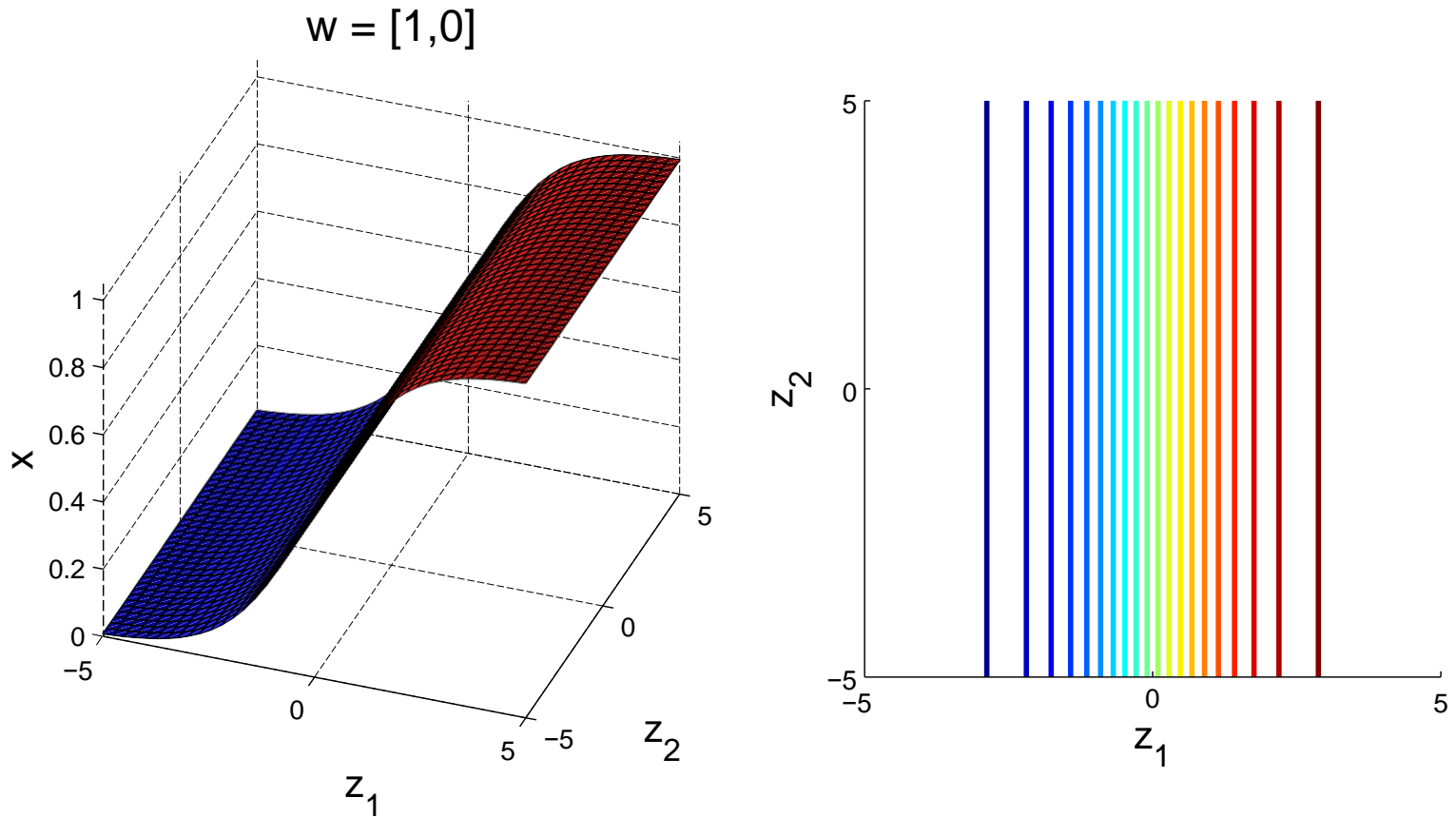
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



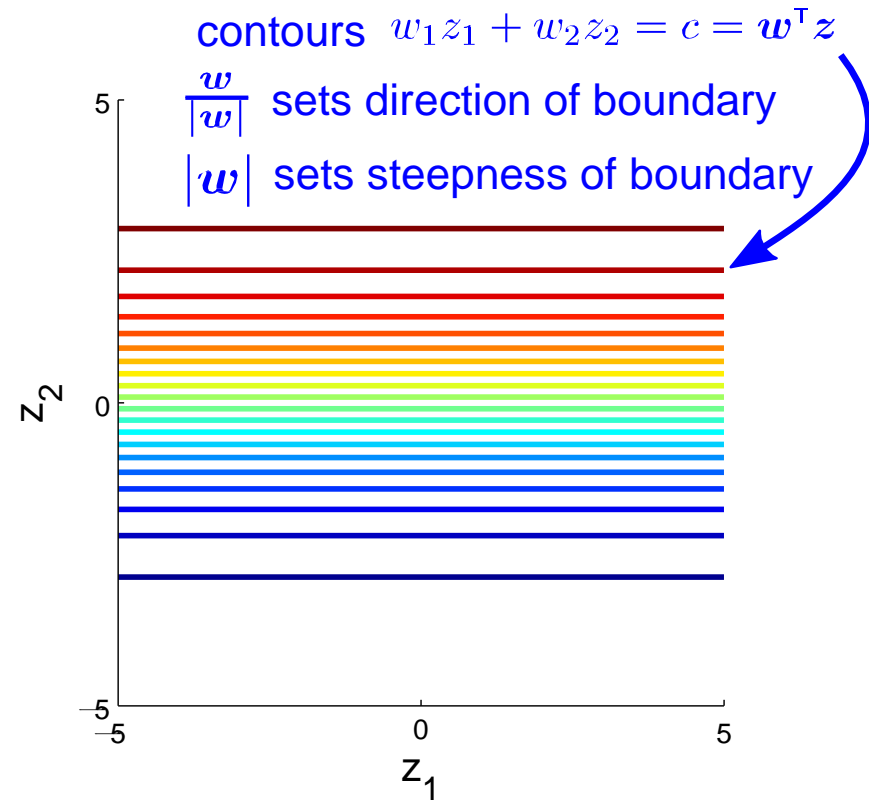
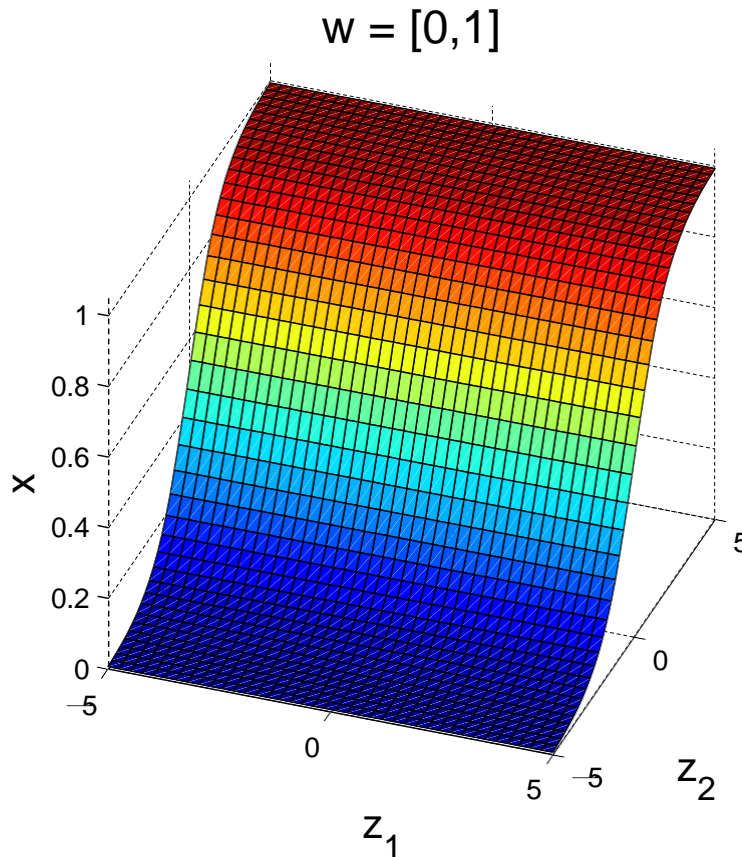
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron



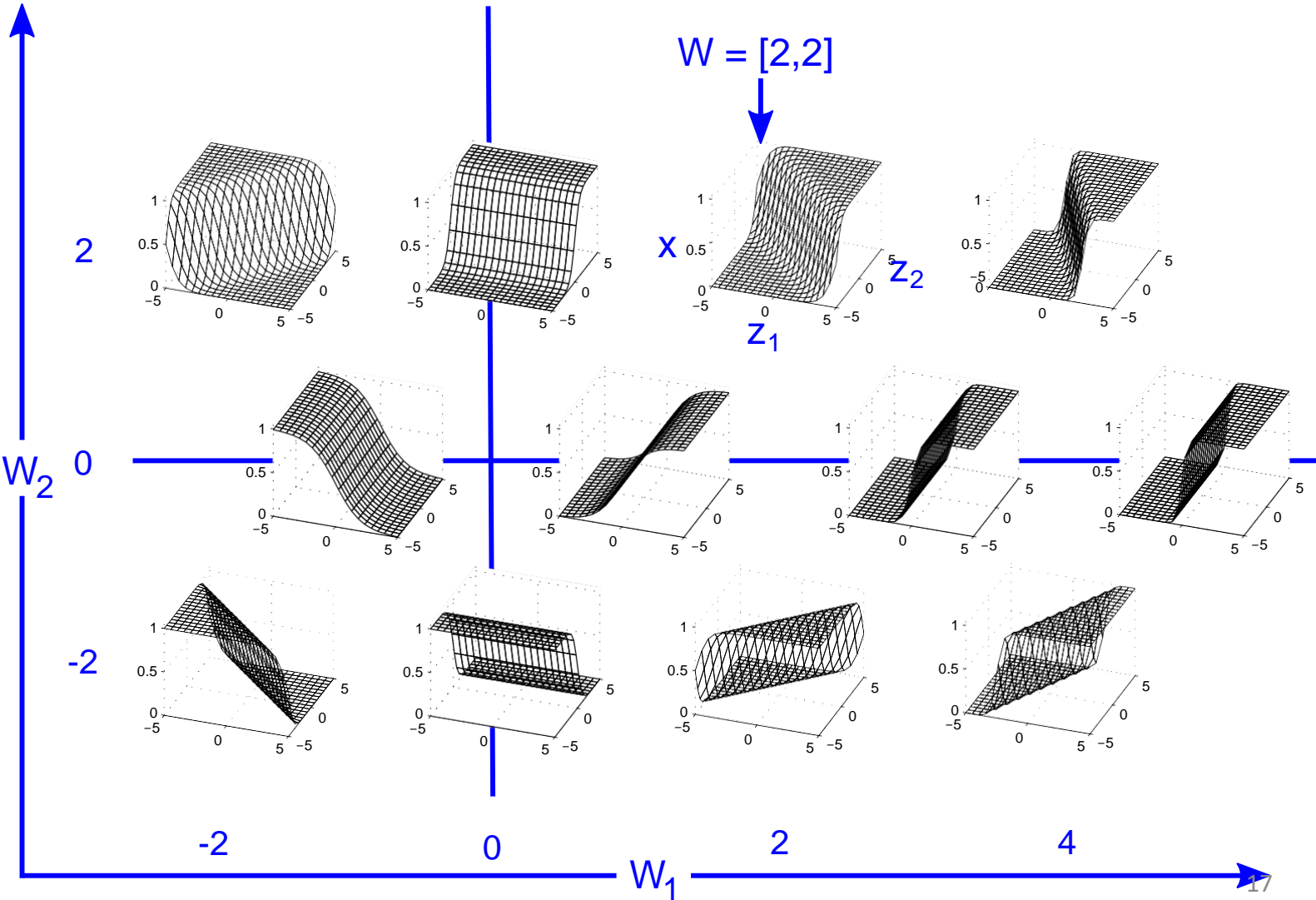
$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Input-Output Function of a Single Neuron (cont'd)

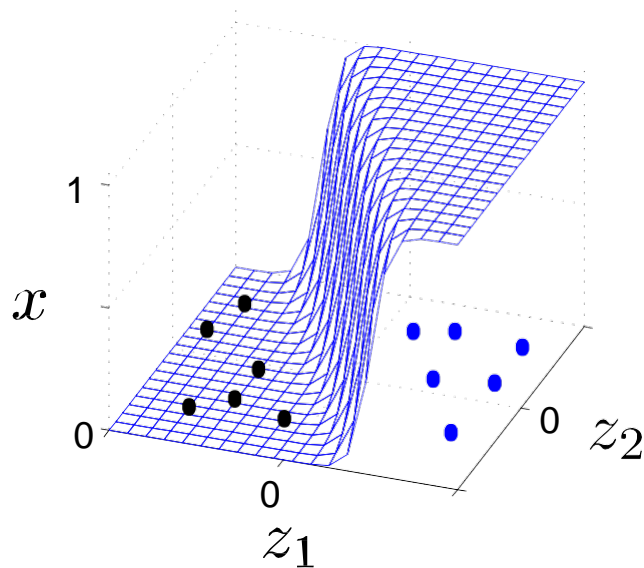


$$x(z_1, z_2) = \frac{1}{1 + \exp(-w_1 z_1 - w_2 z_2)}$$

Weight Space of a Single Neuron



Training a single neuron



training data

$$\{\mathbf{z}^{(n)}\}_{n=1}^N \quad \{t^{(n)}\}_{n=1}^N$$

inputs

class labels

objective function:

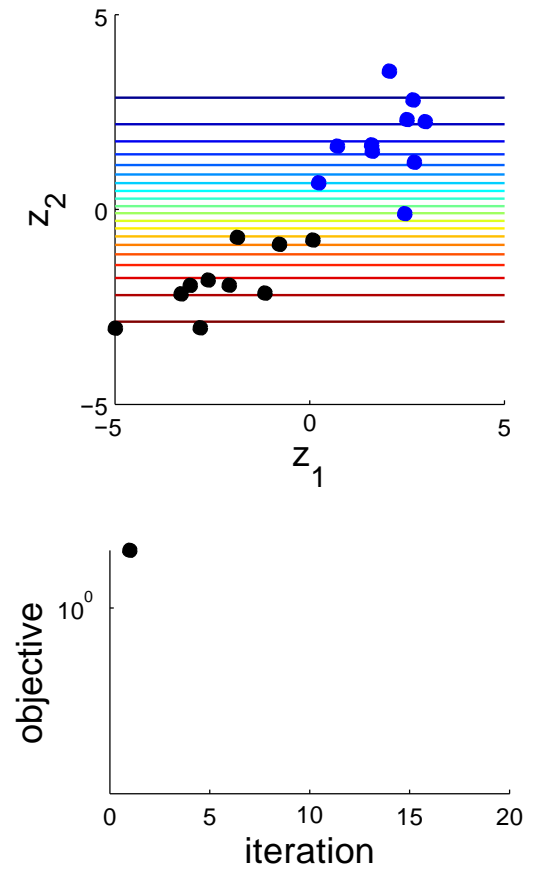
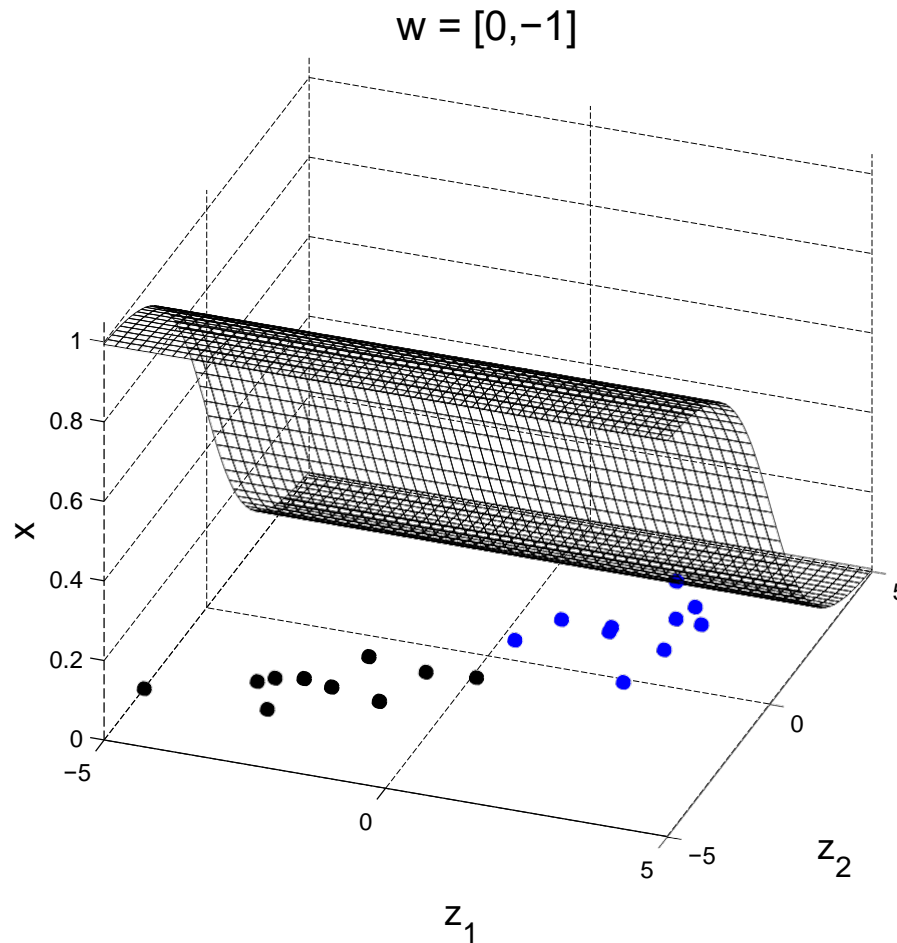
$$G(\mathbf{w}) = - \sum_n [t^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - t^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))]$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_i w_i^2 \quad \text{regulariser discourages the network using extreme weights}$$

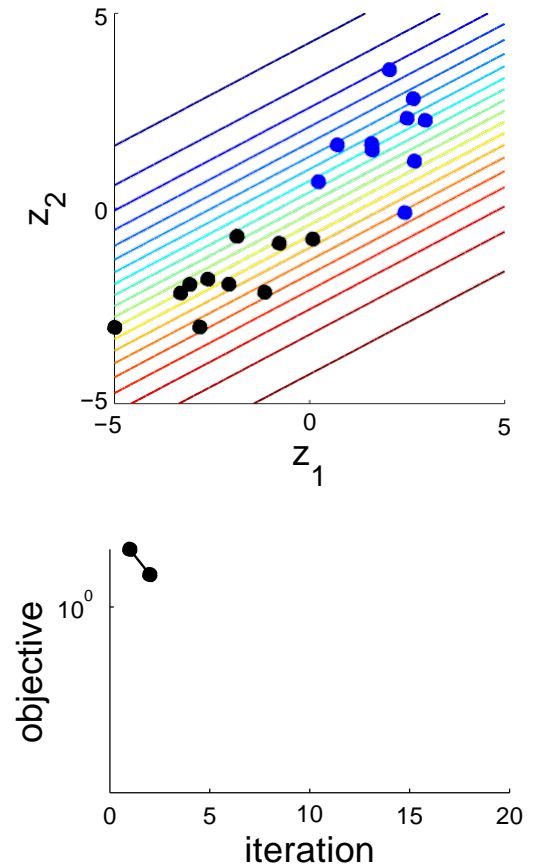
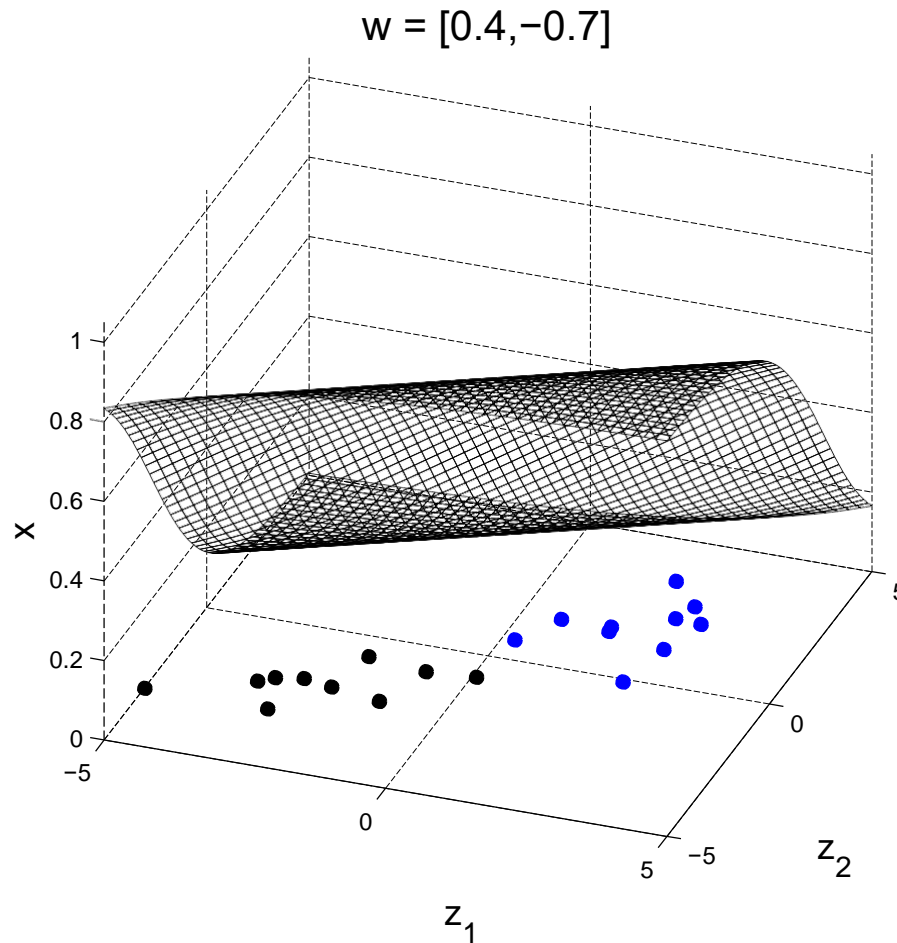
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} M(\mathbf{w}) = \arg \min_{\mathbf{w}} [G(\mathbf{w}) + \alpha E(\mathbf{w})]$$

$$\frac{d}{d\mathbf{w}} M(\mathbf{w}) = - \sum_n (t^{(n)} - x^{(n)}) \mathbf{z}^{(n)} + \alpha \mathbf{w} \quad \text{weight decay - shrinks weights towards zero}$$

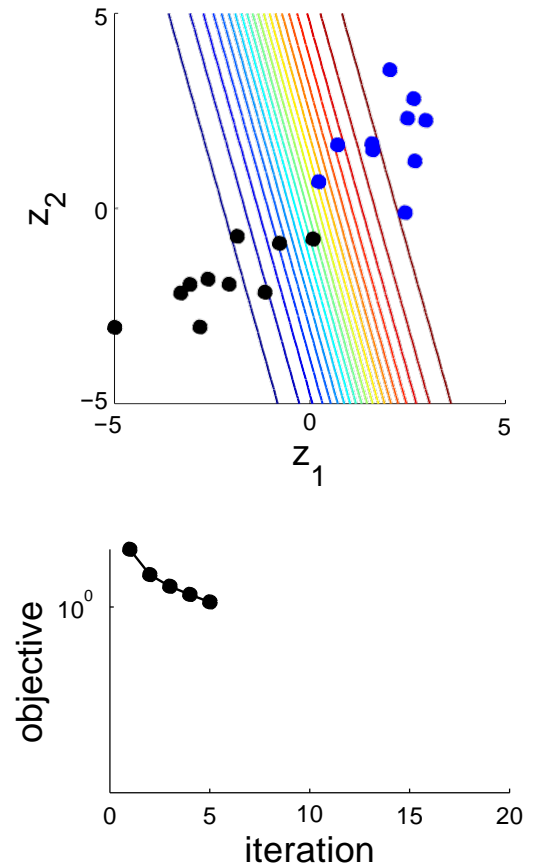
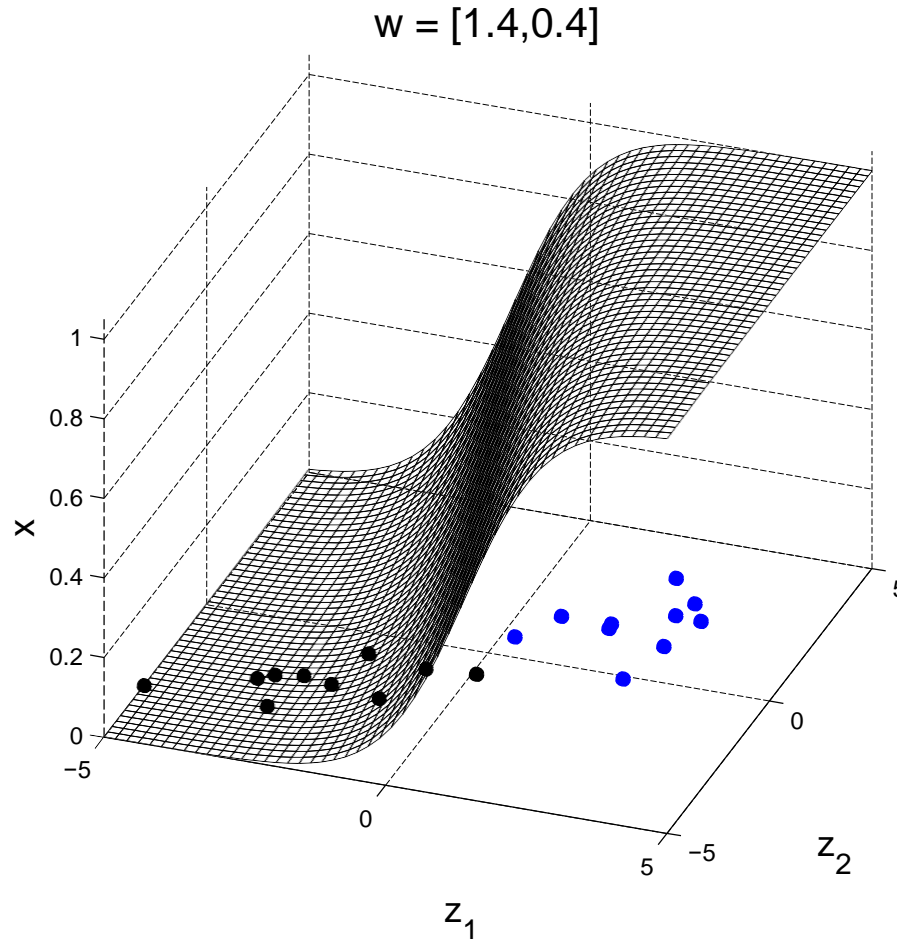
Training a Single Neuron



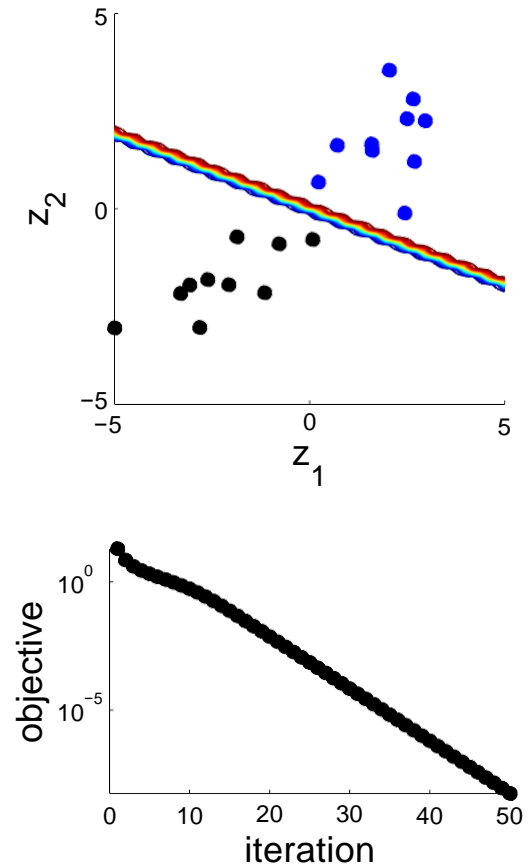
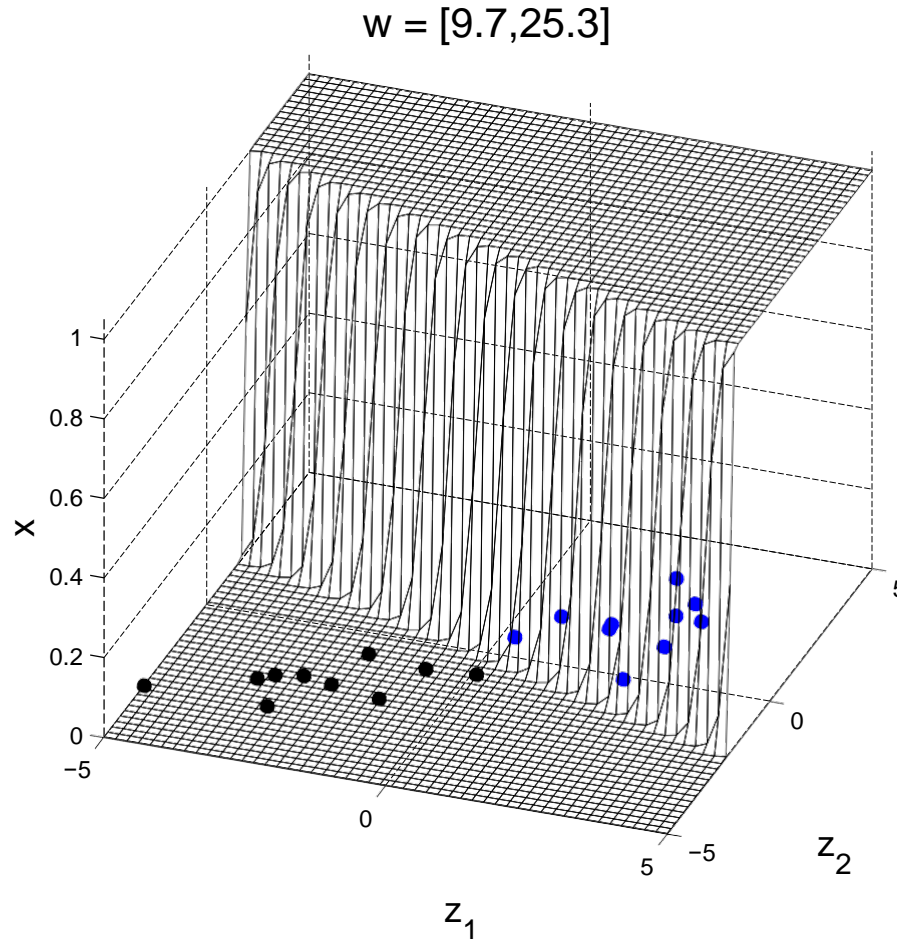
Training a Single Neuron



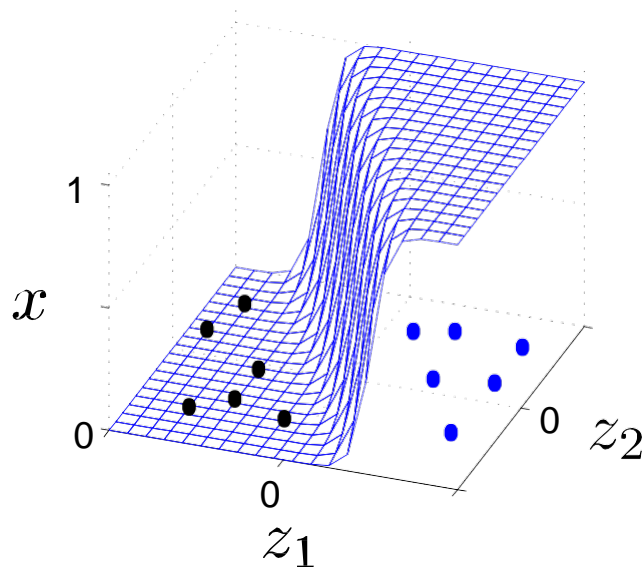
Training a Single Neuron



Training a Single Neuron



Overfitting and Weight Decay



training data

$$\{\mathbf{z}^{(n)}\}_{n=1}^N \quad \{t^{(n)}\}_{n=1}^N$$

inputs

class labels

objective function:

$$G(\mathbf{w}) = - \sum_n [t^{(n)} \log x(\mathbf{z}^{(n)}; \mathbf{w}) + (1 - t^{(n)}) \log (1 - x(\mathbf{z}^{(n)}; \mathbf{w}))]$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_i w_i^2$$

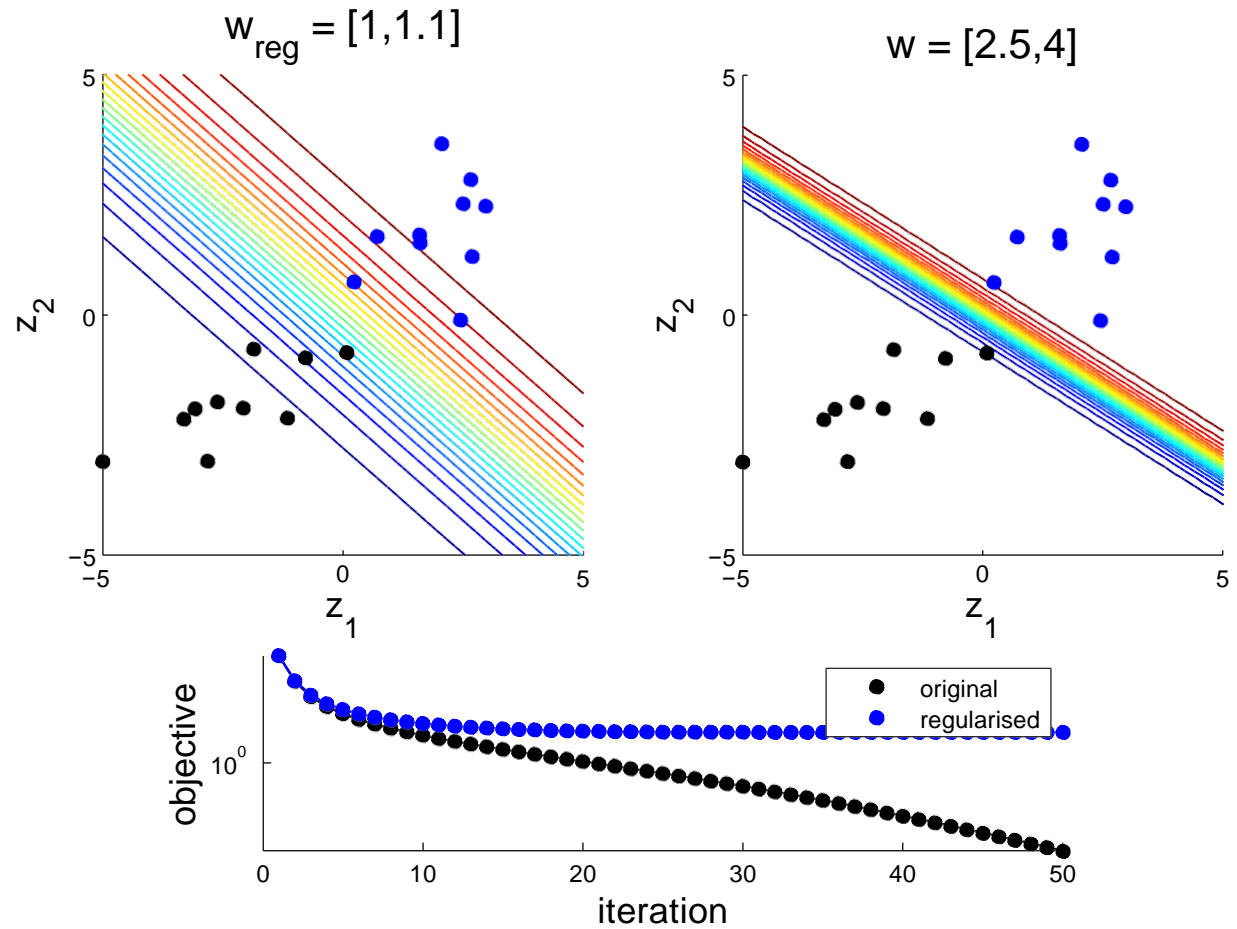
regulariser discourages the network using extreme weights

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} M(\mathbf{w}) = \arg \min_{\mathbf{w}} [G(\mathbf{w}) + \alpha E(\mathbf{w})]$$

$$\frac{d}{d\mathbf{w}} M(\mathbf{w}) = - \sum_n (t^{(n)} - x^{(n)}) \mathbf{z}^{(n)} + \alpha \mathbf{w}$$

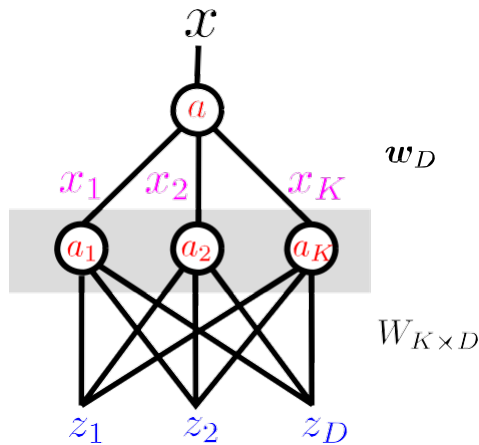
weight decay - shrinks weights towards zero

Training a Single Neuron (cont'd)



Training a Neural Network with Two Hidden Layers

Networks with hidden layers can be fit using gradient descent using an algorithm called **back-propagation**.



$$x(a) = \frac{1}{1 + \exp(-a)}$$

$$a = \sum_{k=1}^K w_k x_k$$

$$x(a_k) = \frac{1}{1 + \exp(-a_k)}$$

$$a_k = \sum_{d=1}^D W_{k,d} z_d$$

objective function:

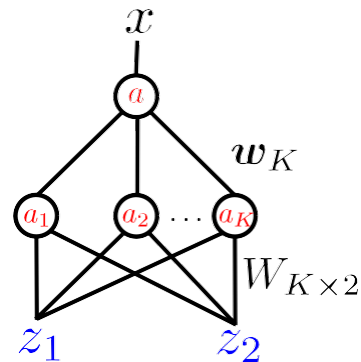
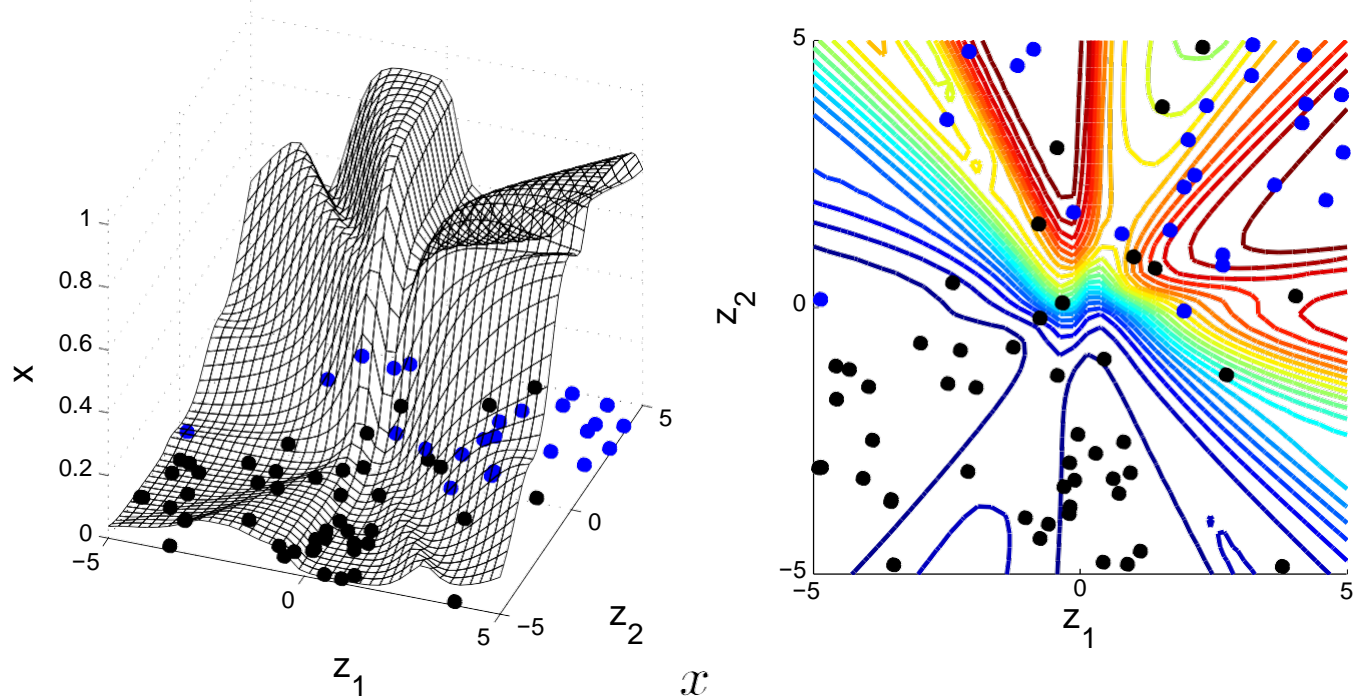
$$G(W, \mathbf{w}) = - \sum_n [t^{(n)} \log x^{(n)} + (1 - t^{(n)}) \log (1 - x^{(n)})] \text{ likelihood same as before}$$

$$E(W, \mathbf{w}) = \frac{1}{2} \sum_i w_i^2 + \frac{1}{2} \sum_{ij} W_{ij}^2 \quad \text{regulariser discourages extreme weights}$$

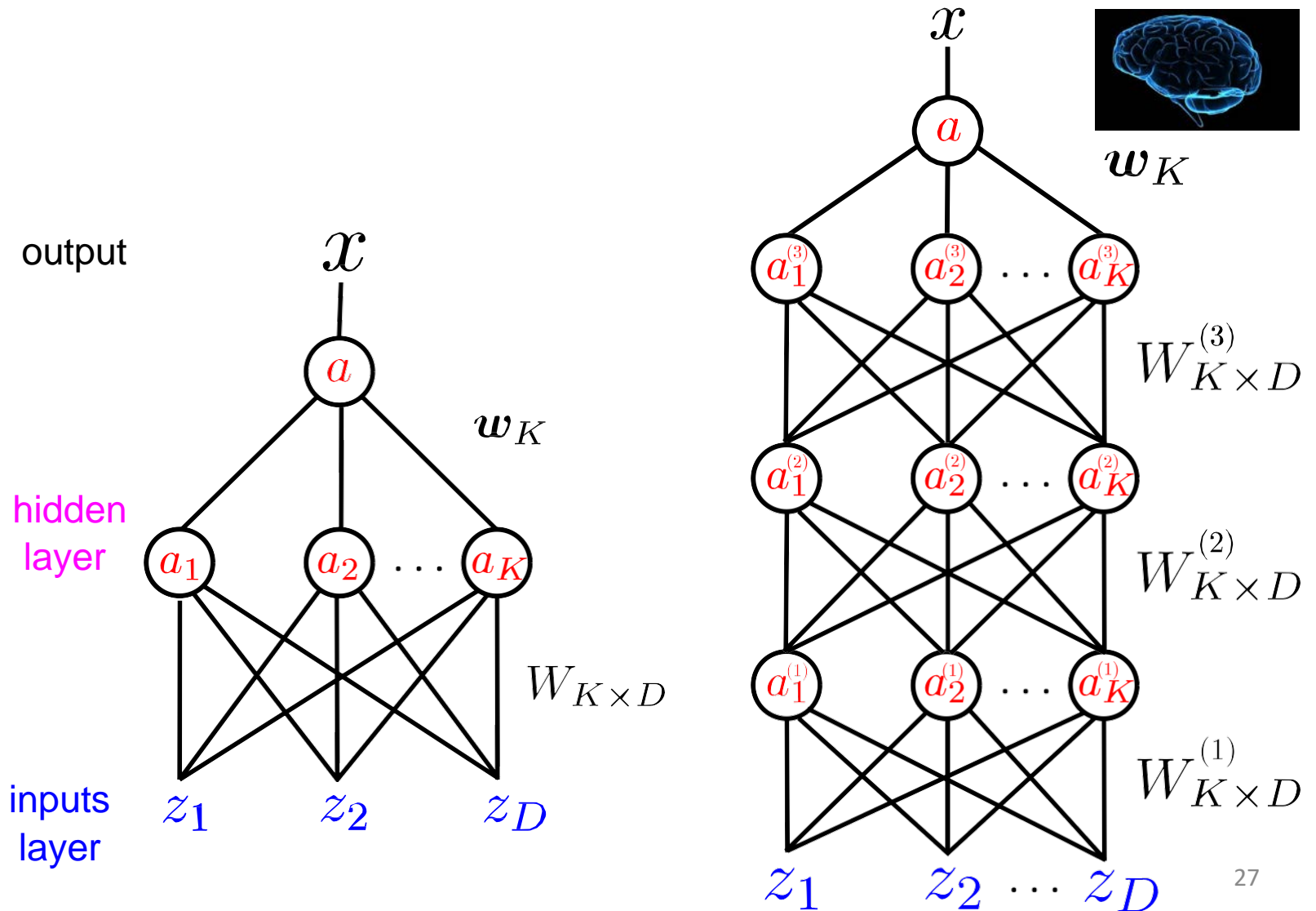
$$\{W, \mathbf{w}^*\} = \arg \min_{W, \mathbf{w}} M(W, \mathbf{w}) = \arg \min_{W, \mathbf{w}} [G(W, \mathbf{w}) + \alpha E(W, \mathbf{w})]$$

$$\begin{aligned} \frac{dG(W, \mathbf{w})}{dW_{ij}} &= \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{dW_{ij}} = \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{da^{(n)}}{dW_{ij}} \\ &= \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{da^{(n)}} \frac{da^{(n)}}{dx_i^{(n)}} \frac{dx_i^{(n)}}{dW_{ij}} = \sum_n \frac{dG(W, \mathbf{w})}{dx^{(n)}} \frac{dx^{(n)}}{da^{(n)}} \frac{da^{(n)}}{dx_i^{(n)}} \frac{da_i^{(n)}}{dW_{ij}} \end{aligned}$$

Training a Neural Network with a Single Hidden Layer



Hierarchical Models with Many Layers

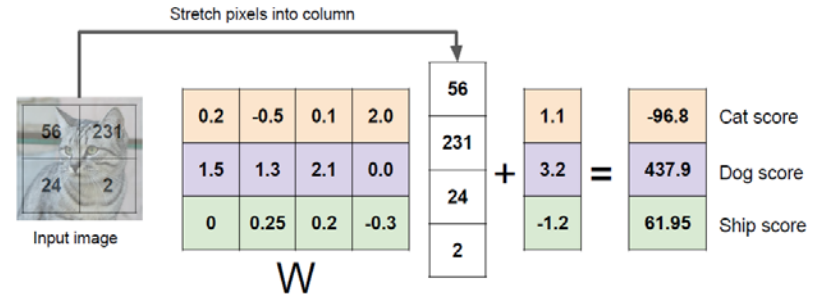
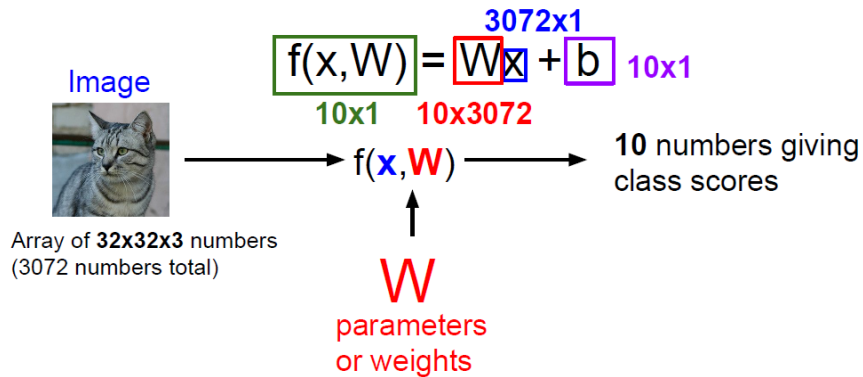


What'd Be Covered in This Crash Course...

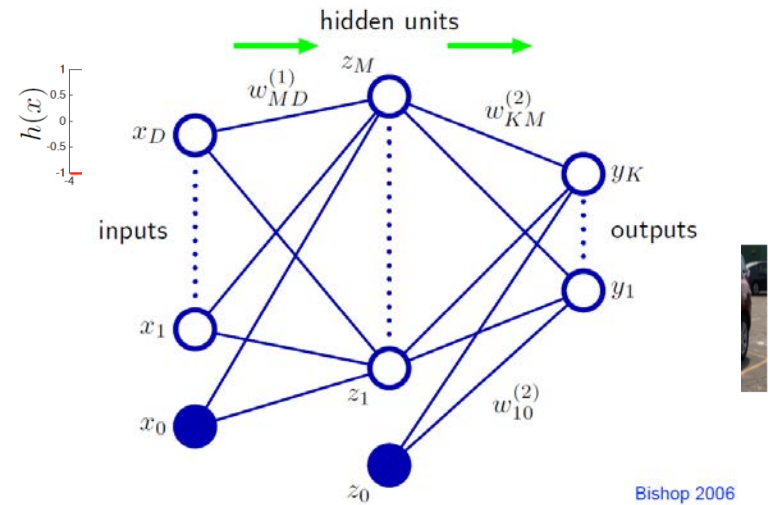
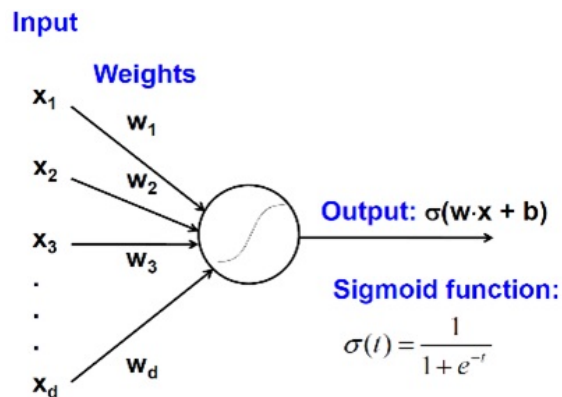
- **Learning-based Computer Vision**
 - From Linear to Non-Linear Classifiers
- **Start of Deep Learning for Computer Vision**
 - Convolutional Neural Networks
 - Self-Supervised Learning
 - Segmentation & Detection
- **Generative Models**
 - Autoencoder, Variational Autoencoder, Generative Adversarial Networks & Diffusion Models
- **Sequence-to-Sequence Learning**
 - Attention is All You Need: Transformer
- **Vision & Language Foundation Models**
 - Image-to-Text vs. Text-to-Image
 - Parameter-Efficient Fine-tuning

Recap: Linear Classification to Neural Nets

- Linear Classifier

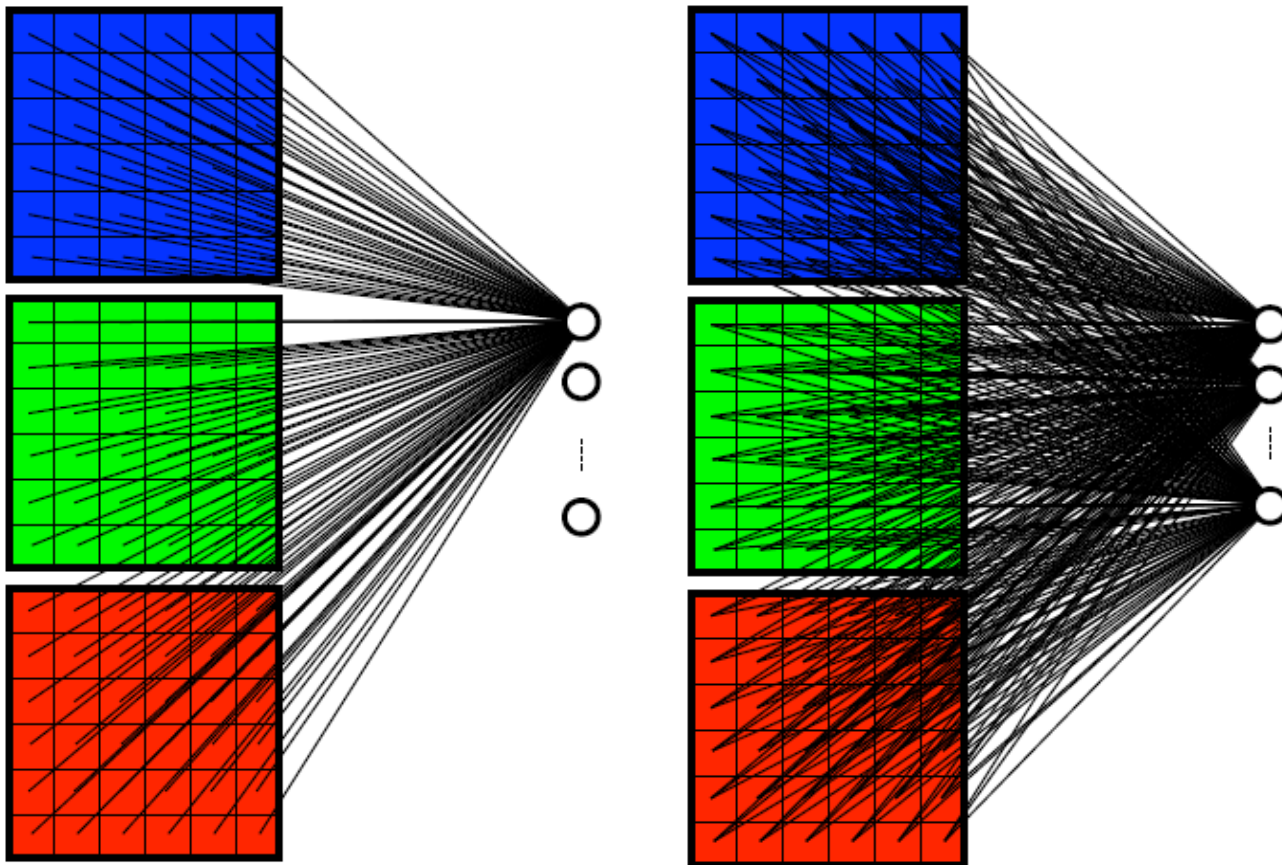


- Neural Network (Multilayer Perceptron)



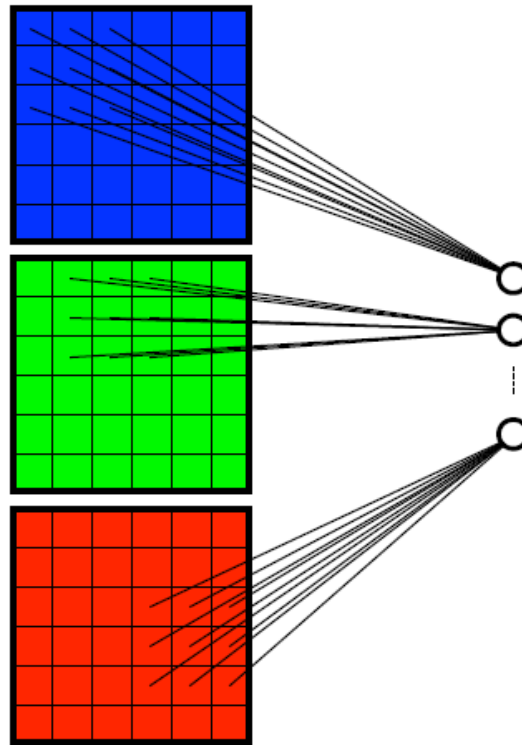
Convolutional Neural Networks

- How many weights for MLPs for images?



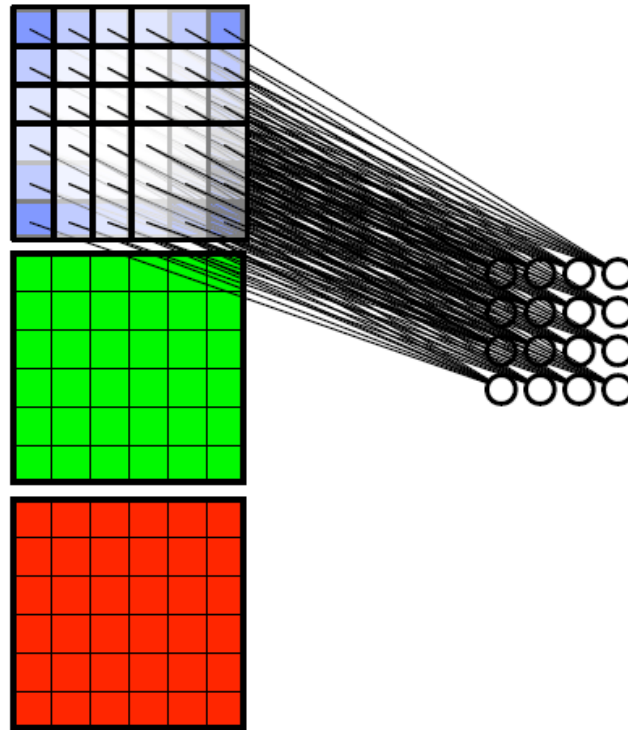
Convolutional Neural Networks

- Property I of CNN: Local Connectivity
 - Each neuron takes info only from a **neighborhood** of pixels.



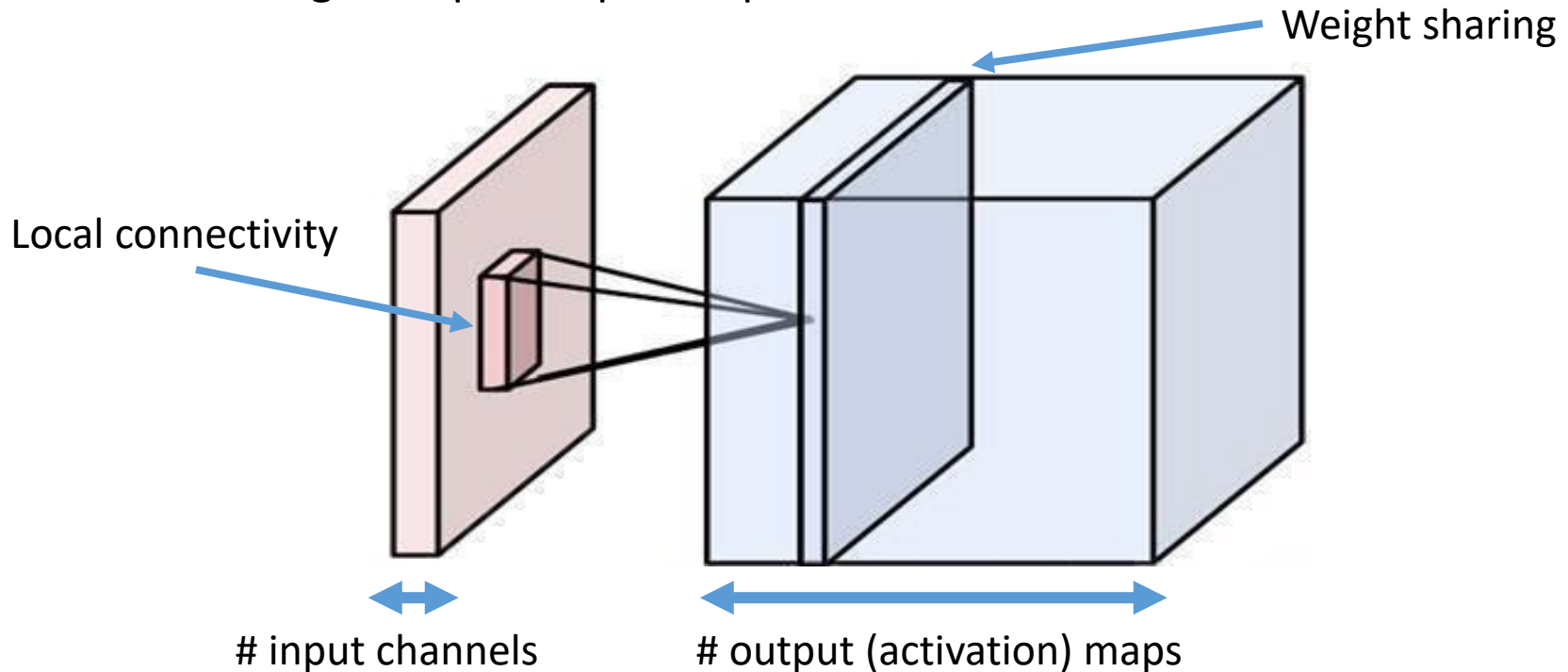
Convolutional Neural Networks

- Property II of CNN: Weight Sharing
 - Neurons connecting all neighborhoods have **identical** weights.

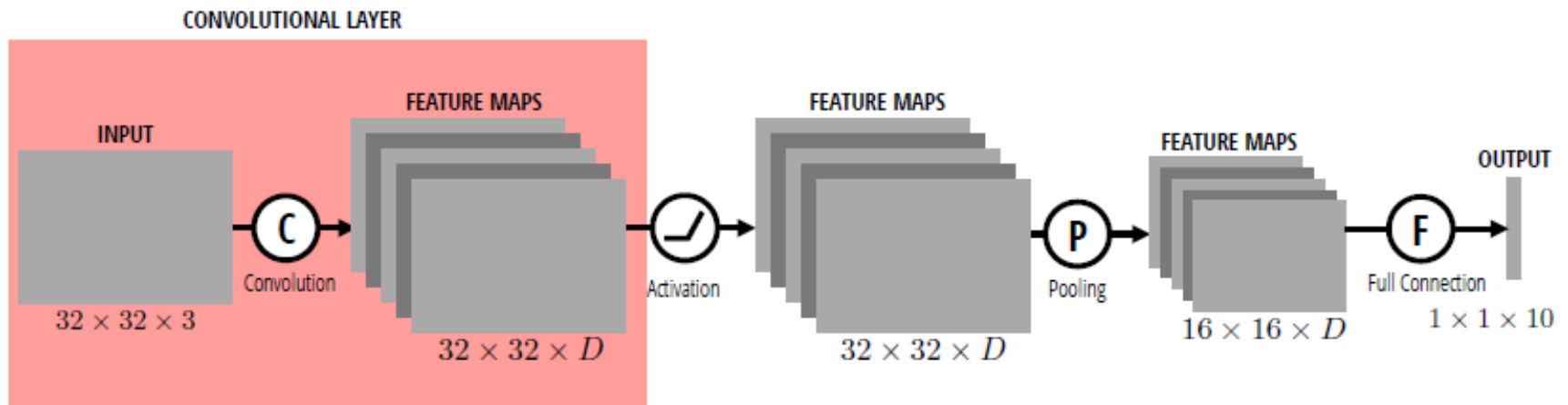


Putting them together → CNN

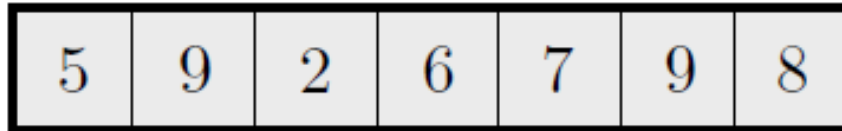
- Local connectivity
- Weight sharing
- Handling multiple input channels
- Handling multiple output maps



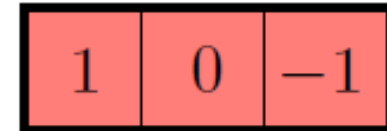
Convolution Layer in CNN



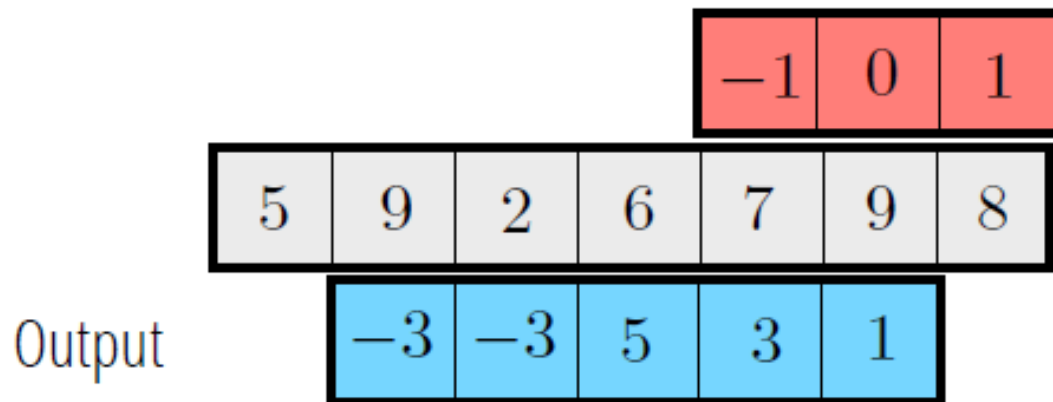
What is a Convolution?



Signal



Filter



Output

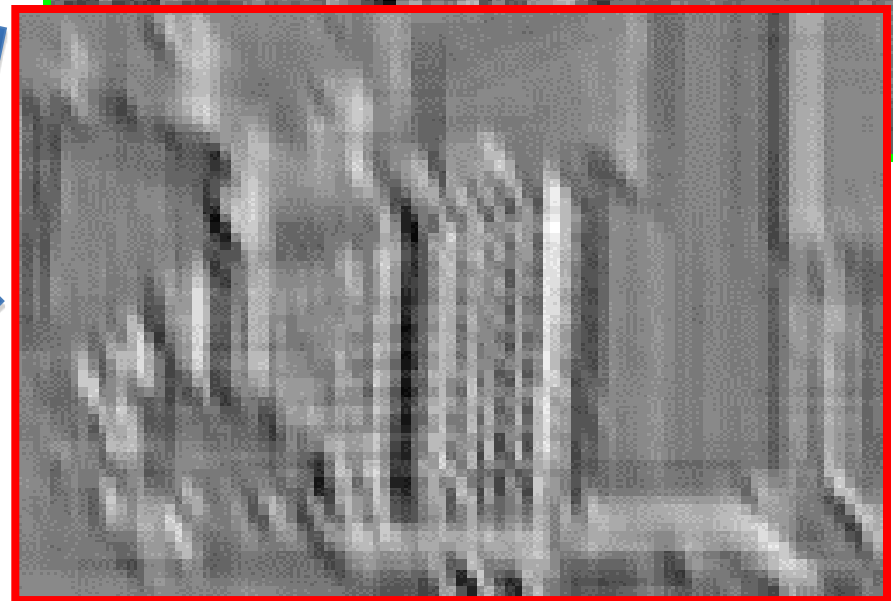
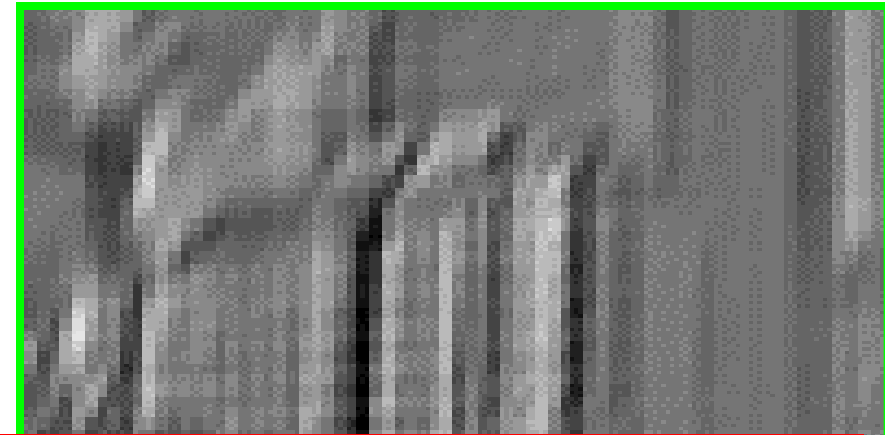
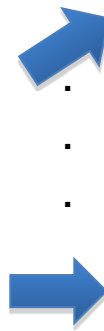
Convolution is a local linear operator

What is a Convolution?

- Weighted moving sum



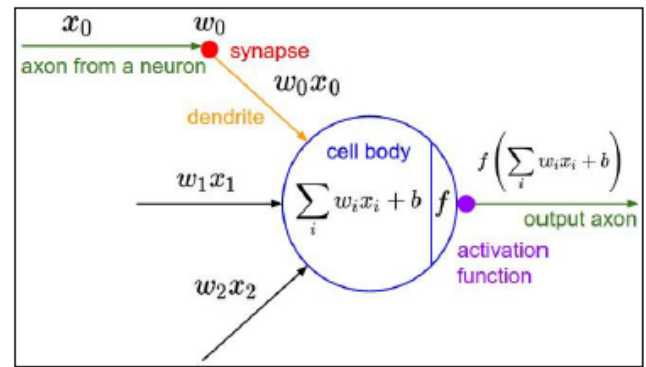
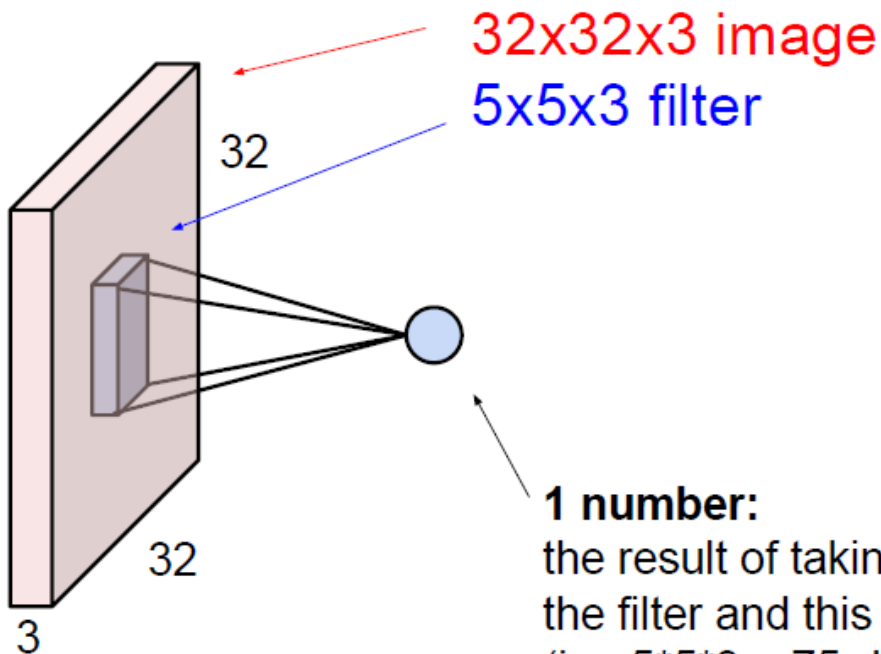
Input



Feature Activation Map

Putting them together

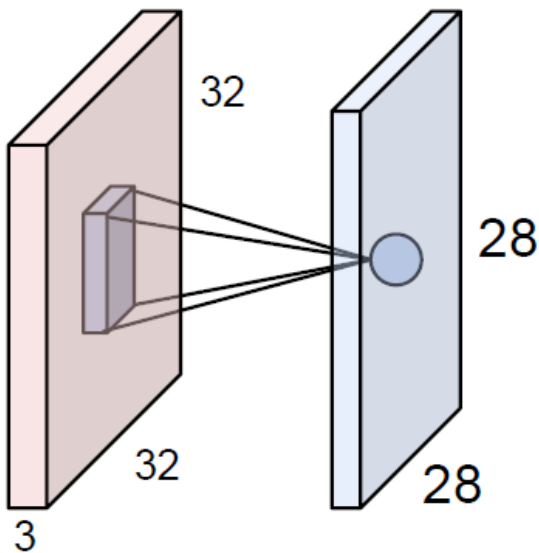
- The brain/neuron view of CONV layer



It's just a neuron with local connectivity...

Putting them together (cont'd)

- The brain/neuron view of CONV layer



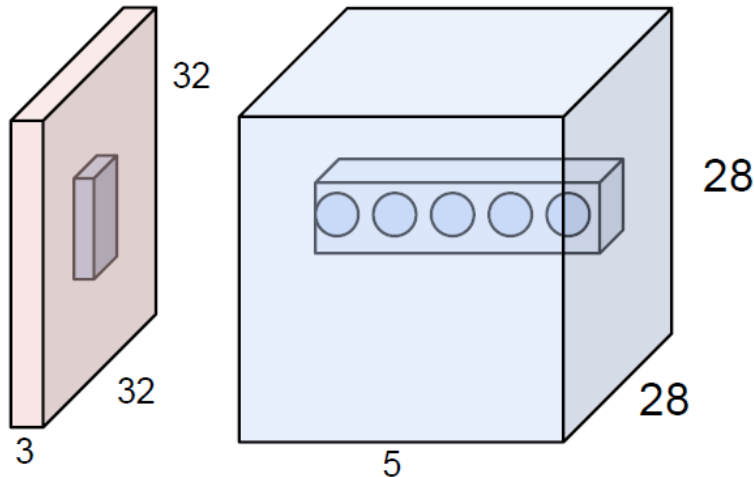
An activation map is a 28x28 sheet of neuron outputs:

1. Each is connected to a small region in the input
2. All of them share parameters

“5x5 filter” -> “5x5 receptive field for each neuron”

Putting them together (cont'd)

- The brain/neuron view of CONV layer

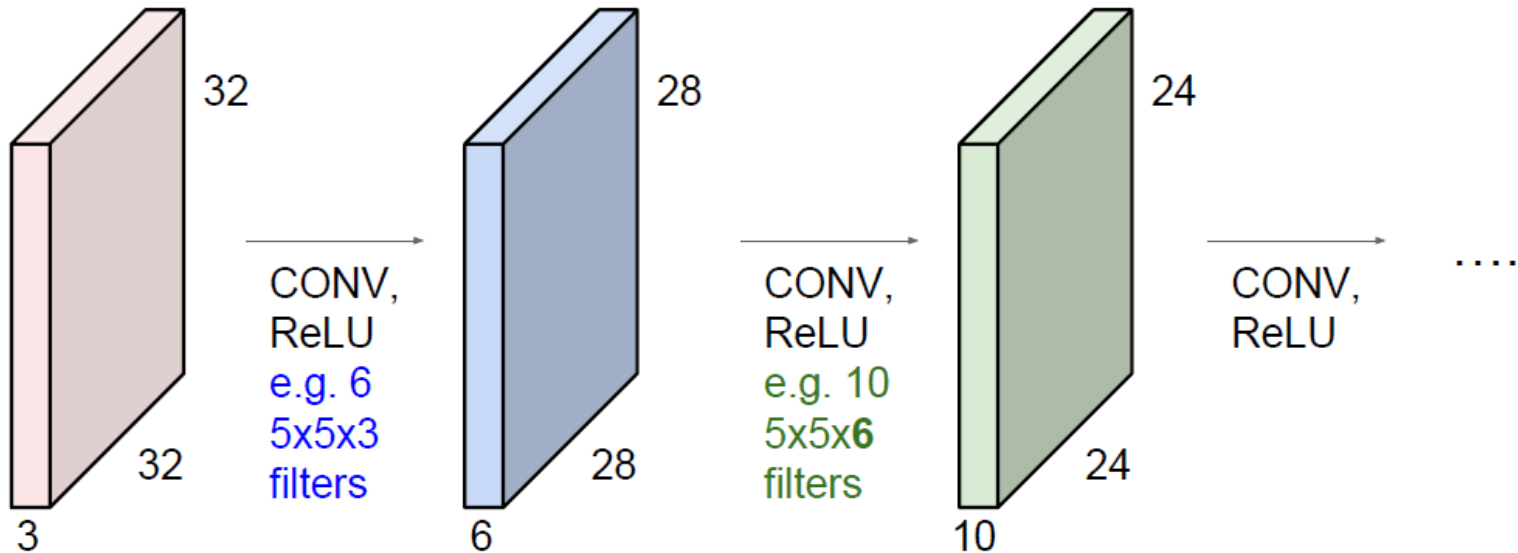


E.g. with 5 filters,
CONV layer consists of
neurons arranged in a 3D grid
(28x28x5)

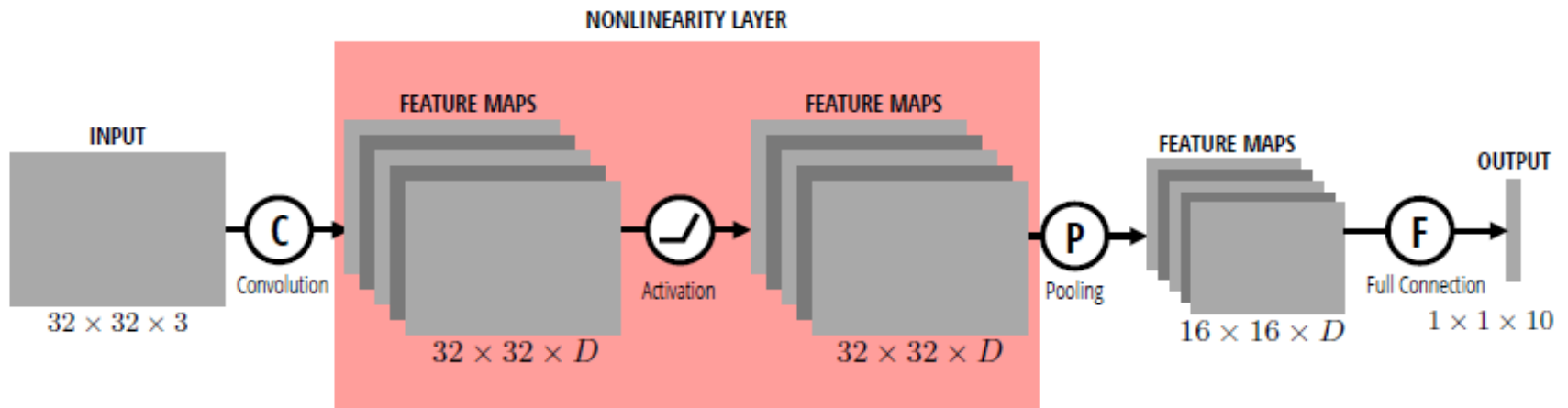
There will be 5 different
neurons all looking at the same
region in the input volume

Putting them together (cont'd)

- Image input with 32 x 32 pixels convolved repeatedly with 5 x 5 x 3 filters shrinks volumes spatially (32 -> 28 -> 24 -> ...).

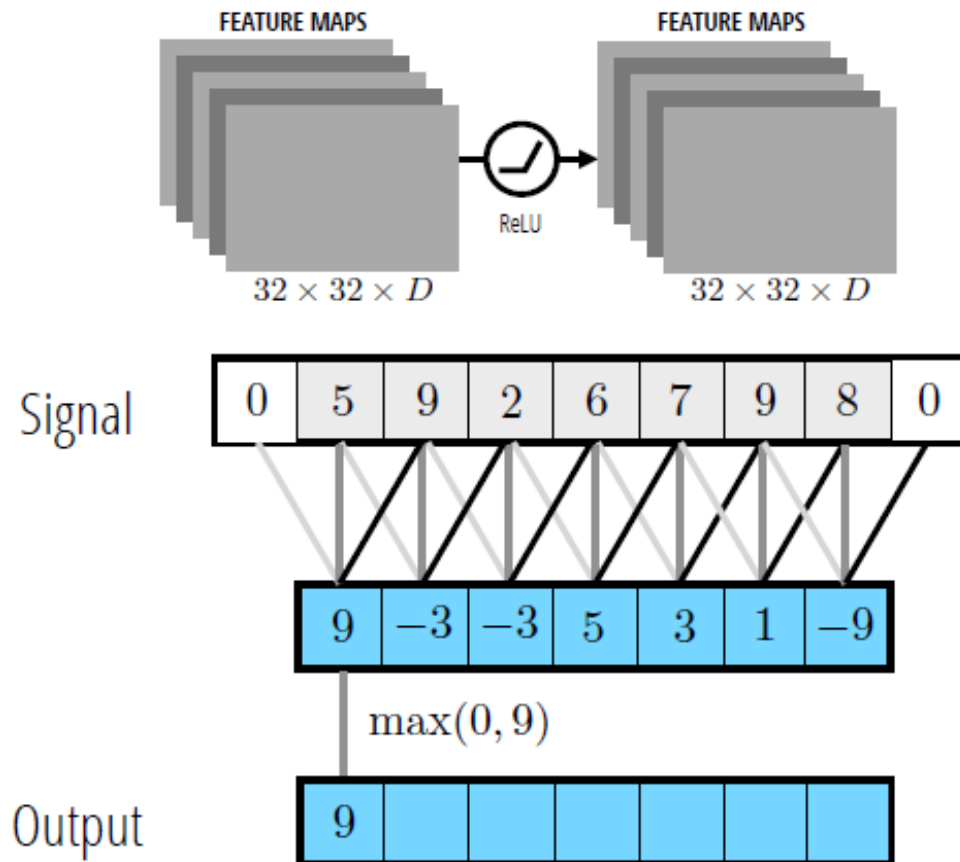


Nonlinearity Layer in CNN



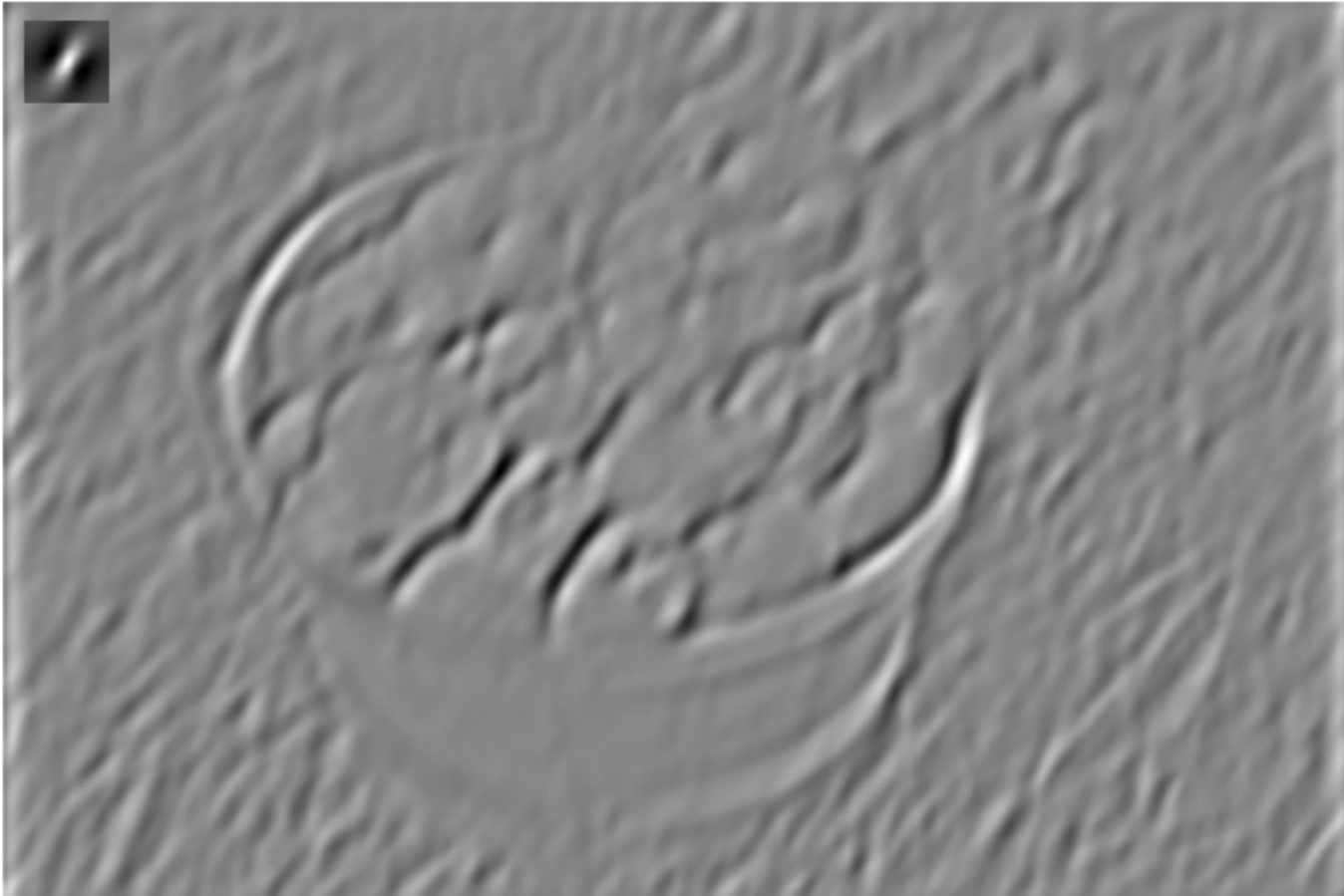
Nonlinearity Layer

- E.g., ReLU (Rectified Linear Unit)
 - Pixel by pixel computation of $\max(0, x)$



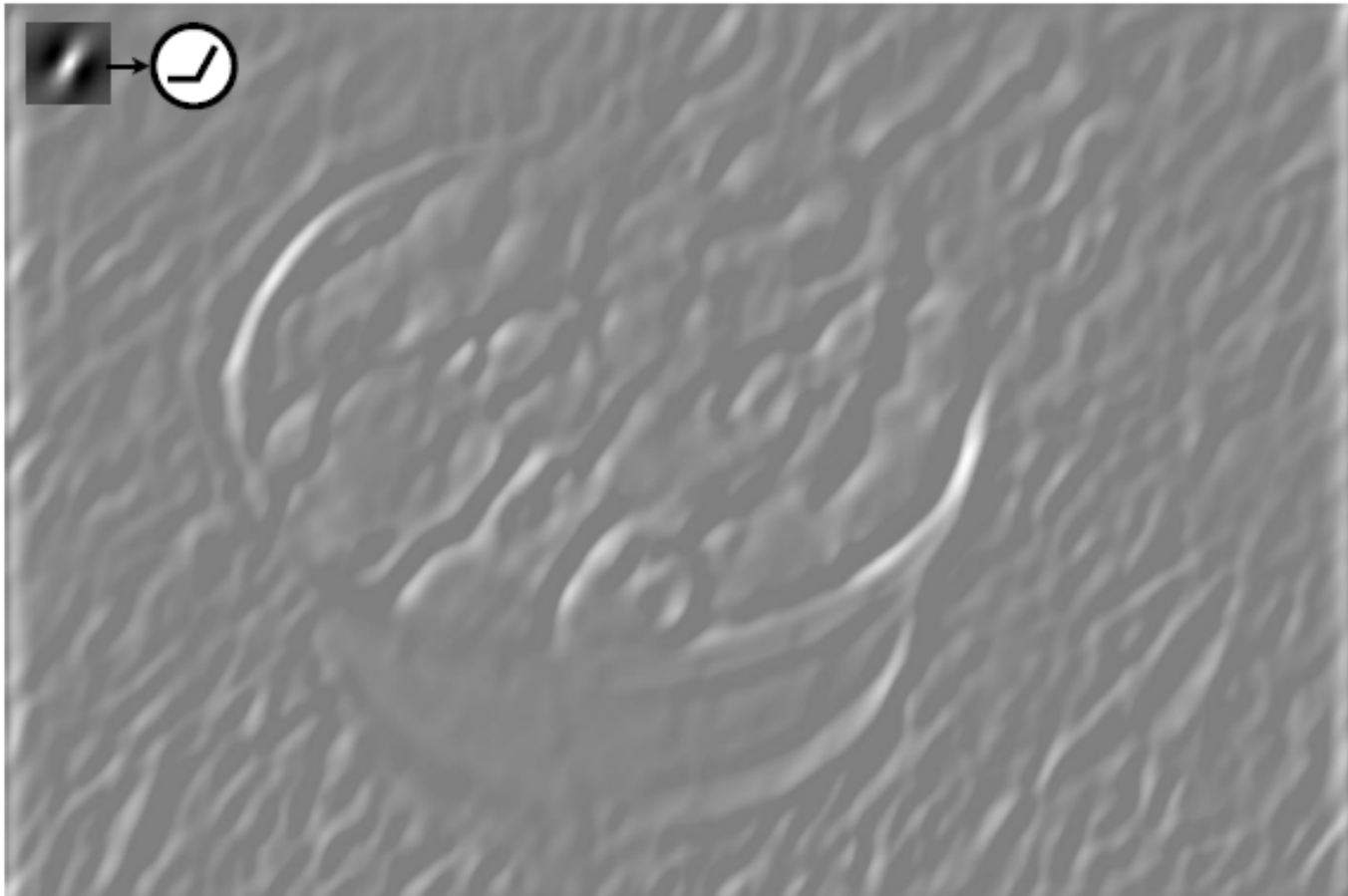
Nonlinearity Layer

- E.g., ReLU (Rectified Linear Unit)
 - Pixel by pixel computation of $\max(0, x)$

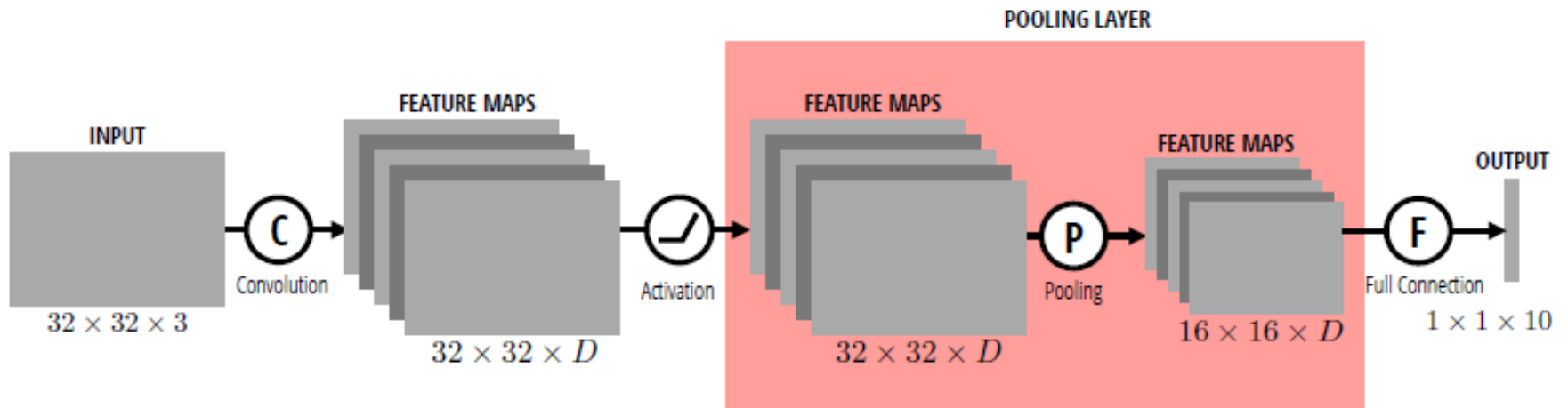


Nonlinearity Layer

- E.g., ReLU (Rectified Linear Unit)
 - Pixel by pixel computation of $\max(0, x)$

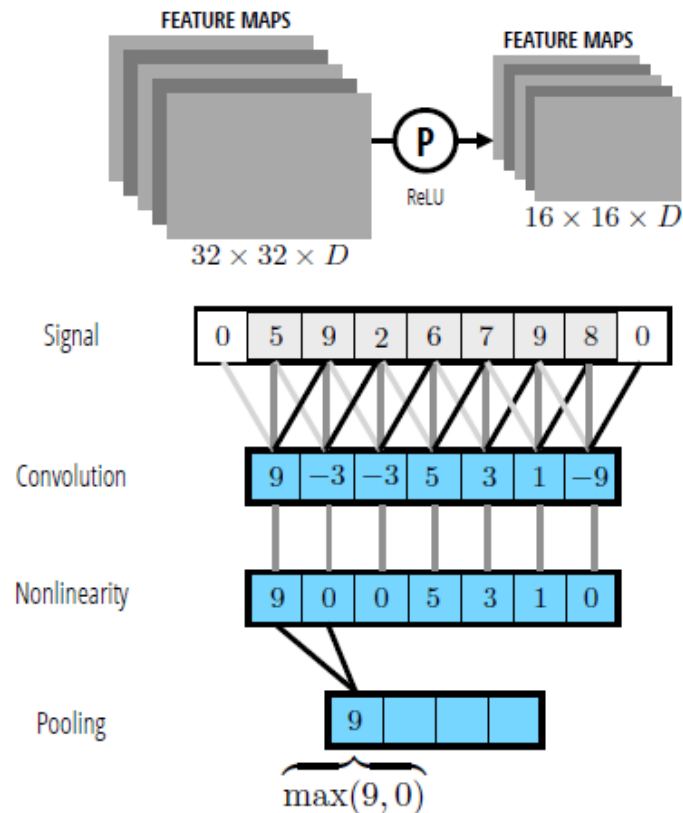


Pooling Layer in CNN



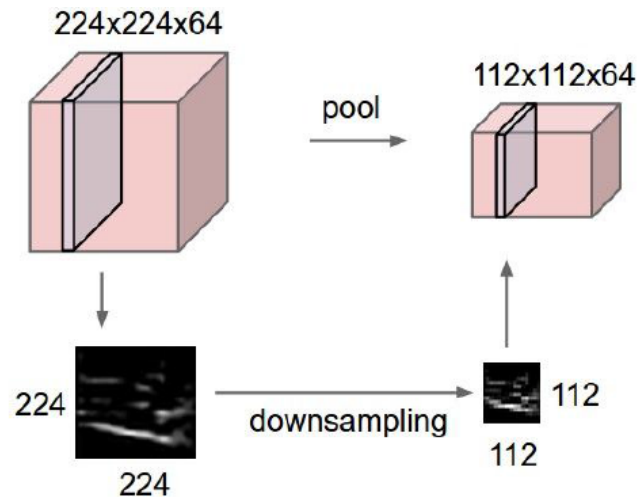
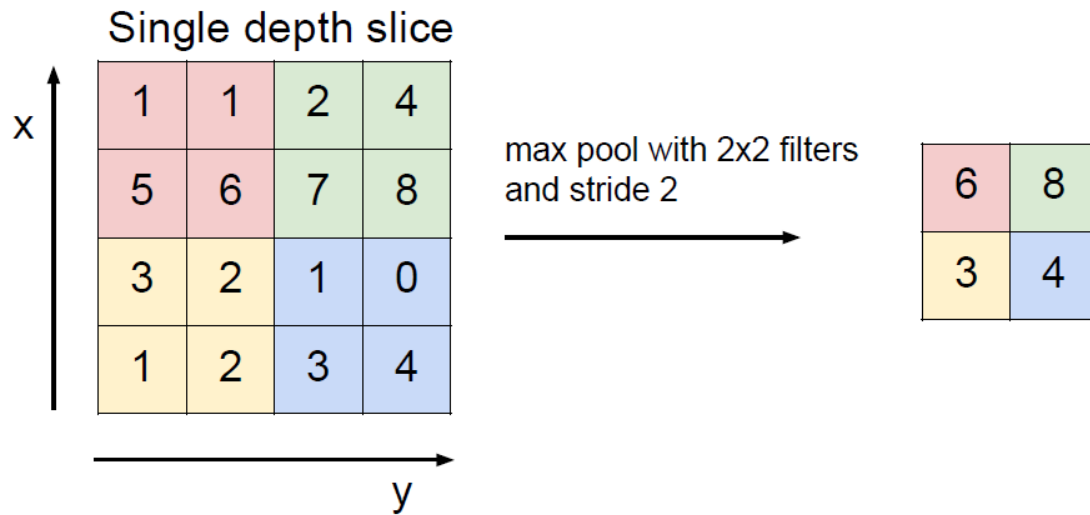
Pooling Layer

- Makes the representations smaller and more manageable
- Operates over each activation map independently
- E.g., Max Pooling

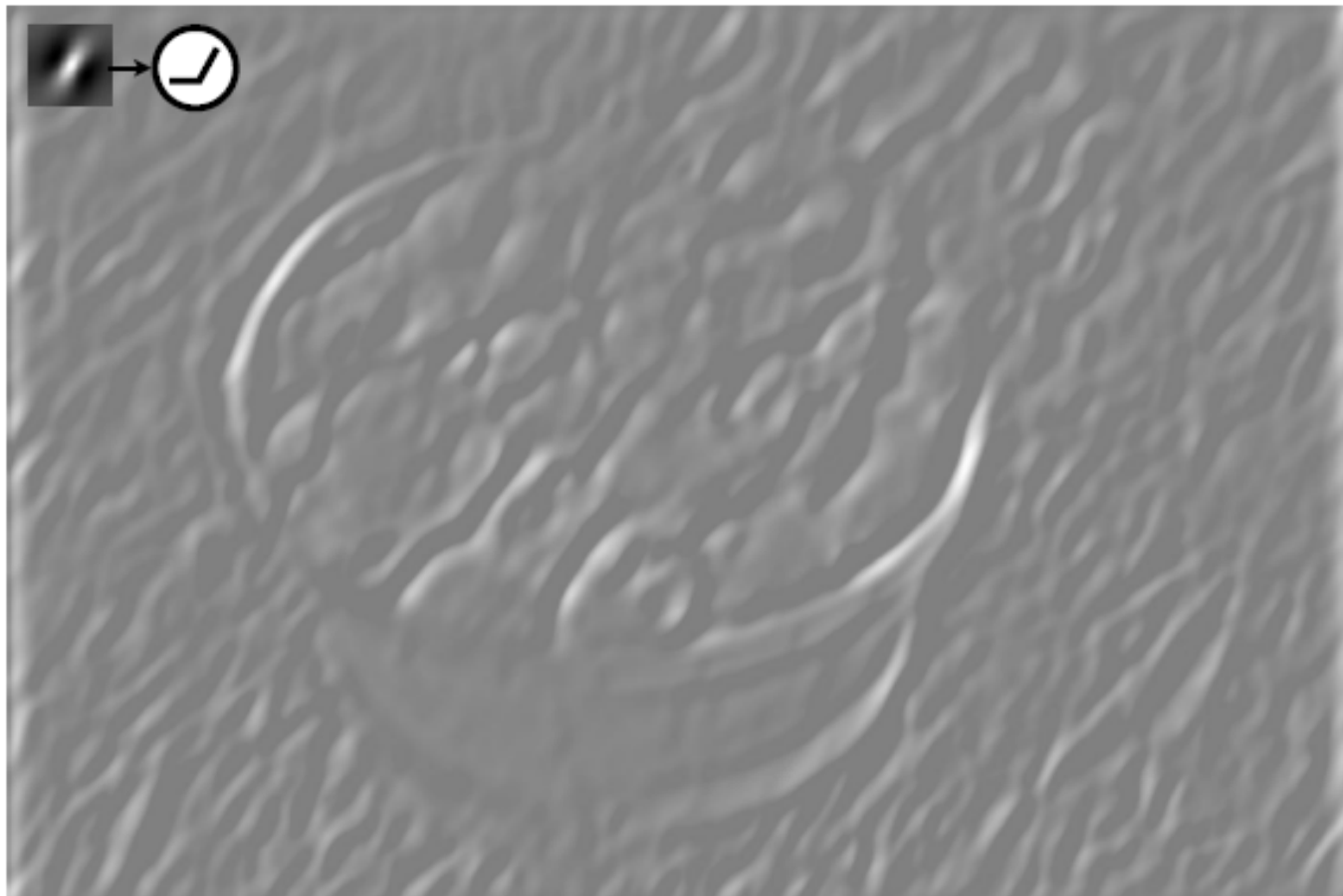


Pooling Layer

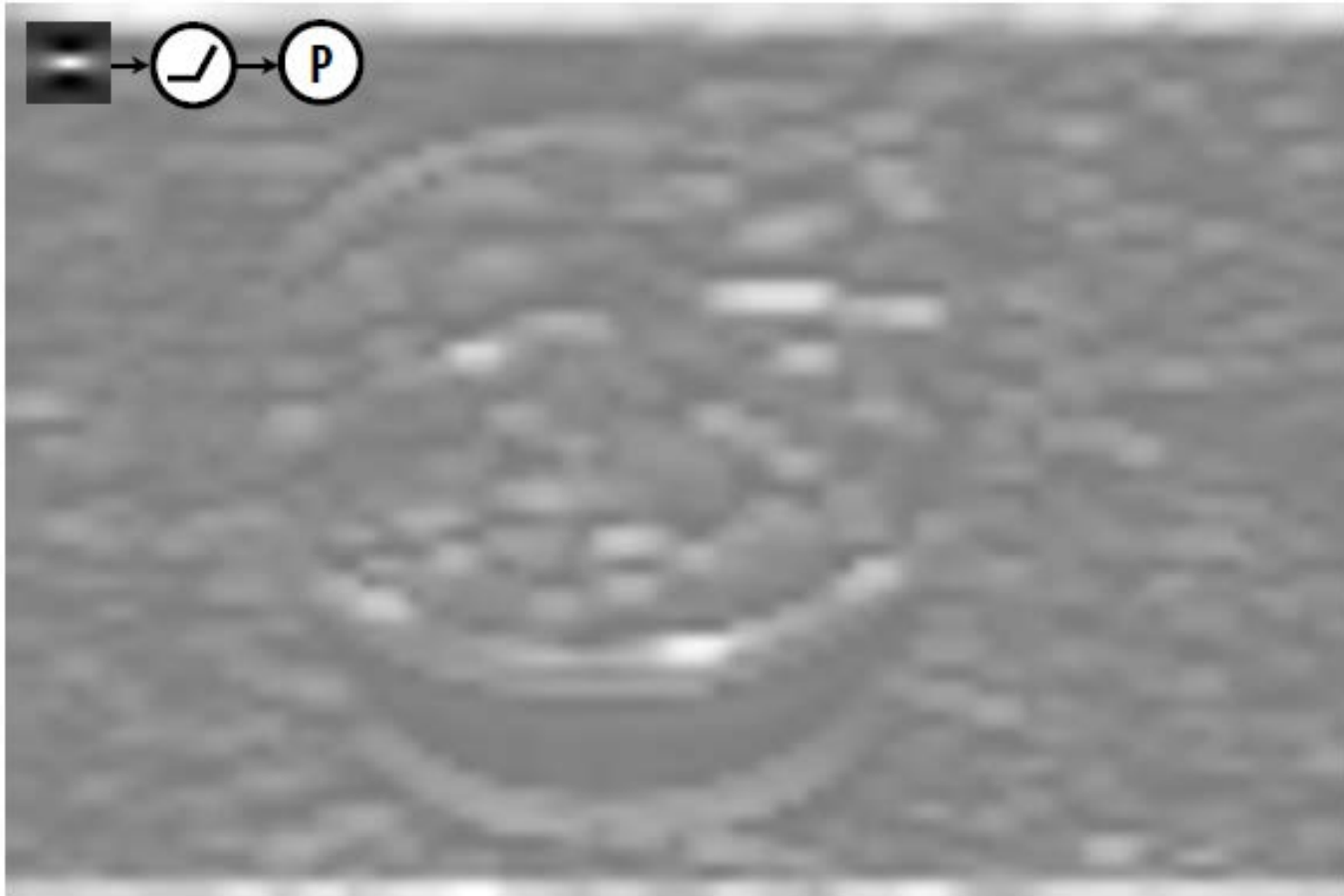
- Reduces the spatial size and provides spatial invariance



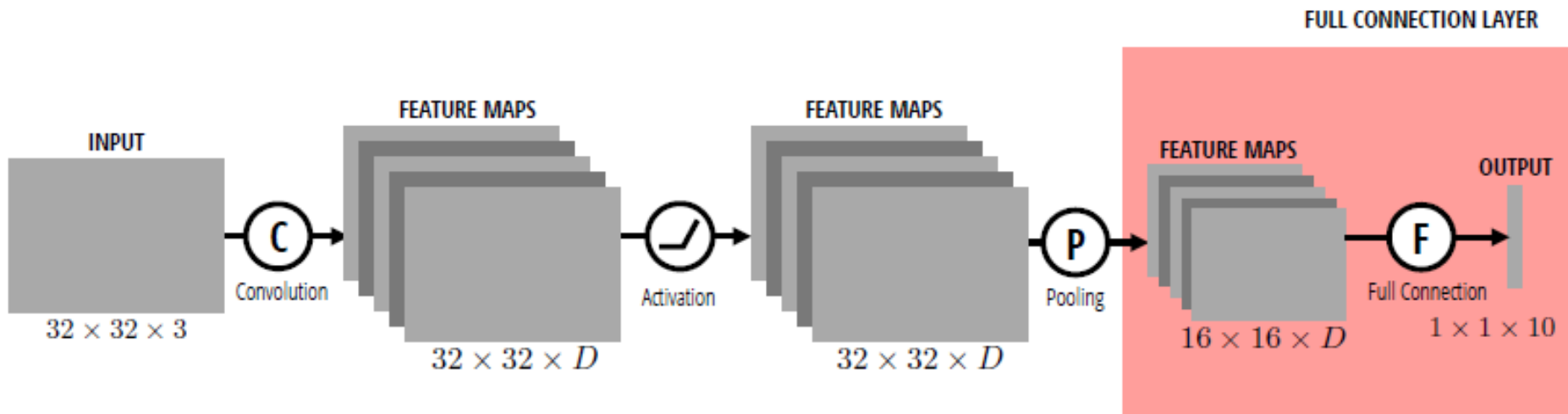
- Example
 - Nonlinearity by ReLU



- Example
 - Max pooling

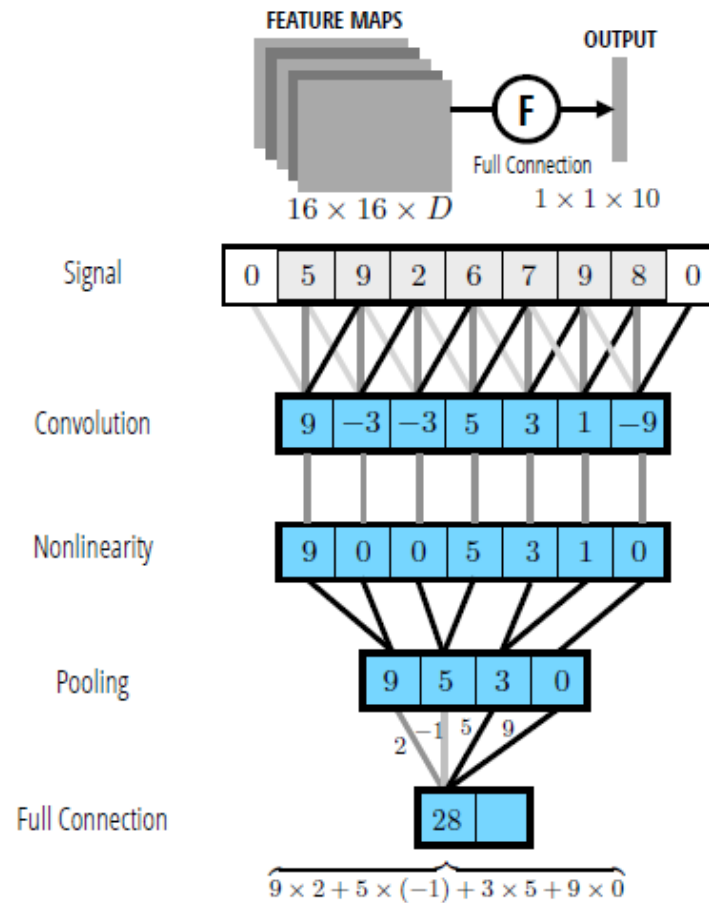


Fully Connected (FC) Layer in CNN



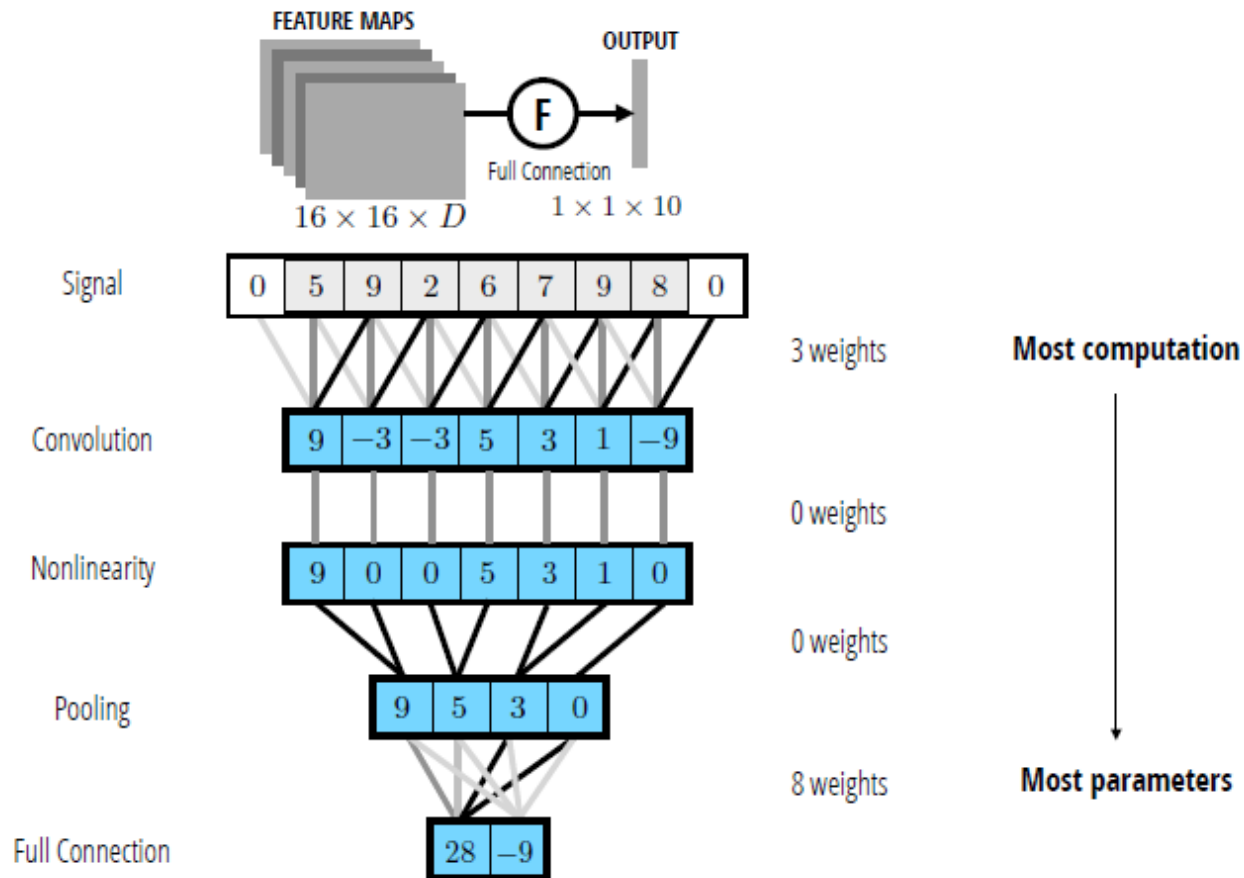
FC Layer

- Contains neurons that connect to the entire input volume, as in ordinary neural networks

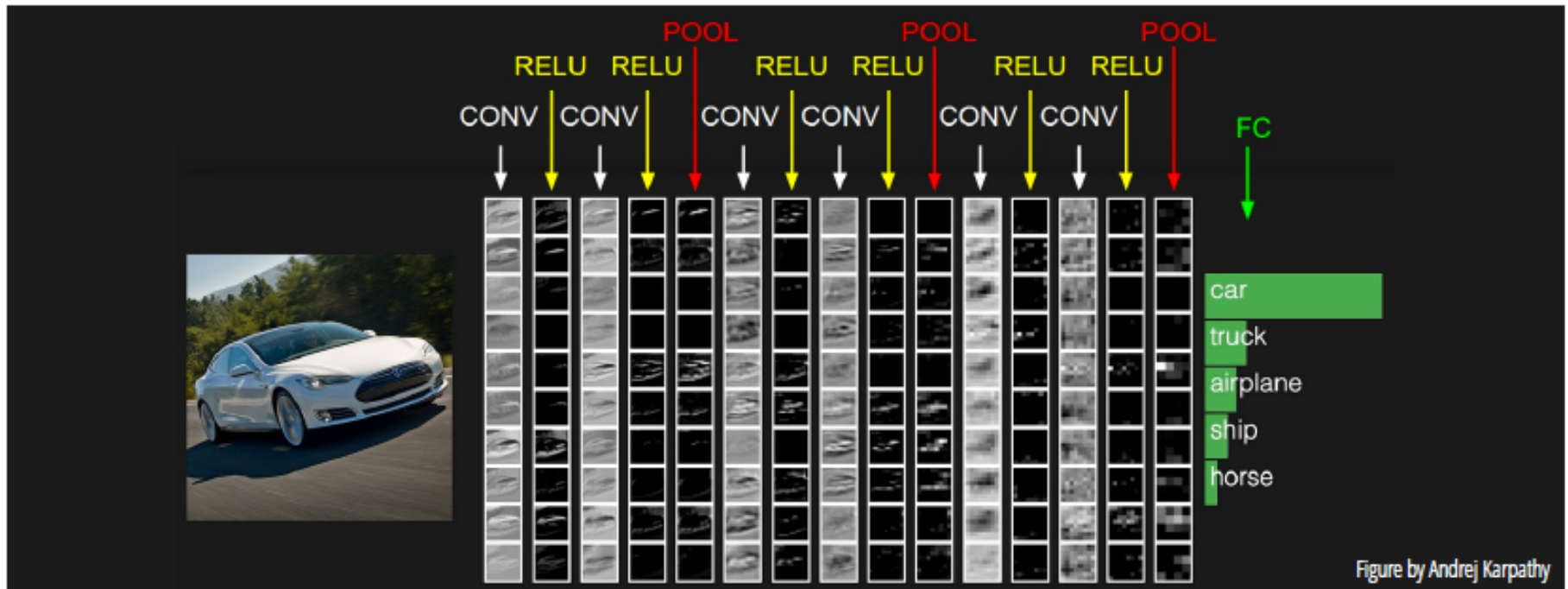
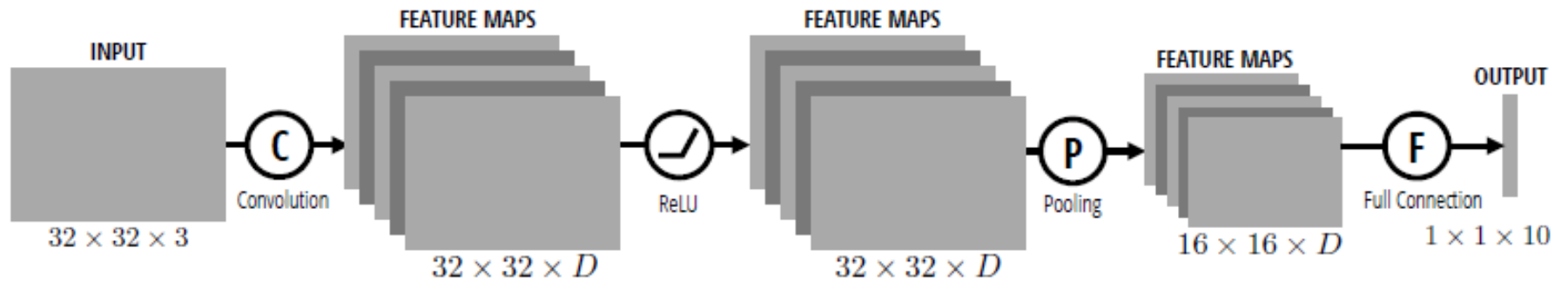


FC Layer

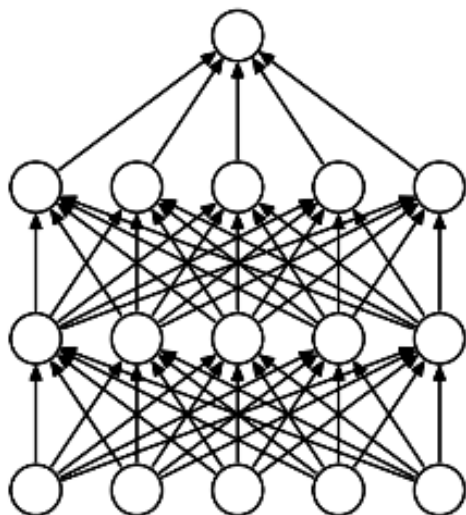
- Contains neurons that connect to the entire input volume, as in ordinary neural networks



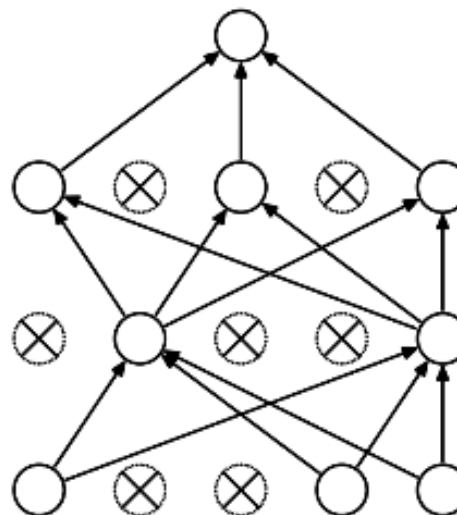
CNN



Training Technique #1: Dropout



(a) Standard Neural Net



(b) After applying dropout.

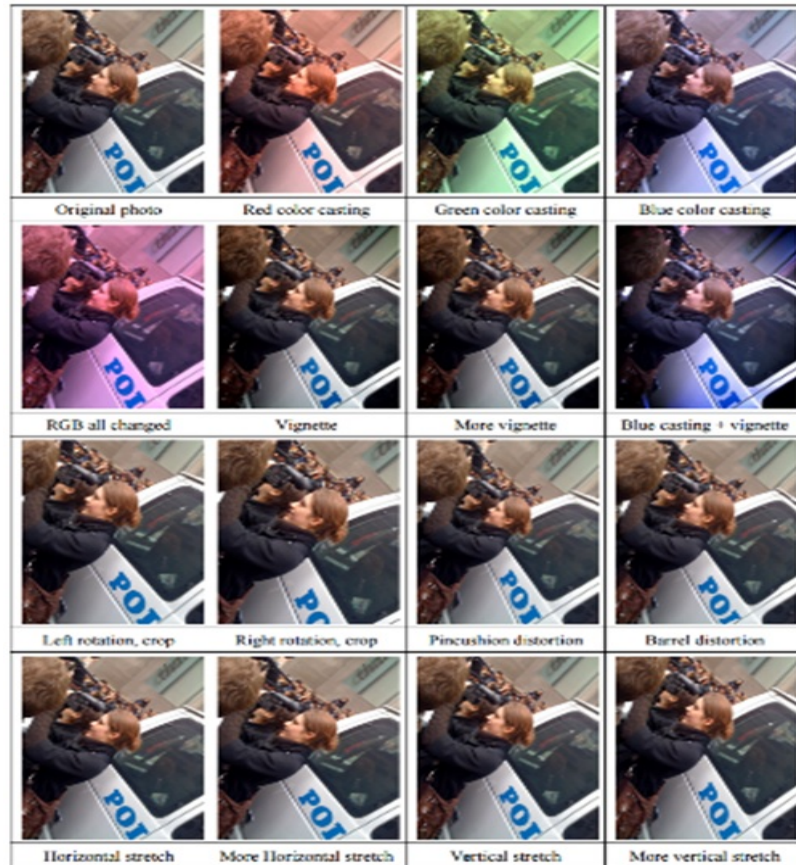
Intuition: successful **conspiracies**

Example: 50 people planning a conspiracy

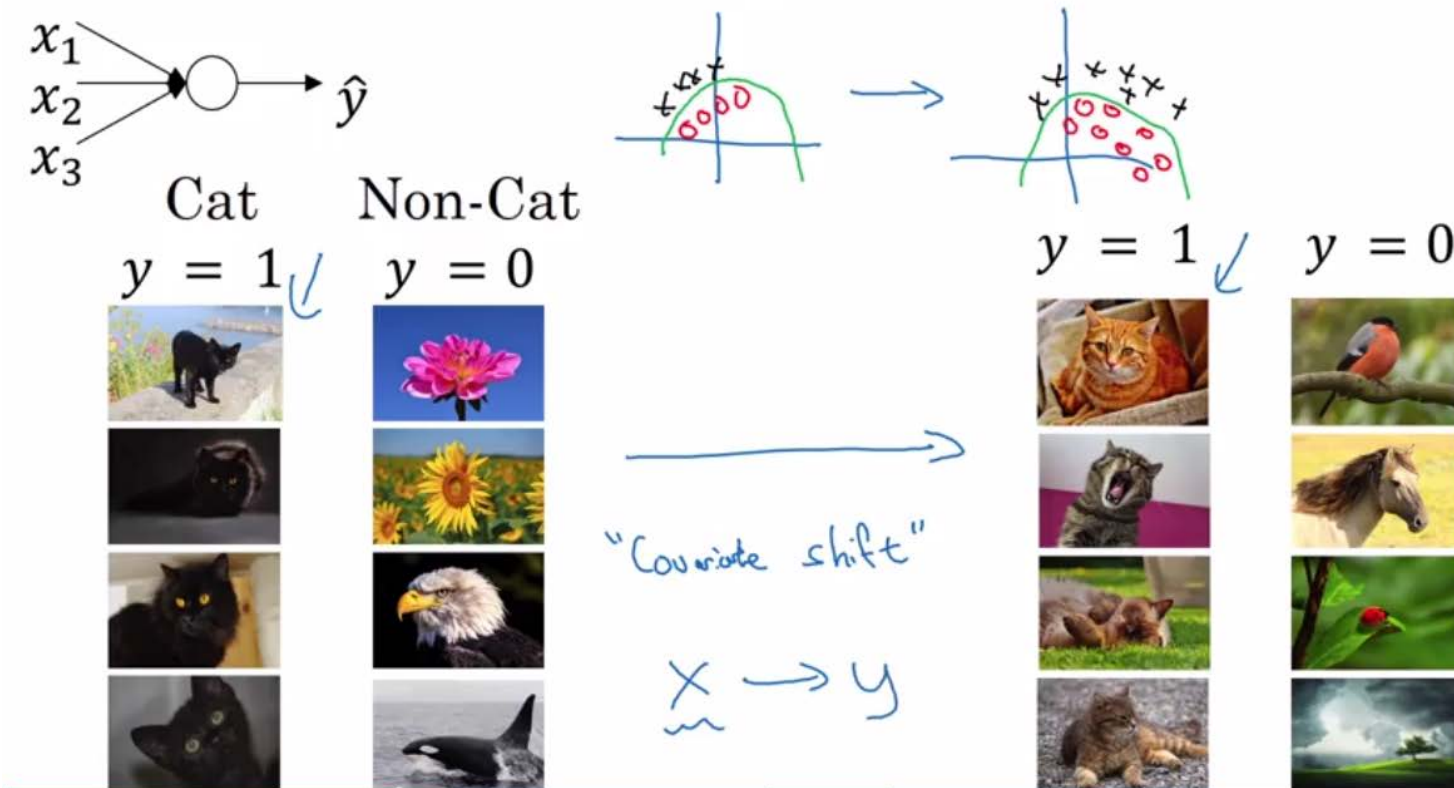
- Strategy A: plan a big conspiracy involving 50 people
 - Likely to fail. 50 people need to play their parts correctly.
- Strategy B: plan 10 conspiracies each involving 5 people
 - Likely to succeed!

Training Technique #2: Data Augmentation

- Create *virtual* training samples
 - Horizontal flip
 - Random crop
 - Color casting
 - Geometric distortion



Training Technique #3: Batch Normalization



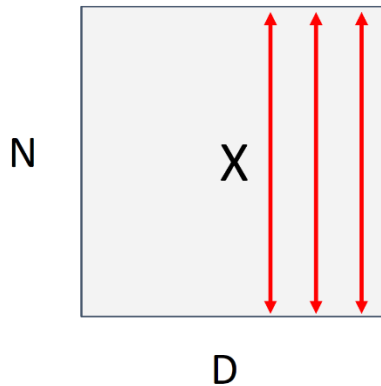
Batch Normalization (cont'd)

- Remarks
 - Differentiable function; back propagation OK

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

- Procedure

Input: $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean
across N samples

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Per-channel std
across N samples

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

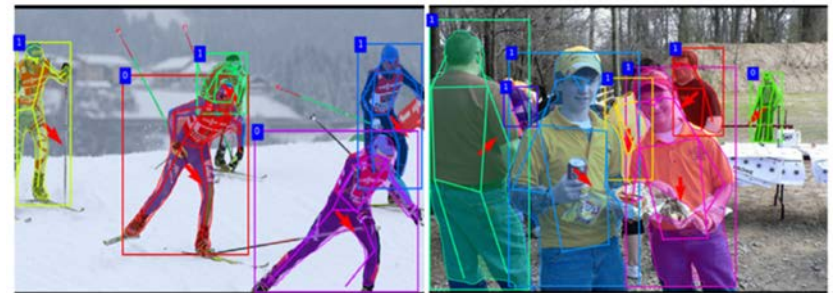
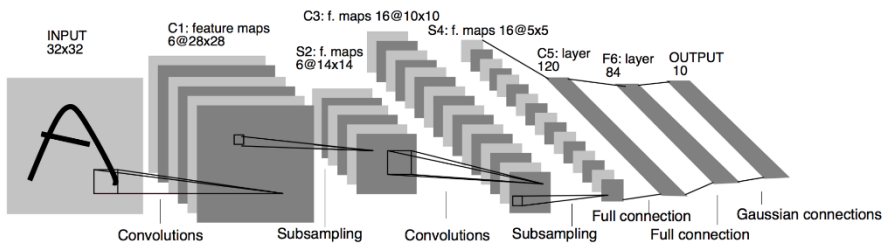
Normalized x,
Shape is N x D

What'd Be Covered in This Crash Course...

- **Learning-based Computer Vision**
 - From Linear to Non-Linear Classifiers
- **Start of Deep Learning for Computer Vision**
 - Convolutional Neural Networks
 - Self-Supervised Learning
 - Segmentation & Detection
- **Generative Models**
 - Autoencoder, Variational Autoencoder, Generative Adversarial Networks & Diffusion Models
- **Sequence-to-Sequence Learning**
 - Attention is All You Need: Transformer
- **Vision & Language Foundation Models**
 - Image-to-Text vs. Text-to-Image
 - Parameter-Efficient Fine-tuning

Supervised Learning

- Deep learning plus supervised learning are rocking the world ...



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.

- In real world scenarios, data-annotation is quite **time-consuming**
- Could one exploit supervised signals from **unlabeled** data?



Self-Supervised Learning (SSL)

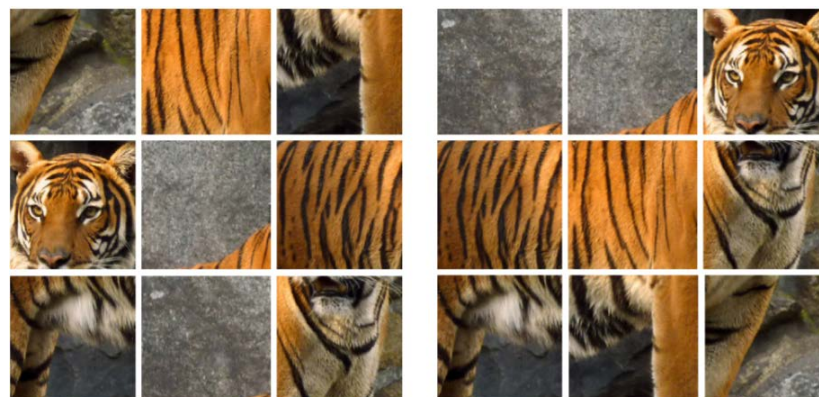
- Learning discriminative representations from **unlabeled** data
- Create self-supervised tasks via **data augmentation**



Colorization



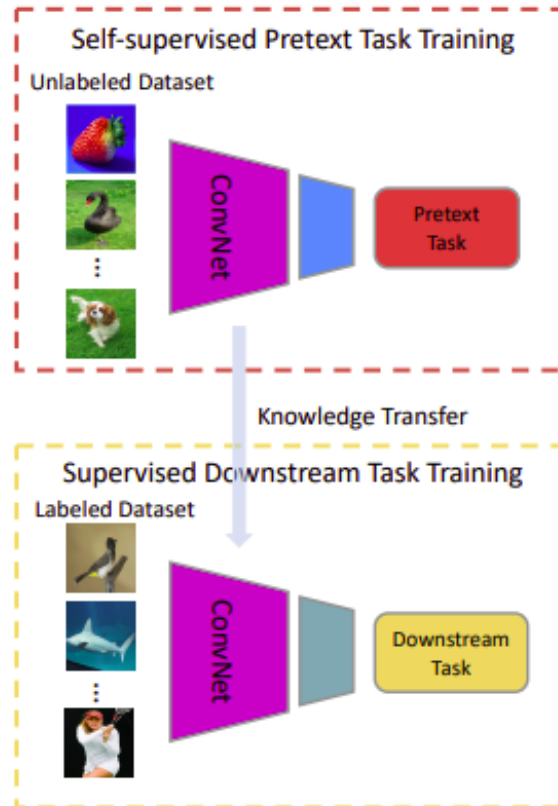
Rotation



Jigsaw Puzzle

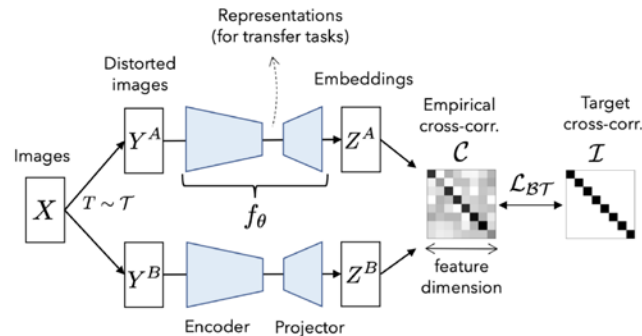
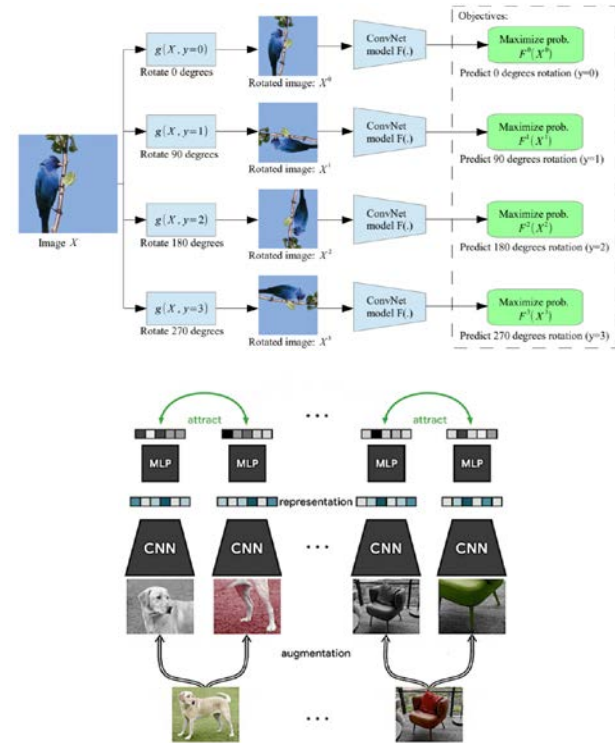
Self-Supervised Learning (SSL)

- Self-Supervised Pretraining
- Supervised Fine-tuning



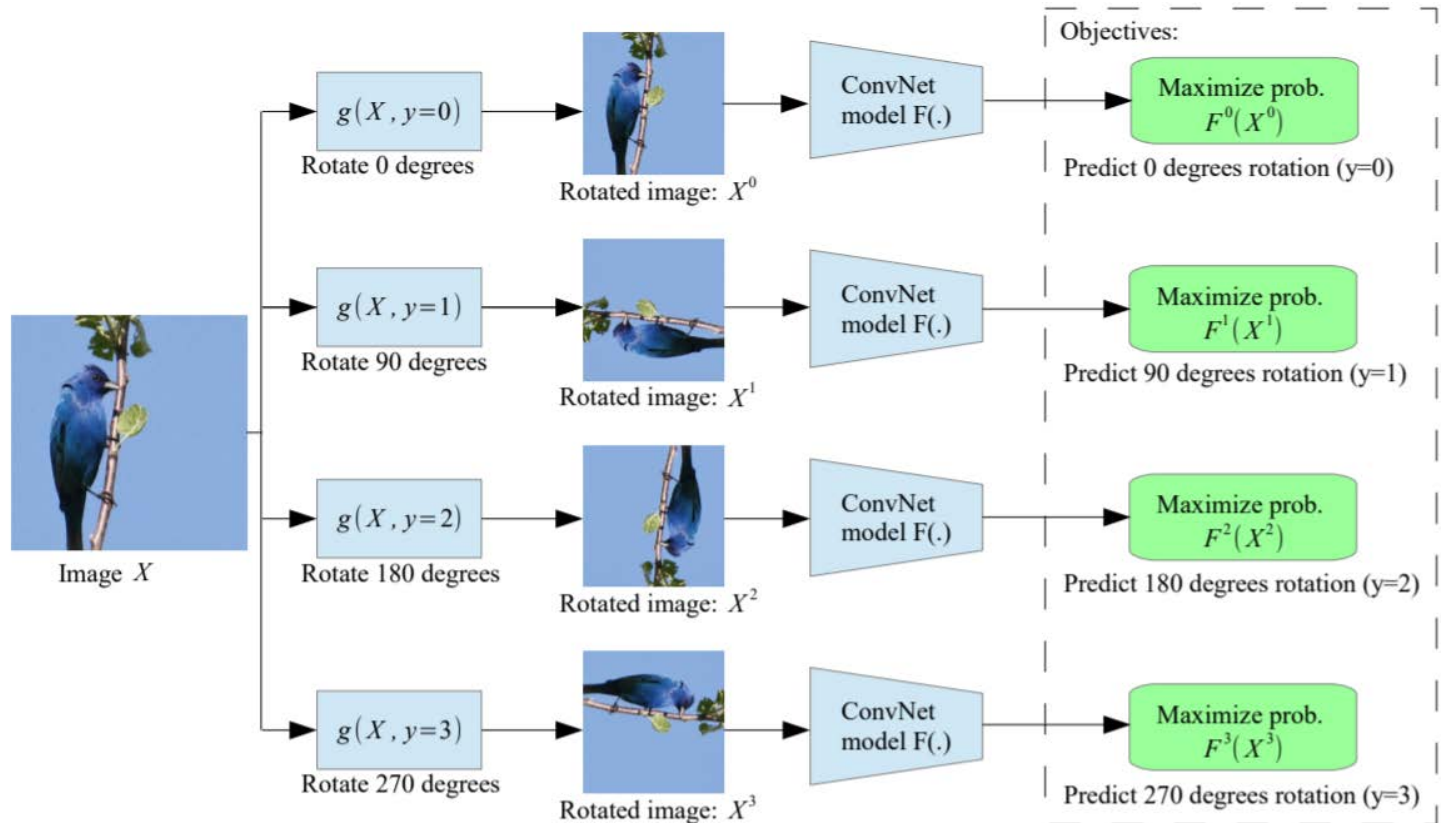
Self-Supervised Learning (SSL)

- Pretext Tasks
 - Jigsaw (ECCV'16)
 - RotNet (ICLR'18)
- Contrastive Learning
 - CPC (ICML'20)
 - SimCLR (ICML'20)
- Learning w/o negative samples
 - BYOL (NeurIPS'20)
 - Barlow Twins (ICML'21)



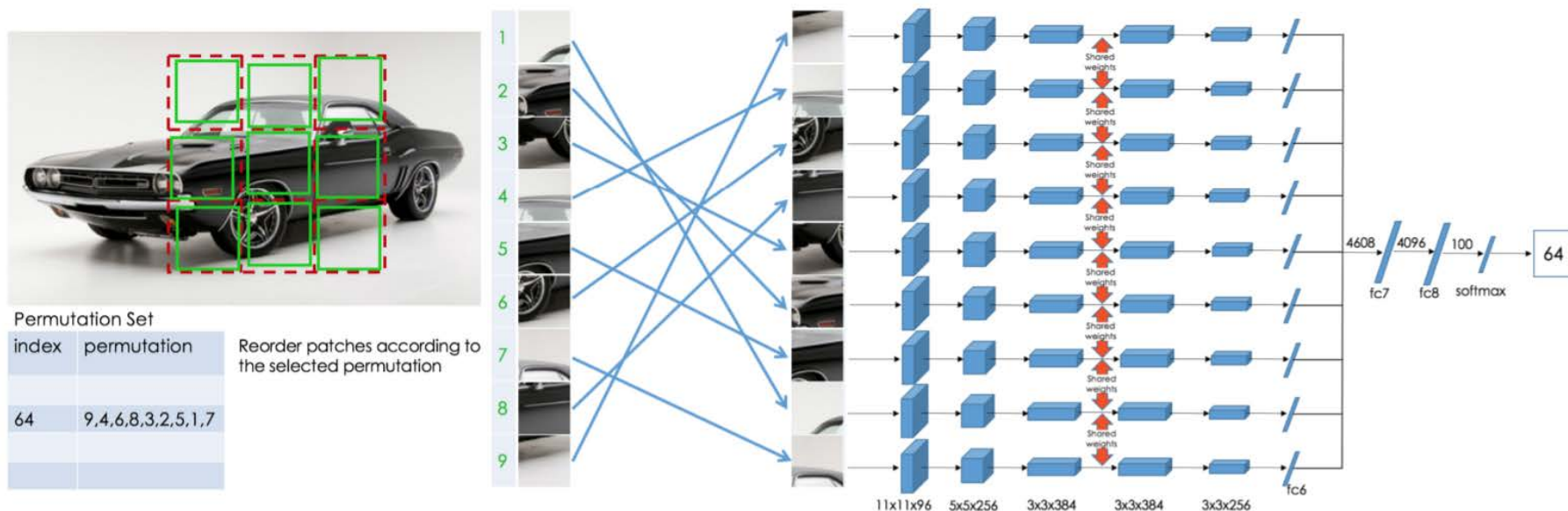
RotNet

- Learning to predict the **rotation** angle



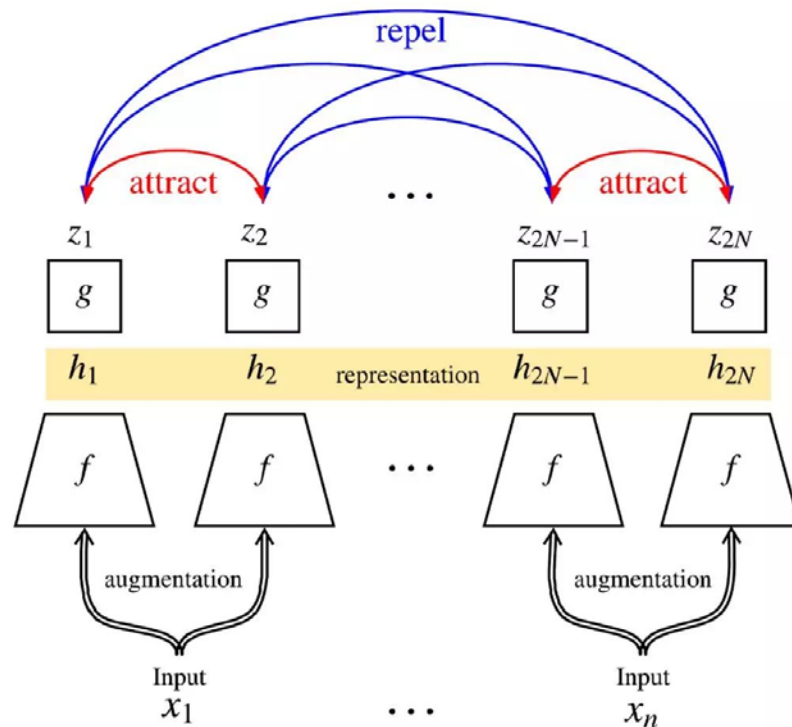
Jigsaw Puzzle

- Assign the **permutation index** and perform augmentation
- Solve jigsaw puzzle by predicting the permutation index



SimCLR

- **Attract** augmented images and **repel** negative samples
- Improve the representation quality with **projection heads** (g)...why?



What'd Be Covered in This Crash Course...

- **Learning-based Computer Vision**
 - From Linear to Non-Linear Classifiers
- **Start of Deep Learning for Computer Vision**
 - Convolutional Neural Networks
 - Self-Supervised Learning
 - Segmentation & Detection
- **Generative Models**
 - Autoencoder, Variational Autoencoder, Generative Adversarial Networks & Diffusion Models
- **Sequence-to-Sequence Learning**
 - Attention is All You Need: Transformer
- **Vision & Language Foundation Models**
 - Image-to-Text vs. Text-to-Image
 - Parameter-Efficient Fine-tuning

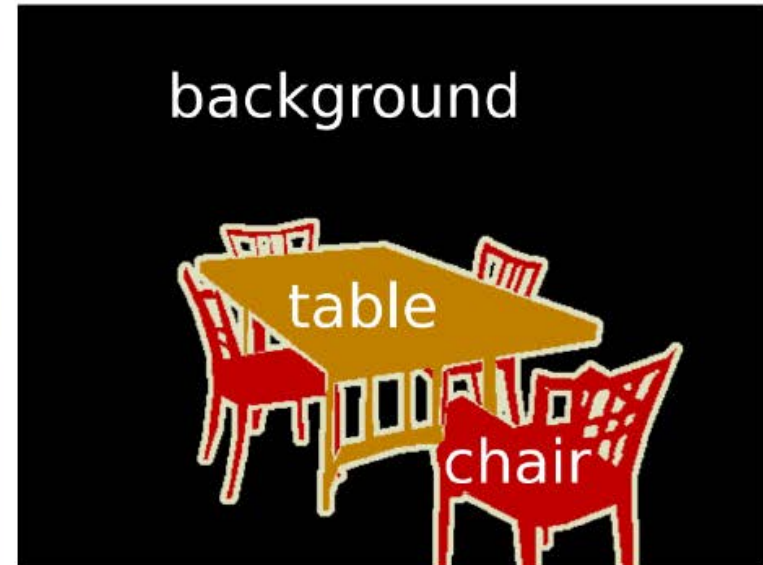
Image Segmentation

- Goal:
Group pixels into meaningful or perceptually similar regions



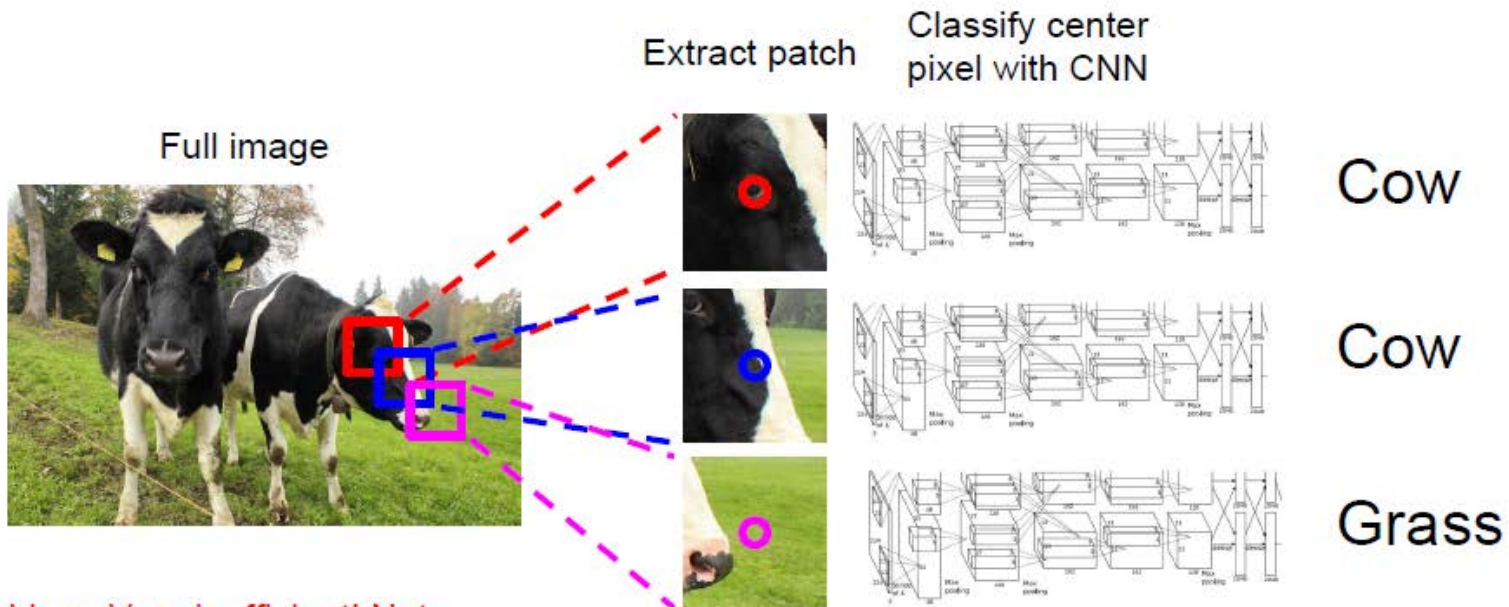
A Practical Segmentation Task

- Semantic Segmentation
 - Supervised learning
 - Assign a class label to each pixel in the input image (i.e., [pixel-level classification](#))
 - Not like instance segmentation, do not differentiate instances; only care about pixel labels



Semantic Segmentation

- Sliding Window



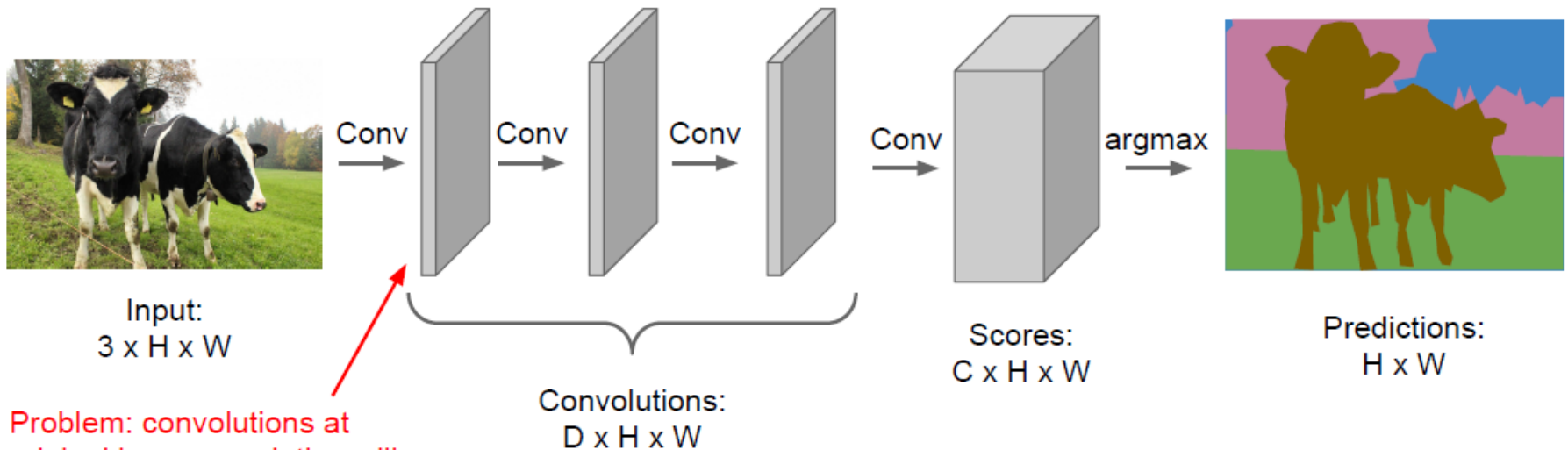
Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation

- Fully Convolutional Nets

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Problem: convolutions at original image resolution will be very expensive ...

Fully Convolutional Networks (FCN)

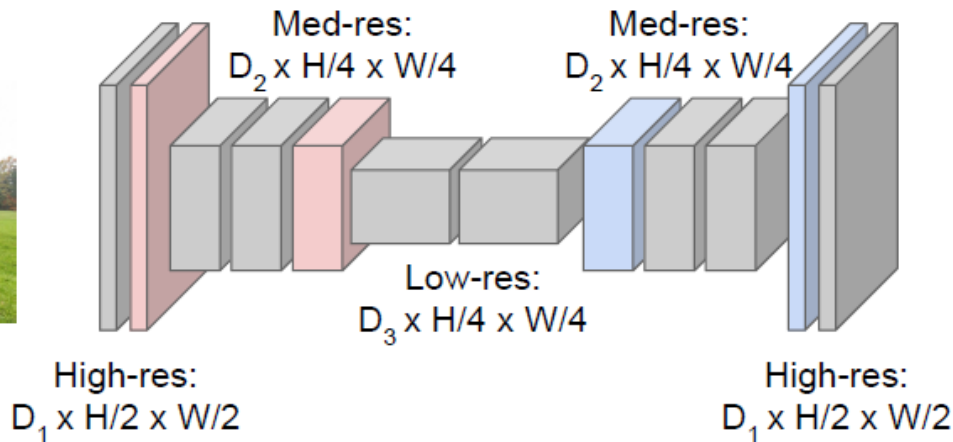
- Remarks
 - All layers are convolutional
 - End-to-end training

Downsampling:
Pooling, strided convolution



Input:
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



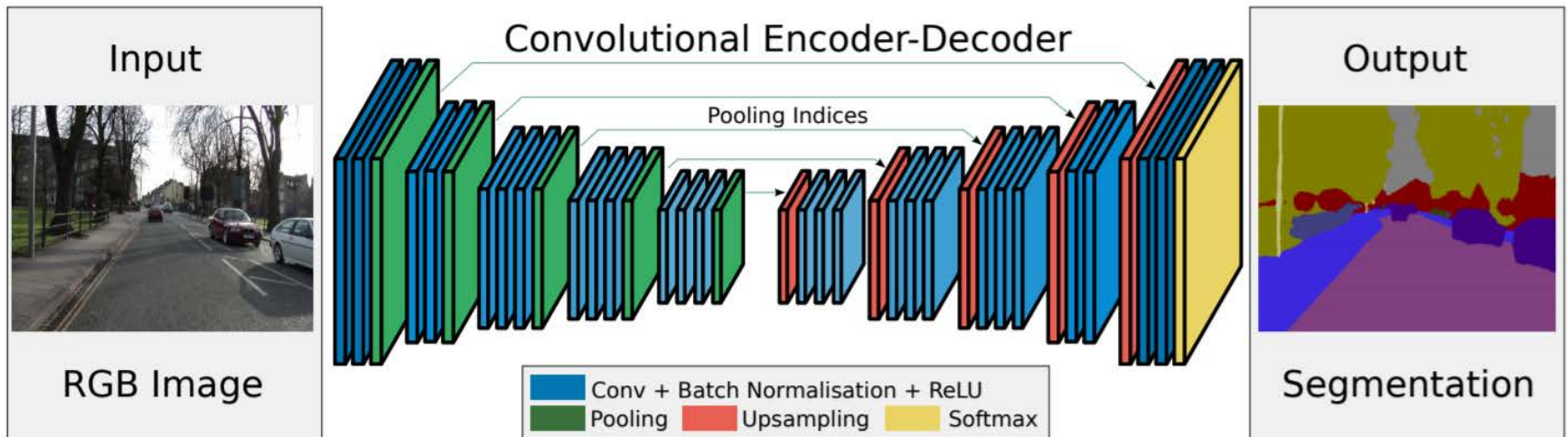
Upsampling:
Unpooling or strided transpose convolution



Predictions:
 $H \times W$

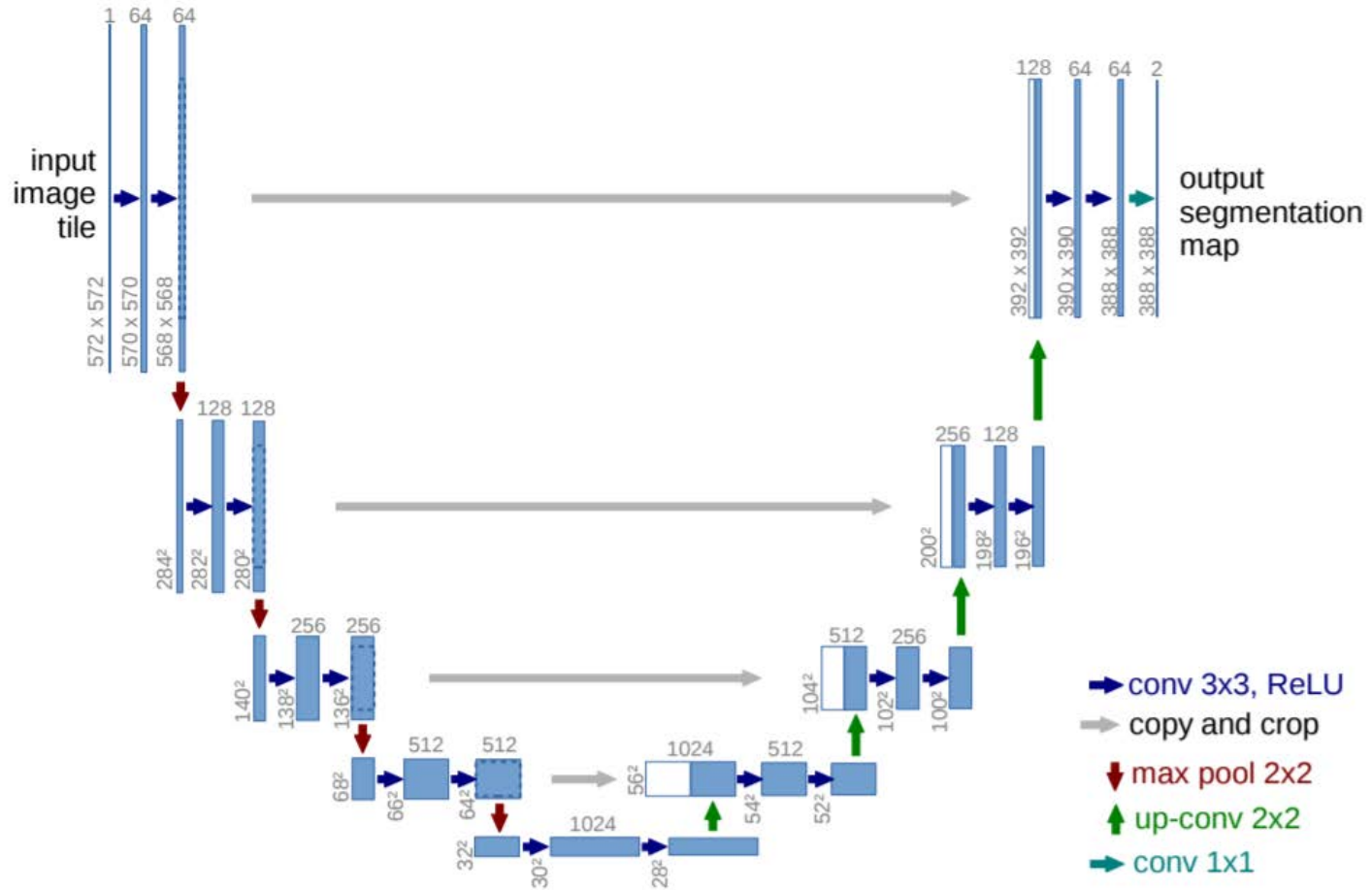
SegNet

- Efficient architecture (memory + computation time)
- Upsampling reusing max-unpooling indices
- Reasonable results without performance boosting addition
- Comparable to FCN



“SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation” [\[link\]](#)

U-Net



U-Net: Convolutional Networks for Biomedical Image Segmentation [\[link\]](#)

Roadmap

Classification



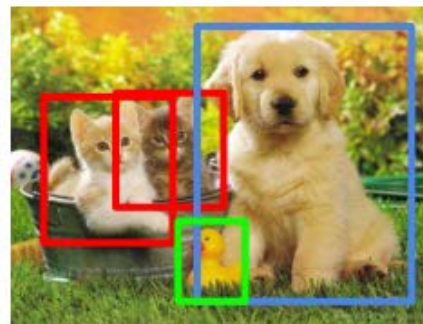
CAT

**Classification
+ Localization**



CAT

Object Detection



CAT, DOG, DUCK

**Instance
Segmentation**



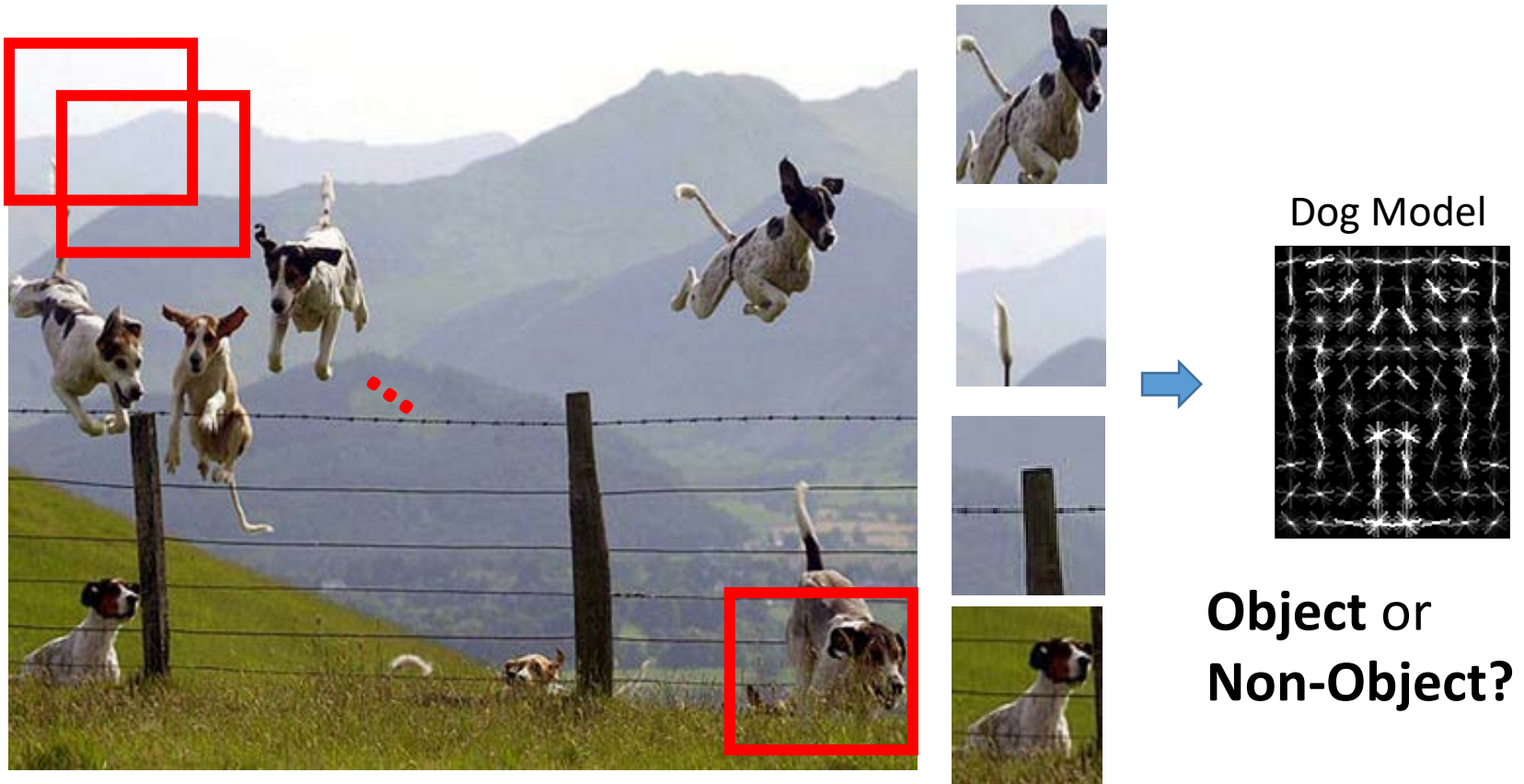
CAT, DOG, DUCK

Single object

Multiple objects

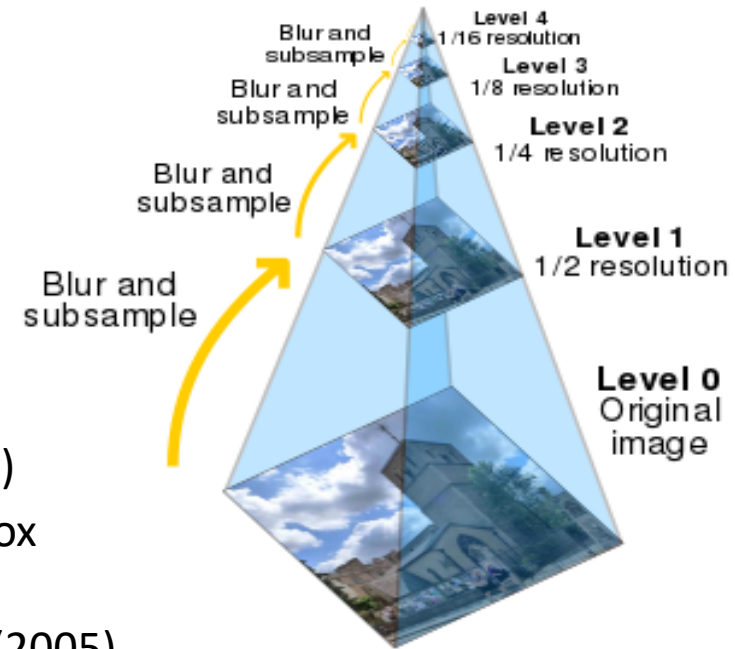
Object Detection

- Focus on object search: “Where is it?”
- Build templates that quickly differentiate object patch from background patch



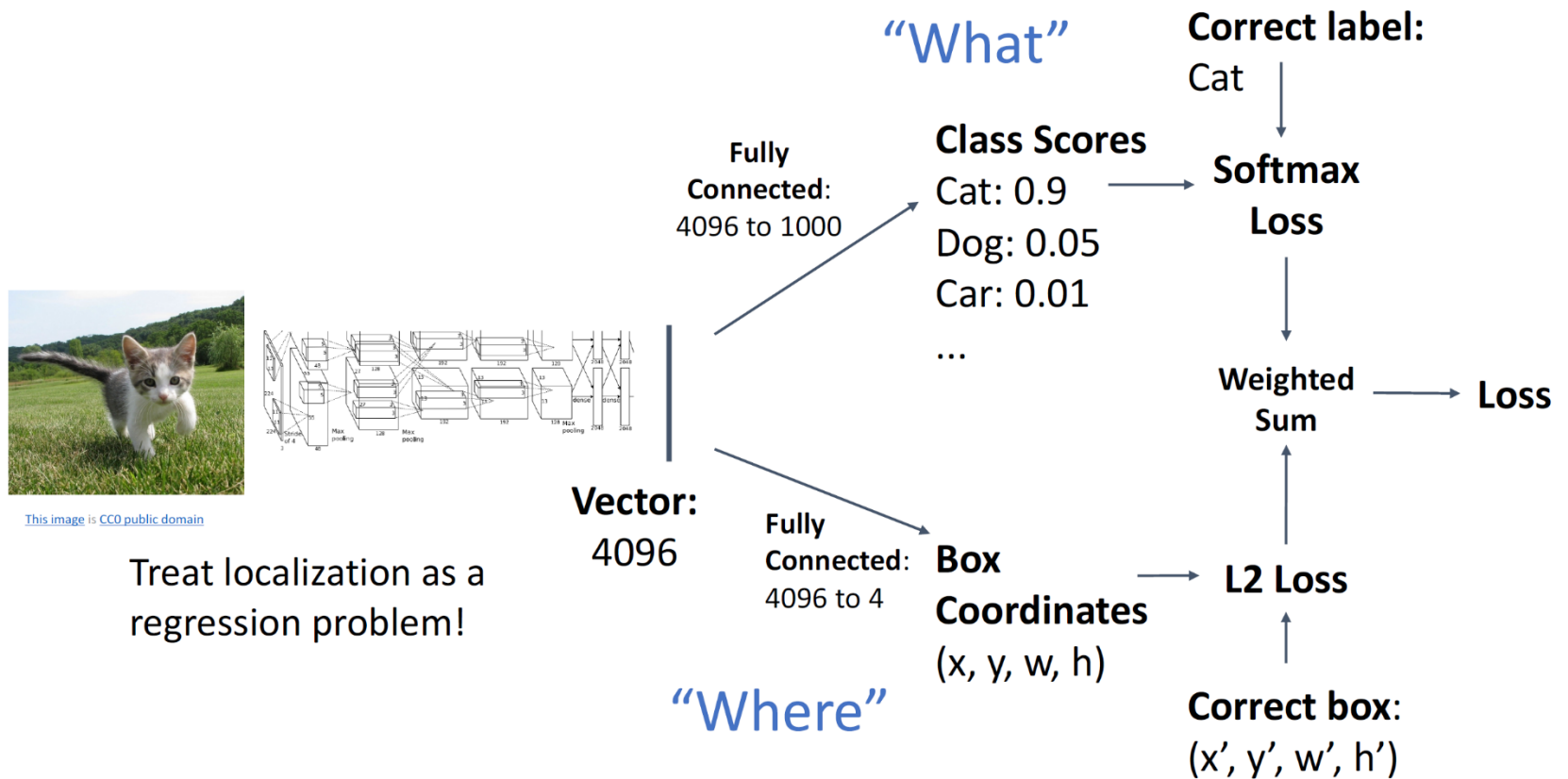
Type of Approaches

- Sliding Windows
 - “Slide” a box around the input image or even across image scales (i.e., image pyramid)
 - Classify each cropped image region inside the box and determine if it’s an object of interest or not
 - E.g., HOG (person) detector by Dalal and Triggs (2005)
Deformable part-based model by Felzenswalb et al. (2010)
Real-time (face) detector by Viola and Jones (2001)
- Region (Object) Proposals
 - Generate region (object) proposals
 - Classify each image region and determine it’s an object or not



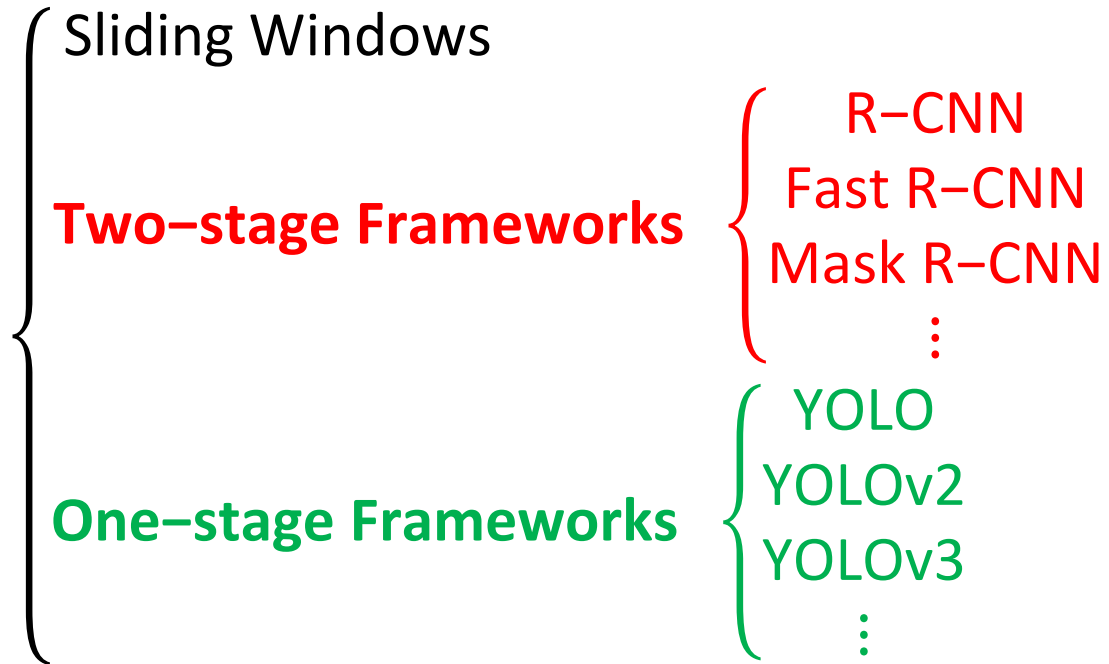
Type of Approaches (cont'd)

- CNN-based Methods



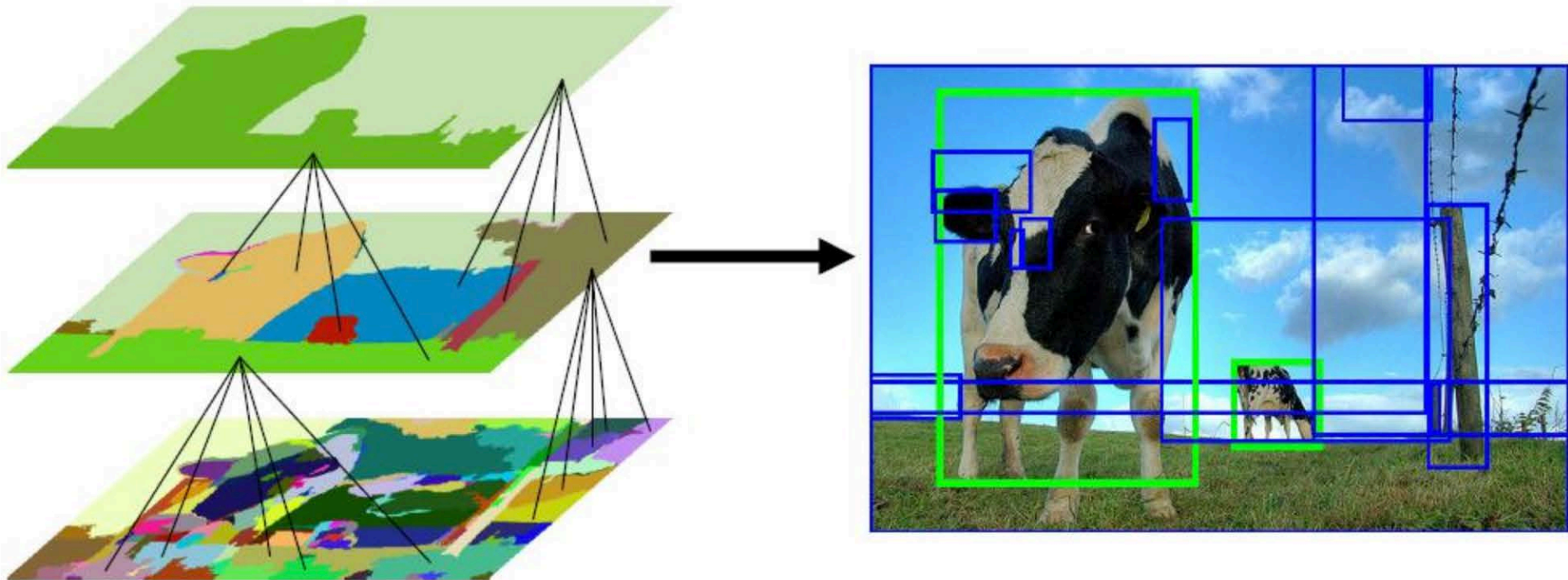
Two-Stage vs. One-Stage Object Detection

Methods

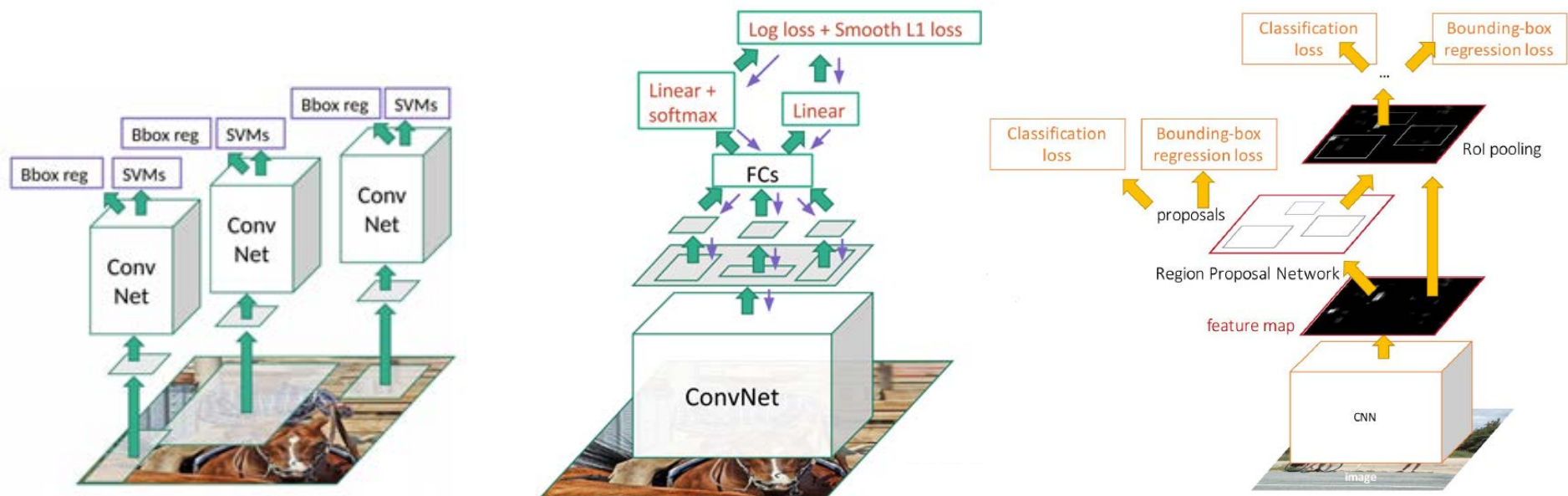


Region Proposal

- Solution
 - Use pre-processing algorithms to filter out some regions first, and feed the regions of interest (i.e., region proposals) into CNN
 - E.g., selective search



R-CNN, Fast R-CNN, & Faster R-CNN

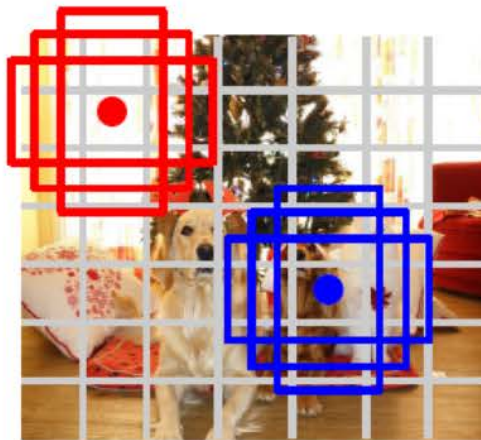


One-Stage Object Detection: Detection without Proposals

Go from input image to tensor of scores with one big convolutional network! →



Input image
 $3 \times H \times W$



Divide image into grid
 7×7

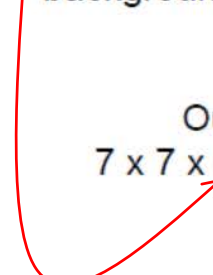
Image a set of **base boxes**
centered at each grid cell
Here $B = 3$



Within each grid cell:

- Regress from each of the B base boxes to a final box with 5 numbers:
($dx, dy, dh, dw, confidence$)
- Predict scores for each of C classes (including background as a class)

Output:
 $7 \times 7 \times (5 * B + C)$

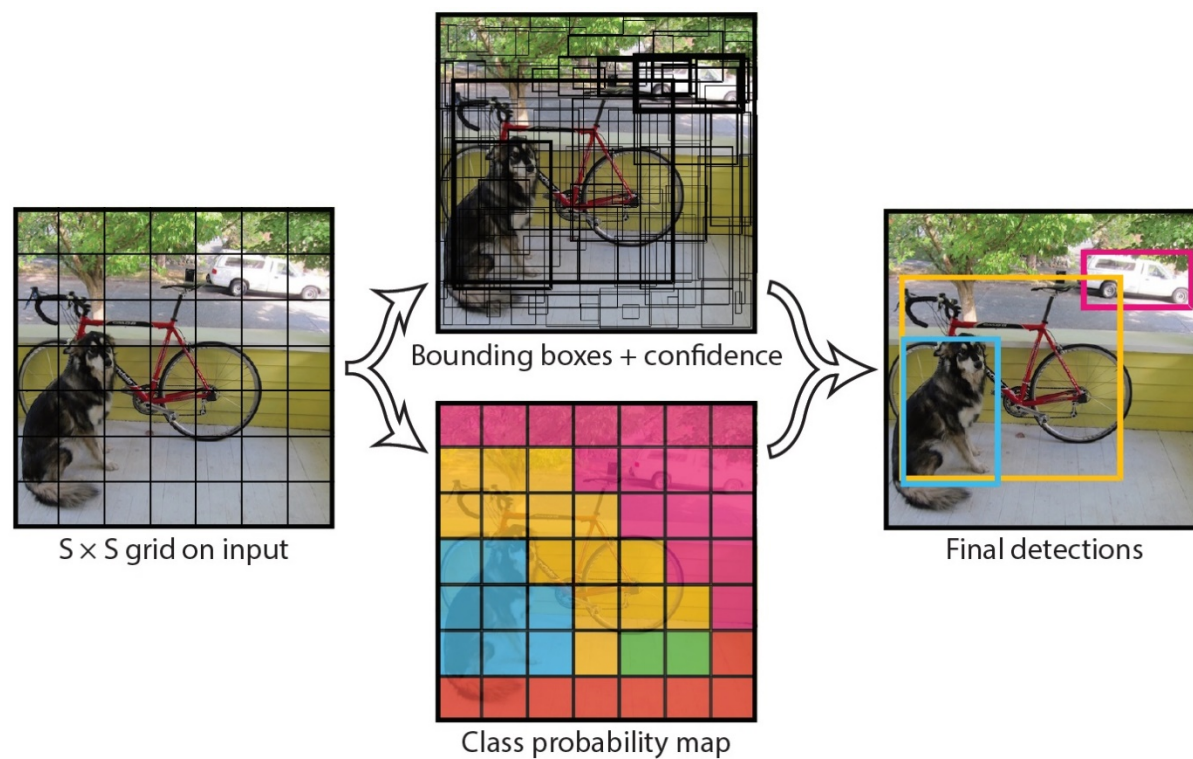


Redmon et al, "You Only Look Once:
Unified, Real-Time Object Detection", CVPR 2016
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

You Only Look Once (YOLO)

Divide the image into an $S \times S$ grid and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities.

These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.

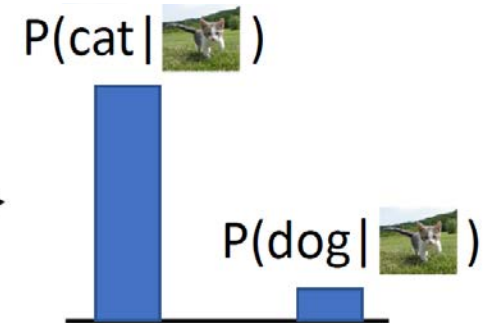


What'd Be Covered in This Crash Course...

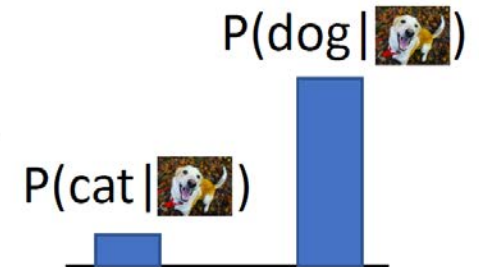
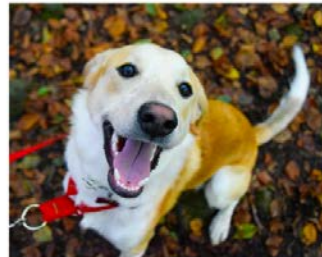
- **Learning-based Computer Vision**
 - From Linear to Non-Linear Classifiers
- **Start of Deep Learning for Computer Vision**
 - Convolutional Neural Networks
 - Self-Supervised Learning
 - Segmentation & Detection
- **Generative Models**
 - Autoencoder, Variational Autoencoder, Generative Adversarial Networks & Diffusion Models
- **Sequence-to-Sequence Learning**
 - Attention is All You Need: Transformer
- **Vision & Language Foundation Models**
 - Image-to-Text vs. Text-to-Image
 - Parameter-Efficient Fine-tuning

Discriminative vs. Generative Models (cont'd)

Discriminative Model:
Learn a probability distribution $p(y|x)$



Generative Model:
Learn a probability distribution $p(x)$

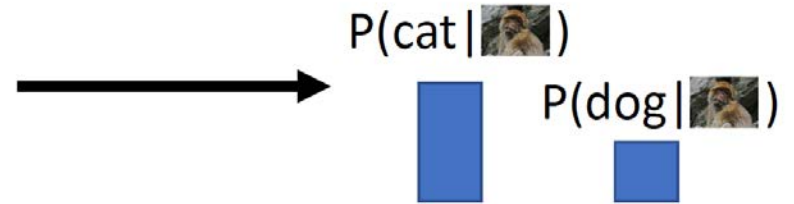


Conditional Generative Model: Learn $p(x|y)$

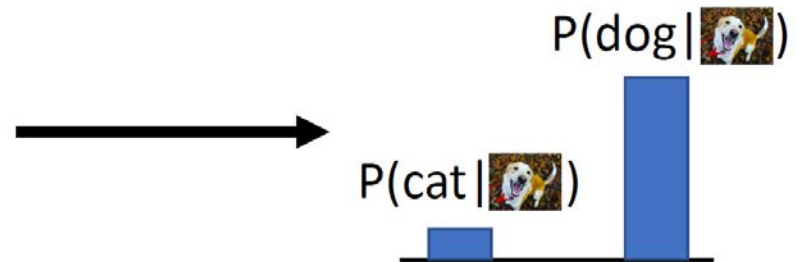
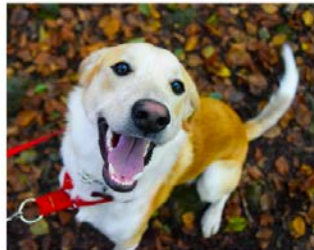
Discriminative model: the possible labels for each input "compete" for probability mass. But no competition between **images**

Discriminative vs. Generative Models (cont'd)

Discriminative Model:
Learn a probability distribution $p(y|x)$



Generative Model:
Learn a probability distribution $p(x)$



Conditional Generative Model: Learn $p(x|y)$

Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

Discriminative vs. Generative Models (cont'd)

Discriminative Model:

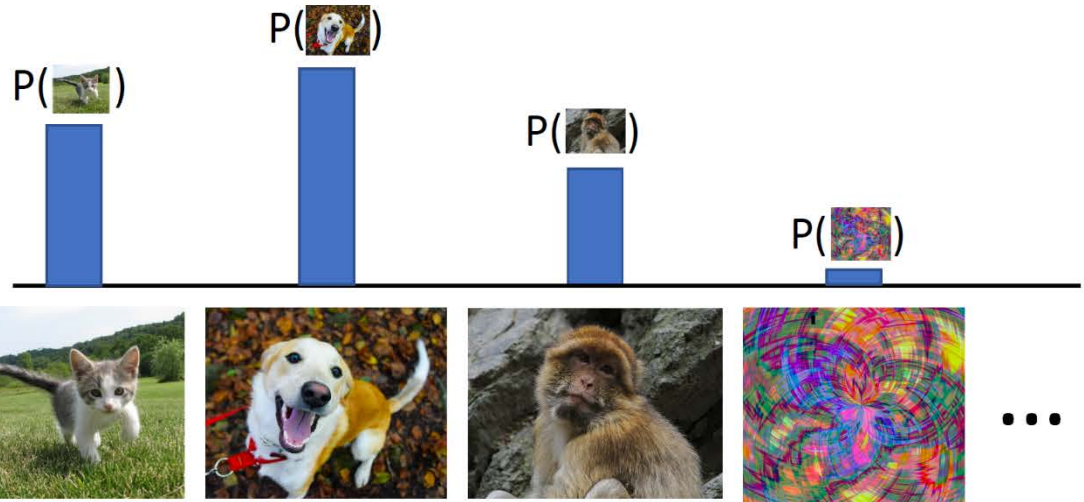
Learn a probability distribution $p(y|x)$

Generative Model:

Learn a probability distribution $p(x)$

Conditional Generative Model:

Learn $p(x|y)$



Generative model: All possible images compete with each other for probability mass

Model can “reject” unreasonable inputs by assigning them small values

Discriminative vs. Generative Models (cont'd)

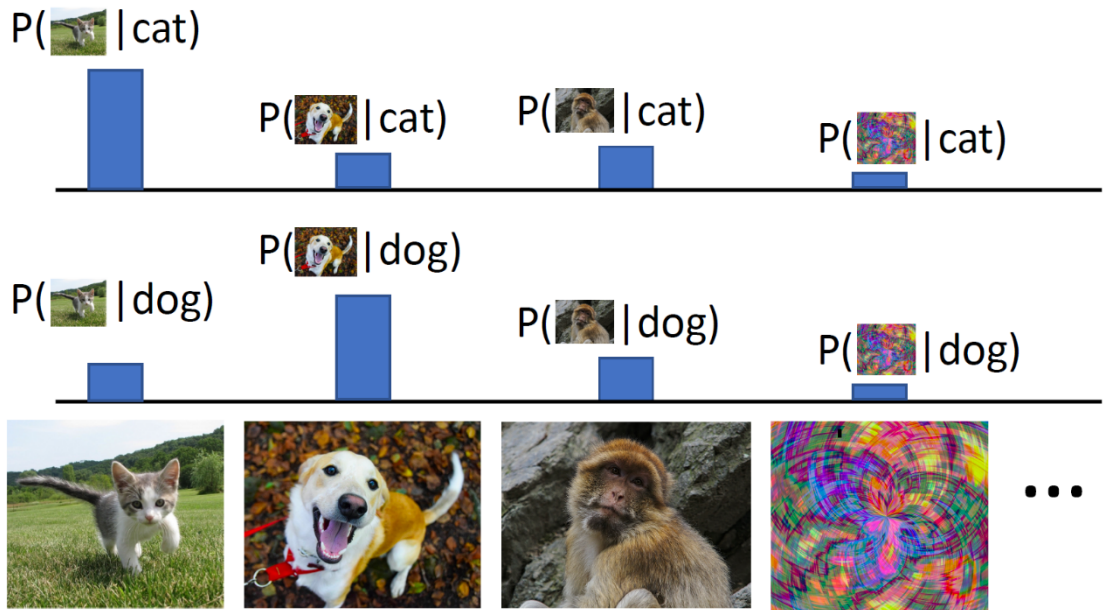
Discriminative Model:

Learn a probability distribution $p(y|x)$

Generative Model:

Learn a probability distribution $p(x)$

Conditional Generative Model: Learn $p(x|y)$

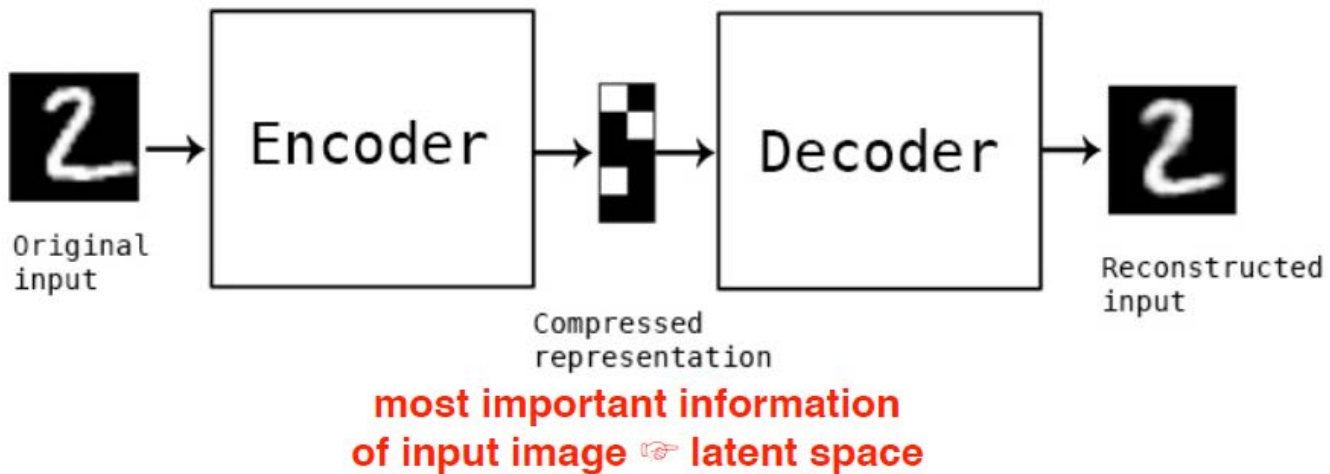


Conditional Generative Model: Each possible label induces a competition among all images

GenAI? Let's Start from Autoencoder

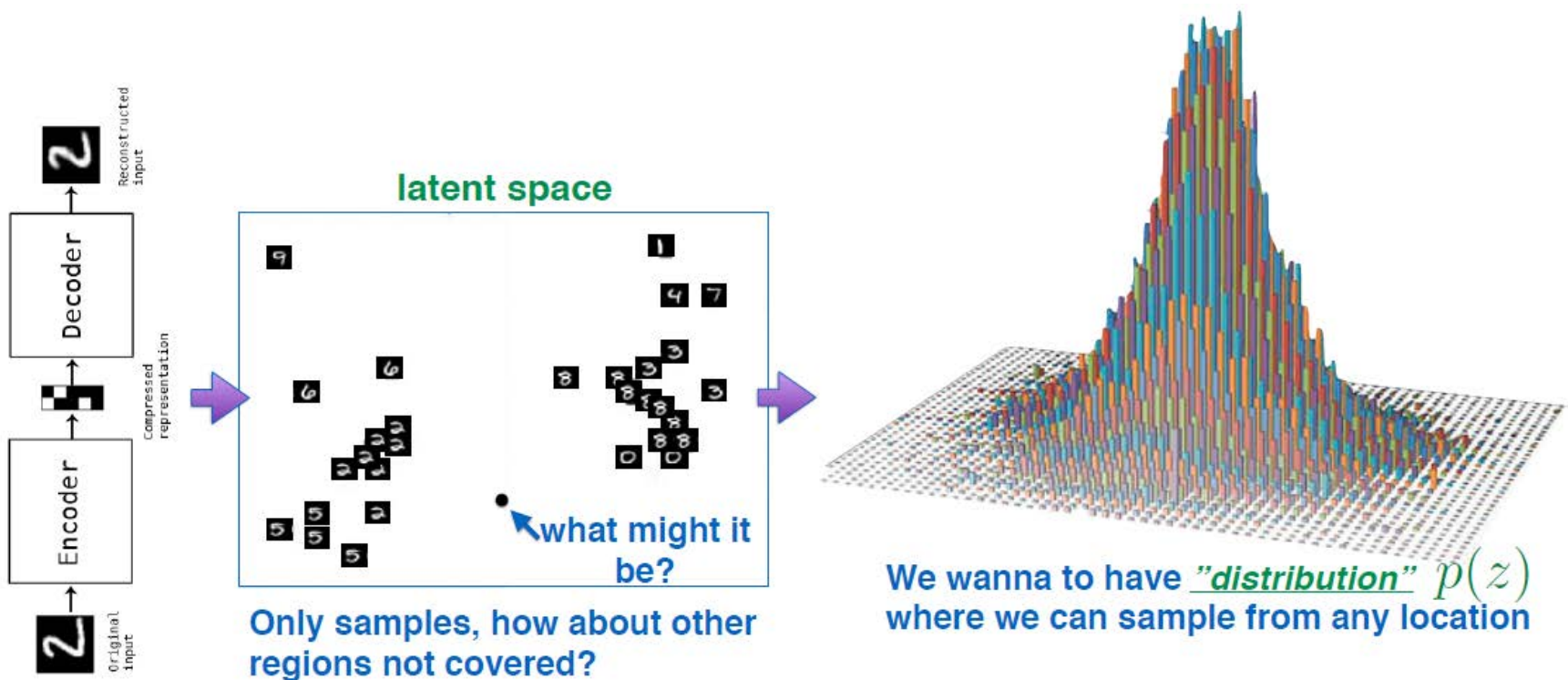
- Autoencoder

- Autoencoding = encoding itself with recovery purposes
- In other words, encode/decode data with reconstruction guarantees
- Latent variables/features as deep representations
- Example objective/loss function at output:
 - L2 norm between input and output, i.e.,



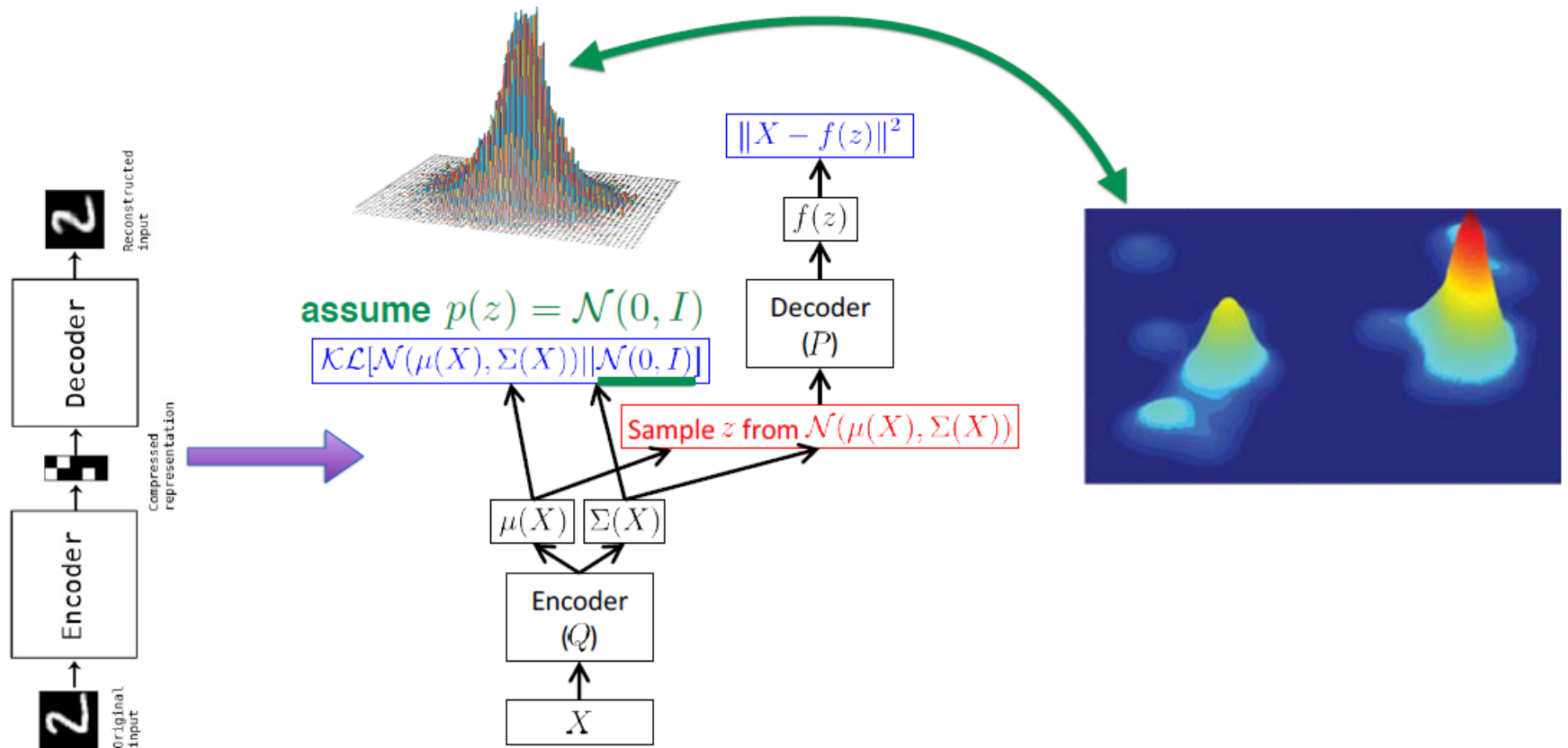
Take a Deep Look to Discover Latent Variables/Representations (cont'd)

- What's the Limitation of Autoencoder?



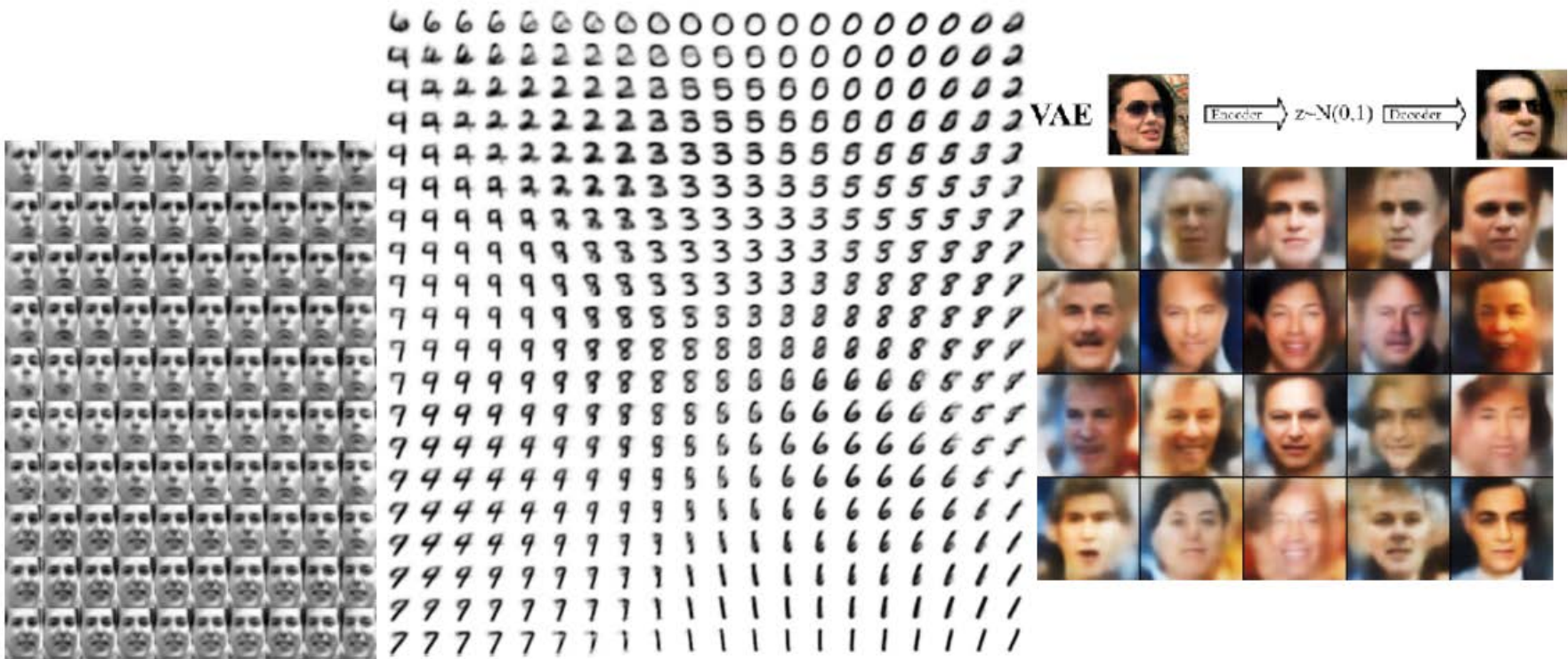
From Autoencoder to Variational Autoencoder

Now is a “distribution”, we can assume it to be a distribution easy to sample from, e.g. Gaussian



From Autoencoder to Variational Autoencoder (cont'd)

- Example Results

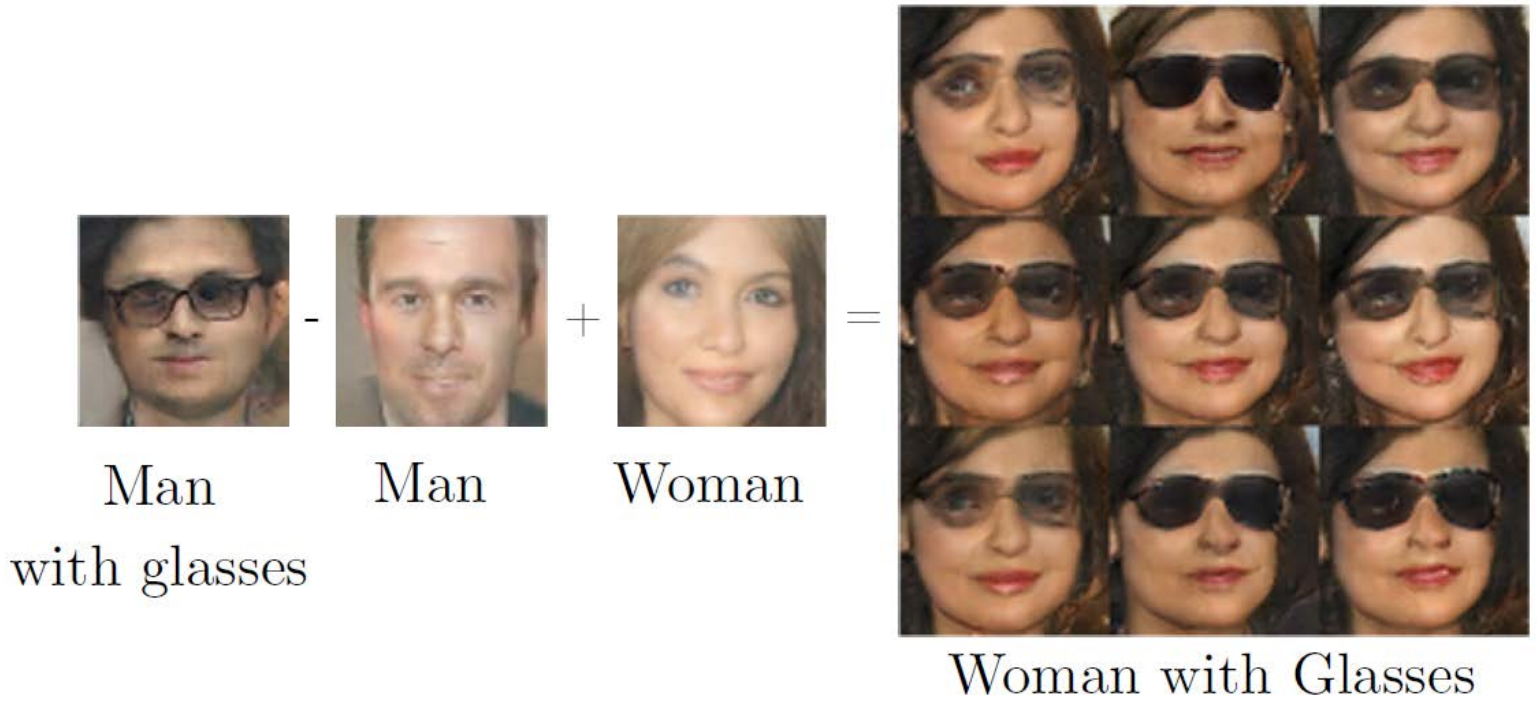


(a) Learned Frey Face manifold

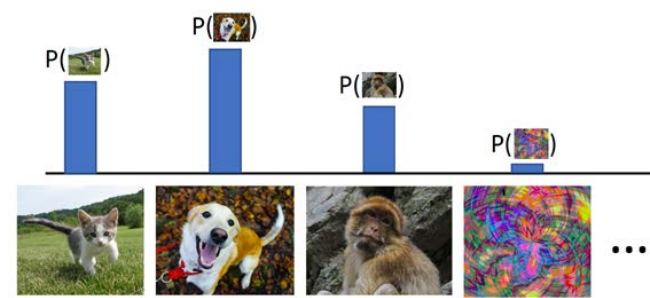
(b) Learned MNIST manifold

From Autoencoder to Variational Autoencoder (cont'd)

- Example Results
 - $A' - A + B = B'$

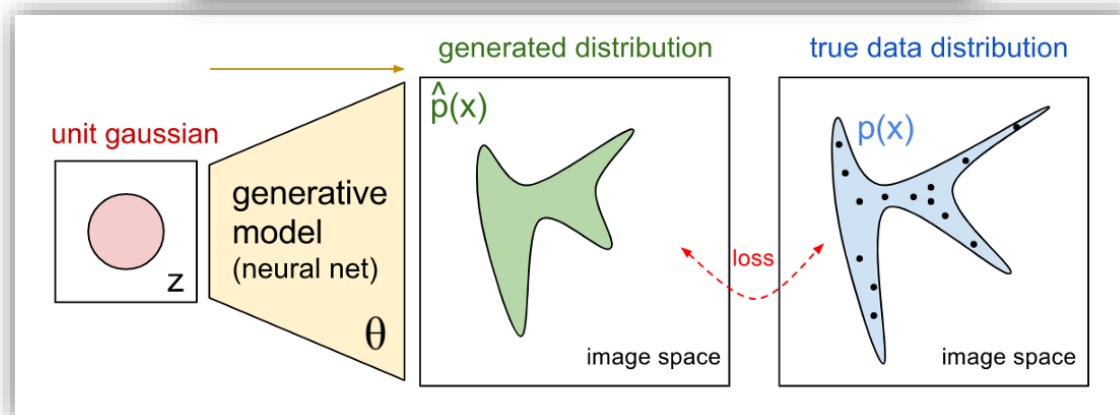
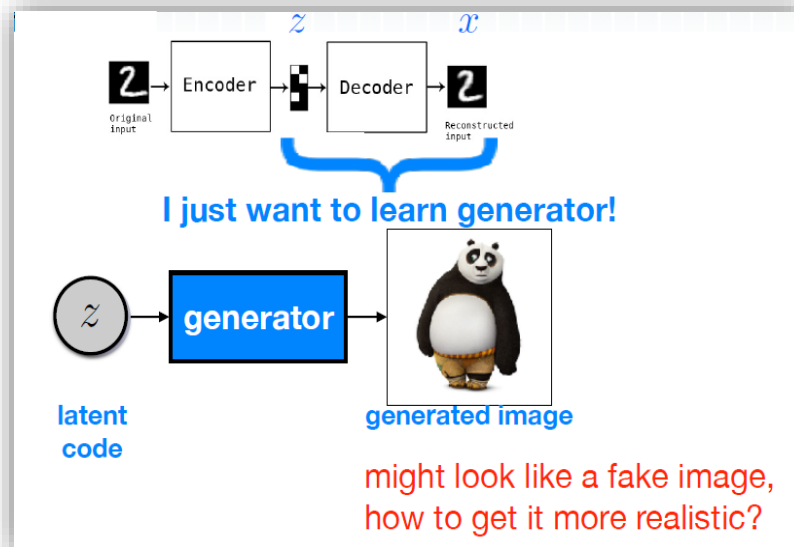


Limitation of VAE?



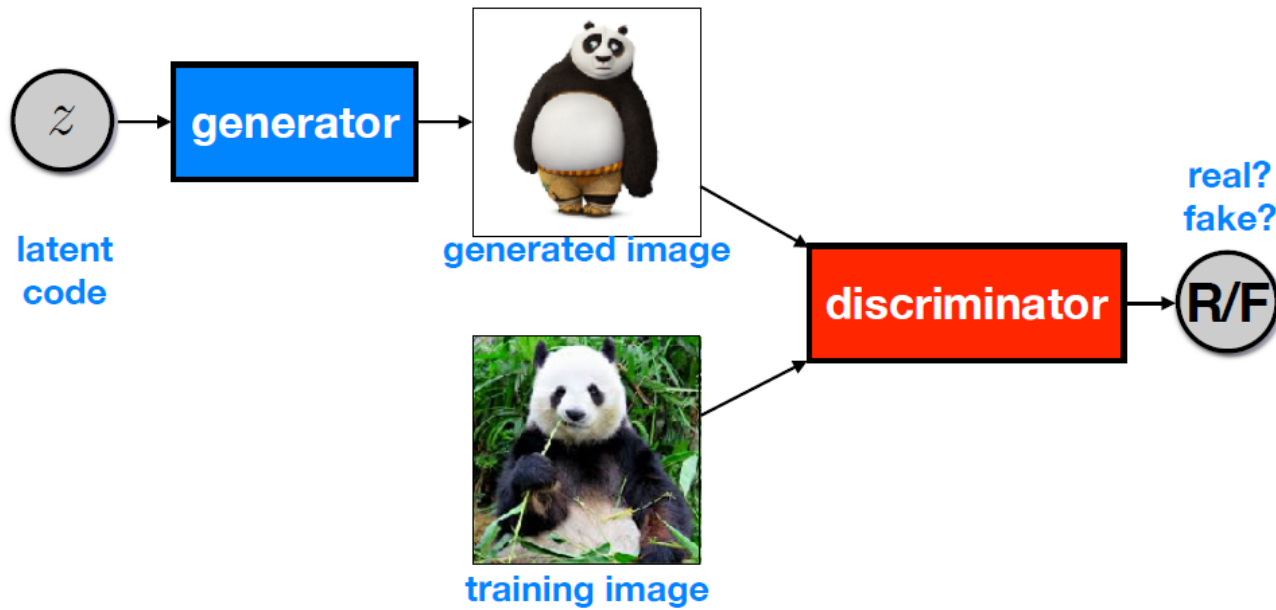
- Remarks

- Why Gaussian distribution is sufficient?
- What if we only need the decoder/generator in practice?
- How do we know if the output images are sufficiently good?



Generative Adversarial Network

- Idea
 - **Generator** to convert a vector z (sampled from P_z) into fake data x (from P_G), while we need $P_G = P_{\text{data}}$
 - **Discriminator** classifies data as real or fake (1/0)
 - How? Impose an **adversarial loss** on the observed data distribution!

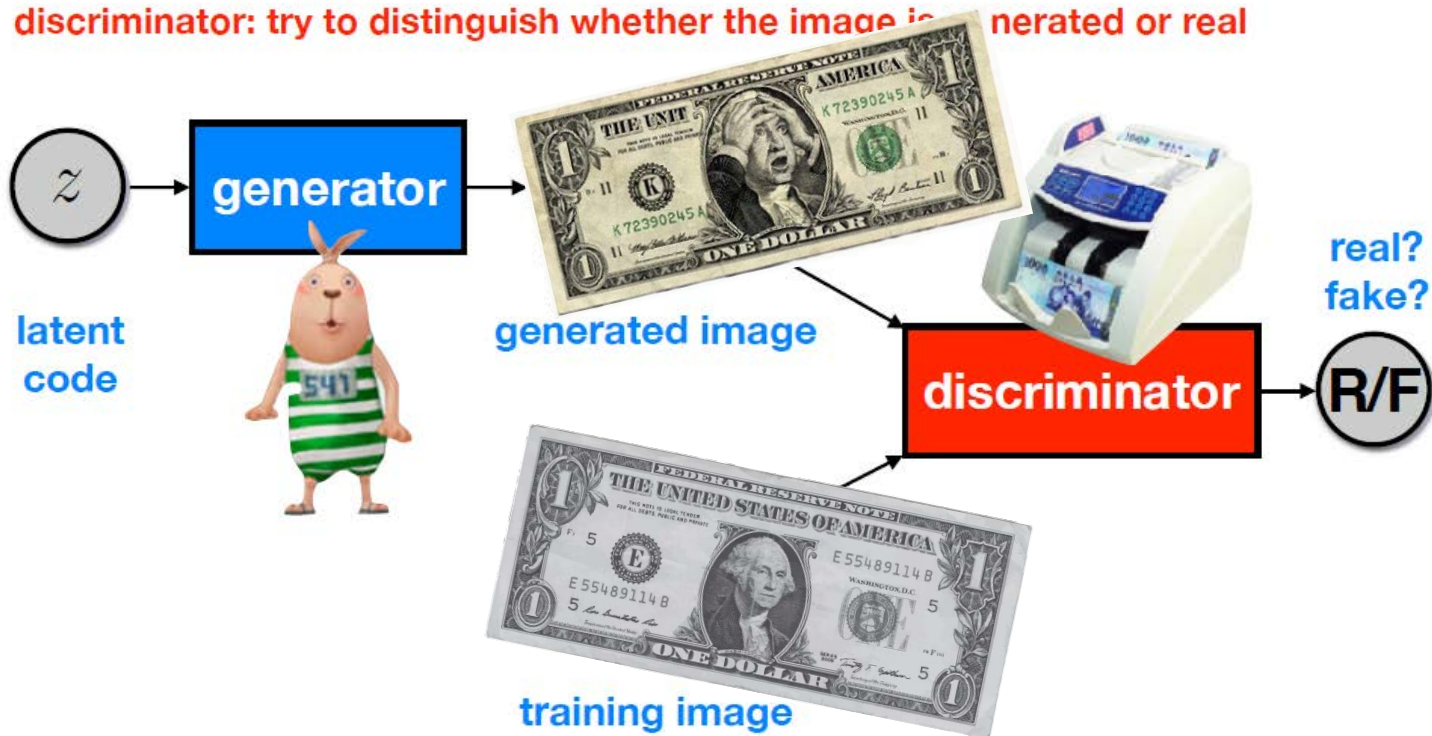


Generative Adversarial Network (cont'd)

- Key idea:
 - Impose *adversarial loss* on data distribution
 - Let's see a practical example...

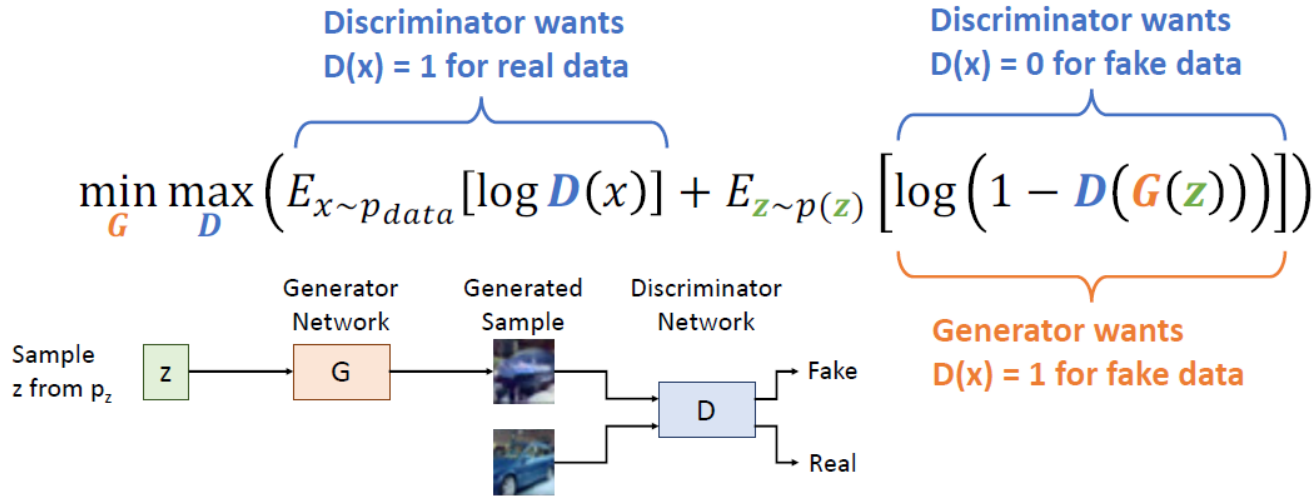
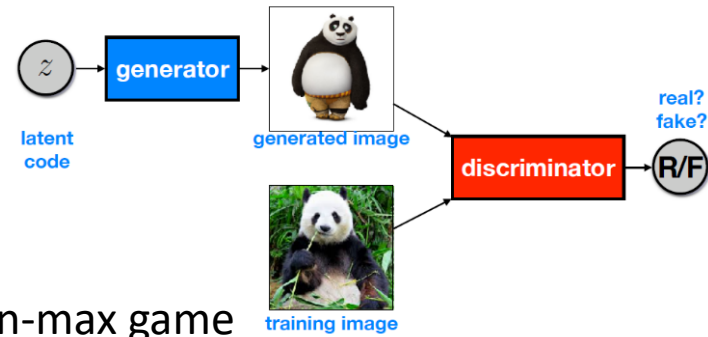
generator: try to generate more realistic images to cheat discriminator

discriminator: try to distinguish whether the image is generated or real



Training Objective of GAN

- Jointly train generator G and discriminator D with a min-max game



- Train G & D with alternating gradient updates

$$\min_G \max_D V(G, D)$$

For t in $1, \dots, T$:

- (Update **D**) $D = D + \alpha_D \frac{\partial V}{\partial D}$
- (Update **G**) $G = G - \alpha_G \frac{\partial V}{\partial G}$

Denoising Diffusion Models

- Recently emerging as powerful **visual** generative models
 - Unconditional image synthesis
 - Conditional image synthesis
 - Outperforms GANs

DALL·E 2

"a teddy bear on a skateboard in times square"



Diffusion Models Beat GANs on Image Synthesis, Dhariwai & Nochol, OpenAI, 2021

Imagen

A group of teddy bears in suit in a corporate office celebrating the birthday of their friend. There is a pizza cake on the desk.

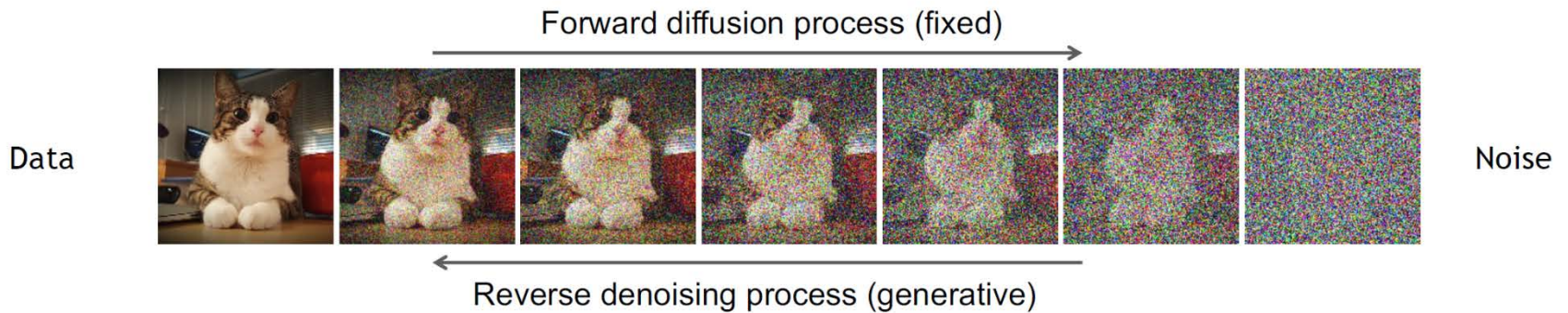
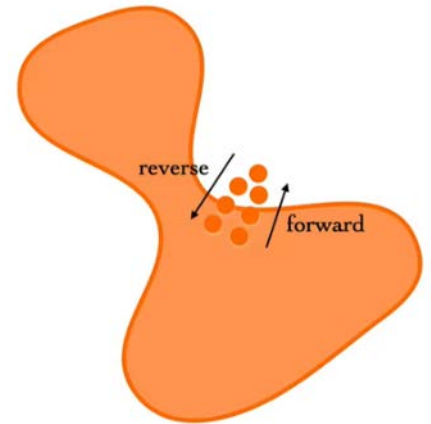


Cascaded Diffusion Models for High Fidelity Image Generation, Ho et al., Google, 2021

Denoising Diffusion Probabilistic Models (DDPM)

Learning to generate by denoising

- 2 processes required for training:
 - Forward diffusion process – gradually add noise to input
 - **Reverse diffusion process**
 - learns to generate/restore data by denoising
 - typically implemented via a **conditional U-net**
 - Comments about noise scheduling (see next slide)



Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020

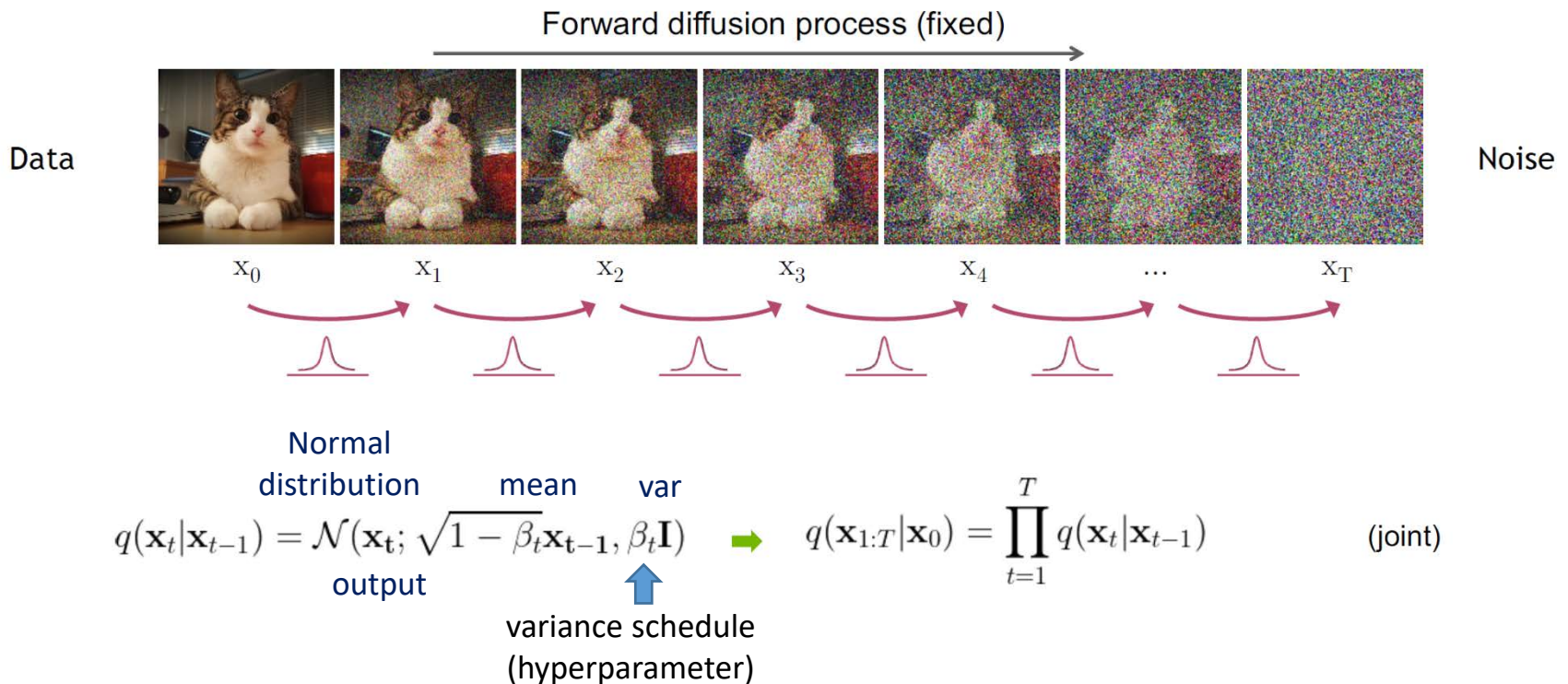
Song et al., Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021

DDPM:

Learning to generate by denoising (cont'd)

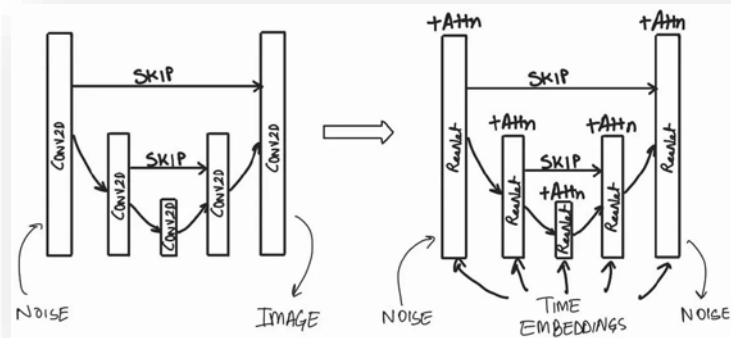
- Forward diffusion process

- Gradually add noise to the input in T steps (e.g., via linear scheduling)
- Recall that x_0 denotes clean input image, and x_T is the final noisy one.
- Comments on $q(x_t|x_{t-1})$



Learning of Diffusion Models

- Summary
 - Training and sample generation



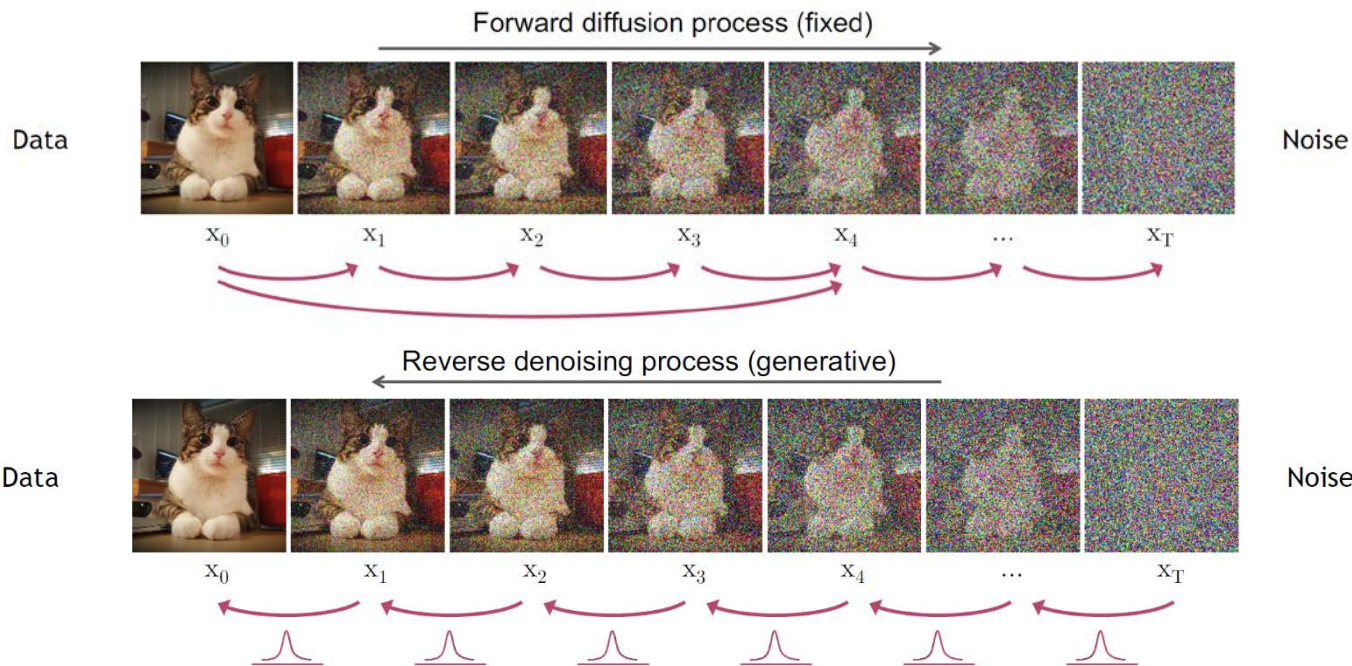
Algorithm 1 Training

- 1: **repeat**
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on $\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$
- 6: **until** converged

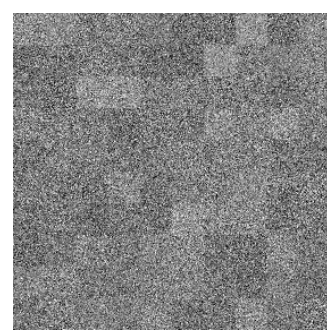
Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for** $t = T, \dots, 1$ **do**
- 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return** \mathbf{x}_0

<https://medium.com/@vedantjumle/class-conditioned-diffusion-models-using-keras-and-tensorflow-9997fa6d958c>



Learning of Diffusion Models



E.g., MNIST

- Summary

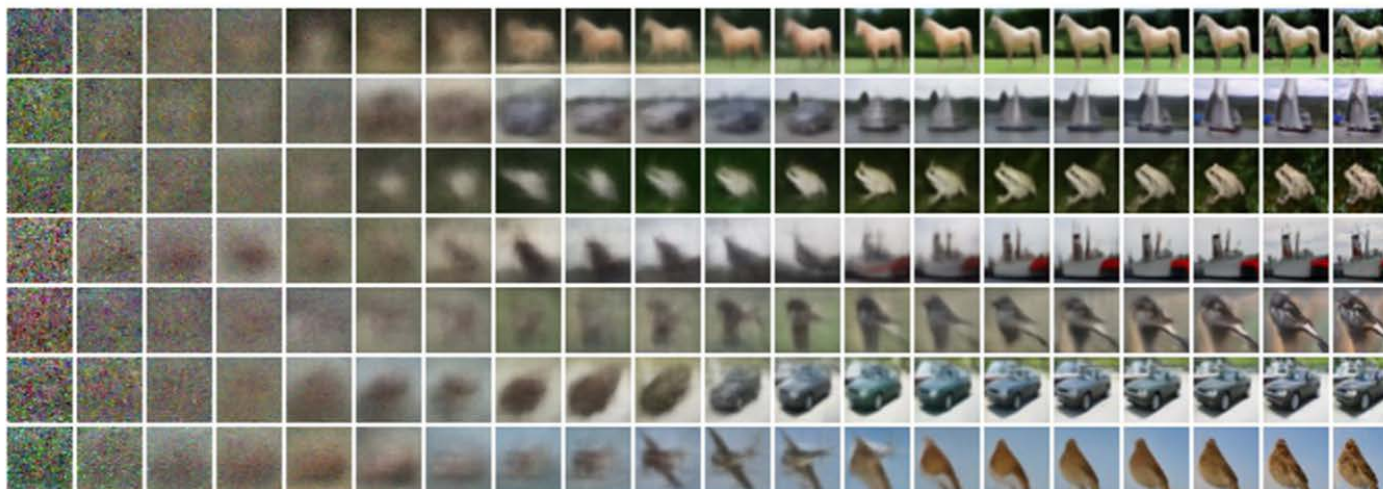
- Training and sample generation

Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on
$$\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon; t)\|^2$$
 - 6: **until** converged
-

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-



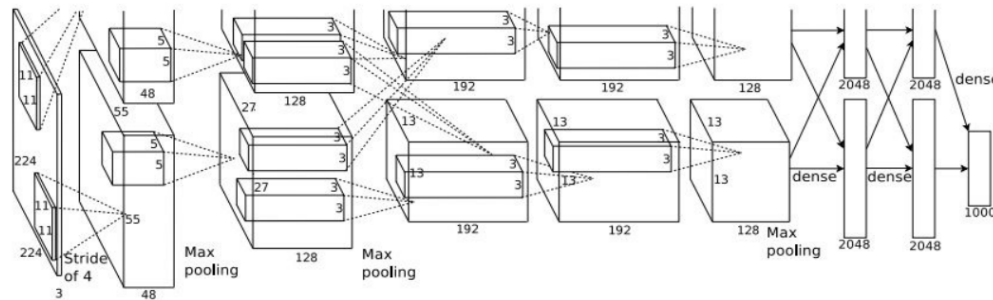
More steps →

What'd Be Covered in This Crash Course...

- **Learning-based Computer Vision**
 - From Linear to Non-Linear Classifiers
- **Start of Deep Learning for Computer Vision**
 - Convolutional Neural Networks
 - Self-Supervised Learning
 - Segmentation & Detection
- **Generative Models**
 - Autoencoder, Variational Autoencoder, Generative Adversarial Networks & Diffusion Models
- **Sequence-to-Sequence Learning**
 - Attention is All You Need: Transformer
- **Vision & Language Foundation Models**
 - Image-to-Text vs. Text-to-Image
 - Parameter-Efficient Fine-tuning

What Are The Limitations of CNN?

- Deal with image data
 - Both input and output are images/vectors
- Simply feed-forward processing



More Applications in Vision

Image Captioning

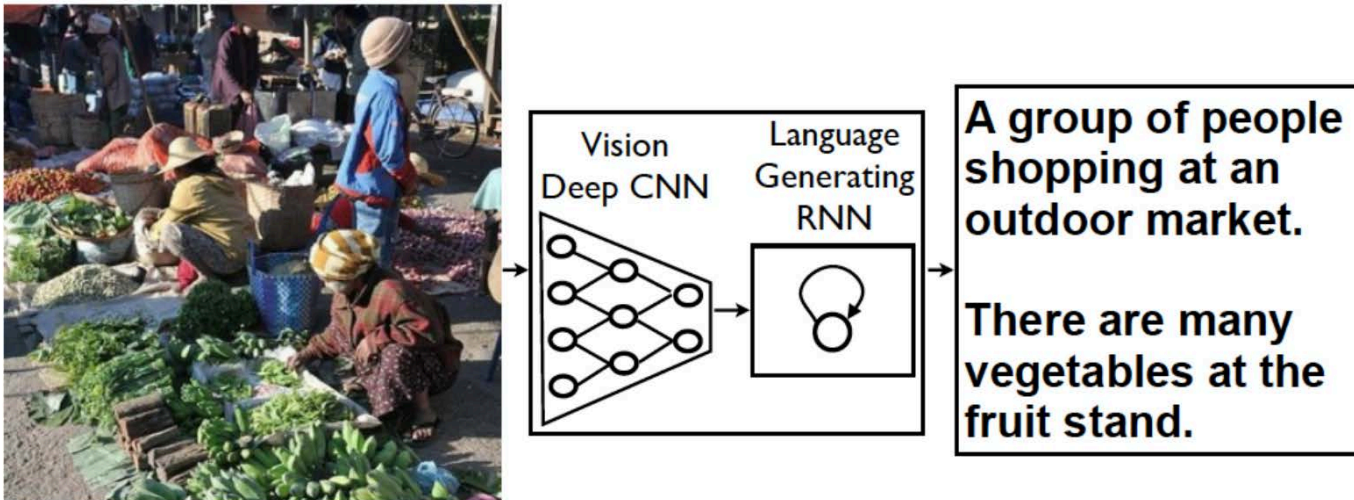
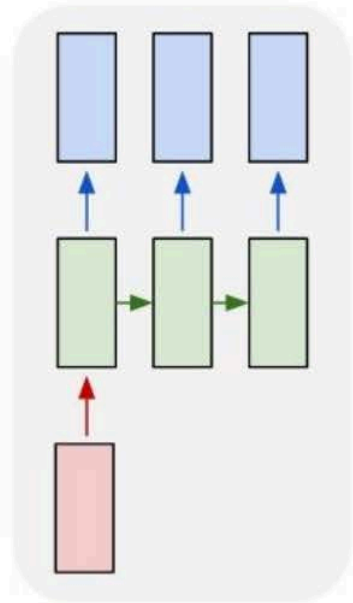


Figure from Vinyals et al, "Show and tell: A neural image caption generator", CVPR 2015

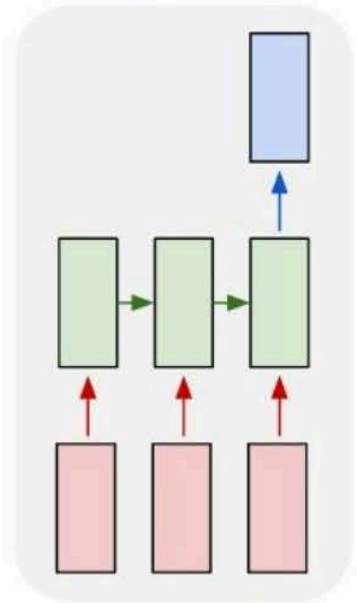
Sequence-to-Sequence Modeling

one to many



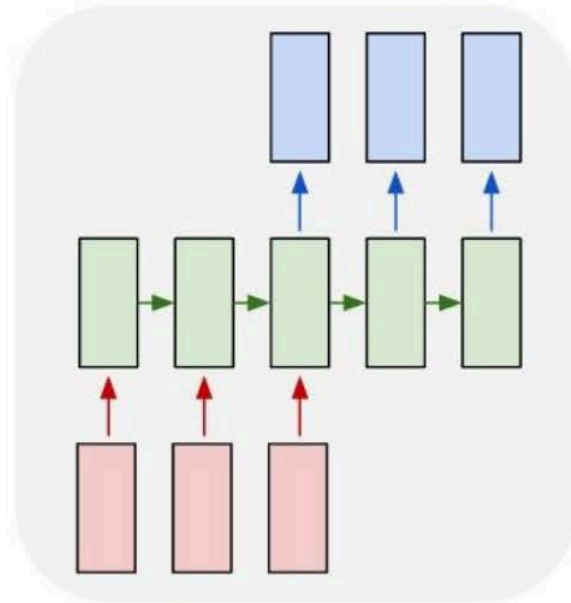
e.g., image caption

many to one



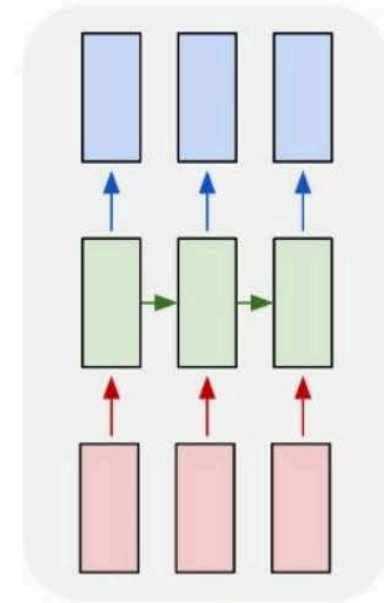
e.g., action recognition

many to many



e.g., video prediction

many to many



e.g., video indexing

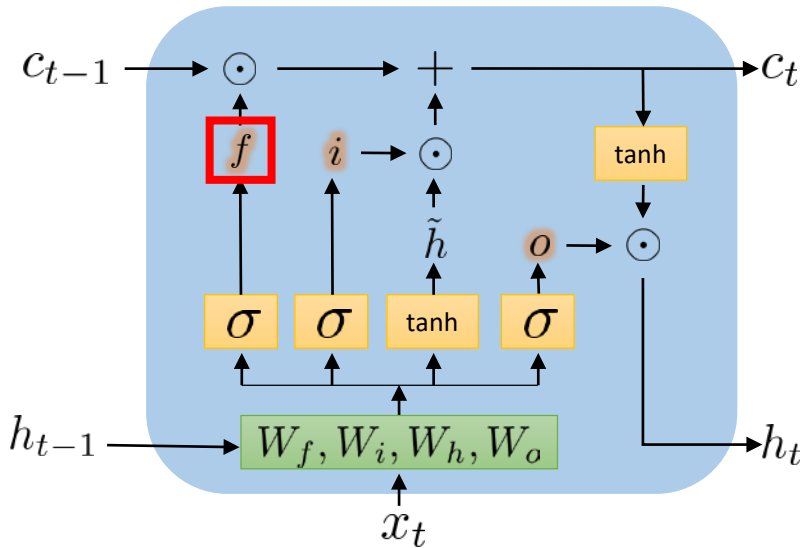
Vanilla RNN, LSTM, & GRU

RNN

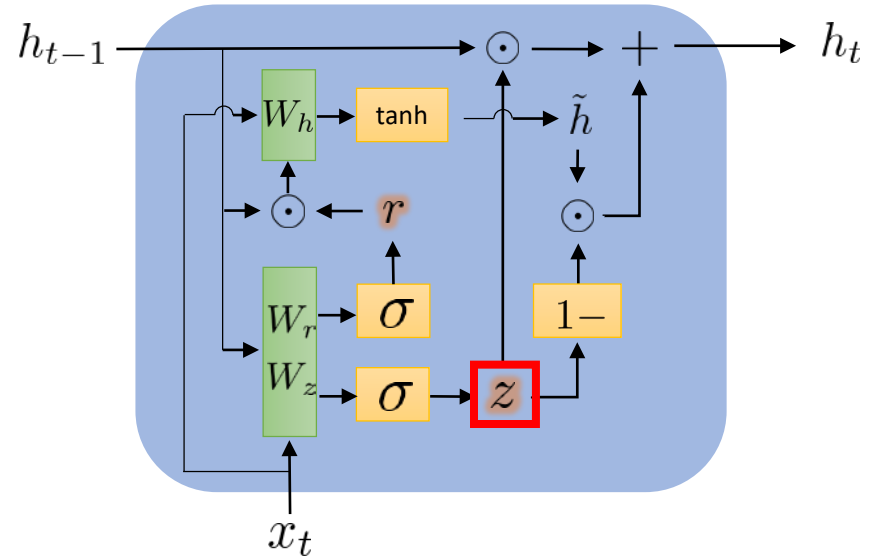
$$h_t = \tanh \left(W \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right)$$

Output in time t

Input in time t



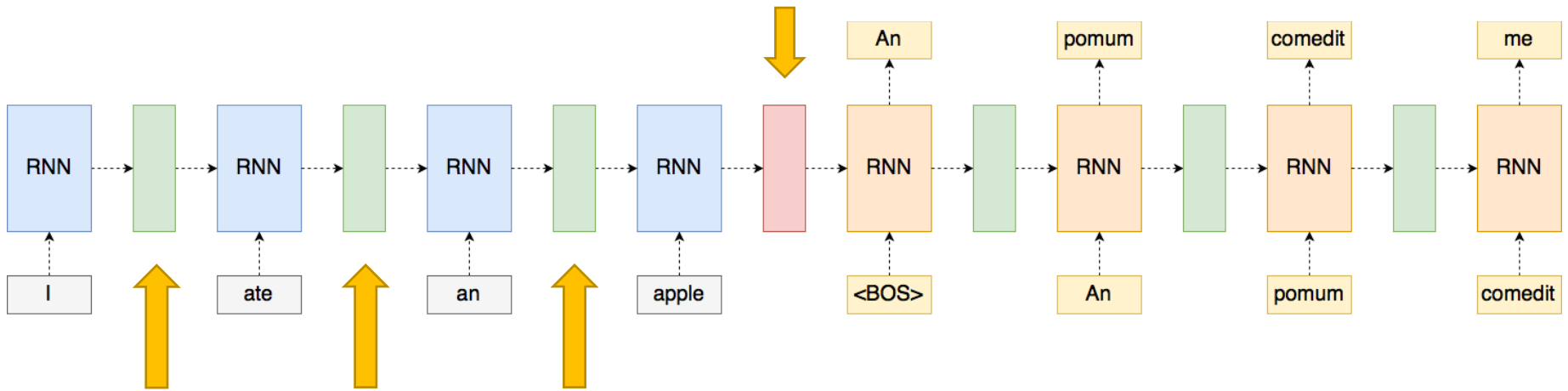
LSTM



GRU

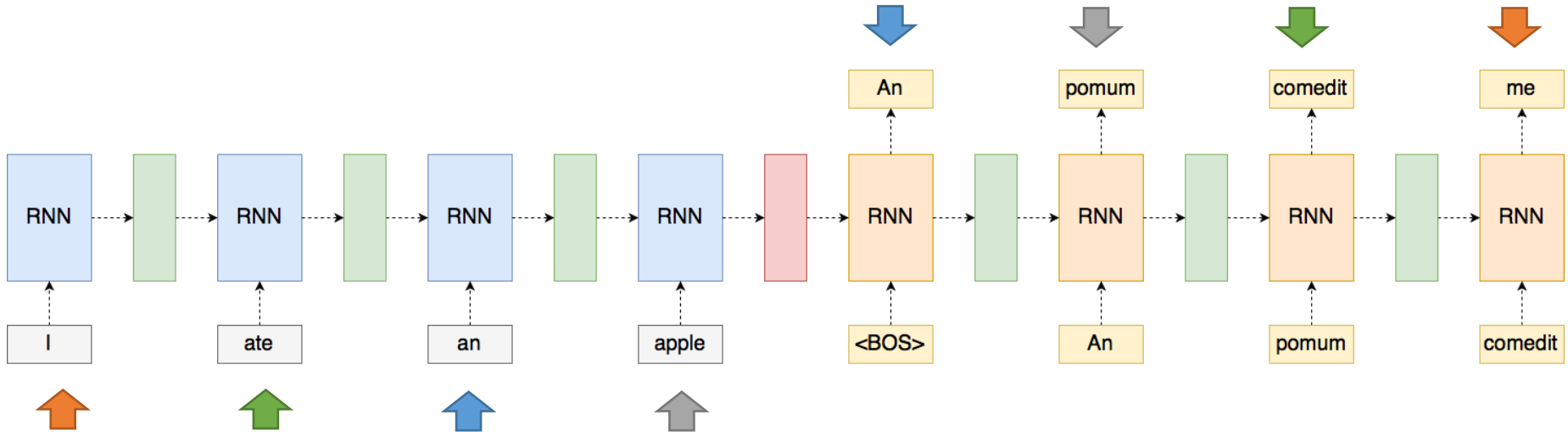
What's the Potential Problem in RNN?

- Each hidden state vector extracts/carries information across time steps (some might be diluted downstream).
- Information of the entire input sequence is embedded into a **single hidden state** vector.



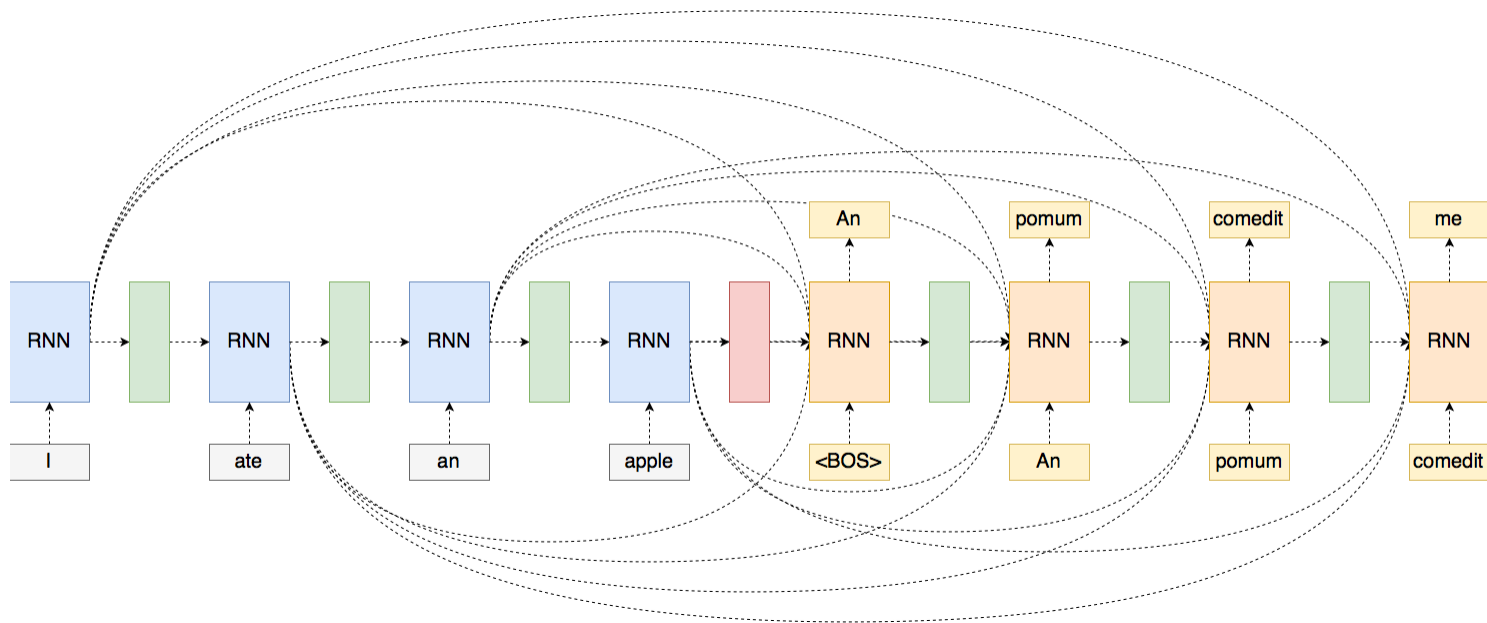
What's the Potential Problem? (cont'd)

- Outputs at different time steps have particular meanings.
- However, **synchrony** between **input** and **output seqs** is **not** required.



What's the Potential Problem? (cont'd)

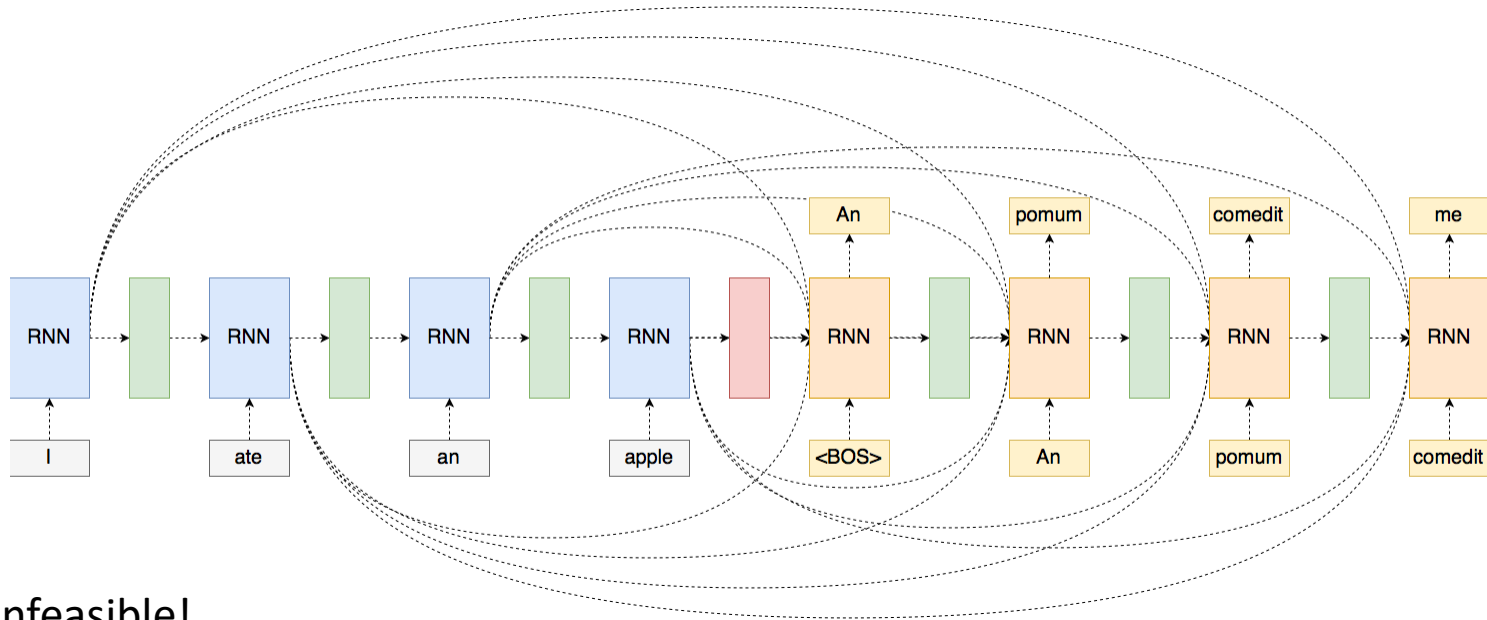
- Connecting every hidden state between encoder and decoder?



- Infeasible!
 - Both inputs and outputs are with varying sizes.
 - **Overparameterized**
 - Possible solution: **attention**

RNN with Attention is Good, But..

- Attention in a pre-defined sequential order
- Information loss due to long sequences...
- Connecting every hidden state between encoder and decoder?

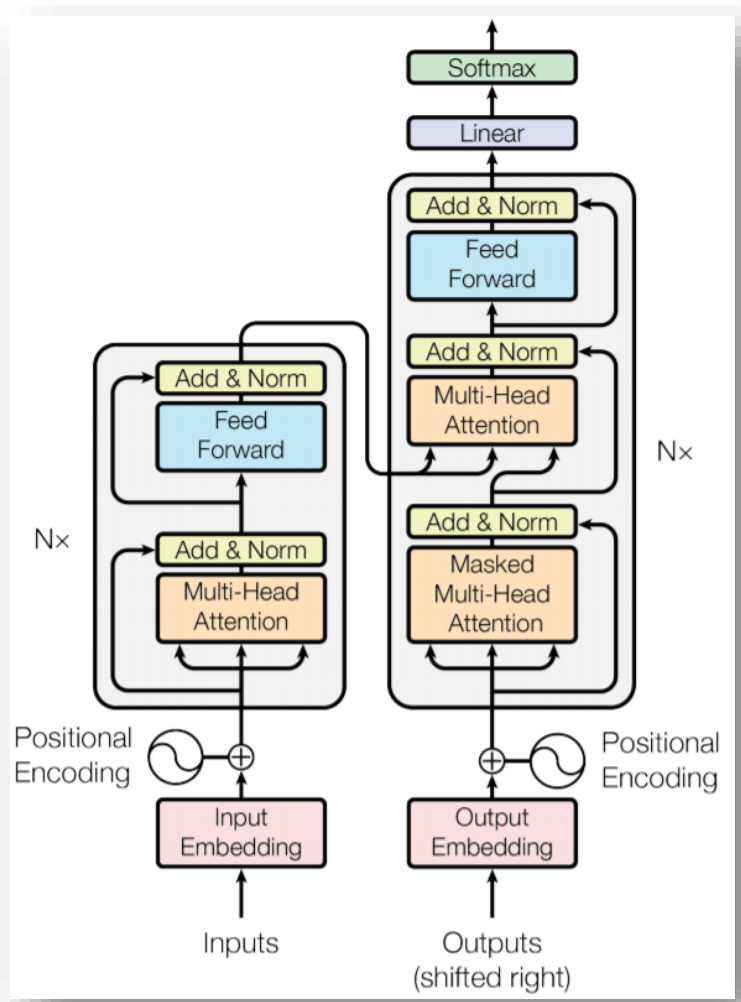
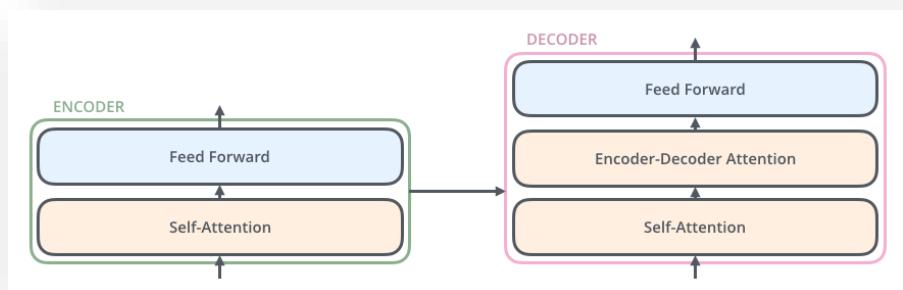


- Infeasible!
 - Both inputs and outputs are with varying sizes.
 - Overparameterized



Transformer

- “Attention is all you need”, NIPS/NeurIPS 2017
- Self-attention for text translation
- Say goodbye to CNN & RNN
- More details available at: <http://jalamar.github.io/illustrated-transformer/>



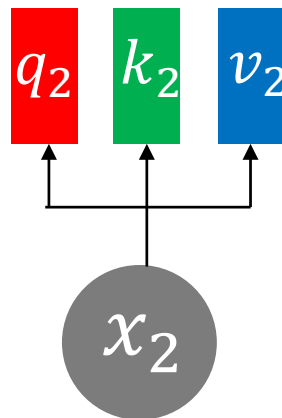
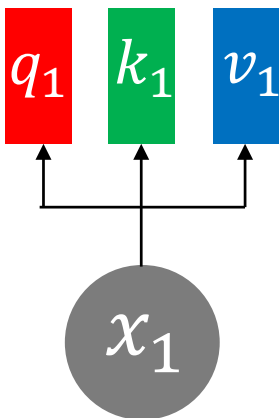
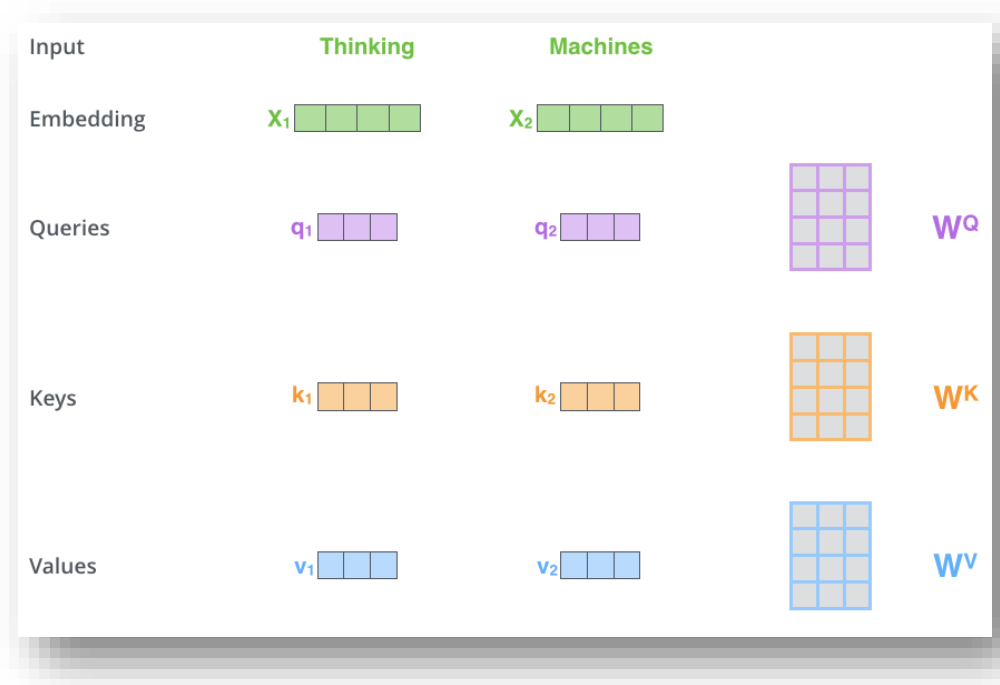
Self-Attention (1/5)

- Query q , key k , value v vectors are learned from each input x

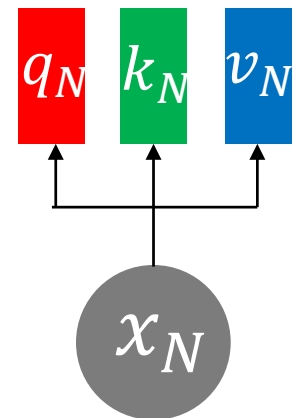
$$q_i = W^Q x_i$$

$$k_i = W^K x_i$$

$$v_i = W^V x_i$$



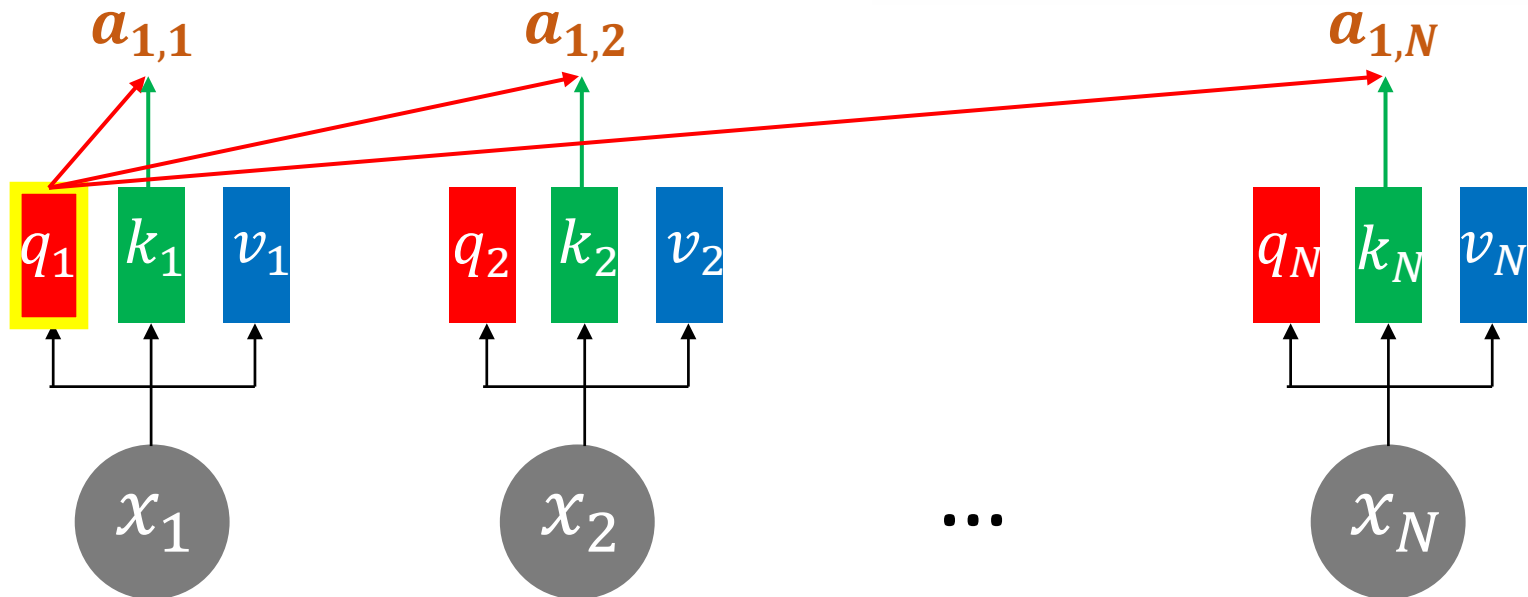
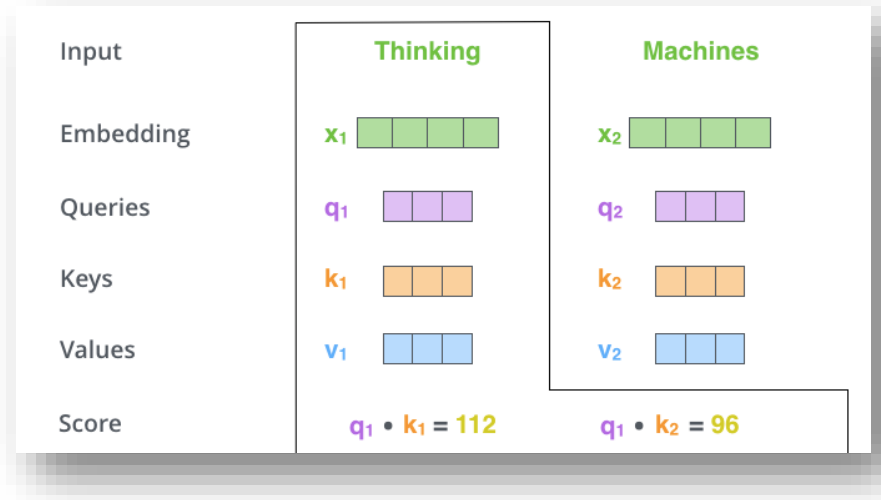
...



Self-Attention (2/5)

- Relation between each input is modeled by inner-product of **query** q and **key** k .

$$a_{1,i} = \frac{q_1 \cdot k_i}{\sqrt{d}}, \text{ where } a \in R, q, k \in R^d$$

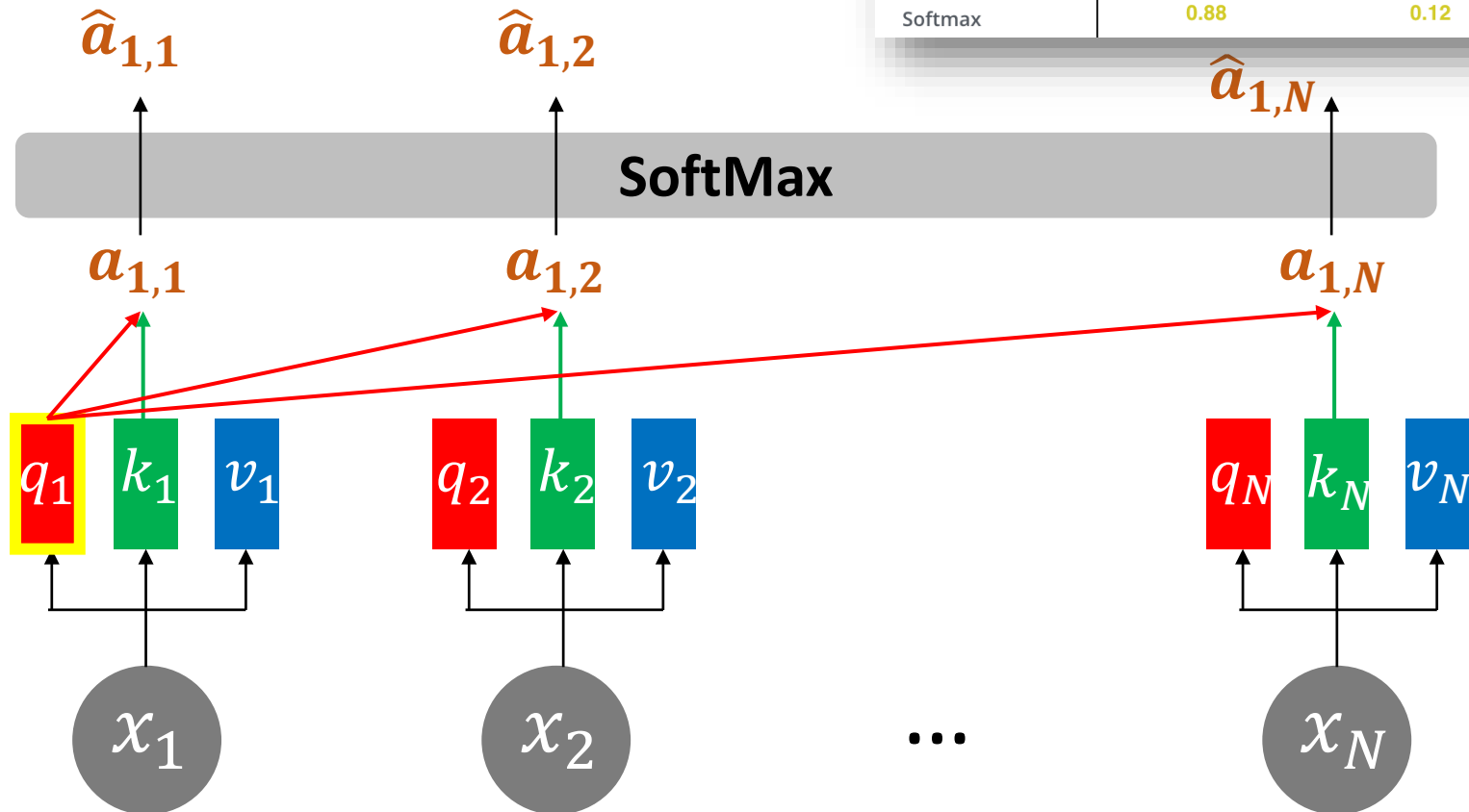


Self-Attention (3/5)

- SoftMax is applied:

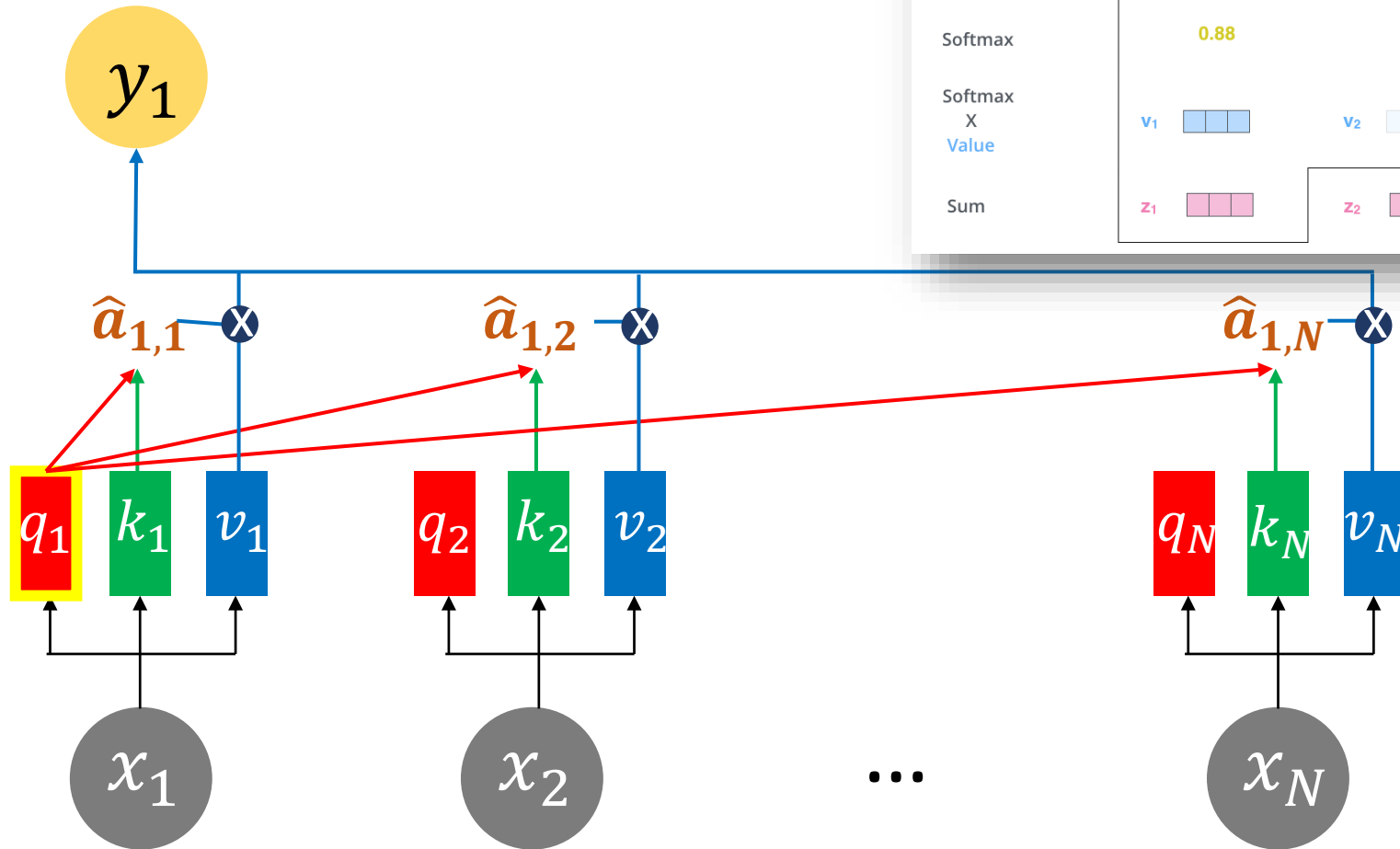
$$0 \leq \hat{a}_i = e^{a_i} / \sum_j^N e^{a_j} \leq 1, \text{ for } i=1, \dots, N$$

Input	Thinking	Machines
Embedding	x_1	x_2
Queries	q_1	q_2
Keys	k_1	k_2
Values	v_1	v_2
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Divide by $8 (\sqrt{d_k})$	14	12
Softmax	0.88	0.12



Self-Attention (4/5)

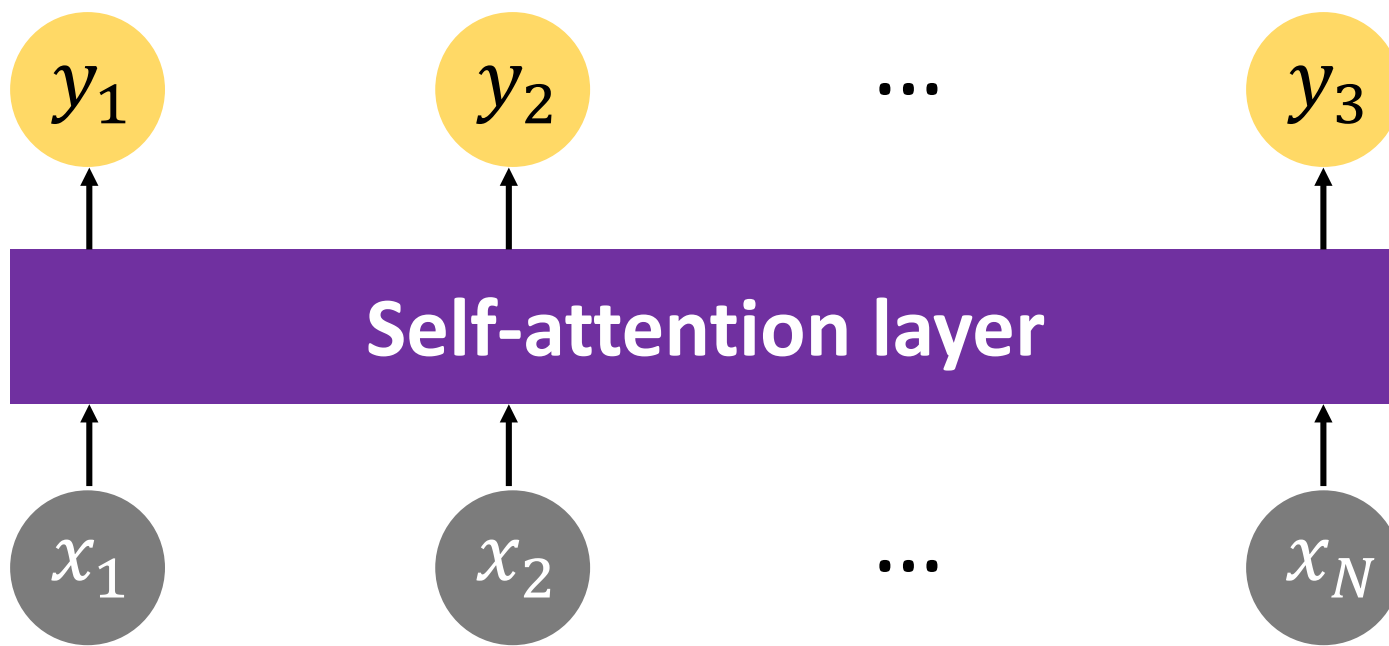
- Value vectors v are aggregated with attention weight \hat{a} , i.e., $y_1 = \sum_i^N \hat{a}_i \cdot v_i$



	Thinking	Machines
Input		
Embedding	x_1 [] [] [] []	x_2 [] [] [] []
Queries	q_1 [] []	q_2 [] []
Keys	k_1 [] []	k_2 [] []
Values	v_1 [] []	v_2 [] []
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Divide by $8 (\sqrt{d_k})$	14	12
Softmax	0.88	0.12
Softmax X Value	v_1 [] []	v_2 [] []
Sum	z_1 [] []	z_2 [] []

Self-Attention (5/5)

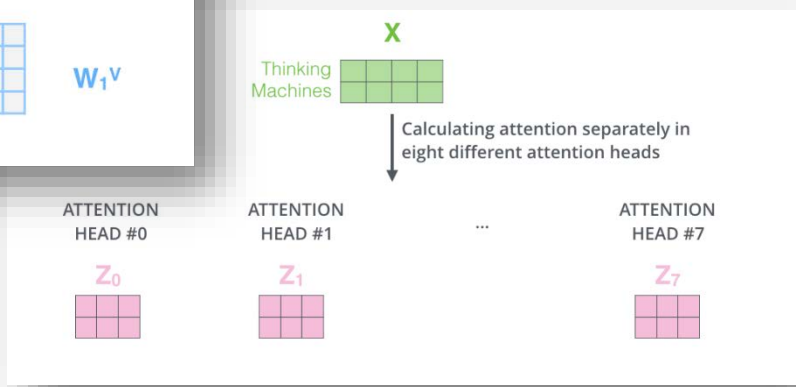
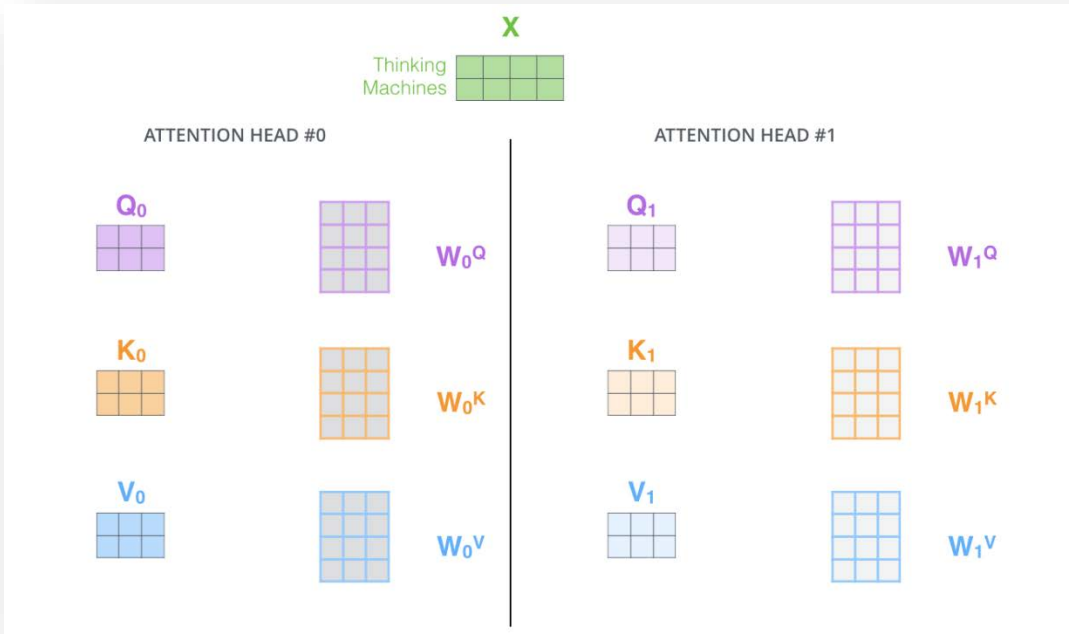
- All y_i can be computed **in parallel**
- Each y_i considers $x_1 \sim x_N$, modeling their **long-distance dependencies**.
- Global feature can be obtained by **average-pooling** over $y_1 \sim y_N$





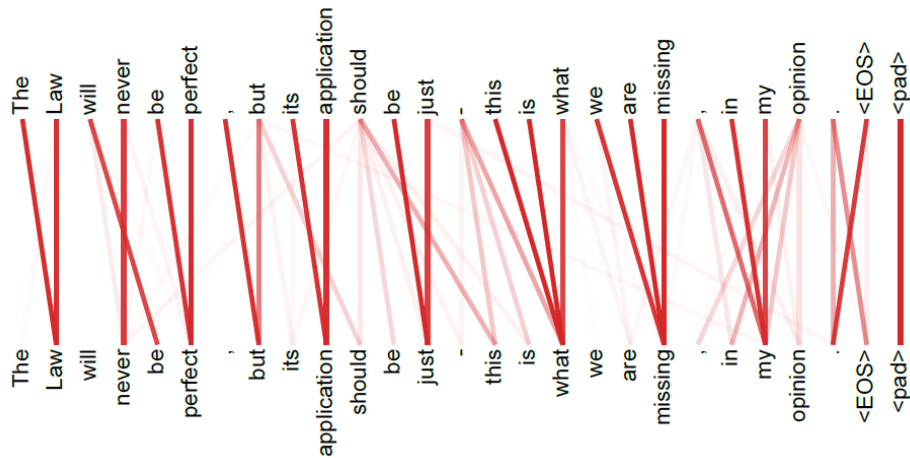
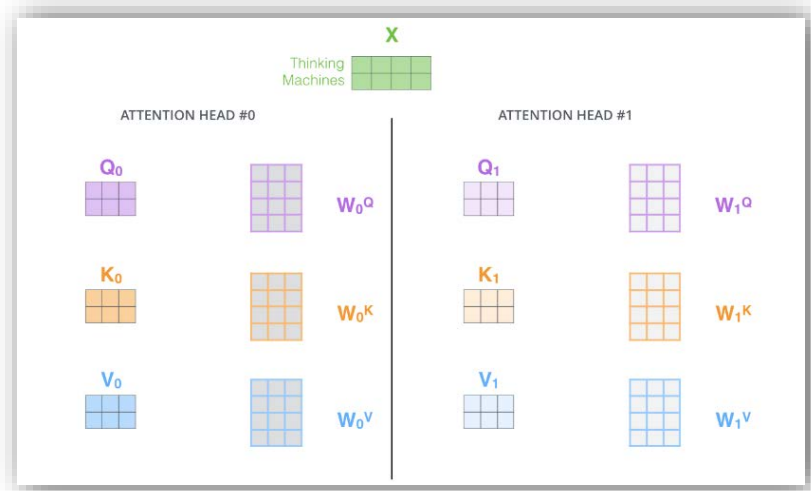
Multi-Head Self-Attention (1/4)

- Perform self-attention at **different subspaces**, implying performing attention over different input feature types (e.g., representations, modalities, positions, etc.)

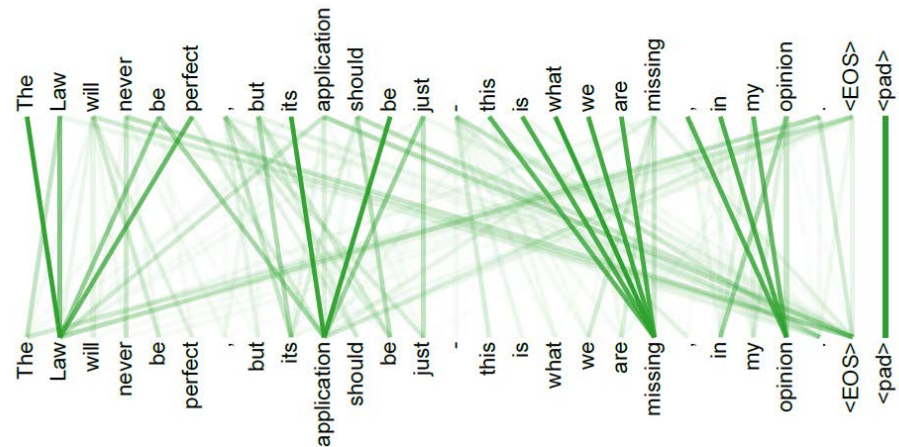


Multi-Head Self-Attention (2/4)

- Perform self-attention at **different subspaces**, implying performing attention over different input feature types
- See example below



Attention weights of Head 1

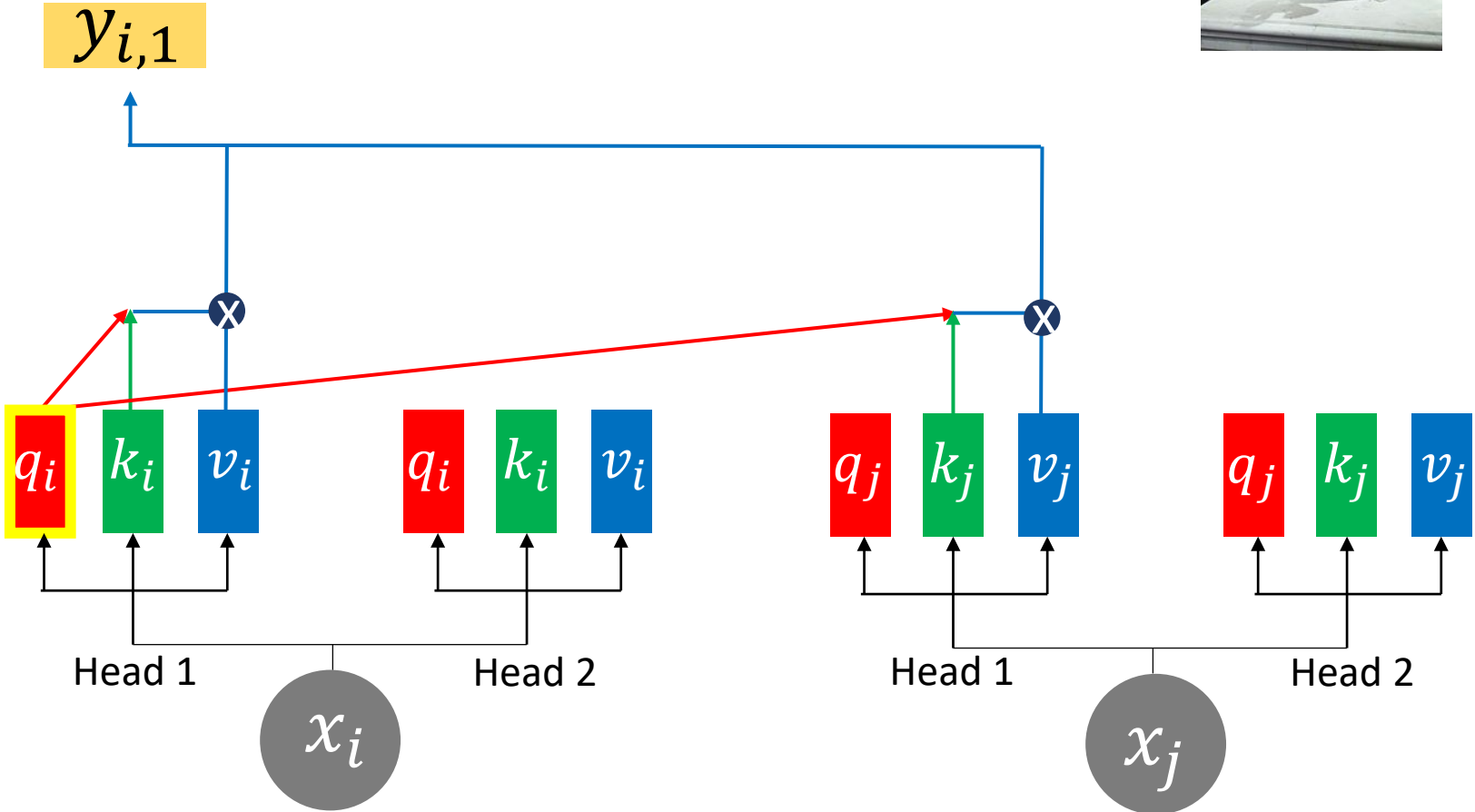


Attention weights of Head 2



Multi-Head Self-Attention (3/4)

- A 2-head example, output of two heads are concatenated.

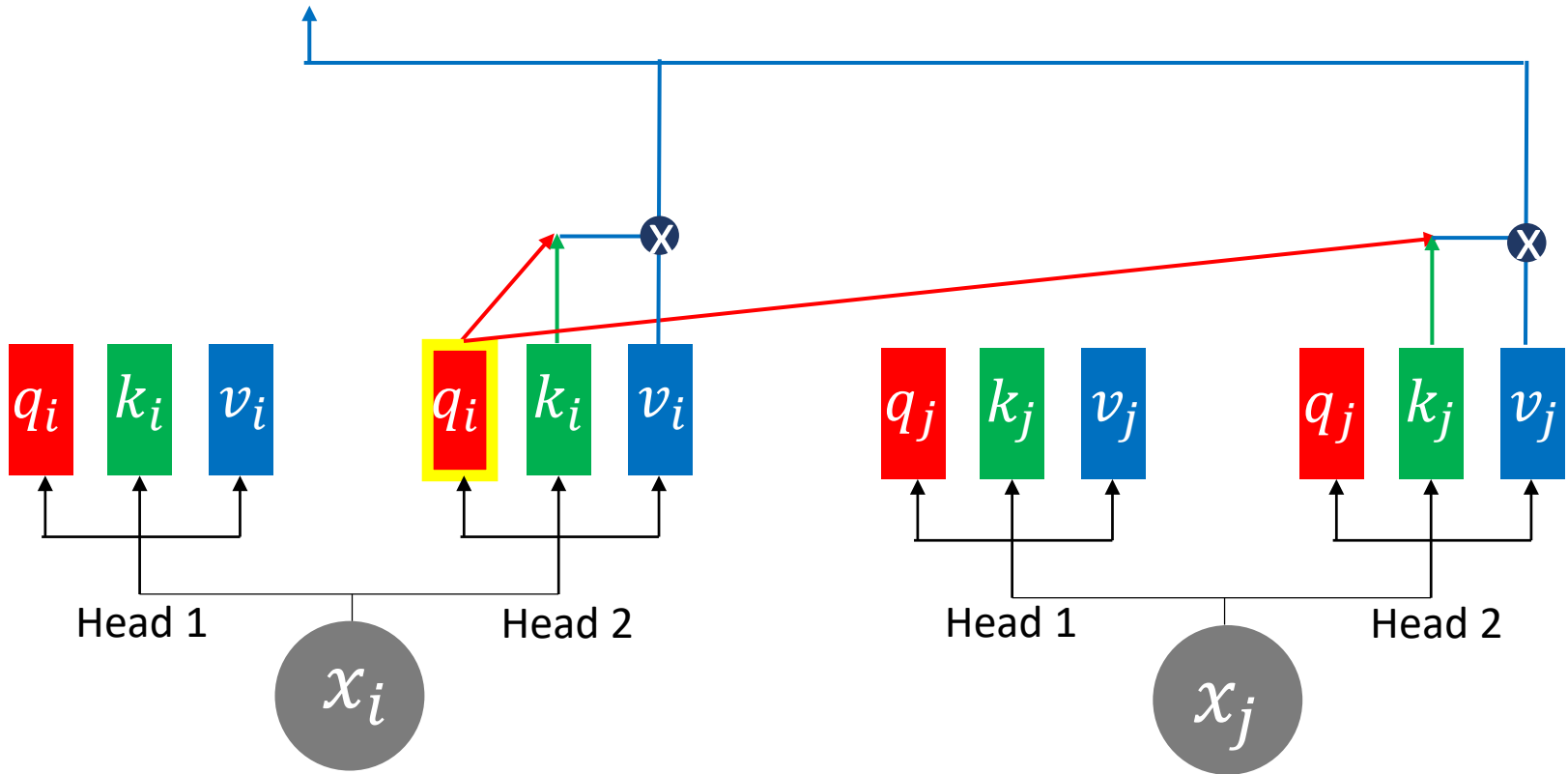




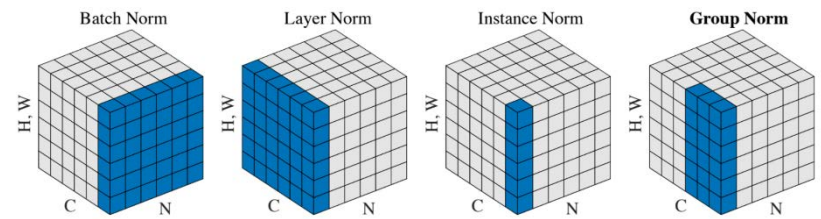
Multi-Head Self-Attention (4/4)

- A 2-head example, output of two heads are concatenated.

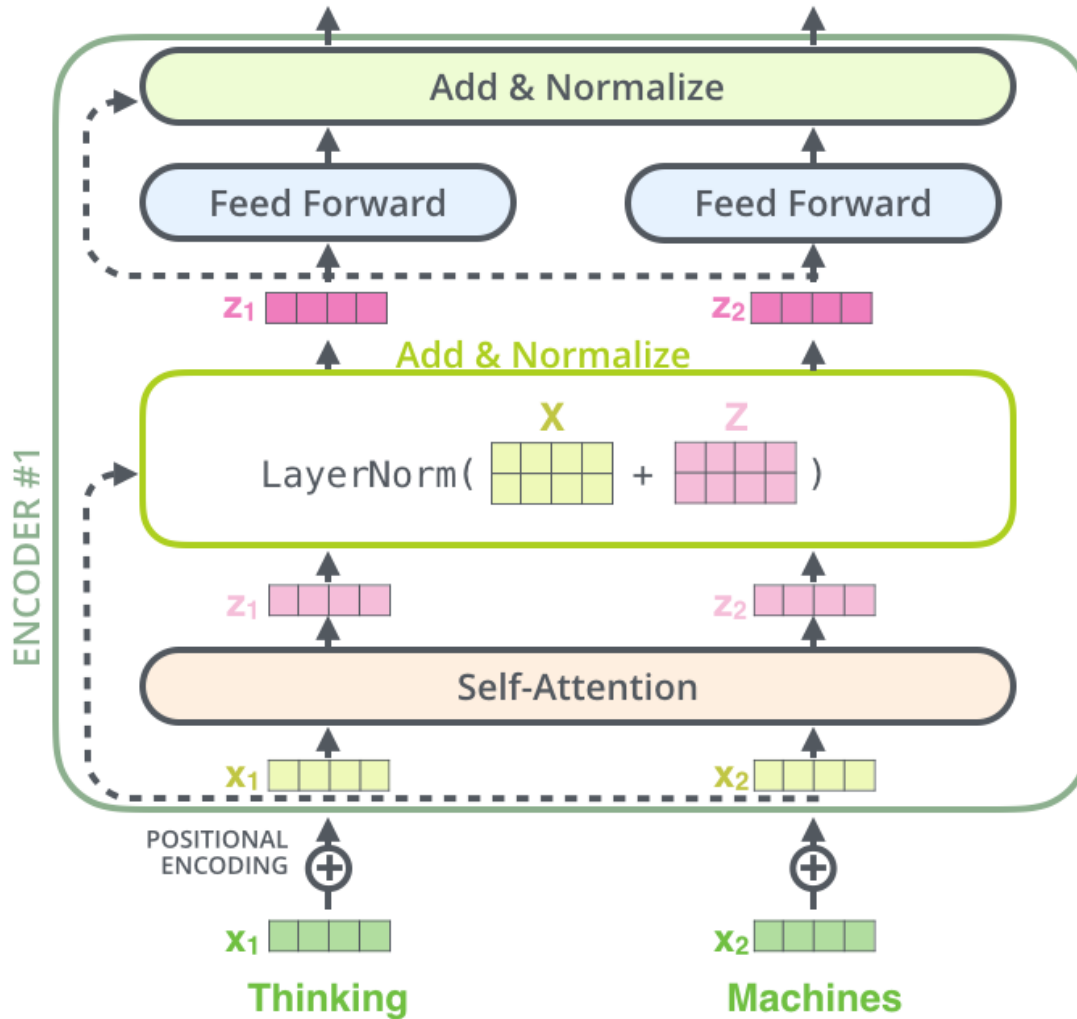
$$y_{i,1} \parallel y_{i,2} = y_i$$



The Residuals

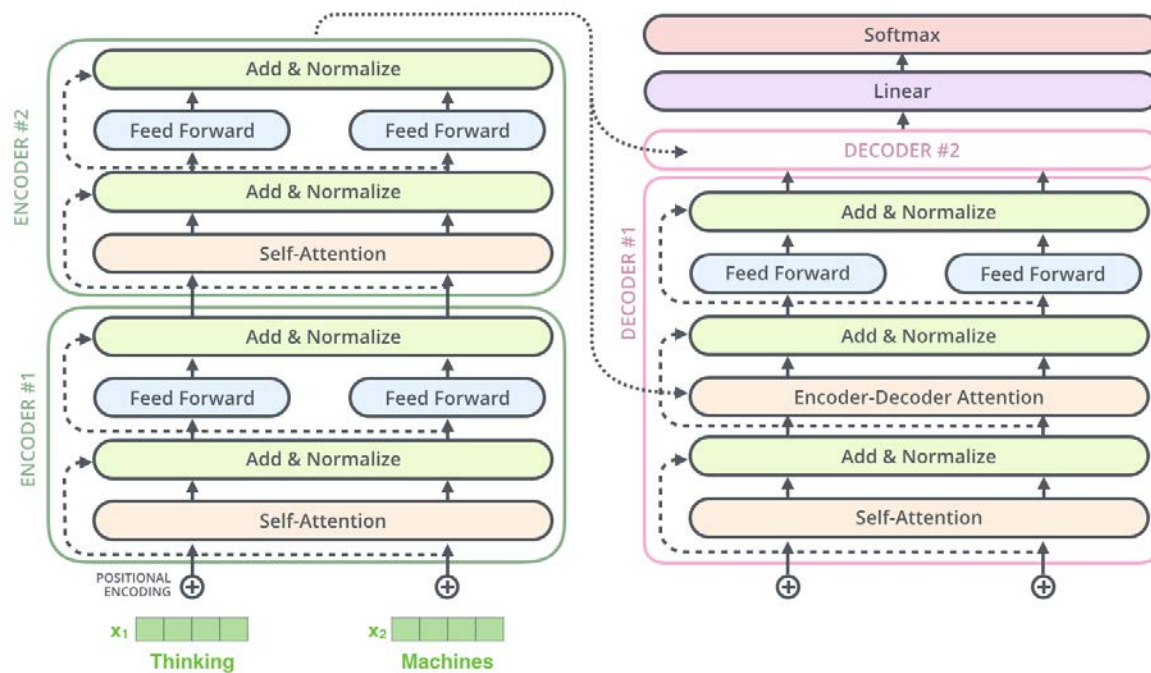


- A residual connection followed by layer normalization



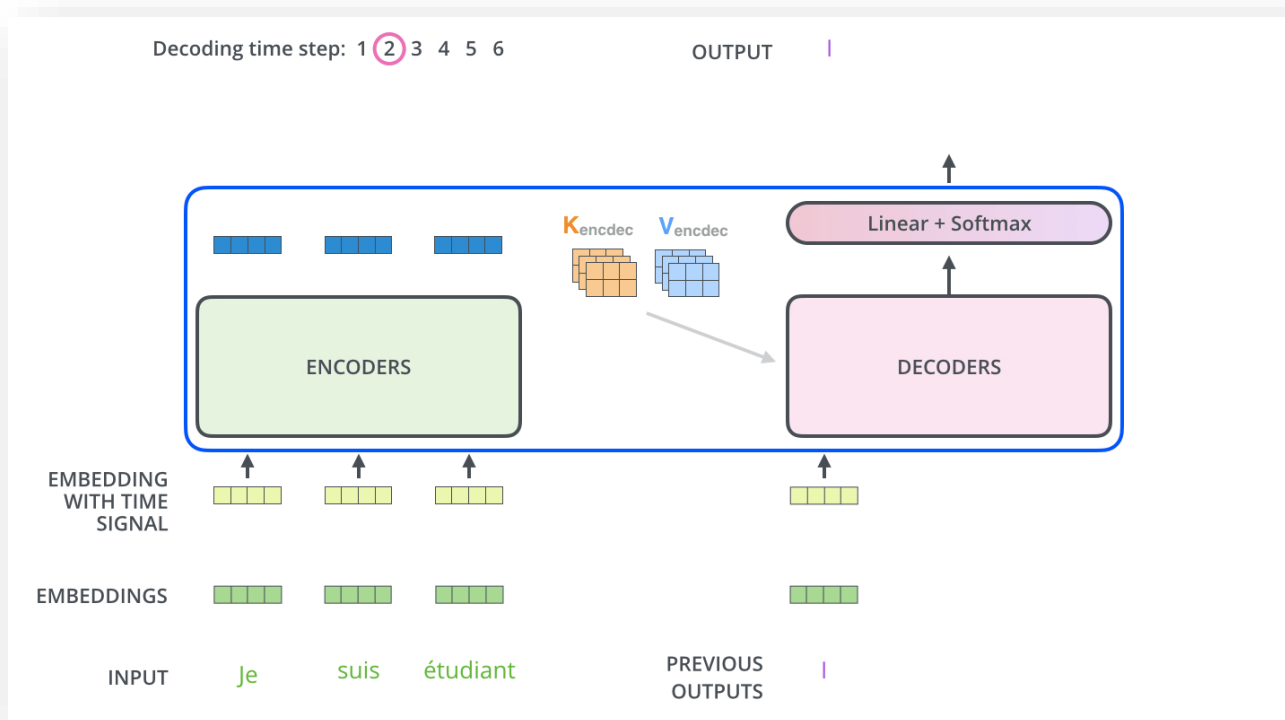
The Decoder in Transformer

- Encoder-decoder attention
 - Q from self-attn in decoder, K & V from encoder outputs
- Masked multi-head attention
 - Design similar to that of encoder, except for decoder #1 which takes additional inputs (of GT/predicted word embeddings).
 - Mask unpredicted tokens during softmax: why?



The Decoder in Transformer (cont'd)

- Encoder-decoder attention
 - Q from self-attn in decoder, K & V from encoder outputs
- Masked multi-head attention
 - Design similar to that of encoder, except for decoder #1 which takes additional inputs (of GT/predicted word embeddings).
 - Mask unpredicted tokens during softmax: why?



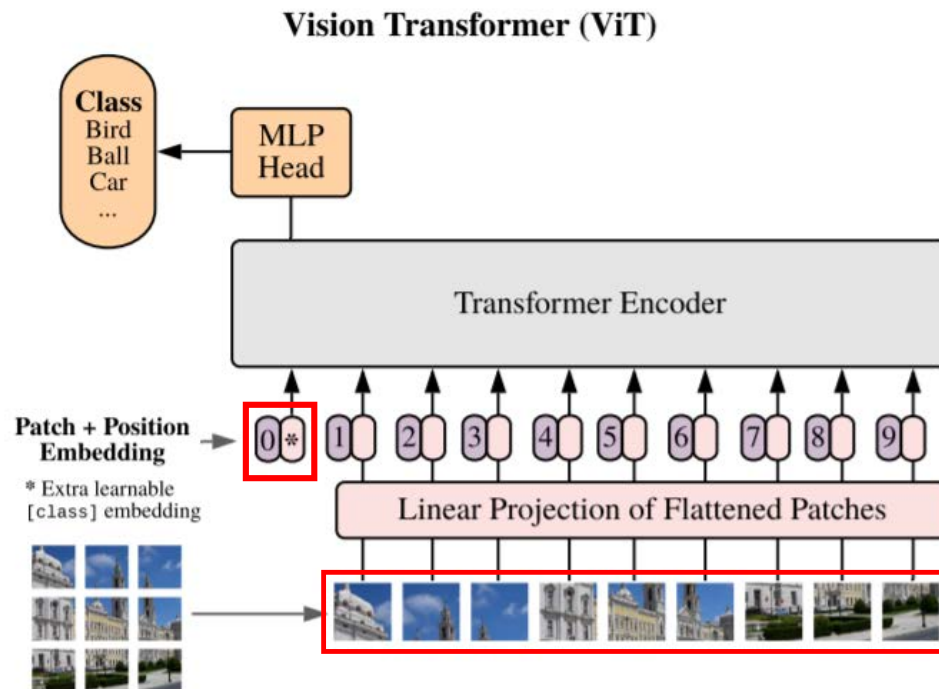
Overview of Decoding in Transformer

- Encoder/Decoder Cross-Attention + Decoder self-attention



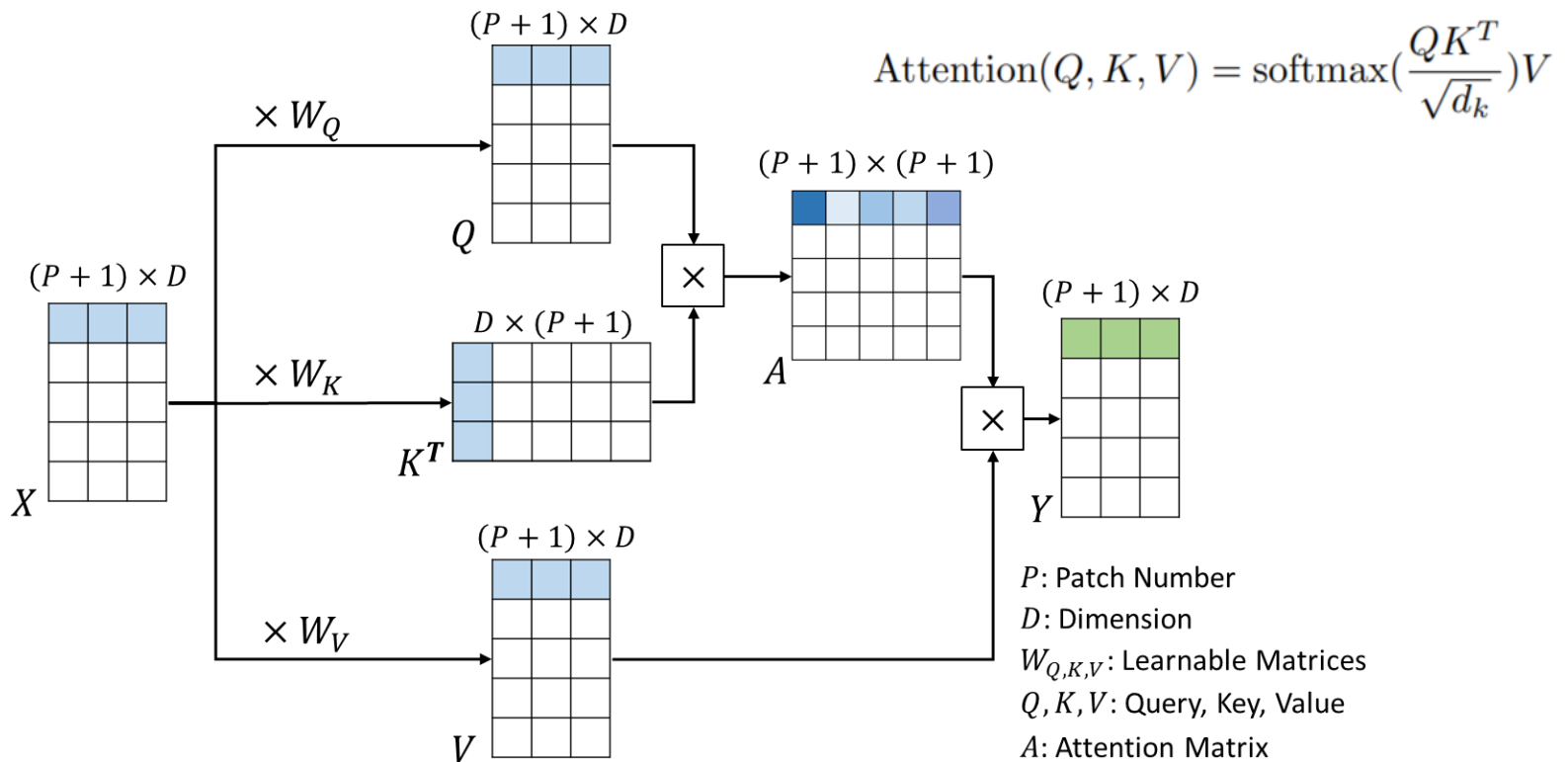
Vision Transformer

- “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”, ICLR, 2021. (Google Research)
- Partition the input image into a **patch sequence**
- An additional **token** (*) is appended to perform attention on patches
- Both the “*” token and positional embeddings (denoted by 0, 1, 2 ...) are **trainable vectors**.



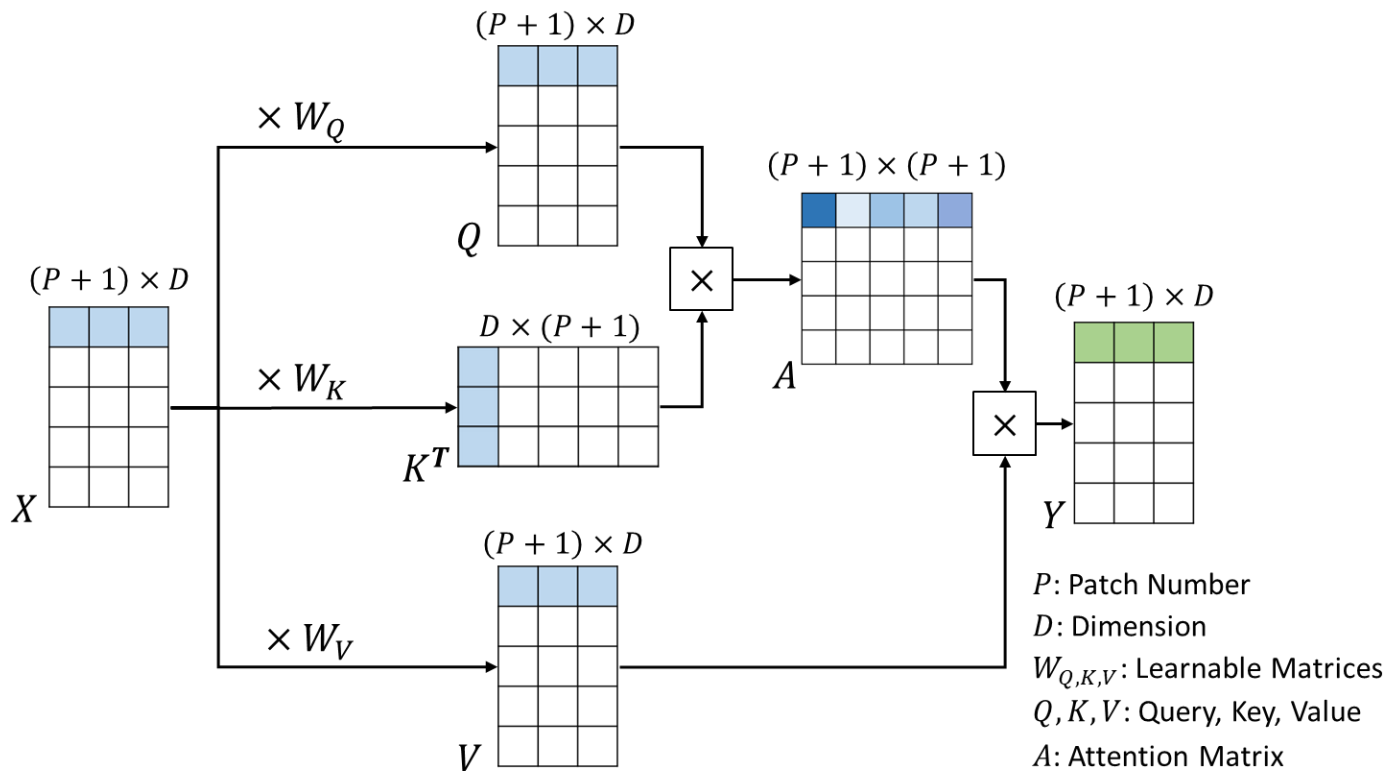
Query-Key-Value Attention in ViT

- E.g., An input image is partitioned into 4 patches, with feature dimension = 3 (i.e., $P=4$ and $D=3$).
- Note that there are $(P+1)$ rows since we have an additional token of *.



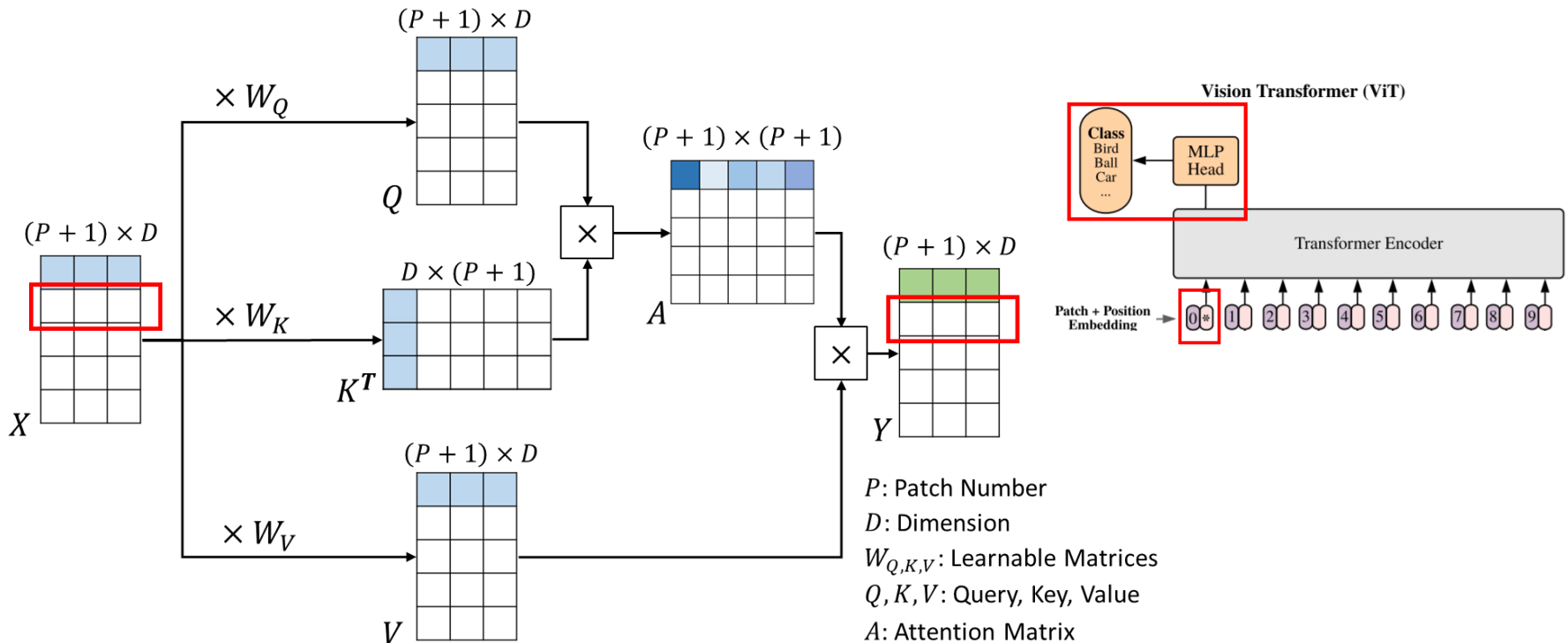
Query-Key-Value Attention in ViT (cont'd)

- By performing attention, the input sequence X (of length $P+1$) is “transformed” into another sequence Y with the same length
- Again, that’s why it is called “**Transformer**” and as a **seq2seq** model.



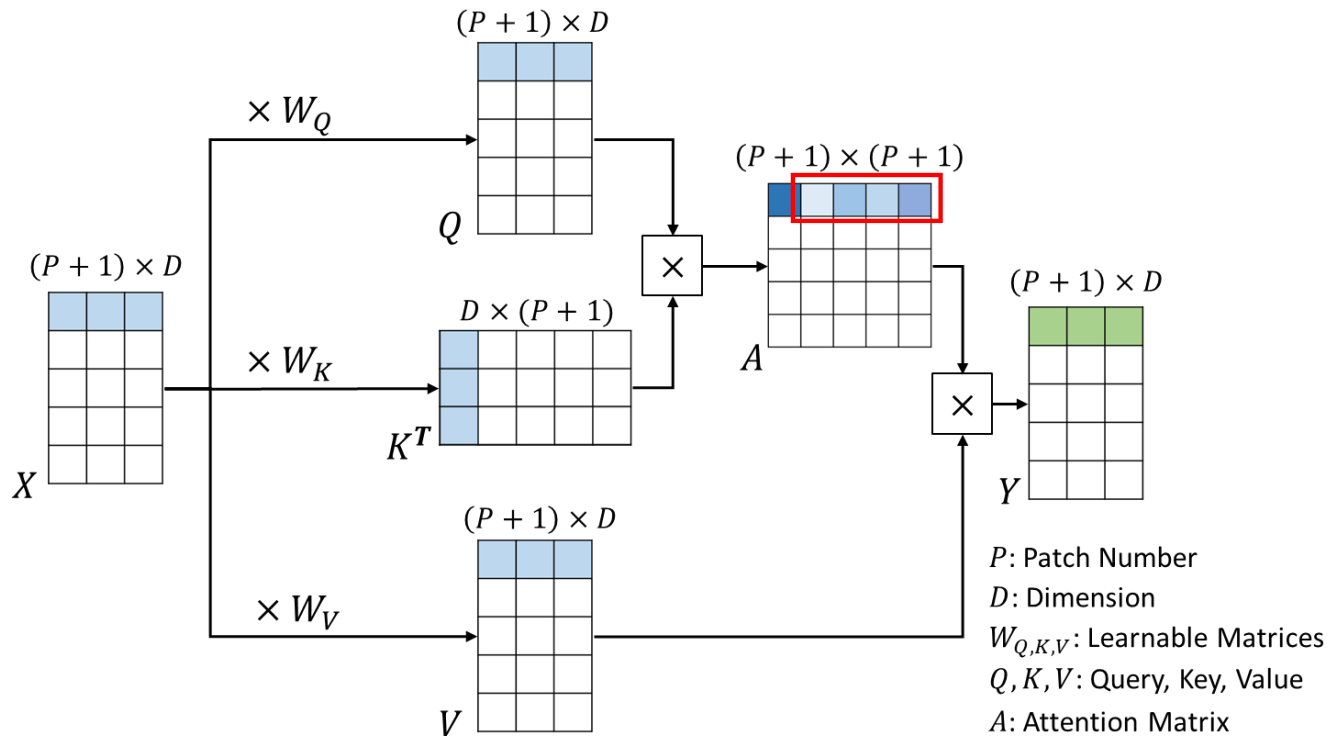
Query-Key-Value Attention in ViT (cont'd)

- In standard vision transformer, we only take the **first output token** of the output sequence (the **first row** of Y) for classification purposes
- This corresponds to the output when **token "0"** serves as query

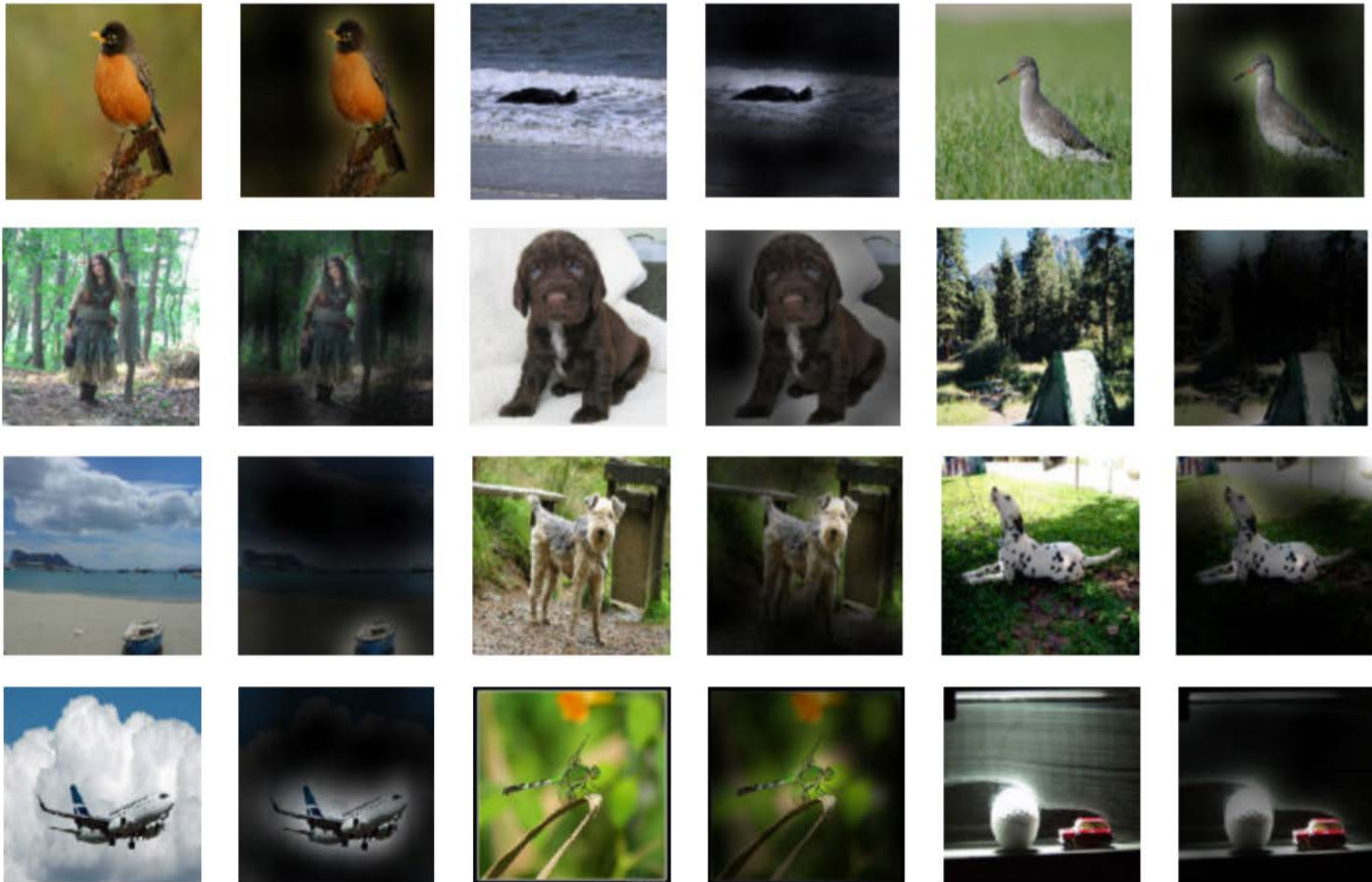


Visualization of ViT

- To visualize the attention maps, we take the attention scores from the **first row** of A (when token “0” serves as query)
- Note the first element is excluded, and thus there are **P scores** corresponding to the P image patches



Example Visualization for Image Classification

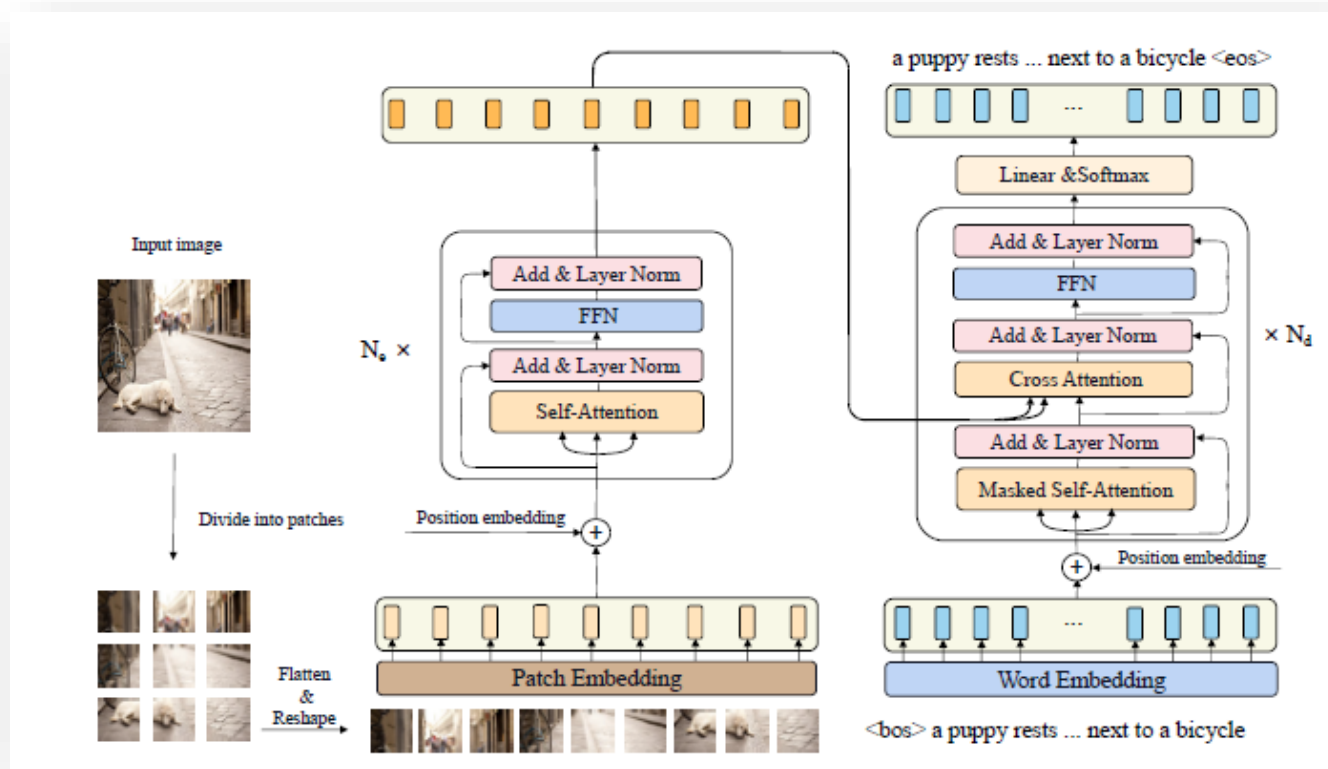


What'd Be Covered in This Crash Course...

- **Learning-based Computer Vision**
 - From Linear to Non-Linear Classifiers
- **Start of Deep Learning for Computer Vision**
 - Convolutional Neural Networks
 - Self-Supervised Learning
 - Segmentation & Detection
- **Generative Models**
 - Autoencoder, Variational Autoencoder, Generative Adversarial Networks & Diffusion Models
- **Sequence-to-Sequence Learning**
 - Attention is All You Need: Transformer
- **Vision & Language Foundation Models**
 - Image-to-Text vs. Text-to-Image
 - Parameter-Efficient Fine-tuning

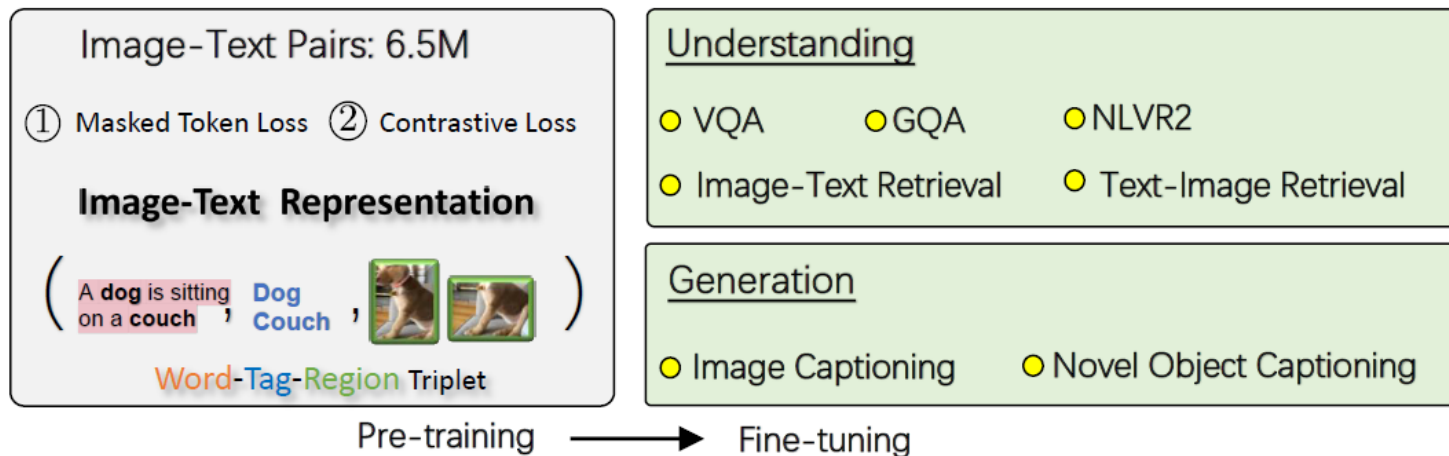
Image Captioning

- **Caption Transformer (CPTR)** -
CPTR: Full Transformer Network for Image Captioning, arxiv 2021
- Motivation: patch translation for image captioning



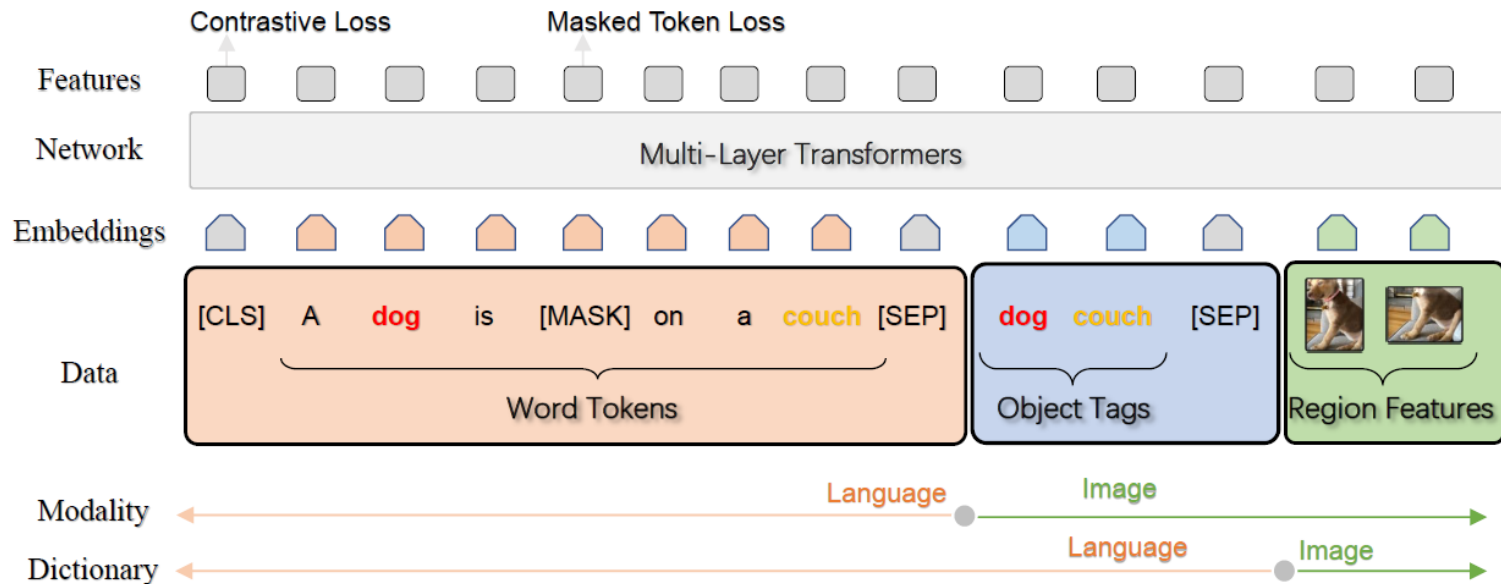
Beyond Image Captioning: Unified Vision & Language Model

- **Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks (ECCV'20)**
 - Training data: triplets of **caption-tag-region**
 - Objectives:
 1. Masked token loss for **words** & **tags**
 2. Contrastive loss **tags** and others
 - Fine-tuning: 5 vision & language tasks (VQA, image-text retrieval, image captioning, NOC, etc.)



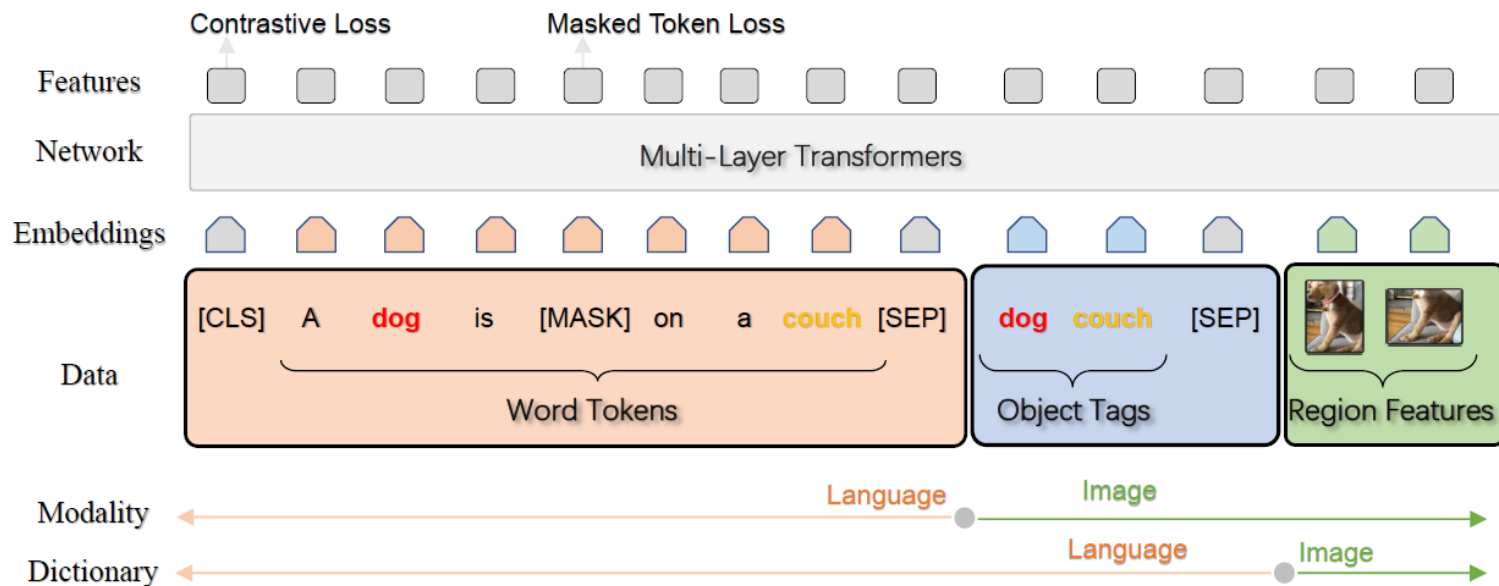
Semantics-Aligned Pre-training for V+L Tasks

- **Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks (ECCV'20)**
 - Training:
 - Inputs: triplets of **caption-tag-region**
 - Objectives: Masked token loss for **words** & **tags** + Contrastive loss **tags** and others
 - Fine-tuning:
5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)

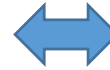


Semantics-Aligned Pre-training for V+L Tasks (cont'd)

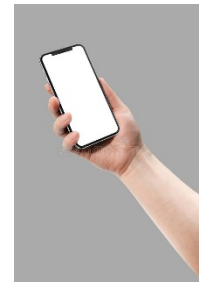
- **Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks (ECCV'20)**
 - Training:
 - Inputs: triplets of word-tag-region
 - Objectives: Masked token loss for words & tags + Contrastive loss tags and others
 - Fine-tuning:
 - 5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)



Holding an apple



or



- **Oscar (cont'd)**

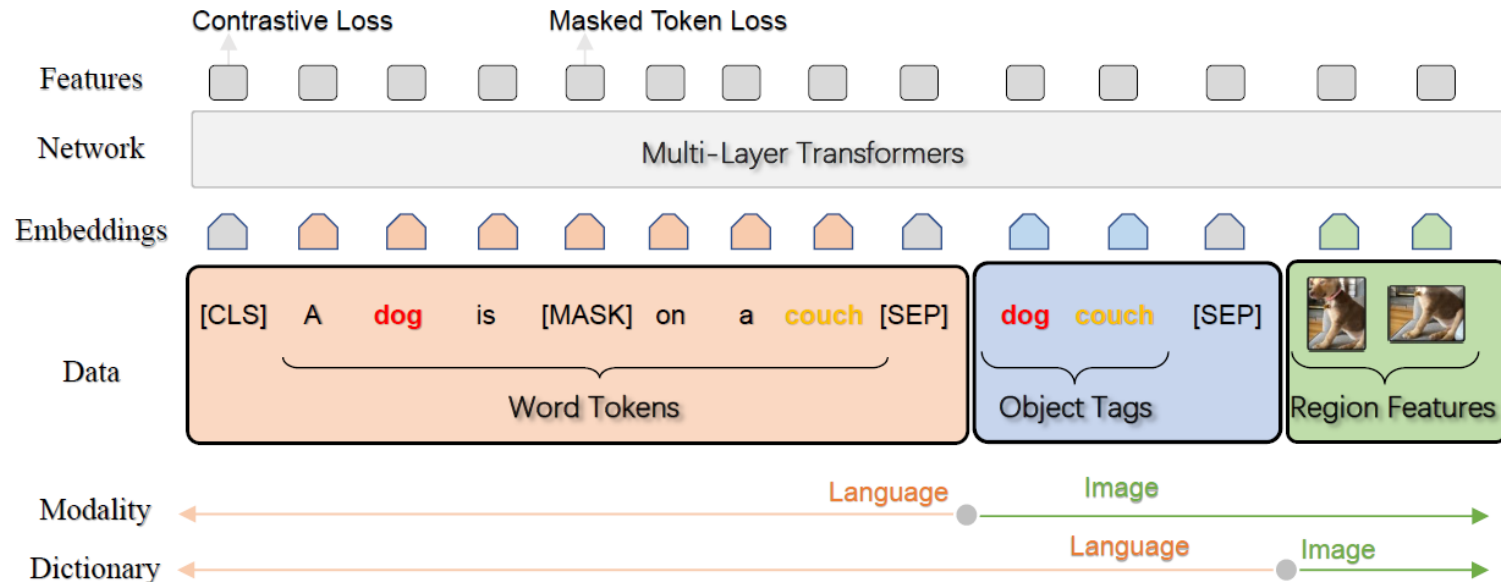
- Fine-tuning:

- 5 vision & language tasks (image captioning, NOC, VQA, image-text retrieval, etc.)

- Take **image-text retrieval** as an example

- Training: aligned/mis-aligned **image-text** pairs as positive/negative input pairs, with **[CLS]** for binary classification (1/0)

- Inference: for either image or text retrieval, calculate classification score of **[CLS]** for the input query



CLIP: Contrastive Language-Image Pretraining

- OpenAI, *Learning Transferable Visual Models From Natural Language Supervision*, NeurIPS WS 2021 (w/ 9000+ citations)
- Why DL/CNN not good enough?
 - Require annotated data for training image classification
 - Domain gap between closed-world and open-world domain data
 - Lack of ability for zero-shot classification



let



geNet V2



ImageNet Rendition



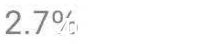
SubjectNet



geNet Sketch



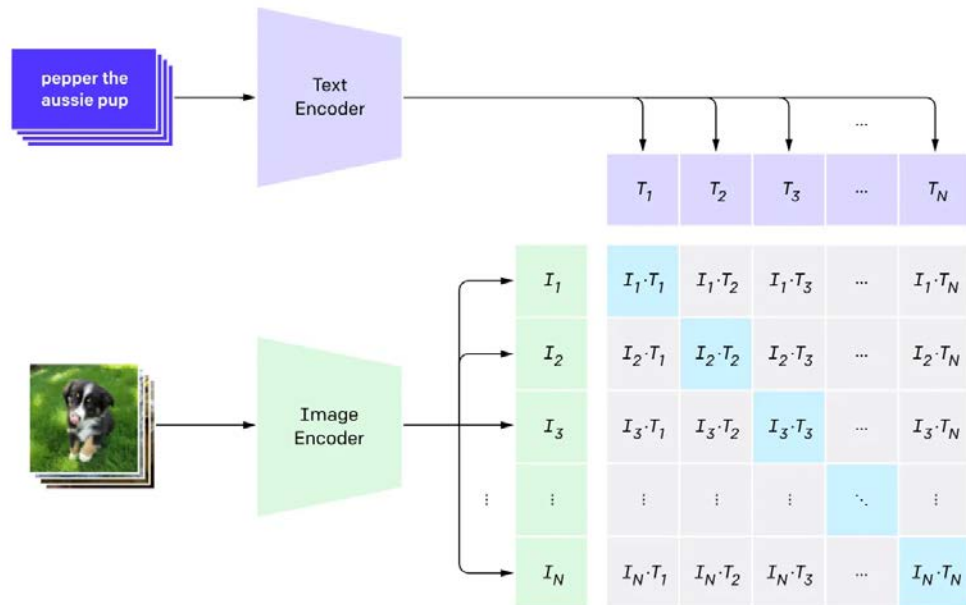
Not Adversarial



CLIP (cont'd)

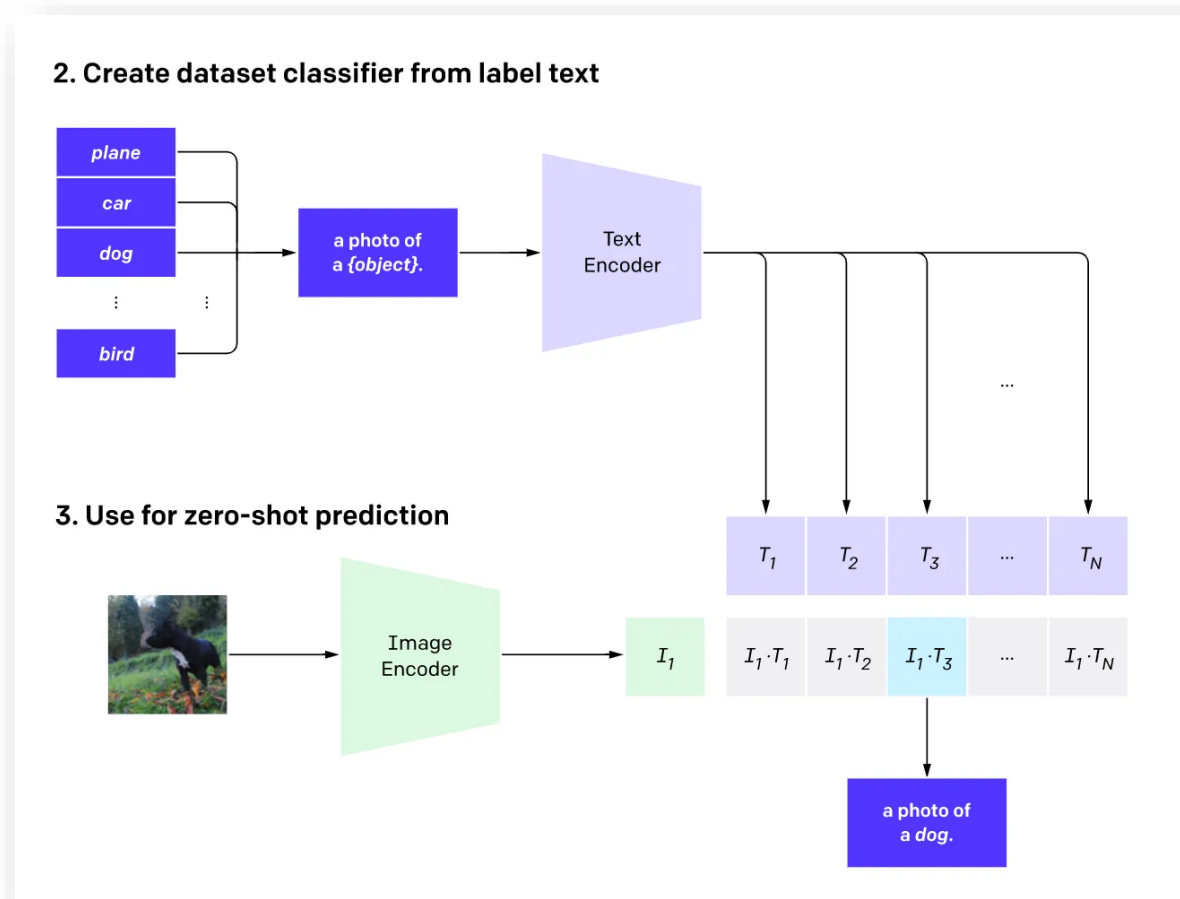
- Why DL/CNN not good enough?
 - Require annotated data for training image classification
 - Domain gap between closed-world and open-world domain data
 - Lack of ability for zero-shot classification
- Motivation/Objectives
 - Cross-domain contrastive learning from large-scale image-language data

1. Contrastive pre-training



CLIP (cont'd)

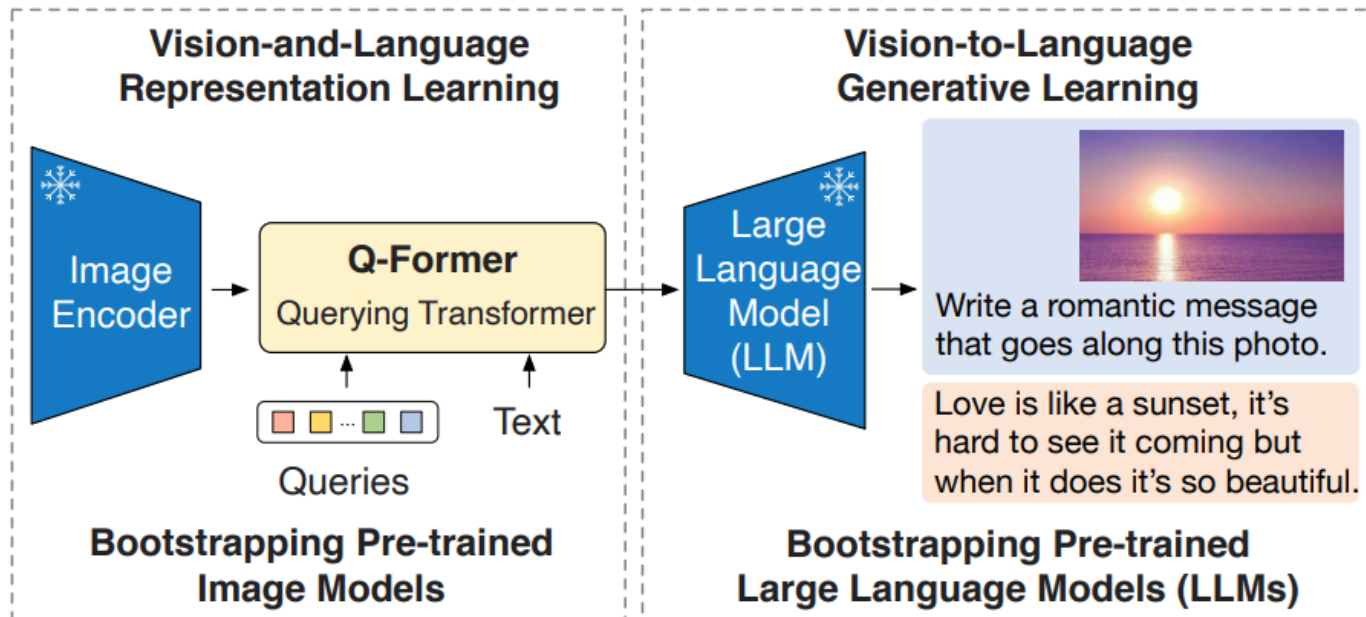
- (Zero-shot) Inference:



- Potential concerns/disadvantages?

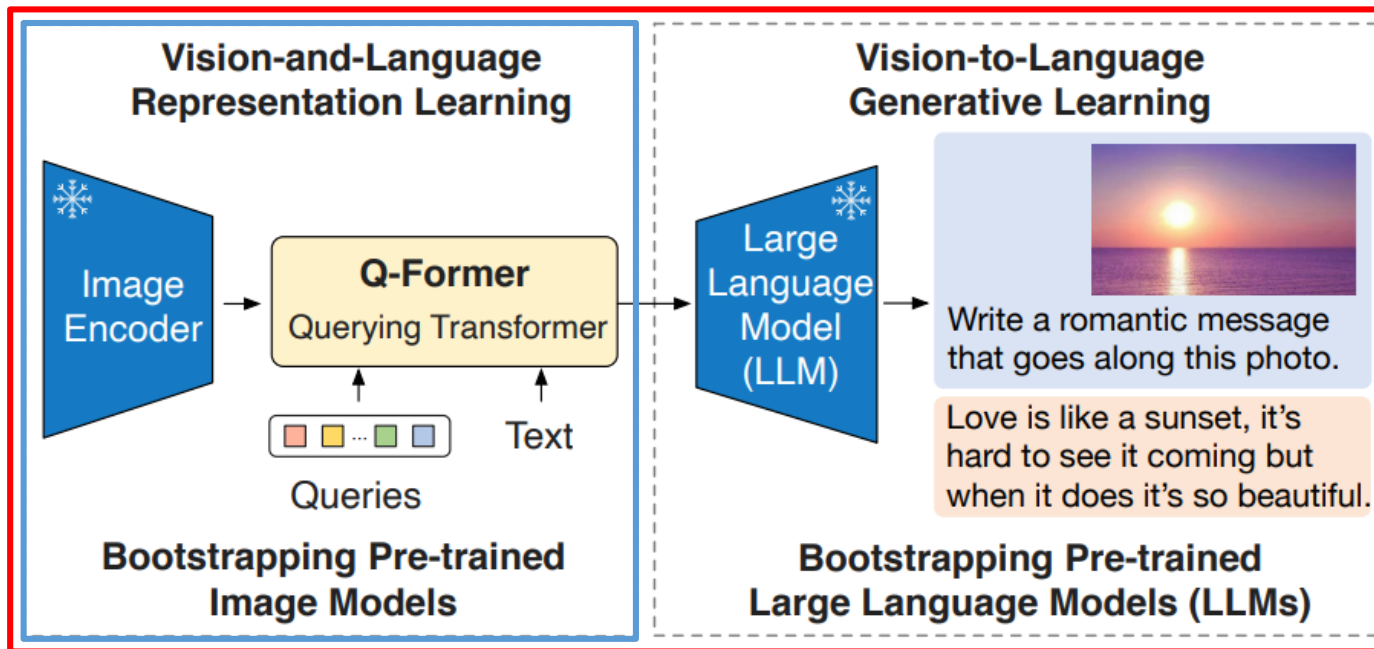
BLIP-2 (ICML'23)

- **BLIP:**
Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation, NeurIPS 2021
- **Goal:**
Bridge the modality gap between off-the-shelf [frozen pre-trained image encoders](#) and [frozen large language models](#) with a lightweight [Querying Transformer \(Q-Former\)](#).
- **Result:**
 1. SOTA performance on various downstream vision-language tasks.
 2. Zero-shot image-to-text generation that can follow natural language instructions.



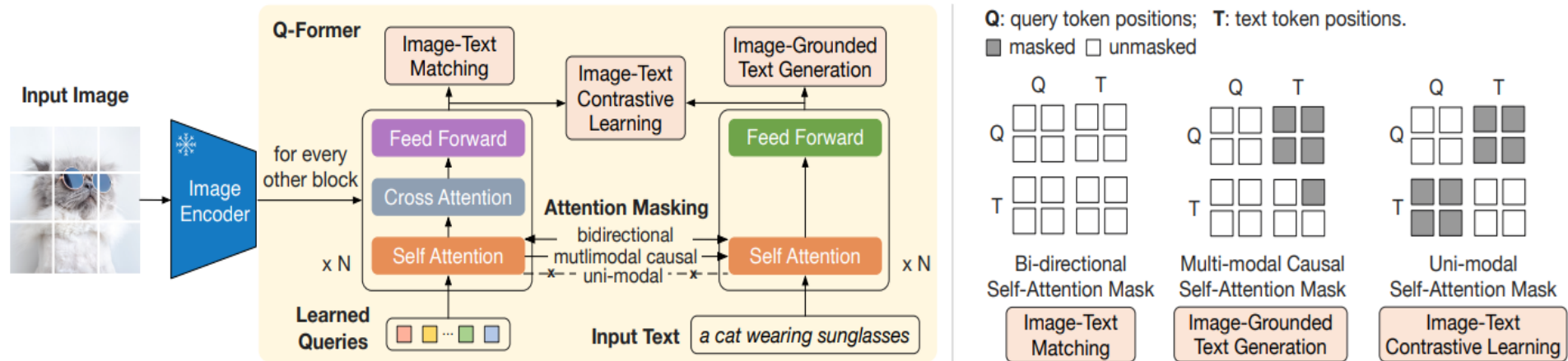
Pre-training

- **Two-stage** Pre-training
 - **Stage 1:**
VL representation learning which enforces the **Q-Former** to learn **visual representation** that is most relevant to the text.
 - **Stage 2:**
VL generative learning makes the output representation of **Q-Former** to be understood by **LLMs**.



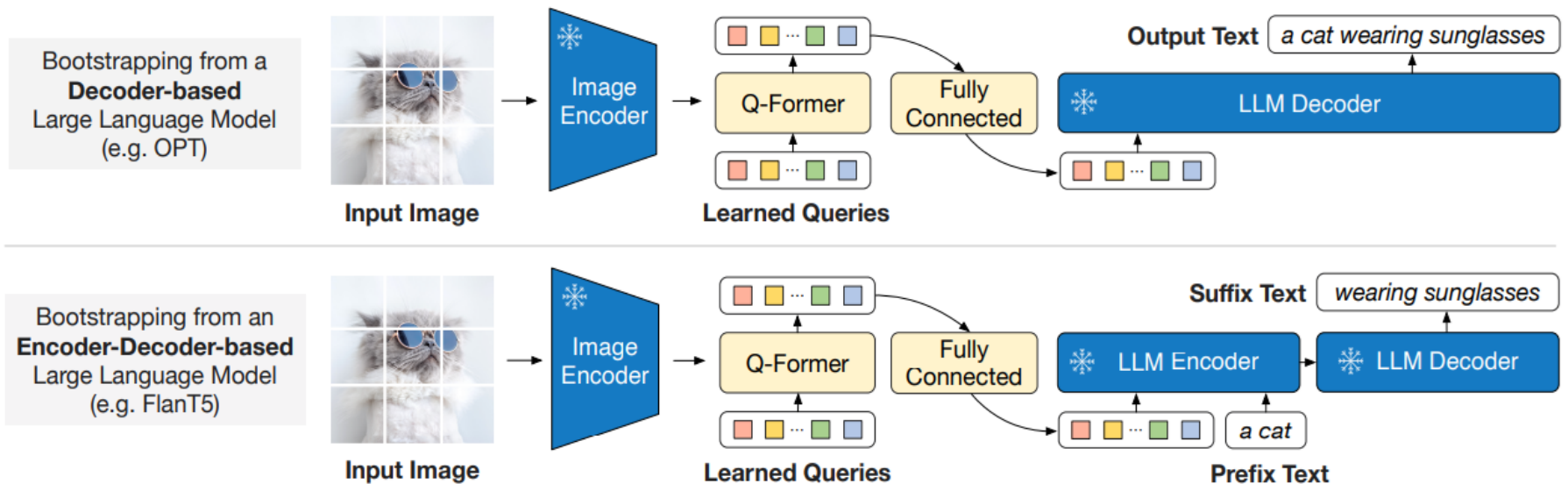
Pre-training Stage 1 - VL Representation Learning

- **Goal:**
enforce the **Q-Former** to extract visual representation relevant to the text.
- **Method:** three pre-training tasks
 - **Image-Text Matching (ITM):**
for each learnable query -> linear classifier for binary decision
 - **Image-grounded Text Generation (ITG):**
self-attn in Q for encoder training; T->Q for image-to-text generation
 - **Image-Text Contrastive Learning (ITC):**
self-attn in Q/T, followed by $\max(\text{sim}(Q, T))$



Pre-training Stage 2 - VL Generative Learning

- **Goal:**
Learning with LLM guidance
i.e., make the output representation of **Q-Former** to be understood by **LLMs**.
- **Method:**
pre-training with Image-grounded Text Generation (ITG)



Parameter-Efficient Fine-Tuning

- **Adapter**

- VL-ADAPTER: Parameter-Efficient Transfer Learning for Vision-and-Language Tasks (CVPR, 2022)

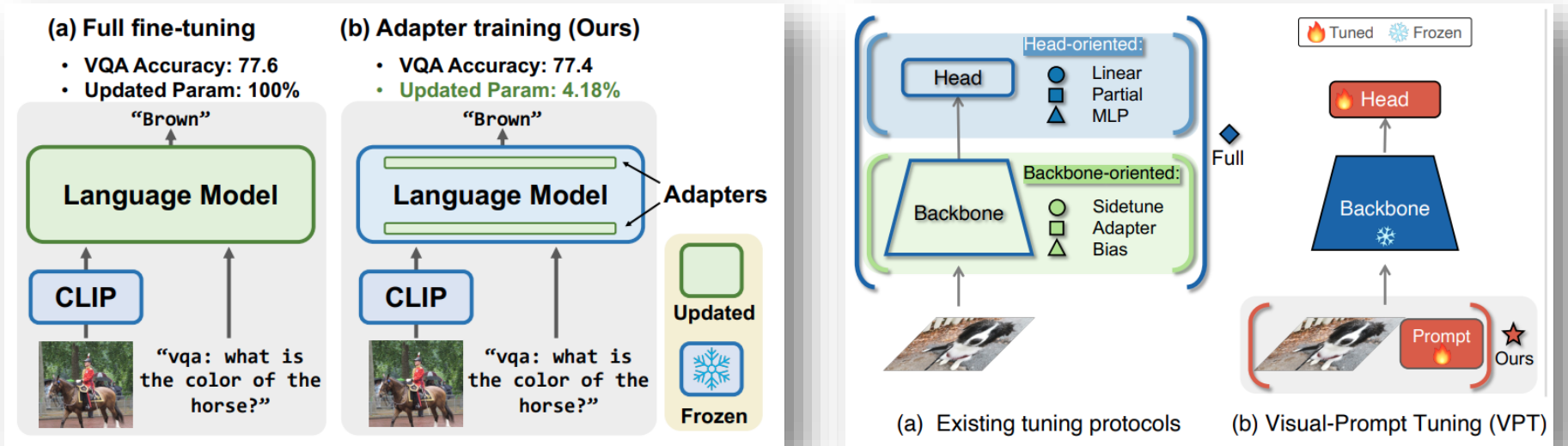
- **Visual Prompt Tuning**

- Visual Prompt Tuning (ECCV, 2022)

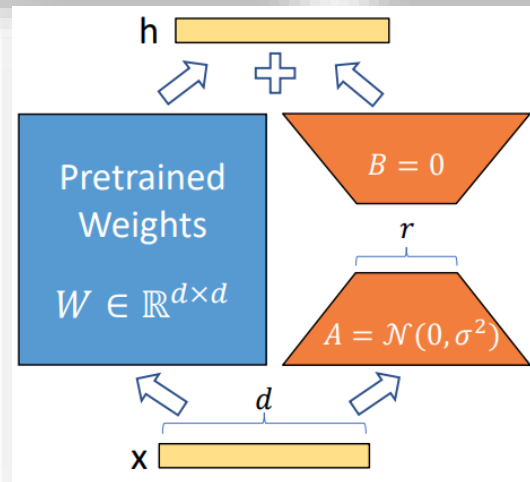
- **LoRA**

- LoRA: Low-Rank Adaptation of Large Language Models (ICLR, 2022)

Parameter Efficient Fine Tuning



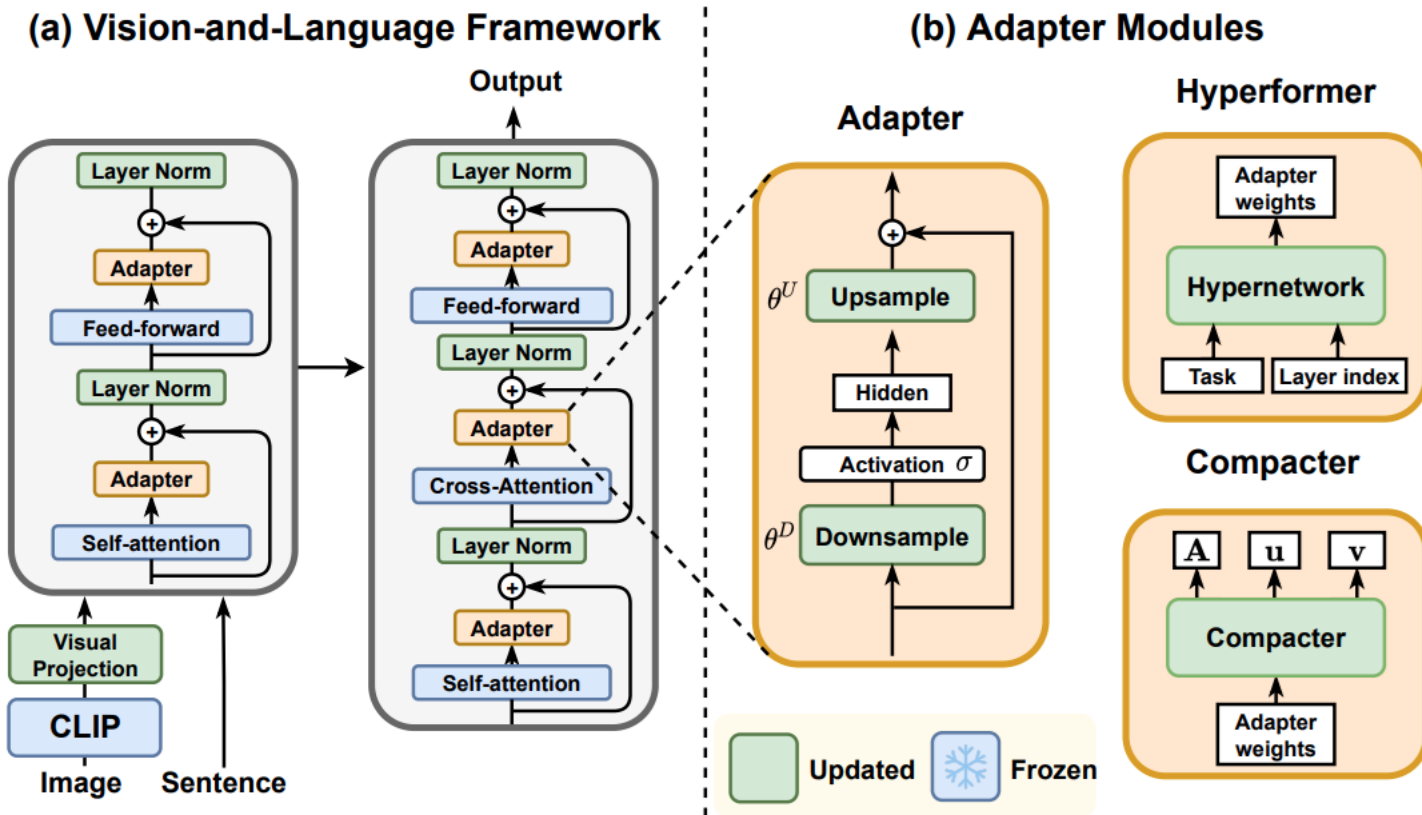
Adapter



LoRA

Prompt Tuning

VL-ADAPTER: Parameter-Efficient Transfer Learning for Vision-and-Language Tasks

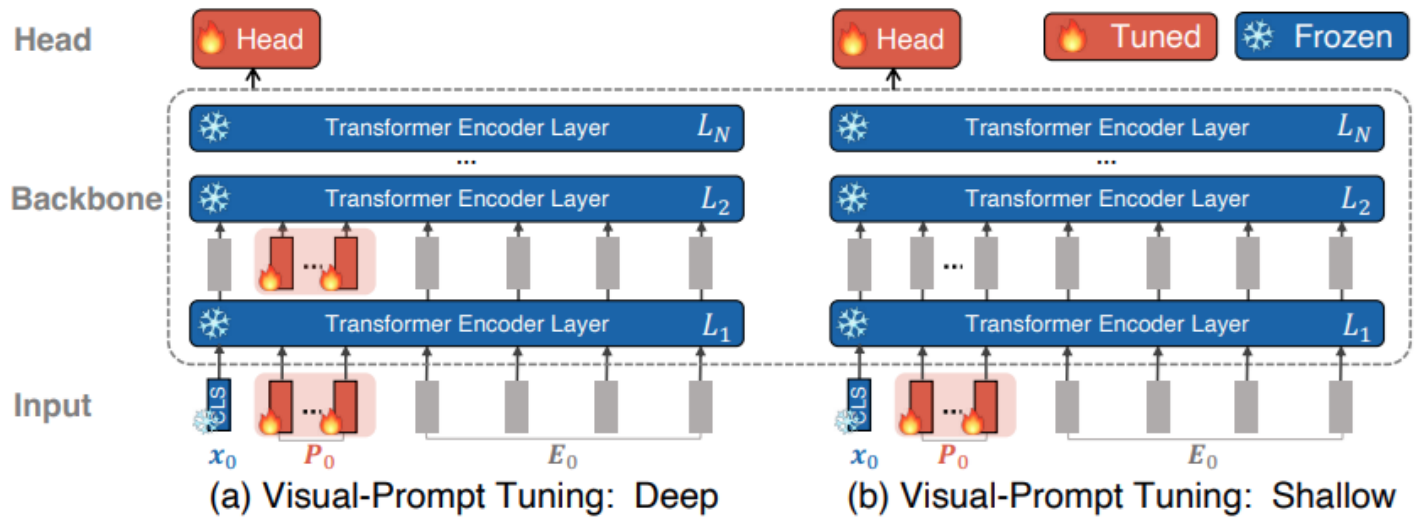
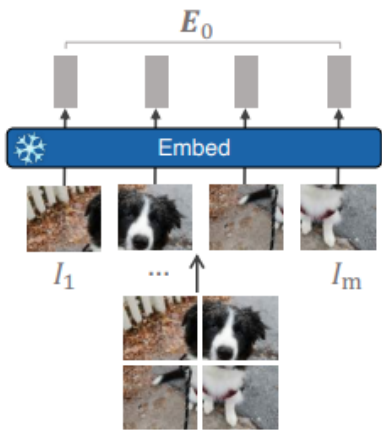


<https://arxiv.org/abs/2112.06825>

Visual Prompt Tuning

- Shallow: $[\mathbf{x}_1, \mathbf{Z}_1, \mathbf{E}_1] = L_1([\mathbf{x}_0, \mathbf{P}, \mathbf{E}_0])$ (4)
- $[\mathbf{x}_i, \mathbf{Z}_i, \mathbf{E}_i] = L_i([\mathbf{x}_{i-1}, \mathbf{Z}_{i-1}, \mathbf{E}_{i-1}])$ $i = 2, 3, \dots, N$ (5)
- $\mathbf{y} = \text{Head}(\mathbf{x}_N)$, (6)

- Deep: $[\mathbf{x}_i, _, \mathbf{E}_i] = L_i([\mathbf{x}_{i-1}, \mathbf{P}_{i-1}, \mathbf{E}_{i-1}])$ $i = 1, 2, \dots, N$ (7)
- $\mathbf{y} = \text{Head}(\mathbf{x}_N)$. (8)



LoRA: Low-Rank Adaptation of Large Language Models

- Previous problems
 - Adapter Layers Introduce Inference Latency
 - Directly Optimizing the Prompt is Hard

- LoRA

$$W_0 \in \mathbb{R}^{d \times k}$$

$$W_0 + \Delta W = W_0 + BA$$

$$B \in \mathbb{R}^{d \times r}, A \in \mathbb{R}^{r \times k}$$

$$\text{rank } r \ll \min(d, k)$$

$$h = W_0x + \Delta Wx = W_0x + BAx$$

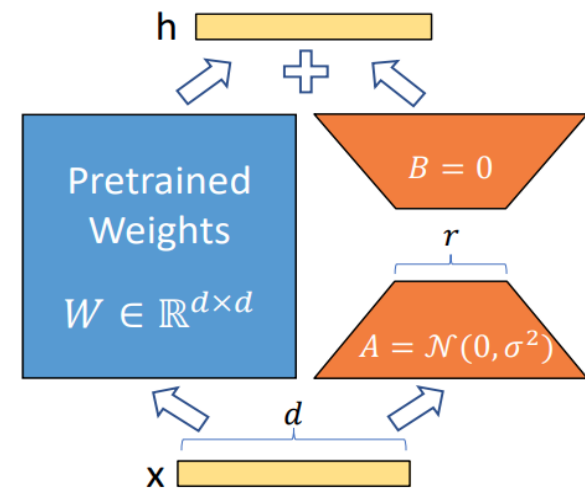


Figure 1: Our reparametrization. We only train A and B .

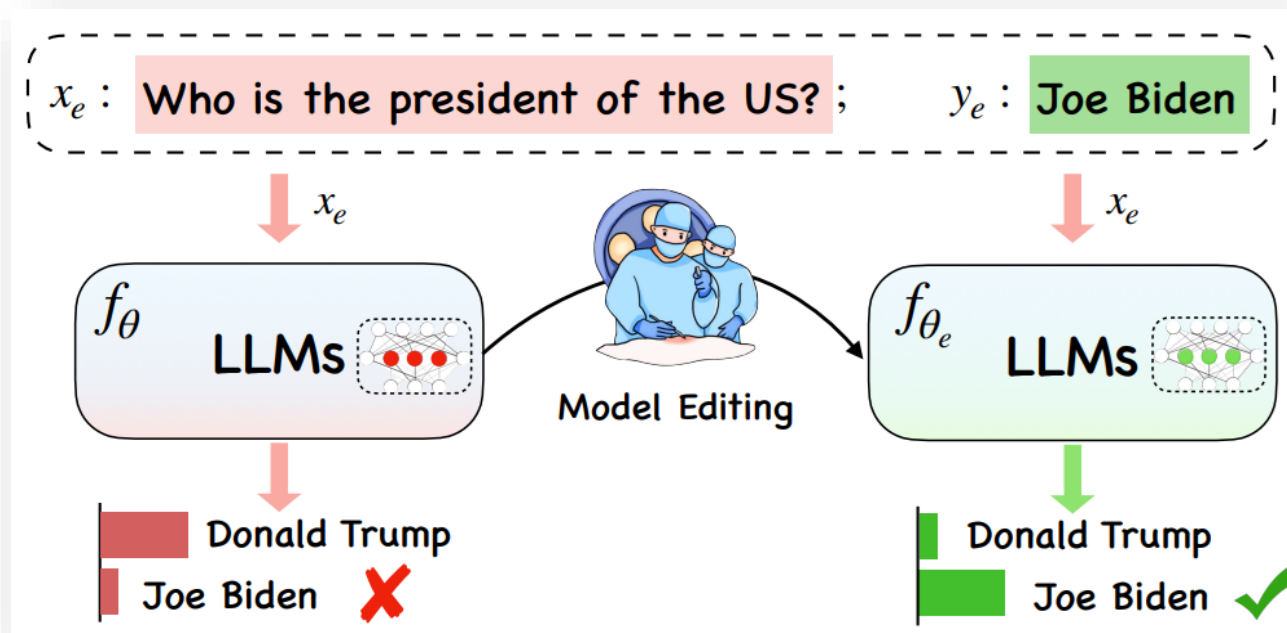
LoRA: Low-Rank Adaptation of LLMs (cont'd)

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
		RoB _{base} (FT)*	125.0M	87.6	94.8	90.2	63.6	92.8	91.9	78.7
RoB _{base} (BitFit)*	0.1M	84.7	93.7	92.7	62.0	91.8	84.0	81.5	90.8	85.2
RoB _{base} (Adpt ^D)*	0.3M	87.1 \pm 0	94.2 \pm 1	88.5 \pm 1.1	60.8 \pm 4	93.1 \pm 1	90.2 \pm 0	71.5 \pm 2.7	89.7 \pm 3	84.4
RoB _{base} (Adpt ^D)*	0.9M	87.3 \pm 1	94.7 \pm 3	88.4 \pm 1	62.6 \pm 9	93.0 \pm 2	90.6 \pm 0	75.9 \pm 2.2	90.3 \pm 1	85.4
RoB _{base} (LoRA)	0.3M	87.5 \pm 3	95.1\pm2	89.7 \pm 7	63.4 \pm 1.2	93.3\pm3	90.8 \pm 1	86.6\pm7	91.5\pm2	87.2
RoB _{large} (FT)*	355.0M	90.2	96.4	90.9	68.0	94.7	92.2	86.6	92.4	88.9
RoB _{large} (LoRA)	0.8M	90.6\pm2	96.2 \pm 5	90.9\pm1.2	68.2\pm1.9	94.9\pm3	91.6 \pm 1	87.4\pm2.5	92.6\pm2	89.0
RoB _{large} (Adpt ^P)†	3.0M	90.2 \pm 3	96.1 \pm 3	90.2 \pm 7	68.3\pm1.0	94.8\pm2	91.9\pm1	83.8 \pm 2.9	92.1 \pm 7	88.4
RoB _{large} (Adpt ^P)†	0.8M	90.5\pm3	96.6\pm2	89.7 \pm 1.2	67.8 \pm 2.5	94.8\pm3	91.7 \pm 2	80.1 \pm 2.9	91.9 \pm 4	87.9
RoB _{large} (Adpt ^H)†	6.0M	89.9 \pm 5	96.2 \pm 3	88.7 \pm 2.9	66.5 \pm 4.4	94.7 \pm 2	92.1 \pm 1	83.4 \pm 1.1	91.0 \pm 1.7	87.8
RoB _{large} (Adpt ^H)†	0.8M	90.3 \pm 3	96.3 \pm 5	87.7 \pm 1.7	66.3 \pm 2.0	94.7 \pm 2	91.5 \pm 1	72.9 \pm 2.9	91.5 \pm 5	86.4
RoB _{large} (LoRA)†	0.8M	90.6\pm2	96.2 \pm 5	90.2\pm1.0	68.2 \pm 1.9	94.8\pm3	91.6 \pm 2	85.2\pm1.1	92.3\pm5	88.6
DeB _{XXL} (FT)*	1500.0M	91.8	97.2	92.0	72.0	96.0	92.7	93.9	92.9	91.1
DeB _{XXL} (LoRA)	4.7M	91.9\pm2	96.9 \pm 2	92.6\pm6	72.4\pm1.1	96.0\pm1	92.9\pm1	94.9\pm4	93.0\pm2	91.3

Model & Method	# Trainable Parameters	E2E NLG Challenge				
		BLEU	NIST	MET	ROUGE-L	CIDEr
GPT-2 M (FT)*	354.92M	68.2	8.62	46.2	71.0	2.47
GPT-2 M (Adapter ^L)*	0.37M	66.3	8.41	45.0	69.8	2.40
GPT-2 M (Adapter ^L)*	11.09M	68.9	8.71	46.1	71.3	2.47
GPT-2 M (Adapter ^H)	11.09M	67.3 \pm 6	8.50 \pm 0.7	46.0 \pm 2	70.7 \pm 2	2.44 \pm 0.1
GPT-2 M (FT ^{Top2})*	25.19M	68.1	8.59	46.0	70.8	2.41
GPT-2 M (PreLayer)*	0.35M	69.7	8.81	46.1	71.4	2.49
GPT-2 M (LoRA)	0.35M	70.4\pm1	8.85\pm0.2	46.8\pm2	71.8\pm1	2.53\pm0.2
GPT-2 L (FT)*	774.03M	68.5	8.78	46.0	69.9	2.45
GPT-2 L (Adapter ^L)	0.88M	69.1 \pm 1	8.68 \pm 0.3	46.3 \pm 0	71.4 \pm 2	2.49\pm0
GPT-2 L (Adapter ^L)	23.00M	68.9 \pm 3	8.70 \pm 0.4	46.1 \pm 1	71.3 \pm 2	2.45 \pm 0.2
GPT-2 L (PreLayer)*	0.77M	70.3	8.85	46.2	71.7	2.47
GPT-2 L (LoRA)	0.77M	70.4\pm1	8.89\pm0.2	46.8\pm2	72.0\pm2	2.47 \pm 0.2

Adv. Topic: Knowledge Editing

- **Motivation:**
Knowledge updates everyday while retraining LLMs is expensive.
- **Goal:**
Propose more efficient & effective solutions to update knowledge in LLM



Adv. Topic: Unlearning



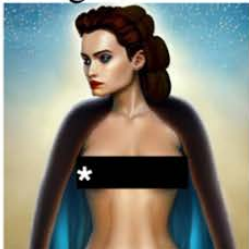
↓ A photo of *Elon Musk* ↓



- **Goal:**
Erase the undesirable visual concepts from Diffusion Models.
Concepts can be abstractive concept, artistic style, object, or personality.
- **Method:**
Diffusion Model **Fine-tuning**.

Erasing Nudity

Original Model



Edited Model



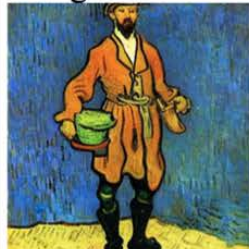
* Added by authors for publication



Erased from model:
“Nudity”

Erasing Artistic Style

Original Model



Edited Model



Erased from model:
“Van Gogh”

Erasing Objects

Original Model



Edited Model



Erased from model:
“Car”

What We (Try to) Cover Today...

- **Learning-based Computer Vision**
 - From Linear to Non-Linear Classifiers
- **Start of Deep Learning for Computer Vision**
 - Convolutional Neural Networks
 - SSL, Segmentation & Detection
- **Generative Models**
 - AE, VAE, GAN, & Diffusion Models
- **Sequence-to-Sequence Learning**
 - Attention is All You Need: Transformer
- **Vision & Language Foundation Models**
 - Image-to-Text vs. Text-to-Image
 - Parameter-Efficient Fine-tuning
- **Lots of research topics we haven't covered...**
 - 3D vision
 - Video-based synthesis & analysis, etc.



Feel free to reach me at ycwang@ntu.edu.tw