# CPSC 425: Computer Vision

**Lecture 23:** Detection, Segmentation

# **Menu** for Today

## **Redings:**

— **Today's** Lecture:  N/A

— **Next** Lecture:      N/A

## **Reminders:**

— **Assignment 6**: Deep Learning is out and due **Thursday**

— **Material** for **Final Prep** will be on Canvas **tonight**

— **Quiz 6** is due **Thursday**

# Computer **Vision Problems**

Categorization



Multi-**class:** Horse
Church
Toothbrush
**Person**

IM**A**GENET

Multi-**label**: **Horse**
Church
Toothbrush
**Person**

# Computer **Vision Problems**

Categorization



Detection



Multi-**class:**   Horse
              Church
              Toothbrush
              **Person**

Horse (x, y, w, h)
Horse (x, y, w, h)
Person (x, y, w, h)
Person (x, y, w, h)

IMAGENET

COCO
Common Objects in Context

Multi-**label:**   **Horse**
              Church
              Toothbrush
              **Person**

# Computer **Vision Problems**

## Categorization
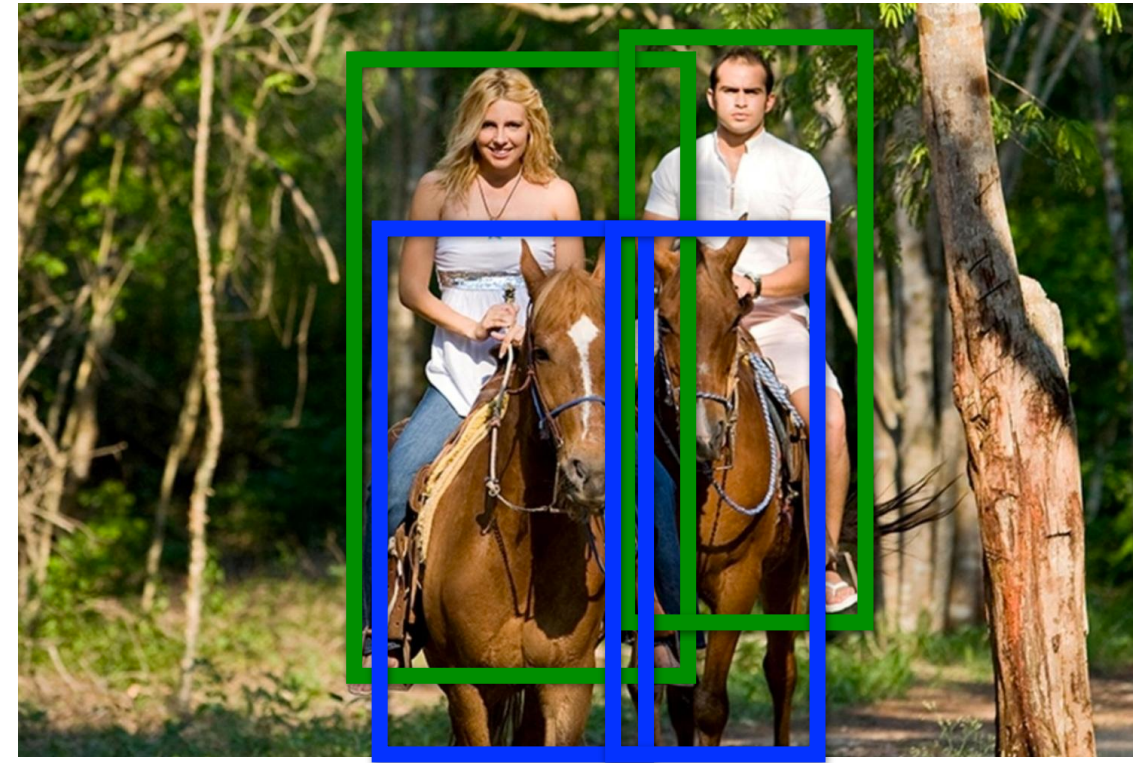


Multi-**class:**　Horse
　　　　　　　Church
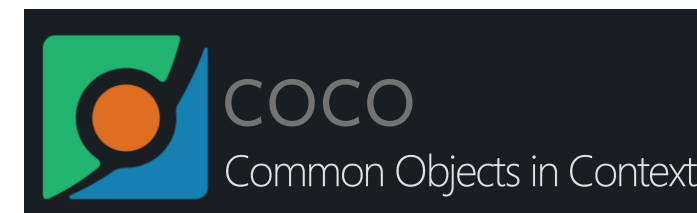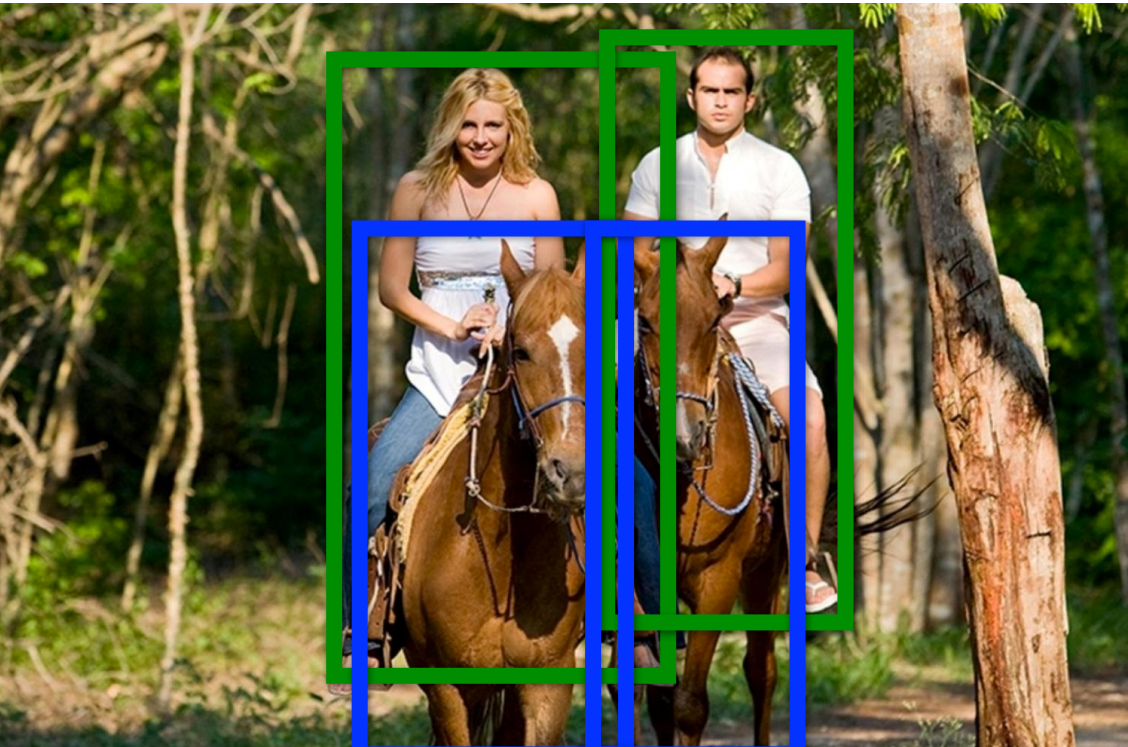　　　　　　　Toothbrush
　　　　　　　**Person**

IM**A**GENET

Multi-**label**:　**Horse**
　　　　　　　Church
　　　　　　　Toothbrush
　　　　　　　**Person**

## Detection



Horse (x, y, w, h)
Horse (x, y, w, h)
Person (x, y, w, h)
Person (x, y, w, h)

COCO
Common Objects in Context

## Segmentation



Horse
Person

COCO
Common Objects in Context

# Computer **Vision Problems**

## Categorization



Multi-**class:**    Horse
                    Church
                    Toothbrush
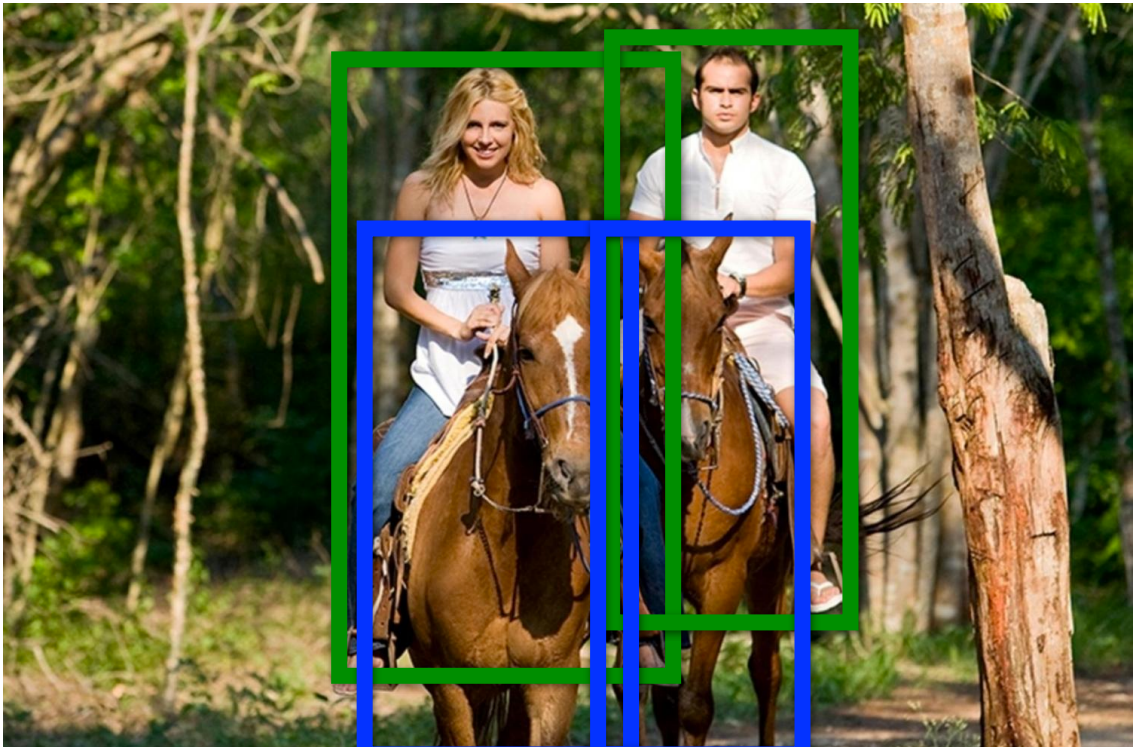                    **Person**

IM**A**GENET

Multi-**label**:    **Horse**
                    Church
                    Toothbrush
                    **Person**

## Detection



Horse (x, y, w, h)
Horse (x, y, w, h)
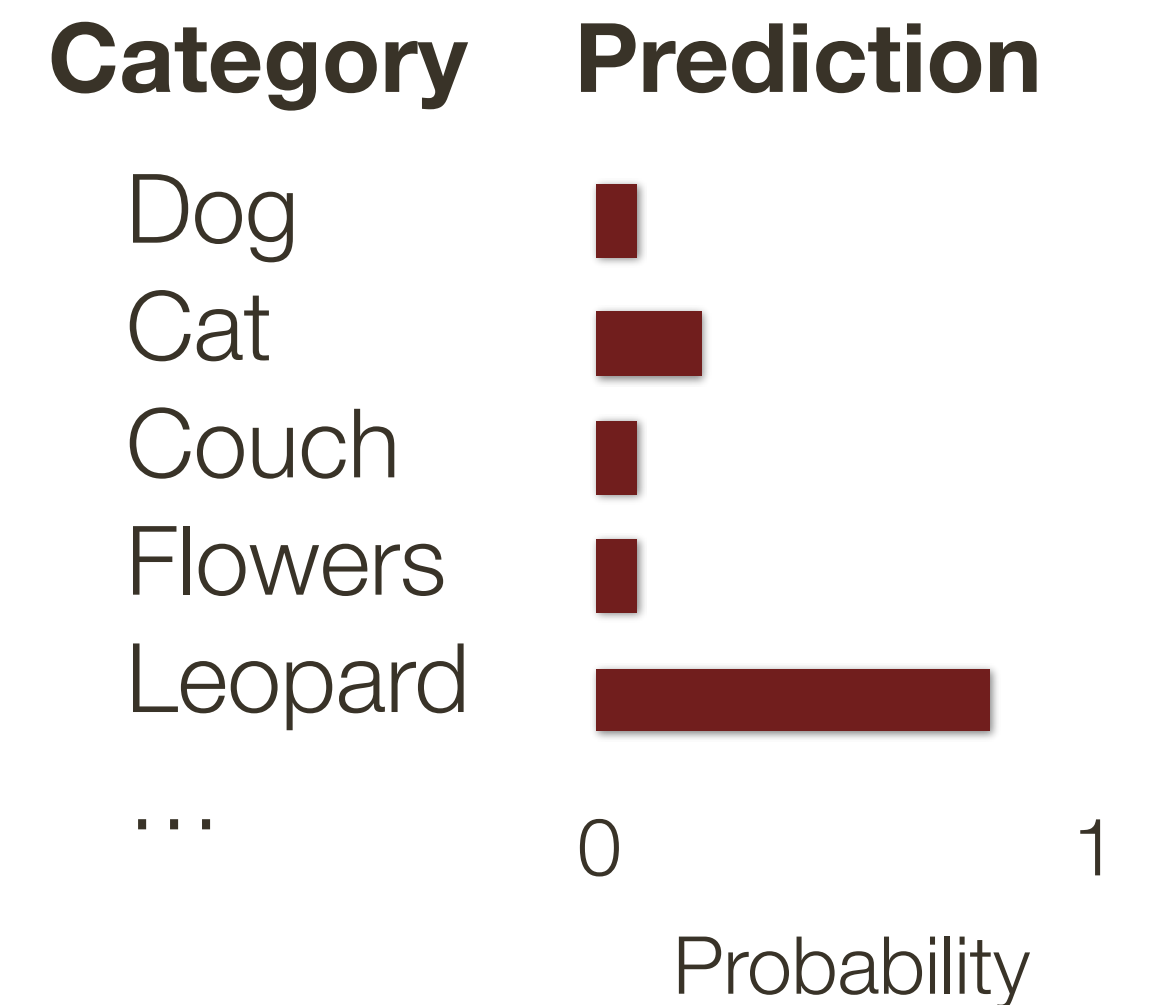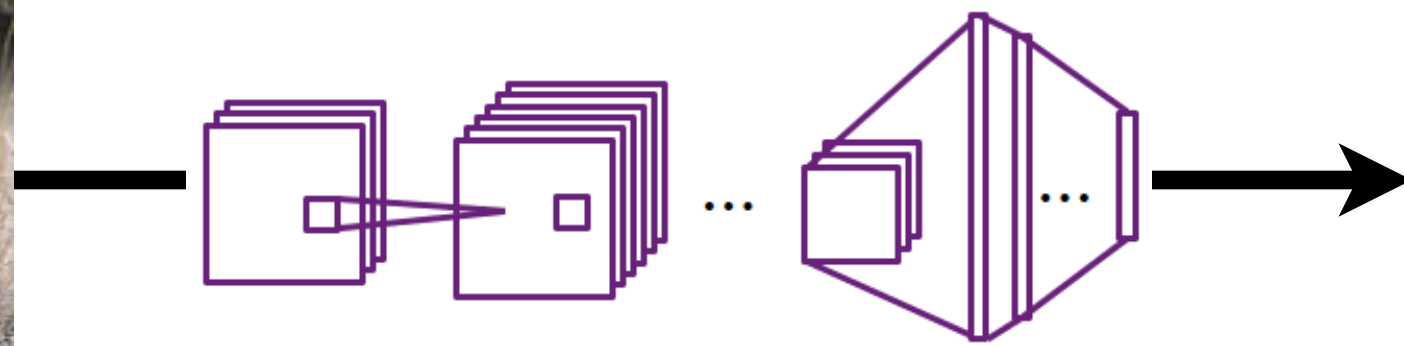Person (x, y, w, h)
Person (x, y, w, h)

COCO
Common Objects in Context

## Segmentation



Horse
Person

COCO
Common Objects in Context

## Instance Segmentation



Horse1
Horse2
Person1
Person2

# Object **Classification**



| Category | Prediction |
|----------|-----------|
| Dog | |
| Cat | |
| Couch | |
| Flowers | $\mathbf{x}^t$ |
| Leopard | |
| … | |

0      1

Probability

**Problem:** For each image predict which category it belongs to out of a fixed set
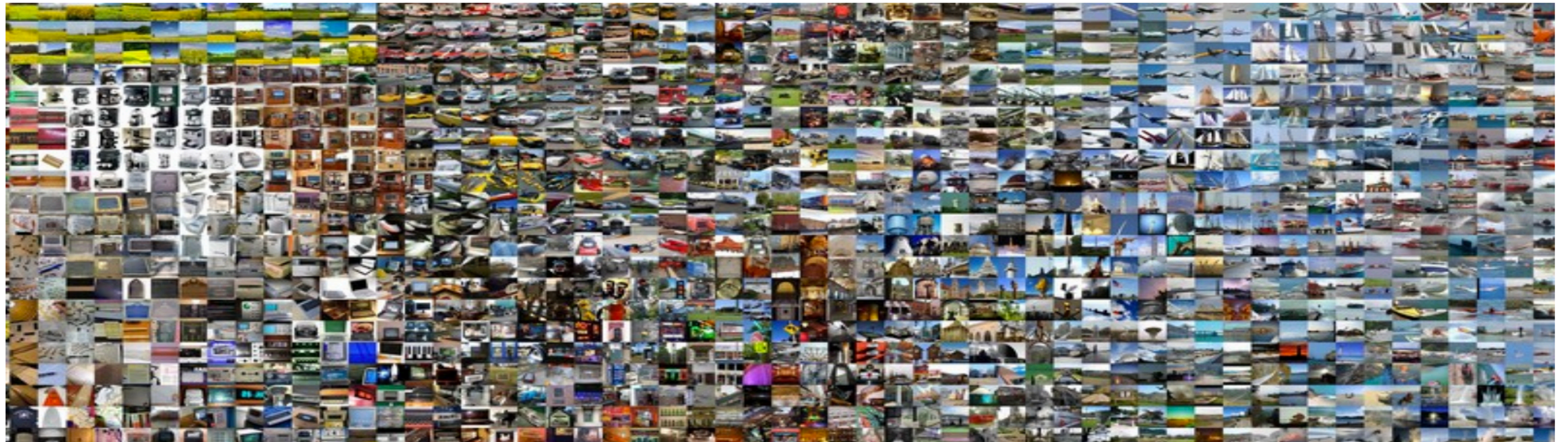
# ImageNet **Competition** (ILSVRC)

Annual competition of image classification at scale

Focuses on a subset of **1K synset** categories

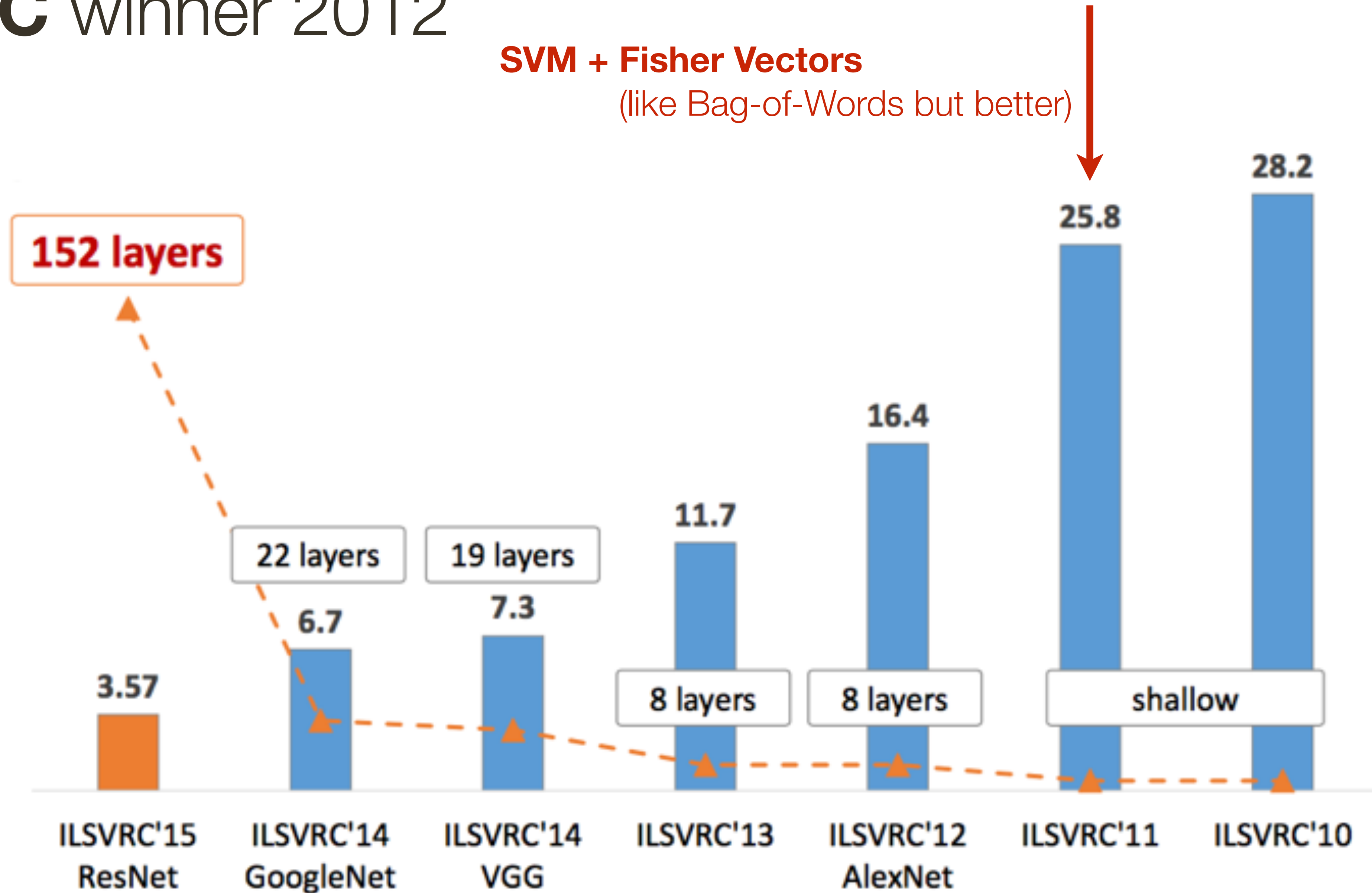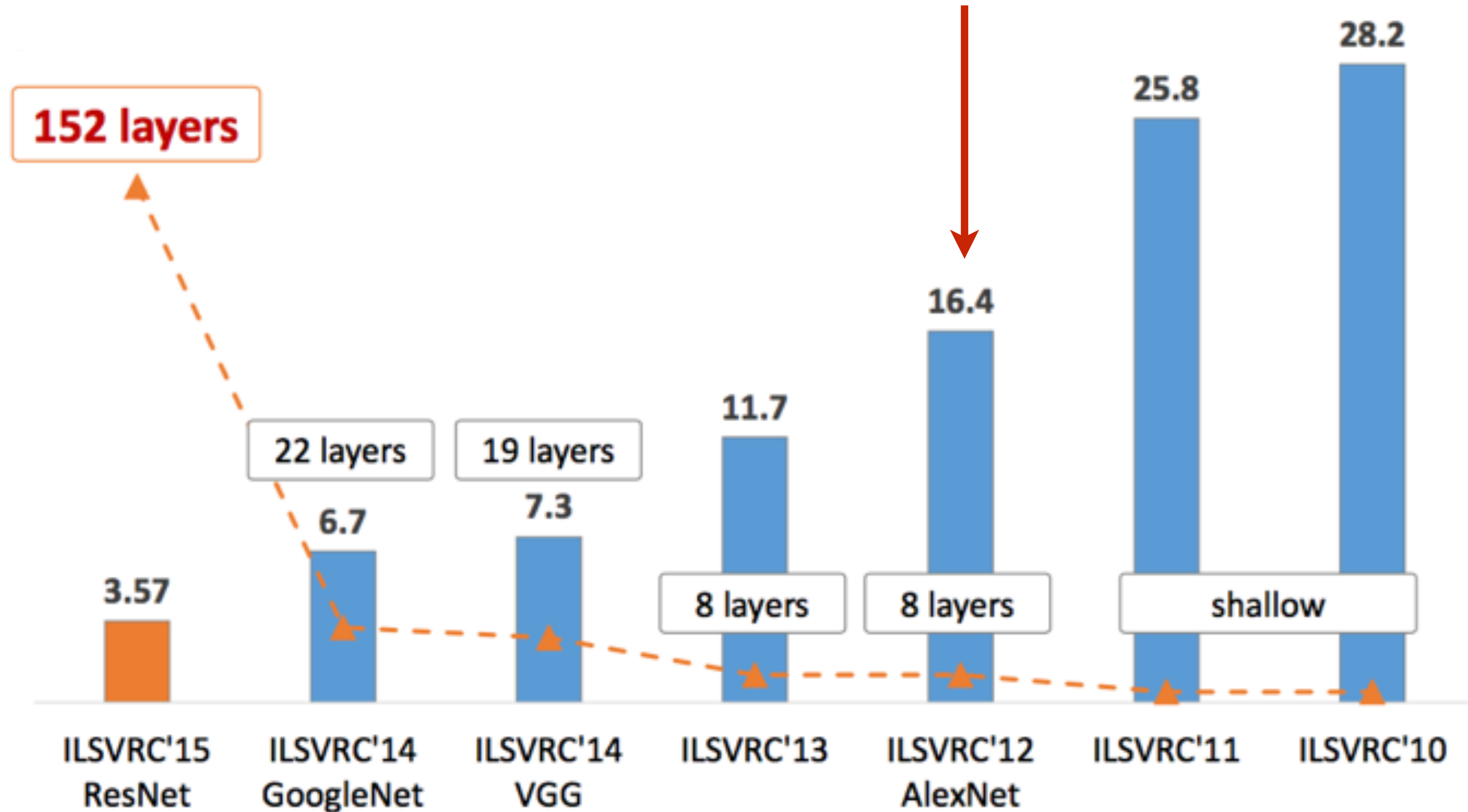**Scoring:** need to predict true label within top K (K=5)

# ILSVRC winner 2012



SVM + Fisher Vectors
(like Bag-of-Words but better)

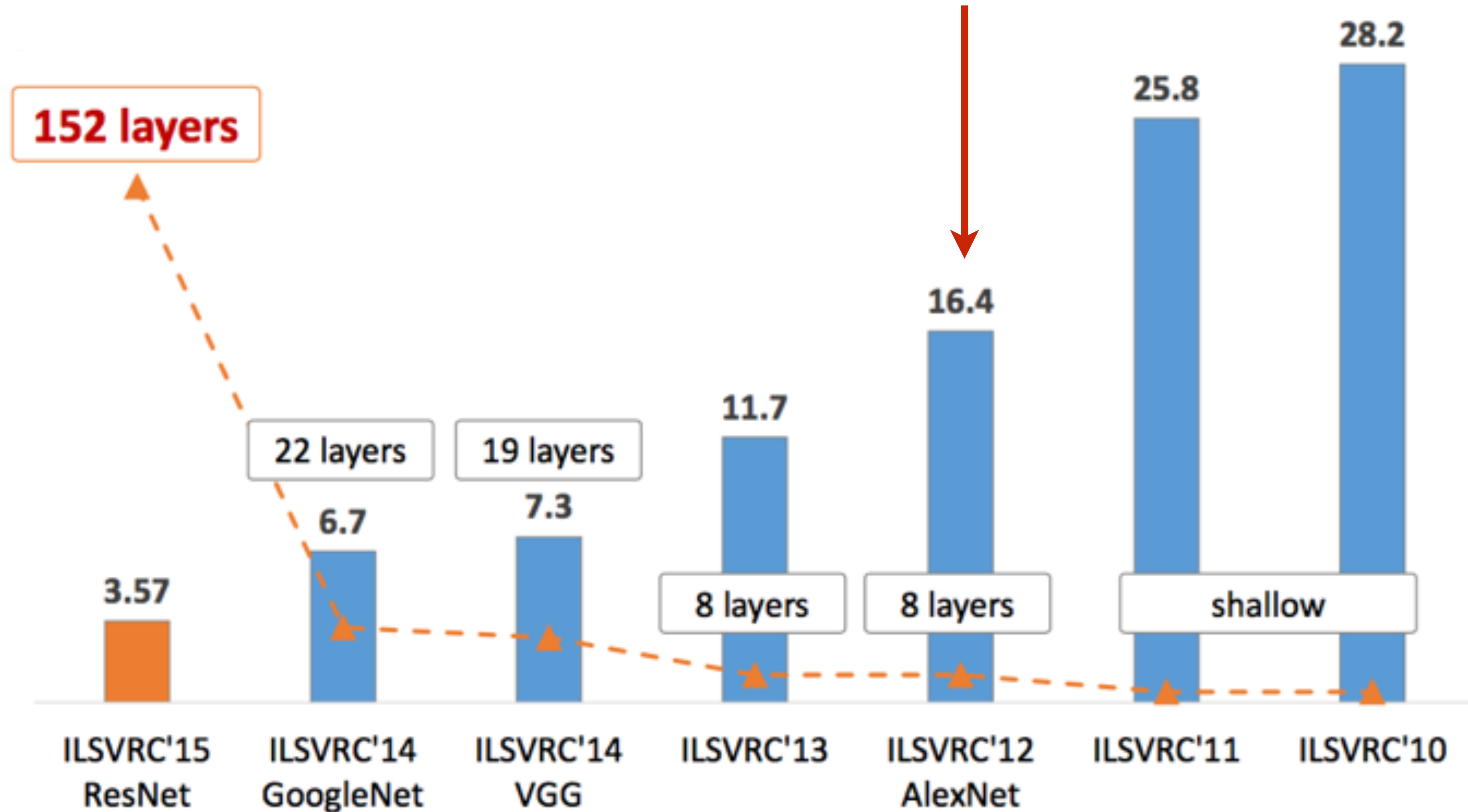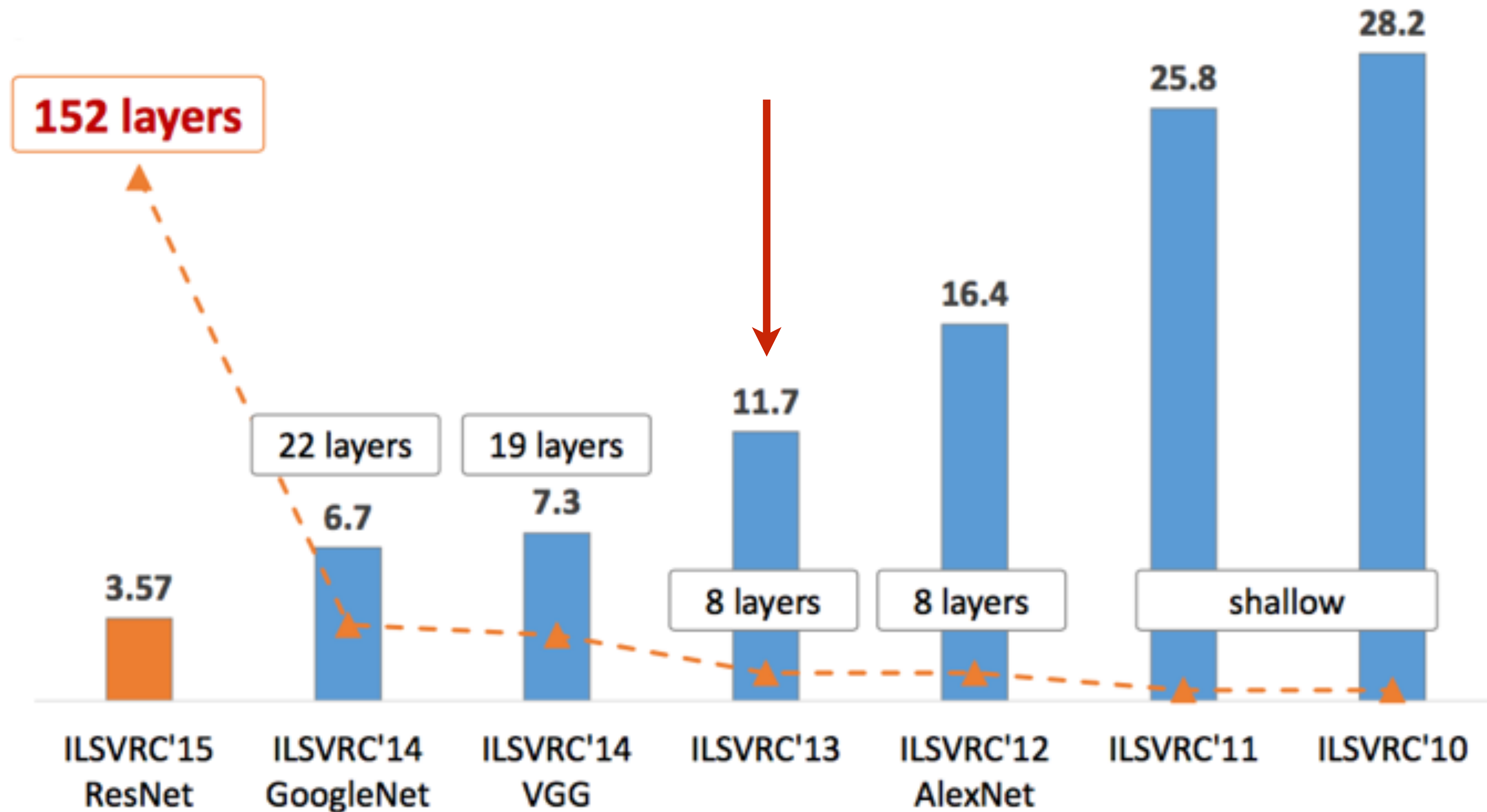* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# ILSVRC winner 2012
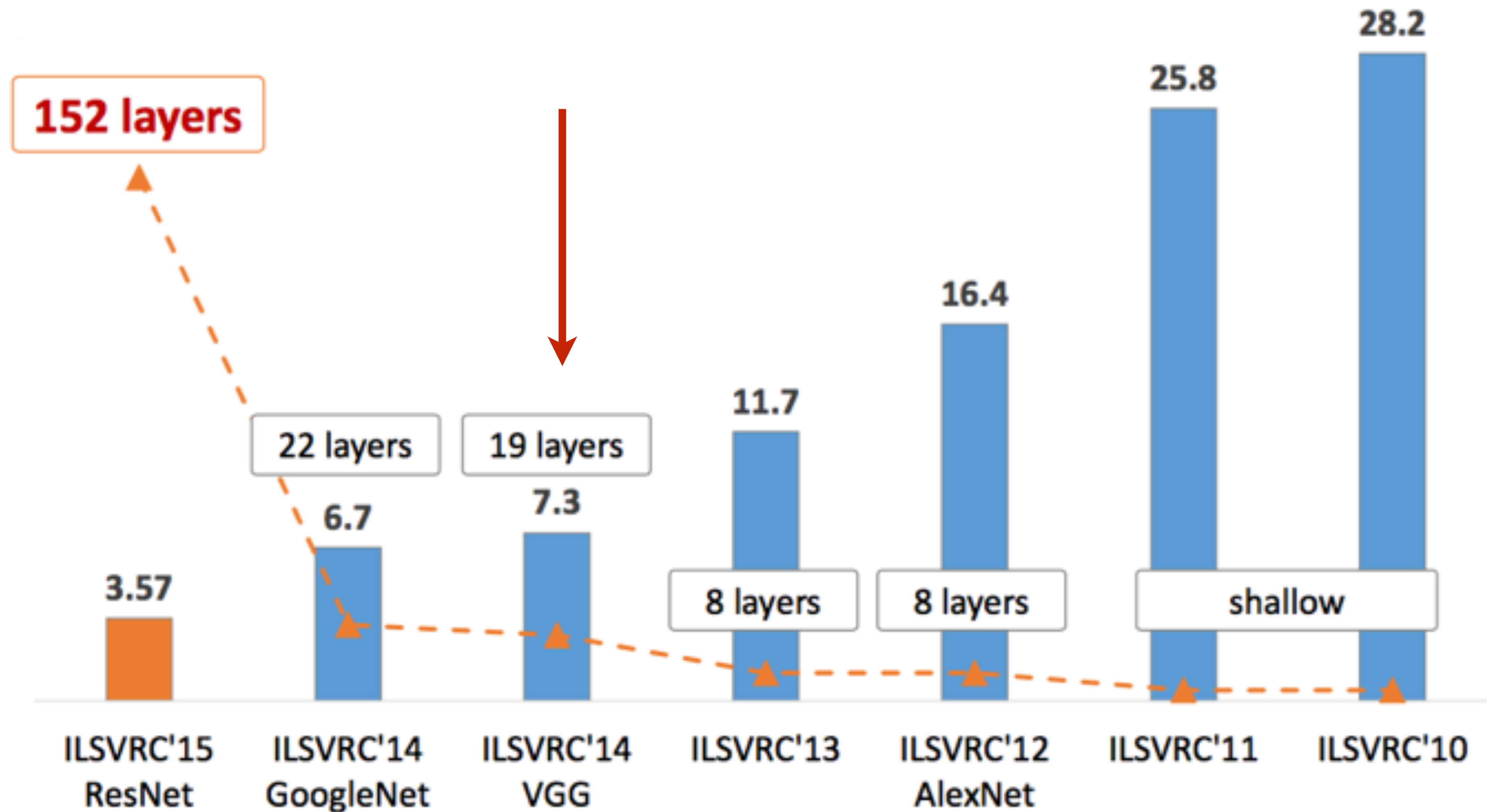
# ILSVRC winner 2012

# ILSVRC winner 2012

# ILSVRC winner 2012

# ILSVRC winner 2012

# ResNet

even deeper — **152 layers**!

using residual connections



F(x) + x

relu

conv

F(x)

relu

conv

X
identity

X
Residual block

Softmax
FC 1000
Pool
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
...
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128 / 2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
Pool
7x7 conv, 64 / 2
Input

# ResNet: Motivation

What happens when we continue to stacking deeper layers on a "plain" CNN

# ResNet: Motivation

What happens when we continue to stacking deeper layers on a "plain" CNN



Whats the **problem**?

# **ResNet:** Motivation

What happens when we continue to stacking deeper layers on a "plain" CNN



## Whats the **problem**?

# Optimizing **Deep** Neural Networks

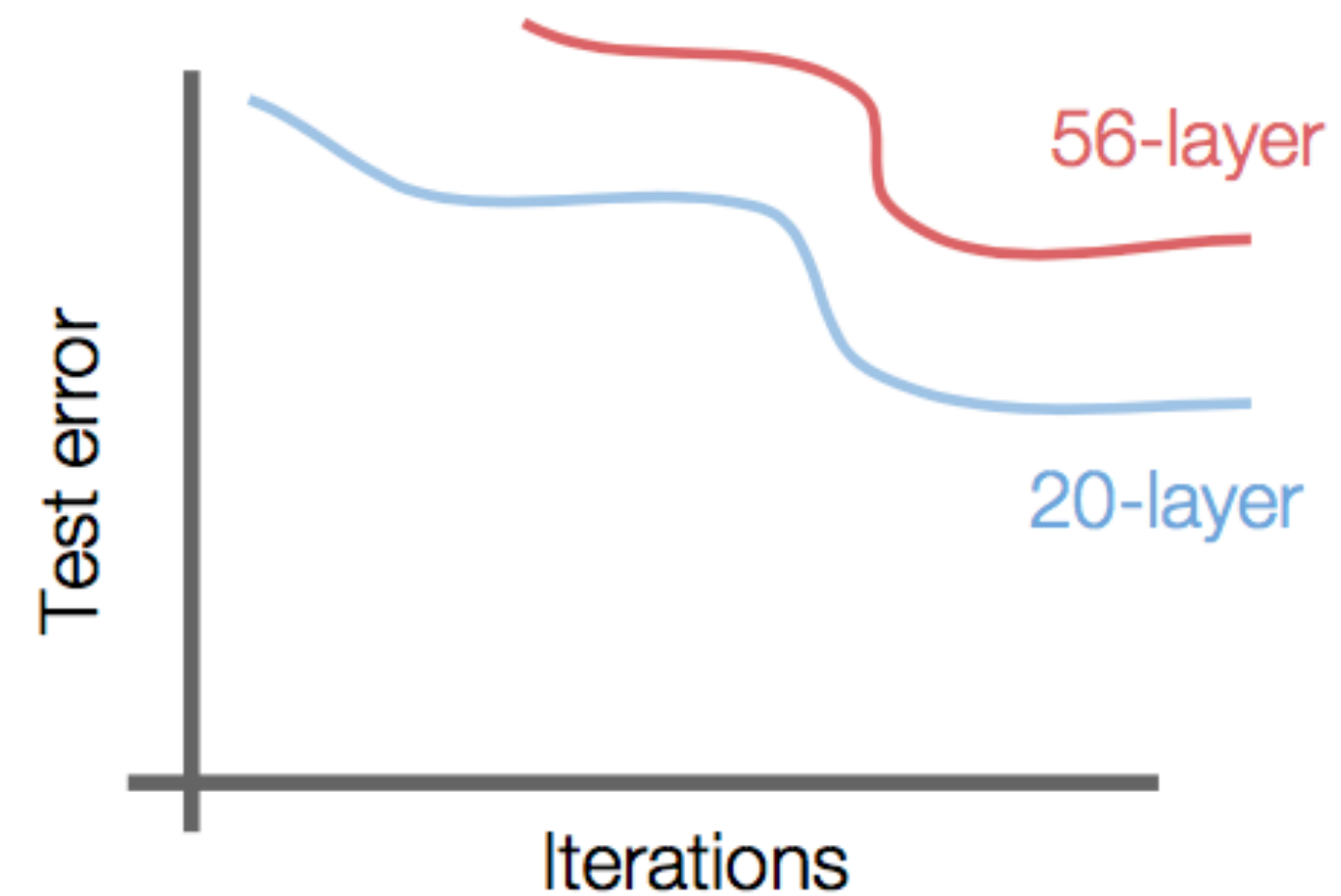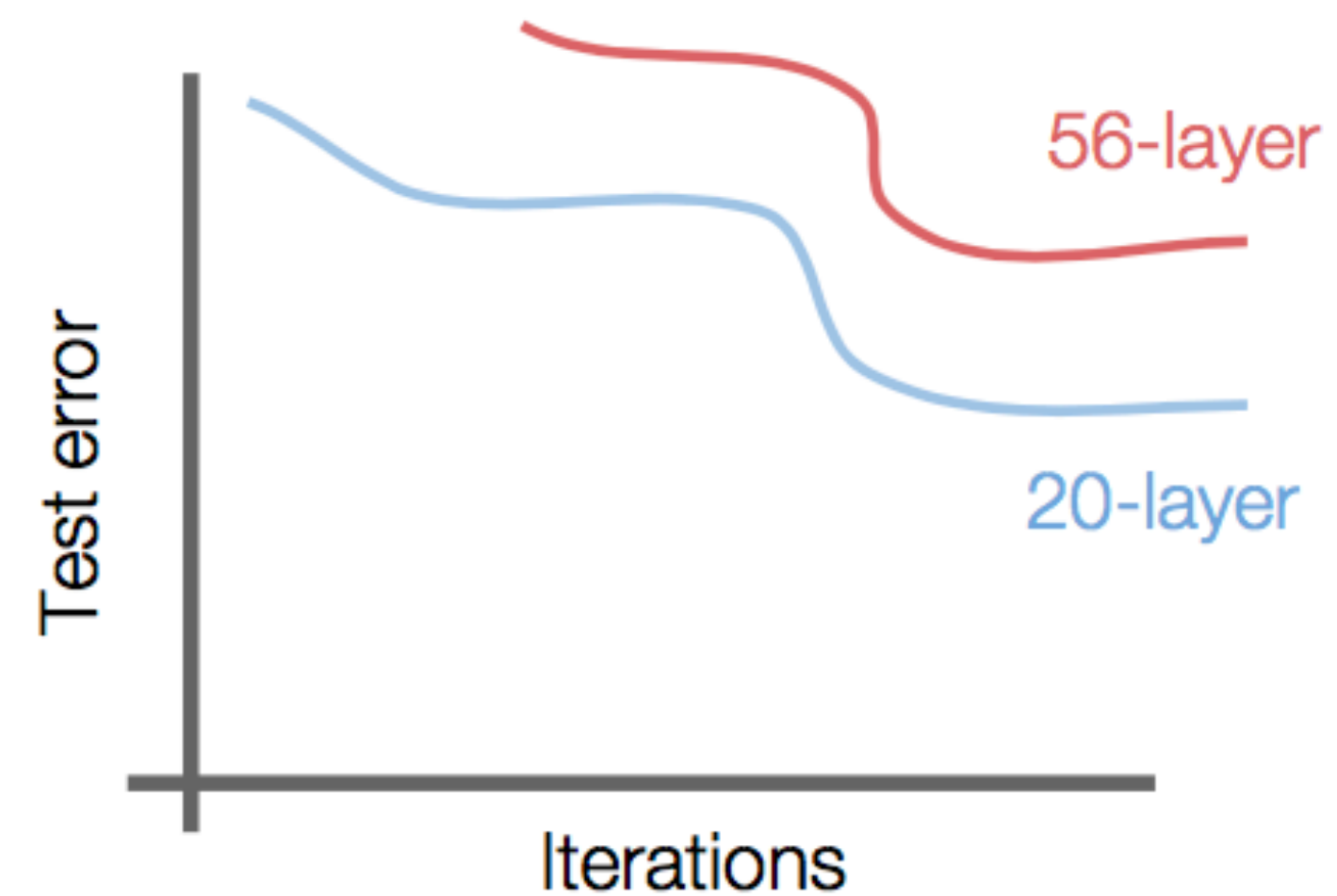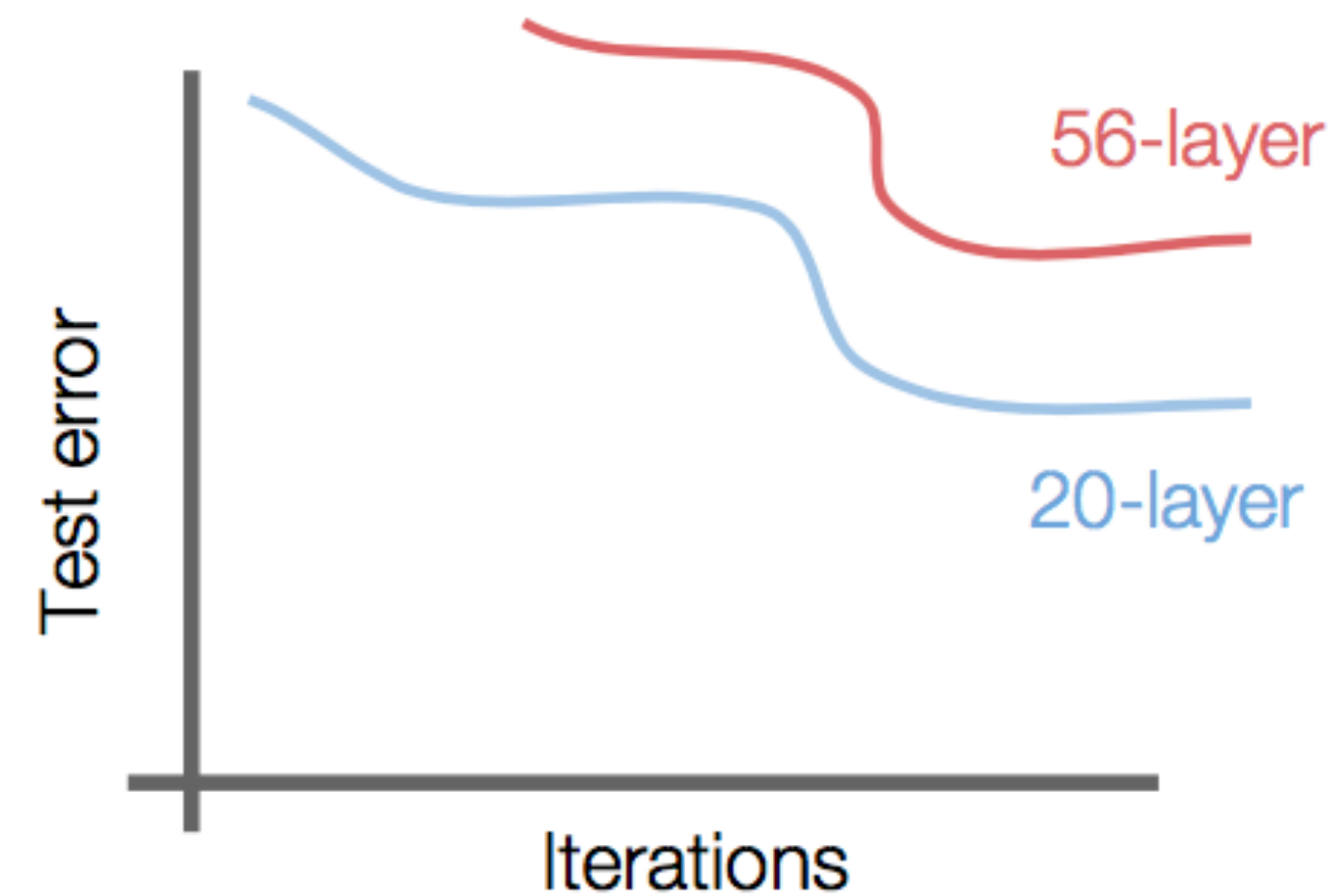$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \overbrace{w_2 \sigma'(z_2)}^{< \frac{1}{4}} \overbrace{w_3 \sigma'(z_3)}^{< \frac{1}{4}} \underbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$

common terms

This is called **vanishing gradient** problem

— makes deep networks hard to train
— later layers learn faster than earlier ones

$$\frac{\partial C}{\partial b_3} = \sigma'(z_3) \overbrace{w_4 \sigma'(z_4) \frac{\partial C}{\partial a_4}}$$

# **ResNet:** Motivation

**Hypothesis:** deeper models are harder to optimize (optimization problem)

# **ResNet:** Motivation

**Hypothesis:** deeper models are harder to optimize (optimization problem)

**Observation:** the deeper model should (conceptually) perform just as well (e.g., take shallower model and use identity for all remaining layers)

# **ResNet:** Motivation

[ He et al., 2015 ]

**Hypothesis:** deeper models are harder to optimize (optimization problem)

**Observation:** the deeper model should (conceptually) perform just as well (e.g., take shallower model and use identity for all remaining layers)

How do we implement this idea in practice

# ResNet

**Solution:** use network to fit residual mapping instead of directly trying to fit a desired underlying mapping

$$H(x) = F(x) + X$$

Use layers to fit **residual**
$F(x) = H(x) - X$ instead of $H(x)$ directly



"Plain" layers

Residual block

# ResNet

## Full details

— Stacked **residual blocks**

— Every residual block consists of **two 3x3 filters**

— Periodically double # of filters and downsample spatially using stride of 2

— Additional convolutional layer in the beginning

— **No FC layers** at the end (only FC to output 1000 classes)

# ILSVRC winner 2012

# **ILSVRC** winner 2012



MSRA @ ILSVRC & COCO 2015 Competitions

• **1st places** in all five main tracks
  • ImageNet Classification: *"Ultra-deep"* (quote Yann) **152-layer** nets
  • ImageNet Detection: **16%** better than 2nd
  • ImageNet Localization: **27%** better than 2nd
  • COCO Detection: **11%** better than 2nd
  • COCO Segmentation: **12%** better than 2nd

**152 layers**

28.2

25.8

11.7

7.3

6.7

3.57

22 layers

19 layers

8 layers

8 layers

shallow

| ILSVRC'15 ResNet | ILSVRC'14 GoogleNet | ILSVRC'14 VGG | ILSVRC'13 | ILSVRC'12 AlexNet | ILSVRC'11 | ILSVRC'10 |

# Regularization: Stochastic Depth

Effectively "dropout" but for layers

Stochastically with some probability **turn off some layer** (for each batch)

Effectively trains a collection of neural networks



Residual Block

$f_\ell(H_{\ell-1})$

$H_{\ell-1}$ | Input → Convolution → Batch Norm. → ReLU → Convolution → Batch Norm. → + → ReLU → Output $H_\ell$

$id(H_{\ell-1})$

# **ResNet**: A little theory

One can view a sequence of outputs from residual layers as a **Dynamical System**



[ Cheng et al., ICLR 2018 ]

# ResNet: A little theory

One can view a sequence of outputs from residual layers as a **Dynamical System**



$$\mathbf{Y}_{j+1} = \mathbf{Y}_j + G(\mathbf{Y}_j, \boldsymbol{\theta}_j)$$

[ Cheng et al., ICLR 2018 ]

# **ResNet**: A little theory

One can view a sequence of outputs from residual layers as a **Dynamical System**



What happens if you take more layers and take smaller steps?

[ Chen et al., NIPS 2018 **best paper** ]

# **ResNet**: A little theory

One can view a sequence of outputs from residual layers as a **Dynamical System**



What happens if you take more layers and take smaller steps?

You can actually treat a neural network as an **ODE**: $\dfrac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \boldsymbol{\theta})$

[ Chen et al., NIPS 2018 **best paper** ]

# An Aside: Neural Network **Cascades**

Cascade

Model 0 → Prediction confident enough?
- Yes → Prediction of model 0 — **Early exit (easy examples)**
- No → Model 1 → Averaged Prediction — **Use more models (hard examples)**

# Comparing **Complexity**



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Computer **Vision Problems** (no language for now)

## Categorization



Multi-**class:**    Horse
                    Church
                    Toothbrush
                    **Person**

IMAGENET

Multi-**label**:    **Horse**
                    Church
                    Toothbrush
                    **Person**

## Detection



Horse (x, y, w, h)
Horse (x, y, w, h)
Person (x, y, w, h)
Person (x, y, w, h)

COCO
Common Objects in Context

# Object **Detection** as Regression Problem



CAT (x, y, w ,h)

# Object **Detection** as Regression Problem



CAT (x, y, w ,h)

DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
…

# Object **Detection** as Regression Problem

CAT (x, y, w ,h)

**Problem:** each image needs a different number of outputs

DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
DUCK (x, y, w ,h)
…

# Object **Detection** as Classification Problem



| Category | Prediction |
|----------|-----------|
| Dog | No |
| Cat | No |
| Couch | No |
| Flowers | No |
| Background | **Yes** |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem

| Category | Prediction |
|----------|------------|
| Dog | No |
| Cat | No |
| Couch | No |
| Flowers | No |
| Background | **Yes** |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem



| Category | Prediction |
|------------|------------|
| Dog | **Yes** |
| Cat | No |
| Couch | No |
| Flowers | No |
| Background | No |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN
classifies each patch as object or background

# Object **Detection** as Classification Problem



| Category | Prediction |
|---|---|
| Dog | **Yes** |
| Cat | No |
| Couch | No |
| Flowers | No |
| Background | No |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem



| Category | Prediction |
|------------|------------|
| Dog | No |
| Cat | **Yes** |
| Couch | No |
| Flowers | No |
| Background | No |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem

**Problem:** Need to apply CNN to **many** patches in each image



| Category | Prediction |
|---|---|
| Dog | No |
| Cat | **Yes** |
| Couch | No |
| Flowers | No |
| Background | No |
| … | … |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# **Region** Proposals (older idea in vision)

[ Alexe et al, TPAMI 2012 ]
[ Uijkings et al, IJCV 2013 ]
[ Cheng et al, CVPR 2014 ]
[ Zitnick and Dollar, ECCV 2014 ]

Find image **regions that are likely contain objects** (any object at all)

- typically works by looking at histogram distributions, region aspect ratio, closed contours, coherent color

Relatively **fast to run** (Selective Search gives 1000 region proposals in a few seconds on a CPU)



**Goal:** Get "true" object regions to be in as few top K proposals as possible

# R-CNN

Input **Image**

# R-CNN

**Regions of Interest** from
a proposal method (~2k)

Input **Image**

* image from Ross Girshick

# R-CNN

**Warped** image regions

**Regions of Interest** from
a proposal method (~2k)

Input **Image**

* image from Ross Girshick

# R-CNN

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

Input **Image**

* image from Ross Girshick

# R-CNN

**Classify** regions with SVM

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

Input **Image**

# R-CNN

**Linear Regression** for bounding box offsets

**Classify** regions with SVM

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

Input **Image**

# R-CNN

R-CNN (Regions with CNN features) algorithm:

— Extract promising candidate regions using an object proposals algorithm

— Resize each proposal window to the size of the input layer of a trained convolutional neural network

— Input each resized image patch to the convolutional neural network

**Implementation detail**: Instead of using the classification scores of the network directly, the output of the final fully-connected layer can be used as an input feature to a trained support vector machine (SVM)

# Fast **R-CNN**

Input **Image**

# Fast **R-CNN**

[ Girshick et al, ICCV 2015 ]



"conv5" feat

Forward wh

ConvNet

Input **Image**

* image from Ross Girshick

# Fast **R-CNN**

"conv5" feature map

Forward prop the **whole image** through CNN

ConvNet

Input **Image**

# Fast **R-CNN**

**Regions of Interest**
from the proposal method



"conv5" feature map

Forward prop the **whole image** through CNN

ConvNet

Input **Image**

# Fast **R-CNN**

**Regions of Interest** from the proposal method

"RoI Pooling" layer

"conv5" feature map

Forward prop the **whole image** through CNN

ConvNet

Input **Image**

Girshick, "Fast R-C
Figure copyright R

* image from Ross Girshick

# RoI Align



15 x 15 pixel Region of Interest
in the original image

Corresponding region in the
Feature Map (2.93 x 2.93)

Original Image: 128 x 128

CNN

Feature Map: 25 x 25

bilinear
interpolation

variable size RoI

# Fast **R-CNN**

Multi-task loss

Log loss + Smooth L1 loss

Object classification

Linear + softmax

Linear — Bounding box regression

FCs

**Regions of Interest** from the proposal method

"RoI Pooling" layer

"conv5" feature map

Forward prop the **whole image** through CNN

ConvNet

Input **Image**

* image from Ross Girshick

# Fast **R-CNN:** Training



Log loss + Smooth L1 loss — Multi-task loss

[ Girshick et al, ICCV 2015 ]

Object classification

Linear + softmax

Linear — Bounding box regression

FCs

**Regions of Interest** from the proposal method

"RoI Pooling" layer

"conv5" feature map

ConvNet

Forward prop the **whole image** through CNN

Input **Image**

* image from Ross Girshick

# R-CNN vs. SPP vs. Fast R-CNN

[ Girshick et al, CVPR 2014 ]

[ Girshick et al, ICCV 2015 ]

[ He et al, ECCV 2014 ]

## Training time (Hours)

| Model | Hours |
|---|---|
| R-CNN | 84 |
| SPP-Net | 25.5 |
| Fast R-CNN | 8.75 |

## Test time (seconds)

| Model | Including Region proposals | Excluding Region Proposals |
|---|---|---|
| R-CNN | 49 | 47 |
| SPP-Net | 4.3 | 2.3 |
| Fast R-CNN | 2.3 | 0.32 |

# R-CNN vs. SPP vs. Fast R-CNN

**Observation:** Performance dominated by the region proposals at this point!

# Faster **R-CNN**

Make CNN do proposals!

Insert **Region Proposal Network (RPN)** to predict proposals from features

Jointly train with 4 losses:
1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



Classification loss

Bounding-box regression loss

RoI pooling

Classification loss

Bounding-box regression loss

proposals

Region Proposal Network

feature map

CNN

image

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

# YOLO: You Only Look Once

Input image
3 x H x W

Divide image into grid
7 x 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 x 7 x (5 * B + C)

# YOLO: You Only Look Once

Input image
3 x H x W

Divide image into grid
7 x 7

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

P(Object)  X  Y  Width  Height  P(Object)  X  Y  Width  Height  P(Cat | Object)  P(Bird | Object)  P(TV | Object)

1st - 5th          6th - 10th         11th - 30th
Box #1             Box #2             Class Probabilities

# YOLO v2

# YOLO v2

http://pjreddie.com/yolo

# YOLO: You Only Look Once

Input image
3 x H x W

Divide image into grid
7 x 7

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

P(Object) X Y Width Height P(Object) X Y Width Height P(Cat | Object) P(Bird | Object) P(TV | Object)

1st - 5th          6th - 10th          11th - 30th
Box #1             Box #2              Class Probabilities

# Computer **Vision Problems** (no language for now)

## Segmentation



Horse
Person

# Semantic **Segmentation**

Label **every pixel** with a category label (without differentiating instances)

# Semantic **Segmentation:** Sliding Window

Extract **patches**          **Classify center pixel** with CNN



Cow

Cow

Grass

# Semantic **Segmentation:** Sliding Window

Extract **patches**          **Classify center pixel** with CNN

Cow

Cow

Grass

**Problem:** VERY inefficient, no reuse of computations for overlapping patches

# Semantic **Segmentation:** Fully Convolutional CNNs

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



**Input Image**

3 x H x W

**Convolutions**

D x H x W

**Class Scores**

C x H x W

**Predicted** Labels

H x W

# Semantic **Segmentation:** Fully Convolutional CNNs

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



| Input **Image** | | **Convolutions** | | Class **Scores** | **Predicted** Labels |
|---|---|---|---|---|---|
| 3 x H x W | | D x H x W | | C x H x W | H x W |

**Problem:** Convolutions at the original image scale will be very expensive

# Semantic **Segmentation:** Fully Convolutional CNNs

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



Input **Image**

3 x H x W

Med-res:
$D_2$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

High-res:
$D_1$ x H/2 x W/2

High-res:
$D_1$ x H/2 x W/2

**Predicted** Labels

H x W

[ Long et al, CVPR 2015 ]
[ Noh et al, ICCV 2015 ]

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# Semantic **Segmentation:** Fully Convolutional CNNs

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



Input **Image**

3 x H x W

Med-res:
$D_2$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

High-res:
$D_1$ x H/2 x W/2

High-res:
$D_1$ x H/2 x W/2

**Predicted** Labels

H x W

**Downsampling** = Pooling

**Upsampling** = ???

[ Long et al, CVPR 2015 ]
[ Noh et al, ICCV 2015 ]

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# In-network **Up Sampling** (a.k.a "Unpooling")

Nearest Neighbor

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

$\longrightarrow$

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

**Input:** 2 x 2          **Output:** 4 x 4

# In-network **Up Sampling** (a.k.a "Unpooling")

**Nearest Neighbor**

$$
\begin{array}{cc}
1 & 2 \\
3 & 4
\end{array}
\longrightarrow
\begin{array}{cc|cc}
1 & 1 & 2 & 2 \\
1 & 1 & 2 & 2 \\
\hline
3 & 3 & 4 & 4 \\
3 & 3 & 4 & 4
\end{array}
$$

**Input:** 2 x 2          **Output:** 4 x 4

**"Bed of Nails"**

$$
\begin{array}{cc}
1 & 2 \\
3 & 4
\end{array}
\longrightarrow
\begin{array}{cc|cc}
1 & 0 & 2 & 0 \\
0 & 0 & 0 & 0 \\
\hline
3 & 0 & 4 & 0 \\
0 & 0 & 0 & 0
\end{array}
$$

**Input:** 2 x 2          **Output:** 4 x 4

# In-network **Up Sampling:** Max Unpooling

**Max Pooling**
Remember which element was max!

**Max Unpooling**
Use positions from pooling layer



**Input:** 4 x 4

**Output:** 2 x 2

Rest of the network

**Input:** 2 x 2

**Output:** 4 x 4

Corresponding pairs of downsampling and upsampling layers

# In-network **Up Sampling:** Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 1 pad 1



Dot product between filter and input

**Input:** 4 x 4

**Output:** 4 x 4

# In-network **Up Sampling:** Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 1 pad 1



**Input:** 4 x 4

Dot product between filter and input

**Output:** 4 x 4

# In-network **Up Sampling:** Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



Dot product between filter and input

**Input:** 4 x 4

**Output:** 2 x 2

# In-network **Up Sampling:** Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



**Input:** 4 x 4

Dot product between filter and input



**Output:** 2 x 2

Filter moves 2 pixels in the **input** for every one pixel in the **output**

Stride gives ratio in movement in input vs output

# In-network **Up Sampling:** Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

**Output:** 4 x 4

**Input:** 2 x 2

# In-network **Up Sampling:** Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

Input gives weight for filter

**Input:** 2 x 2

**Output:** 4 x 4

# In-network **Up Sampling:** Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

Sum where output overlaps

Input gives weight for filter

**Input:** 2 x 2

**Output:** 4 x 4

Filter moves 2 pixels in the **output** for every one pixel in the **input**

Stride gives ratio in movement in output vs input

# Transpose Convolution: 1-D Example

**Input**

**Filter**

**Output**



Output contains copies of the filter weighted multiplied by the input, summing at overlaps in the output

# **U-Net** Architecture

ResNet-like Fully convolutional CNN



[ Ronneberger  et al, CVPR 2015 ]

# Computer **Vision Problems** (no language for now)

Instance Segmentation



Horse1
Horse2
Person1
Person2

# **Mask** R-CNN



[ He et al, 2017 ]

# **Mask** R-CNN



[ He et al, 2017 ]

# Summary

Common types of layers:

1. **Convolutional** Layer
   — Parameters define a set of learnable filters

2. **Pooling** Layer
   — Performs a downsampling along the spatial dimensions

3. **Fully-Connected** Layer
   — As in a regular neural network

Each layer accepts an input 3D volume and transforms it to an output 3D volume through a differentiable function

# Summary

The parameters of a neural network are learned using **backpropagation**, which computes gradients via recursive application of the chain rule

A **convolutional neural network** assumes inputs are images, and constrains the network architecture to reduce the number of parameters

A **convolutional layer** applies a set of learnable filters

A **pooling layer** performs spatial downsampling

A **fully-connected** layer is the same as in a regular neural network

Convolutional neural networks can be seen as learning a hierarchy of filters

# **Attention** Layer

**Inputs**:
**Query vectors**: $\mathbf{Q}$ (Shape: $N_Q \times D_Q$)
**Input vectors**: $\mathbf{X}$ (Shape: $N_X \times D_X$)
**Key matrix**: $\mathbf{W_K}$ (Shape: $D_X \times D_Q$)
**Value matrix**: $\mathbf{W_V}$ (Shape: $D_X \times D_V$)

**Computation**:
**Key vectors**: $\mathbf{K} = \mathbf{X}\mathbf{W_K}$ (Shape: $N_X \times D_Q$)
**Value Vectors**: $\mathbf{V} = \mathbf{X}\mathbf{W_V}$ (Shape: $N_X \times D_V$)
**Similarities**: $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$) $E_{i,j} = (\mathbf{Q_i} \cdot \mathbf{K_j}) / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = A\mathbf{V}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j}\mathbf{V_j}$

$X_1$

$X_2$

$X_3$

$Q_1$ $Q_2$ $Q_3$ $Q_4$

# **Attention** Layer

**Inputs**:
**Query vectors**: $\mathbf{Q}$ (Shape: $N_Q$ x $D_Q$)
**Input vectors**: $\mathbf{X}$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $\mathbf{W_K}$ (Shape: $D_X$ x $D_Q$)
**Value matrix:** $\mathbf{W_V}$ (Shape: $D_X$ x $D_V$)

**Computation**:
**Key vectors**: $\mathbf{K}$ = $\mathbf{X}\mathbf{W_K}$  (Shape: $N_X$ x $D_Q$)
**Value Vectors**: $\mathbf{V}$ = $\mathbf{X}\mathbf{W_V}$ (Shape: $N_X$ x $D_V$)
**Similarities**: $E = \mathbf{Q}\mathbf{K}^\mathsf{T} / \sqrt{D_Q}$ (Shape: $N_Q$ x $N_X$) $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$
**Attention weights**: A = softmax(E, dim=1)  (Shape: $N_Q$ x $N_X$)
**Output vectors**: Y = A$\mathbf{V}$ (Shape: $N_Q$ x $D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

$X_1 \rightarrow K_1$

$X_2 \rightarrow K_2$

$X_3 \rightarrow K_3$

$Q_1$  $Q_2$  $Q_3$  $Q_4$

[ slide from Justin Johnson, U Michigan ]

# **Attention** Layer

**Inputs**:
**Query vectors**: $\textcolor{green}{Q}$ (Shape: $N_Q \times D_Q$)
**Input vectors**: $\textcolor{blue}{X}$ (Shape: $N_X \times D_X$)
**Key matrix**: $\textcolor{orange}{W_K}$ (Shape: $D_X \times D_Q$)
**Value matrix**: $\textcolor{purple}{W_V}$ (Shape: $D_X \times D_V$)

**Computation**:
**Key vectors**: $\textcolor{orange}{K} = \textcolor{blue}{X}\textcolor{orange}{W_K}$ (Shape: $N_X \times D_Q$)
**Value Vectors**: $\textcolor{purple}{V} = \textcolor{blue}{X}\textcolor{purple}{W_V}$ (Shape: $N_X \times D_V$)
**Similarities**: $E = \textcolor{green}{Q}\textcolor{orange}{K}^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$) $E_{i,j} = (\textcolor{green}{Q_i} \cdot \textcolor{orange}{K_j}) / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = A\textcolor{purple}{V}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j}\textcolor{purple}{V_j}$



[ slide from Justin Johnson, U Michigan ]

# **Attention** Layer

**Inputs**:

**Query vectors**: $\mathbf{Q}$ (Shape: $N_Q \times D_Q$)

**Input vectors**: $\mathbf{X}$ (Shape: $N_X \times D_X$)

**Key matrix**: $\mathbf{W_K}$ (Shape: $D_X \times D_Q$)

**Value matrix**: $\mathbf{W_V}$ (Shape: $D_X \times D_V$)

**Computation**:

**Key vectors**: $\mathbf{K} = \mathbf{X}\mathbf{W_K}$ (Shape: $N_X \times D_Q$)

**Value Vectors**: $\mathbf{V} = \mathbf{X}\mathbf{W_V}$ (Shape: $N_X \times D_V$)

**Similarities**: $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$) $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$

**Attention weights**: $A = \text{softmax}(E, \text{dim=1})$ (Shape: $N_Q \times N_X$)

**Output vectors**: $Y = A\mathbf{V}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j}\mathbf{V}_j$

# **Attention** Layer

**Inputs:**
**Query vectors: Q** (Shape: $N_Q \times D_Q$)
**Input vectors: X** (Shape: $N_X \times D_X$)
**Key matrix: $W_K$** (Shape: $D_X \times D_Q$)
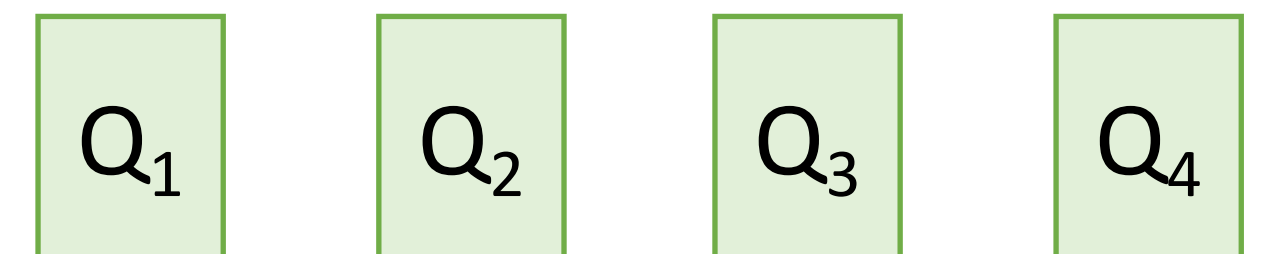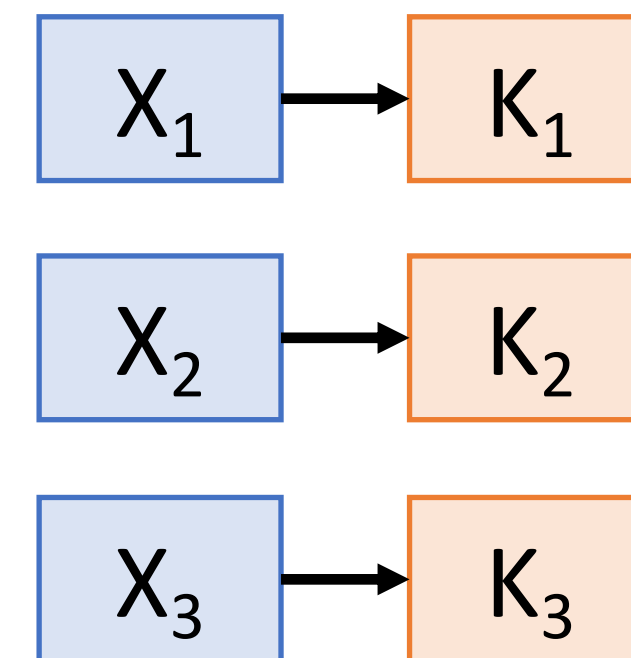**Value matrix: $W_V$** (Shape: $D_X \times D_V$)

**Computation:**
**Key vectors: K** = $XW_K$  (Shape: $N_X \times D_Q$)
**Value Vectors: V** = $XW_V$ (Shape: $N_X \times D_V$)
**Similarities:** E = $QK^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$
**Attention weights:** A = softmax(E, dim=1)  (Shape: $N_Q \times N_X$)
**Output vectors:** Y = A**V** (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} V_j$



Product(→),  Sum(↑)

| $Y_1$ | $Y_2$ | $Y_3$ | $Y_4$ |

| $V_1$ | $A_{1,1}$ | $A_{2,1}$ | $A_{3,1}$ | $A_{4,1}$ |
| $V_2$ | $A_{1,2}$ | $A_{2,2}$ | $A_{3,2}$ | $A_{4,2}$ |
| $V_3$ | $A_{1,3}$ | $A_{2,3}$ | $A_{3,3}$ | $A_{4,3}$ |

Softmax(↑)

| $X_1$ | $K_1$ | $E_{1,1}$ | $E_{2,1}$ | $E_{3,1}$ | $E_{4,1}$ |
| $X_2$ | $K_2$ | $E_{1,2}$ | $E_{2,2}$ | $E_{3,2}$ | $E_{4,2}$ |
| $X_3$ | $K_3$ | $E_{1,3}$ | $E_{2,3}$ | $E_{3,3}$ | $E_{4,3}$ |

| $Q_1$ | $Q_2$ | $Q_3$ | $Q_4$ |

# **Self-attention** Layer

One **query** per **input vector**

**Inputs**:

**Query vectors**: $\mathbf{Q}$ (Shape: $N_Q$ x $D_Q$)

**Input vectors**: $\mathbf{X}$ (Shape: $N_X$ x $D_X$)

**Key matrix**: $\mathbf{W_K}$ (Shape: $D_X$ x $D_Q$)

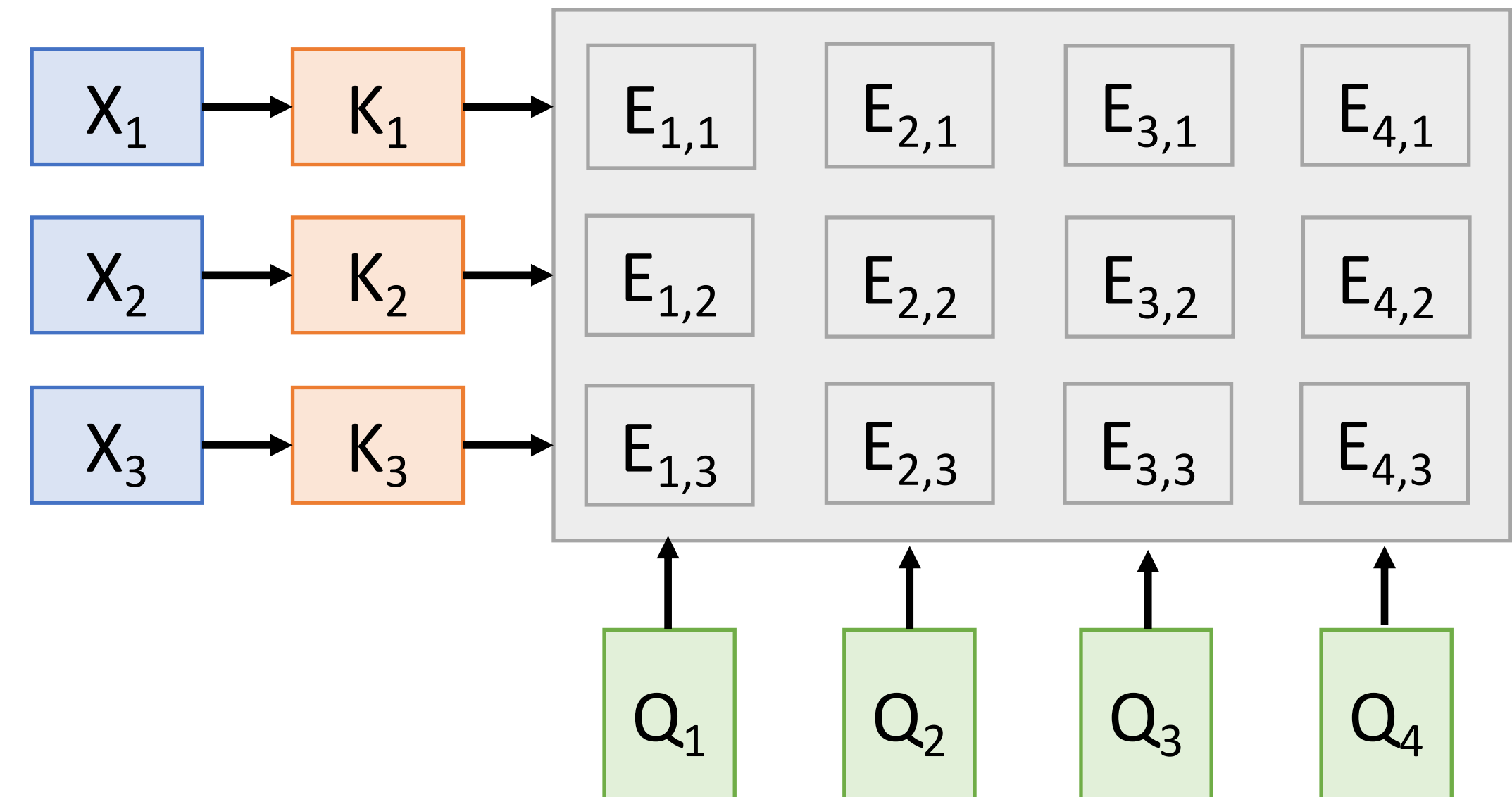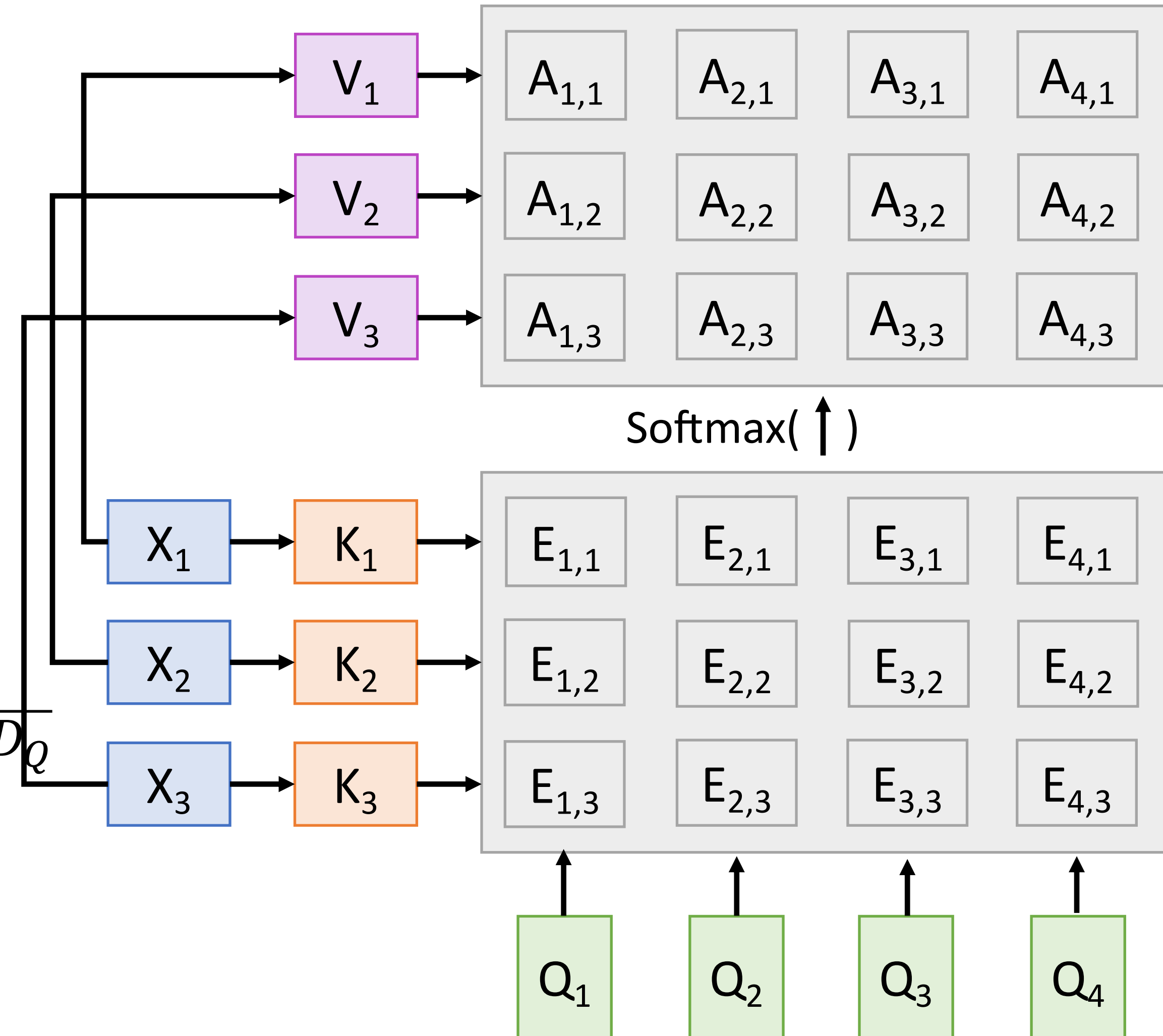**Value matrix:** $\mathbf{W_V}$ (Shape: $D_X$ x $D_V$)

**Computation**:

**Key vectors**: $\mathbf{K} = \mathbf{XW_K}$  (Shape: $N_X$ x $D_Q$)

**Value Vectors**: $\mathbf{V} = \mathbf{XW_V}$ (Shape: $N_X$ x $D_V$)

**Similarities**: $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$ (Shape: $N_Q$ x $N_X$) $E_{i,j} = (\mathbf{Q_i} \cdot \mathbf{K_j}) / \sqrt{D_Q}$

**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$  (Shape: $N_Q$ x $N_X$)

**Output vectors**: $Y = A\mathbf{V}$ (Shape: $N_Q$ x $D_V$) $Y_i = \sum_j A_{i,j}\mathbf{V_j}$

$X_1$    $X_2$    $X_3$

# **Self-attention** Layer

**Inputs**:

**Input vectors**: **X** (Shape: $N_X$ x $D_X$)

**Key matrix**: **$W_K$** (Shape: $D_X$ x $D_Q$)

**Value matrix**: **$W_V$** (Shape: $D_X$ x $D_V$)

**Query matrix**: **$W_Q$** (Shape: $D_X$ x $D_Q$)

**Computation**:

**Query vectors**: **Q** = **$XW_Q$**

**Key vectors**: **K** = **$XW_K$**  (Shape: $N_X$ x $D_Q$)

**Value Vectors**: **V** = **$XW_V$** (Shape: $N_X$ x $D_V$)

**Similarities**: E = **$QK^T$** $/\sqrt{D_Q}$ (Shape: $N_x$ x $N_x$) $E_{i,j}$ = (**$Q_i$** $\cdot$ **$K_j$** ) $/\sqrt{D_Q}$

**Attention weights**: A = softmax(E, dim=1)  (Shape: $N_X$ x $N_X$)

**Output vectors**: Y = A**V** (Shape: $N_X$ x $D_V$) $Y_i$ = $\sum_j A_{i,j}$**$V_j$**

# Self-attention Layer

One **query** per **input vector**

**Inputs**:
**Input vectors**: $\textbf{X}$ (Shape: $N_X \times D_X$)
**Key matrix**: $\textbf{W}_K$ (Shape: $D_X \times D_Q$)
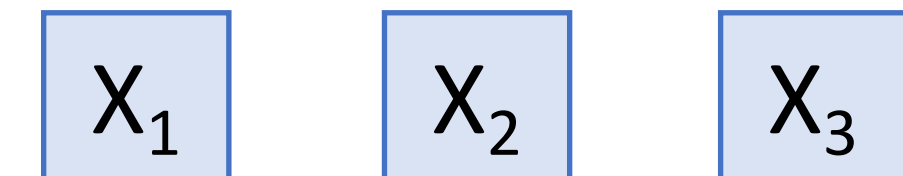**Value matrix**: $\textbf{W}_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $\textbf{W}_Q$ (Shape: $D_X \times D_Q$)

**Computation**:
**Query vectors**: $\textbf{Q} = \textbf{X}\textbf{W}_Q$
**Key vectors**: $\textbf{K} = \textbf{X}\textbf{W}_K$ (Shape: $N_X \times D_Q$)
**Value Vectors**: $\textbf{V} = \textbf{X}\textbf{W}_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = \textbf{Q}\textbf{K}^\mathsf{T} / \sqrt{D_Q}$ (Shape: $N_X \times N_X$) $E_{i,j} = (\textbf{Q}_i \cdot \textbf{K}_j) / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = A\textbf{V}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j}\textbf{V}_j$

# **Self-attention** Layer

One **query** per **input vector**

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix:** $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)
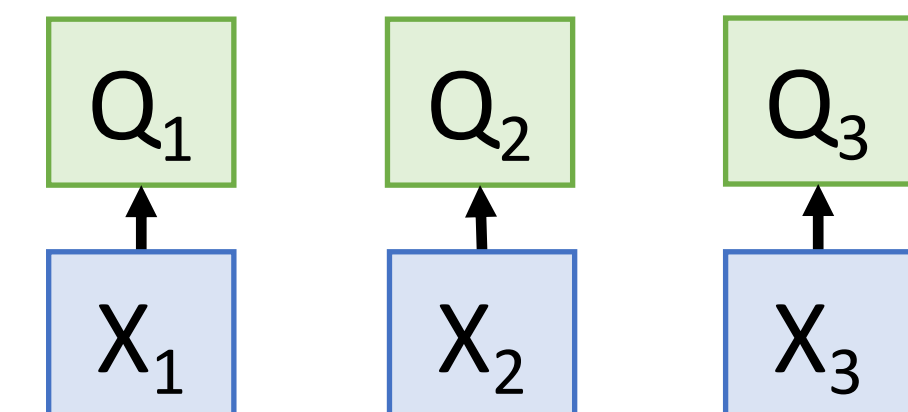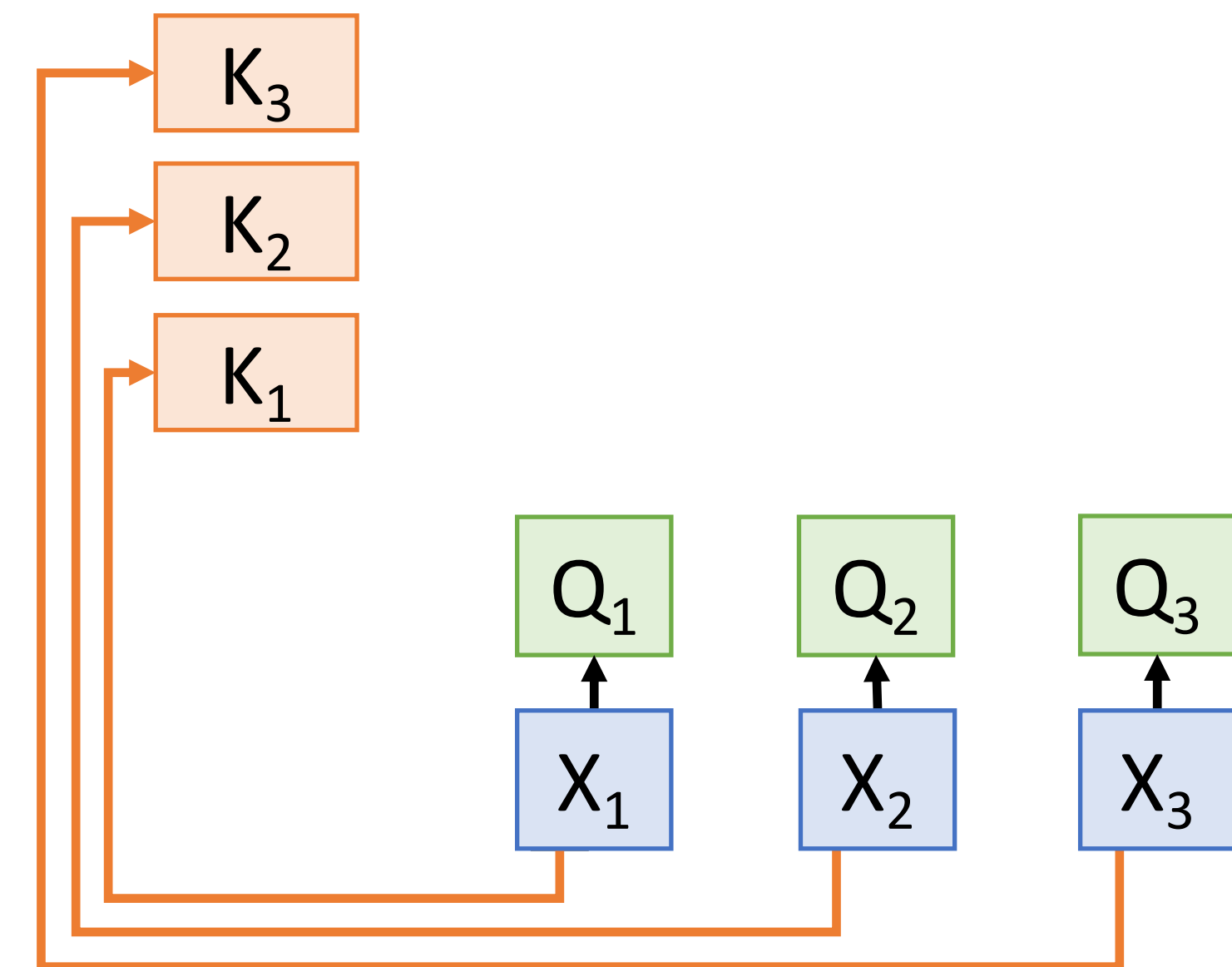
Computation:
**Query vectors**: $Q = XW_Q$
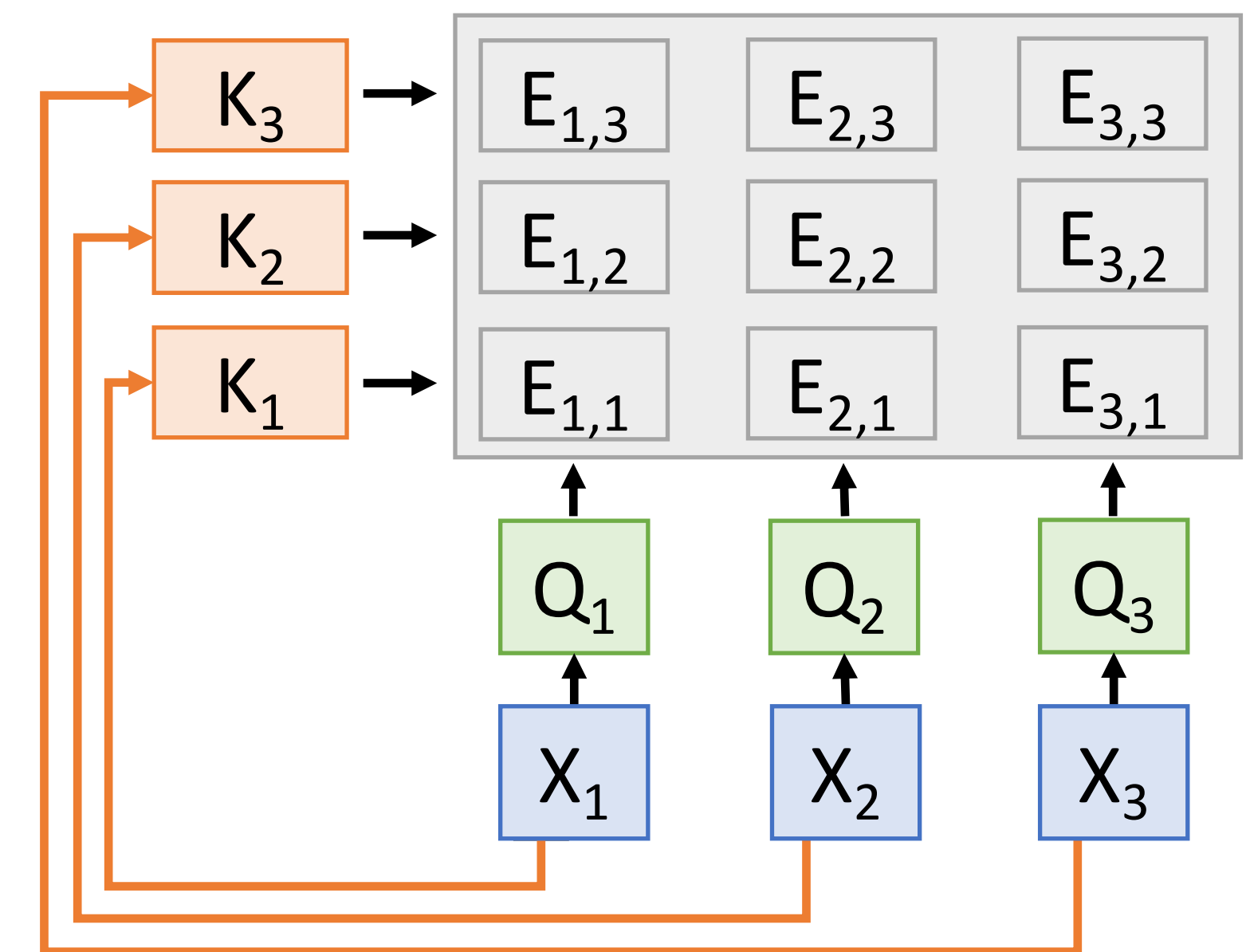**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value Vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_X \times N_X$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

# **Self-attention** Layer

One **query** per **input vector**

**Inputs**:
**Input vectors**: **X** (Shape: $N_X$ x $D_X$)
**Key matrix**: **$W_K$** (Shape: $D_X$ x $D_Q$)
**Value matrix**: **$W_V$** (Shape: $D_X$ x $D_V$)
**Query matrix**: **$W_Q$** (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: **Q** = **X$W_Q$**
**Key vectors**: **K** = **X$W_K$**  (Shape: $N_X$ x $D_Q$)
**Value Vectors**: **V** = **X$W_V$** (Shape: $N_X$ x $D_V$)
**Similarities**: E = **Q$K^T$** $/\sqrt{D_Q}$ (Shape: $N_X$ x $N_X$) $E_{i,j} = (\mathbf{Q_i} \cdot \mathbf{K_j}) /\sqrt{D_Q}$
**Attention weights**: A = softmax(E, dim=1)  (Shape: $N_X$ x $N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V_j}$

# **Self-attention** Layer

One **query** per **input vector**

**Inputs**:
**Input vectors**: **X** (Shape: $N_X \times D_X$)
**Key matrix**: **$W_K$** (Shape: $D_X \times D_Q$)
**Value matrix**: **$W_V$** (Shape: $D_X \times D_V$)
**Query matrix**: **$W_Q$** (Shape: $D_X \times D_Q$)

**Computation**:
**Query vectors**: **Q** = **X$W_Q$**
**Key vectors**: **K** = **X$W_K$** (Shape: $N_X \times D_Q$)
**Value Vectors**: **V** = **X$W_V$** (Shape: $N_X \times D_V$)
**Similarities**: E = **Q$K^T$** $/\sqrt{D_Q}$ (Shape: $N_X \times N_X$) $E_{i,j} = (Q_i \cdot K_j)/\sqrt{D_Q}$
**Attention weights**: A = softmax(E, dim=1) (Shape: $N_X \times N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

# **Self-attention** Layer

One **query** per **input vector**

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)
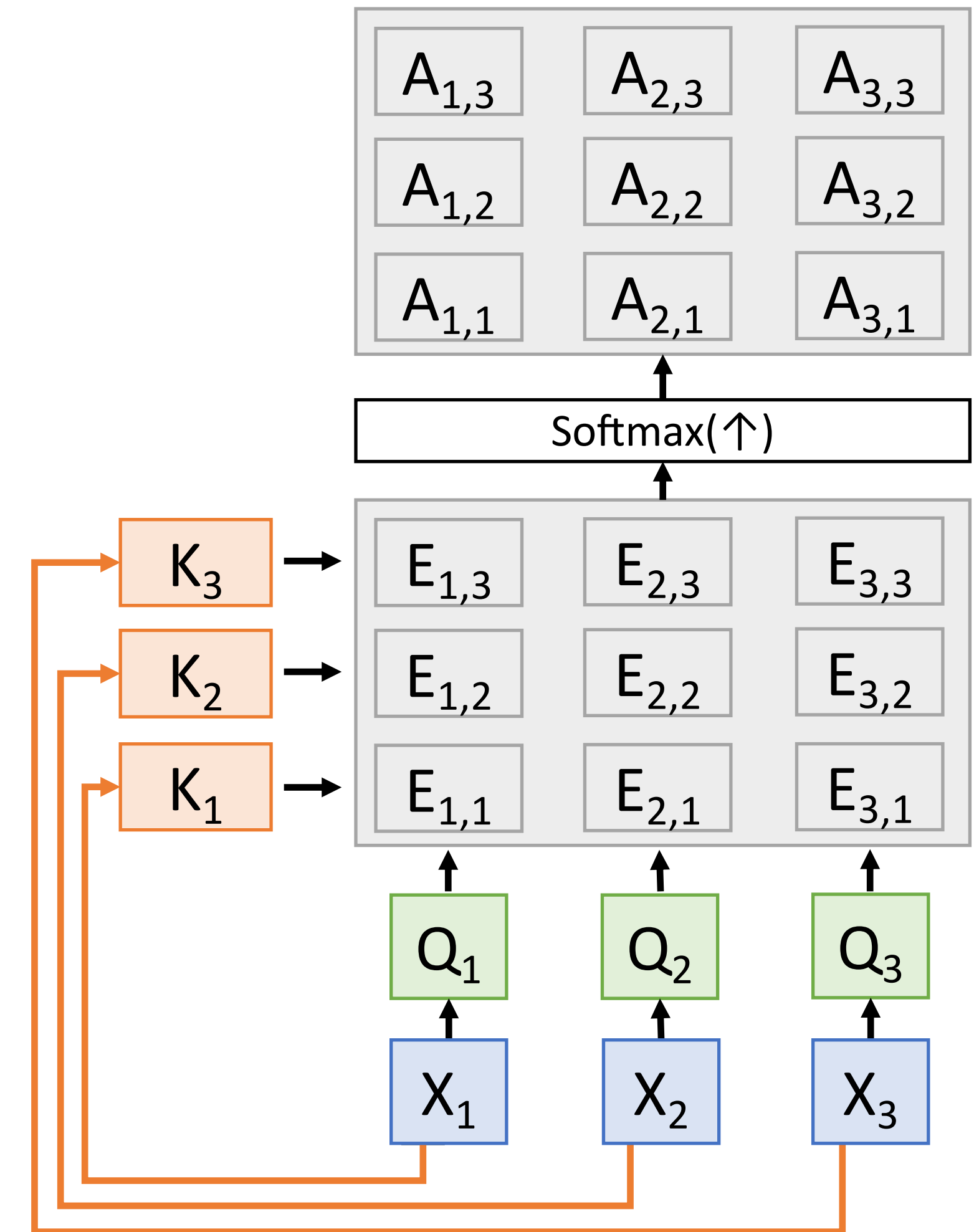
**Computation**:
**Query vectors**: $Q = XW_Q$
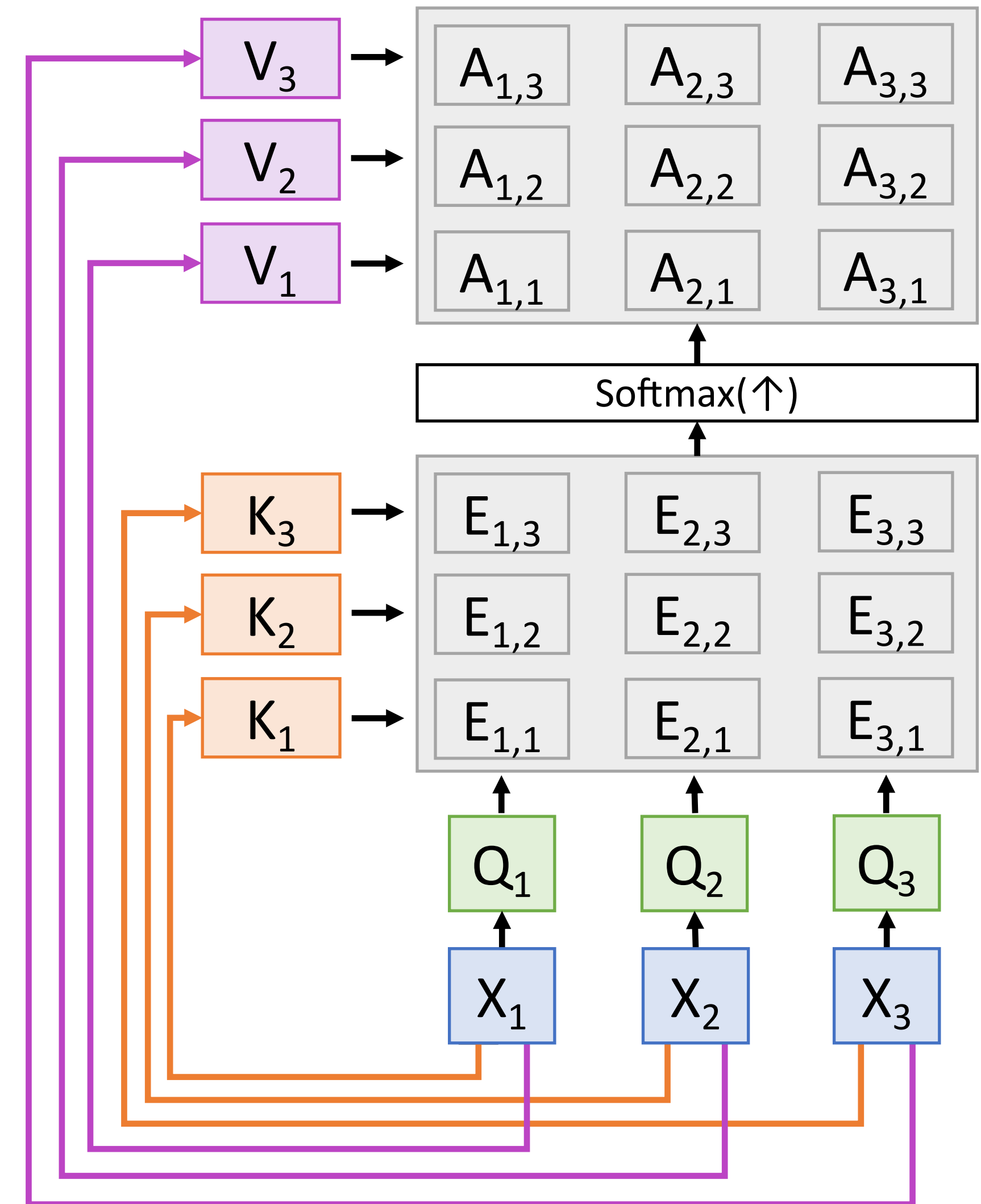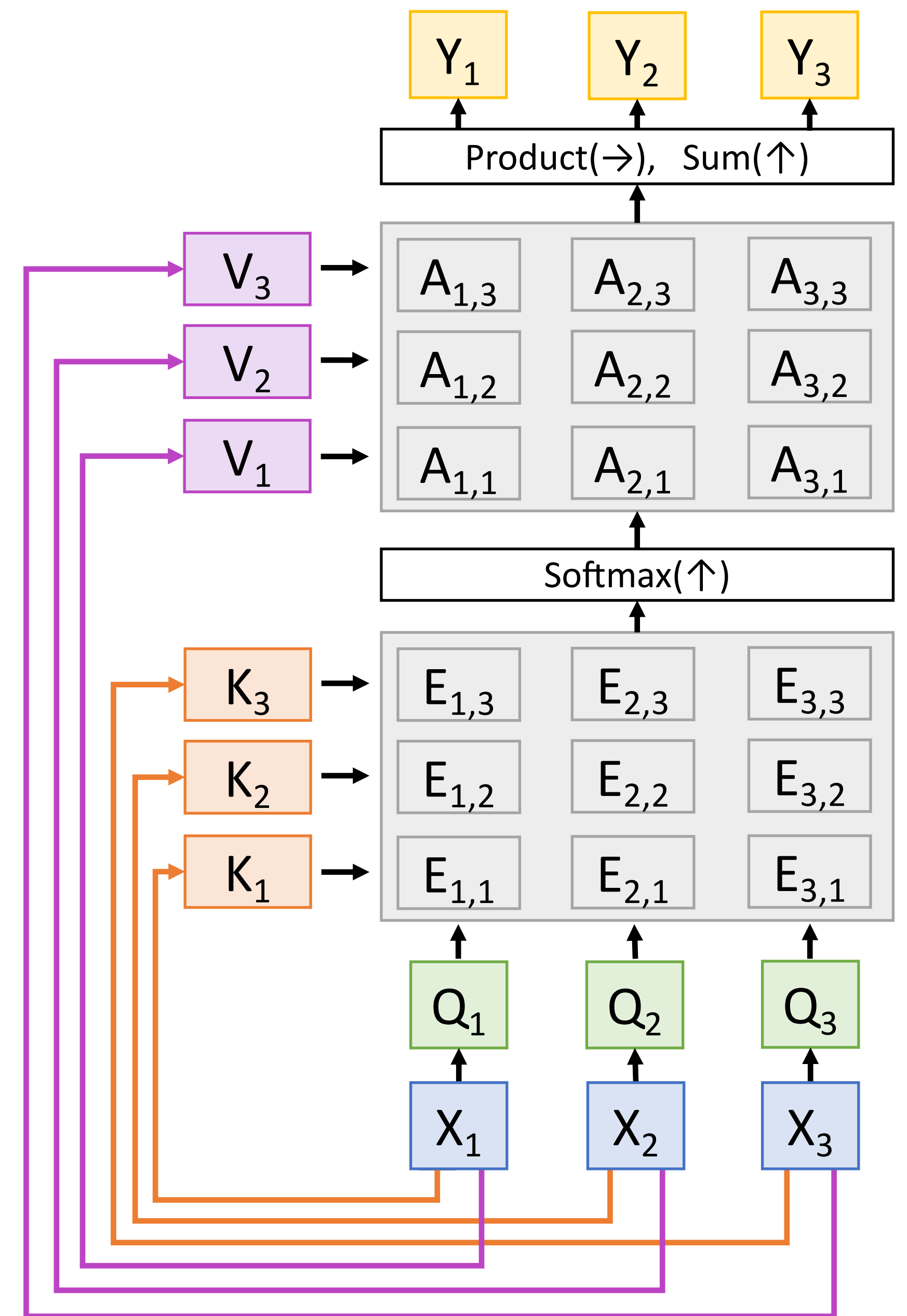**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value Vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_X \times N_X$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$
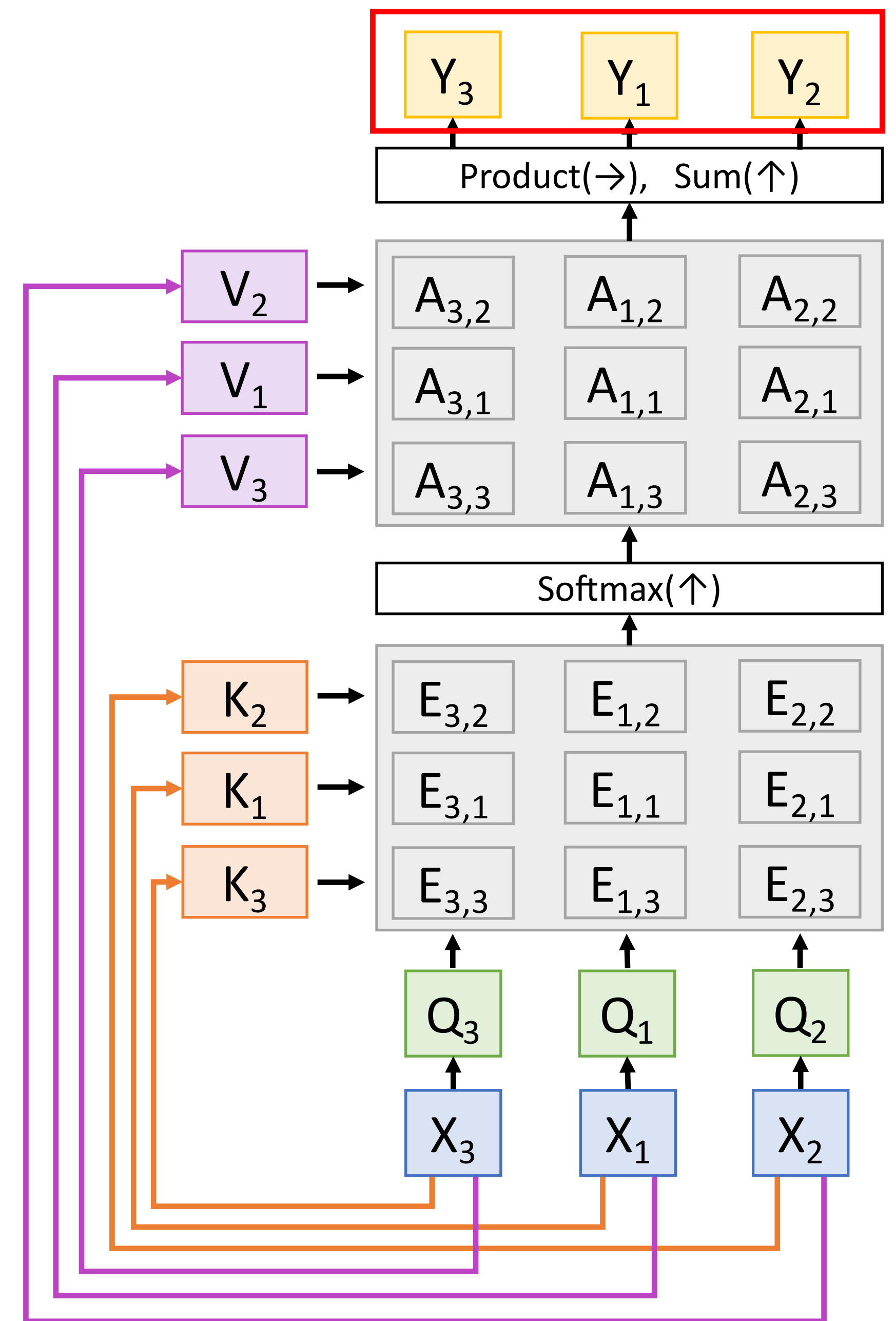**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$



[ slide from Justin Johnson, U Michigan ]

# **Self-attention** Layer

**Inputs**:
**Input vectors**: $\mathbf{X}$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $\mathbf{W_K}$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $\mathbf{W_V}$ (Shape: $D_X$ x $D_V$)
**Query matrix**: $\mathbf{W_Q}$ (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: $\mathbf{Q} = \mathbf{X}\mathbf{W_Q}$
**Key vectors**: $\mathbf{K} = \mathbf{X}\mathbf{W_K}$ (Shape: $N_X$ x $D_Q$)
**Value Vectors**: $\mathbf{V} = \mathbf{X}\mathbf{W_V}$ (Shape: $N_X$ x $D_V$)
**Similarities**: $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$ (Shape: $N_X$ x $N_X$) $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X$ x $N_X$)
**Output vectors**: $Y = A\mathbf{V}$ (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j}\mathbf{V}_j$

Consider **permuting** the input vectors:

Outputs will be the same, but permuted

Self-attention layer is **Permutation Equivariant** $f(s(x)) = s(f(x))$

Self-Attention layer works on **sets** of vectors



[ slide from Justin Johnson, U Michigan ]

# **Multi-head** Self-attention Layer

**Inputs**:

**Input vectors**: $X$ (Shape: $N_X \times D_X$)

**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)

**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)

**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)

Use H independent "Attention Heads" in parallel

**Computation**:

**Query vectors**: $Q = XW_Q$

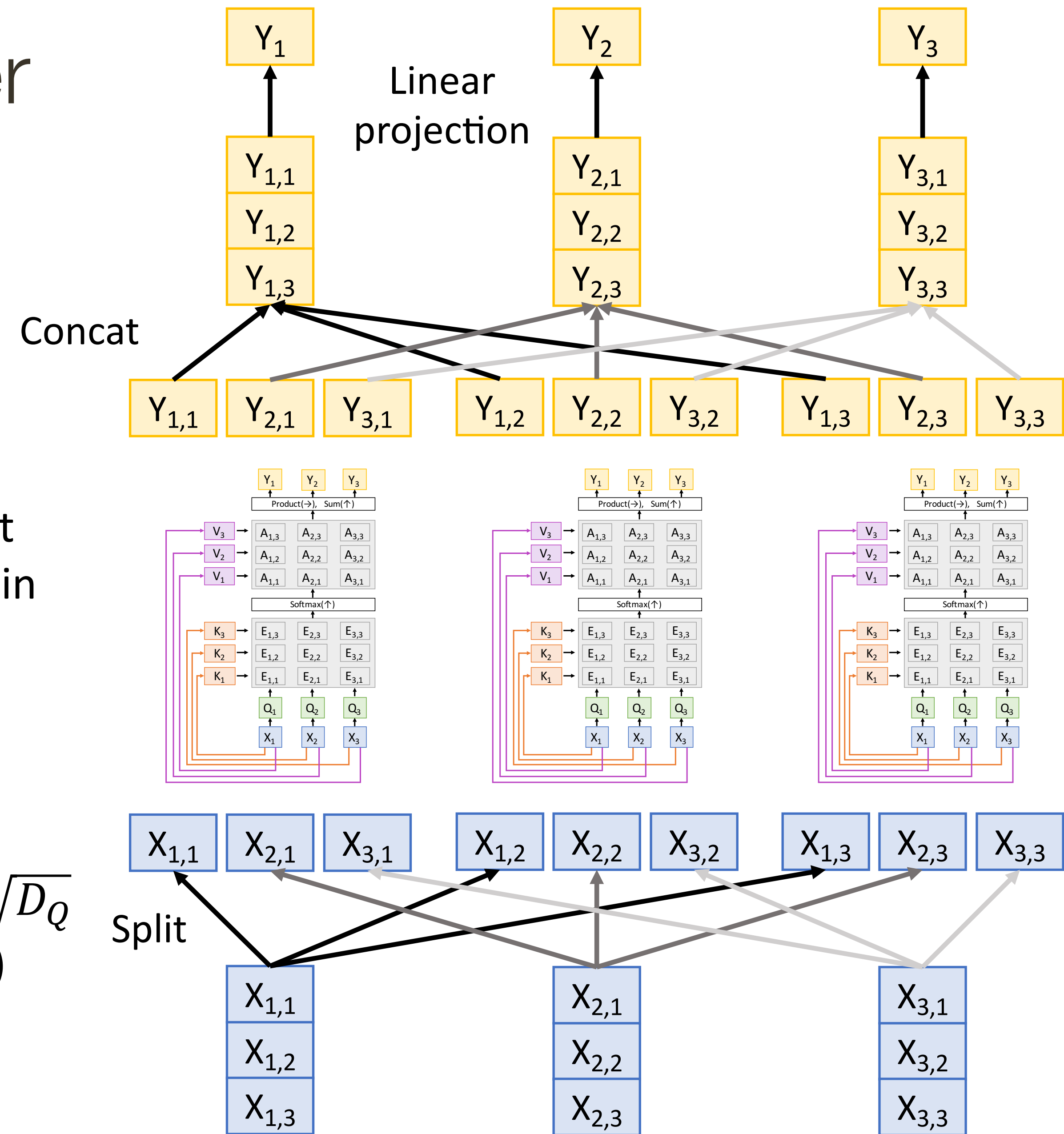**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)

**Value Vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)

**Similarities**: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_X \times N_X$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)

**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$



[ slide from Justin Johnson, U Michigan ]

# CNN with **Self-attention**

Input Image

CNN

Features:
C x H x W

Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

[ slide from Justin Johnson, U Michigan ]

# CNN with Self-attention

Input Image

CNN

Features:
C x H x W

**Queries**:
C' x H x W

1x1 Conv

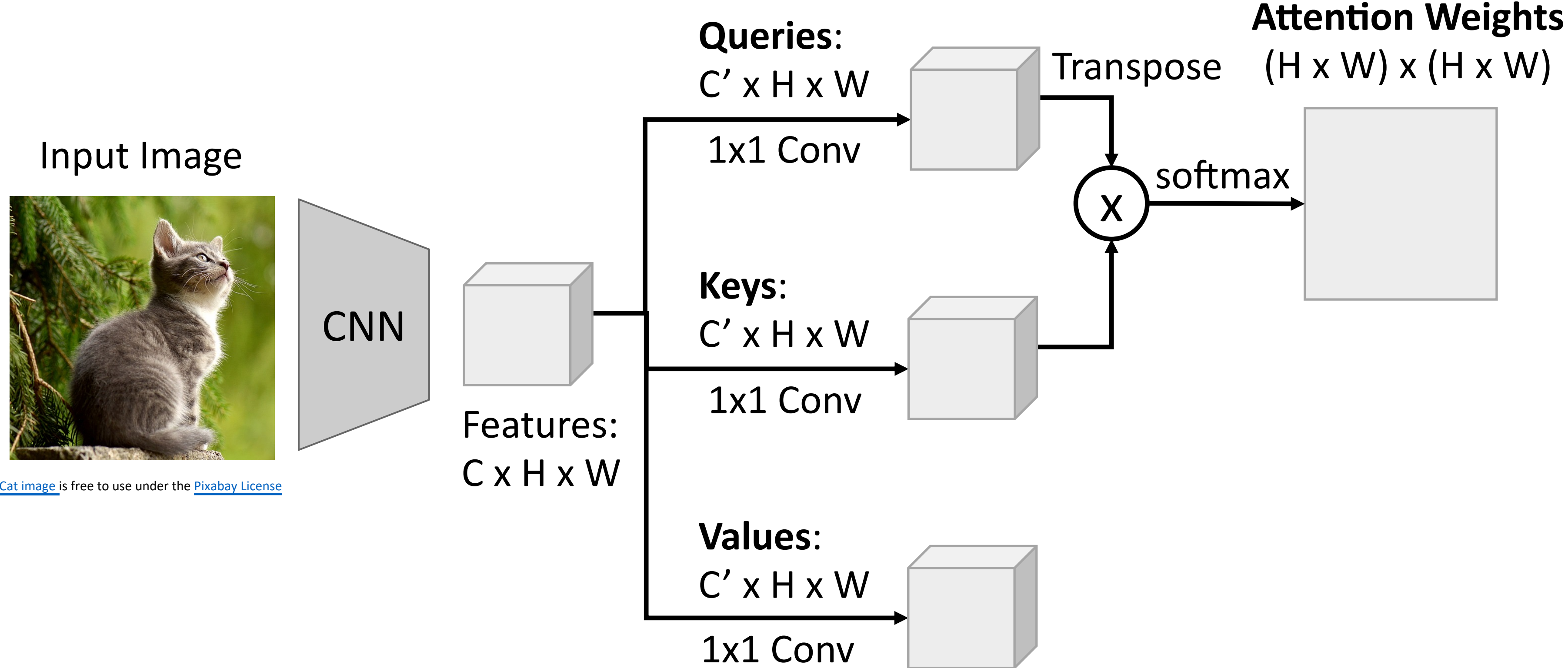**Keys**:
C' x H x W

1x1 Conv

**Values**:
C' x H x W

1x1 Conv

Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

[ slide from Justin Johnson, U Michigan ]

# CNN with Self-attention

Input Image



Cat image is free to use under the Pixabay License

CNN

Features:
C x H x W

**Queries**:
C' x H x W

1x1 Conv

**Keys**:
C' x H x W

1x1 Conv

**Values**:
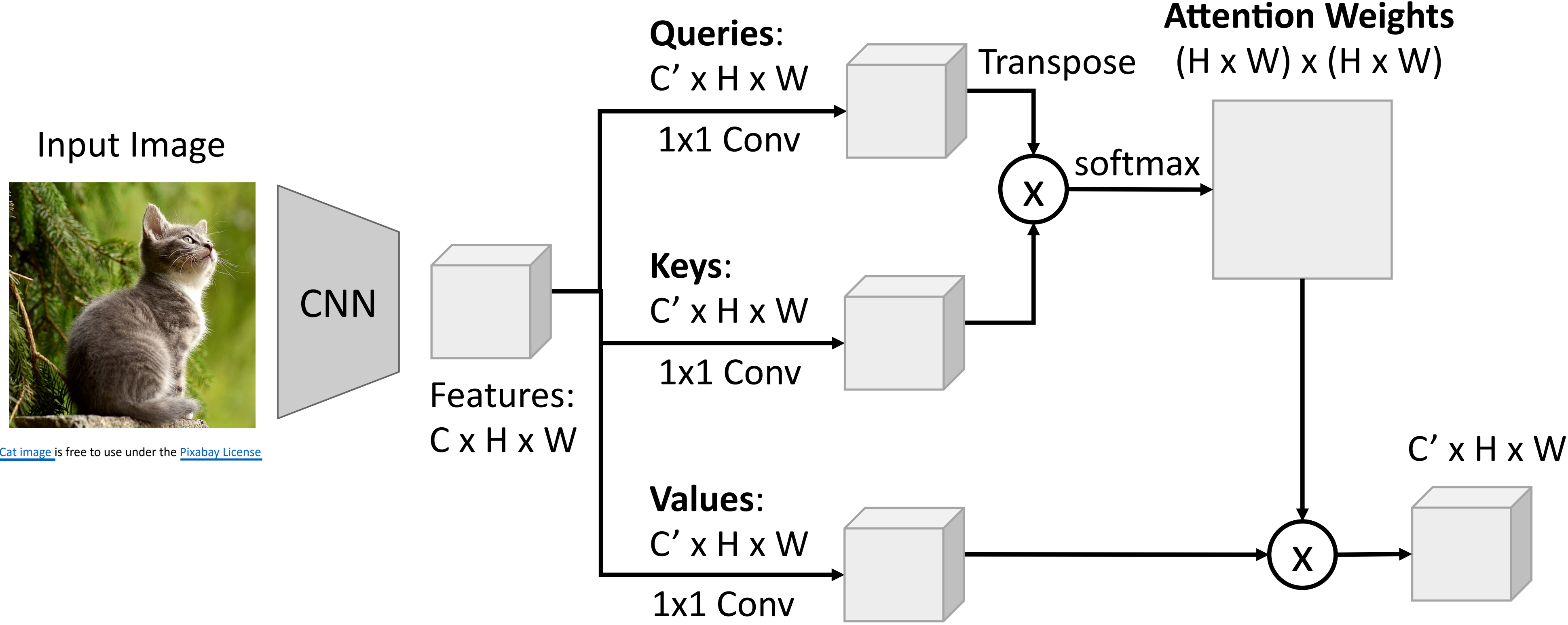C' x H x W

1x1 Conv

Transpose

X   softmax

**Attention Weights**
(H x W) x (H x W)

Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

[ slide from Justin Johnson, U Michigan ]

# CNN with Self-attention

Input Image

CNN

Features:
C x H x W

**Queries**:
C' x H x W

1x1 Conv

**Keys**:
C' x H x W

1x1 Conv

**Values**:
C' x H x W

1x1 Conv

Transpose

X softmax

**Attention Weights**
(H x W) x (H x W)

X

C' x H x W

Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

[ slide from Justin Johnson, U Michigan ]

# CNN with Self-attention

Input Image

CNN

Features:
C x H x W

**Queries**:
C' x H x W

1x1 Conv

**Keys**:
C' x H x W

1x1 Conv

**Values**:
C' x H x W

1x1 Conv

Transpose

X

softmax

**Attention Weights**
(H x W) x (H x W)
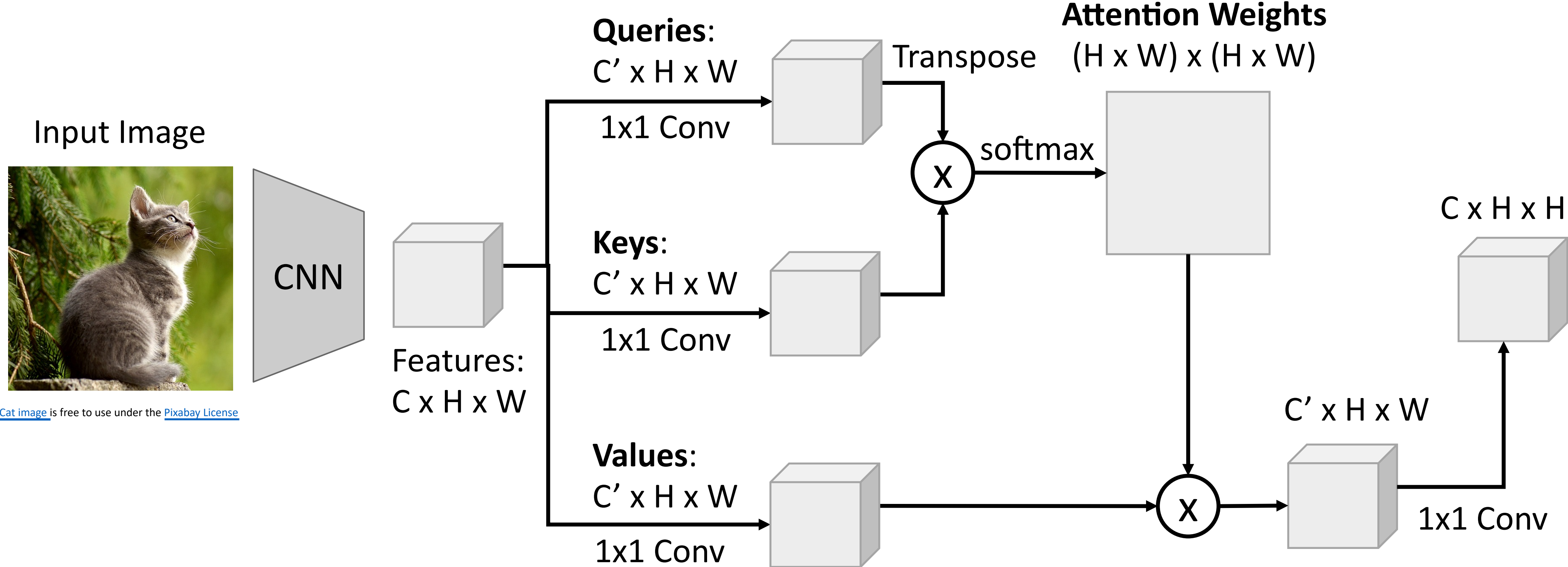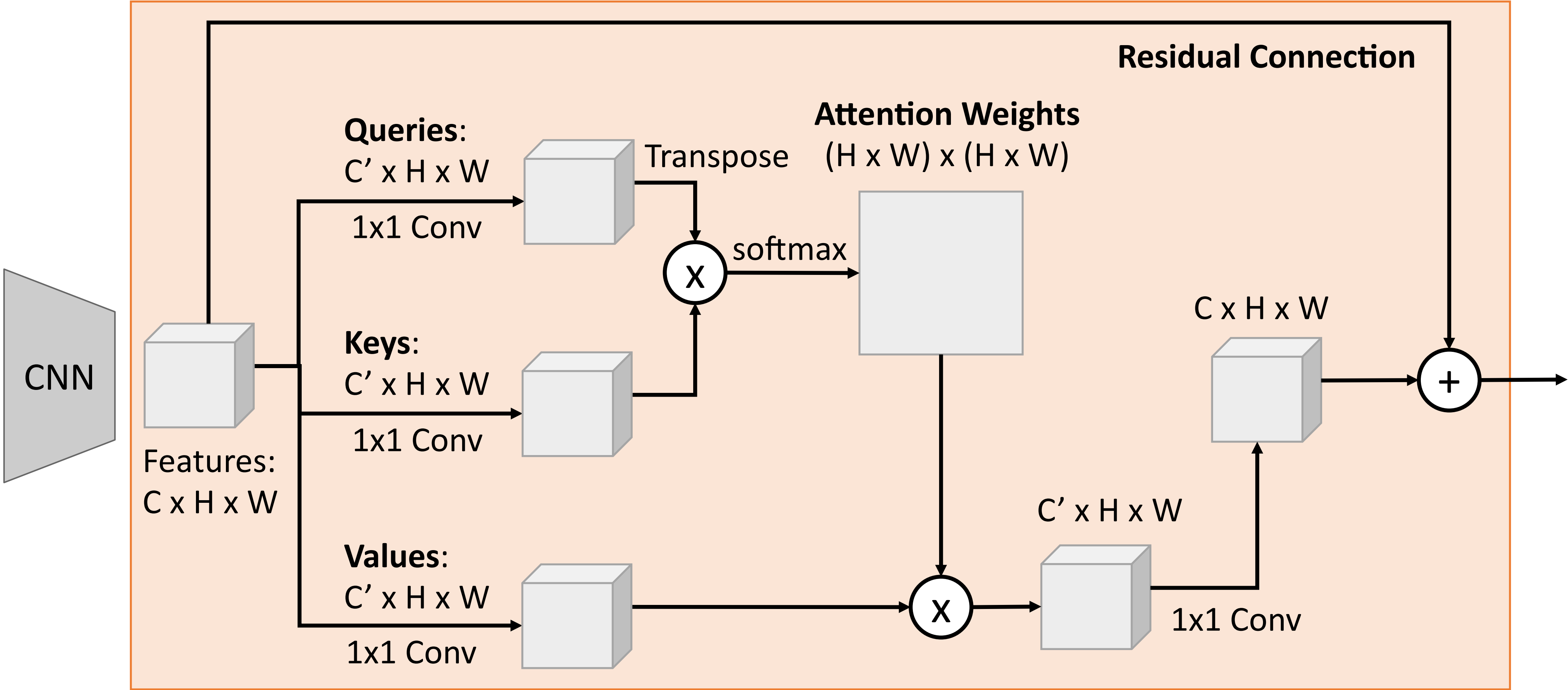
C' x H x W

X

1x1 Conv

C x H x H

Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018
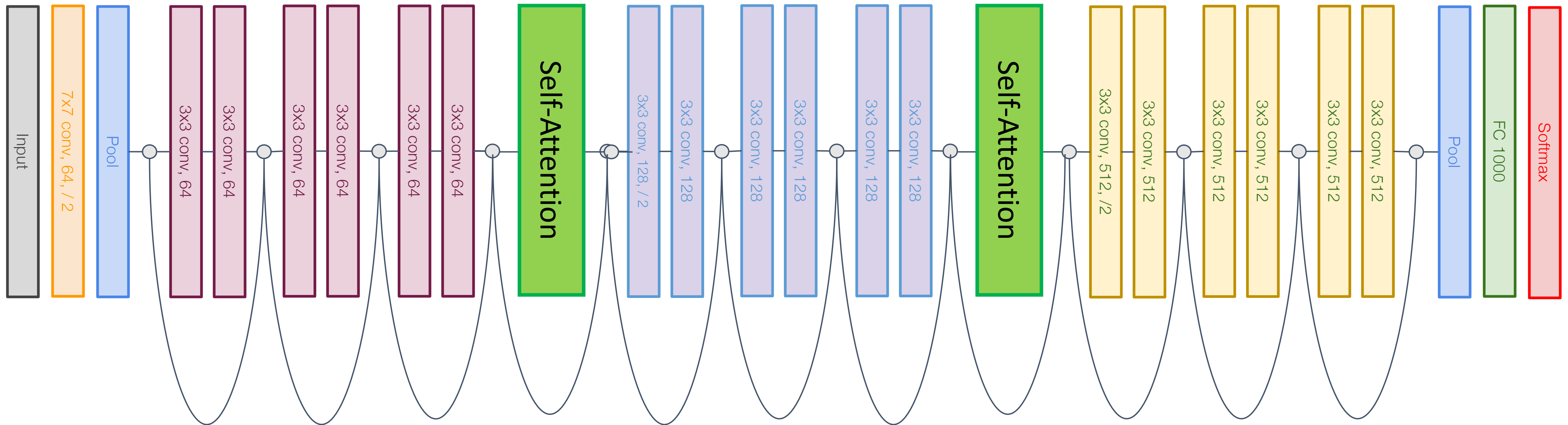
[ slide from Justin Johnson, U Michigan ]

# CNN with Self-attention



Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

[ slide from Justin Johnson, U Michigan ]

# **Attention** with Existing CNNs



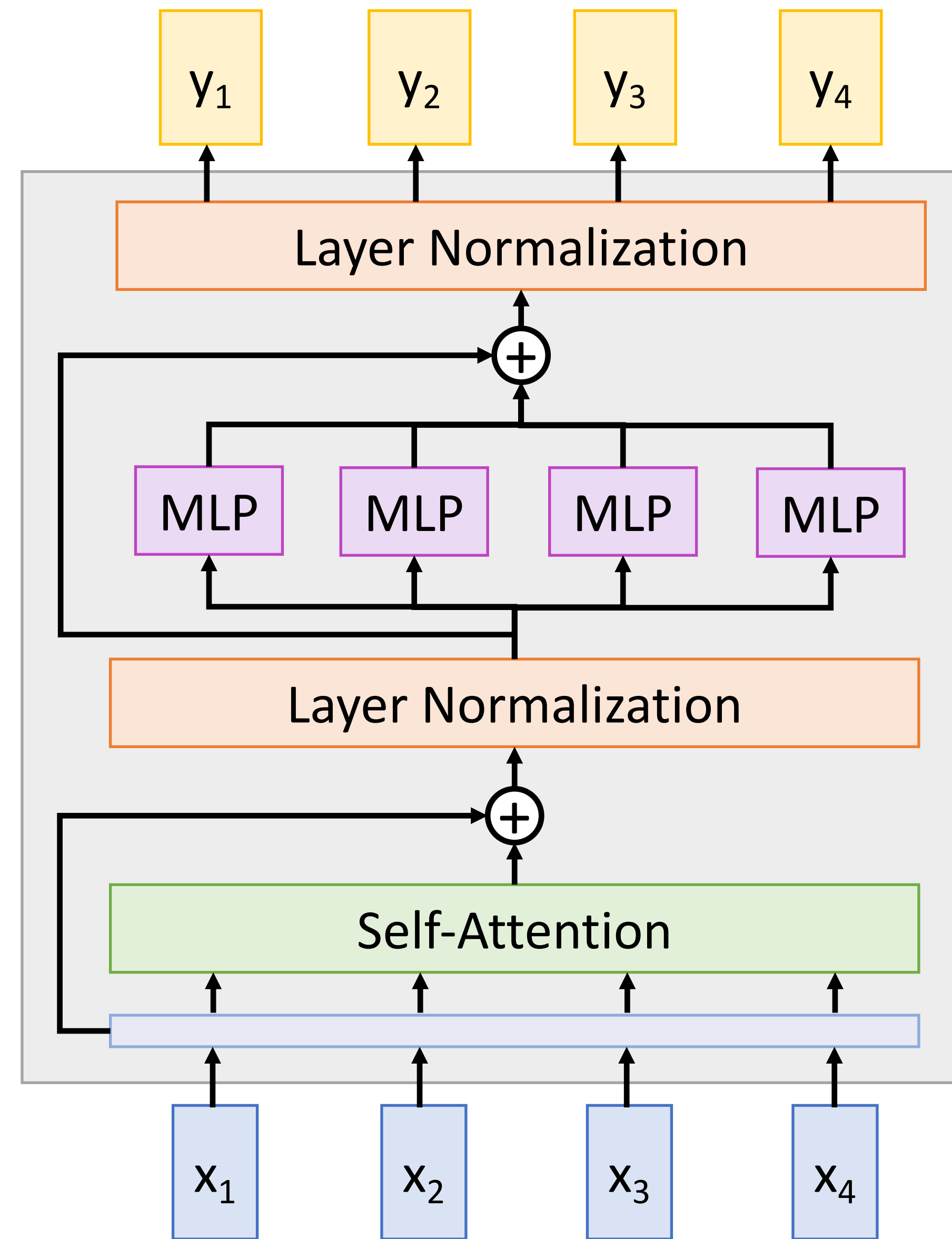Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018
Wang et al, "Non-local Neural Networks", CVPR 2018

[ slide from Justin Johnson, U Michigan ]

# Transformer

Transfomer block inputs a set of vectors, outputs a set of vectors.

Vectors only communicate via (multiheaded) self-attention



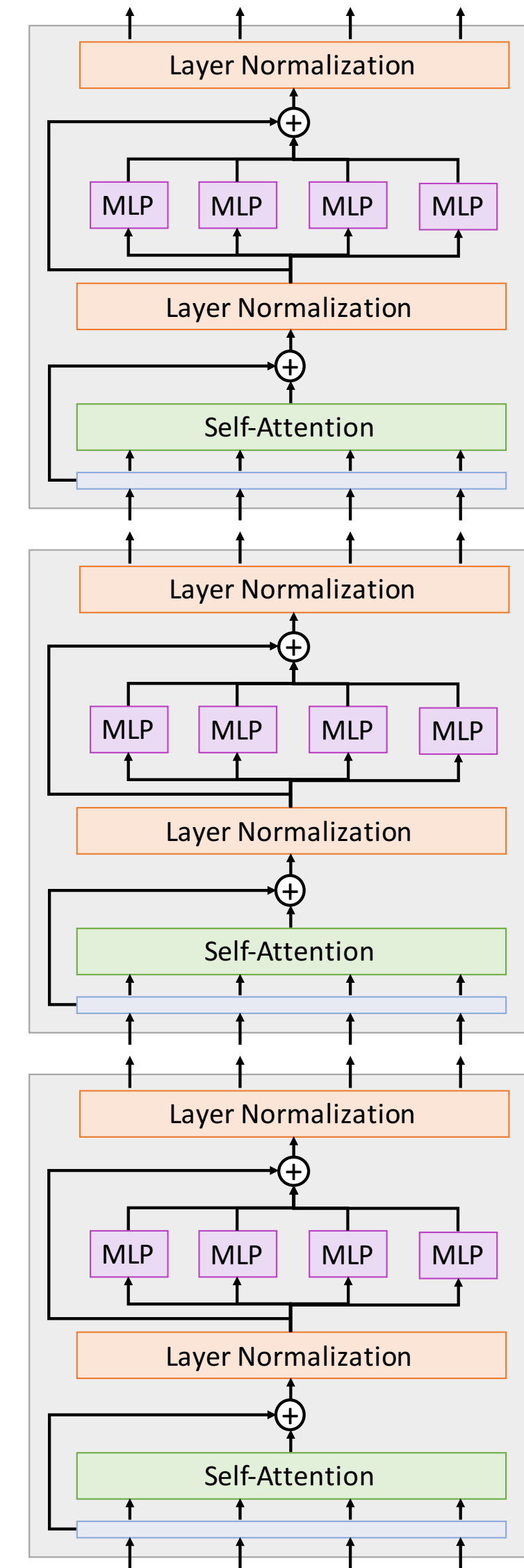Vaswani et al, "Attention is all you need", NeurIPS 2017

# Transformer

**Transformer Block:**

**Input**: Set of vectors x

**Output**: Set of vectors y

Hyperparameters:
- Number of blocks
- Number of heads per block
- Width (channels per head, FFN width)



Vaswani et al, "Attention is all you need", NeurIPS 2017

[ slide from Justin Johnson, U Michigan ]

# **Transformer** on Image Patches

N input patches, each
of shape 3x16x16



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021
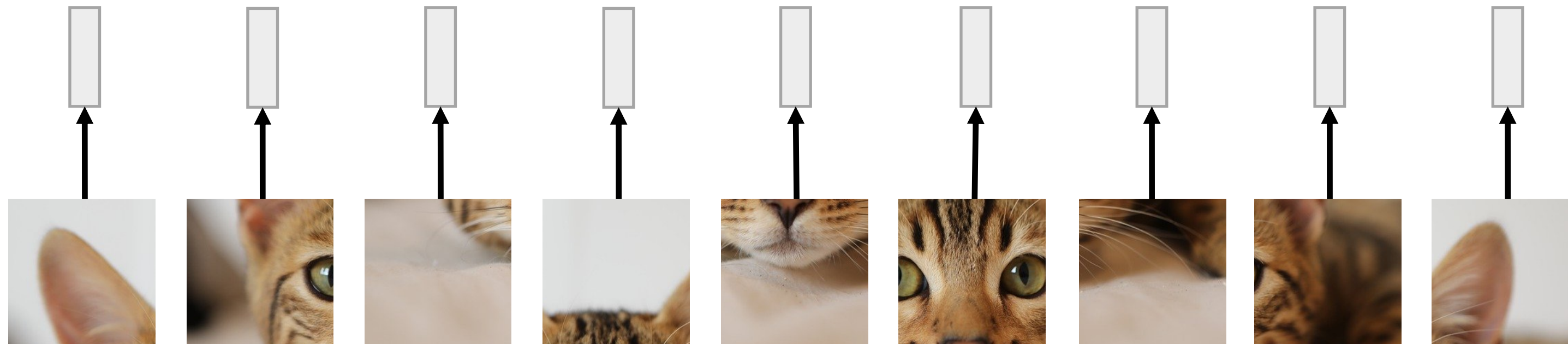
[ slide from Justin Johnson, U Michigan ]

# **Transformer** on Image Patches

Linear projection to
D-dimensional vector

N input patches, each
of shape 3x16x16

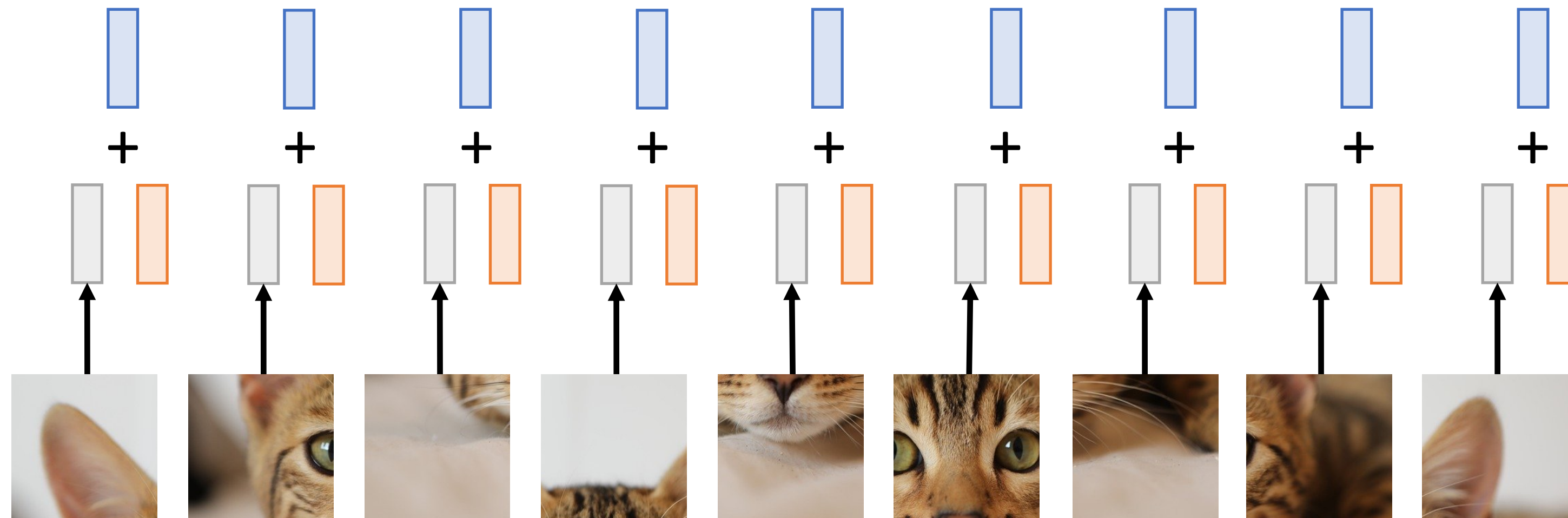Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[ slide from Justin Johnson, U Michigan ]

# **Transformer** on Image Patches

Add positional embedding: learned D-dim vector per position

Linear projection to D-dimensional vector

N input patches, each of shape 3x16x16

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[ slide from Justin Johnson, U Michigan ]

# **Transformer** on Image Patches



Output vectors

Exact same as
NLP Transformer!

Add positional
embedding: learned D-
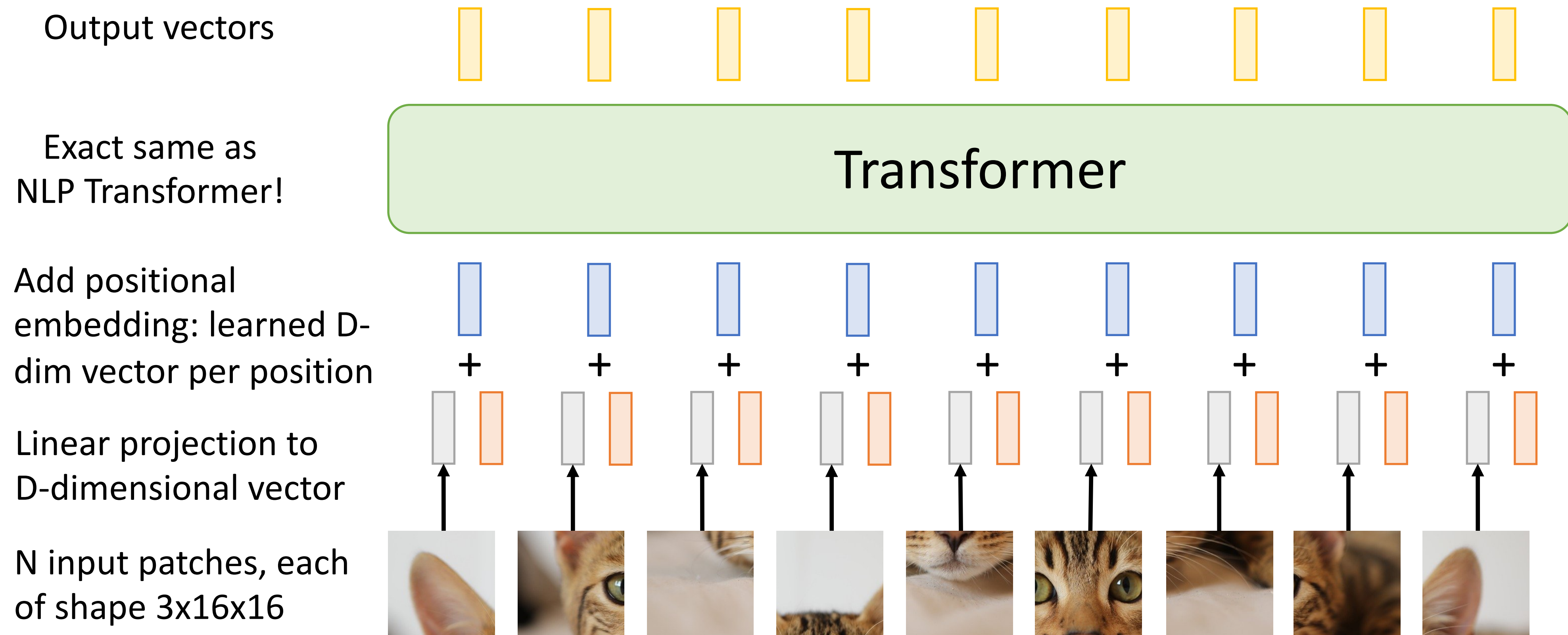dim vector per position

Linear projection to
D-dimensional vector

N input patches, each
of shape 3x16x16

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[ slide from Justin Johnson, U Michigan ]

# **Transformer** on Image Patches



Output vectors

Exact same as NLP Transformer!

Transformer

Add positional embedding: learned D-dim vector per position

Special extra input: **classification token** (D dims, learned)

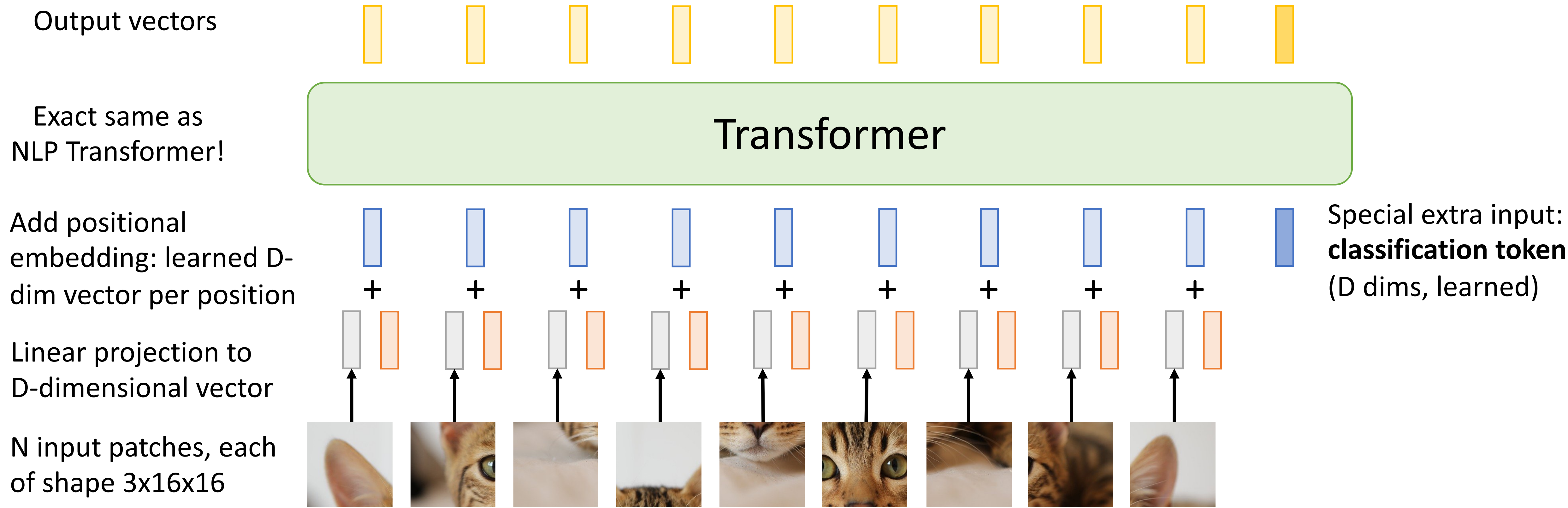Linear projection to D-dimensional vector

N input patches, each of shape 3x16x16

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

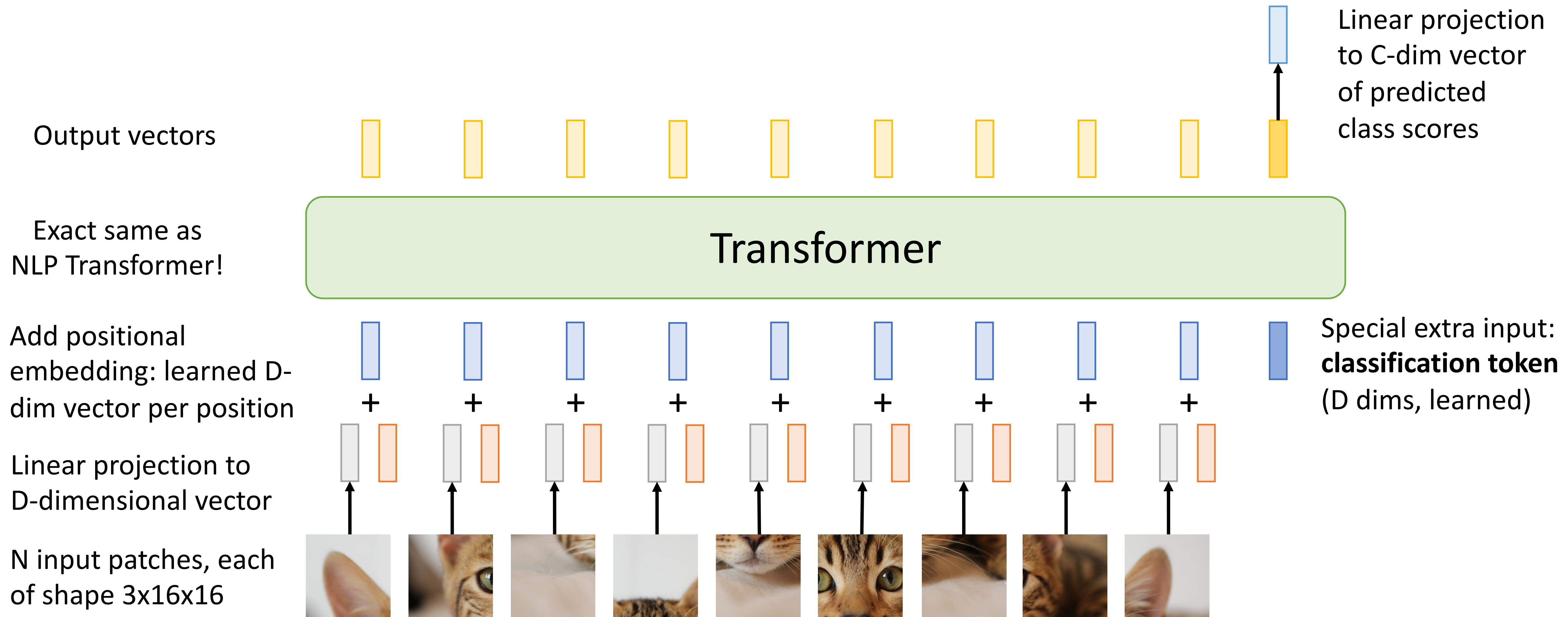Cat image is free for commercial use under a Pixabay license

[ slide from Justin Johnson, U Michigan ]

# **Transformer** on Image Patches

Output vectors

Exact same as
NLP Transformer!

Add positional
embedding: learned D-
dim vector per position

Linear projection to
D-dimensional vector

N input patches, each
of shape 3x16x16



Linear projection
to C-dim vector
of predicted
class scores

Transformer

Special extra input:
**classification token**
(D dims, learned)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[ slide from Justin Johnson, U Michigan ]

# **Vision Transformer** (ViT)

Computer vision model
with no convolutions!

Linear projection
to C-dim vector
of predicted
class scores

Output vectors

Transformer
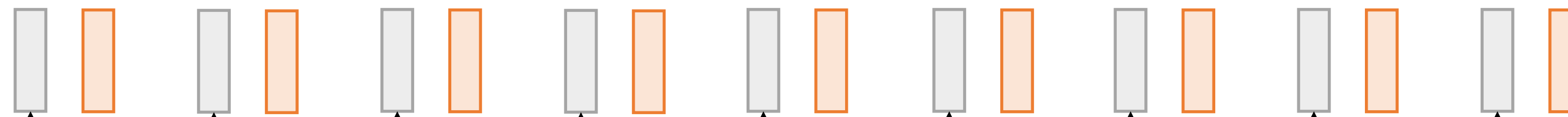
Exact same as
NLP Transformer!

Add positional
embedding: learned D-
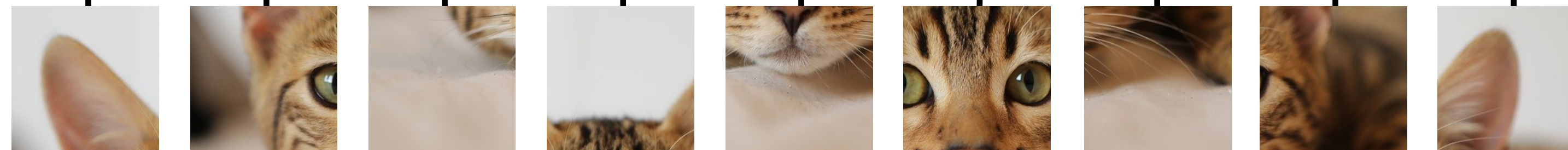dim vector per position

+ + + + + + + + +

Special extra input:
**classification token**
(D dims, learned)

Linear projection to
D-dimensional vector

N input patches, each
of shape 3x16x16

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021
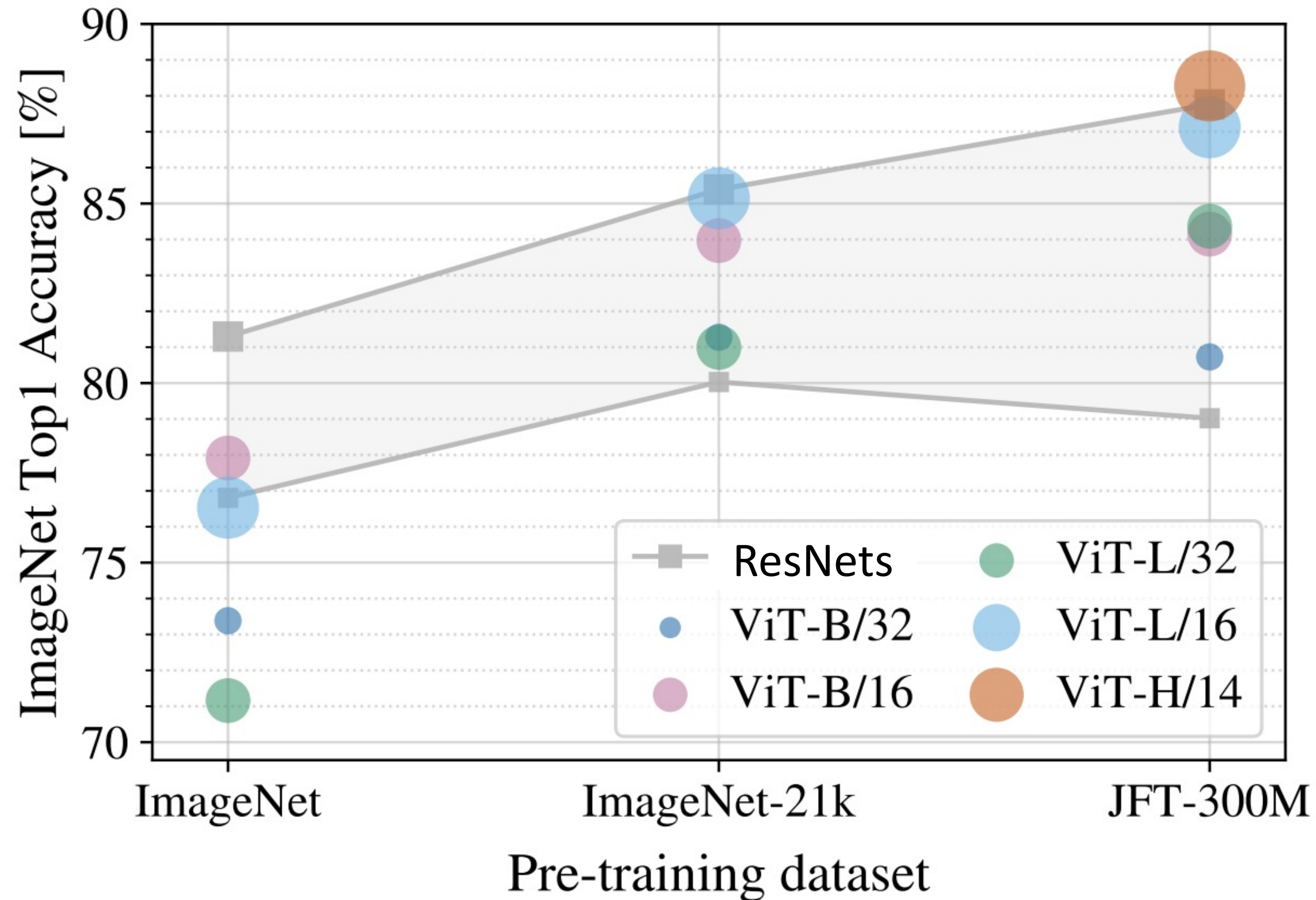
Cat image is free for commercial
use under a Pixabay license

[ slide from Justin Johnson, U Michigan ]

# Vision Transformer (ViT) vs. ResNet

JFT-300M is an internal Google dataset with 300M labeled images

If you pretrain on JFT and finetune on ImageNet, large ViTs outperform large ResNets
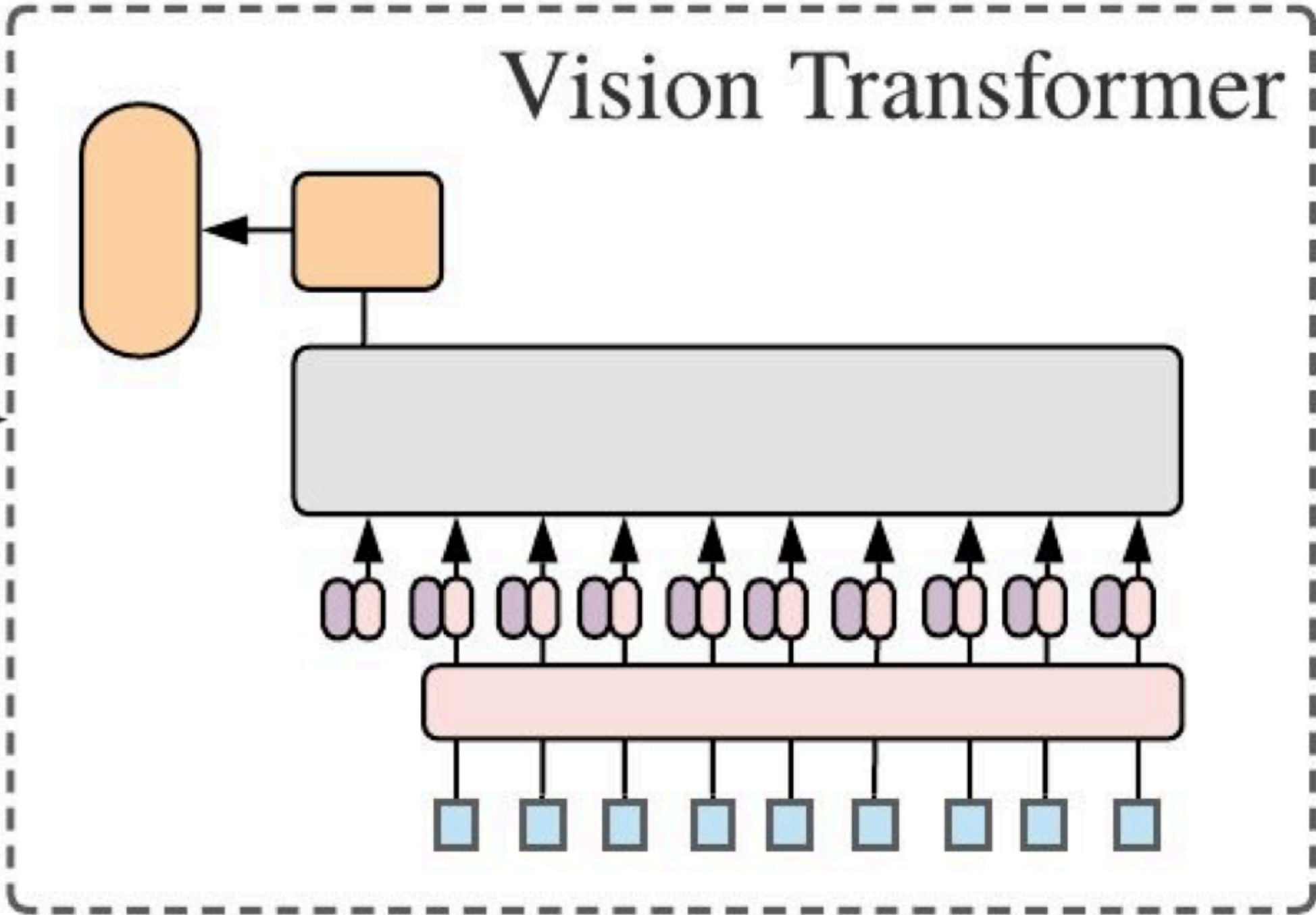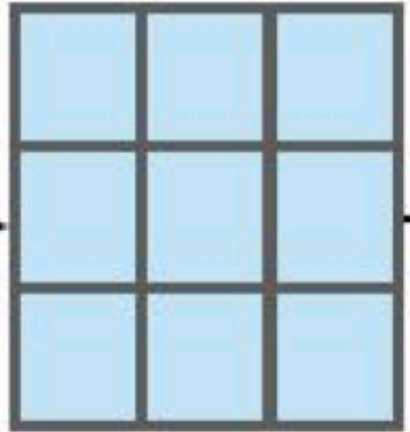


B = Base
L = Large
H = Huge

/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021
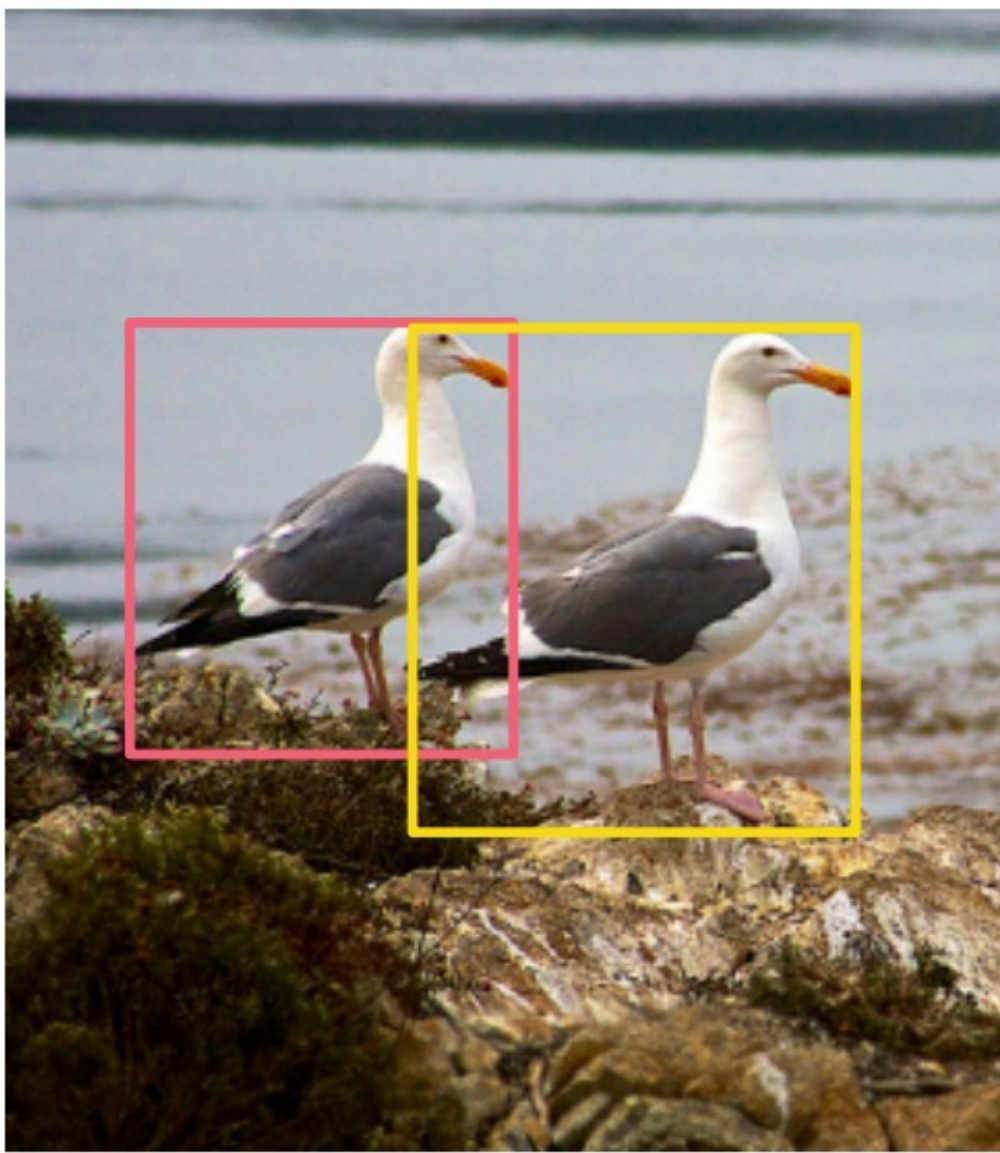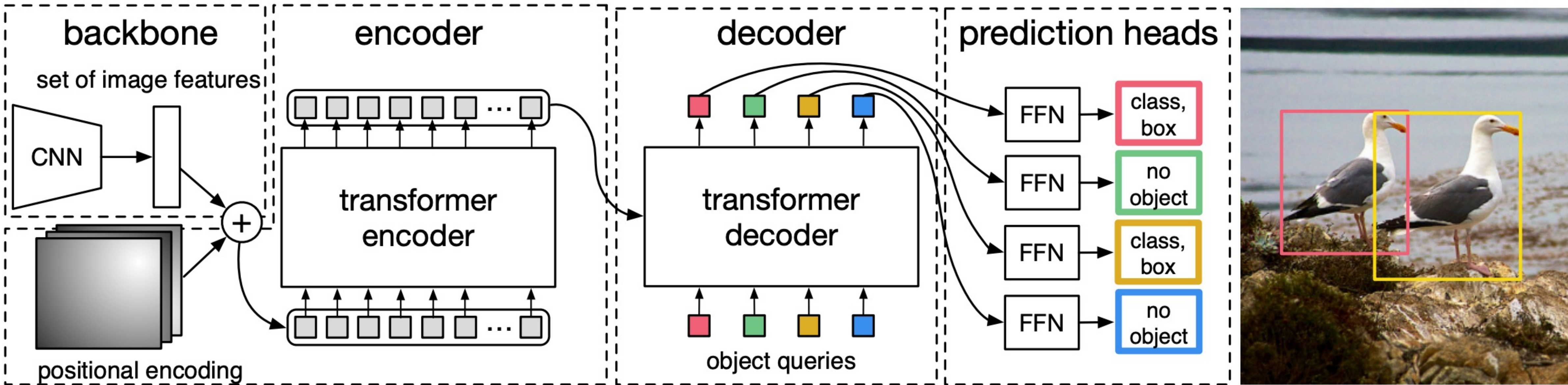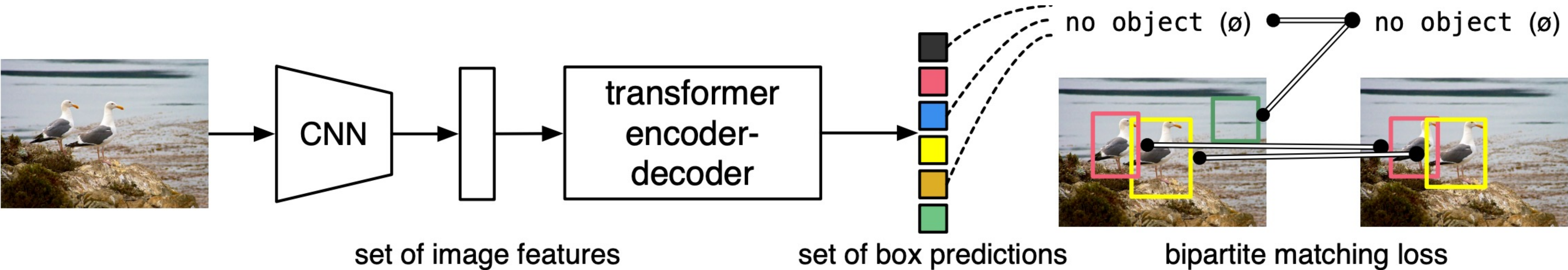
[ slide from Justin Johnson, U Michigan ]

# **ResNet-ViT** Hybrid
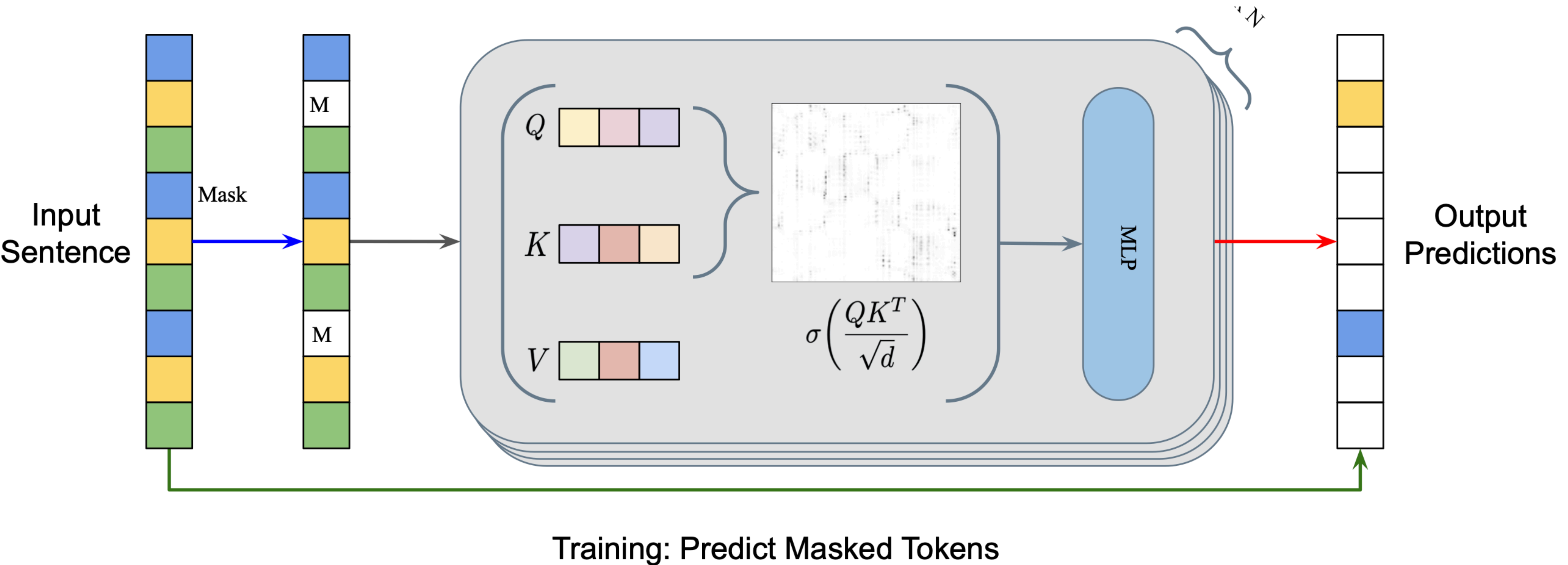
# Object Detection with Transformers: DETR Model



set of image features

set of box predictions

bipartite matching loss

no object (ø)     no object (ø)

backbone

set of image features

CNN

positional encoding

encoder

transformer encoder

decoder

transformer decoder

object queries

prediction heads

FFN → class, box

FFN → no object

FFN → class, box

FFN → no object

Carion et al, "End-to-End Object Detection with Transformers", ECCV 2020

# **Masked Modeling** with Transformers (BERT, GPT, etc.)



Training: Predict Masked Tokens

$$\mathcal{L}_{\mathrm{MLM}}(X;\theta) = \underset{x \sim X}{\mathbb{E}} \underset{\mathrm{mask}}{\mathbb{E}} \sum_{i \in \mathrm{mask}} \log p(x_i | x_{j \notin \mathrm{mask}}; \theta)$$

(mask 15% at a time)

# **Masked** Self-Attention Layer

Don't let vectors "look ahead" in the sequence
Used for language modeling (predict next word)

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X$ x $D_Q$)
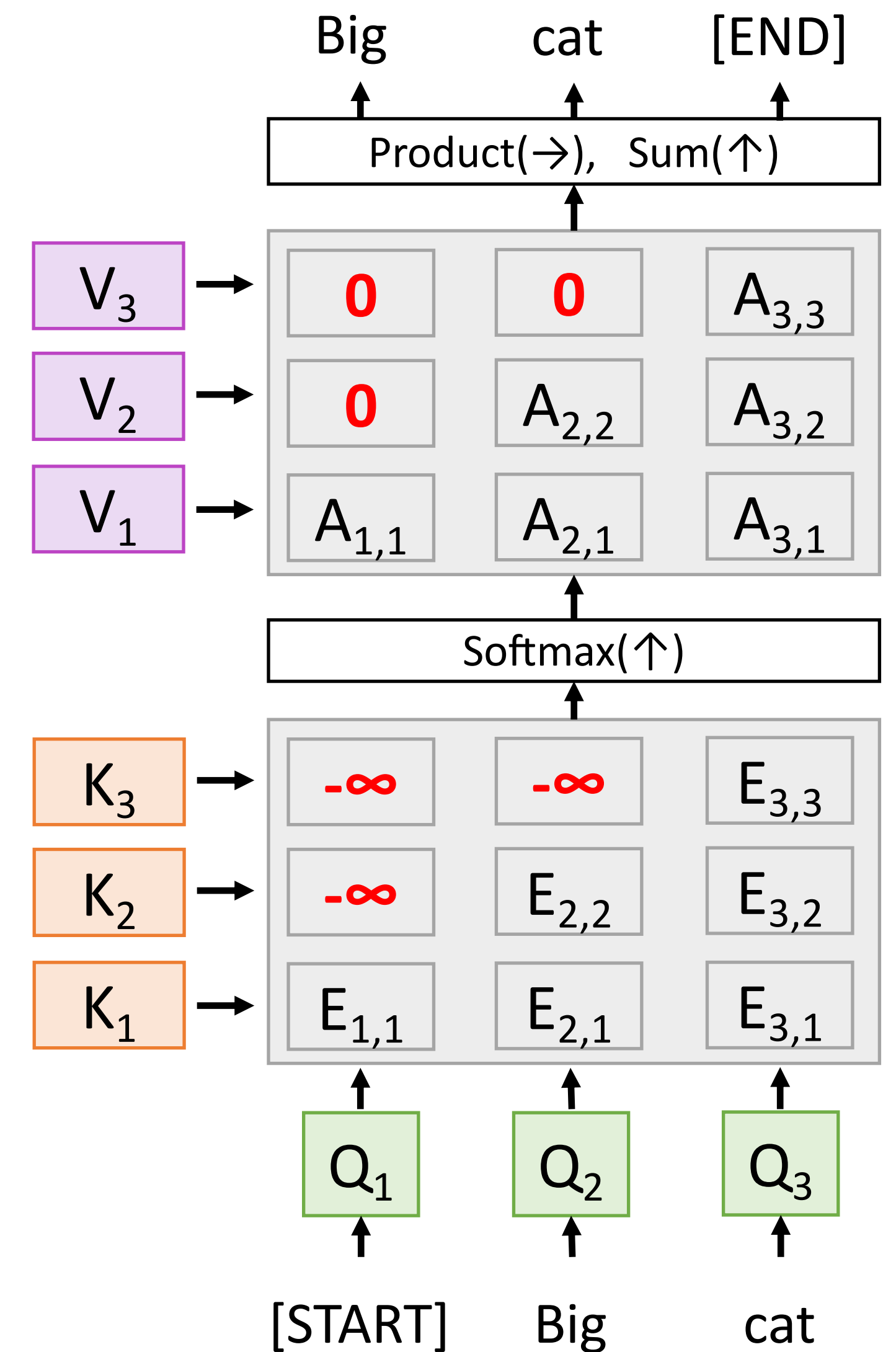
**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X$ x $D_Q$)
**Value Vectors**: $V = XW_V$ (Shape: $N_X$ x $D_V$)
**Similarities**: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_X$ x $N_X$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X$ x $N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$

Please fill out
**Student Evaluations**
(on Canvas)