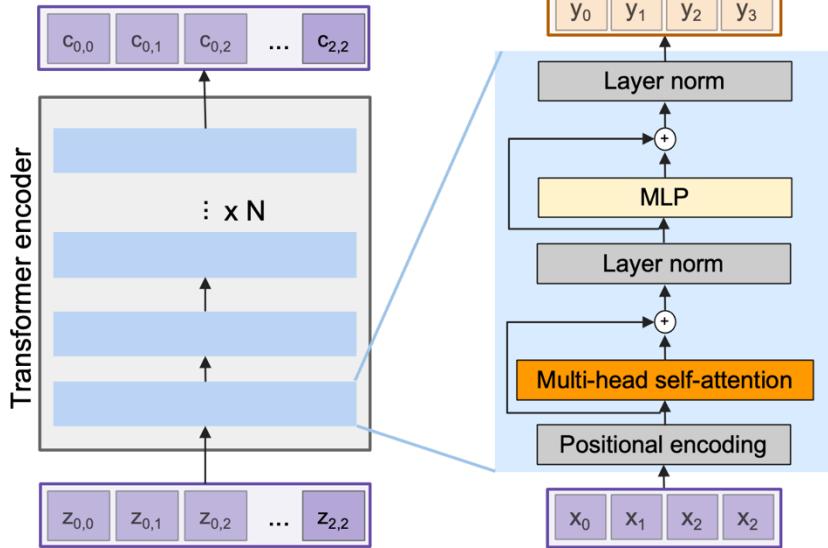
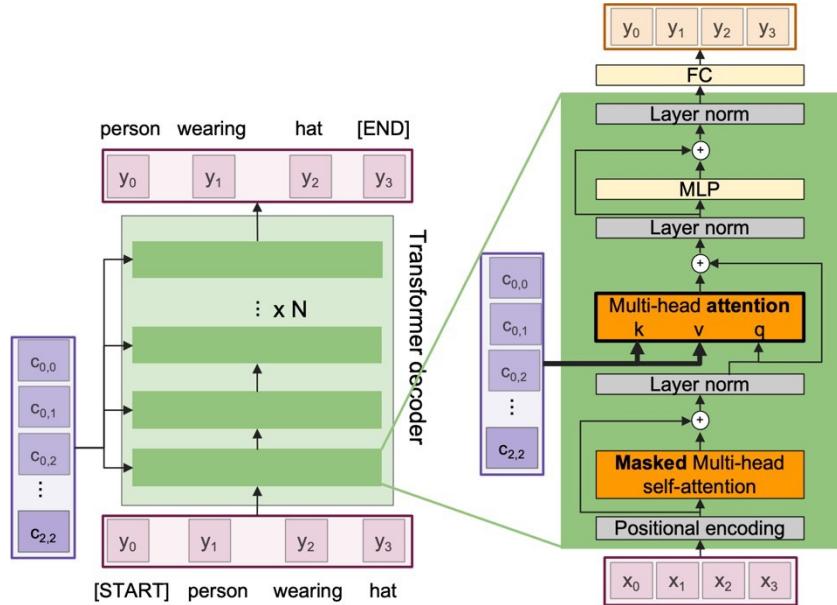


# Lecture 9: Detection, Segmentation, Visualization, and Understanding

# Last time: Transformer



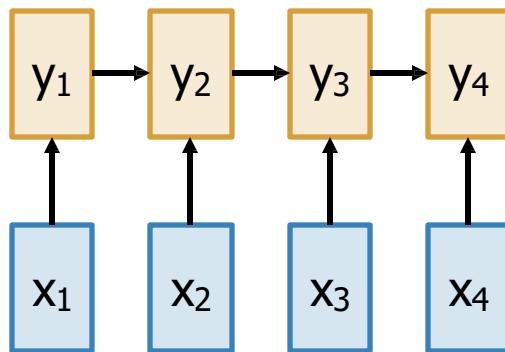
Encoder



Decoder

# Three Ways of Processing Sequences

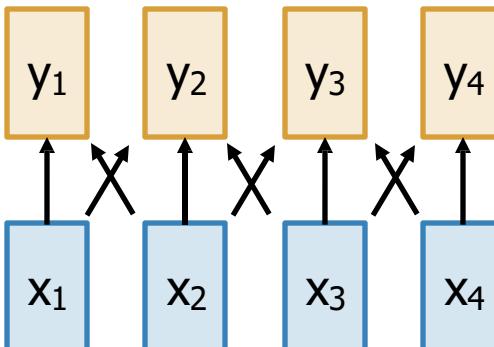
Recurrent Neural Network



Works on 1D ordered sequences

- (+) Theoretically good at long sequences:  $O(N)$  compute and memory for a sequence of length  $N$
- (-) Not parallelizable. Need to compute hidden states sequentially

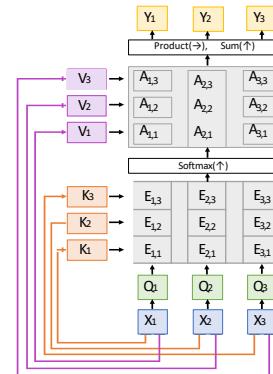
Convolution



Works on N-dimensional grids

- (-) Bad for long sequences: need to stack many layers to build up large receptive fields
- (+) Parallelizable, outputs can be computed in parallel

Self-Attention



Works on sets of vectors

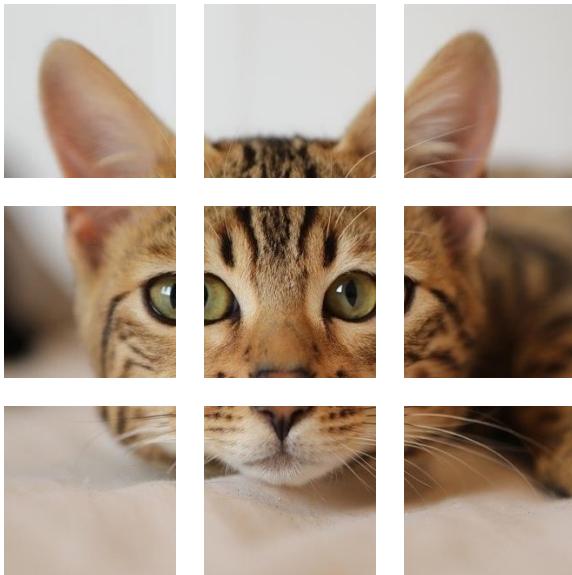
- (+) Great for long sequences; each output depends directly on all inputs
- (+) Highly parallel, it's just 4 matmuls
- (-) Expensive:  $O(N^2)$  compute,  $O(N)$  memory for sequence of length  $N$

# Vision Transformers (ViT)



Input image:  
e.g. 224x224x3

# Vision Transformers (ViT)



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](#)

# Vision Transformers (ViT)

N input patches, each of  
shape 3x16x16



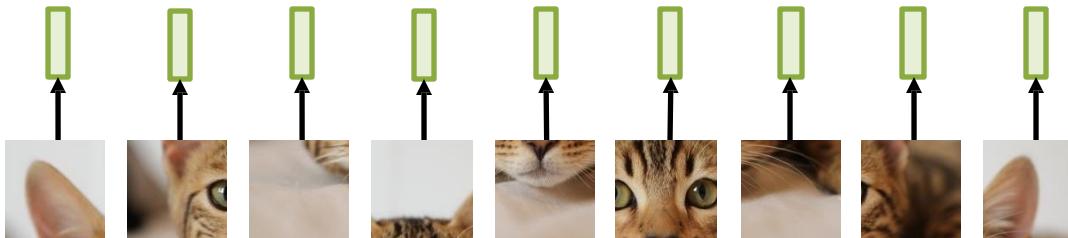
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

# Vision Transformers (ViT)

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

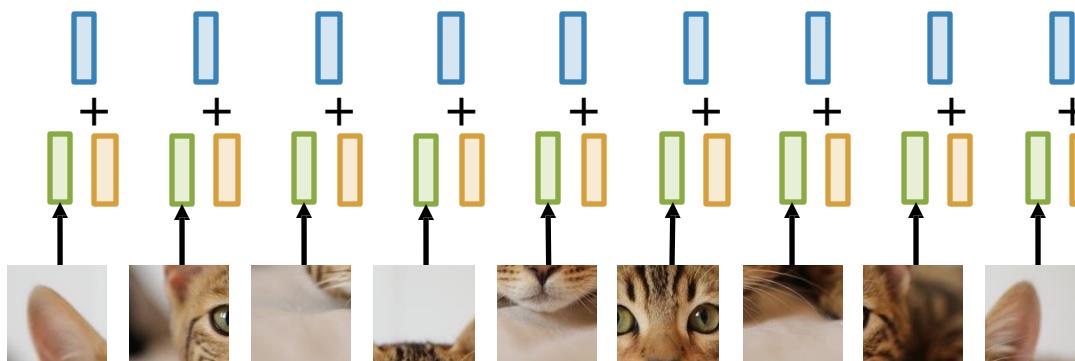
[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformers (ViT)

Add positional embedding:  
learned D-dim vector per  
position

Linear projection to D-  
dimensional vector

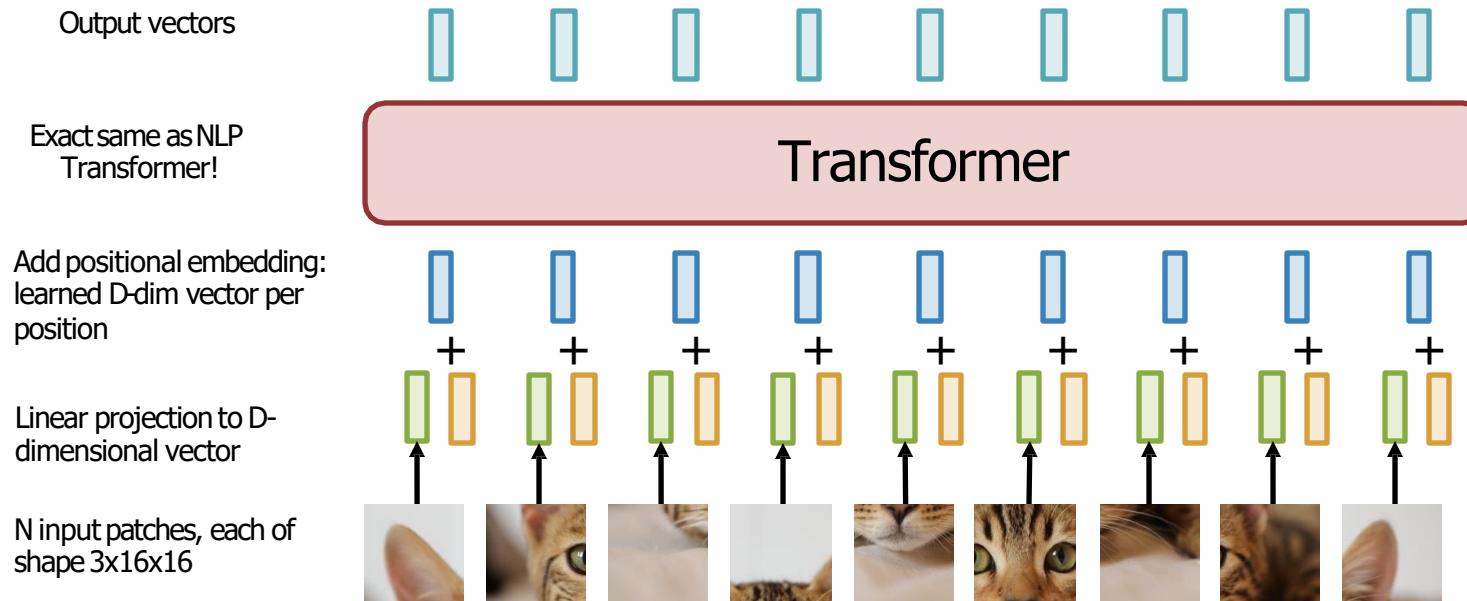
N input patches, each of  
shape 3x16x16



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

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

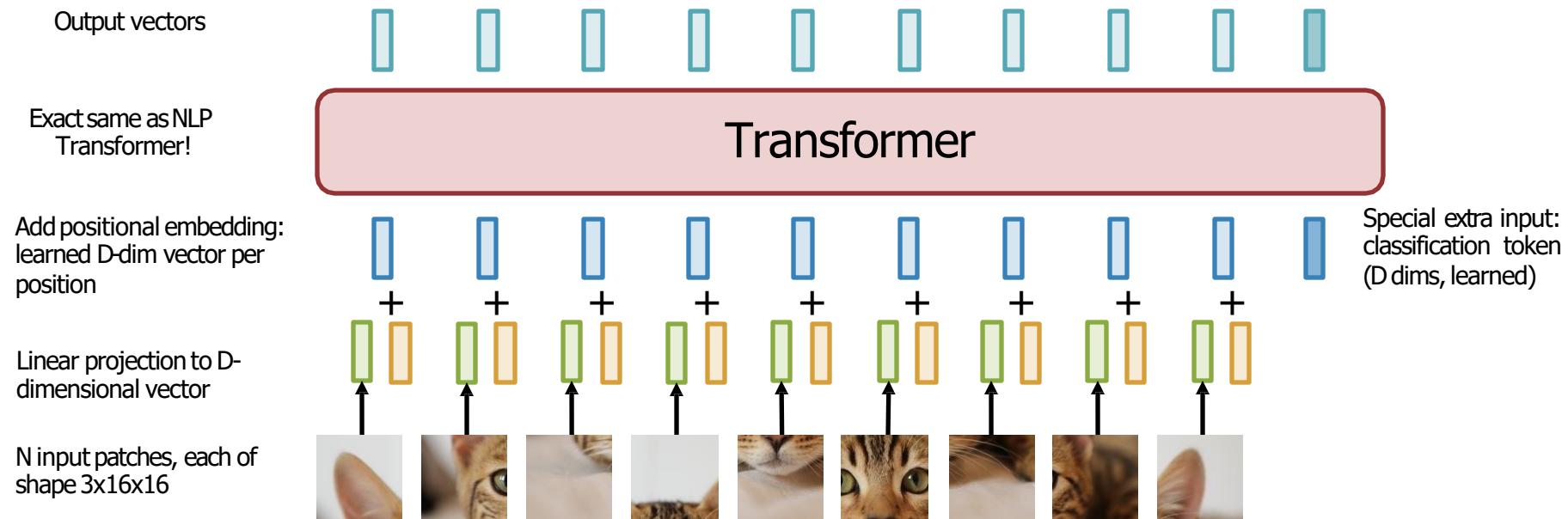
# Vision Transformers (ViT)



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](#)

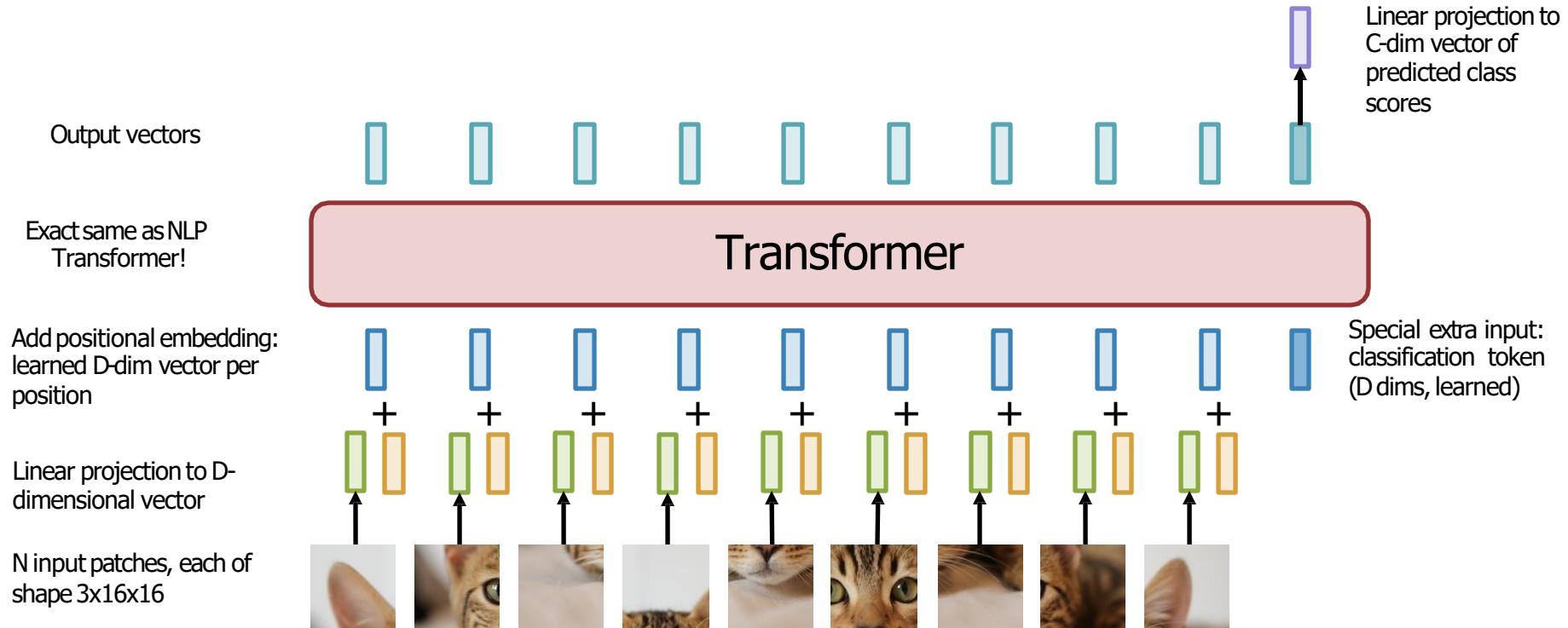
# Vision Transformers (ViT)



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](#)

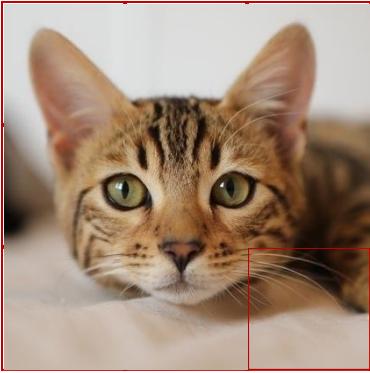
# Vision Transformers (ViT)



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](#)

# Vision Transformers (ViT) – a similar approach (different classifier)



Input image:  
e.g. 224x224x3

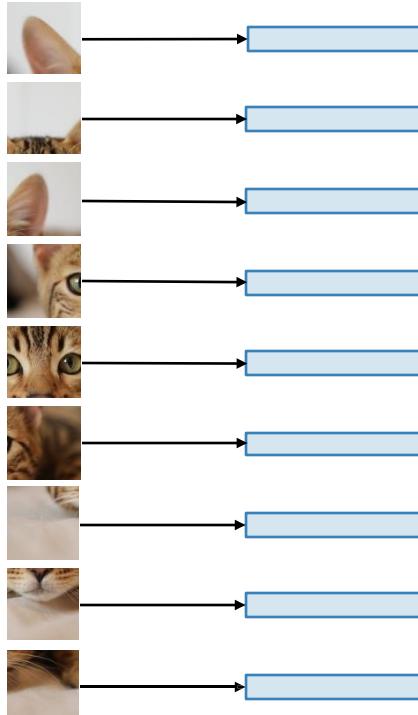


Break into patches  
e.g. 16x16x3

# Vision Transformers (ViT)



Input image:  
e.g. 224x224x3



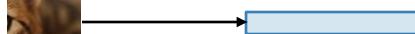
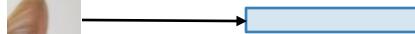
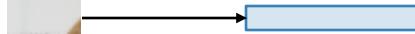
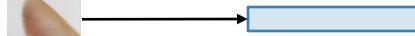
Break into patches  
e.g. 16x16x3

Flatten and apply a linear  
transform  $768 \Rightarrow D$

# Vision Transformers (ViT)

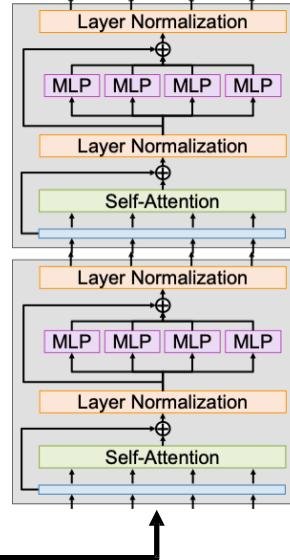
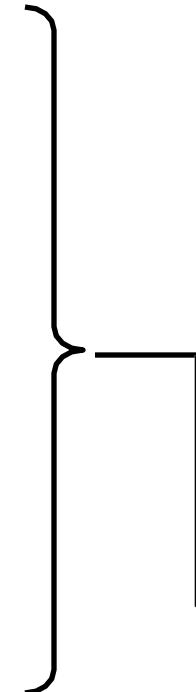


Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

Flatten and apply a linear  
transform  $768 \Rightarrow D$

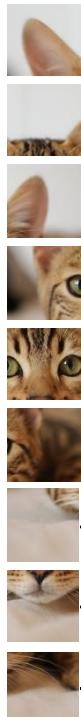


D-dim vector per patch are  
the input vectors to the  
Transformer

# Vision Transformers (ViT)

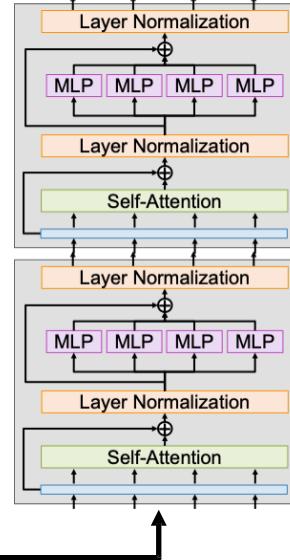
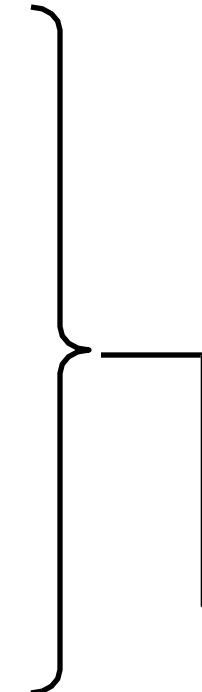


Input image:  
e.g. 224x224x3



Break into patches  
e.g. 16x16x3

Flatten and apply a linear  
transform  $768 \Rightarrow D$



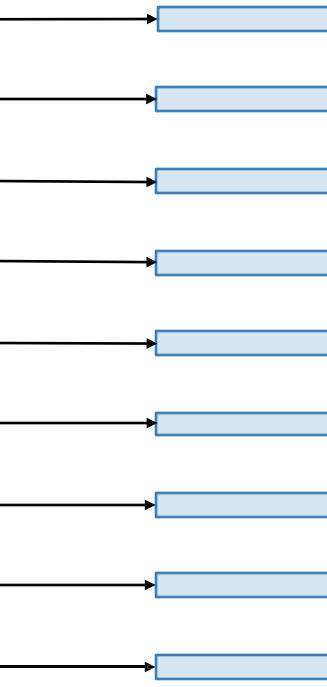
D-dim vector per patch are  
the input vectors to the  
Transformer

Use positional  
encoding to tell  
the transformer  
the 2D position  
of each patch

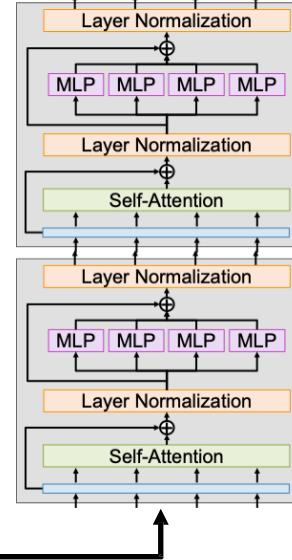
# Vision Transformers (ViT)



Input image:  
e.g. 224x224x3



Flatten and apply a linear  
transform  $768 \Rightarrow D$



D-dim vector per patch are  
the input vectors to the  
Transformer

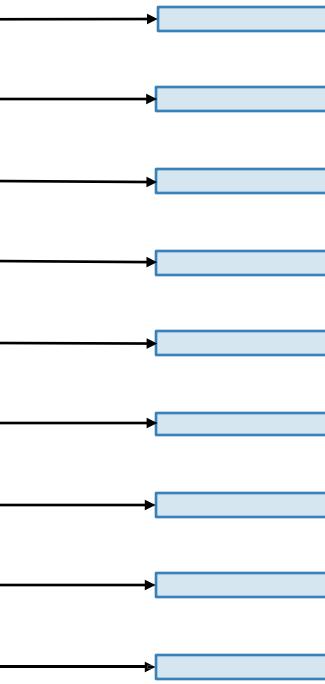
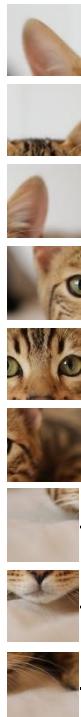
Don't use any  
masking; each  
image patch can  
look at all other  
image patches

Use positional  
encoding to tell  
the transformer  
the 2D position  
of each patch

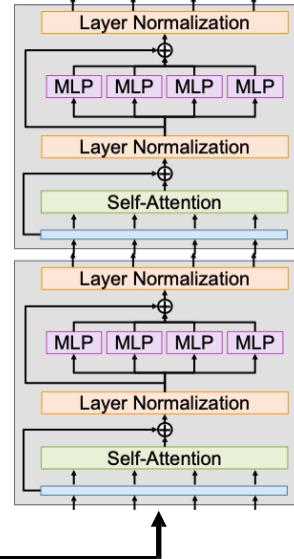
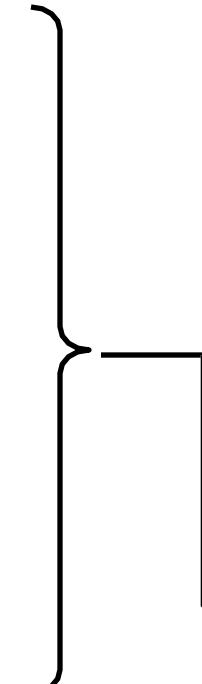
# Vision Transformers (ViT)



Input image:  
e.g. 224x224x3



Flatten and apply a linear  
transform  $768 \Rightarrow D$



D-dim vector per patch are  
the input vectors to the  
Transformer

Transformer  
gives an output  
vector per patch

Don't use any  
masking; each  
image patch can  
look at all other  
image patches

Use positional  
encoding to tell  
the transformer  
the 2D position  
of each patch

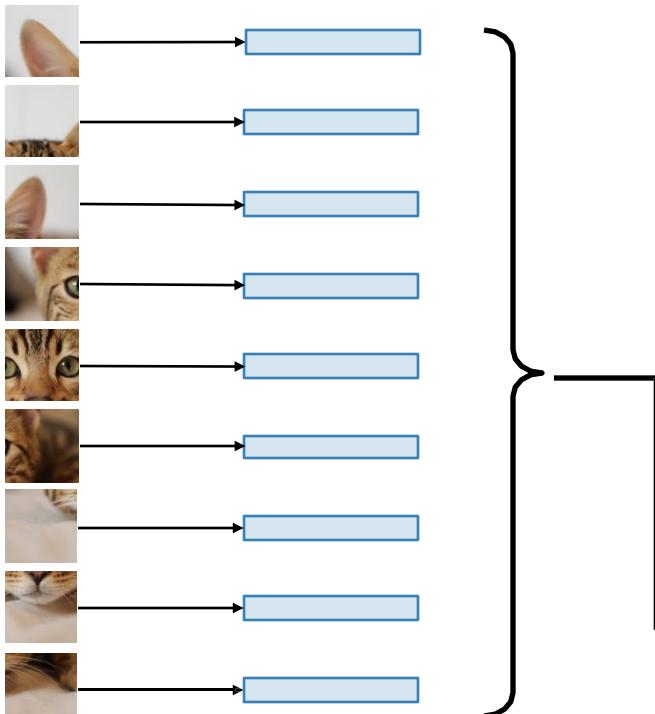
# Vision Transformers (ViT)



Input image:  
e.g. 224x224x3

Break into patches  
e.g. 16x16x3

Flatten and apply a linear transform  $768 \Rightarrow D$



Average pool  $N \times D$  vectors to  
 $1 \times D$ , apply a linear layer  $D \Rightarrow C$  to  
predict class scores

Transformer  
gives an output  
vector per patch

Don't use any masking; each image patch can look at all other image patches

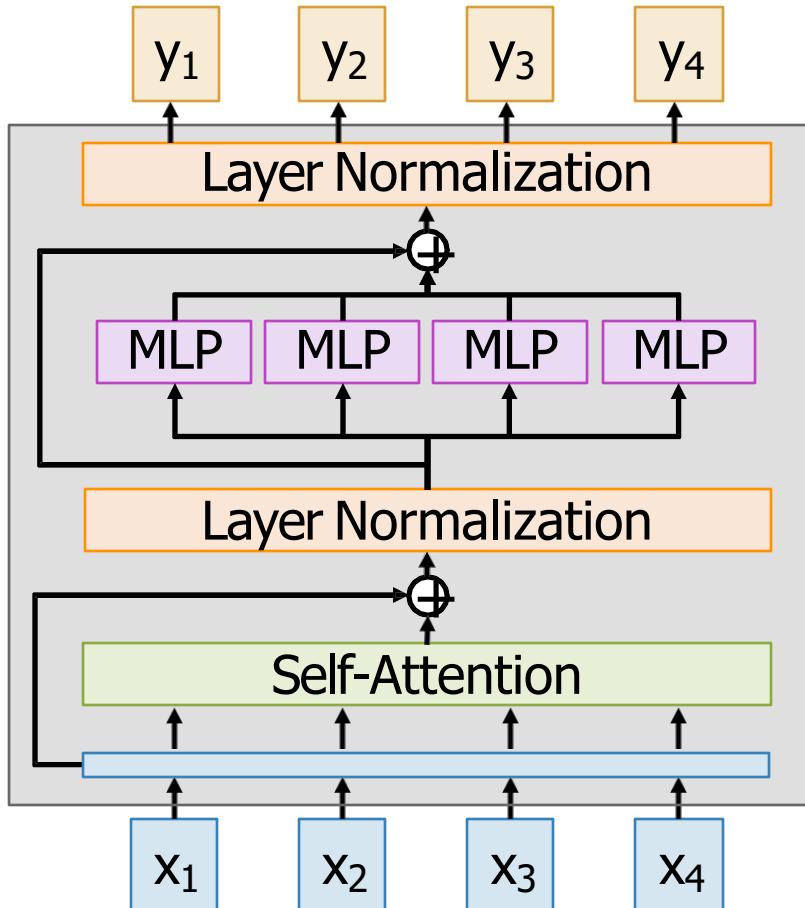
Use positional encoding to tell the transformer the 2D position of each patch

D-dim vector per patch are  
the input vectors to the  
Transformer

# Tweaking Transformers

The Transformer architecture has not changed much since 2017.

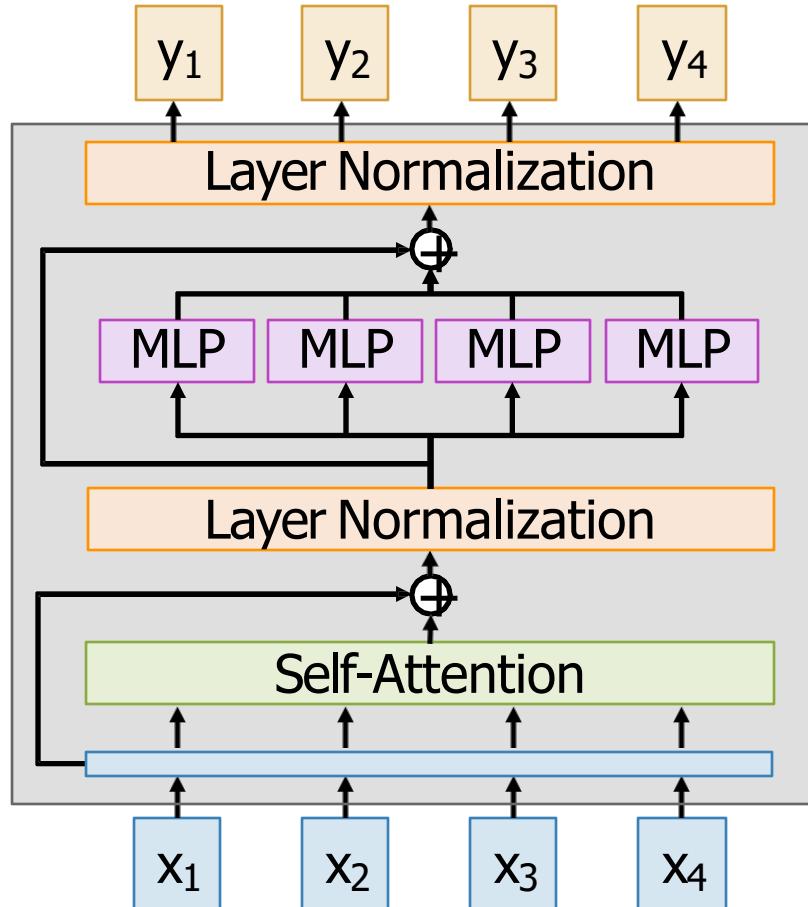
But a few changes have become common:



# Pre-Norm Transformer

Layer normalization is outside the residual connections

Kind of weird, the model can't actually learn the identity function



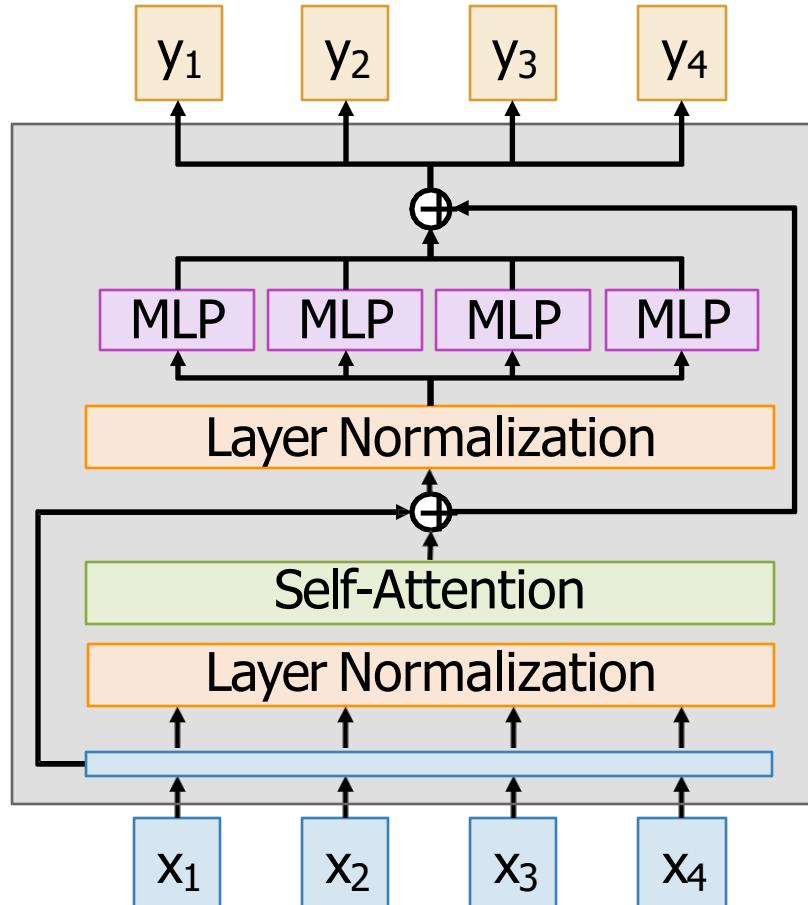
# Pre-Norm Transformer

Layer normalization is outside the residual connections

Kind of weird, the model can't actually learn the identity function

Solution: Move layer normalization before the Self-Attention and MLP, inside the residual connections.

Training is more stable.



# RMSNorm

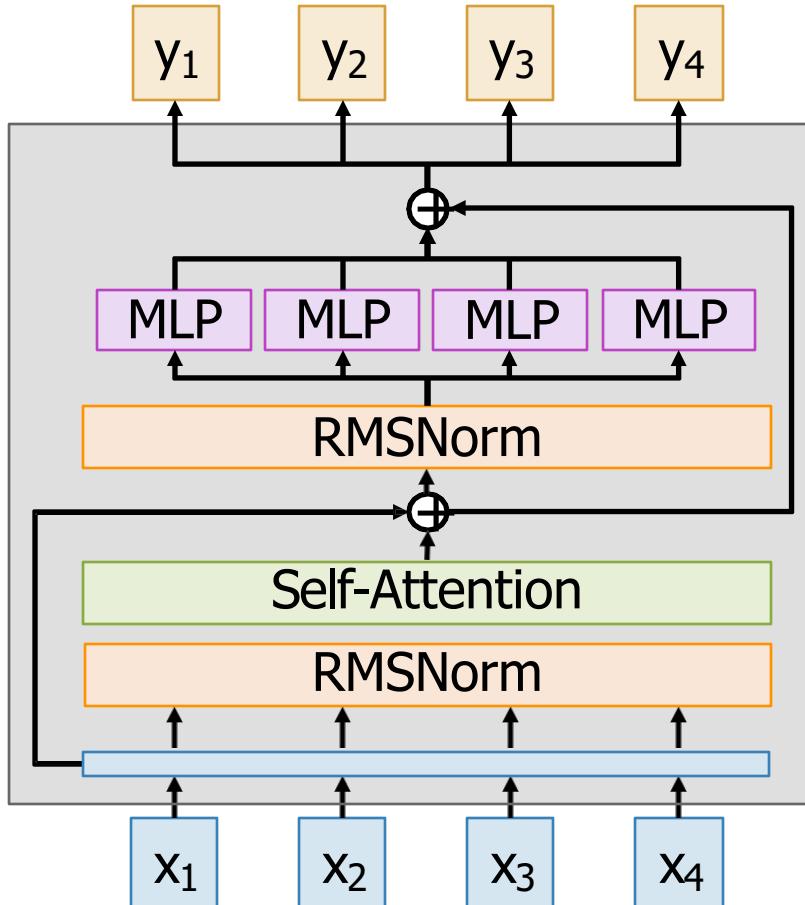
Replace Layer Normalization  
with Root-Mean-Square  
Normalization (RMSNorm)

Input:  $x$  [shape D]  
Output:  $y$  [shape D]  
Weight:  $\gamma$  [shape D]

$$y_i = \frac{x_i}{RMS(x)} * \gamma_i$$

$$RMS(x) = \sqrt{\varepsilon + \frac{1}{N} \sum_{i=1}^N x_i^2}$$

Training is a bit more stable



# SwiGLU MLP

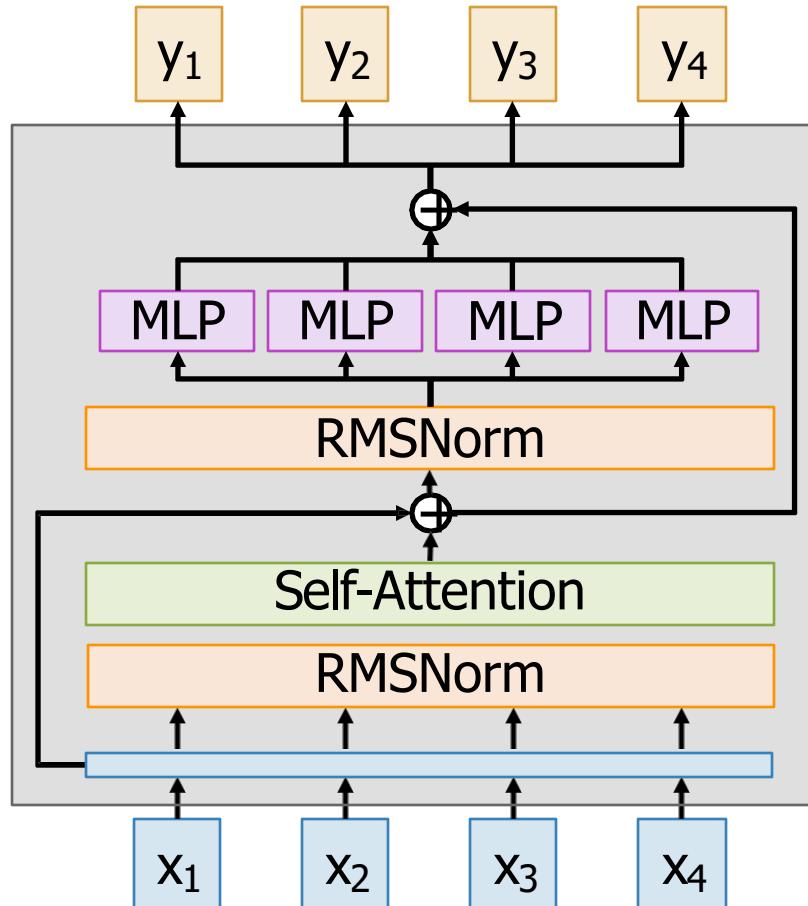
## Classic MLP:

Input:  $X [N \times D]$

Weights:  $W_1 [D \times 4D]$

$W_2 [4D \times D]$

Output:  $Y = \sigma(XW_1)W_2 [N \times D]$



# SwiGLU MLP

## Classic MLP:

Input:  $X [N \times D]$

Weights:  $W_1 [D \times 4D]$

$W_2 [4D \times D]$

Output:  $Y = \sigma(XW_1)W_2 [N \times D]$

## SwiGLU MLP:

Input:  $X [N \times D]$

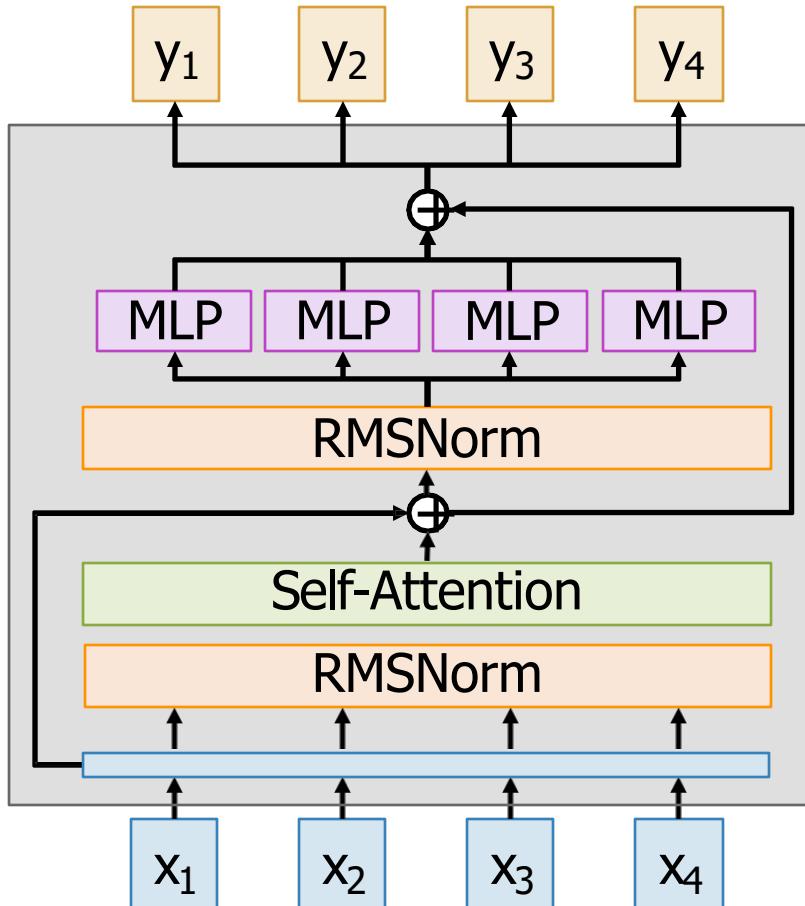
Weights:  $W_1, W_2 [D \times H]$

$W_3 [H \times D]$

Output:

$$Y = (\sigma(XW_1) \odot XW_2)W_3$$

Setting  $H = 8D/3$  keeps  
same total params

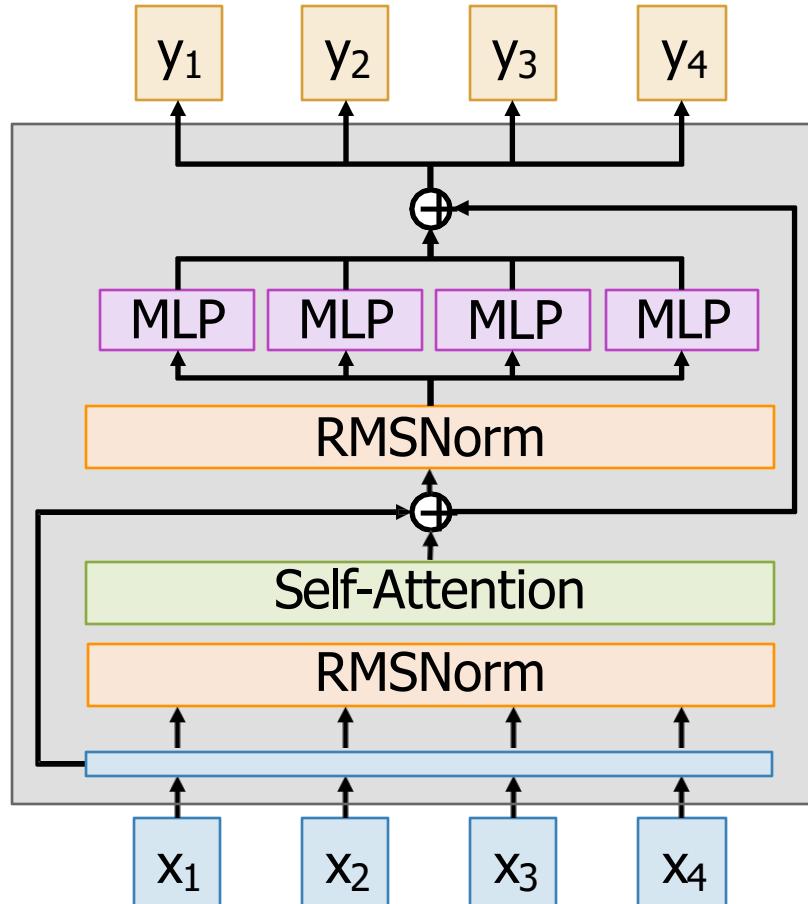


# Mixture of Experts (MoE)

Learn E separate sets of MLP weights in each block; each MLP is an expert

$$W_1: [D \times 4D] \Rightarrow [E \times D \times 4D]$$

$$W_2: [4D \times D] \Rightarrow [E \times 4D \times D]$$



# Mixture of Experts (MoE)

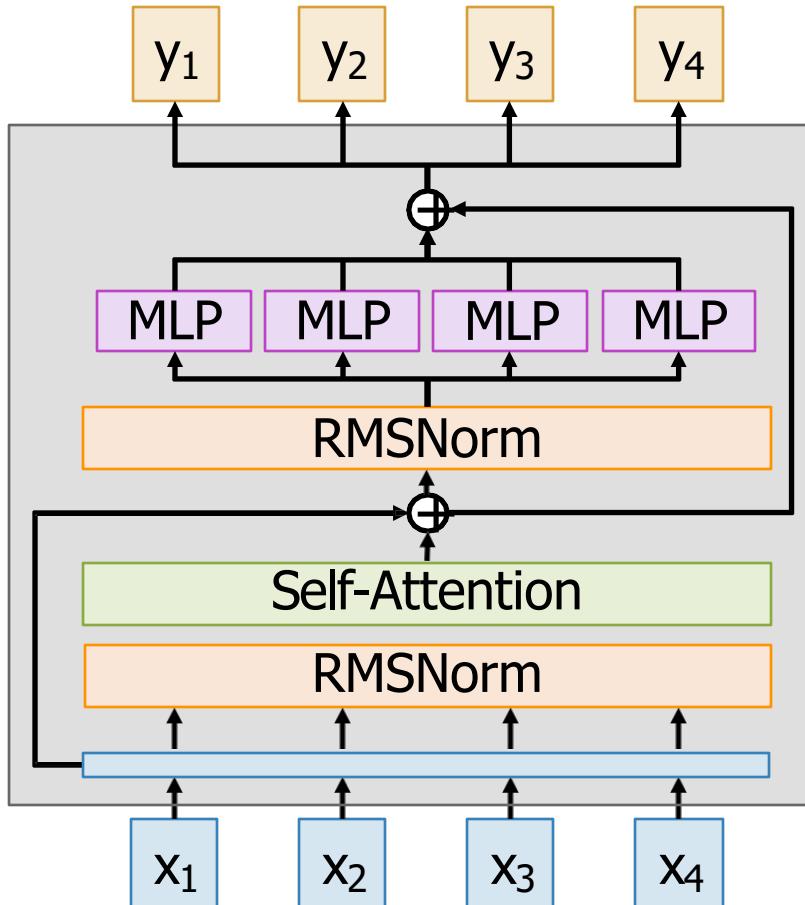
Learn E separate sets of MLP weights in each block; each MLP is an expert

$$W_1: [D \times 4D] \Rightarrow [E \times D \times 4D]$$

$$W_2: [4D \times D] \Rightarrow [E \times 4D \times D]$$

Each token gets routed to A < E of the experts. These are the active experts.

Increases params by E,  
But only increases compute by A



# Mixture of Experts (MoE)

Learn E separate sets of MLP weights in each block; each MLP is an expert

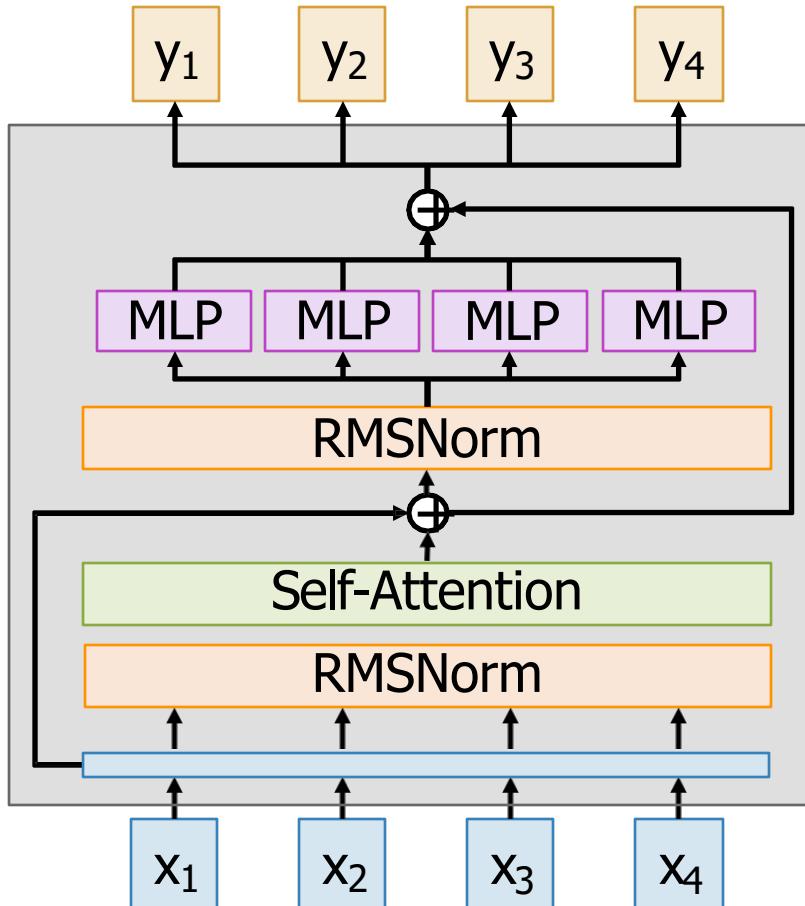
$$W_1: [D \times 4D] \Rightarrow [E \times D \times 4D]$$

$$W_2: [4D \times D] \Rightarrow [E \times 4D \times D]$$

Each token gets routed to A < E of the experts. These are the active experts.

Increases params by E,  
But only increases compute by A

All of the biggest LLMs today (e.g. GPT4o, GPT4.5, Claude 3.7, Gemini 2.5 Pro, etc) almost certainly use MoE and have >1T params; but they don't publish details anymore

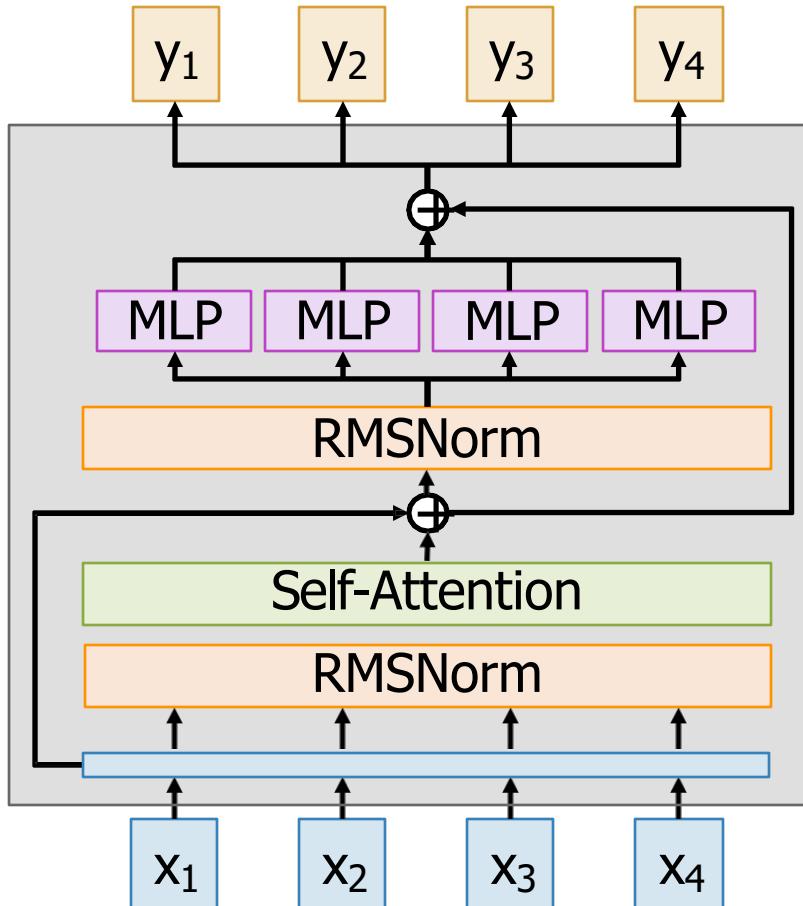


# Tweaking Transformers

The Transformer architecture has not changed much since 2017.

But a few changes have become common:

- Pre-Norm: Move normalization inside residual
- RMSNorm: Different normalization layer
- SwiGLU: Different MLP architecture
- Mixture of Experts (MoE): Learn E different MLPs, use A < E of them per token. Massively increase params, modest increase to compute cost.



# Today

- Transformers Recap
- **Computer Vision Tasks**
  - Semantic Segmentation
  - Object Detection
  - Instance Segmentation
- Visualization & Understanding
  - Model Layers Visualization
  - Saliency Maps
  - CAM & Grad-CAM

# Computer Vision Tasks

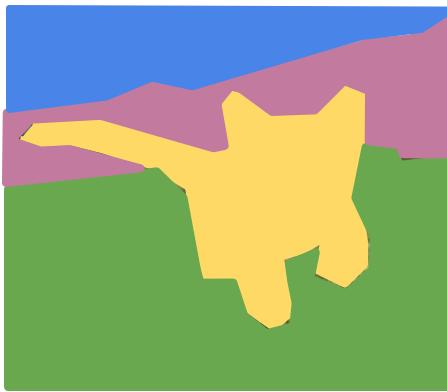
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

[This image](#) is CC0 public domain

# Recall: Image Classification: A core task in Computer Vision



(assume given a set of possible labels)  
{dog, cat, truck, plane, ...}



cat

This image by [Nikita](#) is  
licensed under [CC-BY 2.0](#)

# Semantic Segmentation

Classification



CAT

No spatial extent

## Semantic Segmentation



No objects, just pixels

Object  
Detection



DOG, DOG, CAT

Multiple Object

Instance  
Segmentation



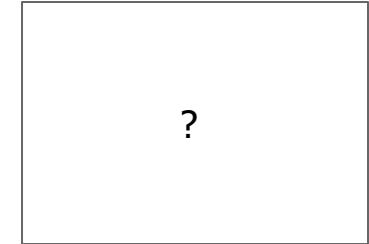
DOG, DOG, CAT

# Semantic Segmentation: The Problem



GRASS, CAT, TREE,  
SKY, ...

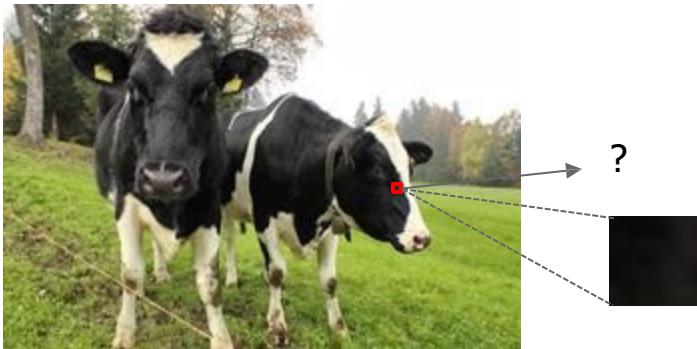
Paired training data: for each training image,  
each pixel is labeled with a semantic category.



At test time, classify each pixel of a new image.

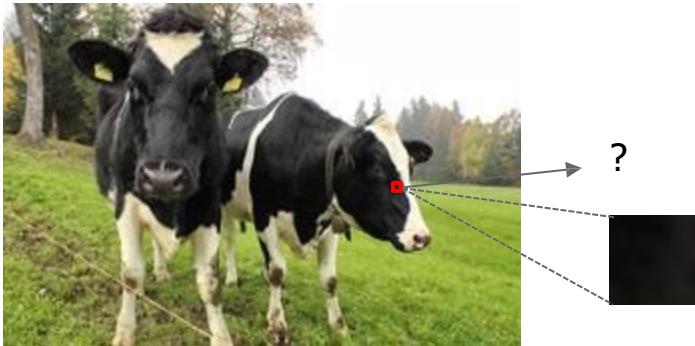
# Semantic Segmentation Idea: Sliding Window

Full image



# Semantic Segmentation Idea: Sliding Window

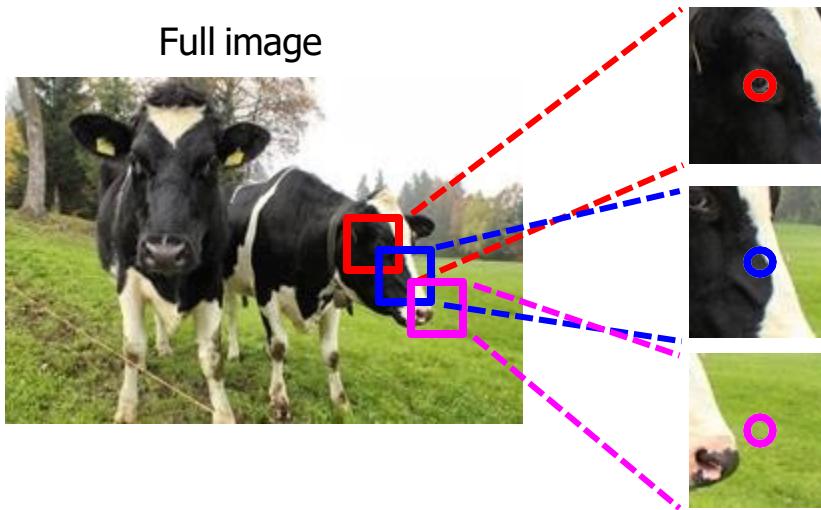
Full image



Impossible to classify without context

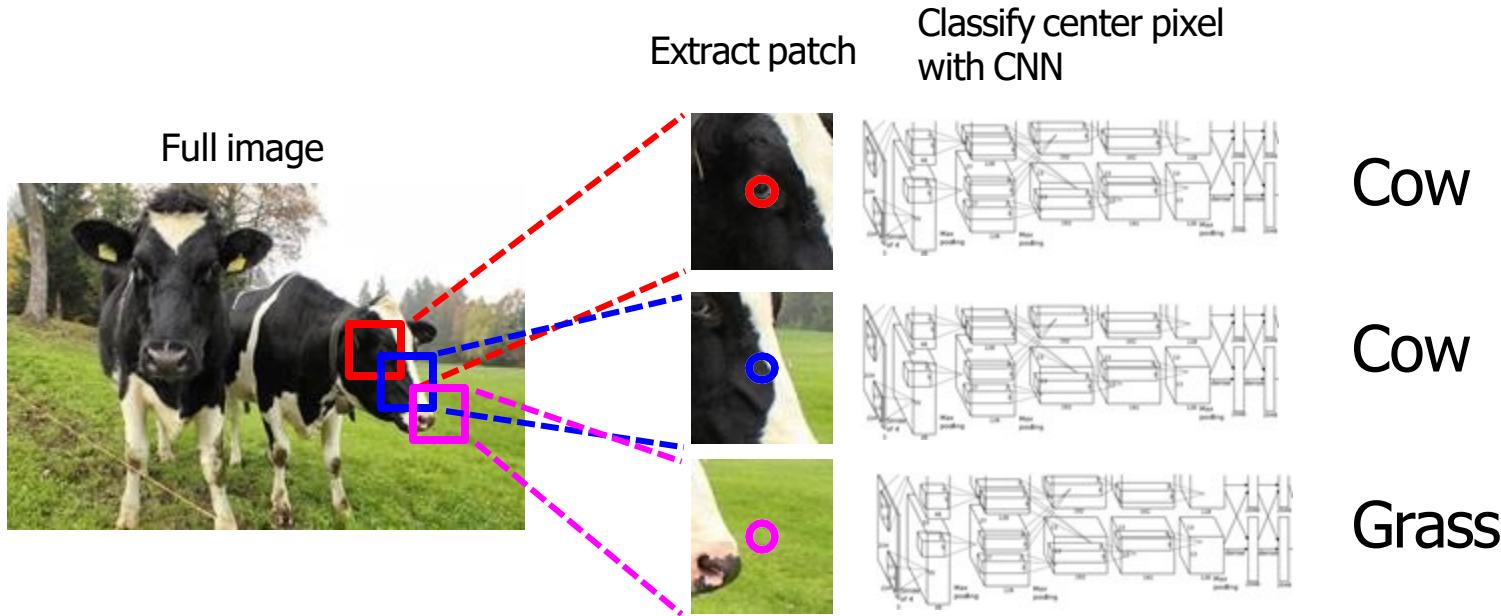
Q: how do we include context?

# Semantic Segmentation Idea: Sliding Window



Q: how do we model this?

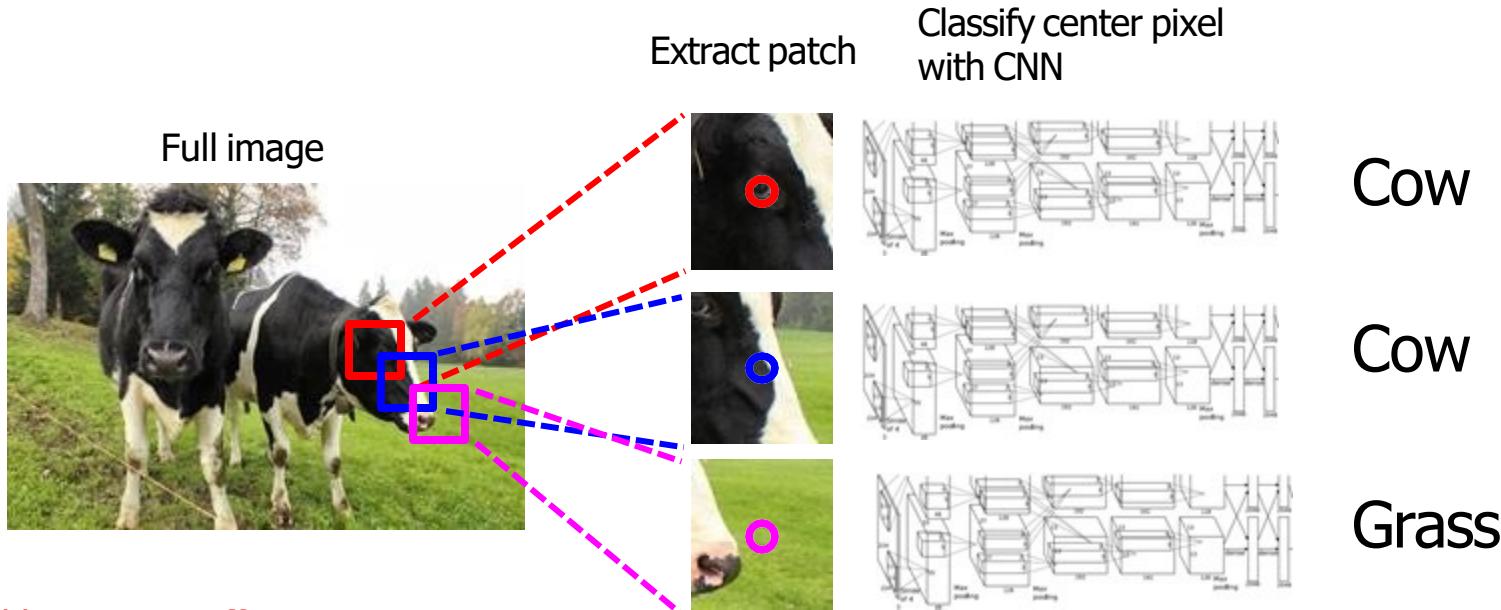
# Semantic Segmentation Idea: Sliding Window



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 Idea: Sliding Window



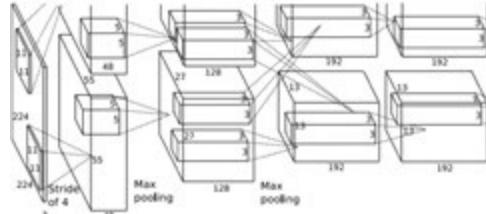
Problem: Very inefficient! Not reusing shared features between overlapping patches

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 Idea: Convolution

Full image

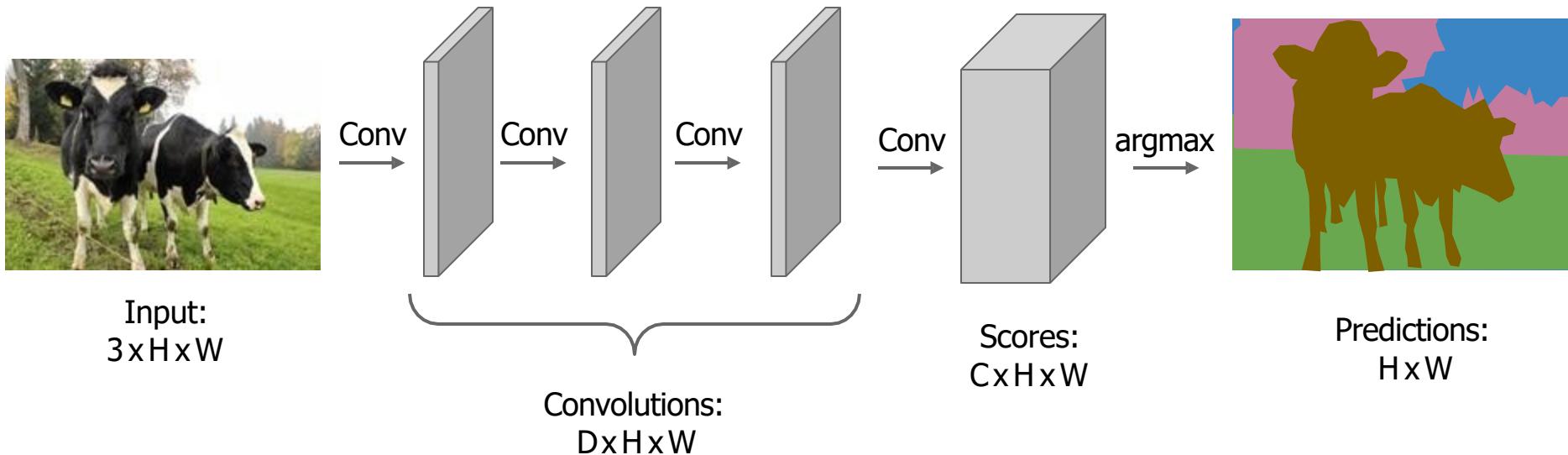


An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

Problem: classification architectures often reduce feature spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as input size.

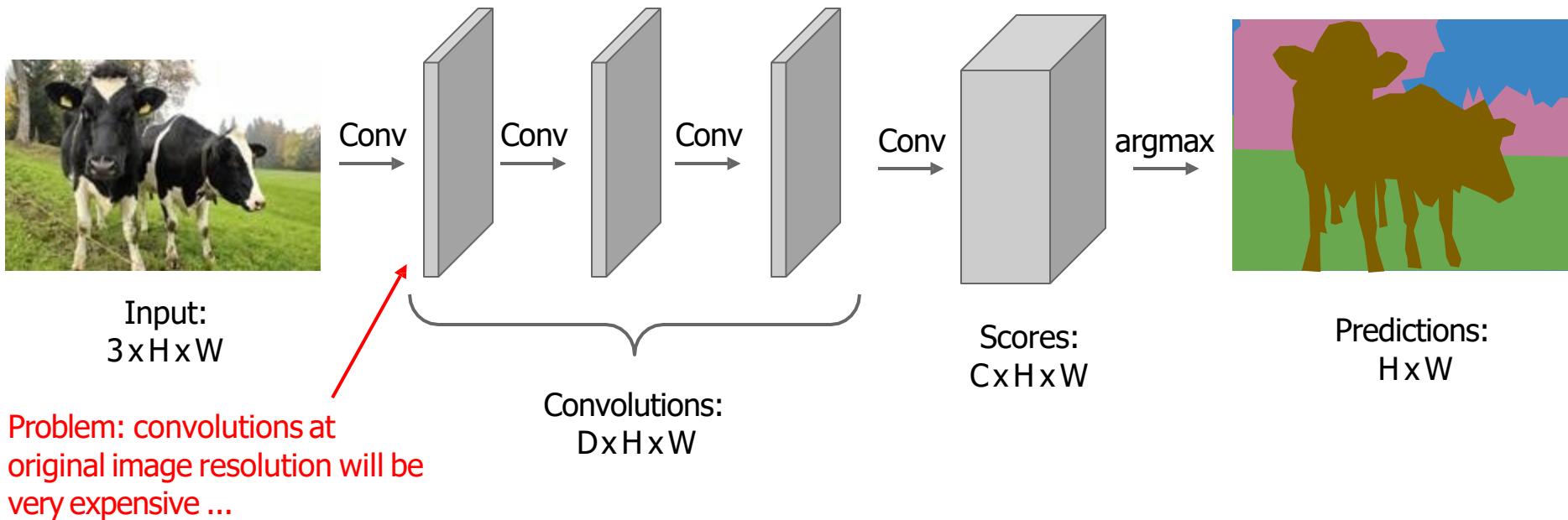
# Semantic Segmentation Idea: Fully Convolutional

Design a network with only convolutional layers  
without downsampling operators to make predictions  
for pixels all at once!



# Semantic Segmentation Idea: Fully Convolutional

Design a network with only convolutional layers  
without downsampling operators to make predictions  
for pixels all at once!

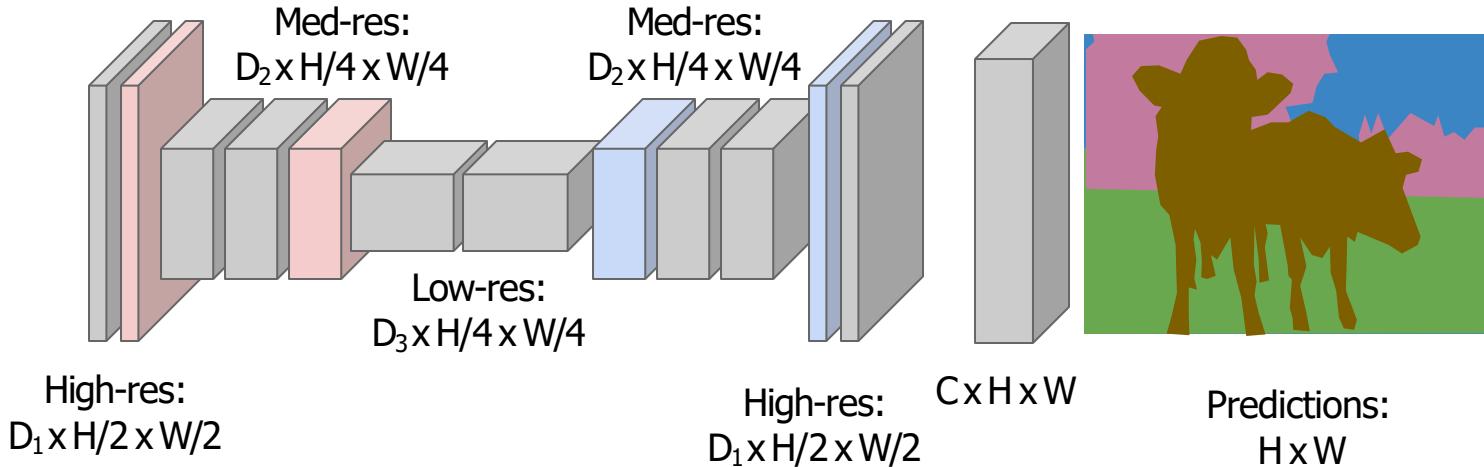


# Semantic Segmentation Idea: Fully Convolutional

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



Input:  
 $3 \times H \times W$



Predictions:  
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

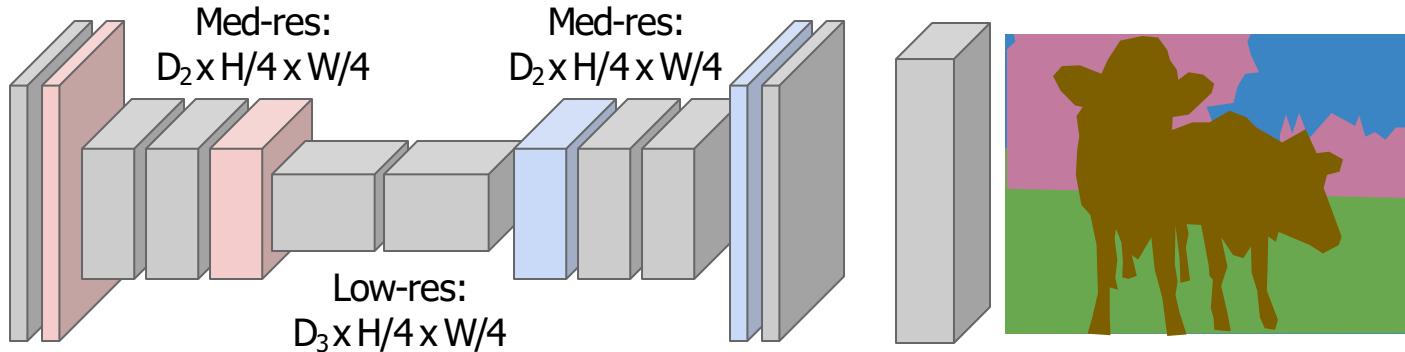
# Semantic Segmentation Idea: Fully Convolutional

Downsampling:  
Pooling, strided  
convolution

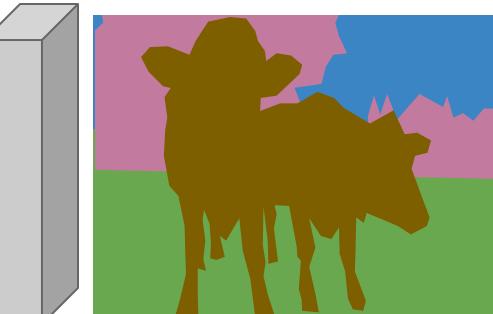


Input:  
 $3 \times H \times W$

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



Upsampling:  
???



Predictions:  
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# In-Network upsampling: “Unpooling”

Nearest Neighbor

1	2
3	4



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Input: 2 x 2

Output: 4 x 4

“Bed of Nails”

1	2
3	4



1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Output: 4 x 4

# In-Network upsampling: “Max Unpooling”

Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4 x 4

5	6
7	8

Output: 2 x 2

Max Unpooling

Use positions from  
pooling layer

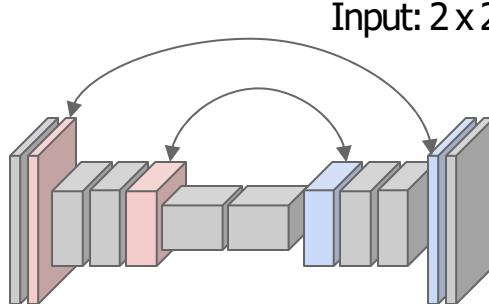
1	2
3	4

Rest of the network

0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

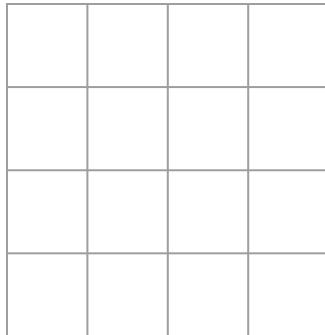
Output: 4 x 4

Corresponding pairs of  
downsampling and  
upsampling layers

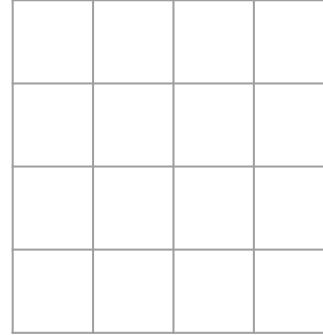


# Learnable Upsampling

Recall: Normal  $3 \times 3$  convolution, stride 1 pad 1



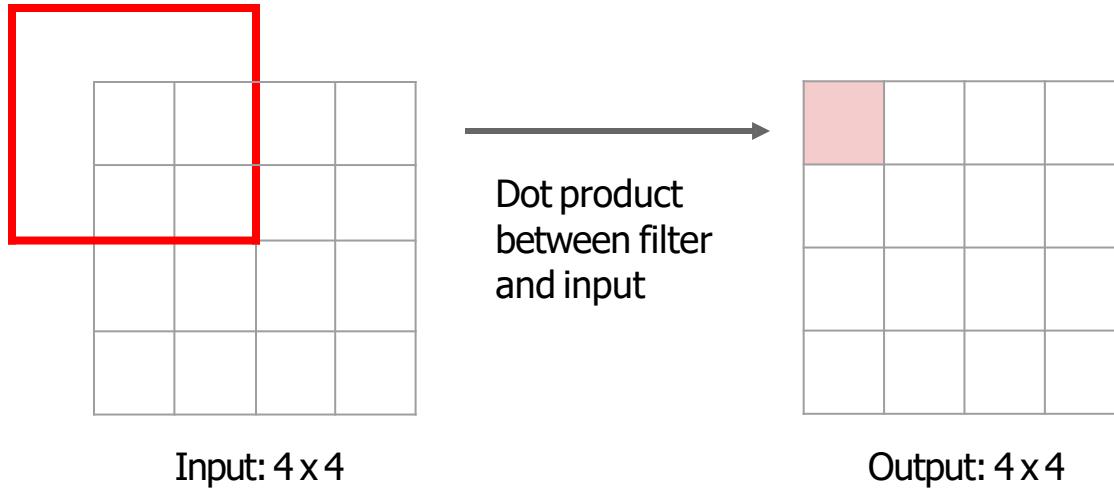
Input:  $4 \times 4$



Output:  $4 \times 4$

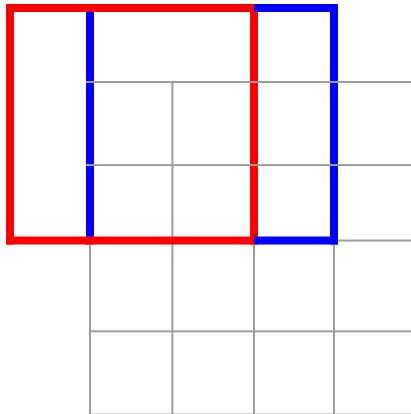
# Learnable Upsampling

Recall: Normal  $3 \times 3$  convolution, stride 1 pad 1



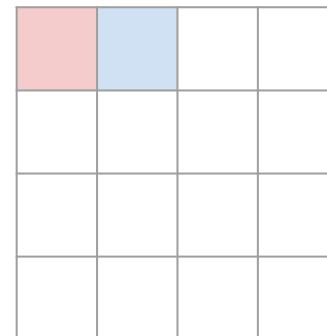
# Learnable Upsampling

Recall: Normal  $3 \times 3$  convolution, stride 1 pad 1



Input:  $4 \times 4$

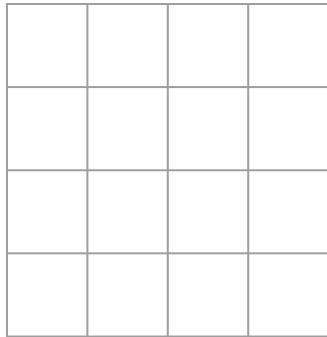
Dot product  
between filter  
and input



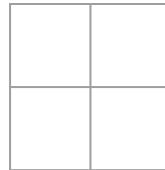
Output:  $4 \times 4$

# Learnable Upsampling

Recall: Normal  $3 \times 3$  convolution, stride 2 pad 1



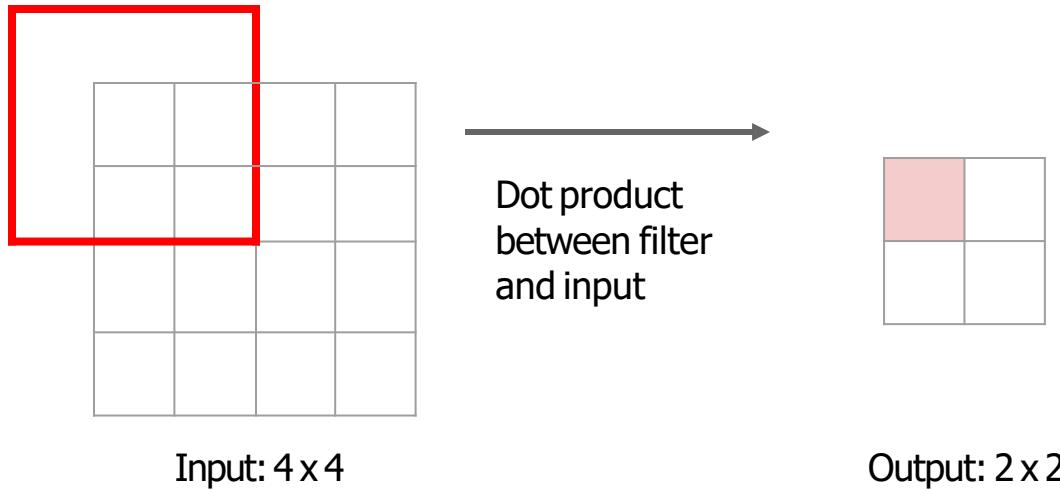
Input:  $4 \times 4$



Output:  $2 \times 2$

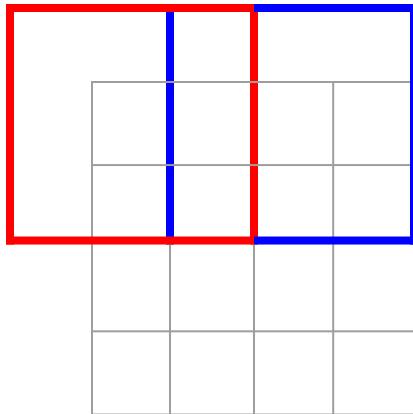
# Learnable Upsampling

Recall: Normal  $3 \times 3$  convolution, stride 2 pad 1



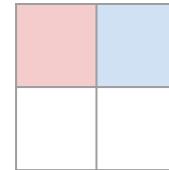
# Learnable Upsampling

Recall: Normal  $3 \times 3$  convolution, stride 2 pad 1



Input:  $4 \times 4$

Dot product  
between filter  
and input



Output:  $2 \times 2$

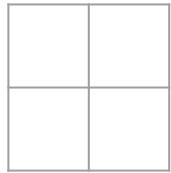
Filter moves 2 pixels in the input for every one pixel in the output

Stride gives ratio between movement in input and output

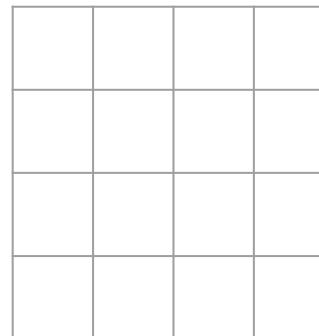
We can interpret strided convolution as “learnable downsampling”.

# Learnable Upsampling: Transposed Convolution

$3 \times 3$  transposed convolution, stride 2 pad 1



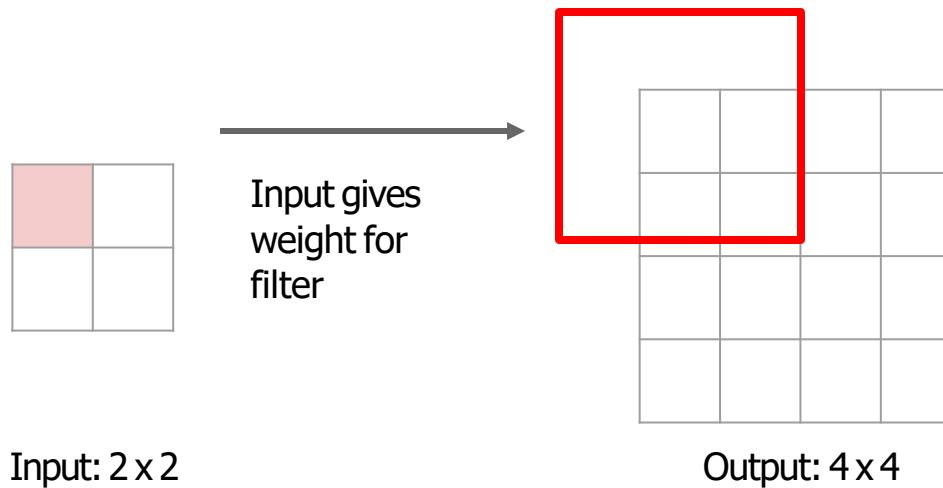
Input:  $2 \times 2$



Output:  $4 \times 4$

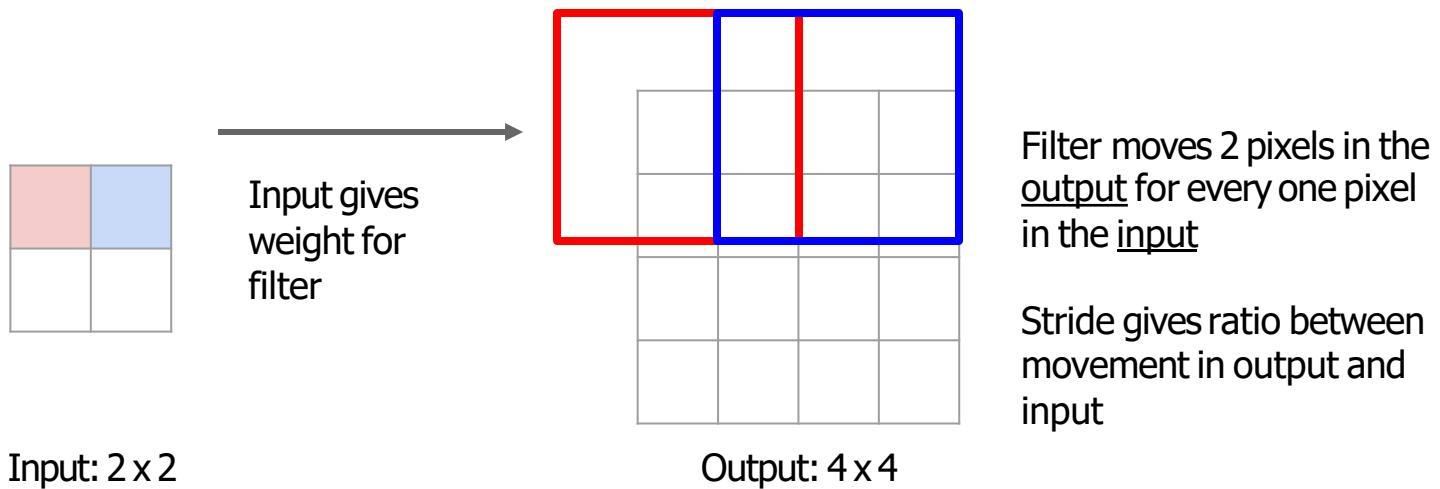
# Learnable Upsampling: Transposed Convolution

3 x 3 transposed convolution, stride 2 pad 1



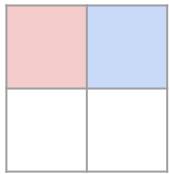
# Learnable Upsampling: Transposed Convolution

3 x 3 transposed convolution, stride 2 pad 1



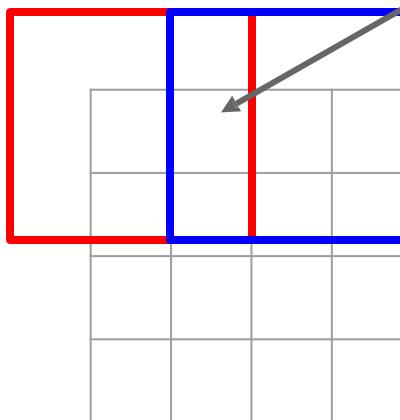
# Learnable Upsampling: Transposed Convolution

3 x 3 transposed convolution, stride 2 pad 1



Input: 2 x 2

Input gives weight for filter



Output: 4 x 4

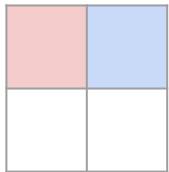
Sum where output overlaps

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

Stride gives ratio between movement in output and input

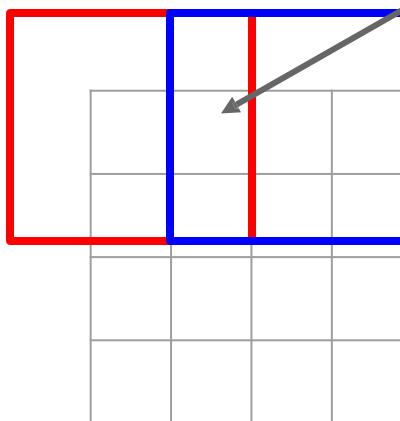
# Learnable Upsampling: Transposed Convolution

3 x 3 transposed convolution, stride 2 pad 1



Input: 2 x 2

Input gives weight for filter



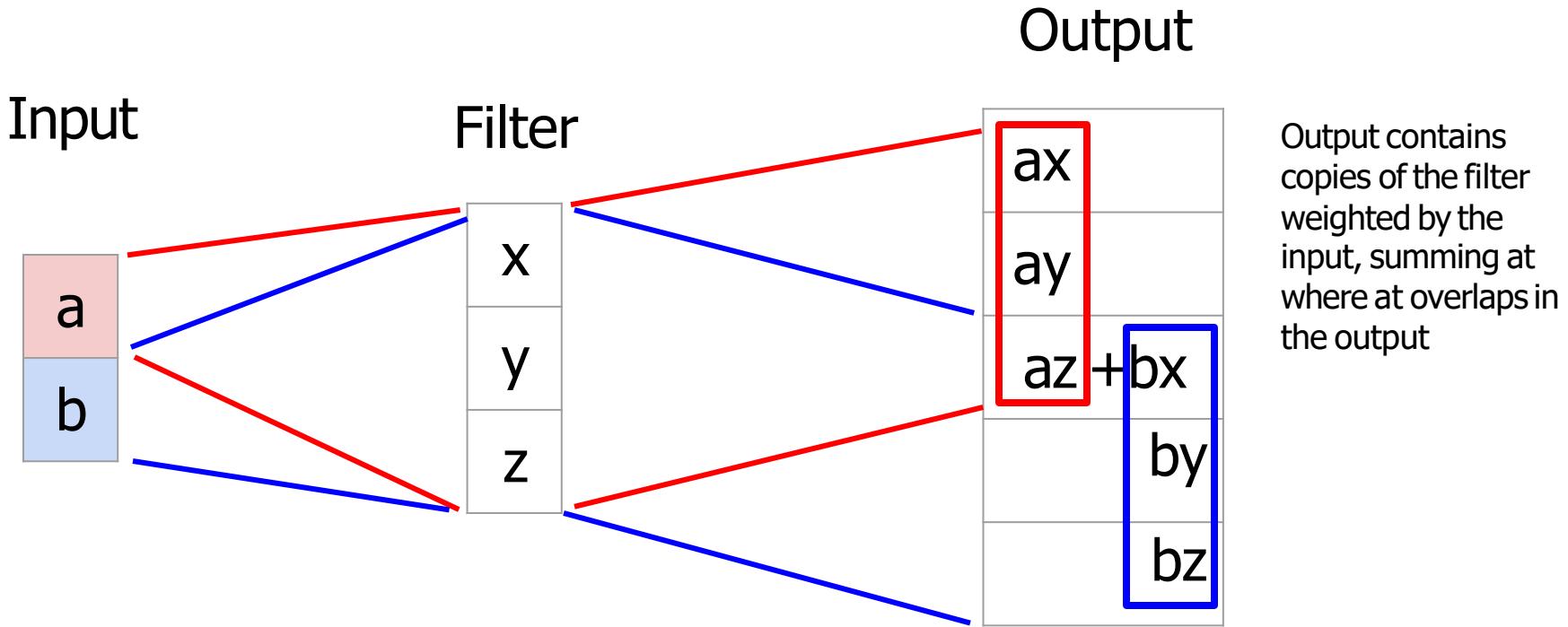
Output: 4 x 4

Sum where output overlaps

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

Stride gives ratio between movement in output and input

# Learnable Upsampling: 1D Example



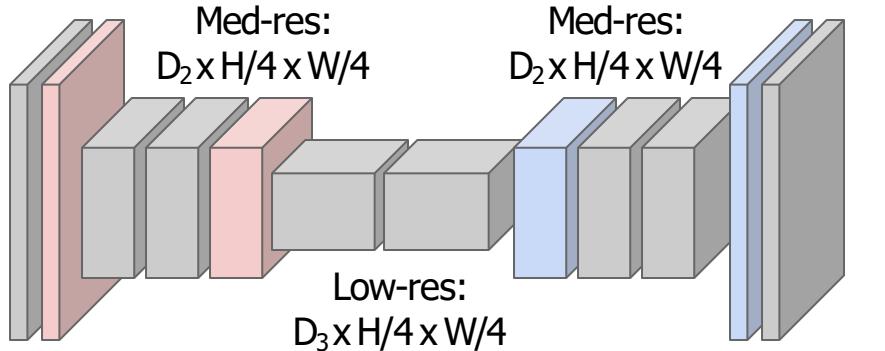
# Semantic Segmentation Idea: Fully Convolutional

Downsampling:  
Pooling, strided  
convolution

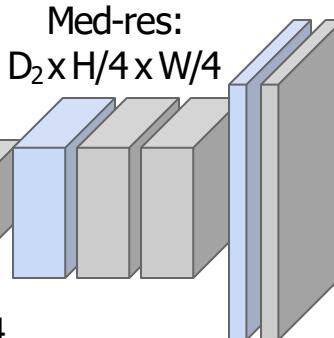


Input:  
 $3 \times H \times W$

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



High-res:  
 $D_1 \times H/2 \times W/2$



High-res:  
 $D_1 \times H/2 \times W/2$

Upsampling:  
Unpooling or strided  
transposed convolution

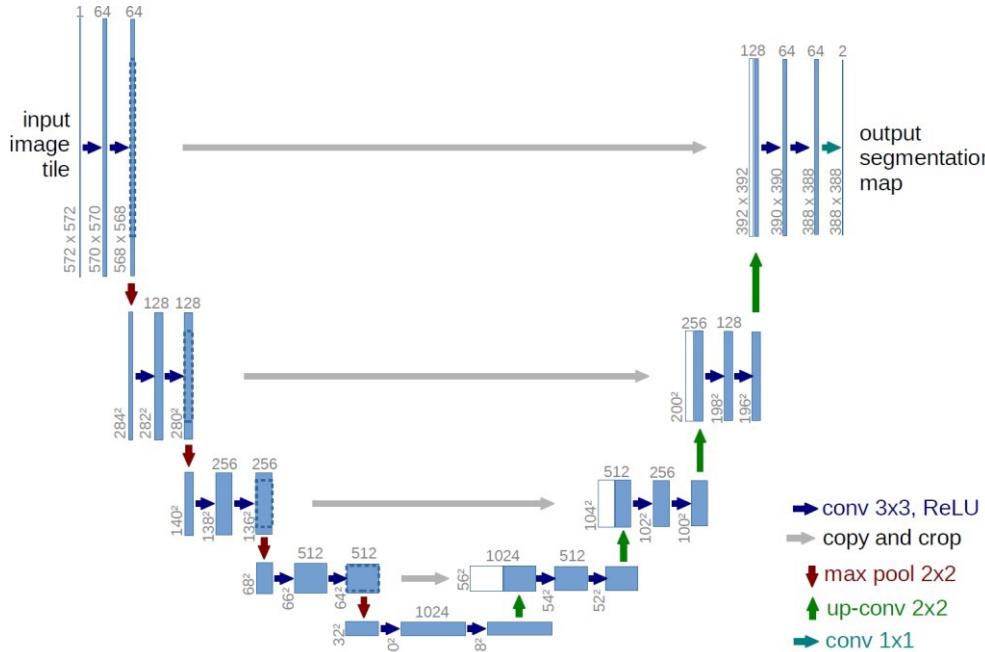


Predictions:  
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

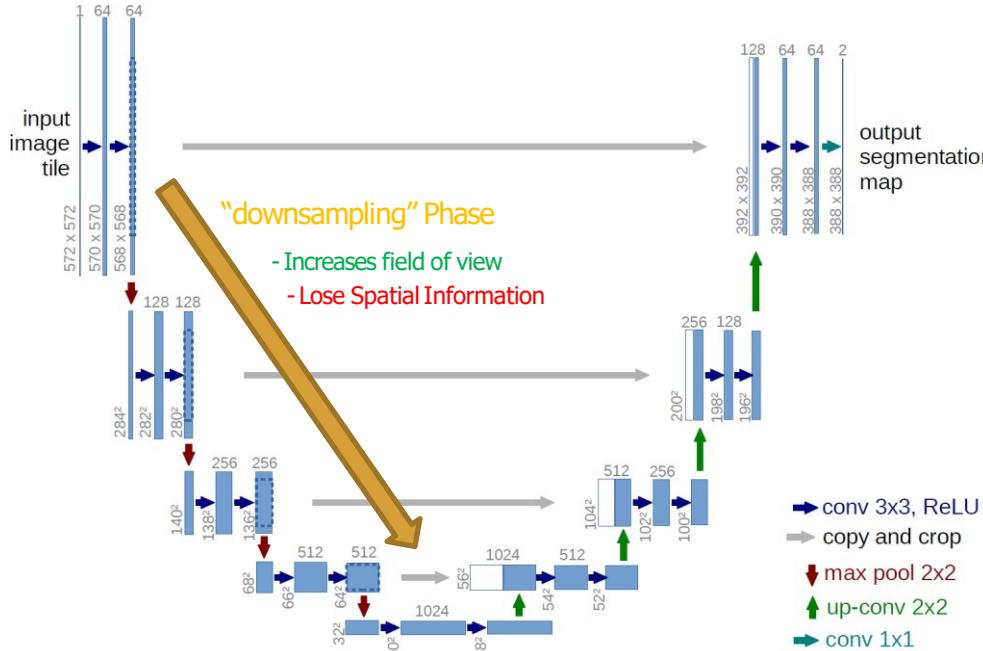
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# U-Net



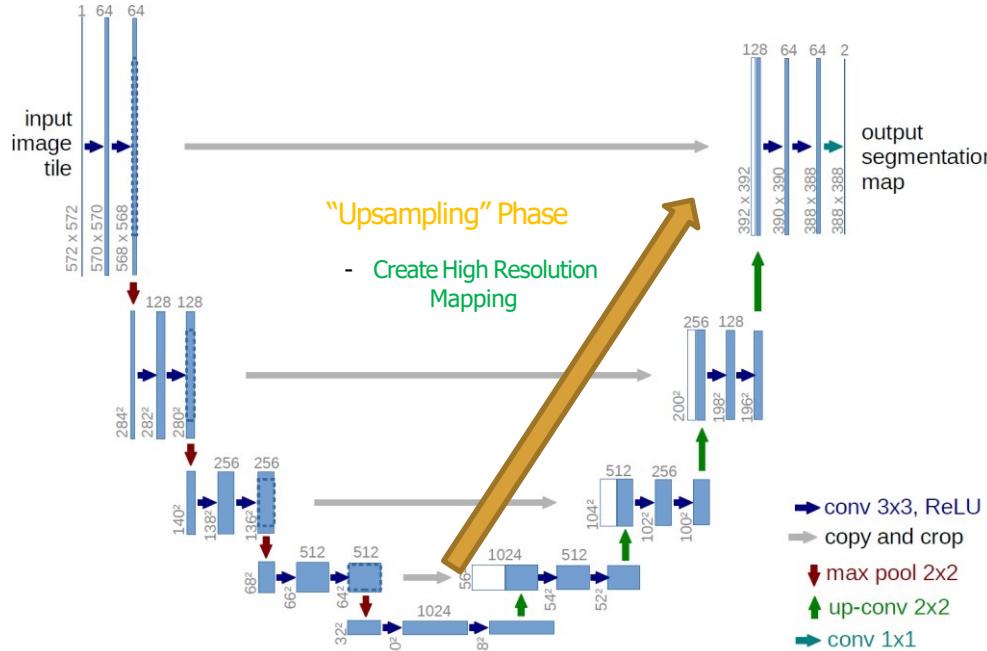
Ronneberger et al. (2015) U-net Architecture

# U-Net

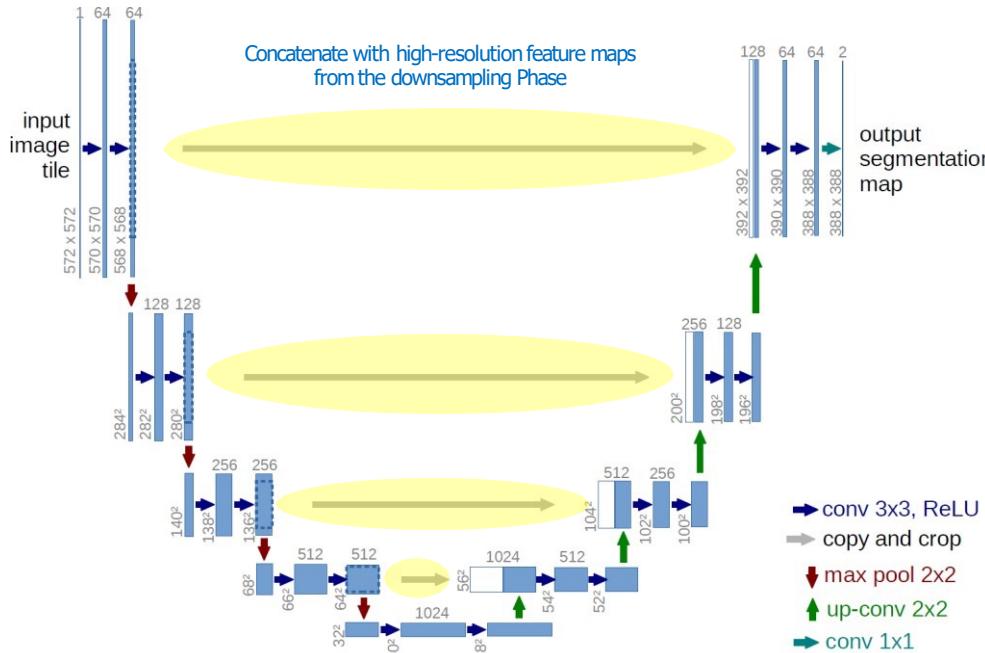


Ronneberger et al. (2015) U-net Architecture

## U-Net

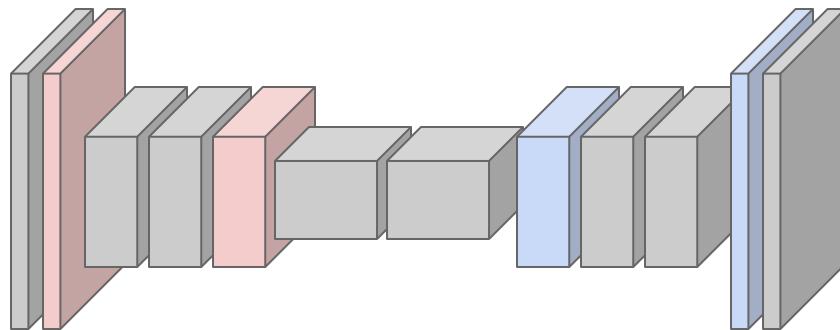
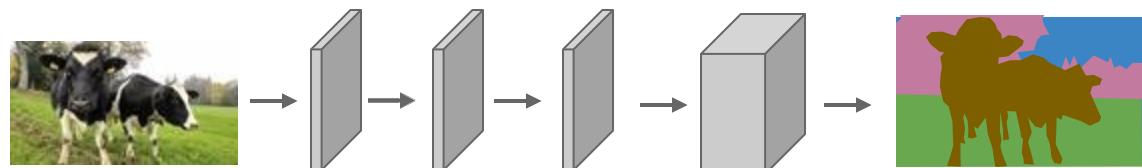
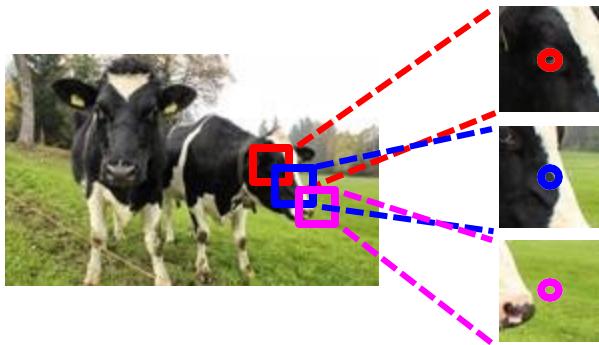


# U-Net



Ronneberger et al. (2015) U-net Architecture

# Semantic Segmentation: Summary



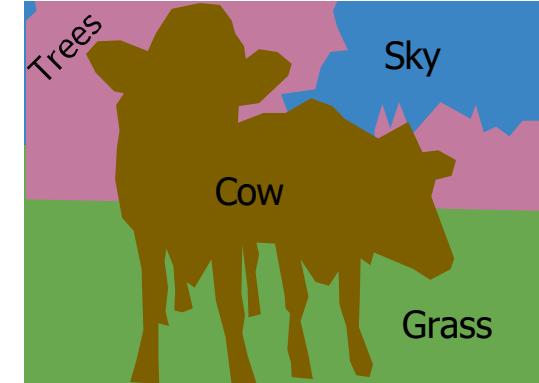
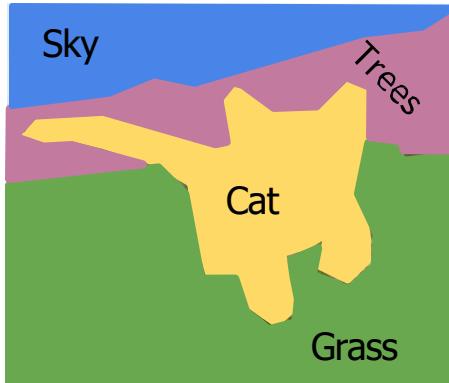
# Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



[This image is CC0 public domain](#)



# Object Detection

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object  
Detection



DOG, DOG, CAT

Multiple Object

Instance  
Segmentation



DOG, DOG, CAT

# Object Detection

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object  
Detection



DOG, DOG, CAT

Multiple Object

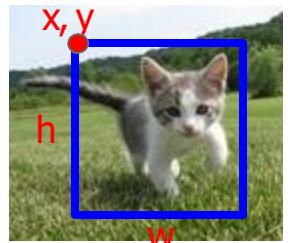
Instance  
Segmentation



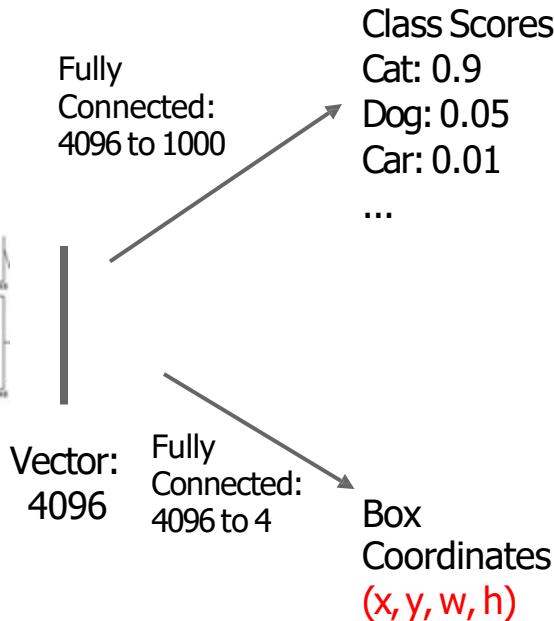
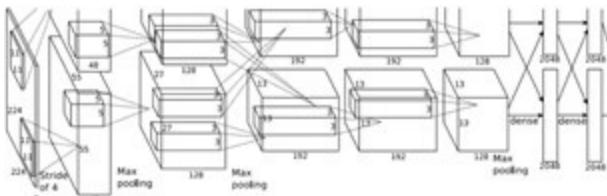
DOG, DOG, CAT

# Object Detection: Single Object

(Classification + Localization)

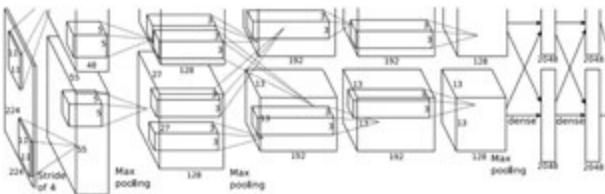
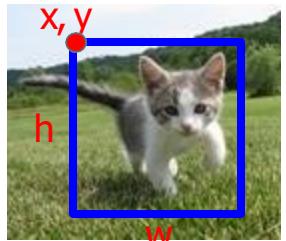


This image is CC0 public domain



# Object Detection: Single Object

(Classification + Localization)



Vector:  
4096

Fully  
Connected:  
4096 to 1000

Class Scores  
Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

Correct label:  
Cat

Softmax  
Loss

Fully  
Connected:  
4096 to 4

Box  
Coordinates  
( $x, y, w, h$ )

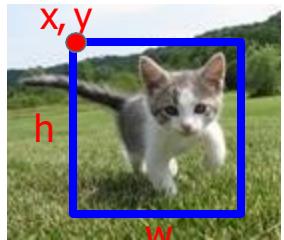
L2 Loss

Correct box:  
( $x', y', w', h'$ )

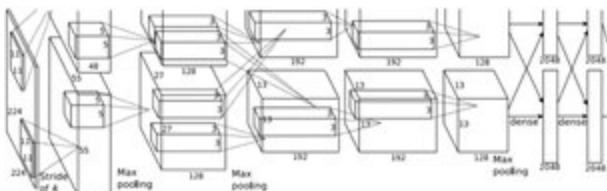
Treat localization as a  
regression problem!

# Object Detection: Single Object

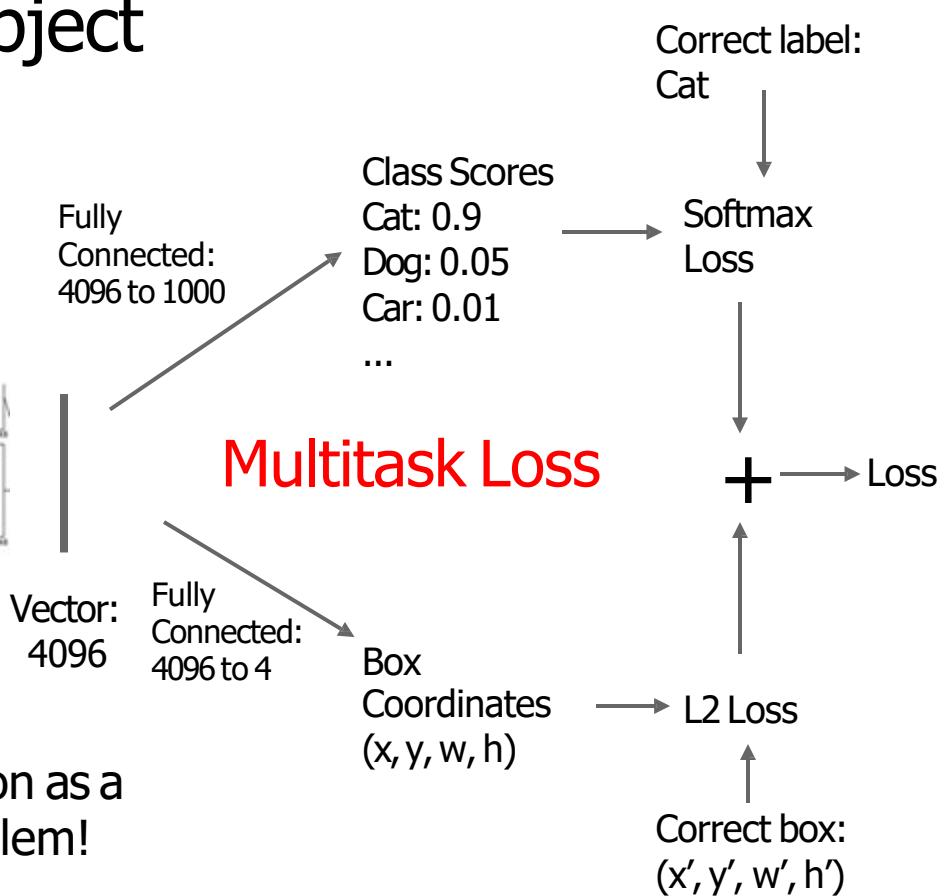
(Classification + Localization)



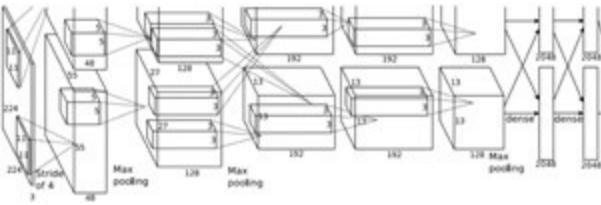
This image is CC0 public domain



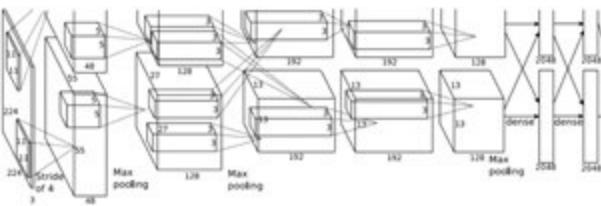
Treat localization as a  
regression problem!



# Object Detection: Multiple Objects



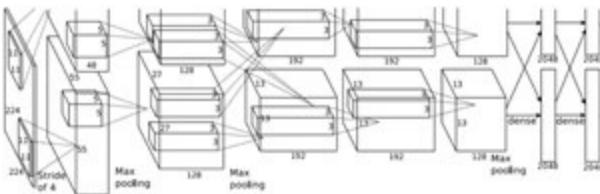
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



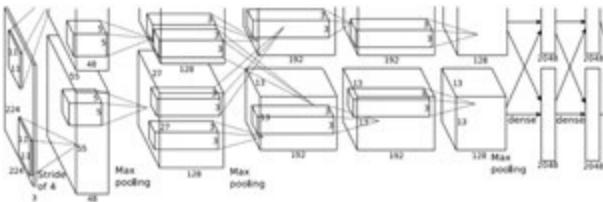
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

...

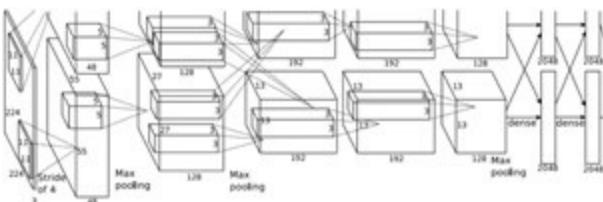
# Object Detection: Multiple Objects

Each image needs a different number of outputs!



CAT:  $(x, y, w, h)$

4 numbers

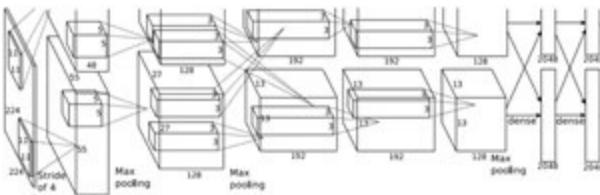


DOG:  $(x, y, w, h)$

DOG:  $(x, y, w, h)$

CAT:  $(x, y, w, h)$

12 numbers



DUCK:  $(x, y, w, h)$

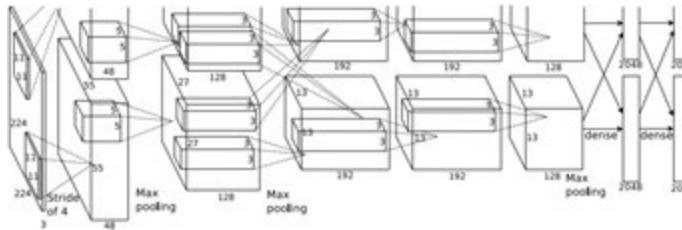
DUCK:  $(x, y, w, h)$

Many numbers!

...

# Object Detection: Multiple Objects

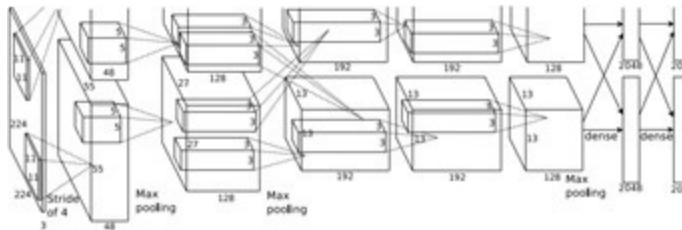
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? NO  
Background? YES

# Object Detection: Multiple Objects

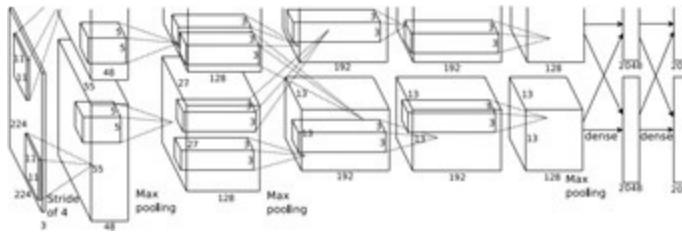
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection: Multiple Objects

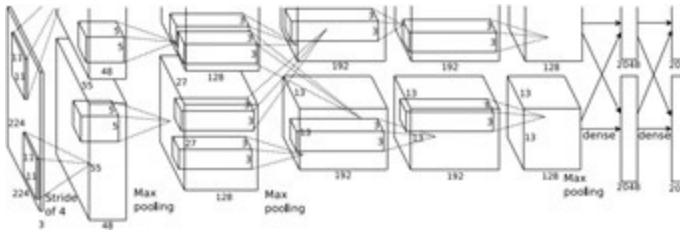
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

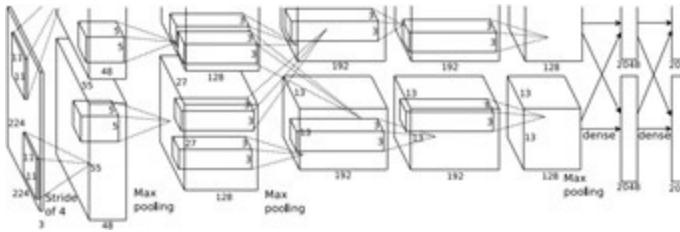
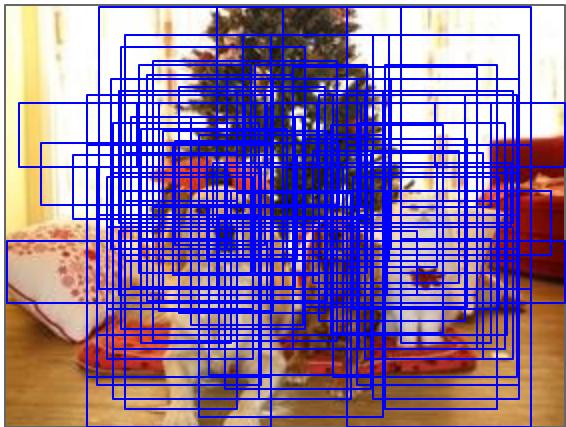


Dog? NO  
Cat? YES  
Background? NO

Q: What's the problem with this approach?

# Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

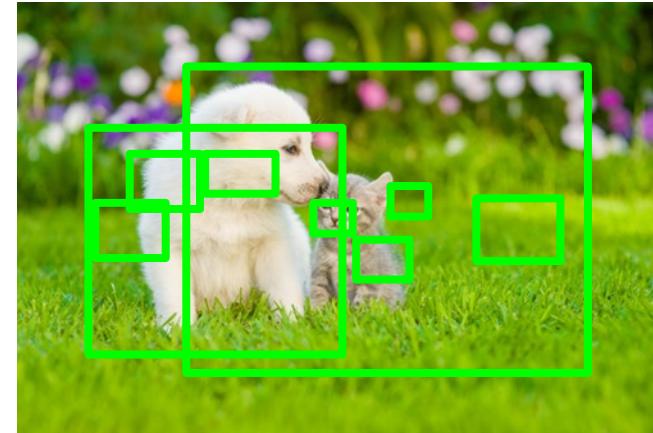


Dog? NO  
Cat? YES  
Background? NO

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!

# Region Proposals: Selective Search

- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



Alex et al, "Measuring the objectness of image windows", TPAMI 2012

Uijlings et al, "Selective Search for Object Recognition", IJCV 2013

Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014

Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

# R-CNN

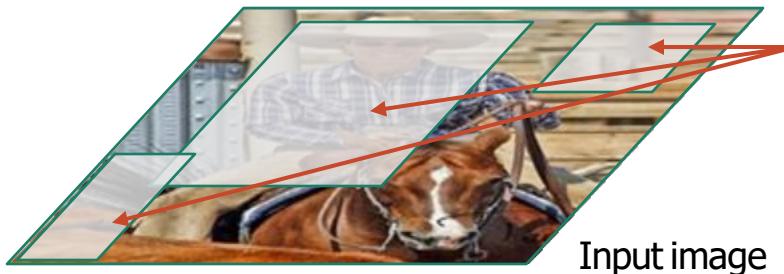


Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# R-CNN



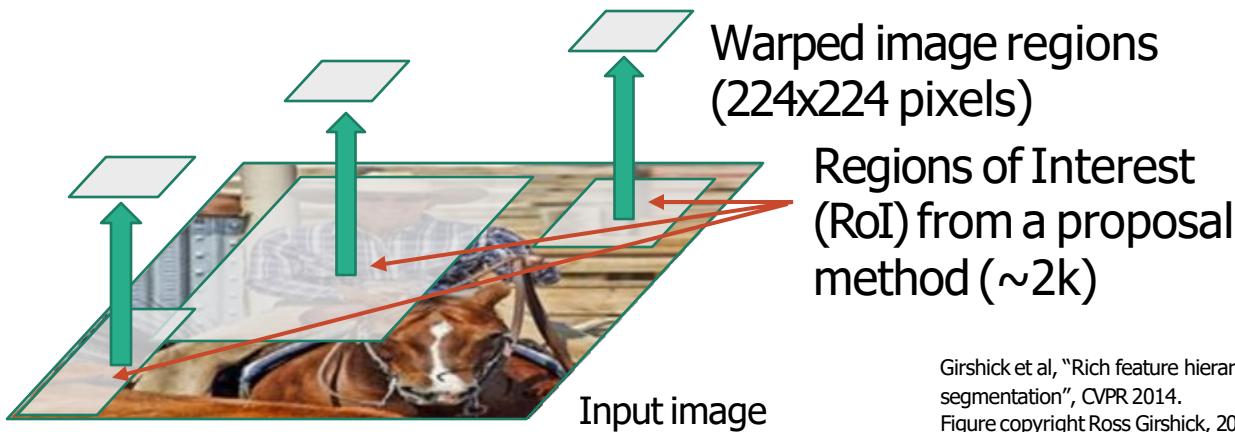
Input image

Regions of Interest  
(RoI) from a proposal  
method ( $\sim 2k$ )

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

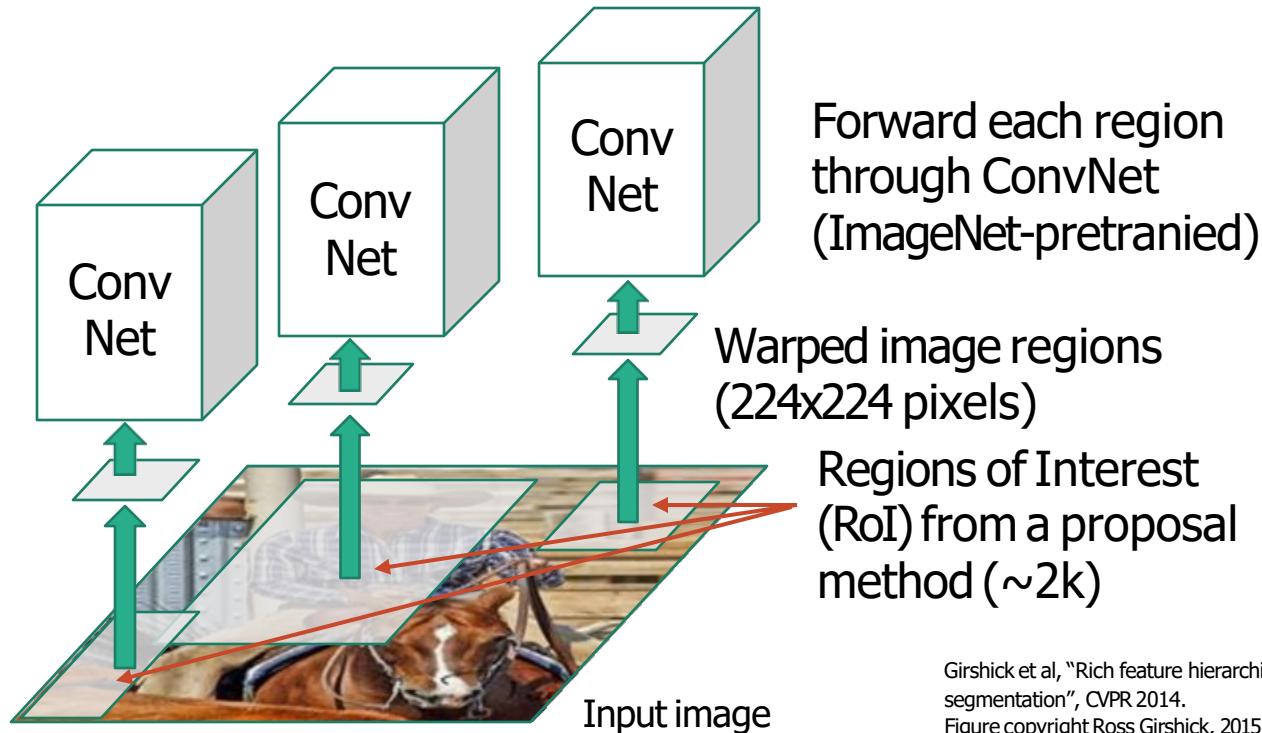
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# R-CNN



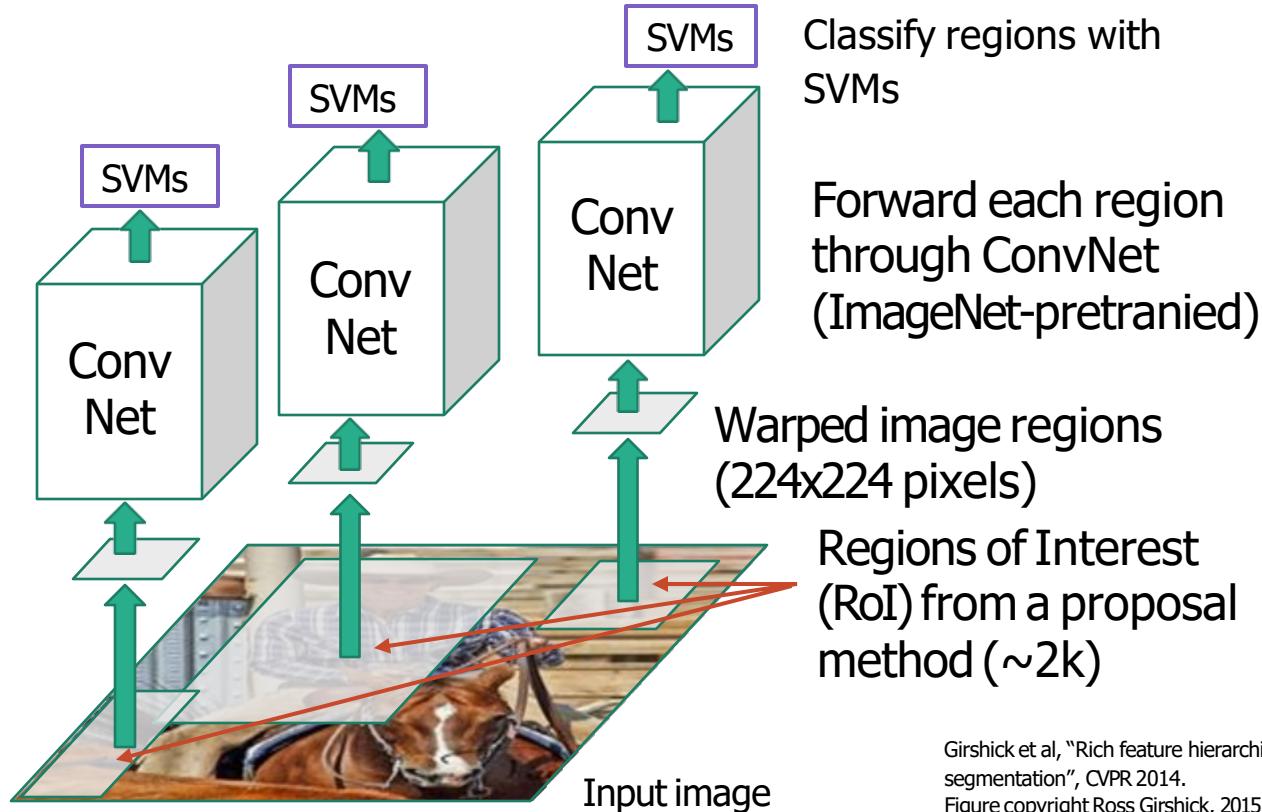
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

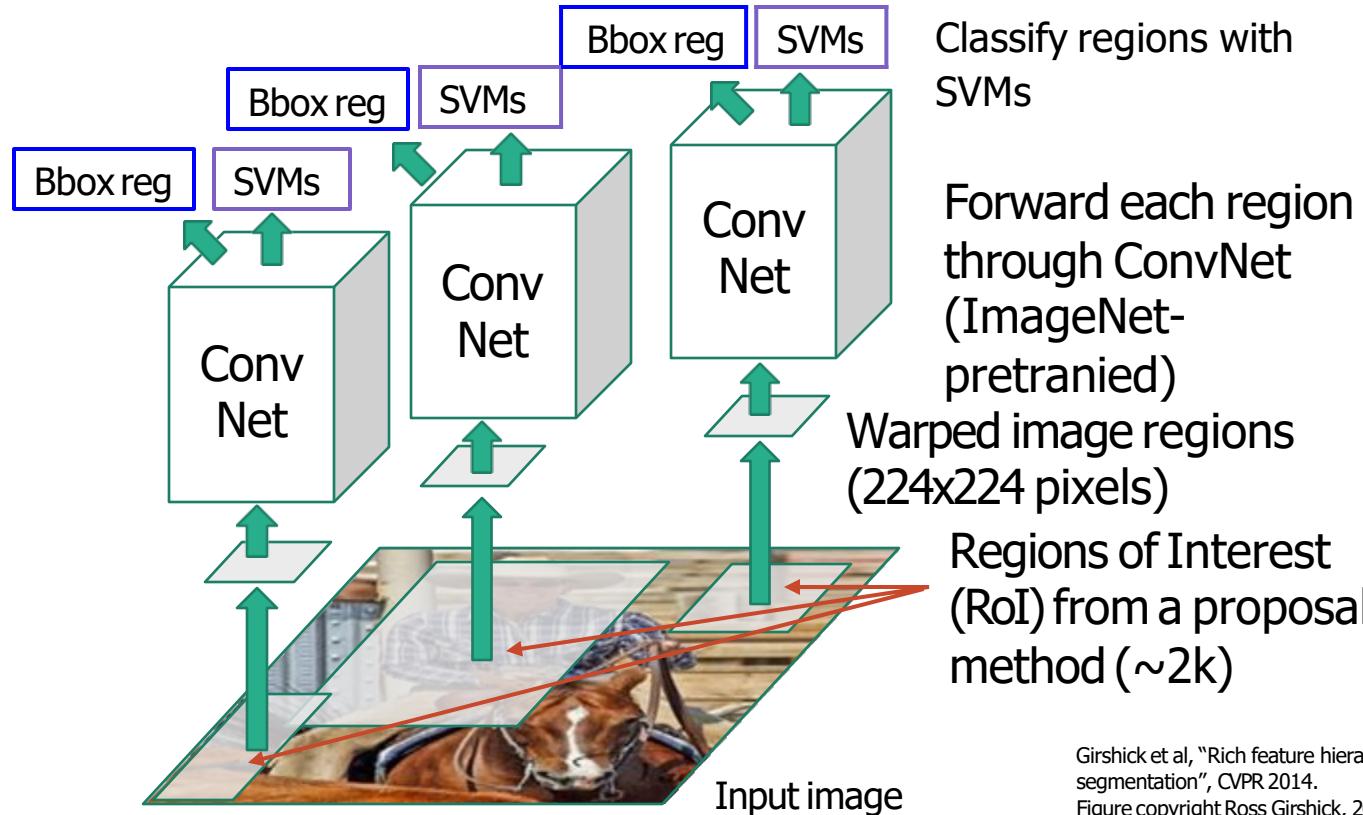
# R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

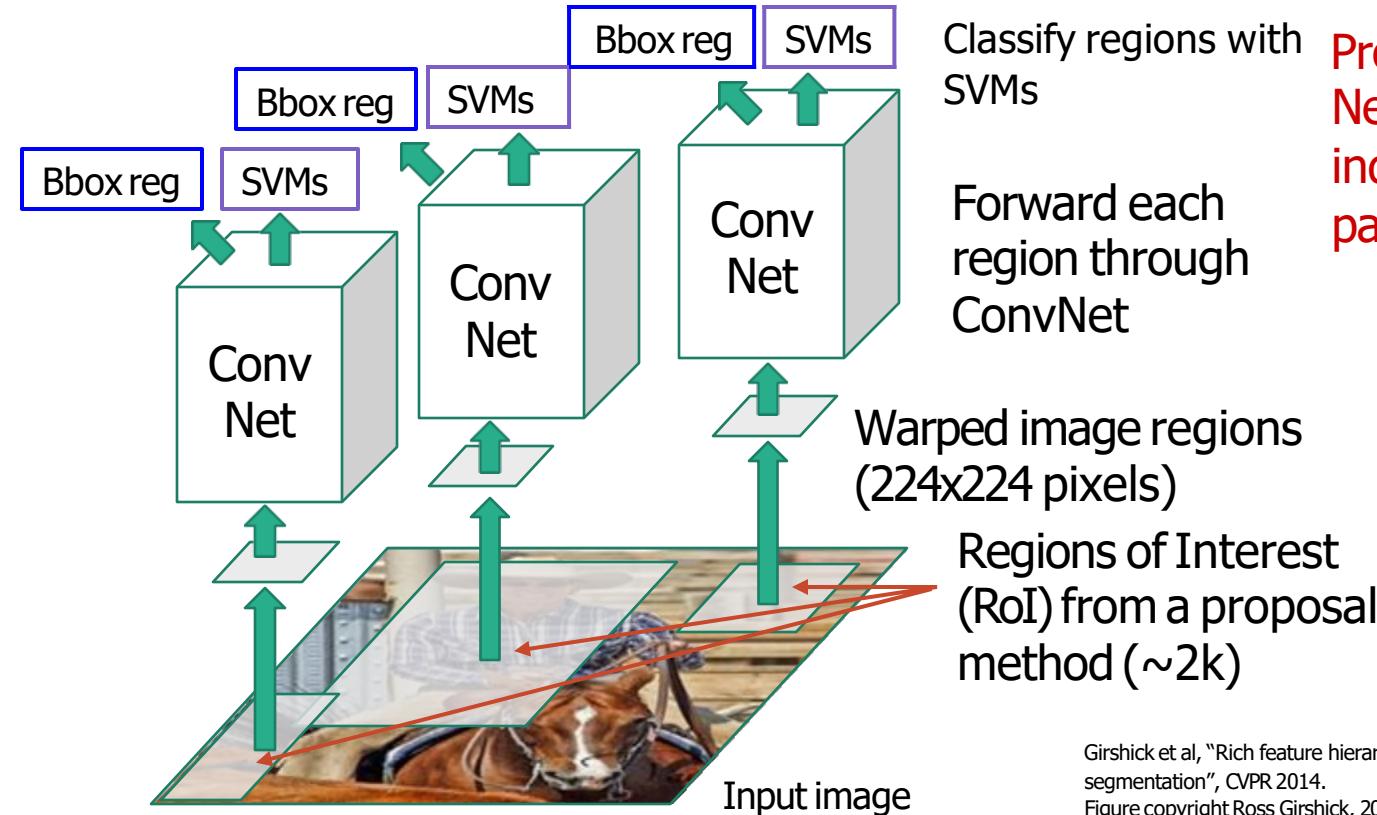
# R-CNN

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# R-CNN



Predict “corrections” to the RoI: 4 numbers: ( $dx, dy, dw, dh$ )

Classify regions with SVMs

Forward each region through ConvNet

Warped image regions (224x224 pixels)

Regions of Interest (RoI) from a proposal method ( $\sim 2k$ )

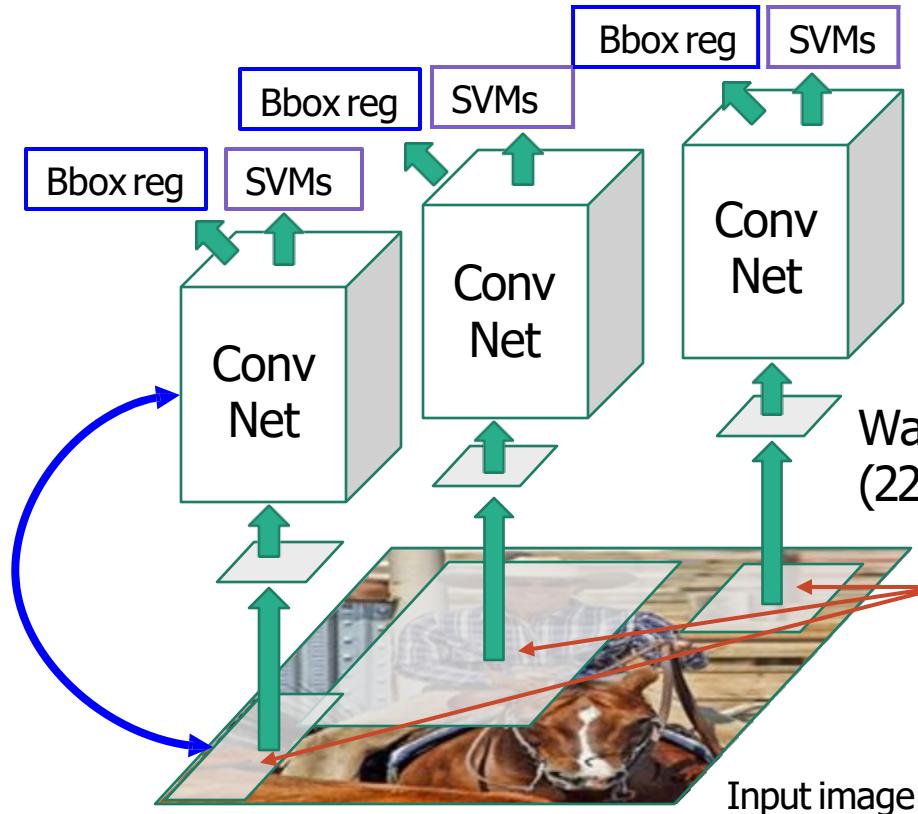
**Problem: Very slow!**  
Need to do  $\sim 2k$  independent forward passes for each image!

Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# “Slow” R-CNN

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



Classify regions with SVMs

Forward each region through ConvNet

Warped image regions (224x224 pixels)

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

**Problem:** Very slow!  
Need to do ~2k independent forward passes for each image!

Idea: Pass the image through convnet before cropping! Crop the conv feature instead!

Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

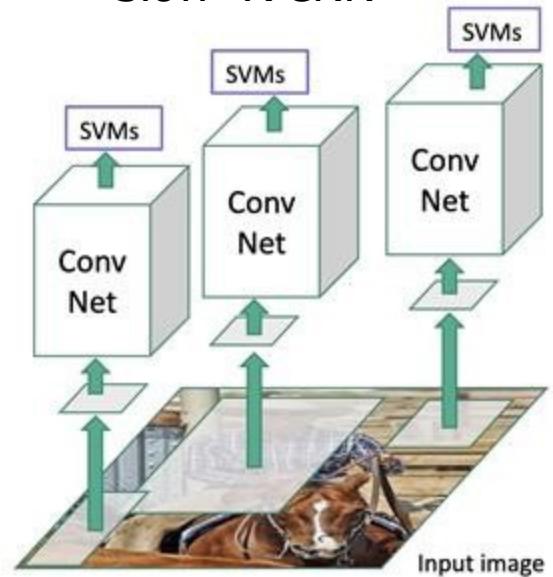
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN



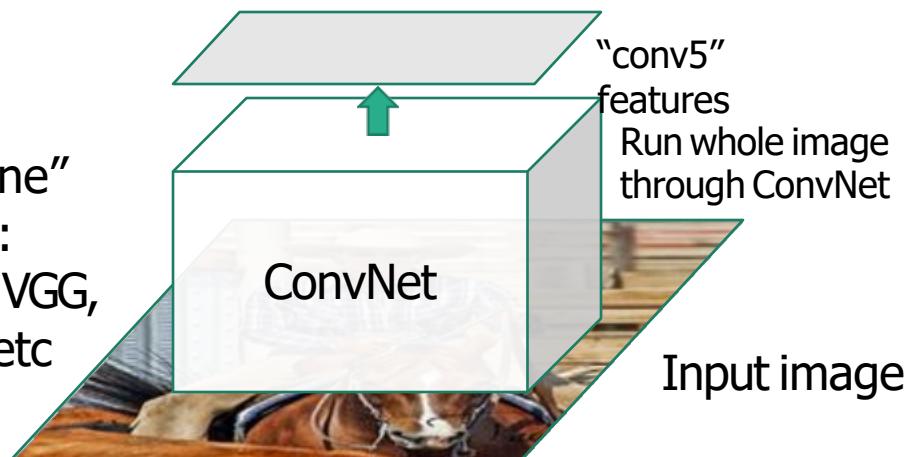
Input image

“Slow” R-CNN

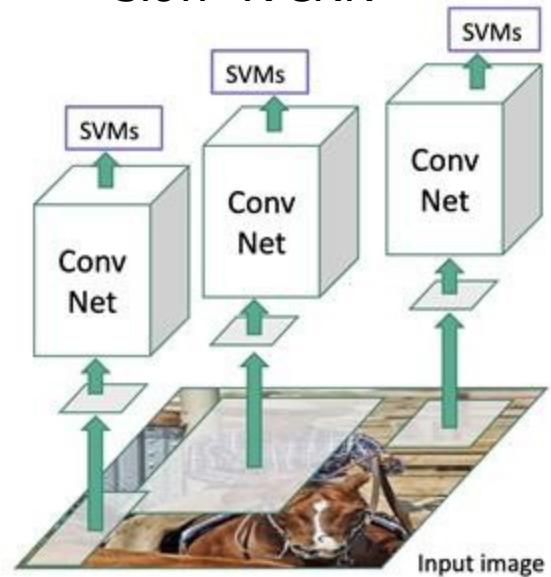


# Fast R-CNN

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc



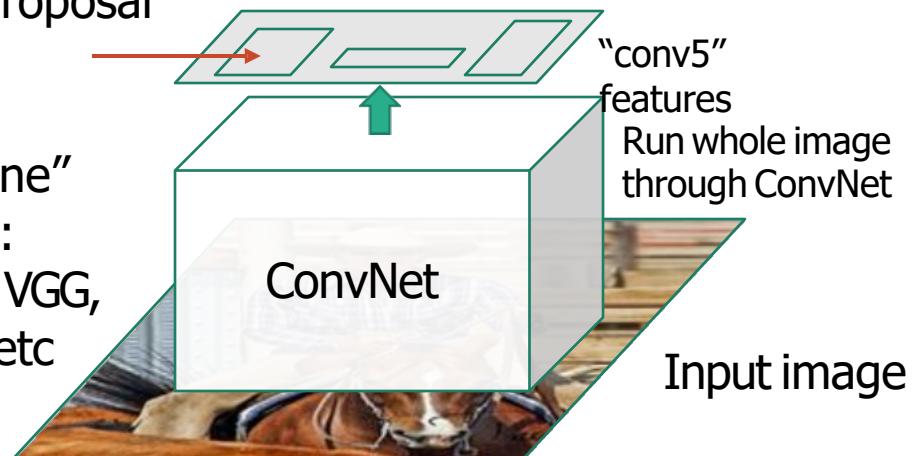
## “Slow” R-CNN



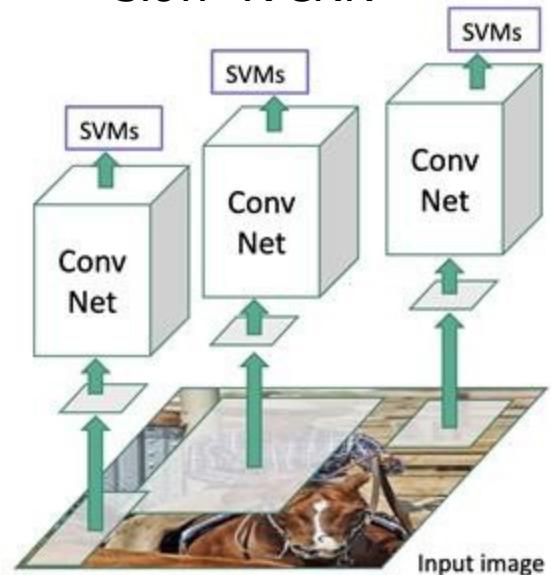
# Fast R-CNN

Regions of Interest (RoIs)  
from a proposal  
method

“Backbone”  
network:  
AlexNet, VGG,  
ResNet, etc



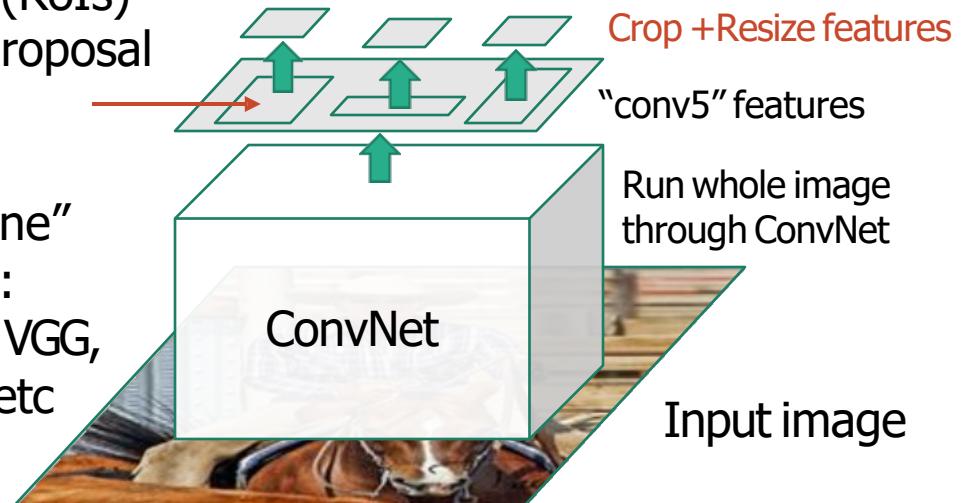
“Slow” R-CNN



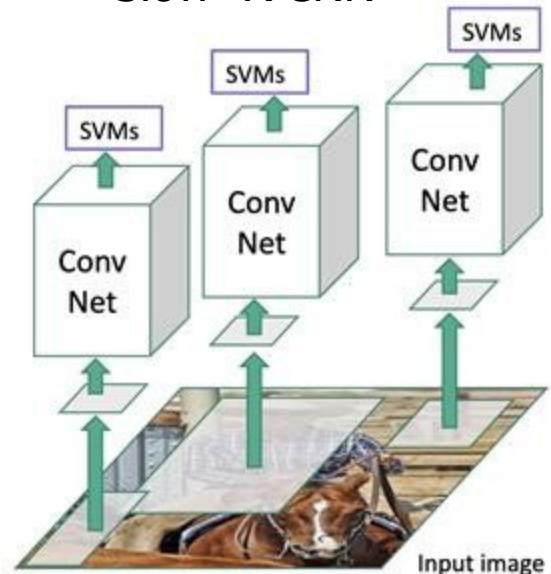
# Fast R-CNN

Regions of Interest (RoIs) from a proposal method

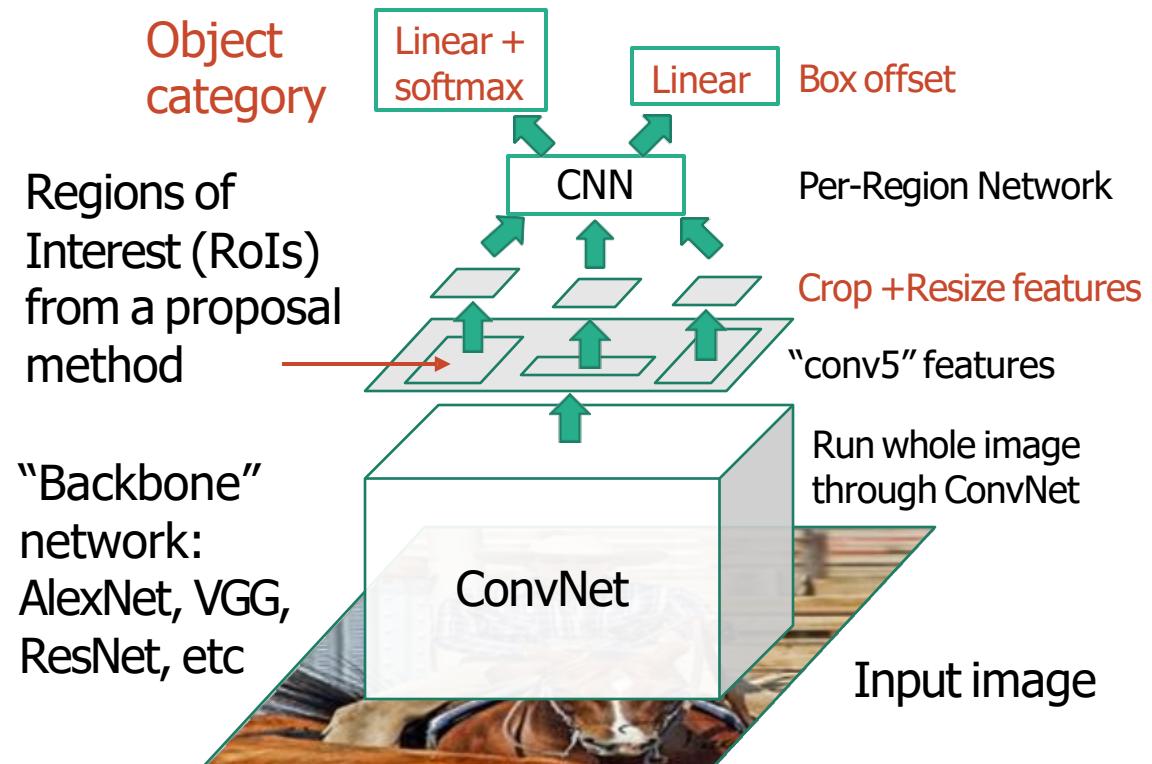
“Backbone” network:  
AlexNet, VGG,  
ResNet, etc



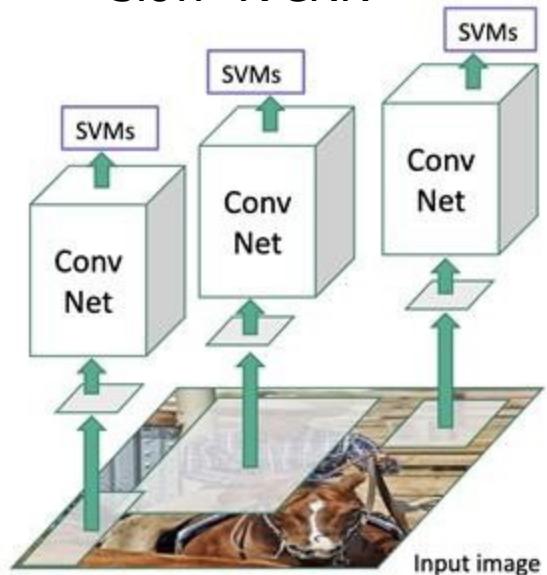
## “Slow” R-CNN



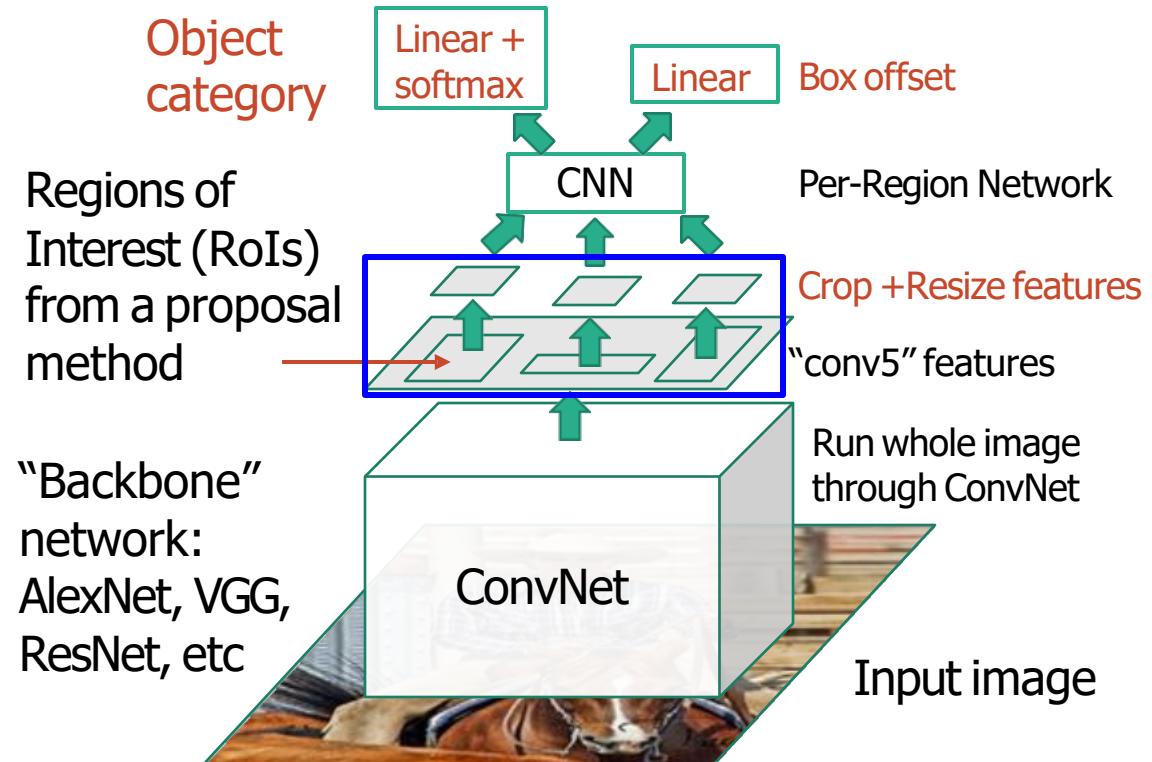
# Fast R-CNN



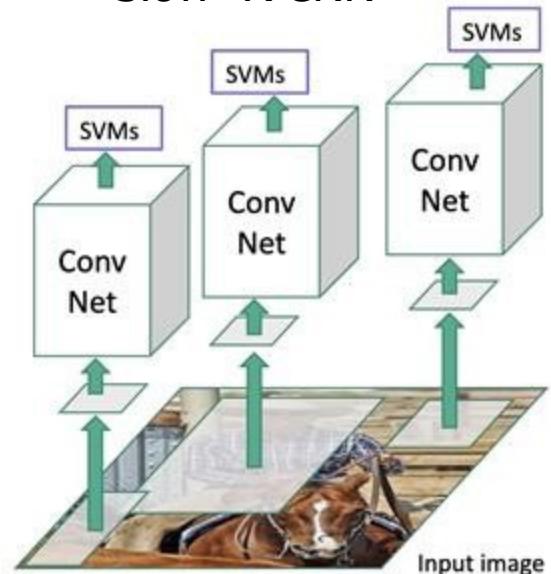
“Slow” R-CNN



# Fast R-CNN



## “Slow” R-CNN



# Region Proposal Network



Input Image  
(e.g.  $3 \times 640 \times 480$ )

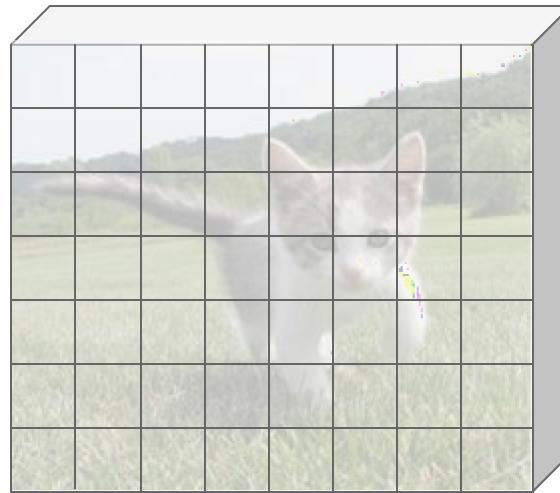
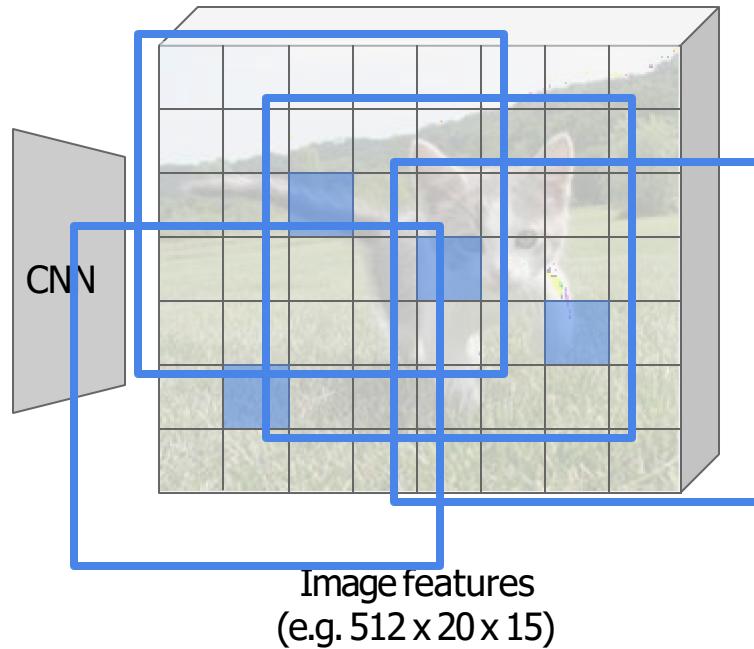


Image features  
(e.g.  $512 \times 20 \times 15$ )

# Region Proposal Network



Input Image  
(e.g.  $3 \times 640 \times 480$ )

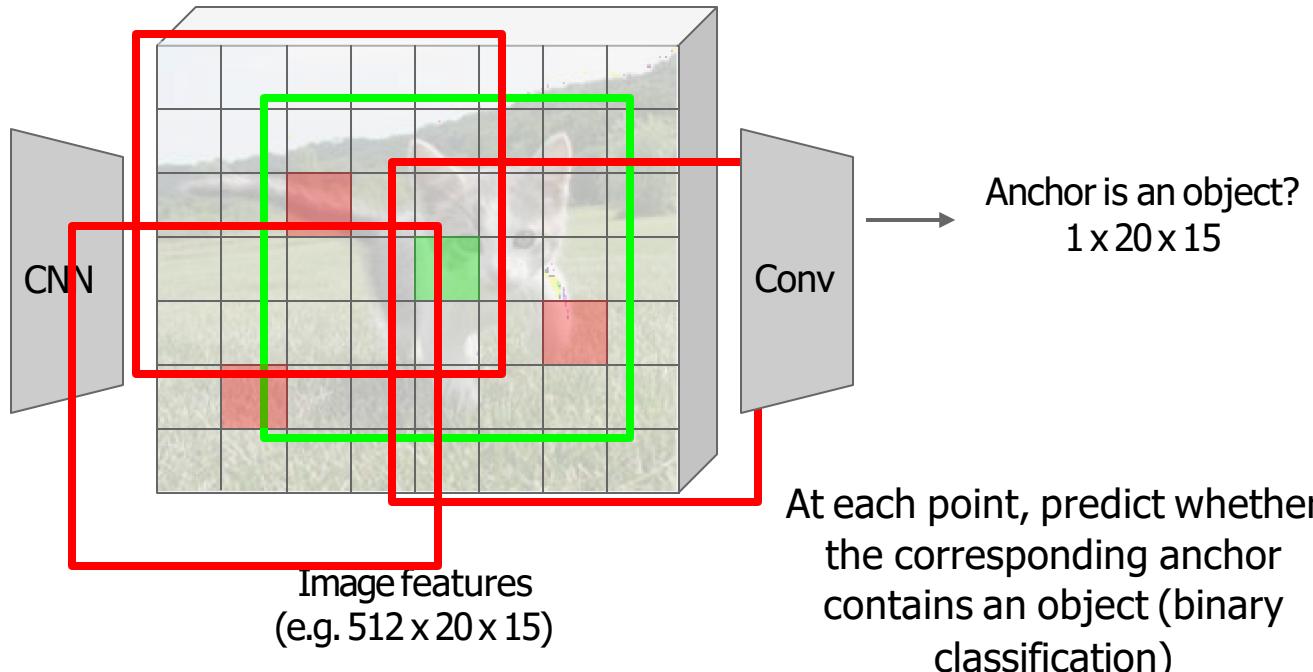


Imagine an anchor box of fixed size at each point in the feature map

# Region Proposal Network



Input Image  
(e.g.  $3 \times 640 \times 480$ )



# Region Proposal Network



Input Image  
(e.g.  $3 \times 640 \times 480$ )

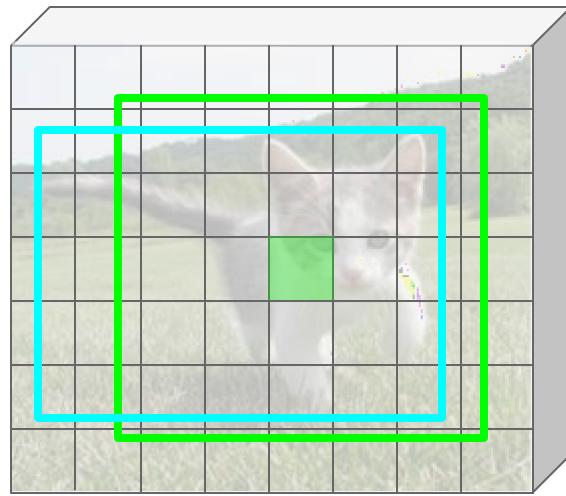


Image features  
(e.g.  $512 \times 20 \times 15$ )



Imagine an anchor box of fixed size at each point in the feature map

Anchor is an object?  
 $1 \times 20 \times 15$

Box corrections  
 $4 \times 20 \times 15$

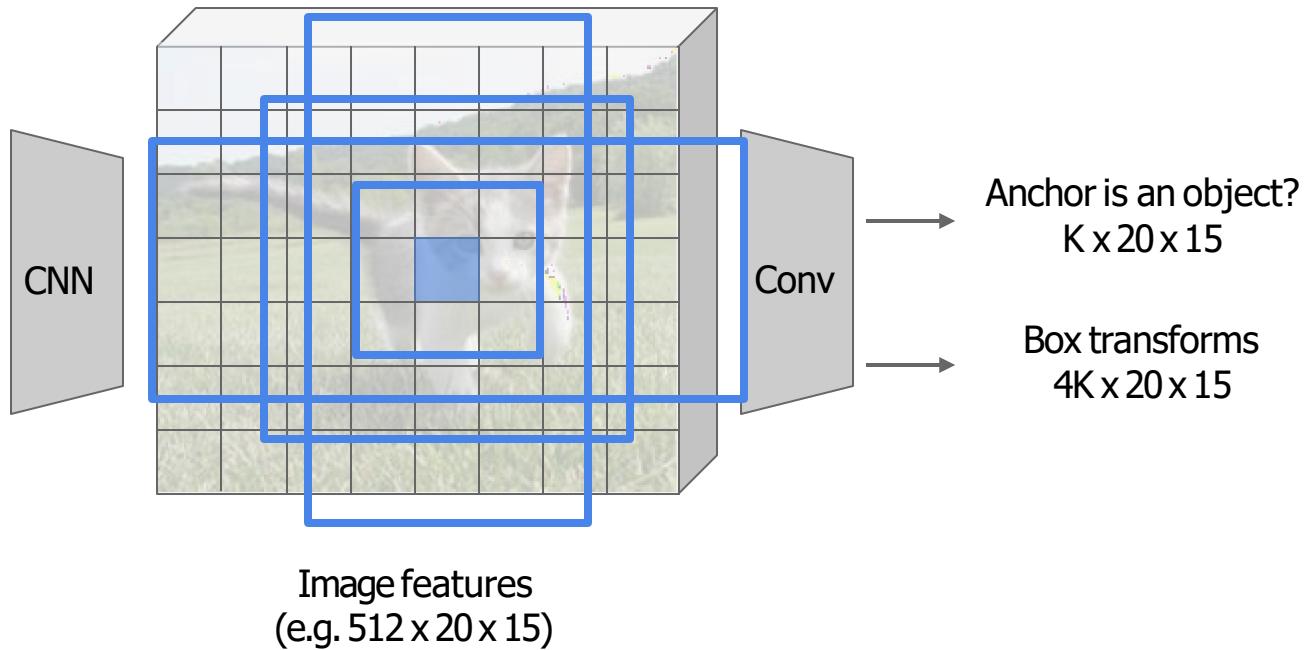
For positive boxes, also predict a corrections from the anchor to the ground-truth box (regress 4 numbers per pixel)

# Region Proposal Network

In practice use K different anchor boxes of different size / scale at each point



Input Image  
(e.g.  $3 \times 640 \times 480$ )

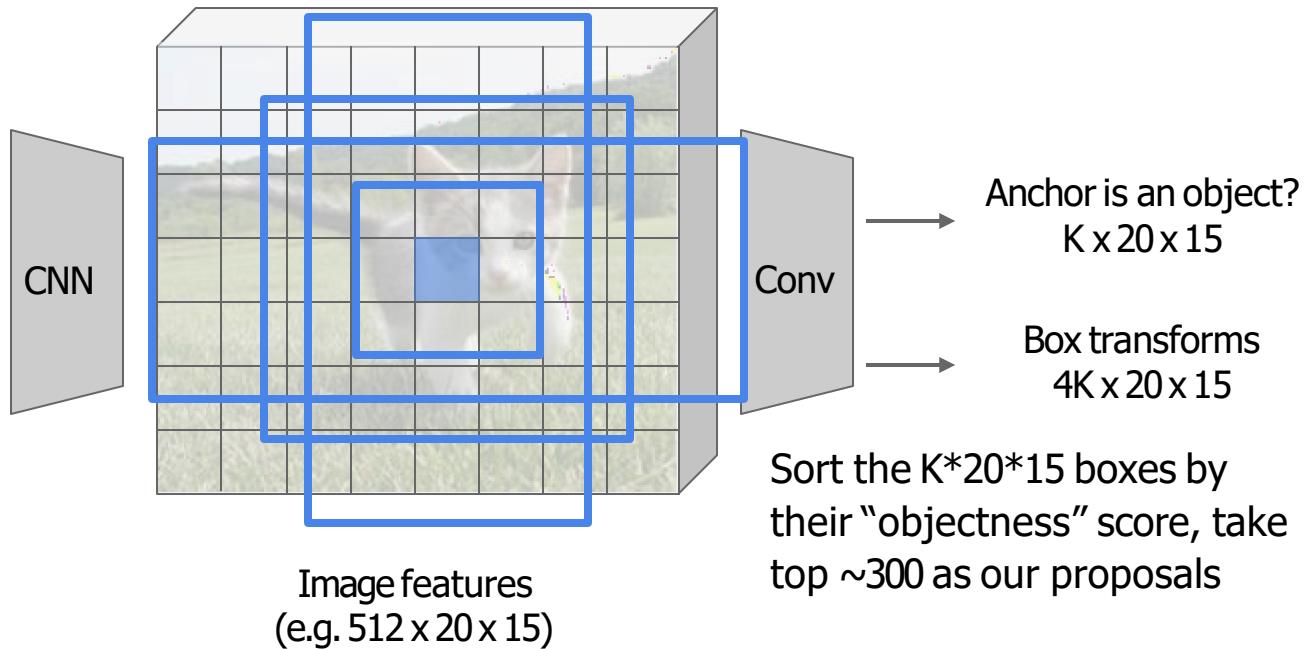


# Region Proposal Network

In practice use K different anchor boxes of different size / scale at each point



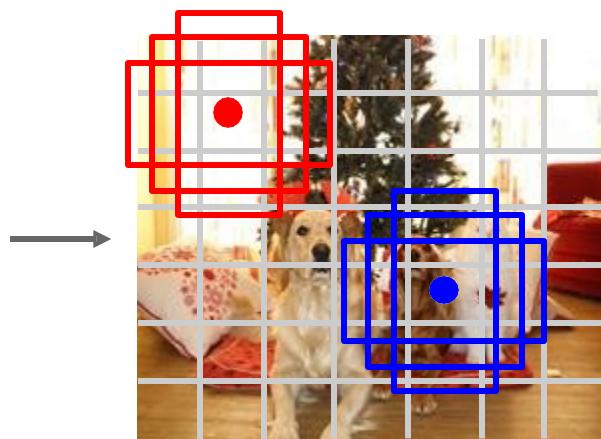
Input Image  
(e.g.  $3 \times 640 \times 480$ )



# Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

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

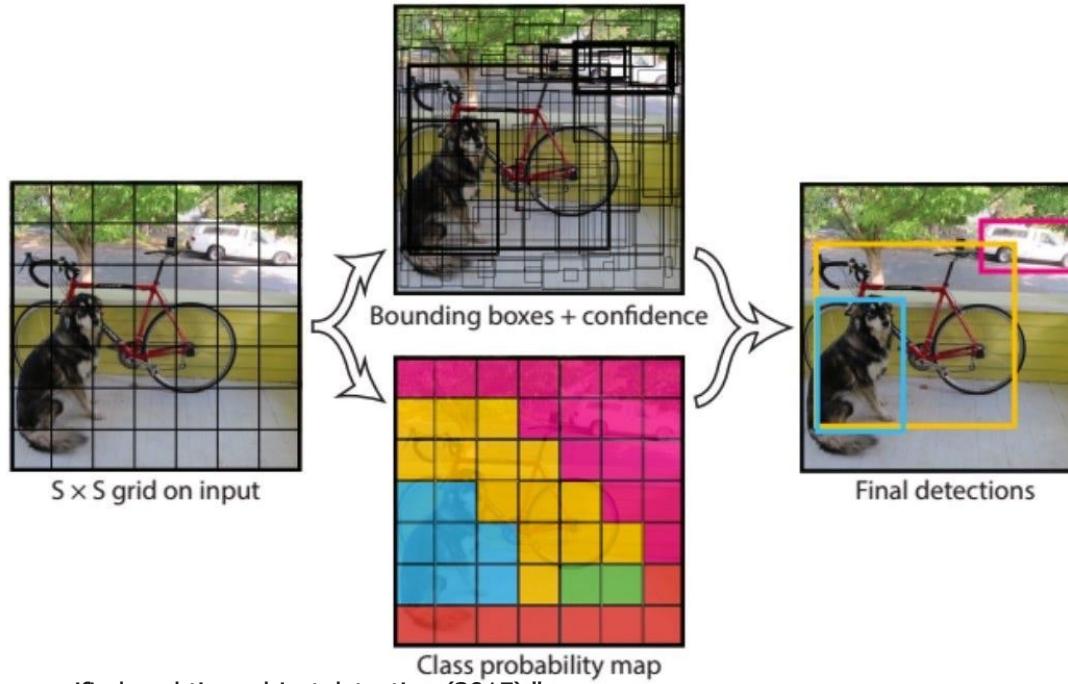
Within each grid cell:

- 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)
- Looks a lot like RPN, but category-specific!

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  
Lin et al, "Focal Loss for Dense Object Detection", ICCV 2017

# YOLO (You Only Look Once) real-time object detection



Redmon et al. "You only look once: unified, real-time object detection (2015)."

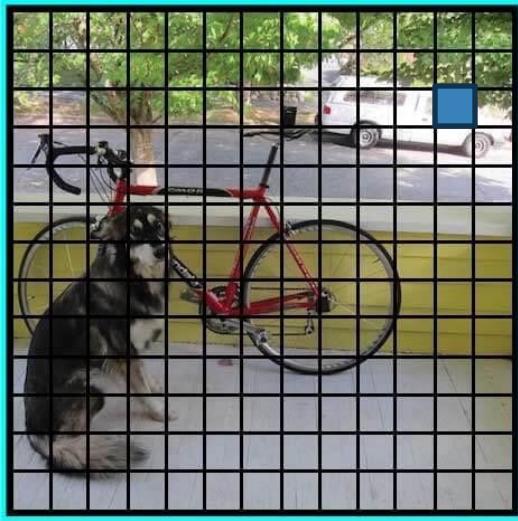
# YOLO



SxS Grid

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO



SxS Grid

For each box output:

- $P(\text{object})$ : probability that the box contains an object
- $B$  bounding boxes ( $x, y, h, w$ )
- $P(\text{class})$ : probability of belonging to a class

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO



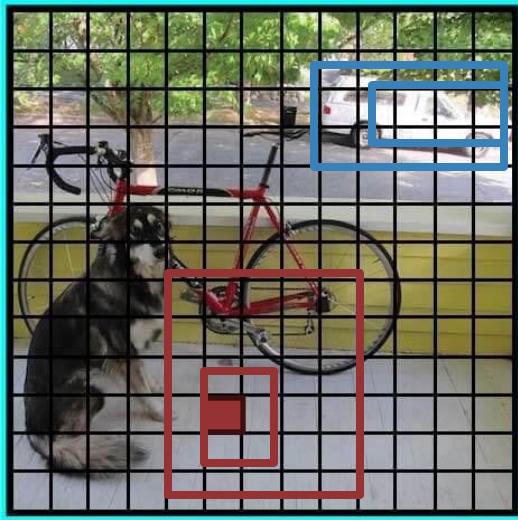
SxS Grid

For each box output:

- $P(\text{object})$ : probability that the box contains an object
- $B$  bounding boxes ( $x, y, h, w$ )
- $P(\text{class})$ : probability of belonging to a class

$$B=2$$

# YOLO



For each box output:

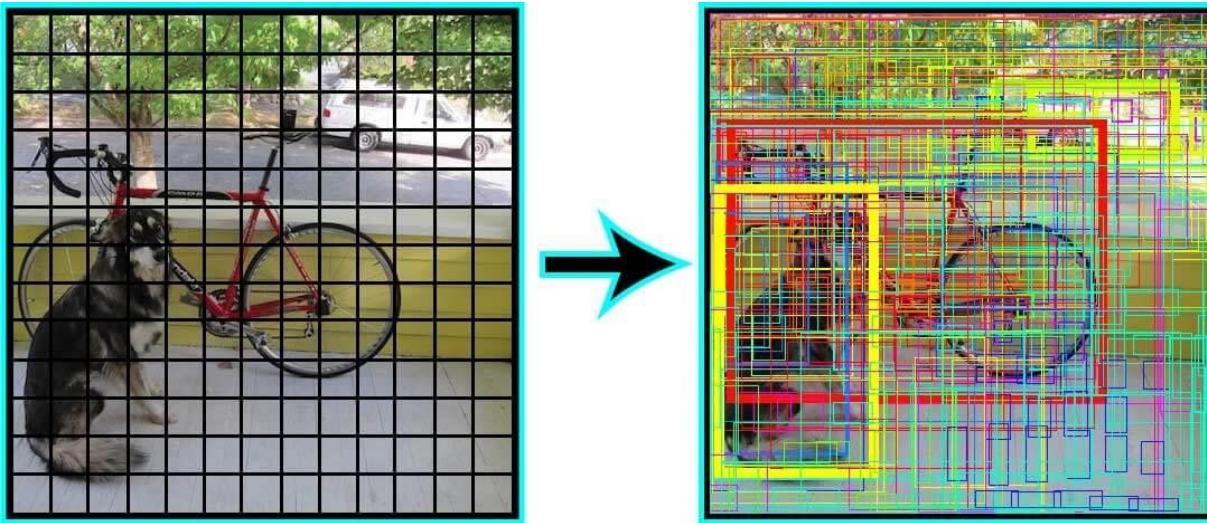
- $P(\text{object})$ : probability that the box contains an object
- $B$  bounding boxes ( $x, y, h, w$ )
- $P(\text{class})$ : probability of belonging to a class

$$B=2$$

SxS Grid

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO

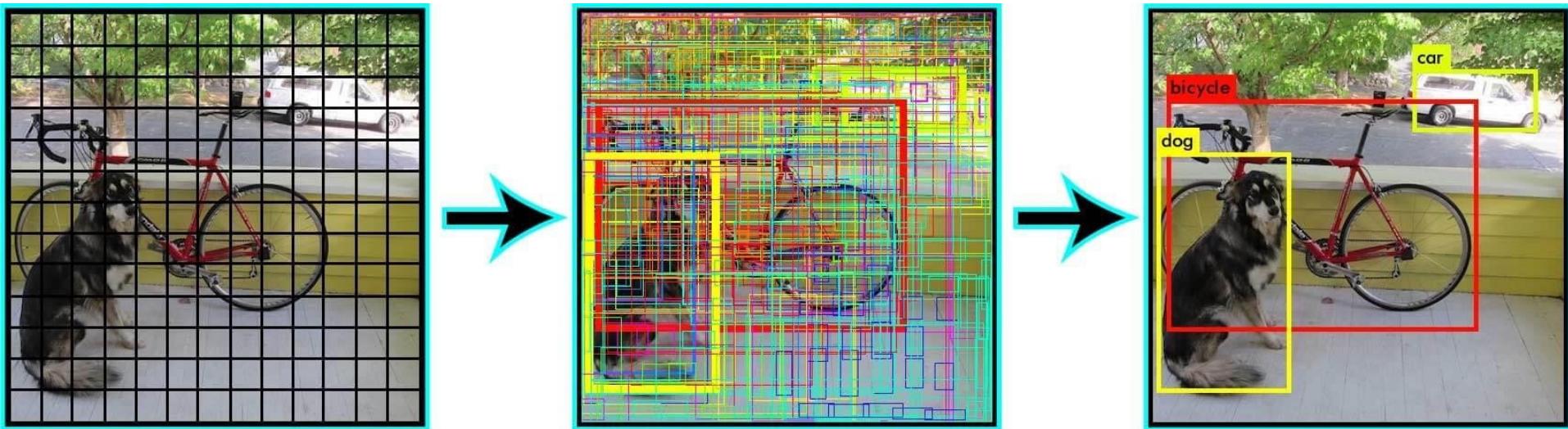


SxS Grid

Many bounding boxes with  
different object probabilities

Redmon et al. "You only look once: unified, real-time object detection (2015)."

# YOLO



$S \times S$  Grid

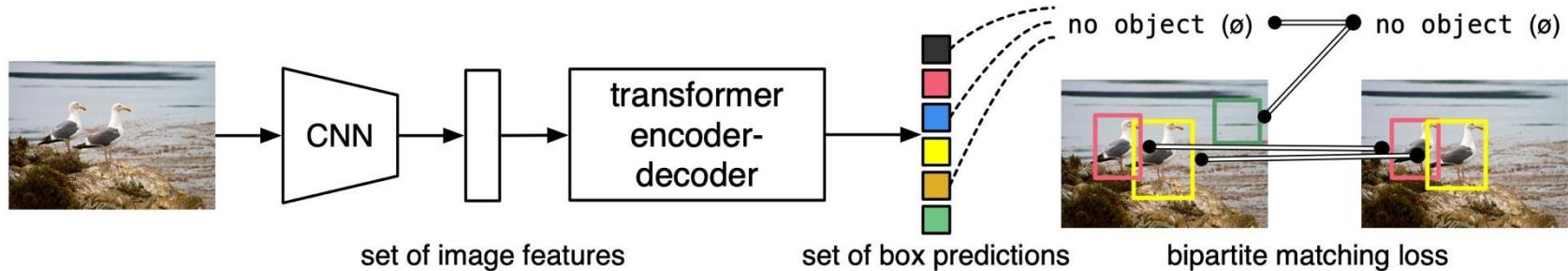
Redmon et al. "You only look once: unified, real-time object detection (2015)."

# Object Detection with Transformers: DETR

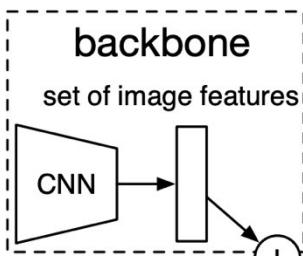
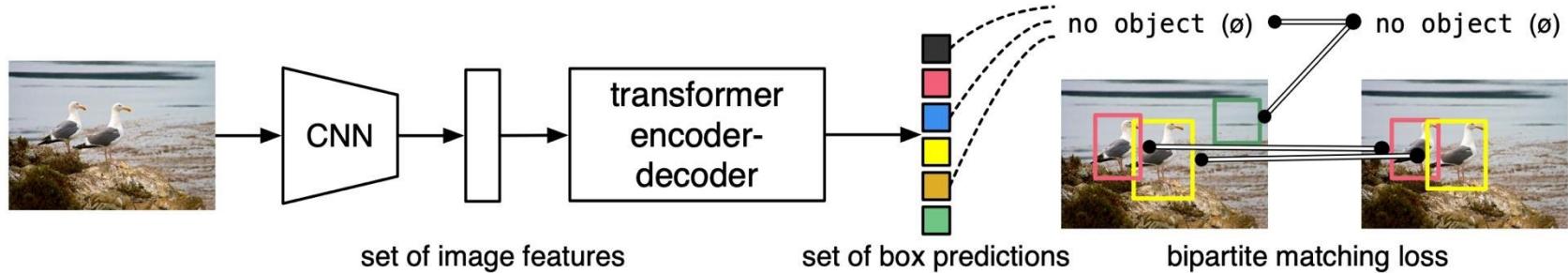
Simple object detection pipeline: directly output a set of boxes from a Transformer

No anchors, no regression of box transforms

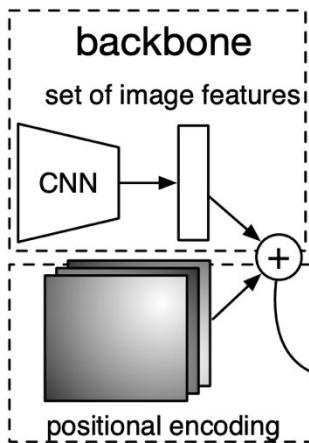
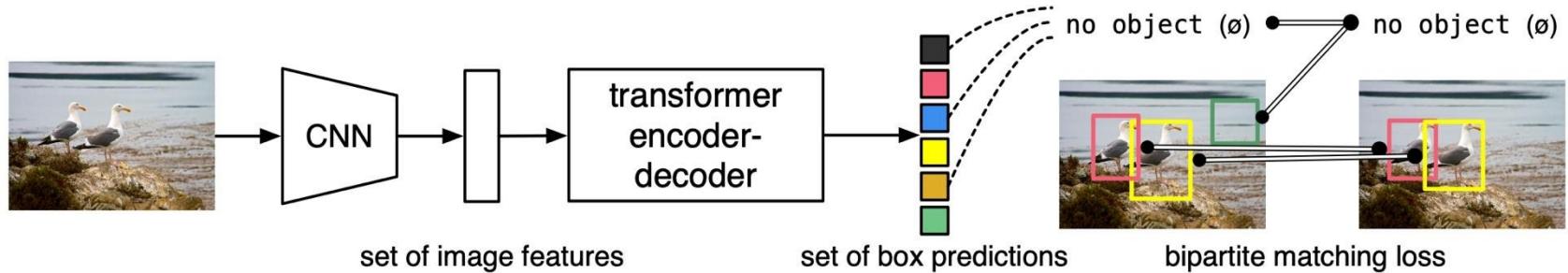
Match predicted boxes to GT boxes with bipartite matching; train to regress box coordinates



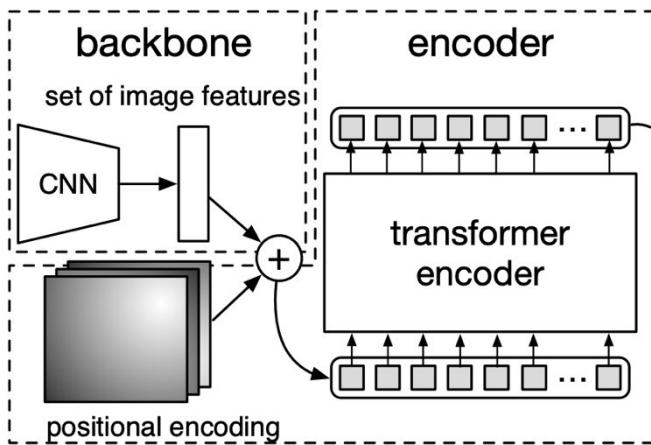
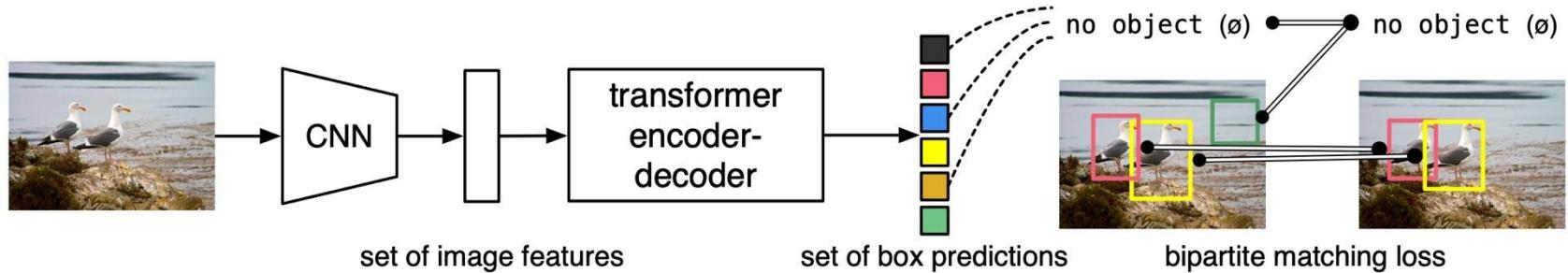
# Object Detection with Transformers: DETR



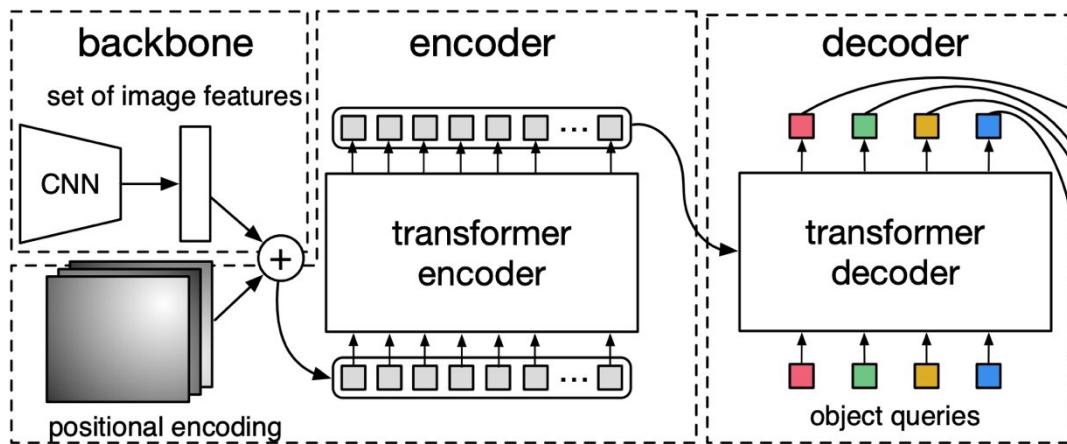
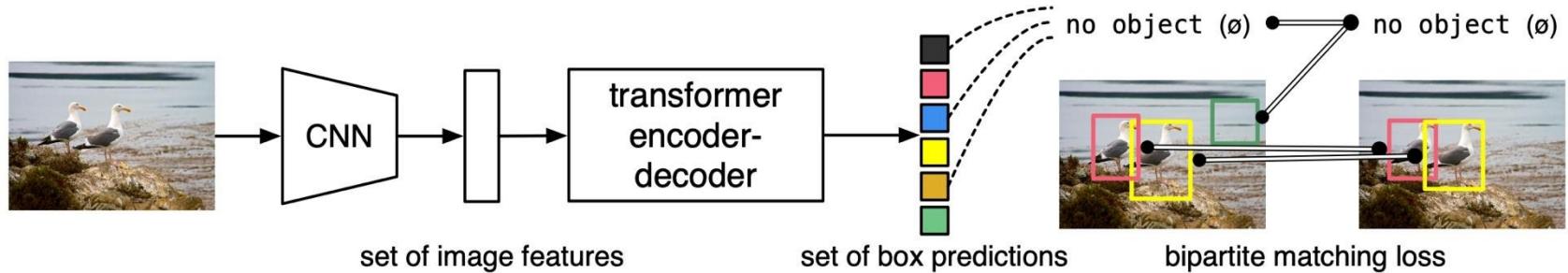
# Object Detection with Transformers: DETR



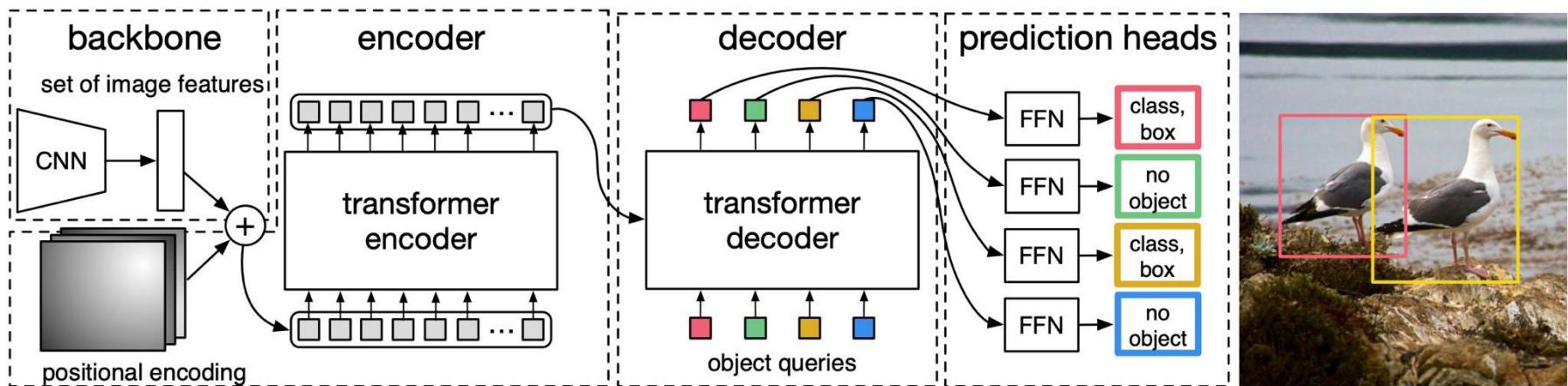
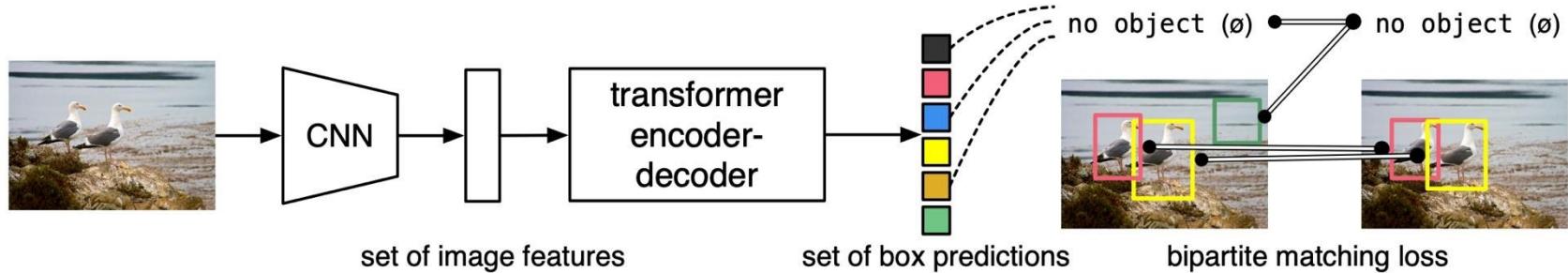
# Object Detection with Transformers: DETR



# Object Detection with Transformers: DETR



# Object Detection with Transformers: DETR



# Instance Segmentation

Classification



CAT

No spatial extent

Semantic  
Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object  
Detection



DOG, DOG, CAT

Multiple Object

Instance  
Segmentation



DOG, DOG, CAT

# Object Detection: Faster R-CNN

**Object  
Detection**



DOG, DOG, CAT

**Instance  
Segmentation**



DOG, DOG, CAT

Classification  
loss

Bounding-box  
regression loss

Classification  
loss      Bounding-box  
regression loss  
...

proposals

Region Proposal Network

feature map

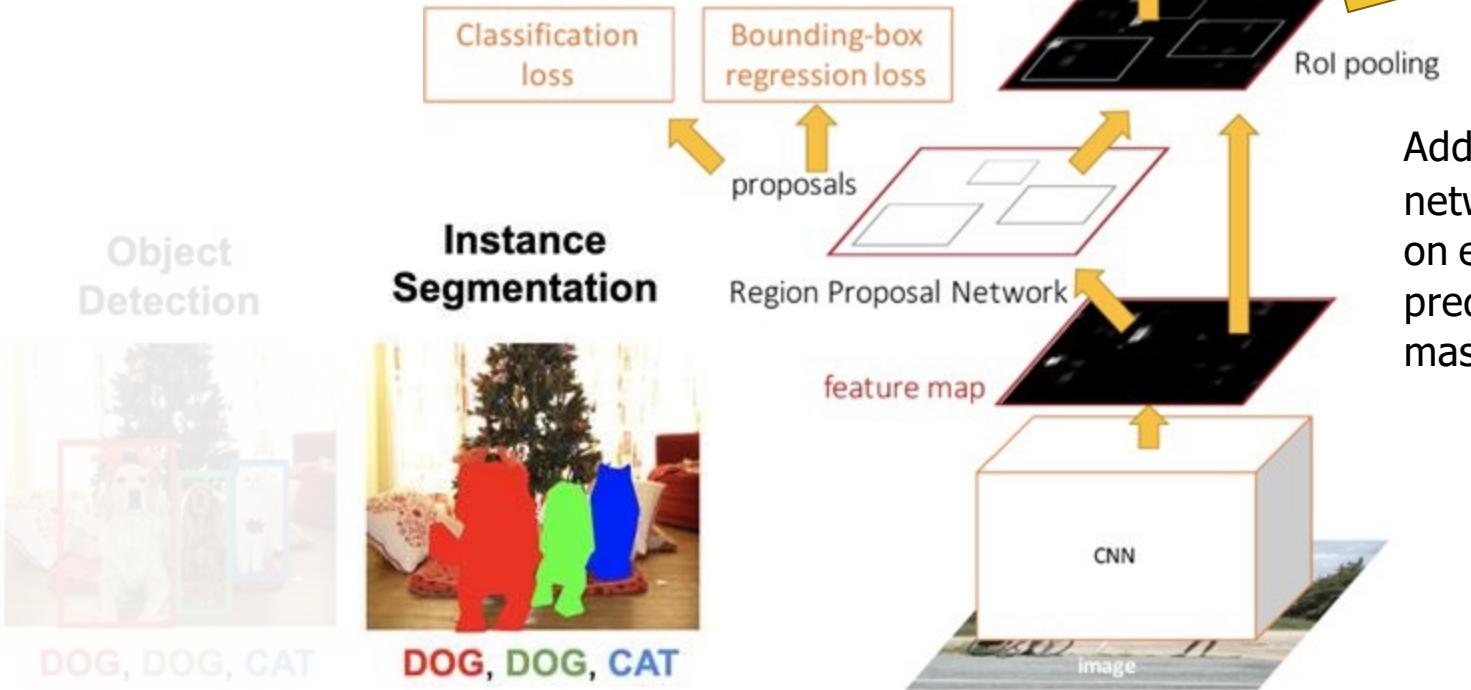
CNN

Image

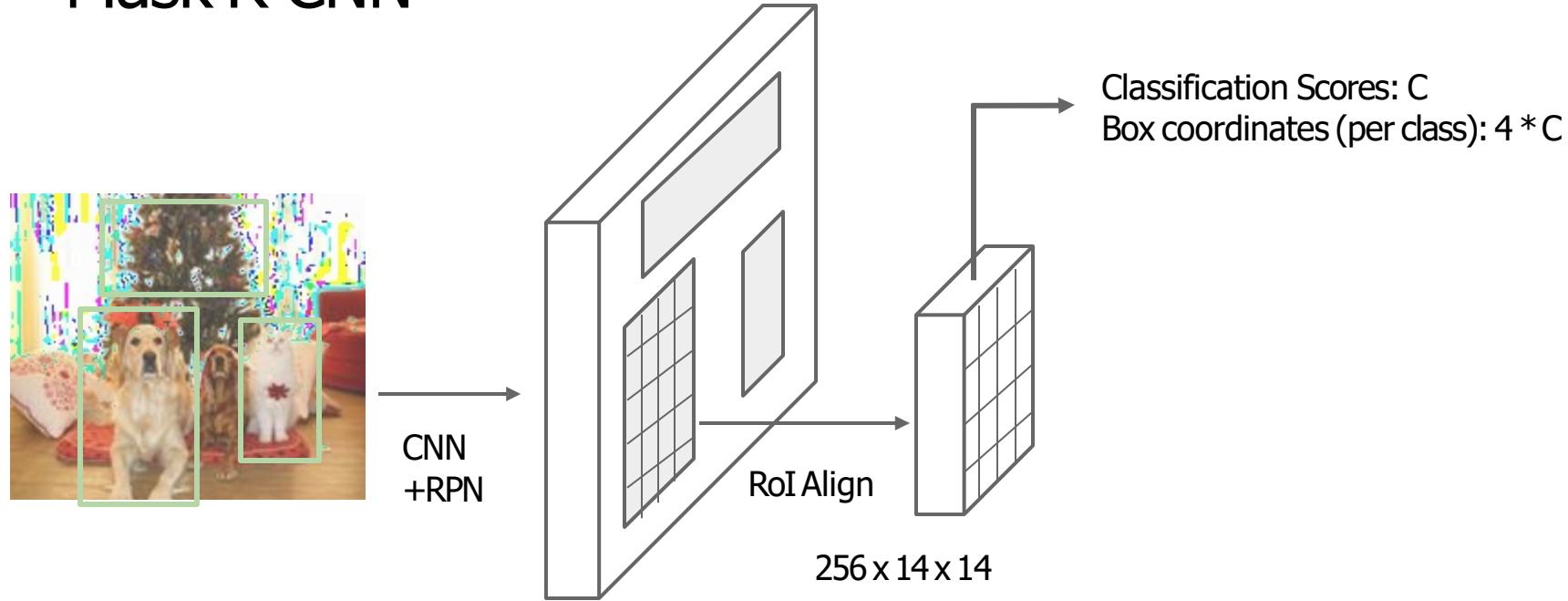
Classification  
loss      Bounding-box  
regression loss  
...

RoI pooling

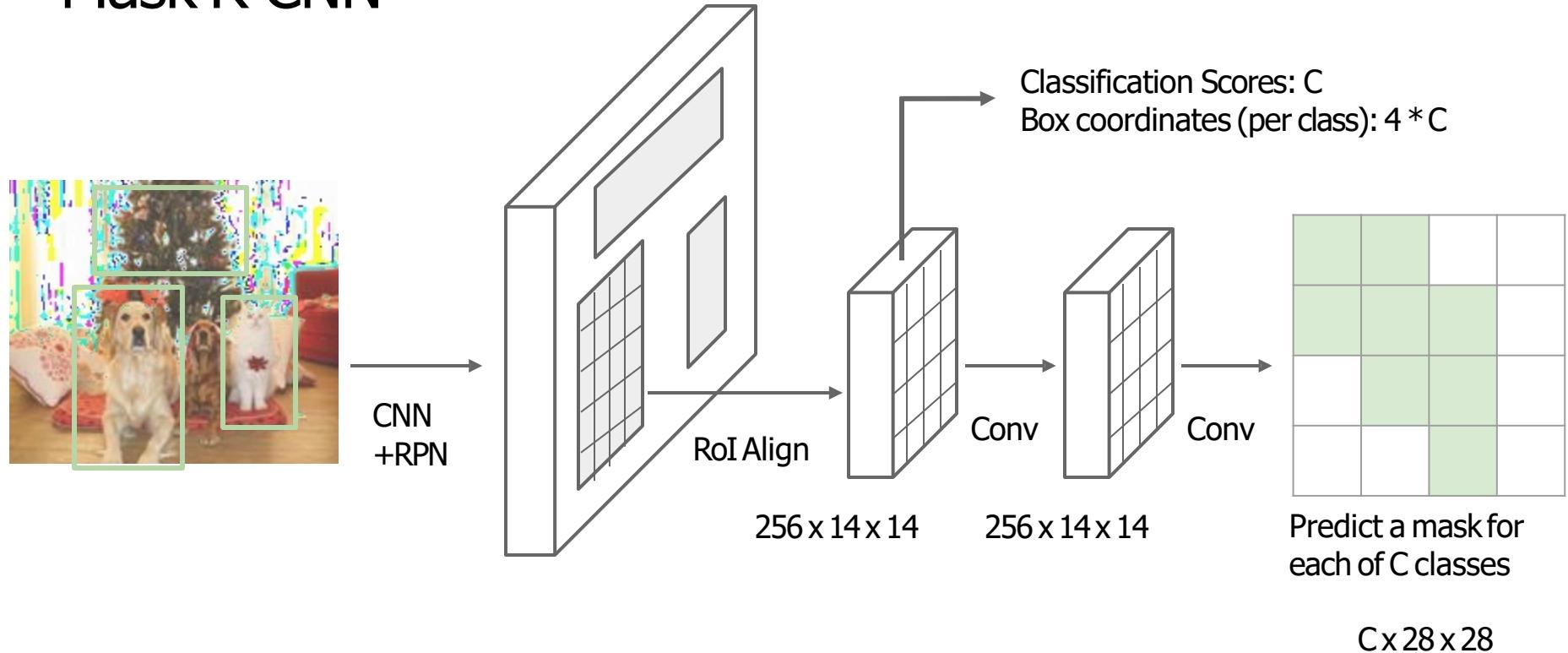
# Instance Segmentation: Mask R-CNN



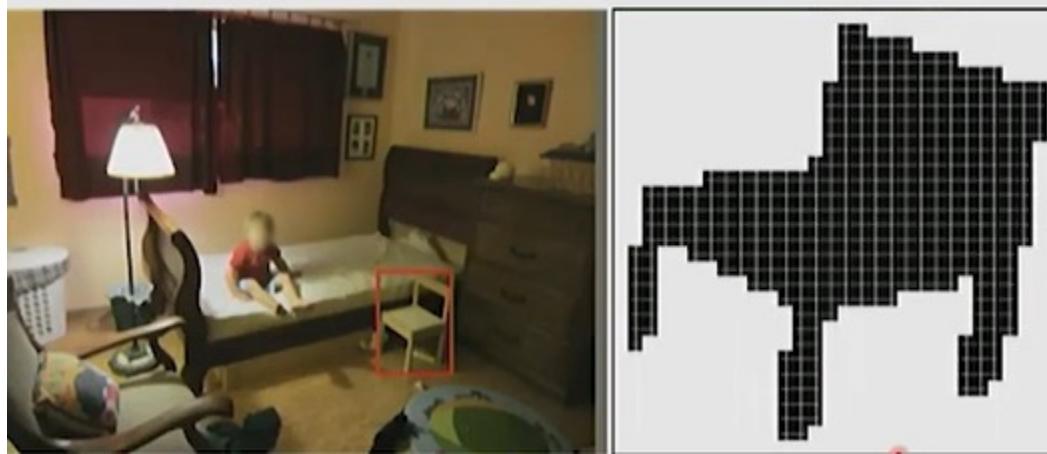
# Mask R-CNN



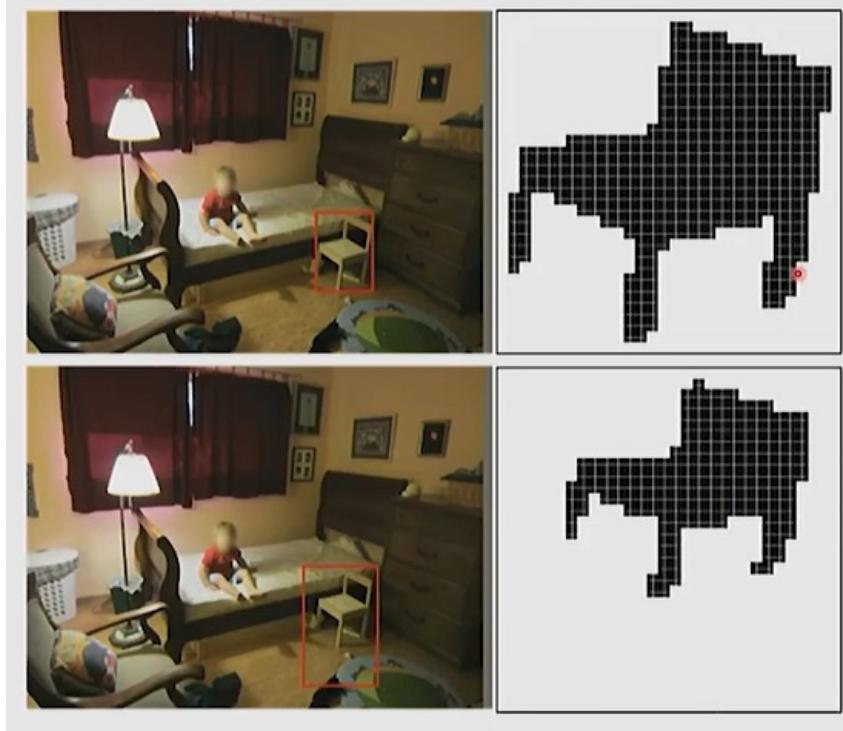
# Mask R-CNN



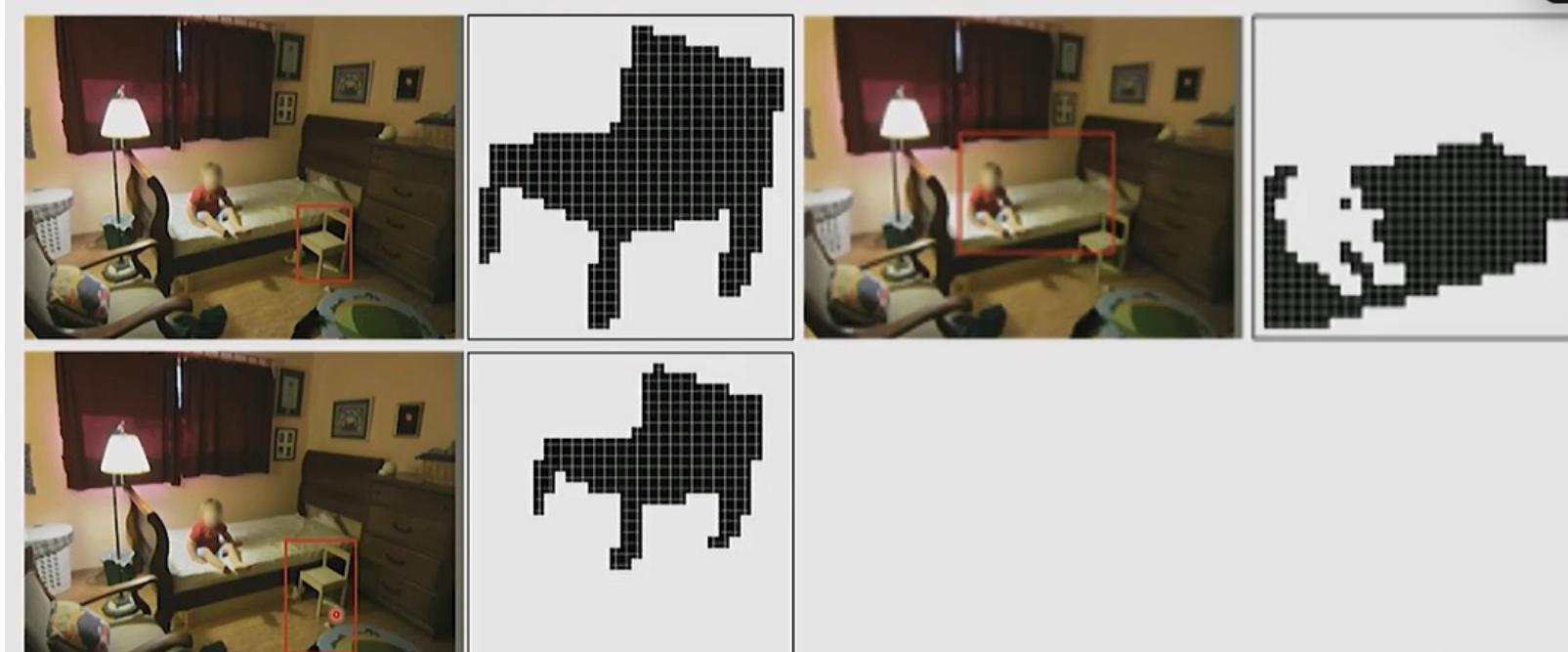
# Mask R-CNN: Example Mask Training Targets



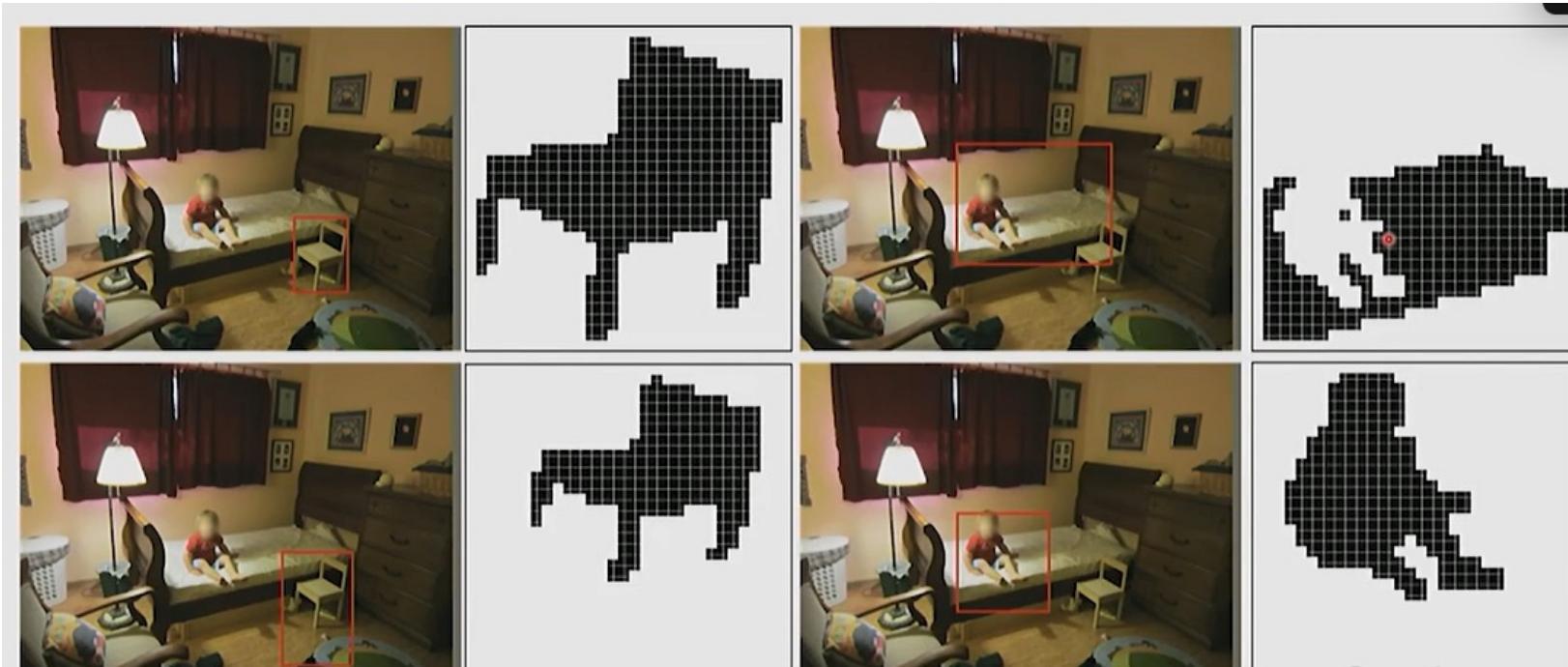
# Mask R-CNN: Example Mask Training Targets



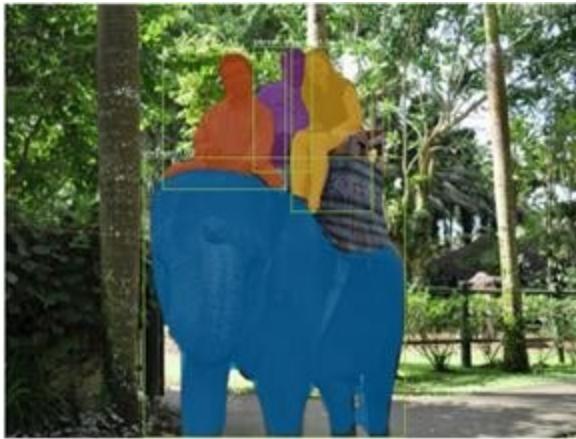
# Mask R-CNN: Example Mask Training Targets



# Mask R-CNN: Example Mask Training Targets

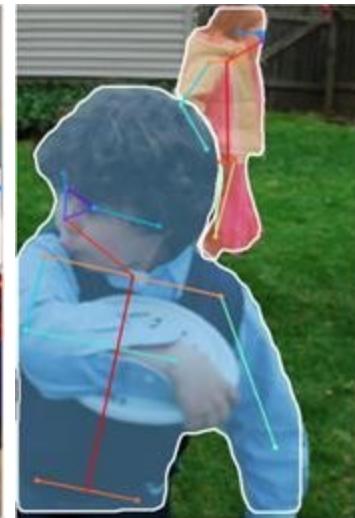


# Mask R-CNN: Very Good Results!



# Mask R-CNN

## Also does pose



# Open Source Frameworks

Lots of good implementations on GitHub!

TensorFlow Detection API:

[https://github.com/tensorflow/models/tree/master/research/object\\_detection](https://github.com/tensorflow/models/tree/master/research/object_detection)

Faster RCNN, SSD, RFCN, Mask R-CNN, ...

Detectron2 (PyTorch)

<https://github.com/facebookresearch/detectron2>

Mask R-CNN, RetinaNet, Faster R-CNN, RPN, Fast R-CNN, R-FCN, ...

Finetune on your own dataset with pre-trained models

# Recap: Lots of computer vision tasks!

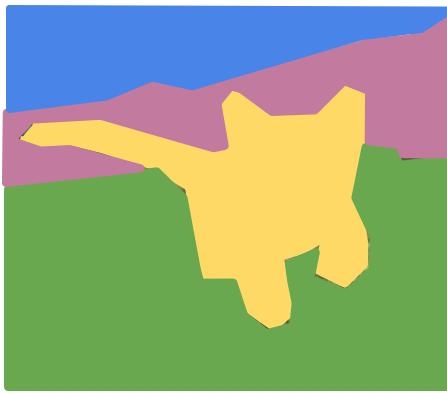
Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

[This image](#) is CC0 public domain

# Today

- Transformers Recap
- **Computer Vision Tasks**
  - Semantic Segmentation
  - Object Detection
  - Instance Segmentation
- Visualization & Understanding
  - Model Layers Visualization
  - Saliency Maps
  - CAM & Grad-CAM

# Interpreting a Linear Classifier: Visual Viewpoint

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



plane



car



bird



cat



deer



dog



frog



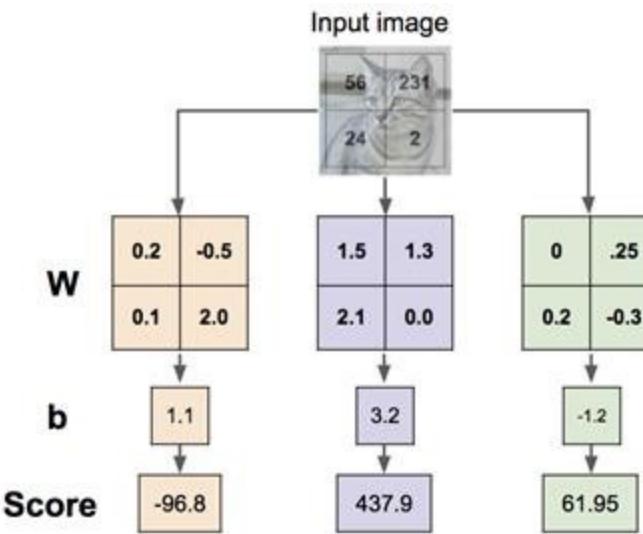
horse



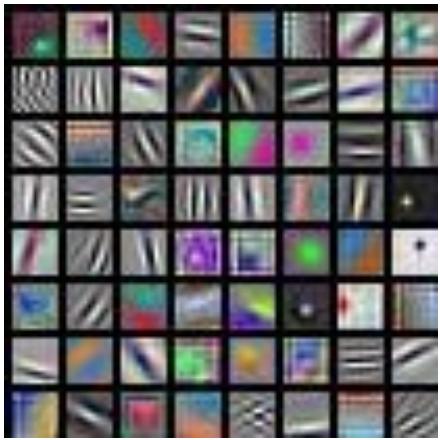
ship



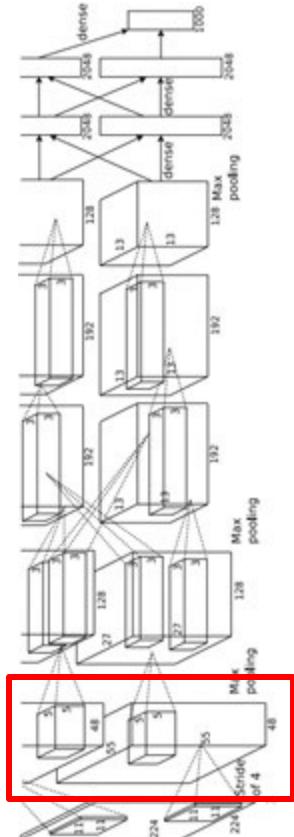
truck



# First Layer: Visualize Filters



AlexNet:  
 $64 \times 3 \times 11 \times 11$

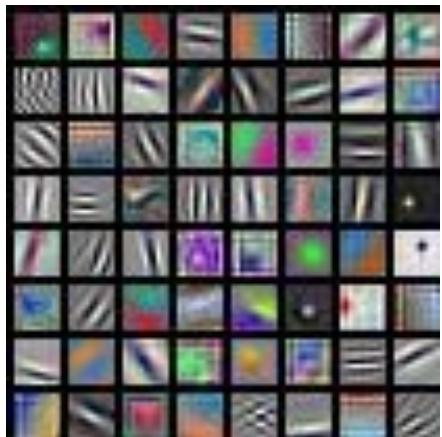


Krizhevsky, "On a weird trick for parallelizing convolutional neural networks", arXiv 2014

He et al, "Deep Residual Learning for Image Recognition", CVPR 2016

Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

# First Layer: Visualize Filters



## AlexNet: 64x3x11x11



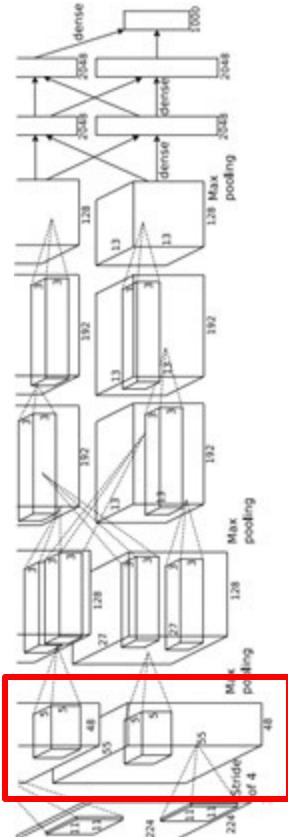
ResNet-18:  
64x3x7x7



# ResNet-101: 64 x 3 x 7 x 7

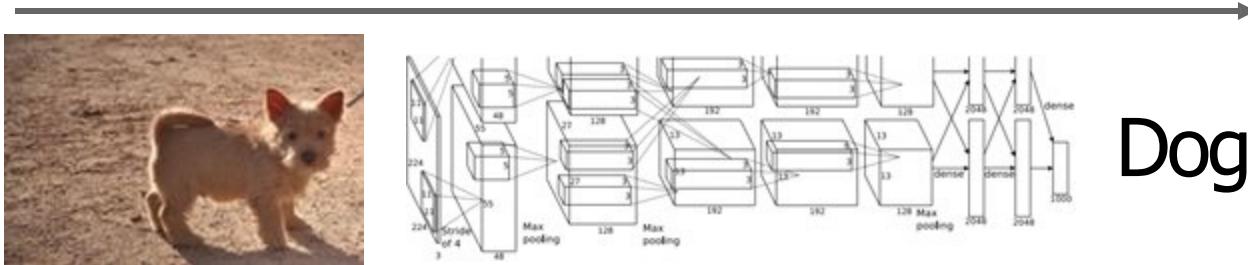


# DenseNet-121: 64 x 3 x 7 x 7



# Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities

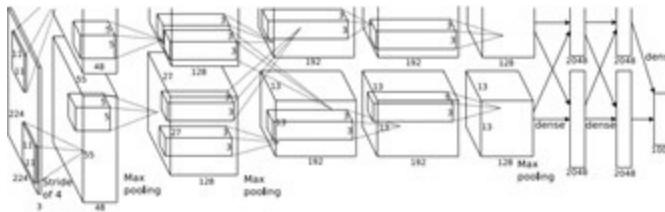


Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

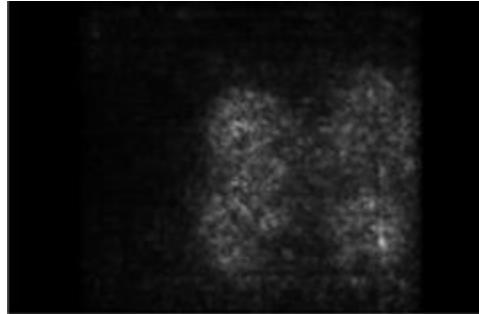
# Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities

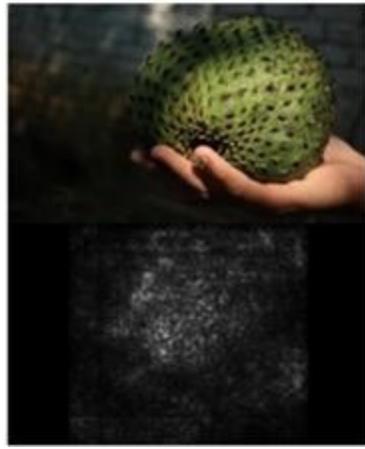


Dog

Compute gradient of (unnormalized) class score  
with respect to image pixels, take absolute value  
and max over RGB channels



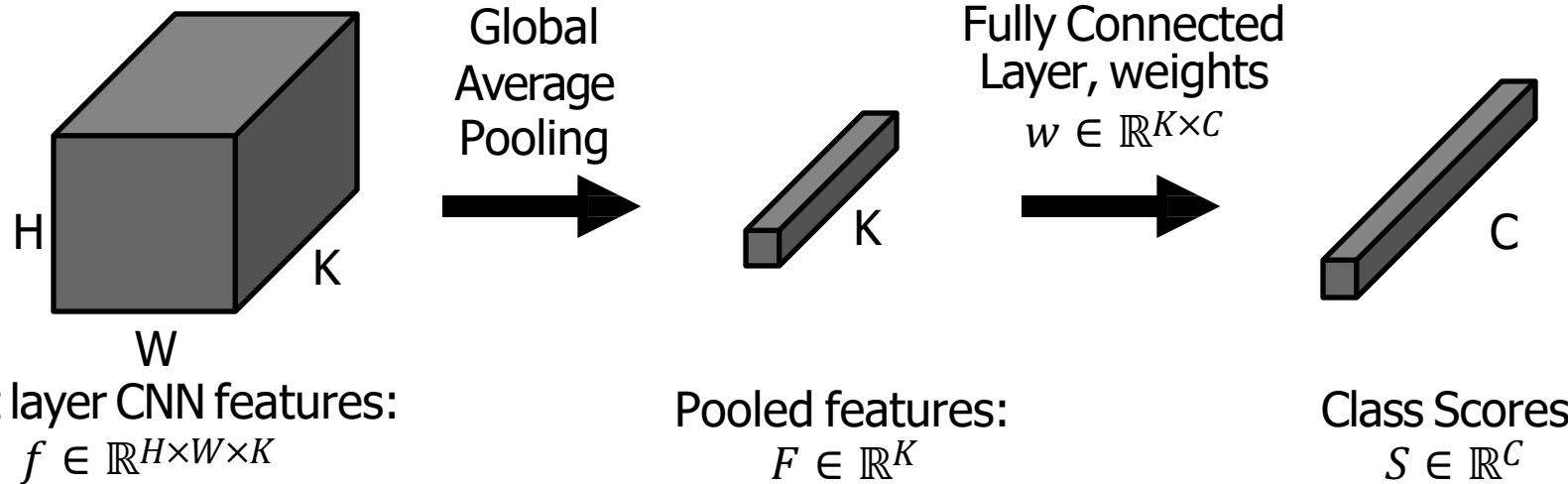
# Saliency Maps



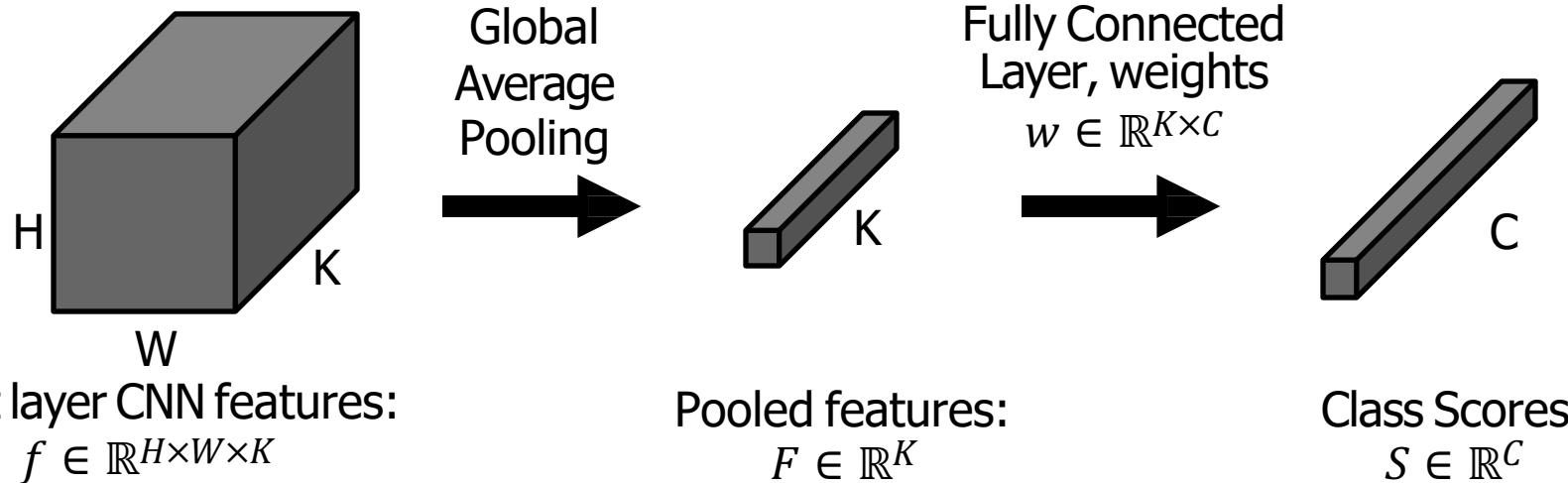
Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

# Class Activation Mapping (CAM)

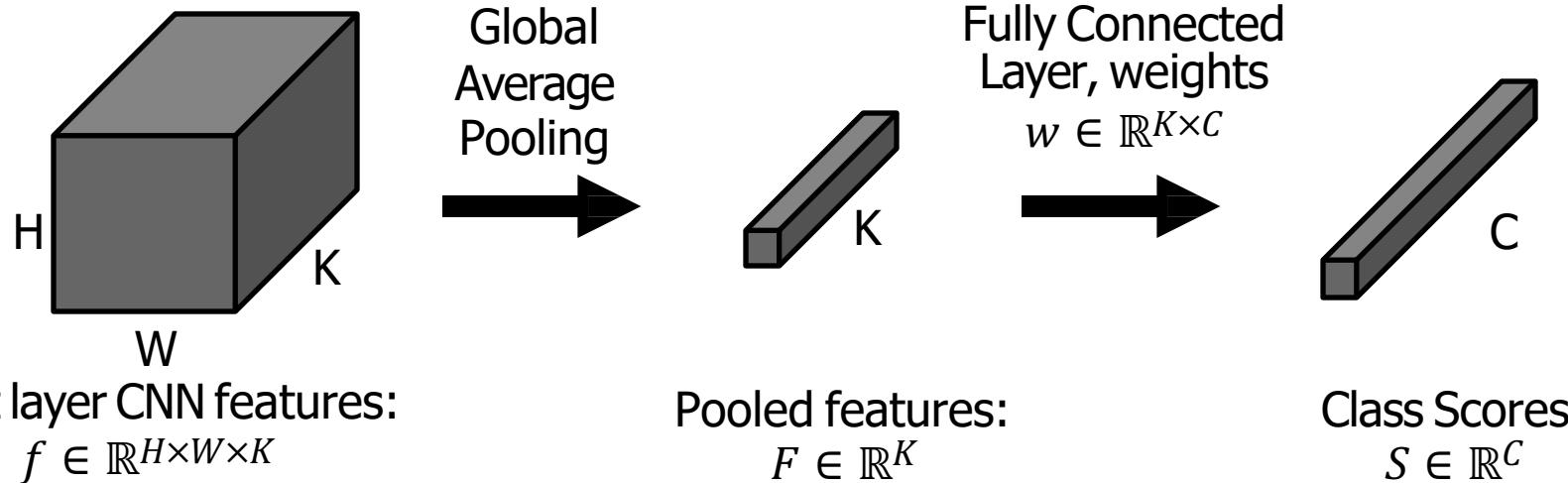


# Class Activation Mapping (CAM)



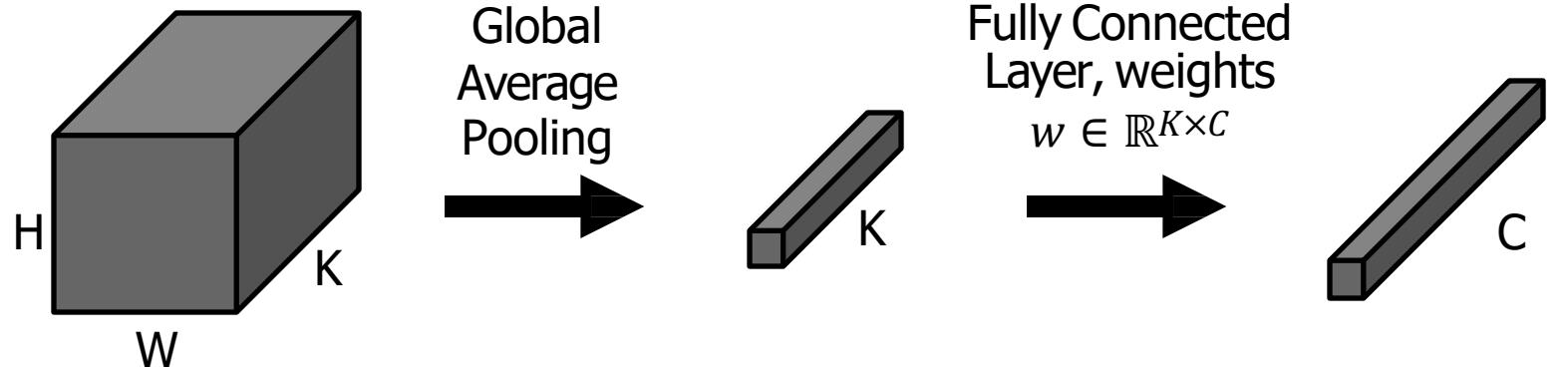
$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k}$$

# Class Activation Mapping (CAM)



$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k} \quad S_c = \sum_k w_{k,c} F_k$$

# Class Activation Mapping (CAM)



Last layer CNN features:

$$f \in \mathbb{R}^{H \times W \times K}$$

Pooled features:

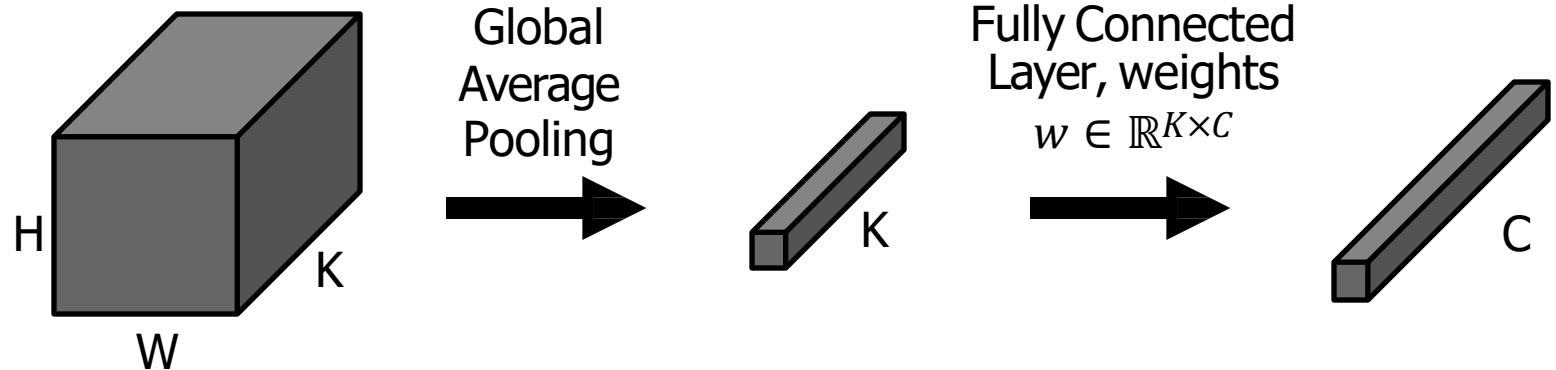
$$F \in \mathbb{R}^K$$

Class Scores:

$$S \in \mathbb{R}^C$$

$$F_k = \frac{1}{HW} \sum_{h,w} f_{h,w,k} \quad S_c = \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_k w_{k,c} \sum_{h,w} f_{h,w,k}$$

# Class Activation Mapping (CAM)



Last layer CNN features:

$$f \in \mathbb{R}^{H \times W \times K}$$

Pooled features:

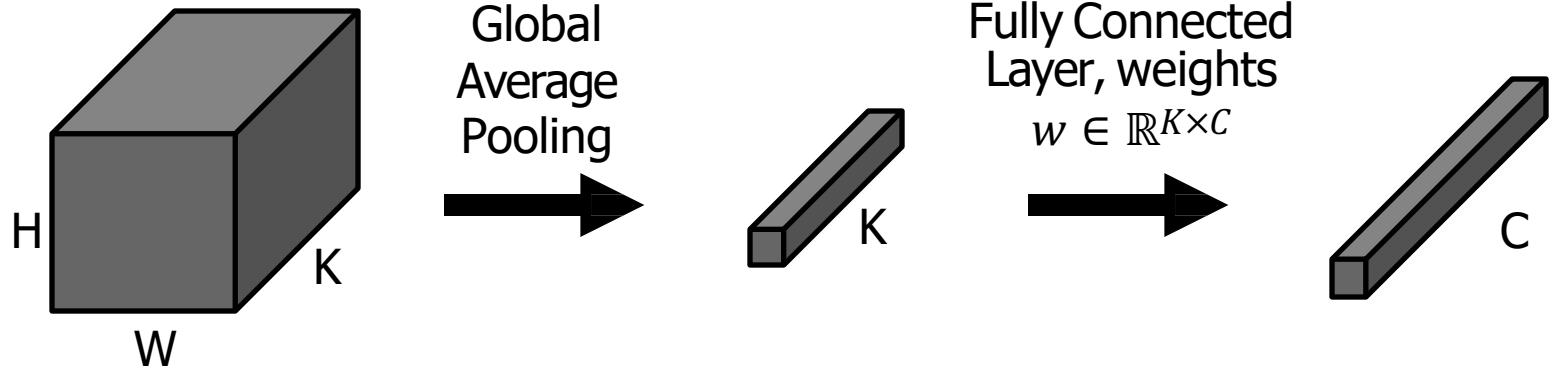
$$F \in \mathbb{R}^K$$

Class Scores:

$$S \in \mathbb{R}^C$$

$$\begin{aligned} F_k &= \frac{1}{HW} \sum_{h,w} f_{h,w,k} & S_c &= \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_k w_{k,c} \sum_{h,w} f_{h,w,k} \\ &&&= \frac{1}{HW} \sum_{h,w} \sum_k w_{k,c} f_{h,w,k} \end{aligned}$$

# Class Activation Mapping (CAM)



Last layer CNN features:

$$f \in \mathbb{R}^{H \times W \times K}$$

Pooled features:

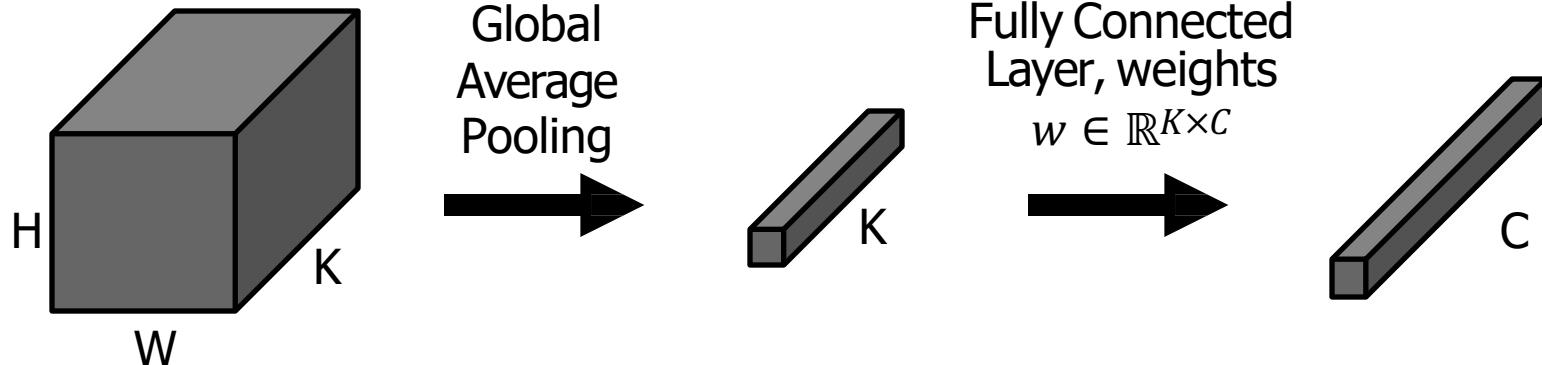
$$F \in \mathbb{R}^K$$

Class Scores:

$$S \in \mathbb{R}^C$$

$$\begin{aligned} F_k &= \frac{1}{HW} \sum_{h,w} f_{h,w,k} & S_c &= \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_k w_{k,c} \sum_{h,w} f_{h,w,k} \\ & & &= \frac{1}{HW} \sum_{h,w} \underbrace{\sum_k w_{k,c} f_{h,w,k}}_{\textcolor{blue}{w_{k,c} f_{h,w,k}}} \end{aligned}$$

# Class Activation Mapping (CAM)



Last layer CNN features:

$$f \in \mathbb{R}^{H \times W \times K}$$

Pooled features:

$$F \in \mathbb{R}^K$$

Class Scores:

$$S \in \mathbb{R}^C$$

$$\begin{aligned} F_k &= \frac{1}{HW} \sum_{h,w} f_{h,w,k} & S_c &= \sum_k w_{k,c} F_k = \frac{1}{HW} \sum_k w_{k,c} \sum_{h,w} f_{h,w,k} \\ & & &= \frac{1}{HW} \sum_{h,w} \sum_k w_{k,c} f_{h,w,k} \end{aligned}$$

Class Activation Maps:

$$M \in \mathbb{R}^{C,H,W}$$

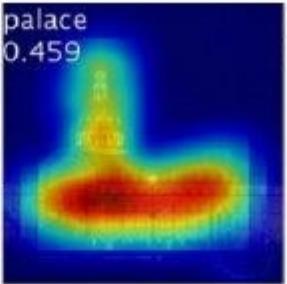
$$M_{c,h,w} = \sum_k w_{k,c} f_{h,w,k}$$

# Class Activation Mapping (CAM)

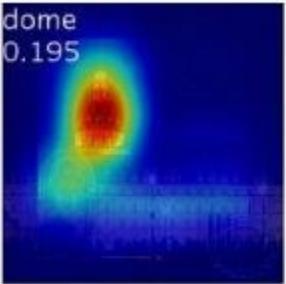
dome



palace  
0.459



dome  
0.195

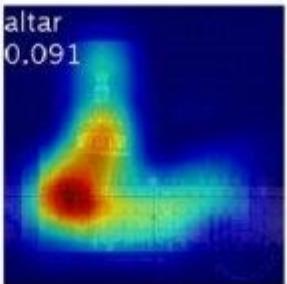


church

0.146



altar  
0.091



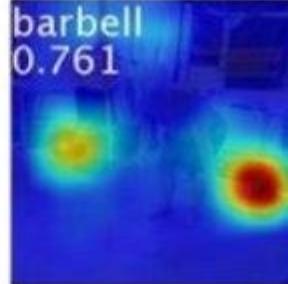
monastery  
0.051



Class activation maps of top 5 predictions



barbell  
0.761



barbell  
0.447



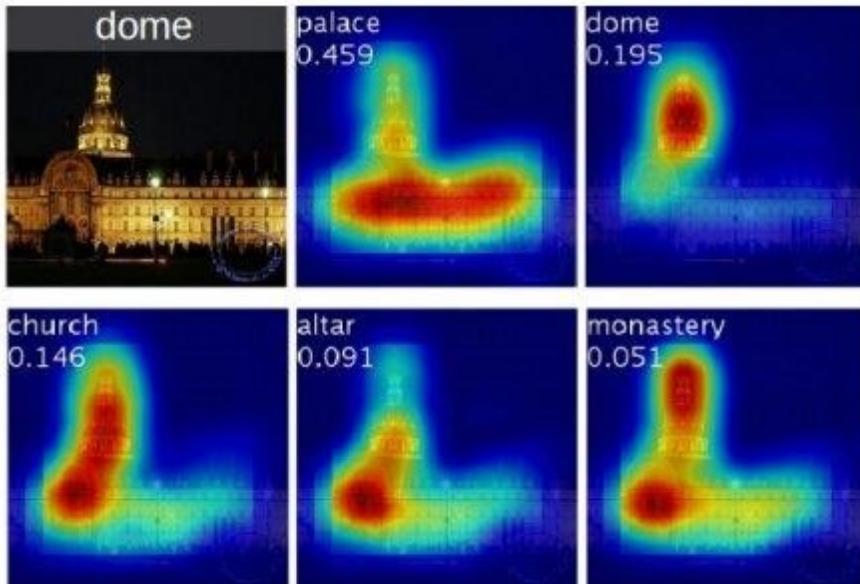
barbell  
0.999



Class activation maps for one object class

# Class Activation Mapping (CAM)

Problem: Can only apply to last conv layer



Class activation maps of top 5 predictions



Class activation maps for one object class

# Gradient-Weighted Class Activation Mapping (Grad-CAM)

1. Pick any layer, with activations  $A \in \mathbb{R}^{H \times W \times K}$

# Gradient-Weighted Class Activation Mapping (Grad-CAM)

1. Pick any layer, with activations  $A \in \mathbb{R}^{H \times W \times K}$
2. Compute gradient of class score  $S_c$  with respect to A:

$$\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$$

# Gradient-Weighted Class Activation Mapping (Grad-CAM)

1. Pick any layer, with activations  $A \in \mathbb{R}^{H \times W \times K}$
2. Compute gradient of class score  $S_c$  with respect to A:

$$\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$$

3. Global Average Pool the gradients to get weights  $\alpha \in \mathbb{R}^K$ :

$$\alpha_k = \frac{1}{HW} \underset{h,w}{\square} \frac{\partial S_c}{\partial A_{h,w,k}}$$

# Gradient-Weighted Class Activation Mapping (Grad-CAM)

1. Pick any layer, with activations  $A \in \mathbb{R}^{H \times W \times K}$
2. Compute gradient of class score  $S_c$  with respect to A:

$$\frac{\partial S_c}{\partial A} \in \mathbb{R}^{H \times W \times K}$$

3. Global Average Pool the gradients to get weights  $\alpha \in \mathbb{R}^K$ :

$$\alpha_k = \frac{1}{HW} \underset{h,w}{\mathbb{E}} \frac{\partial S_c}{\partial A_{h,w,k}}$$

4. Compute activation map  $M^c \in \mathbb{R}^{H,W}$ :

$$M_{h,w}^c = \text{ReLU} \left( \underset{k}{\mathbb{E}} \alpha_k A_{h,w,k} \right)$$

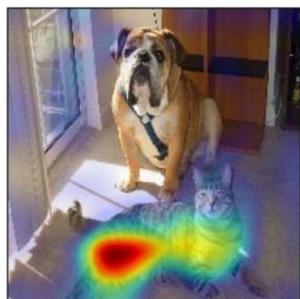
# Gradient-Weighted Class Activation Mapping (Grad-CAM)



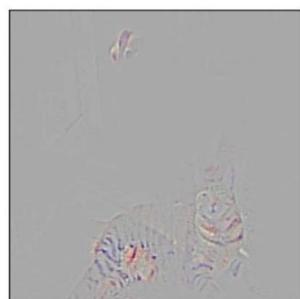
(a) Original Image



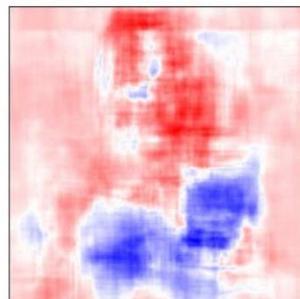
(b) Guided Backprop ‘Cat’



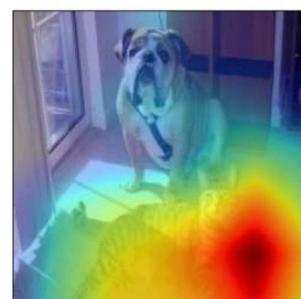
(c) Grad-CAM ‘Cat’



(d) Guided Grad-CAM ‘Cat’



(e) Occlusion map for ‘Cat’



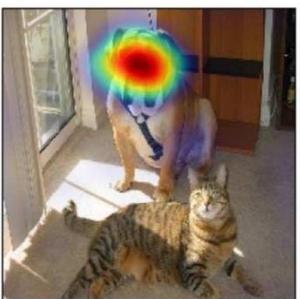
(f) ResNet Grad-CAM ‘Cat’



(g) Original Image



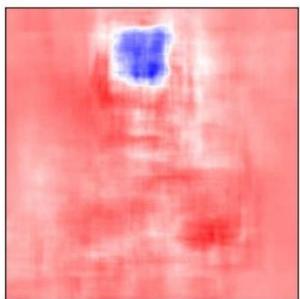
(h) Guided Backprop ‘Dog’



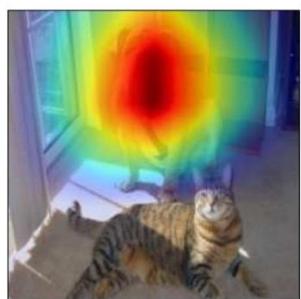
(i) Grad-CAM ‘Dog’



(j) Guided Grad-CAM ‘Dog’

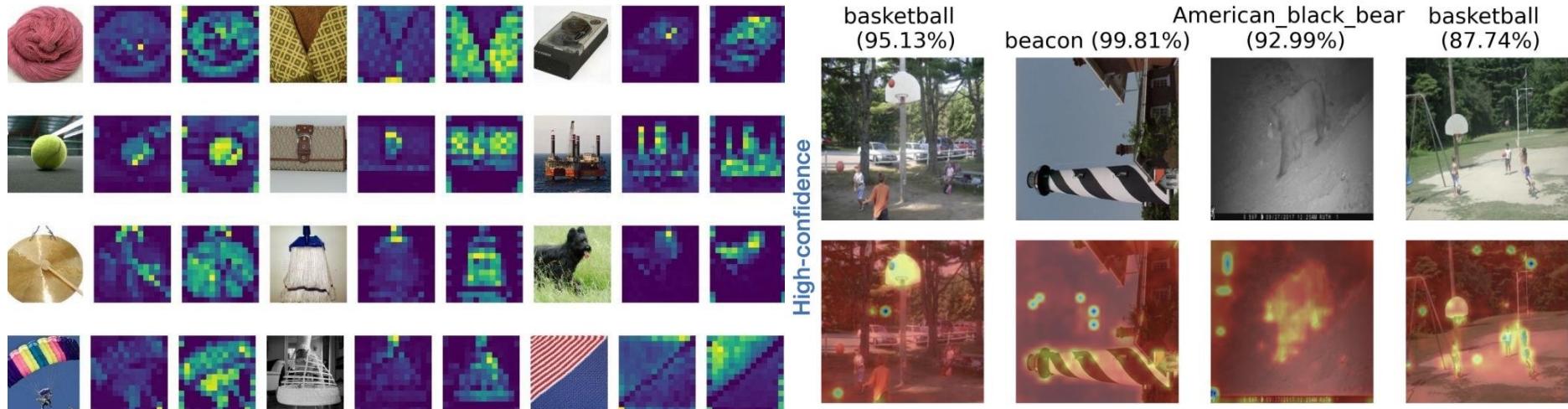


(k) Occlusion map for ‘Dog’



(l) ResNet Grad-CAM ‘Dog’

# Visualizing ViT features



Chen et al., When Vision Transformers OutperformResnets Without Pre-training Or Strong Data Augmentations, ICLR 2022; Paul and Chen, Vision Transformers are Robust Learners, AAAI 2022. Reproduced for educational purposes.

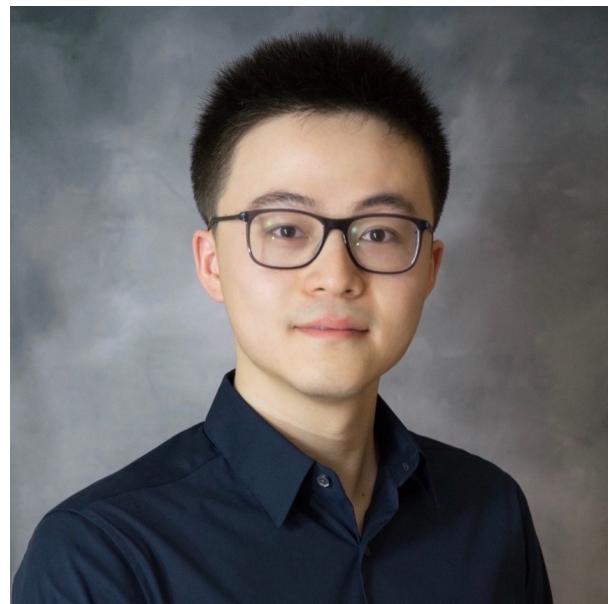
# Today

- Transformers Recap
- **Computer Vision Tasks**
  - Semantic Segmentation
  - Object Detection
  - Instance Segmentation
- Visualization & Understanding
  - Model Layers Visualization
  - Saliency Maps
  - CAM & Grad-CAM

# Lecture 10: Video Understanding



# Instructor today



Ruohan Gao

<https://ruohangao.github.io/>

I taught CS231N at Stanford from 2021-2023

Teaching multimodal  
compute vision now.

Ph.D. at UT Austin

Postdoc at Stanford

Some time at Meta

University of Maryland,  
College Park



# Recall: (2D) Image classification



This image by [Nikita](#) is  
licensed under [CC-BY 2.0](#)

(assume given a set of possible labels)  
{dog, cat, truck, plane, ...}

cat

# Last Lecture: (2D) Detection and Segmentation

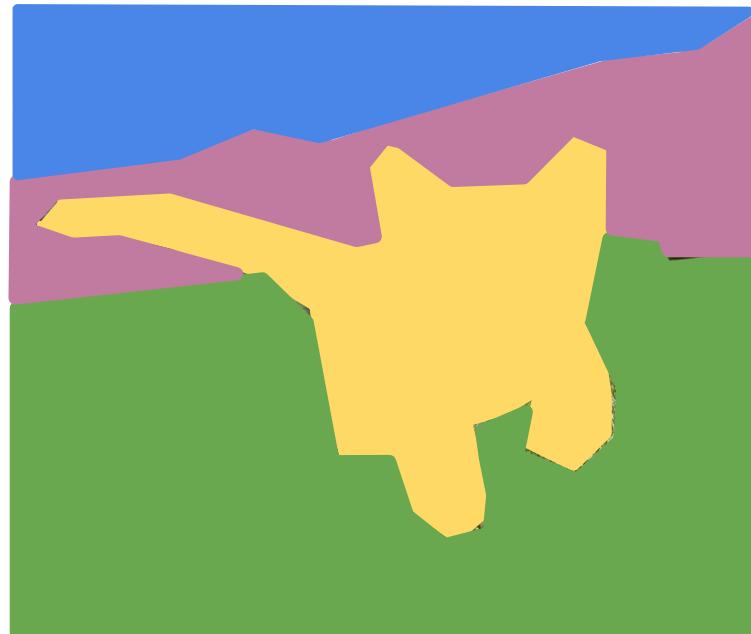
Classification



CAT

No spatial extent

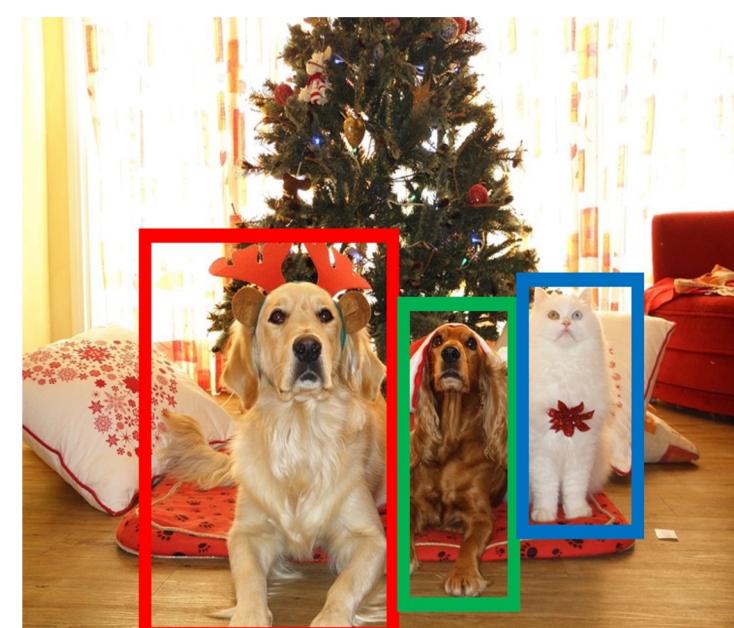
Semantic  
Segmentation



GRASS, CAT, TREE,  
SKY

No objects, just pixels

Object  
Detection



DOG, DOG, CAT

Multiple Objects

Instance  
Segmentation



DOG, DOG, CAT

[This image is CC0 public domain](#)

Living room

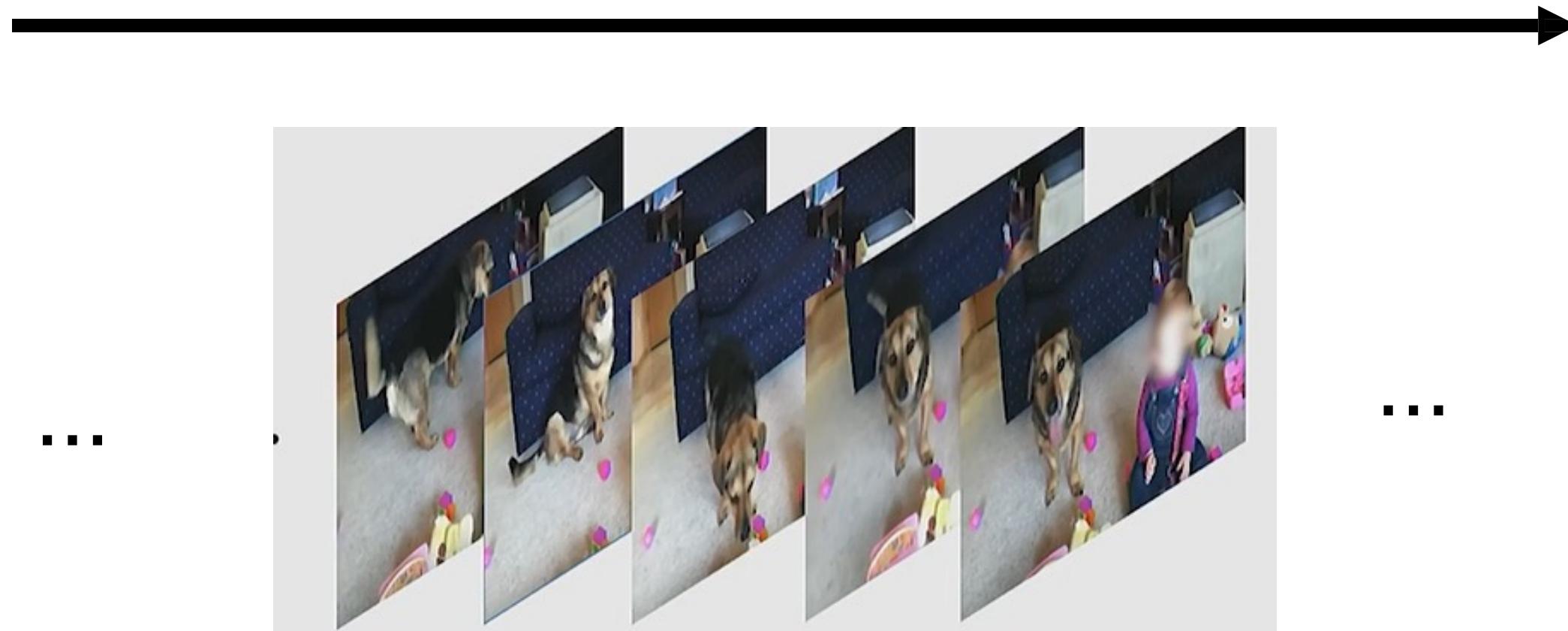
Dog

Baby

# Today: Video = 2D + Time

A video is a sequence of images

4D tensor:  $T \times 3 \times H \times W$   
(or  $3 \times T \times H \times W$ )



This image is CC0 public domain

# Example task: Video Classification



Input video:  
 $T \times 3 \times H \times W$

Swimming  
Running  
Jumping  
Eating  
Standing

# Example task: Video Classification



Images: Recognize objects



Dog  
Cat  
Fish  
Truck



Videos: Recognize actions



Swimming  
Running  
Jumping  
Eating  
Standing

# Problem: Videos are big!

Videos are ~30 frames per second (fps)



Size of uncompressed video  
(3 bytes per pixel):

SD (640 x 480): ~1.5 GB per minute

HD (1920 x 1080): ~10 GB per minute

Input video:

$T \times 3 \times H \times W$

# Problem: Videos are big!

Videos are ~30 frames per second (fps)



Input video:  
 $T \times 3 \times H \times W$

Size of uncompressed video  
(3 bytes per pixel):

SD (640 x 480): ~1.5 GB per minute

HD (1920 x 1080): ~10 GB per minute

Solution: Train on short clips: low  
fps and low spatial resolution  
e.g.  $T = 16$ ,  $H=W=112$   
(3.2 seconds at 5 fps, 588 KB)

# Training on Clips

Raw video: Long, high FPS

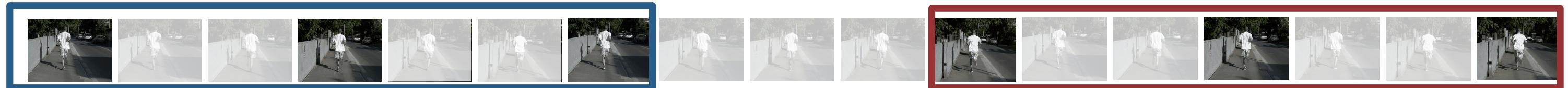


# Training on Clips

Raw video: Long, high FPS



Training: Train model to classify short clips with low FPS

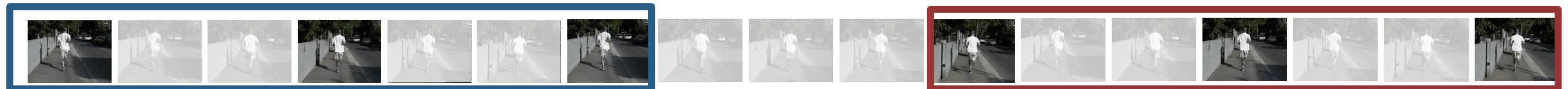


# Training on Clips

Raw video: Long, high FPS



Training: Train model to classify short clips with low FPS



Testing: Run model on different clips, average predictions

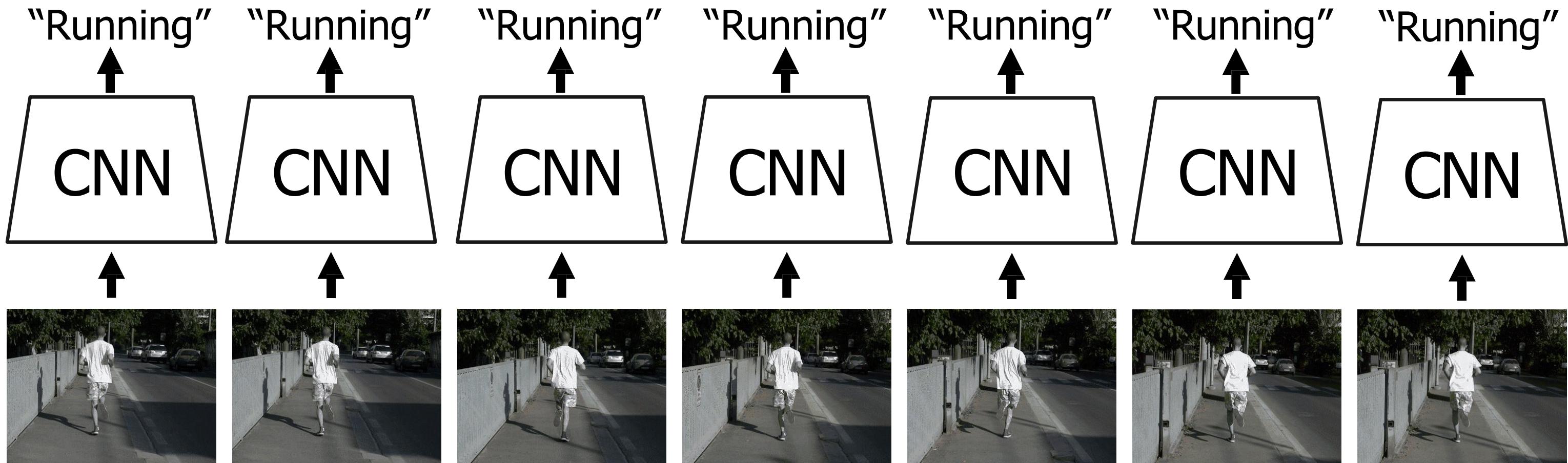


# Video Classification: Single-Frame CNN

**Simple idea:** train normal 2D CNN to classify video frames independently!

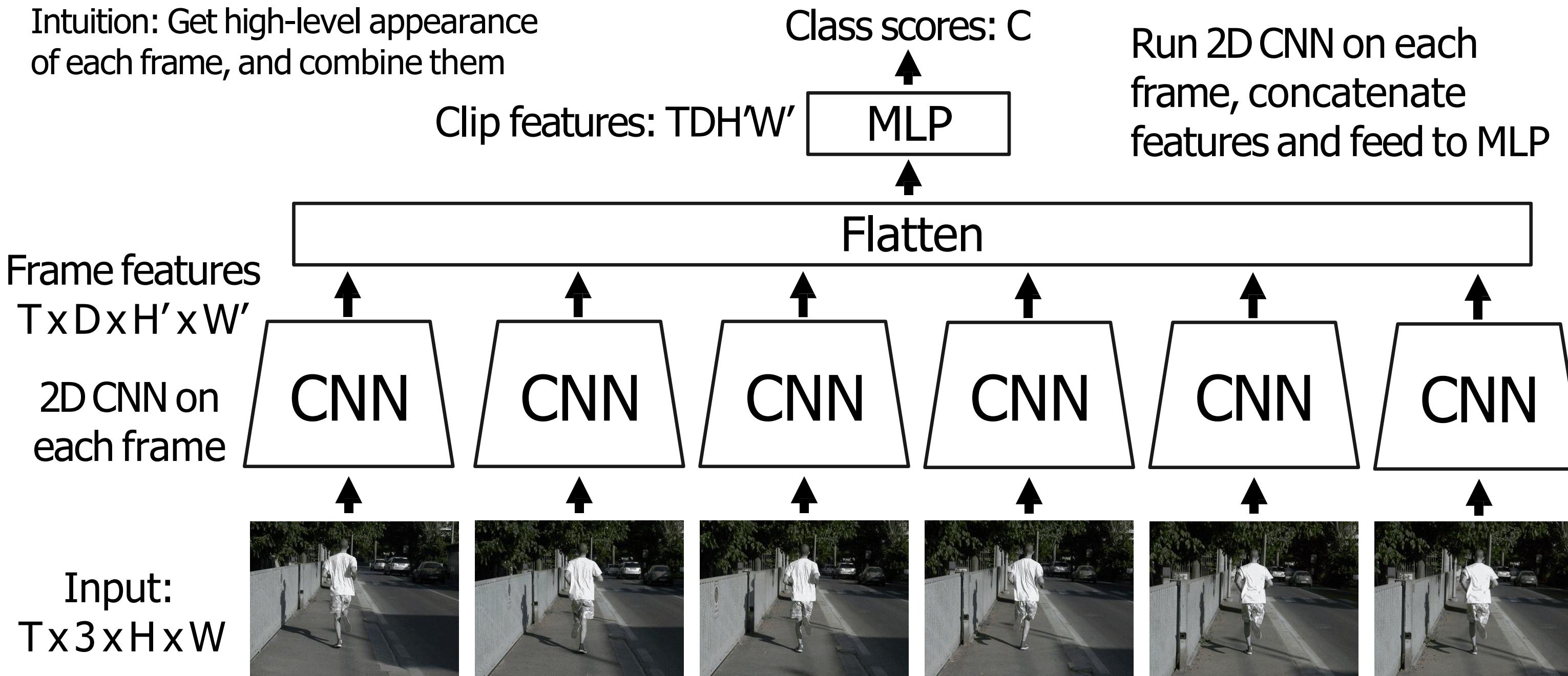
(Average predicted probs at test-time)

Often a very strong baseline for video classification



# Video Classification: Late Fusion (with FC layers)

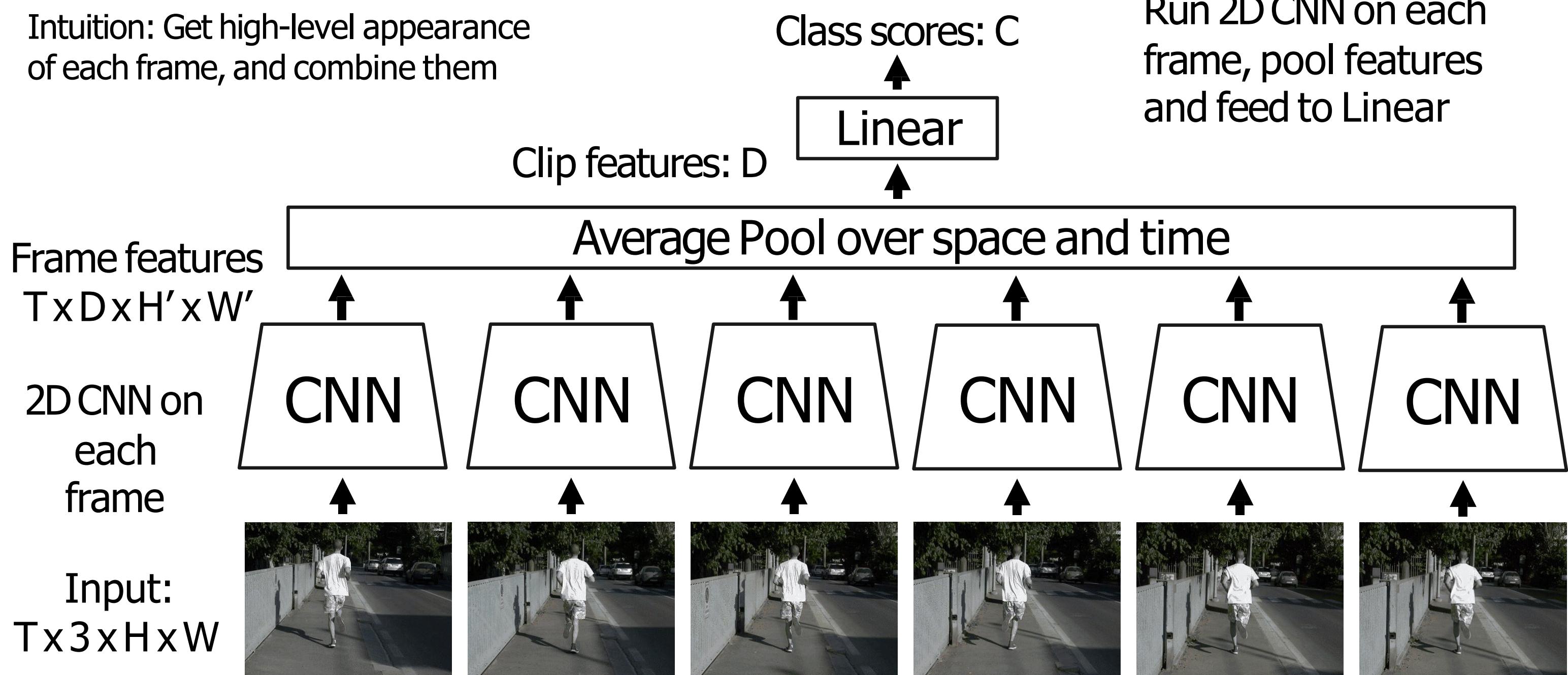
Intuition: Get high-level appearance of each frame, and combine them



Karpathy et al, "Large-scale Video Classification with Convolutional Neural Networks", CVPR 2014

# Video Classification: Late Fusion (with pooling)

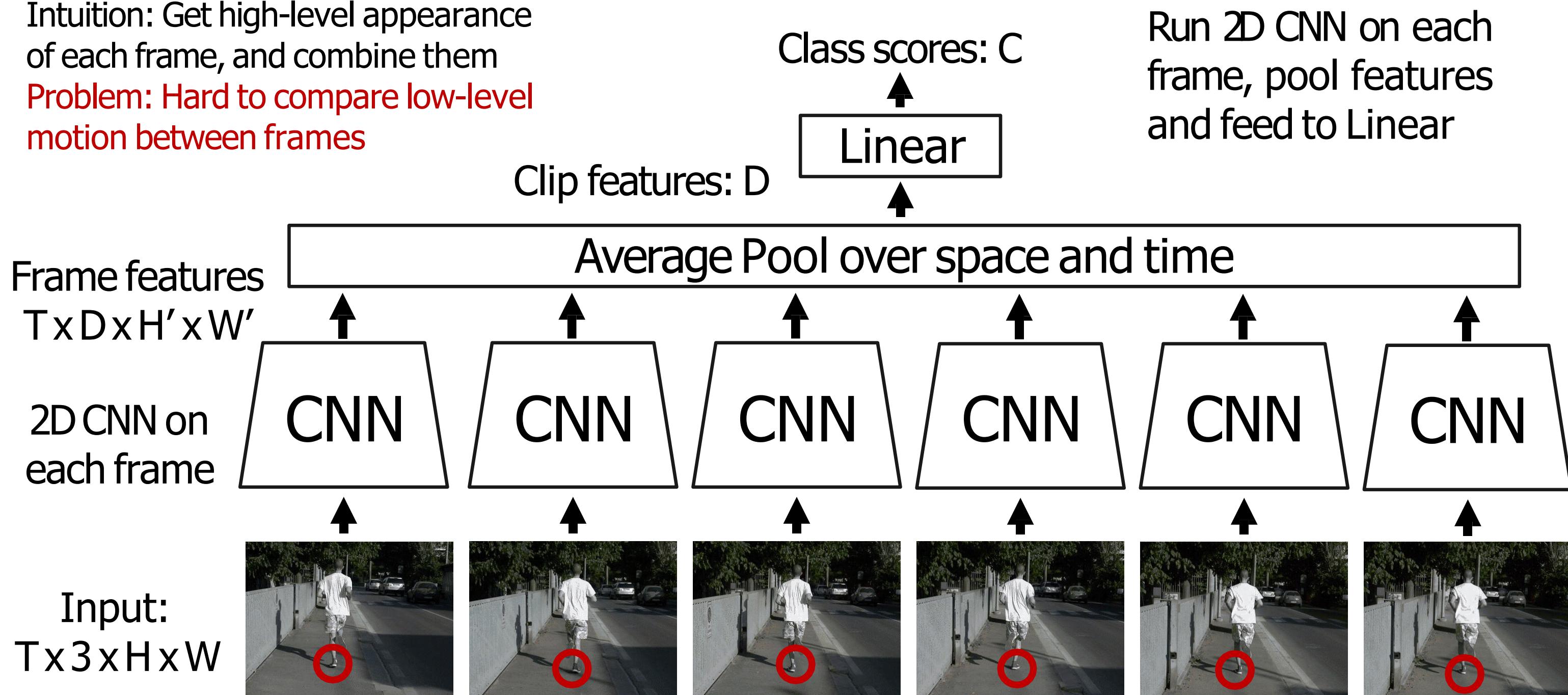
Intuition: Get high-level appearance of each frame, and combine them



# Video Classification: Late Fusion (with pooling)

Intuition: Get high-level appearance of each frame, and combine them

Problem: Hard to compare low-level motion between frames



# Video Classification: Early Fusion

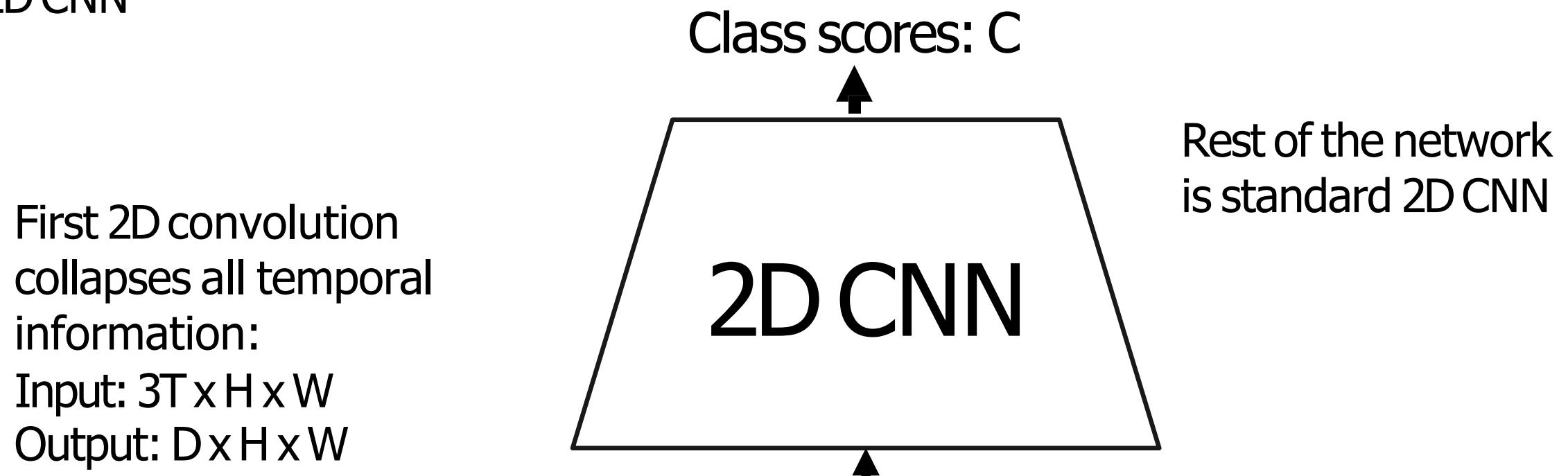
Intuition: Compare frames  
with very first conv layer, after  
that normal 2D CNN

Reshape:  
 $3T \times H \times W$

Input:  
 $T \times 3 \times H \times W$



Karpathy et al, "Large-scale Video Classification with Convolutional Neural Networks", CVPR 2014



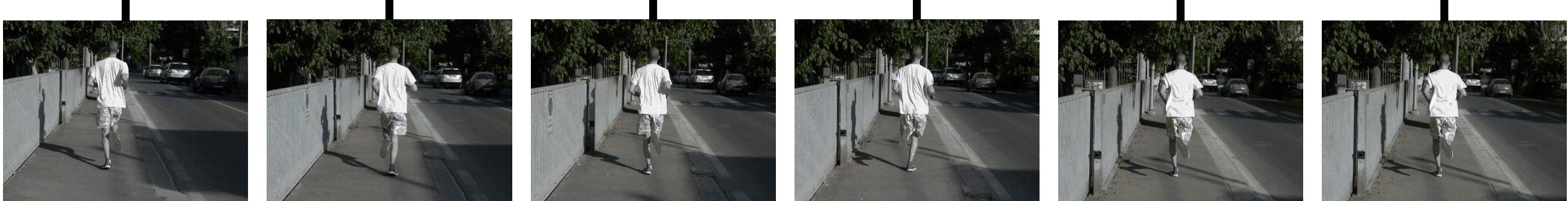
# Video Classification: Early Fusion

Intuition: Compare frames with very first conv layer, after that normal 2D CNN

Problem: One layer of temporal processing may not be enough!

Reshape:  
 $3T \times H \times W$

Input:  
 $T \times 3 \times H \times W$



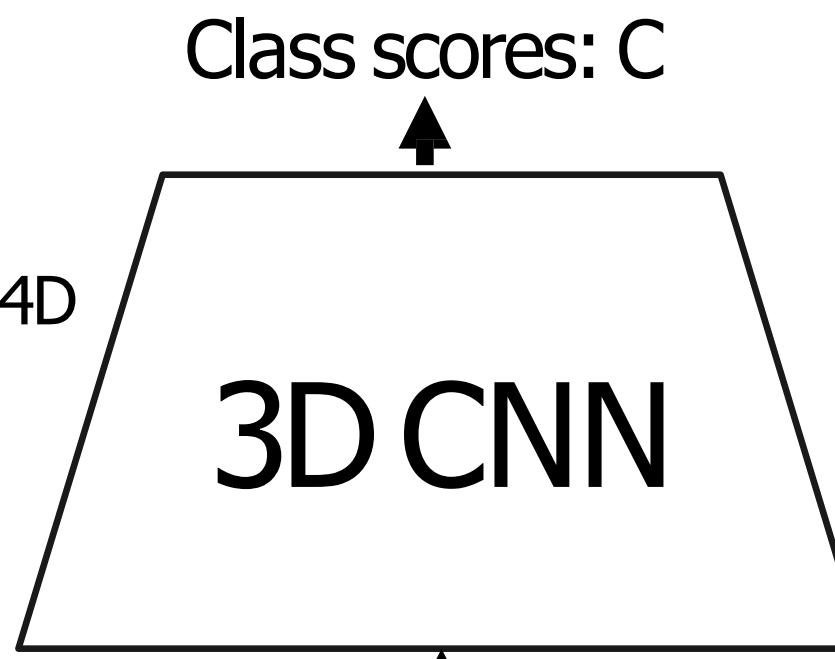
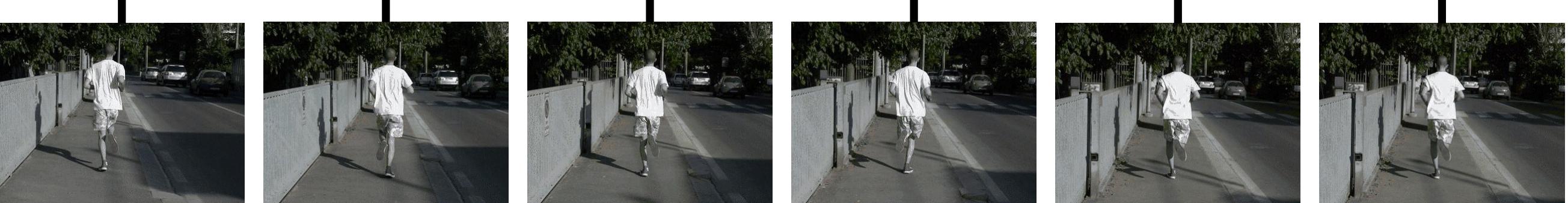
Karpathy et al, "Large-scale Video Classification with Convolutional Neural Networks", CVPR 2014

# Video Classification: 3D CNN

Intuition: Use 3D versions of convolution and pooling to slowly fuse temporal information over the course of the network

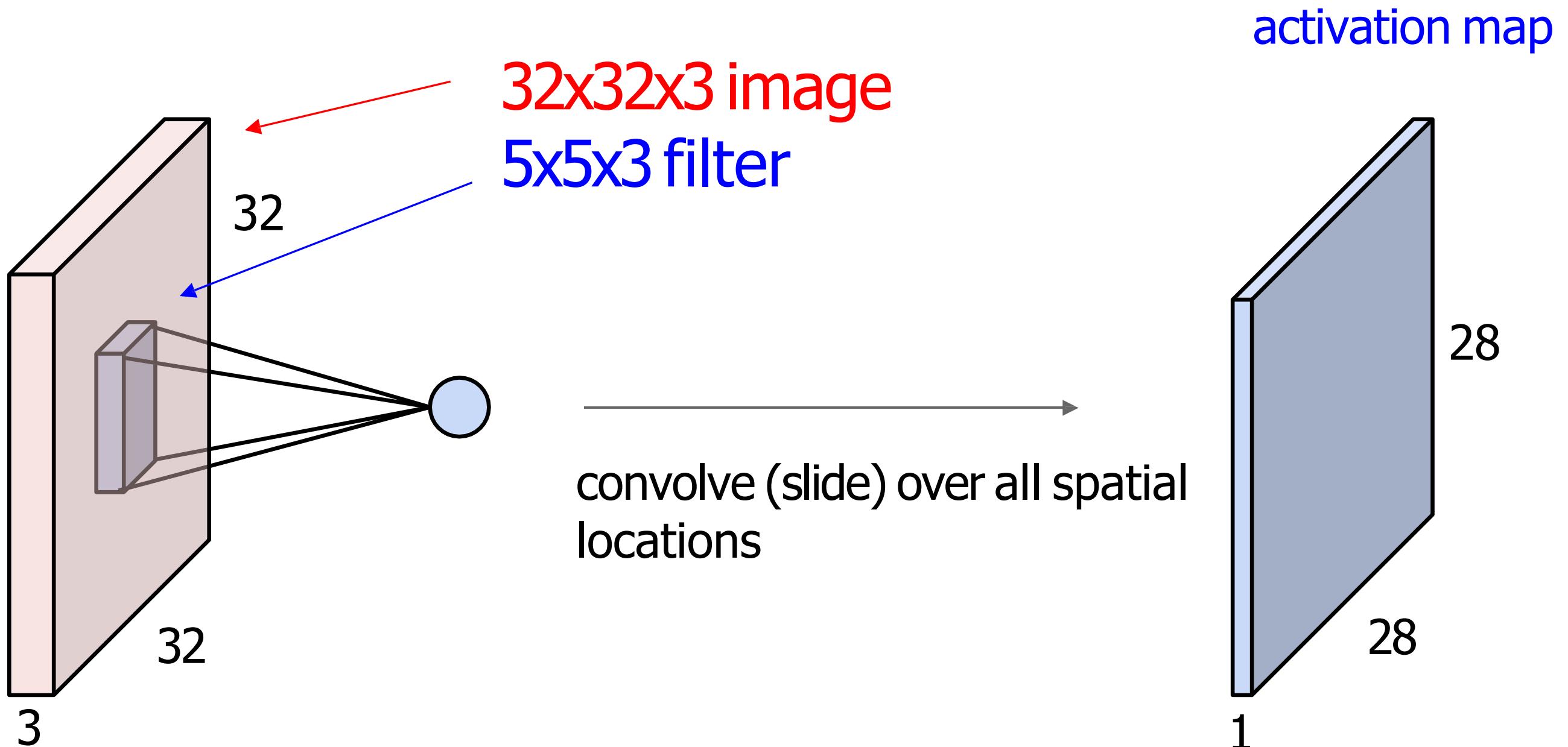
Each layer in the network is a 4D tensor:  $D \times T \times H \times W$   
Use 3D conv and 3D pooling operations

Input:  
 $3 \times T \times H \times W$

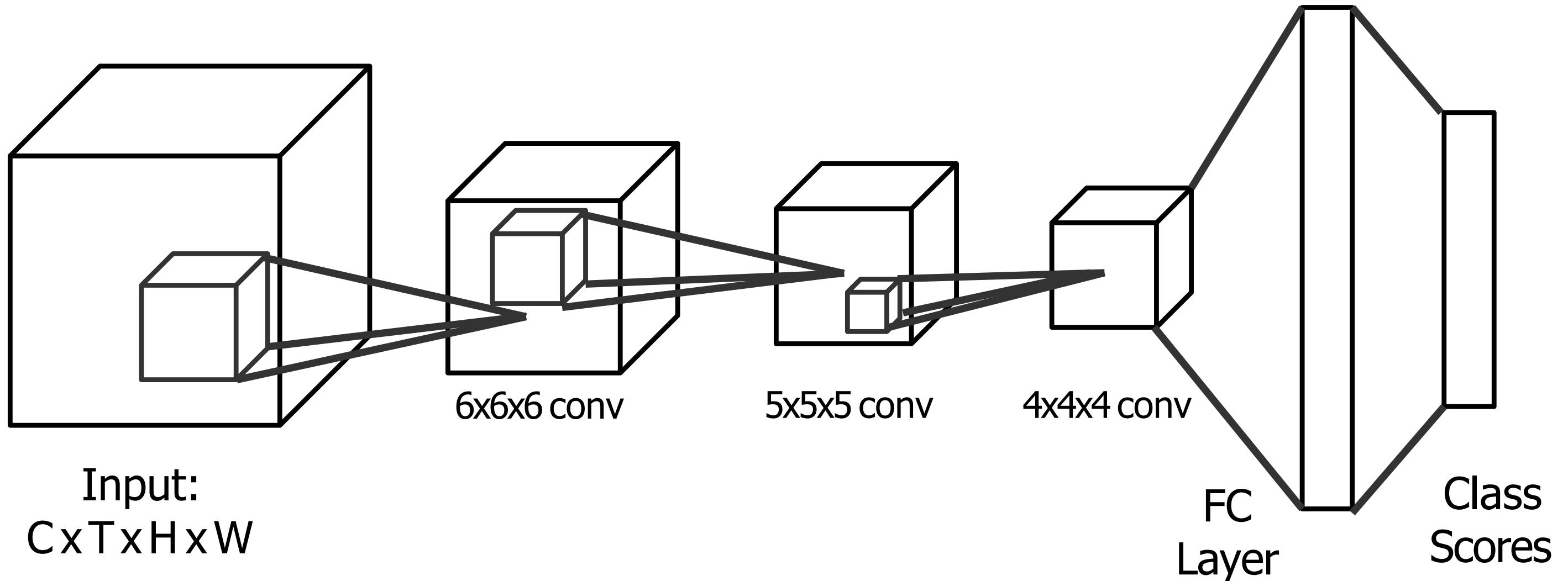


Ji et al, "3D Convolutional Neural Networks for Human Action Recognition", TPAMI 2010 ; Karpathy et al, "Large-scale Video Classification with Convolutional Neural Networks", CVPR 2014

# Convolution Layer



# 3D Convolution



# Early Fusion vs Late Fusion vs 3D CNN

Late  
Fusion

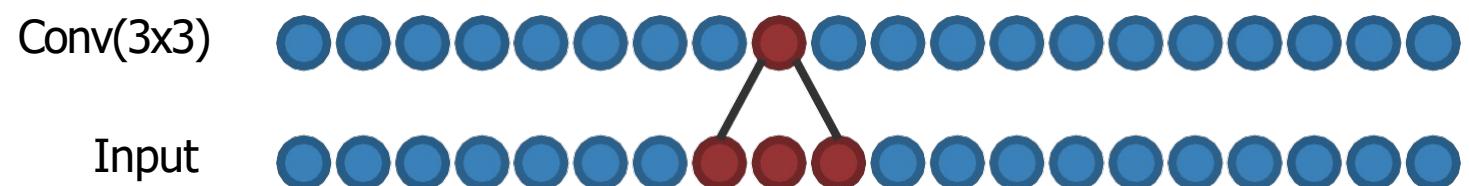
Layer	Size $(C \times T \times H \times W)$	Receptive Field $(T \times H \times W)$
Input	$3 \times 20 \times 64 \times 64$	
Conv2D( $3 \times 3$ , $3 \rightarrow 12$ )	$12 \times 20 \times 64 \times 64$	$1 \times 3 \times 3$

(Small example  
architectures, in  
practice much bigger)

# Early Fusion vs Late Fusion vs 3D CNN

Late  
Fusion

Layer	Size $(C \times T \times H \times W)$	Receptive Field $(T \times H \times W)$
Input	$3 \times 20 \times 64 \times 64$	
Conv2D( $3 \times 3$ , $3 \rightarrow 12$ )	$12 \times 20 \times 64 \times 64$	$1 \times 3 \times 3$

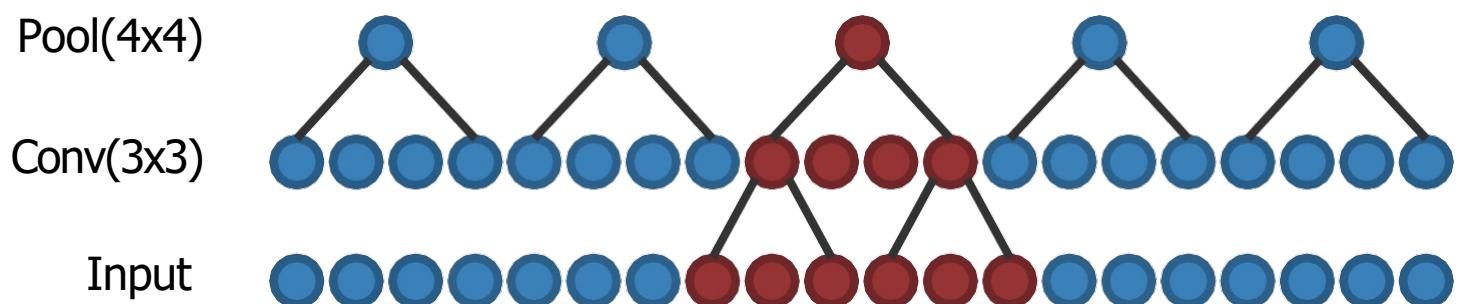


(Small example  
architectures, in  
practice much bigger)

# Early Fusion vs Late Fusion vs 3D CNN

Late  
Fusion

Layer	Size $(C \times T \times H \times W)$	Receptive Field $(T \times H \times W)$
Input	$3 \times 20 \times 64 \times 64$	
Conv2D( $3 \times 3$ , $3 > 12$ )	$12 \times 20 \times 64 \times 64$	$1 \times 3 \times 3$
Pool2D( $4 \times 4$ )	$12 \times 20 \times 16 \times 16$	$1 \times 6 \times 6$



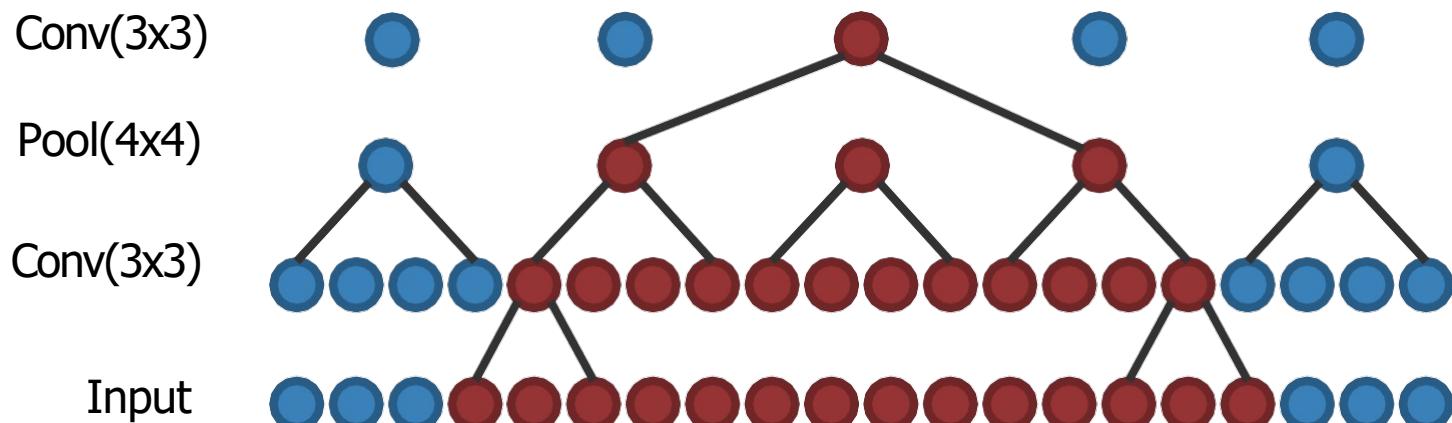
(Small example  
architectures, in  
practice much bigger)

# Early Fusion vs Late Fusion vs 3D CNN

Late  
Fusion

Layer	Size $(C \times T \times H \times W)$	Receptive Field $(T \times H \times W)$
Input	$3 \times 20 \times 64 \times 64$	
Conv2D(3x3, 3->12)	$12 \times 20 \times 64 \times 64$	$1 \times 3 \times 3$
Pool2D(4x4)	$12 \times 20 \times 16 \times 16$	$1 \times 6 \times 6$
Conv2D(3x3, 12->24)	$24 \times 20 \times 16 \times 16$	$1 \times 14 \times 14$

Build slowly in space



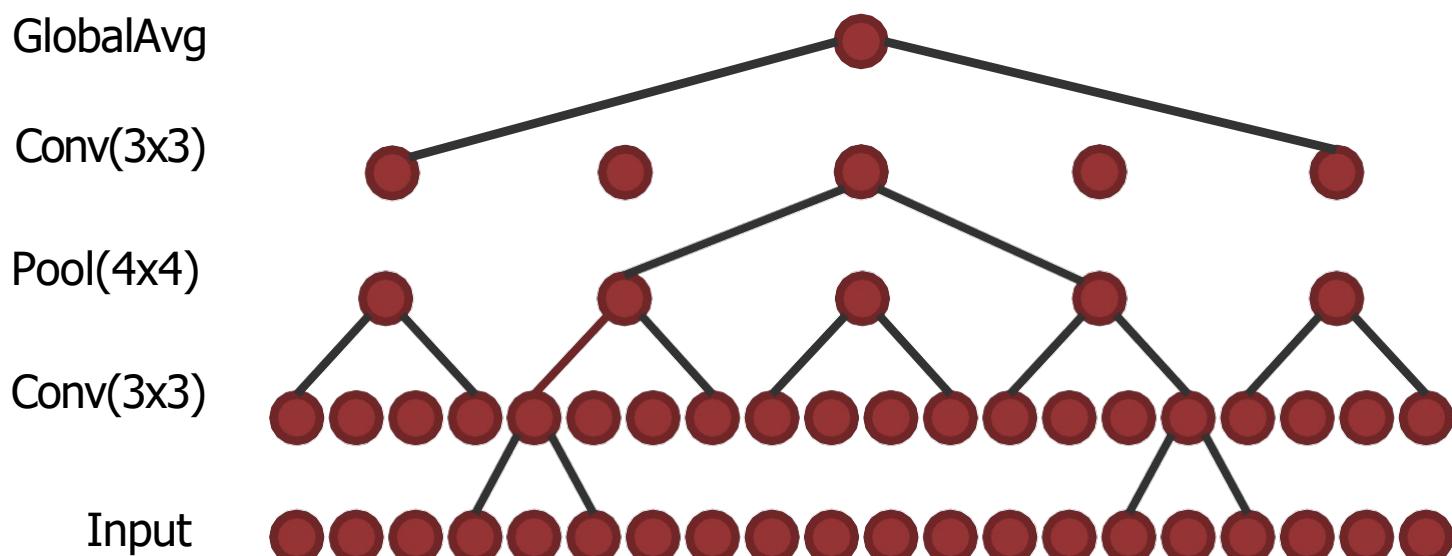
(Small example  
architectures, in  
practice much bigger)

# Early Fusion vs Late Fusion vs 3D CNN

Late  
Fusion

Layer	Size $(C \times T \times H \times W)$	Receptive Field $(T \times H \times W)$
Input	$3 \times 20 \times 64 \times 64$	
Conv2D(3x3, 3->12)	$12 \times 20 \times 64 \times 64$	$1 \times 3 \times 3$
Pool2D(4x4)	$12 \times 20 \times 16 \times 16$	$1 \times 6 \times 6$
Conv2D(3x3, 12->24)	$24 \times 20 \times 16 \times 16$	$1 \times 14 \times 14$
GlobalAvgPool	$24 \times 1 \times 1 \times 1$	$20 \times 64 \times 64$

Build slowly in space,  
All-at-once in time at end



(Small example  
architectures, in  
practice much bigger)

# Early Fusion vs Late Fusion vs 3D CNN

Late  
Fusion

Early  
Fusion

Layer	Size $(C \times T \times H \times W)$	Receptive Field $(T \times H \times W)$
Input	$3 \times 20 \times 64 \times 64$	
Conv2D(3x3, 3->12)	$12 \times 20 \times 64 \times 64$	$1 \times 3 \times 3$
Pool2D(4x4)	$12 \times 20 \times 16 \times 16$	$1 \times 6 \times 6$
Conv2D(3x3, 12->24)	$24 \times 20 \times 16 \times 16$	$1 \times 14 \times 14$
GlobalAvgPool	$24 \times 1 \times 1 \times 1$	$20 \times 64 \times 64$
Input	$3 \times 20 \times 64 \times 64$	
Conv2D(3x3, 3*20->12)	$12 \times 64 \times 64$	$20 \times 3 \times 3$
Pool2D(4x4)	$12 \times 16 \times 16$	$20 \times 6 \times 6$
Conv2D(3x3, 12->24)	$24 \times 16 \times 16$	$20 \times 14 \times 14$
GlobalAvgPool	$24 \times 1 \times 1$	$20 \times 64 \times 64$

Build slowly in space,  
All-at-once in time at end

Build slowly in space,  
All-at-once in time at start

(Small example  
architectures, in  
practice much bigger)

# Early Fusion vs Late Fusion vs 3D CNN

Late  
Fusion

Early  
Fusion

3D CNN

	Size <b>(C x T x H x W)</b>	Receptive Field <b>(T x H x W)</b>
Layer		
	Input	$3 \times 20 \times 64 \times 64$
	Conv2D(3x3, 3->12)	$12 \times 20 \times 64 \times 64$
	Pool2D(4x4)	$12 \times 20 \times 16 \times 16$
	Conv2D(3x3, 12->24)	$24 \times 20 \times 16 \times 16$
Layer	GlobalAvgPool	$24 \times 1 \times 1 \times 1$
		$20 \times 64 \times 64$
	Input	$3 \times 20 \times 64 \times 64$
	Conv2D(3x3, 3*20->12)	$12 \times 64 \times 64$
	Pool2D(4x4)	$12 \times 16 \times 16$
Layer	Conv2D(3x3, 12->24)	$24 \times 16 \times 16$
	GlobalAvgPool	$24 \times 1 \times 1$
		$20 \times 64 \times 64$
	Input	$3 \times 20 \times 64 \times 64$
	Conv3D(3x3x3, 3->12)	$12 \times 20 \times 64 \times 64$
Layer	Pool3D(4x4x4)	$12 \times 5 \times 16 \times 16$
	Conv3D(3x3x3, 12->24)	$24 \times 5 \times 16 \times 16$
	GlobalAvgPool	$24 \times 1 \times 1$
		$20 \times 64 \times 64$

Build slowly in space,  
All-at-once in time at end

Build slowly in space,  
All-at-once in time at start

Build slowly in space,  
Build slowly in time  
"Slow Fusion"

(Small example  
architectures, in  
practice much bigger)

# Early Fusion vs Late Fusion vs 3D CNN

Late  
Fusion

Early  
Fusion

3D  
CNN

	Size <b>(C x T x H x W)</b>	Receptive Field <b>(T x H x W)</b>
Layer		
	Input	$3 \times 20 \times 64 \times 64$
	Conv2D(3x3, 3->12)	$12 \times 20 \times 64 \times 64$
	Pool2D(4x4)	$12 \times 20 \times 16 \times 16$
	Conv2D(3x3, 12->24)	$24 \times 20 \times 16 \times 16$
Layer	GlobalAvgPool	$24 \times 1 \times 1 \times 1$
	Input	$3 \times 20 \times 64 \times 64$
	Conv2D(3x3, 3*20->12)	$12 \times 64 \times 64$
	Pool2D(4x4)	$12 \times 16 \times 16$
	Conv2D(3x3, 12->24)	$24 \times 16 \times 16$
Layer	GlobalAvgPool	$24 \times 1 \times 1$
	Input	$3 \times 20 \times 64 \times 64$
	Conv3D(3x3x3, 3->12)	$12 \times 20 \times 64 \times 64$
	Pool3D(4x4x4)	$12 \times 5 \times 16 \times 16$
	Conv3D(3x3x3, 12->24)	$24 \times 5 \times 16 \times 16$
Layer	GlobalAvgPool	$24 \times 1 \times 1$

Build slowly in space,  
All-at-once in time at end

Build slowly in space,  
All-at-once in time at start

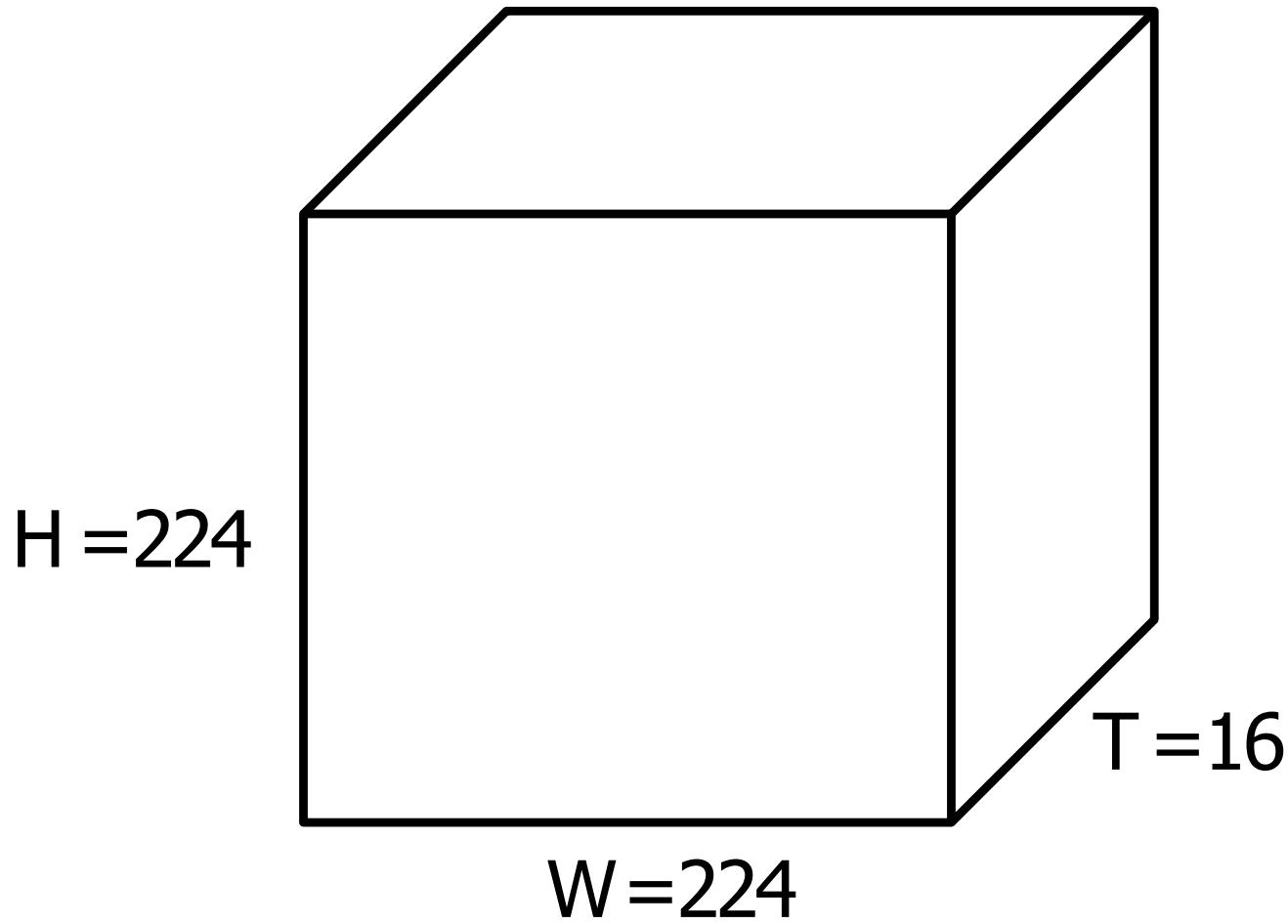
Build slowly in space,  
Build slowly in time  
"Slow Fusion"

What is the  
difference?

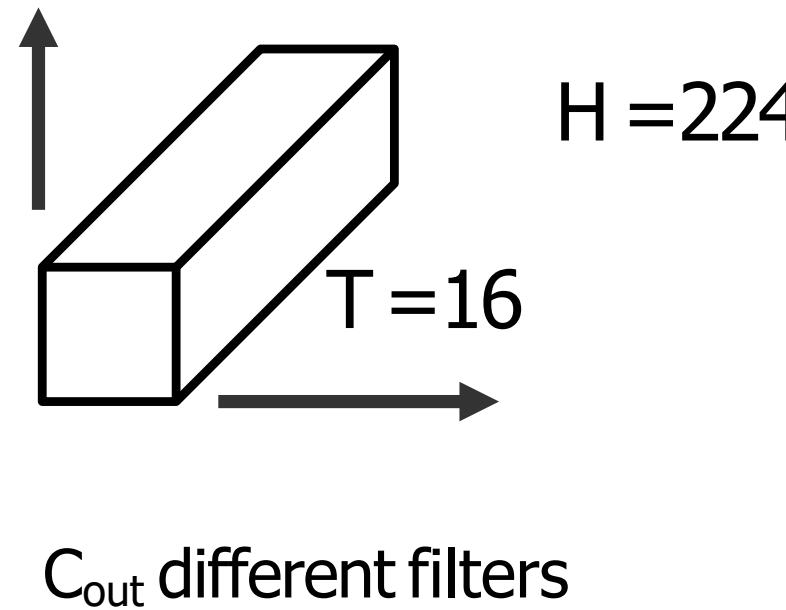
(Small example  
architectures, in  
practice much bigger)

# 2D Conv (Early Fusion) vs 3D Conv (3D CNN)

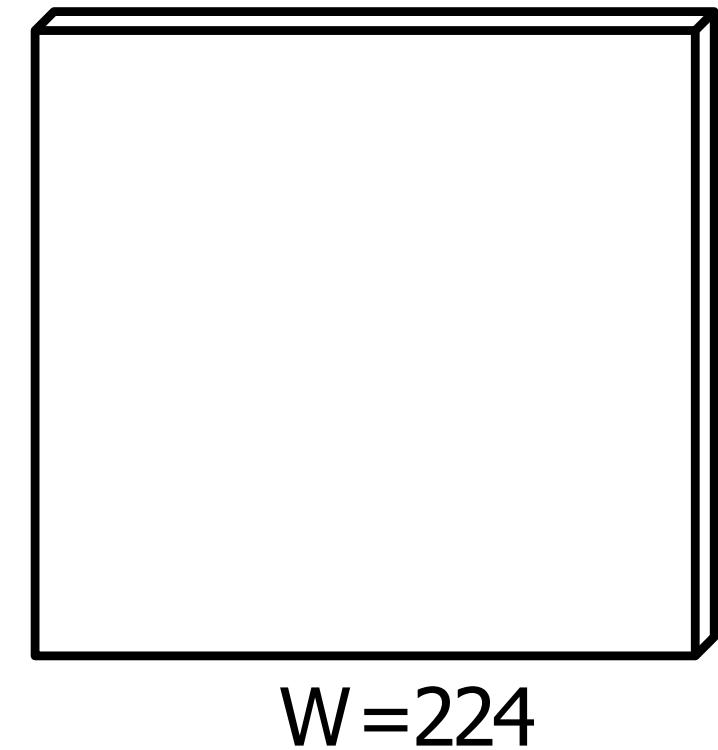
Input:  $C_{in} \times T \times H \times W$   
(3D grid with  $C_{in}$ -dim  
feat at each point)



Weight:  
 $C_{out} \times C_{in} \times T \times 3 \times 3$   
Slide over x and y

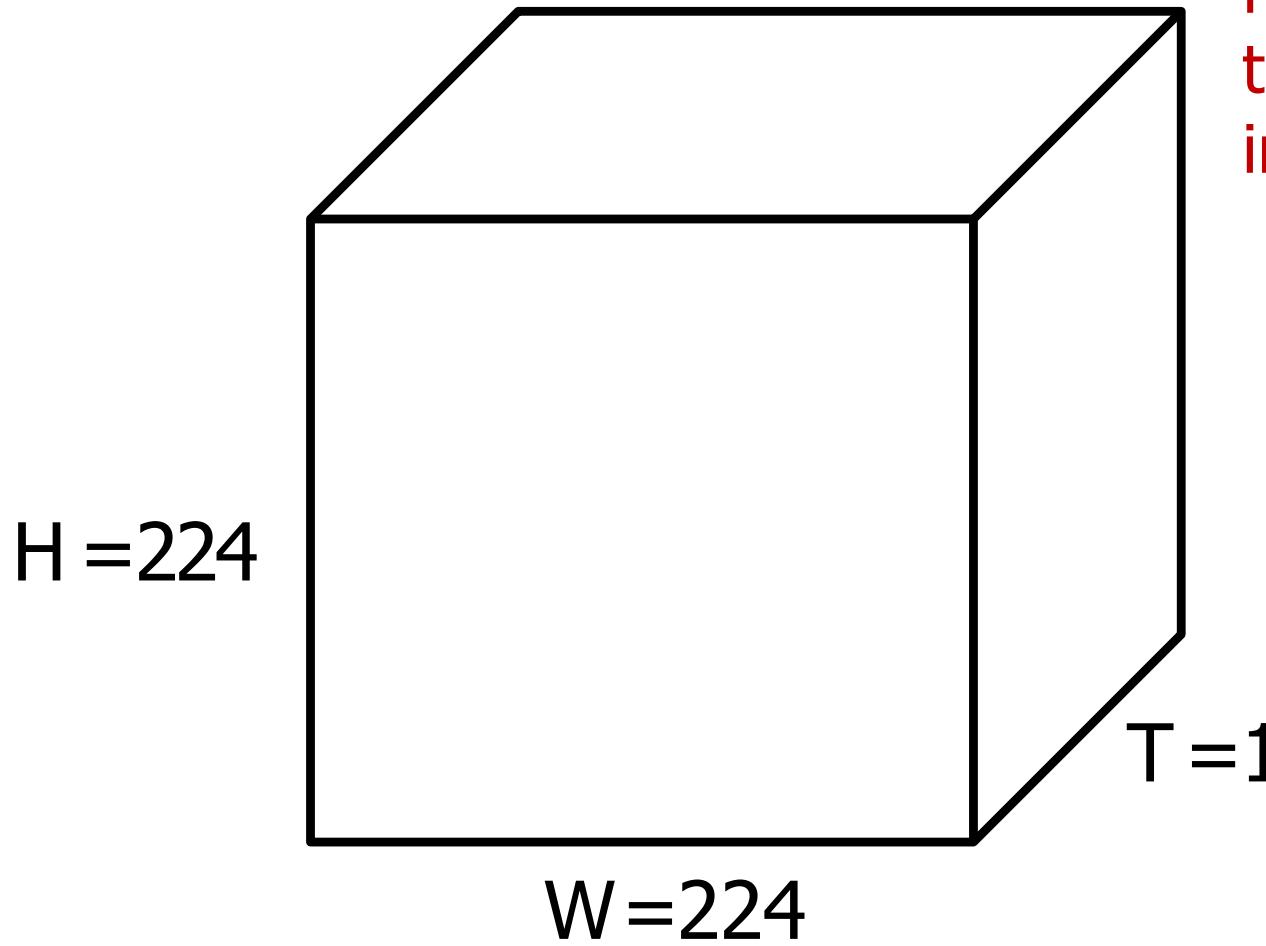


Output:  
 $C_{out} \times H \times W$   
2D grid with  $C_{out}$ -dim  
feat at each point



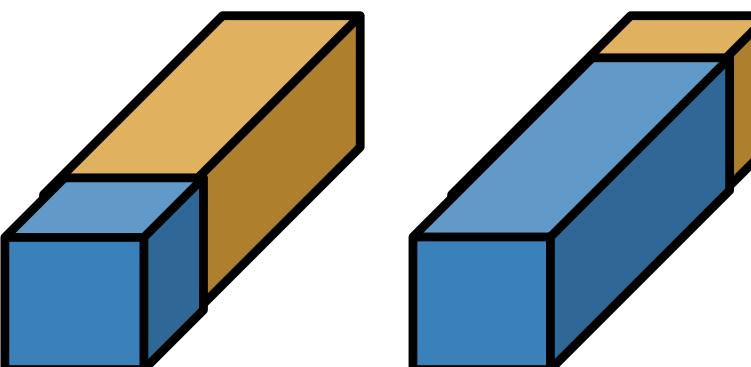
# 2D Conv (Early Fusion) vs 3D Conv (3D CNN)

Input:  $C_{in} \times T \times H \times W$   
(3D grid with  $C_{in}$ -dim feat  
at each point)



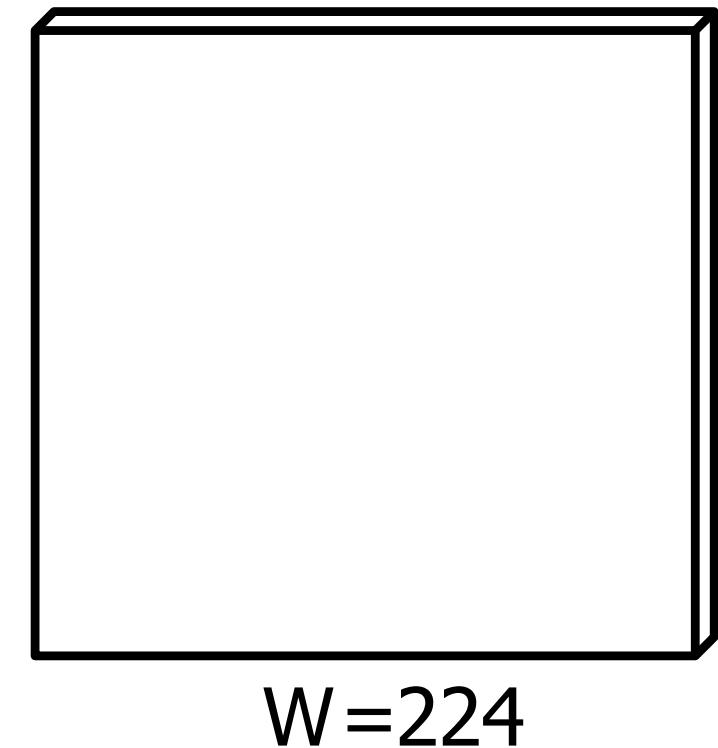
Weight:  
 $C_{out} \times C_{in} \times T \times 3 \times 3$   
Slide over x and y

No temporal shift-invariance!  
Needs to learn separate filters for  
the same motion at different times  
in the clip



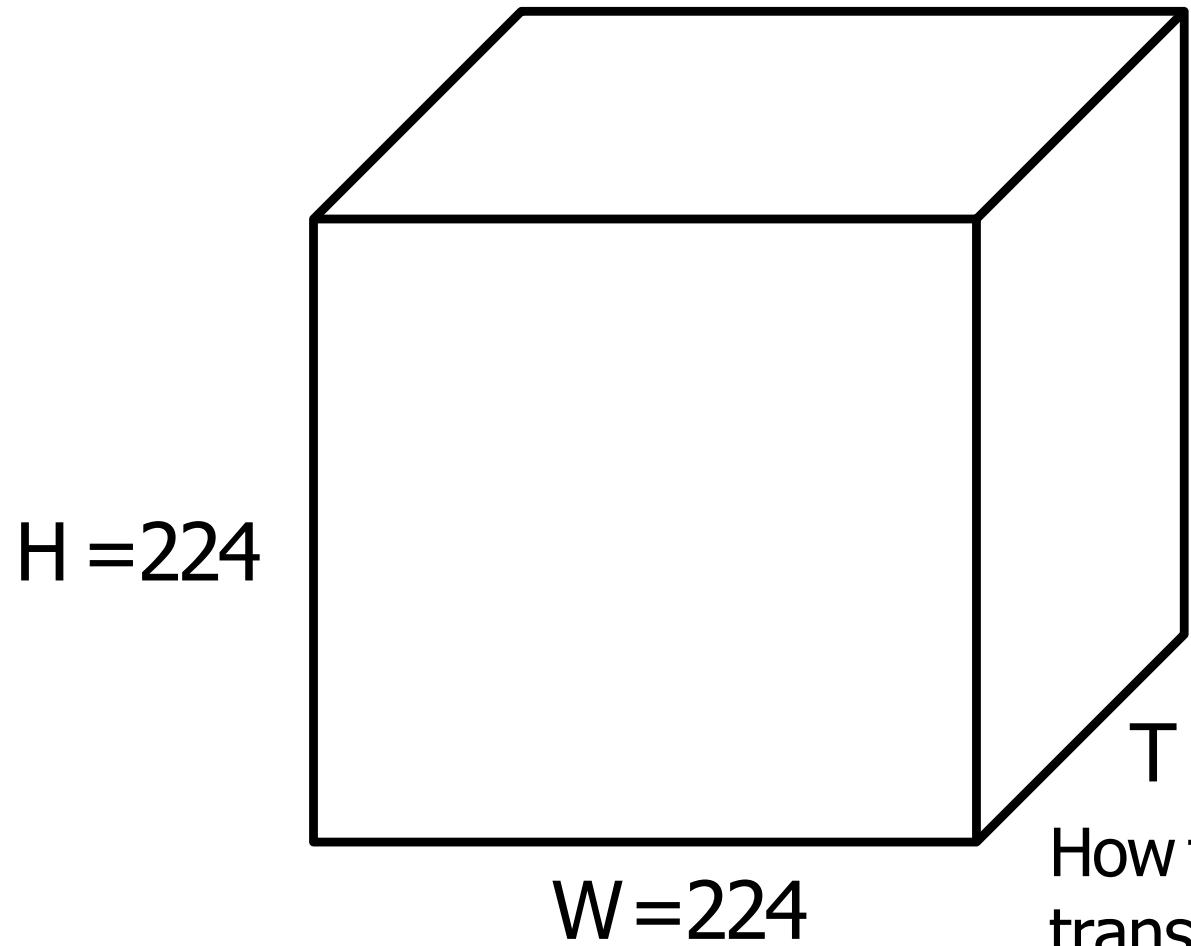
$C_{out}$  different filters

Output:  
 $C_{out} \times H \times W$   
2D grid with  $C_{out}$ -dim  
feat at each point



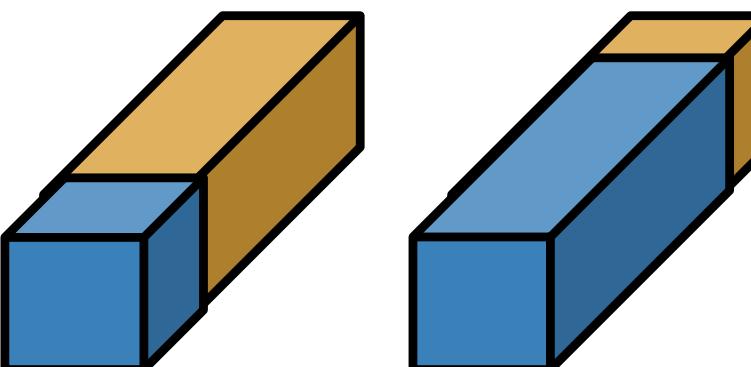
# 2D Conv (Early Fusion) vs 3D Conv (3D CNN)

Input:  $C_{in} \times T \times H \times W$   
(3D grid with  $C_{in}$ -dim feat  
at each point)



Weight:  
 $C_{out} \times C_{in} \times T \times 3 \times 3$   
Slide over x and y

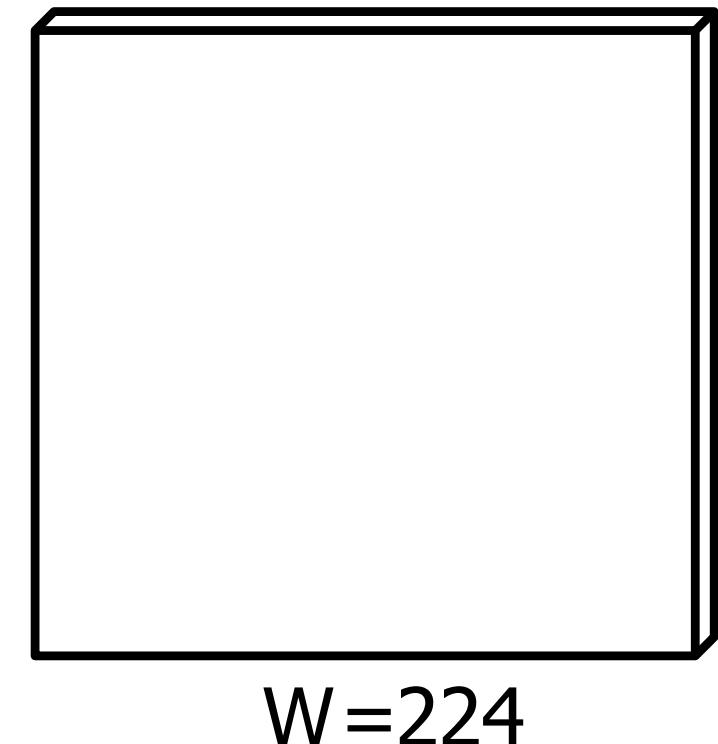
No temporal shift-invariance!  
Needs to learn separate filters for  
the same motion at different times  
in the clip



$C_{out}$  different filters

How to recognize blue to orange  
transitions anywhere in space and time?

Output:  
 $C_{out} \times H \times W$   
2D grid with  $C_{out}$ -dim  
feat at each point

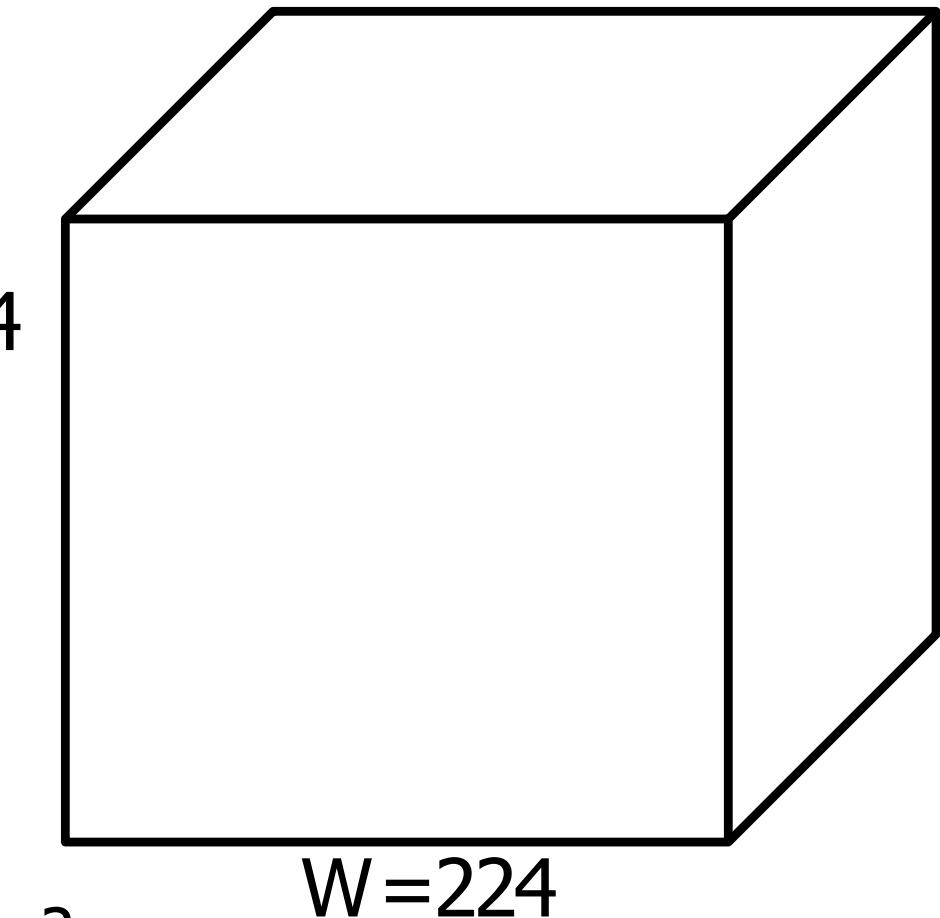
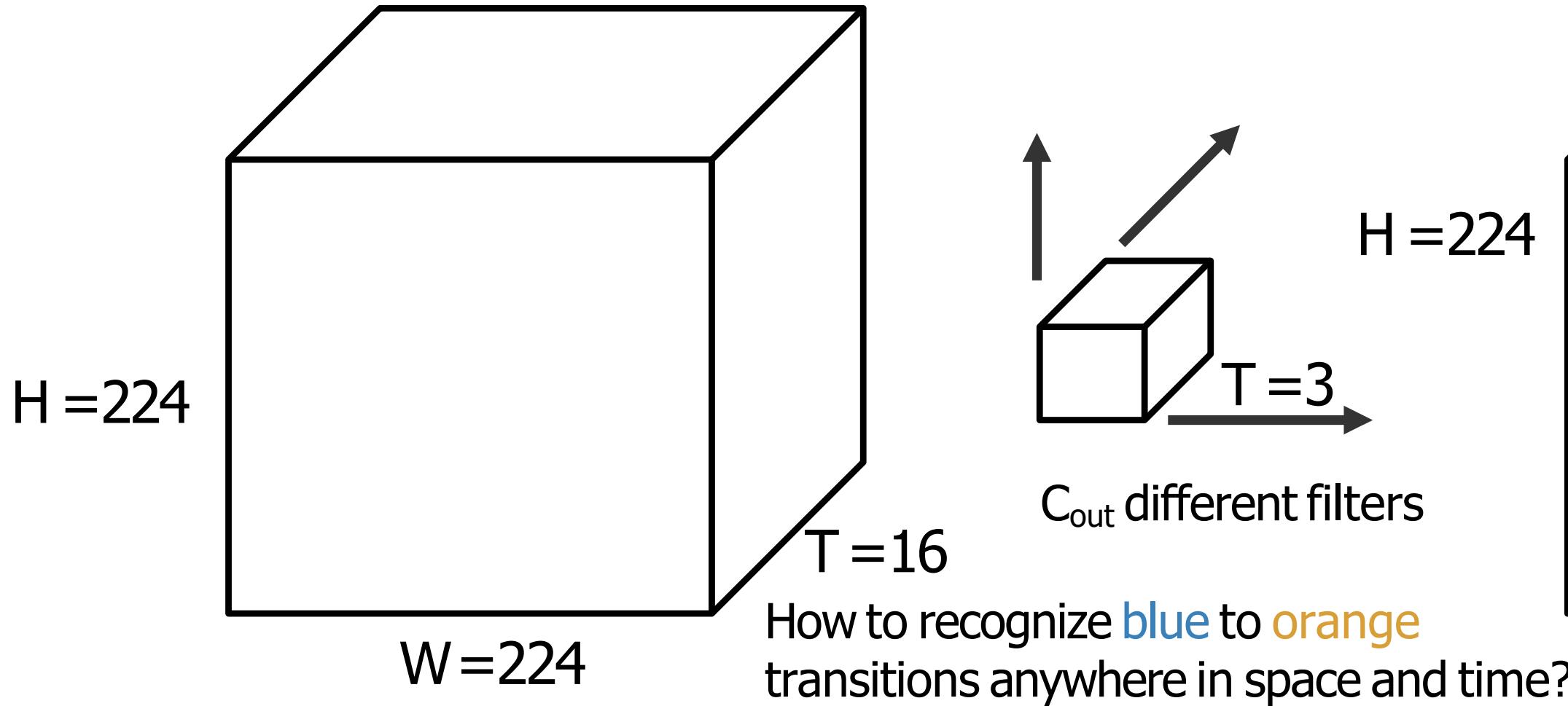


# 2D Conv (Early Fusion) vs 3D Conv (3D CNN)

Input:  $C_{in} \times T \times H \times W$   
(3D grid with  $C_{in}$ -dim  
feat at each point)

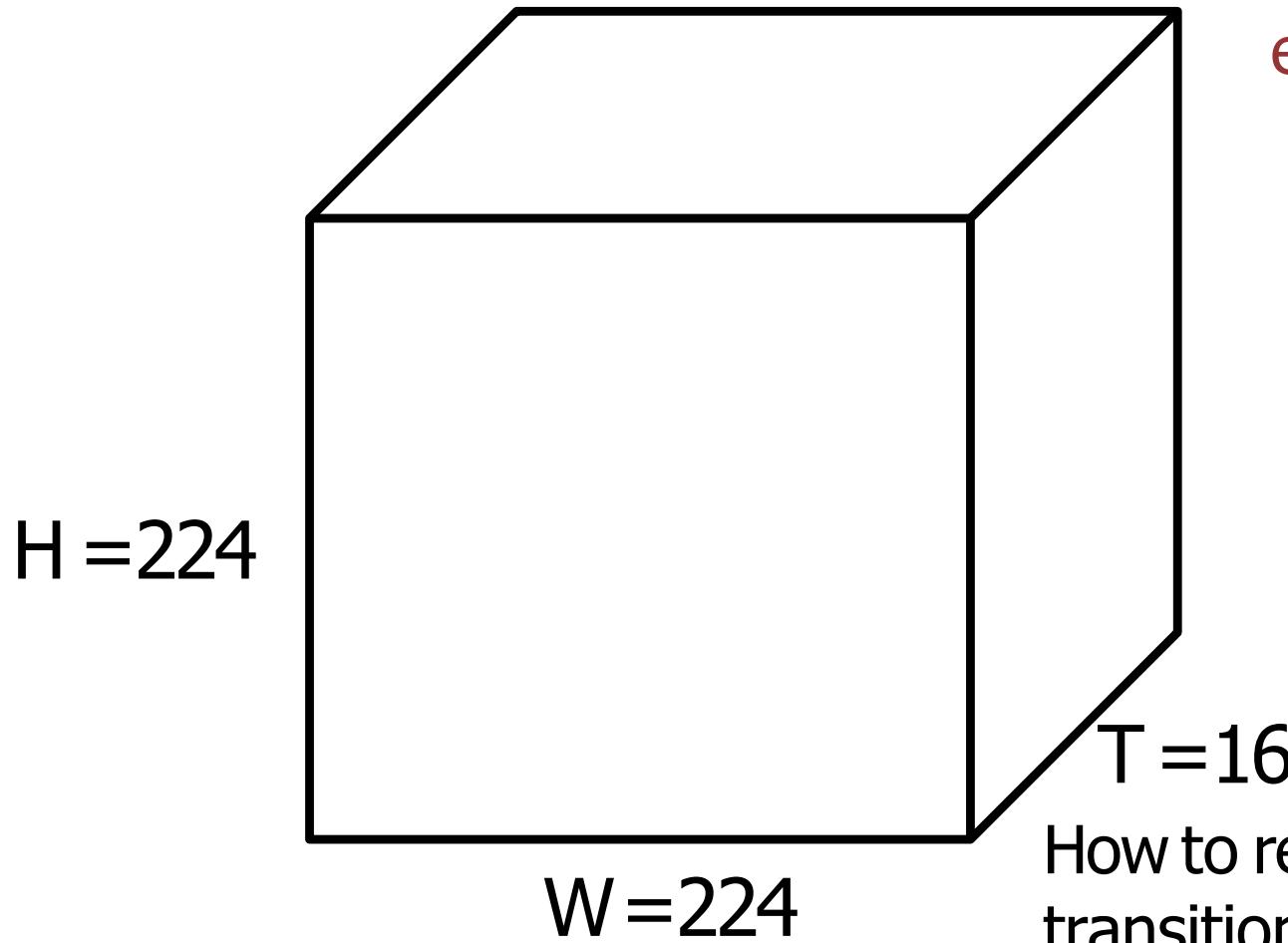
Weight:  
 $C_{out} \times C_{in} \times 3 \times 3 \times 3$   
Slide over x and y

Output:  
 $C_{out} \times T \times H \times W$   
3D grid with  $C_{out}$ -dim  
feat at each point



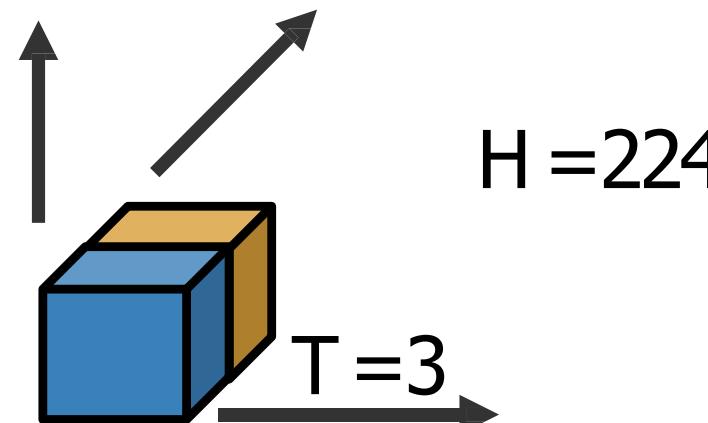
# 2D Conv (Early Fusion) vs 3D Conv (3D CNN)

Input:  $C_{in} \times T \times H \times W$   
(3D grid with  $C_{in}$ -dim  
feat at each point)



Weight:  
 $C_{out} \times C_{in} \times 3 \times 3 \times 3$   
Slide over x and y

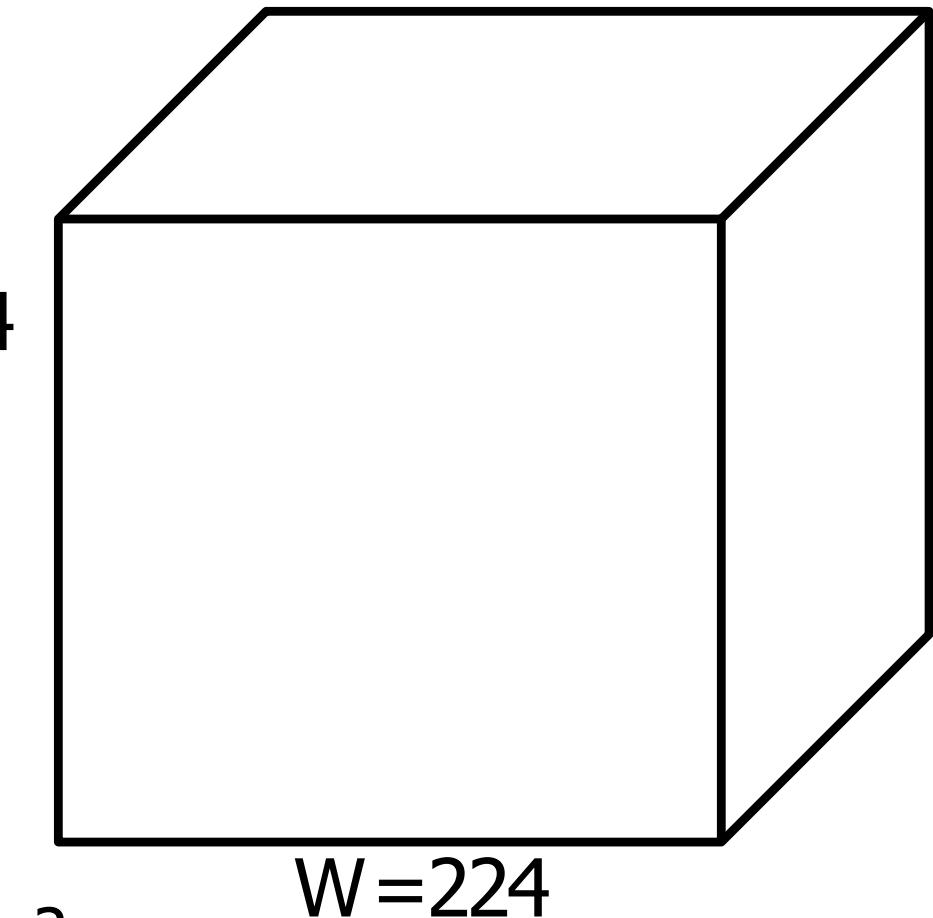
Temporal shift-invariant since  
each filter slides over time!



$C_{out}$  different filters

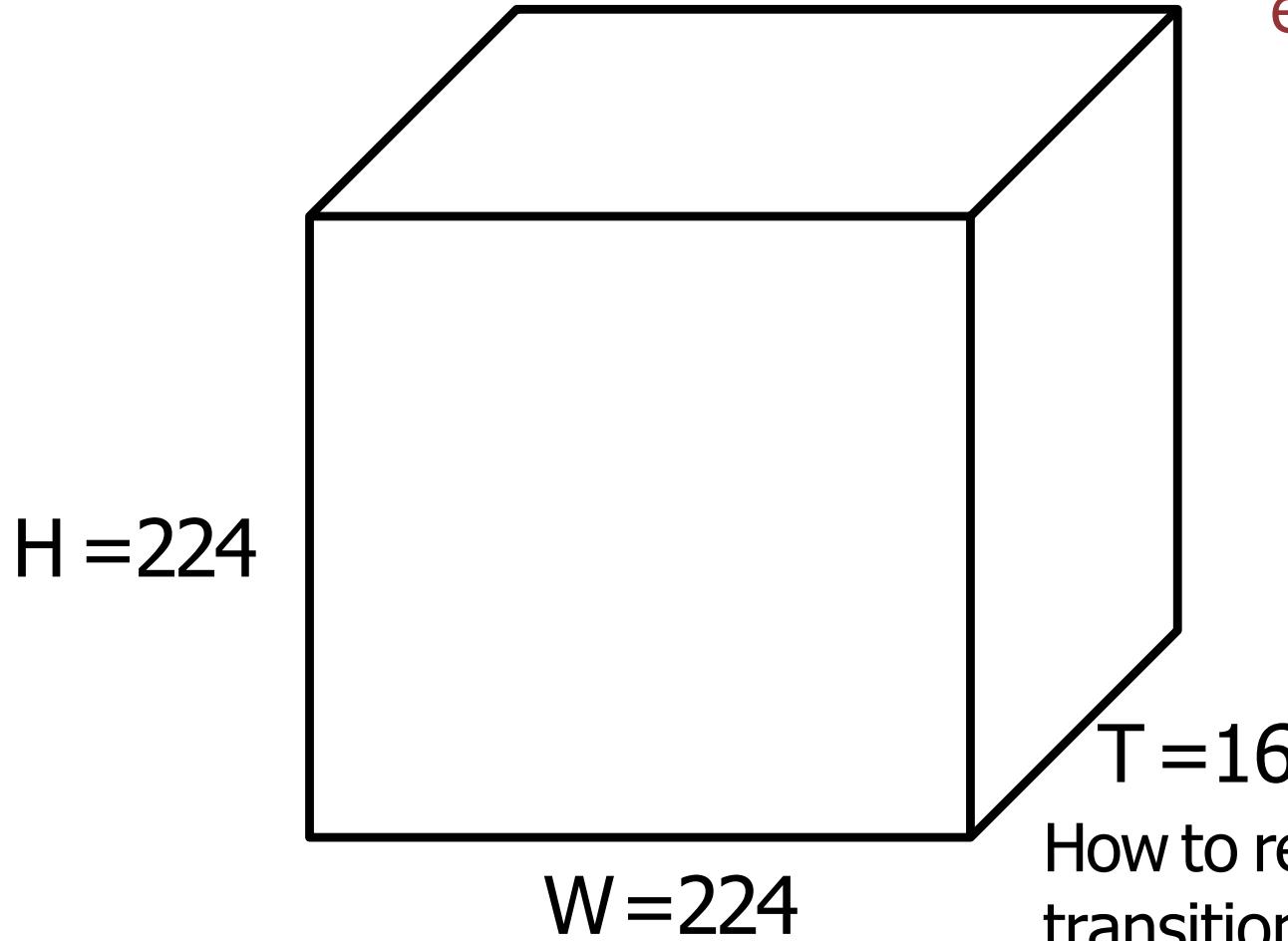
How to recognize **blue** to **orange**  
transitions anywhere in space and time?

Output:  
 $C_{out} \times T \times H \times W$   
3D grid with  $C_{out}$ -dim  
feat at each point



# 2D Conv (Early Fusion) vs 3D Conv (3D CNN)

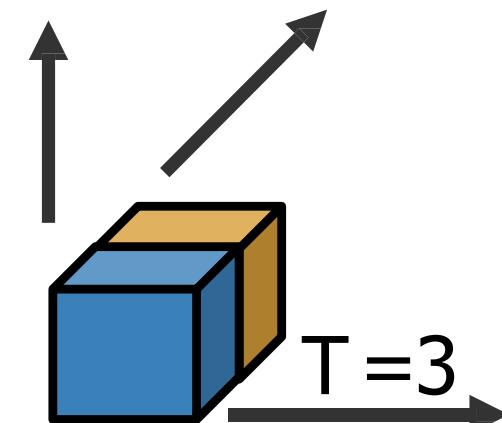
Input:  $C_{in} \times T \times H \times W$   
(3D grid with  $C_{in}$ -dim  
feat at each point)



How to recognize **blue** to **orange**  
transitions anywhere in space and time?

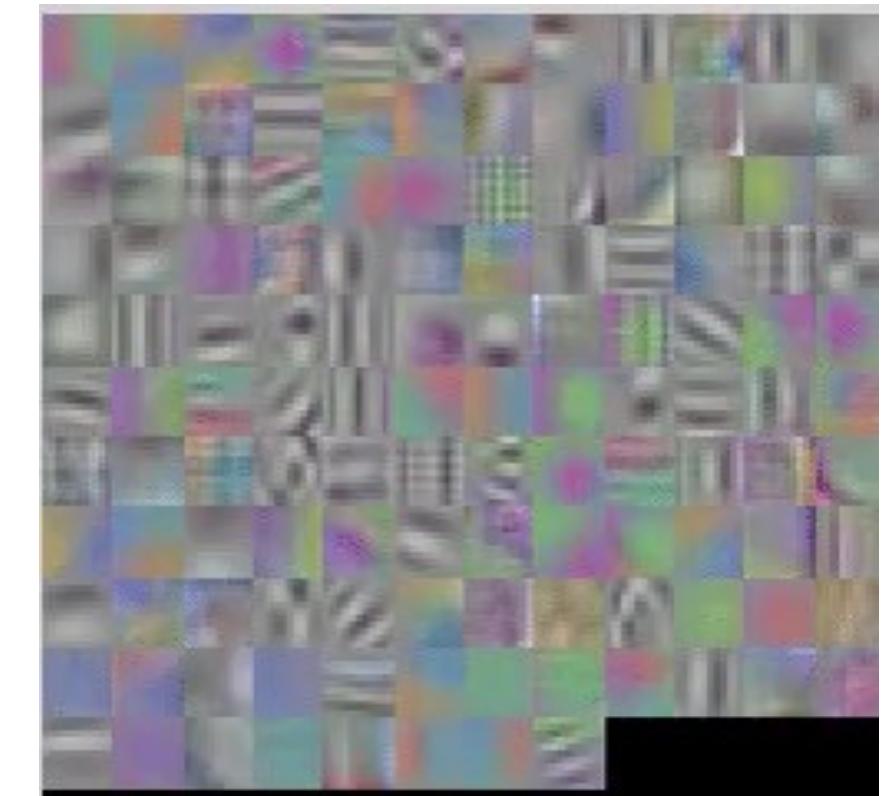
Weight:  
 $C_{out} \times C_{in} \times 3 \times 3 \times 3$   
Slide over x and y

Temporal shift-invariant since  
each filter slides over time!



$C_{out}$  different filters

First-layer filters have shape  
3 (RGB)  $\times$  4 (frames)  $\times$  5  $\times$  5  
(space)  
Can visualize as video clips!



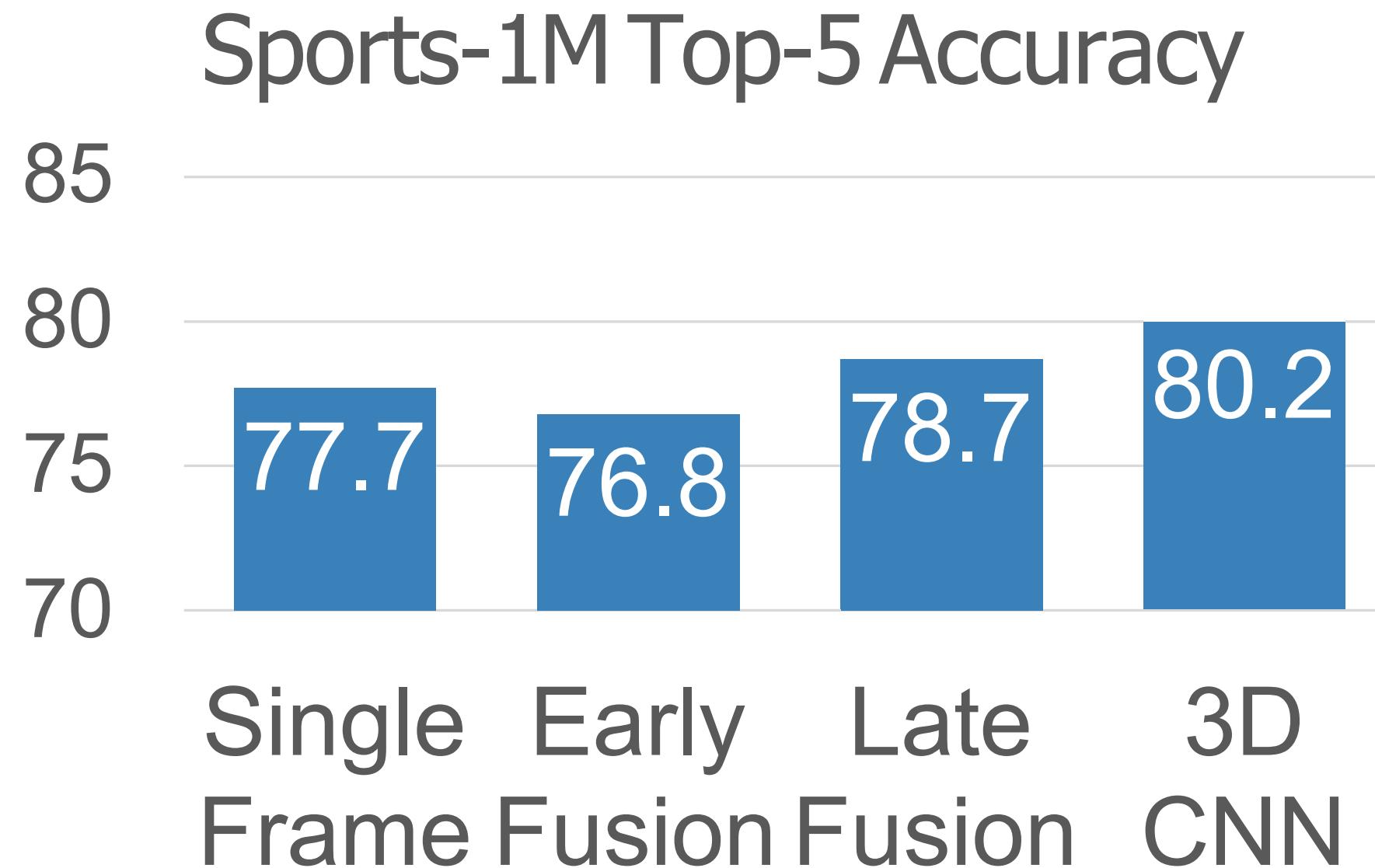
# Example Video Dataset: Sports-1M



1 million YouTube videos  
annotated with labels for 487  
different types of sports

Ground Truth  
Correct prediction  
Incorrect prediction

# Early Fusion vs Late Fusion vs 3D CNN



# C3D: The VGG of 3D CNNs

3D CNN that uses all 3x3x3 conv and  
2x2x2 pooling  
(except Pool1 which is 1x2x2)

Released model pretrained on  
Sports-1M: Many people used this as  
a video feature extractor

Tran et al, "Learning Spatiotemporal Features with 3D Convolutional Networks", ICCV 2015

Layer	Size
Input	3 x 16 x 112 x 112
Conv1 (3x3x3)	64 x 16 x 112 x 112
Pool1 (1x2x2)	64 x 16 x 56 x 56
Conv2 (3x3x3)	128 x 16 x 56 x 56
Pool2 (2x2x2)	128 x 8 x 28 x 28
Conv3a (3x3x3)	256 x 8 x 28 x 28
Conv3b (3x3x3)	256 x 8 x 28 x 28
Pool3 (2x2x2)	256 x 4 x 14 x 14
Conv4a (3x3x3)	512 x 4 x 14 x 14
Conv4b (3x3x3)	512 x 4 x 14 x 14
Pool4 (2x2x2)	512 x 2 x 7 x 7
Conv5a (3x3x3)	512 x 2 x 7 x 7
Conv5b (3x3x3)	512 x 2 x 7 x 7
Pool5	512 x 1 x 3 x 3
FC6	4096
FC7	4096
FC8	C

# C3D: The VGG of 3D CNNs

3D CNN that uses all  $3 \times 3 \times 3$  conv and  
 $2 \times 2 \times 2$  pooling  
(except Pool1 which is  $1 \times 2 \times 2$ )

Released model pretrained on  
Sports-1M: Many people used this as  
a video feature extractor

Problem:  $3 \times 3 \times 3$  conv is very  
expensive!

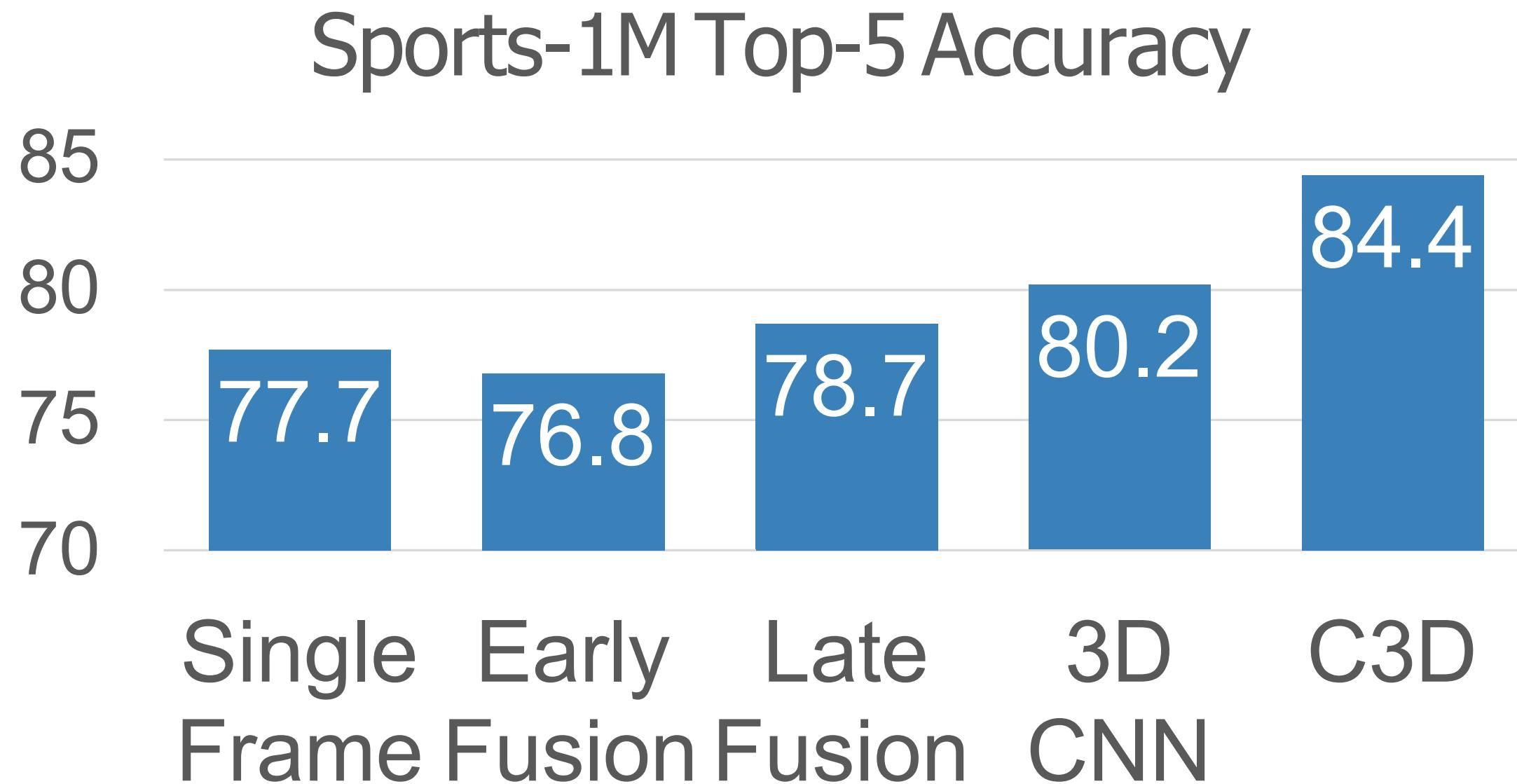
AlexNet: 0.7 GFLOP

VGG-16: 13.6 GFLOP

C3D: 39.5 GFLOP (2.9x VGG!)

Layer	Size	MFLOPs
Input	$3 \times 16 \times 112 \times 112$	
Conv1 ( $3 \times 3 \times 3$ )	$64 \times 16 \times 112 \times 112$	1.04
Pool1 ( $1 \times 2 \times 2$ )	$64 \times 16 \times 56 \times 56$	
Conv2 ( $3 \times 3 \times 3$ )	$128 \times 16 \times 56 \times 56$	11.10
Pool2 ( $2 \times 2 \times 2$ )	$128 \times 8 \times 28 \times 28$	
Conv3a ( $3 \times 3 \times 3$ )	$256 \times 8 \times 28 \times 28$	5.55
Conv3b ( $3 \times 3 \times 3$ )	$256 \times 8 \times 28 \times 28$	11.10
Pool3 ( $2 \times 2 \times 2$ )	$256 \times 4 \times 14 \times 14$	
Conv4a ( $3 \times 3 \times 3$ )	$512 \times 4 \times 14 \times 14$	2.77
Conv4b ( $3 \times 3 \times 3$ )	$512 \times 4 \times 14 \times 14$	5.55
Pool4 ( $2 \times 2 \times 2$ )	$512 \times 2 \times 7 \times 7$	
Conv5a ( $3 \times 3 \times 3$ )	$512 \times 2 \times 7 \times 7$	0.69
Conv5b ( $3 \times 3 \times 3$ )	$512 \times 2 \times 7 \times 7$	0.69
Pool5	$512 \times 1 \times 3 \times 3$	
FC6	4096	0.51
FC7	4096	0.45
FC8	C	0.05

# Early Fusion vs Late Fusion vs 3D CNN



Karpathy et al, "Large-scale Video Classification with Convolutional Neural Networks", CVPR 2014  
Tran et al, "Learning Spatiotemporal Features with 3D Convolutional Networks", ICCV 2015

# Recognizing Actions from Motion

We can easily recognize actions using only motion information



Johansson, "Visual perception of biological motion and a model for its analysis." Perception & Psychophysics. 14(2):201-211. 1973.

# Measuring Motion: Optical Flow

Image at frame t



Image at frame t+1

Simonyan and Zisserman, "Two-stream convolutional networks for action recognition in videos", NeurIPS 2014

# Measuring Motion: Optical Flow

Image at frame t

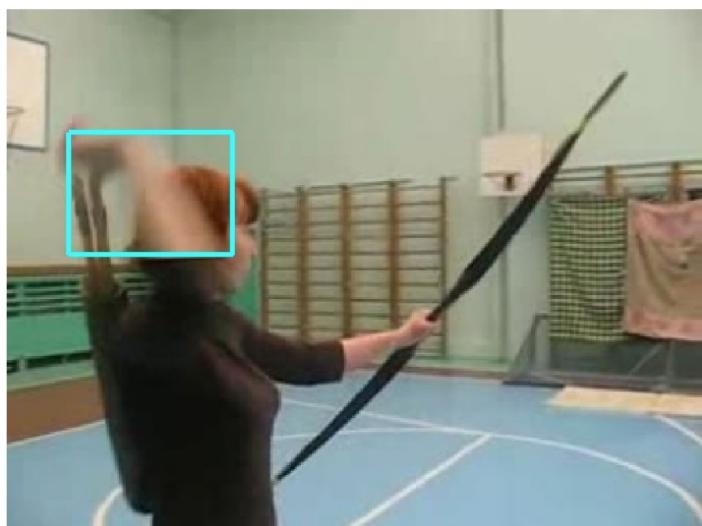
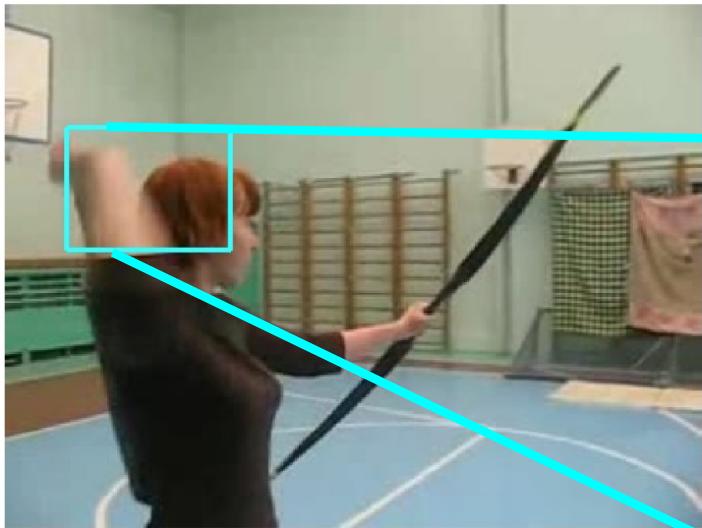
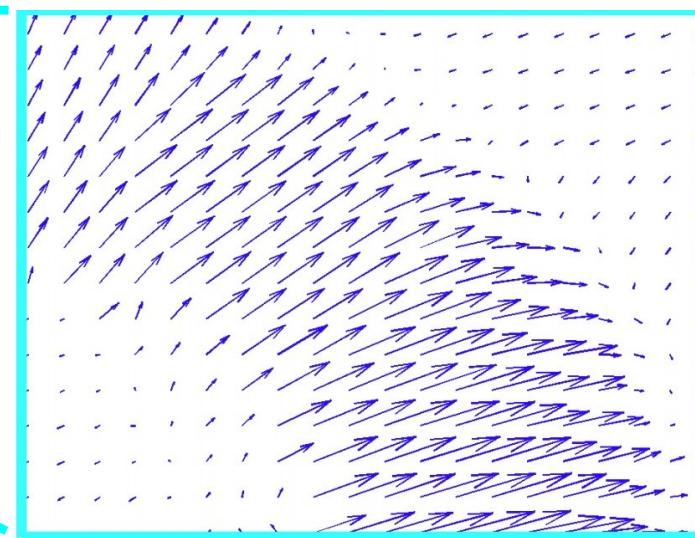


Image at frame t+1

Optical flow gives a displacement field  $F$  between images  $I_t$  and  $I_{t+1}$



Tells where each pixel will move in the next frame:  
 $F(x, y) = (dx, dy)$   
 $I_{t+1}(x+dx, y+dy) = I_t(x, y)$

# Measuring Motion: Optical Flow

Image at frame t

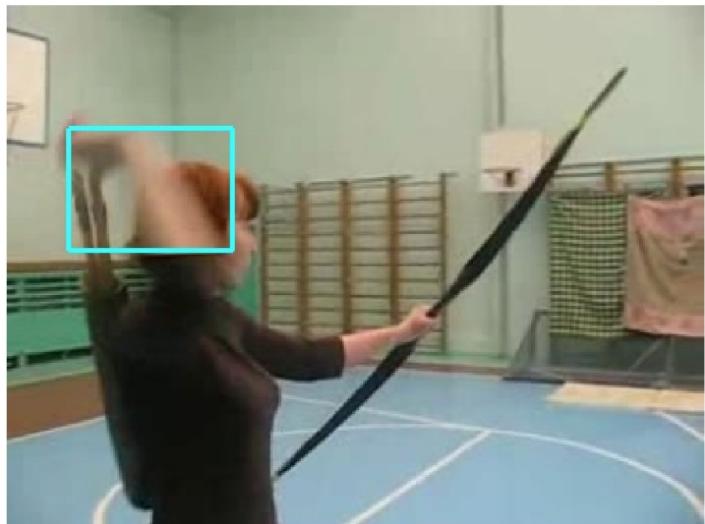
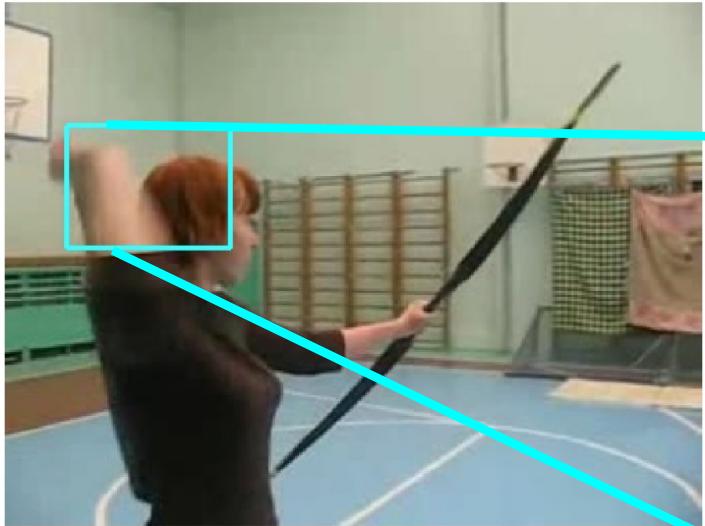
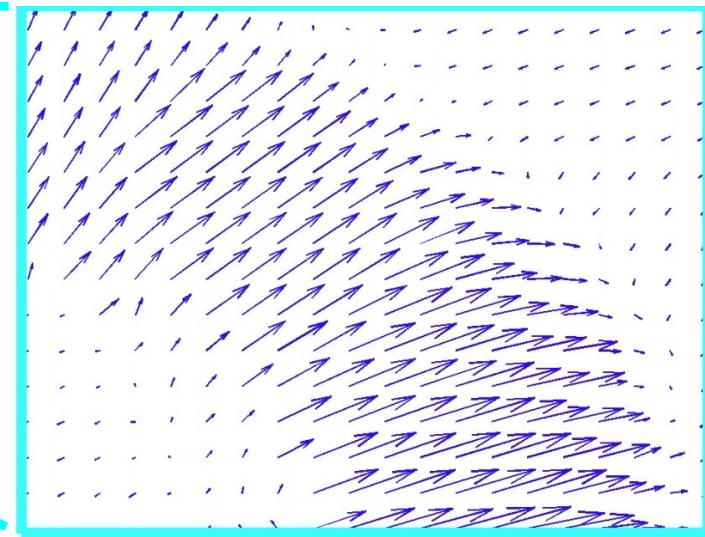


Image at frame  $t+1$

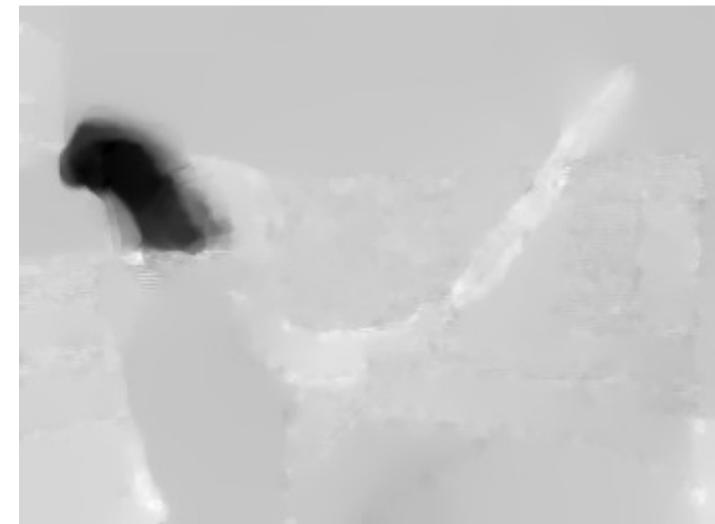
Optical flow gives a displacement field  $F$  between images  $I_t$  and  $I_{t+1}$



Tells where each pixel will move in the next frame:  
 $F(x, y) = (dx, dy)$   
 $I_{t+1}(x+dx, y+dy) = I_t(x, y)$

Optical Flow highlights local motion

Horizontal flow  $dx$

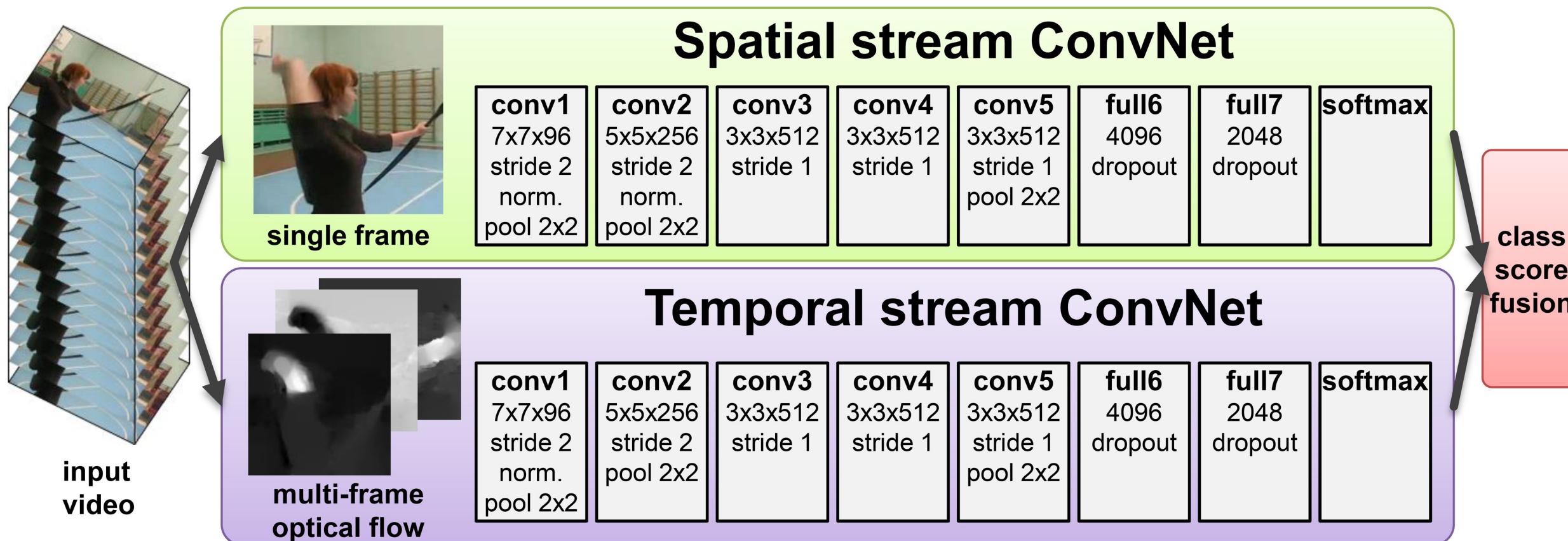


Vertical Flow  $dy$

# Separating Motion and Appearance: Two-Stream Networks

Input: Single Image

$3 \times H \times W$



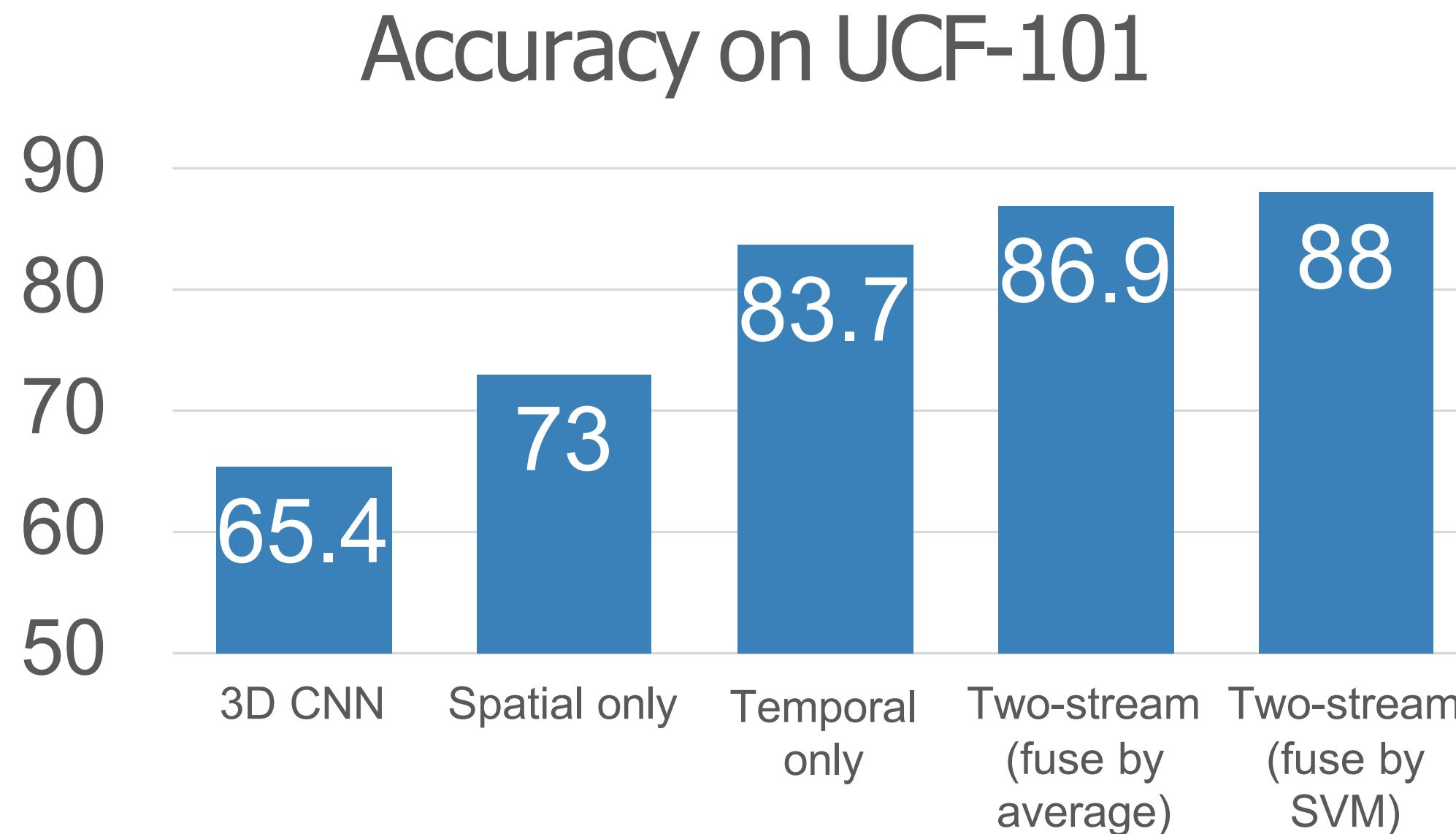
Input: Stack of optical flow:

$[2*(T-1)] \times H \times W$

Early fusion: First 2D conv  
processes all flow images

Simonyan and Zisserman, "Two-stream convolutional networks for action recognition in videos", NeurIPS 2014

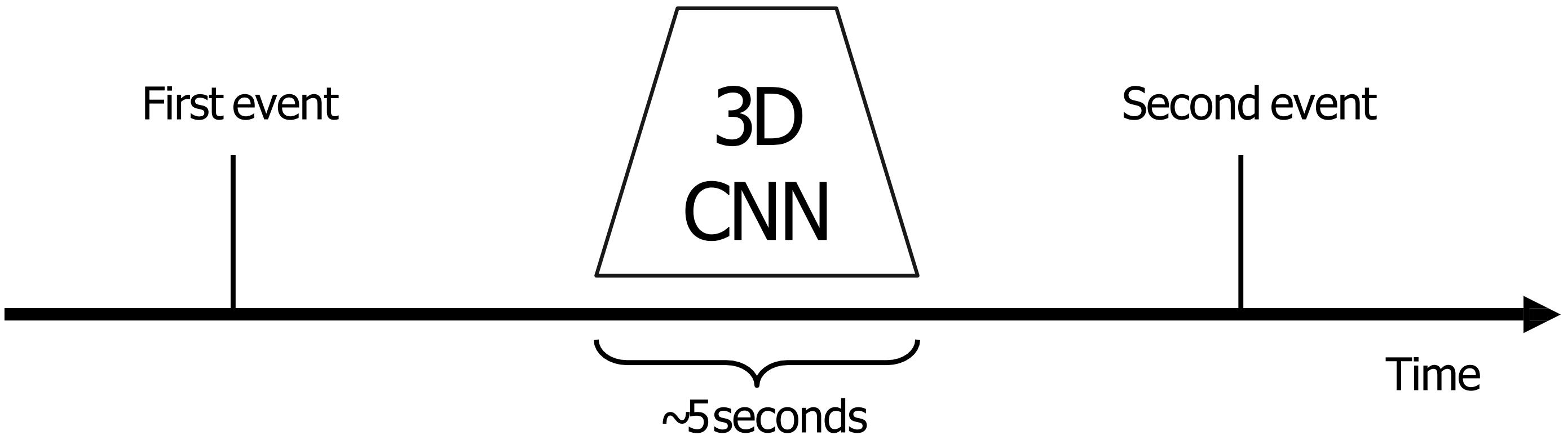
# Separating Motion and Appearance: Two-Stream Networks



Simonyan and Zisserman, "Two-stream convolutional networks for action recognition in videos", NeurIPS 2014

# Modeling long-term temporal structure

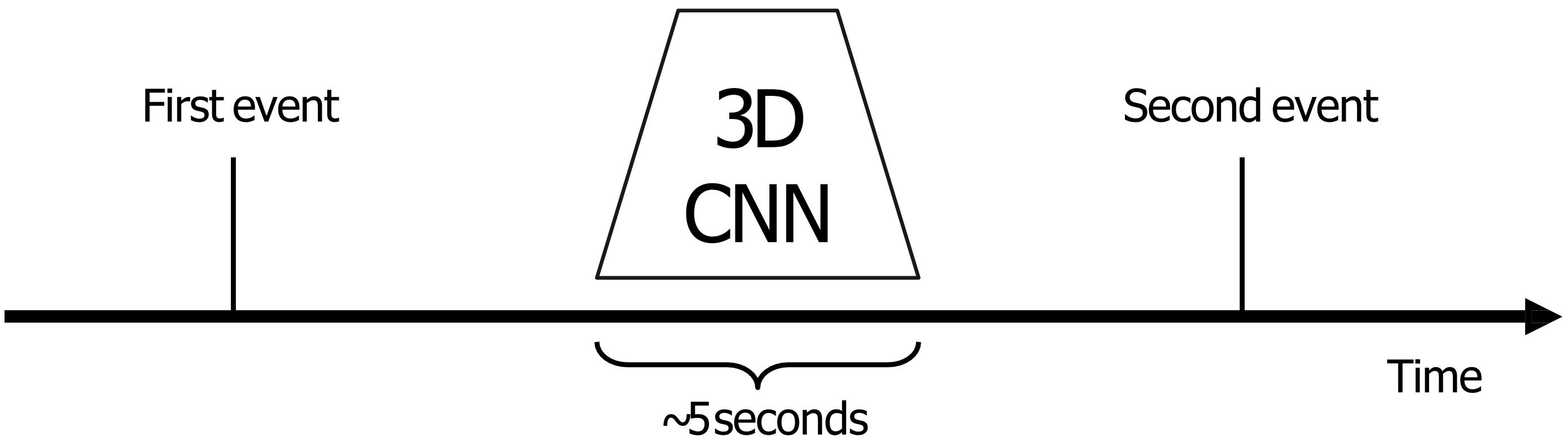
So far all our temporal CNNs only model local motion between frames in very short clips of ~2-5 seconds. What about long-term structure?



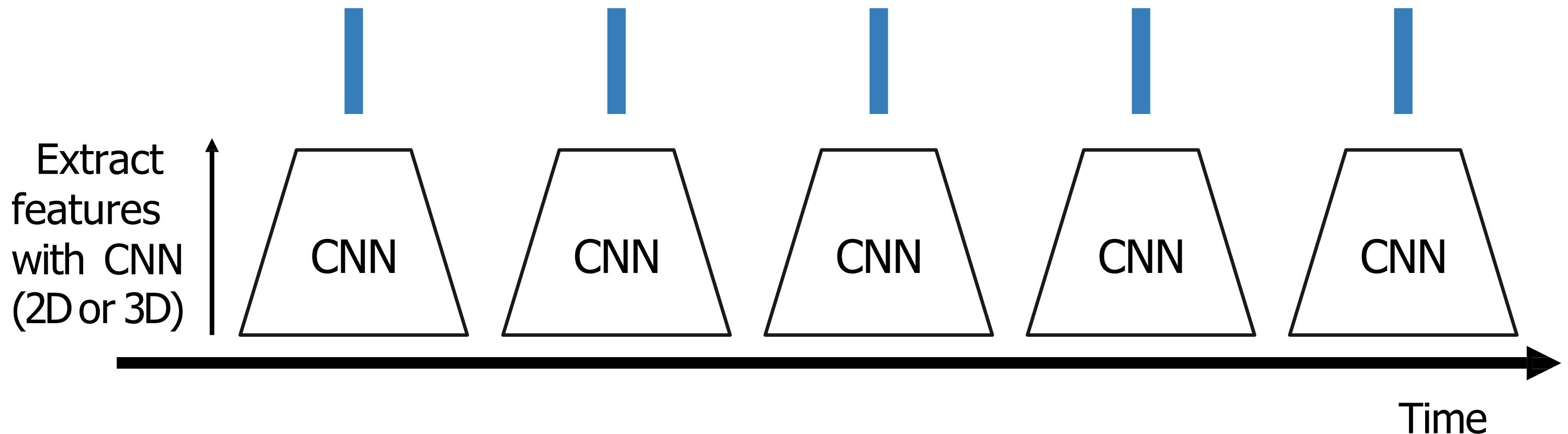
# Modeling long-term temporal structure

So far all our temporal CNNs only model local motion between frames in very short clips of ~2-5 seconds. What about long-term structure?

We know how to handle sequences! How about recurrent networks?

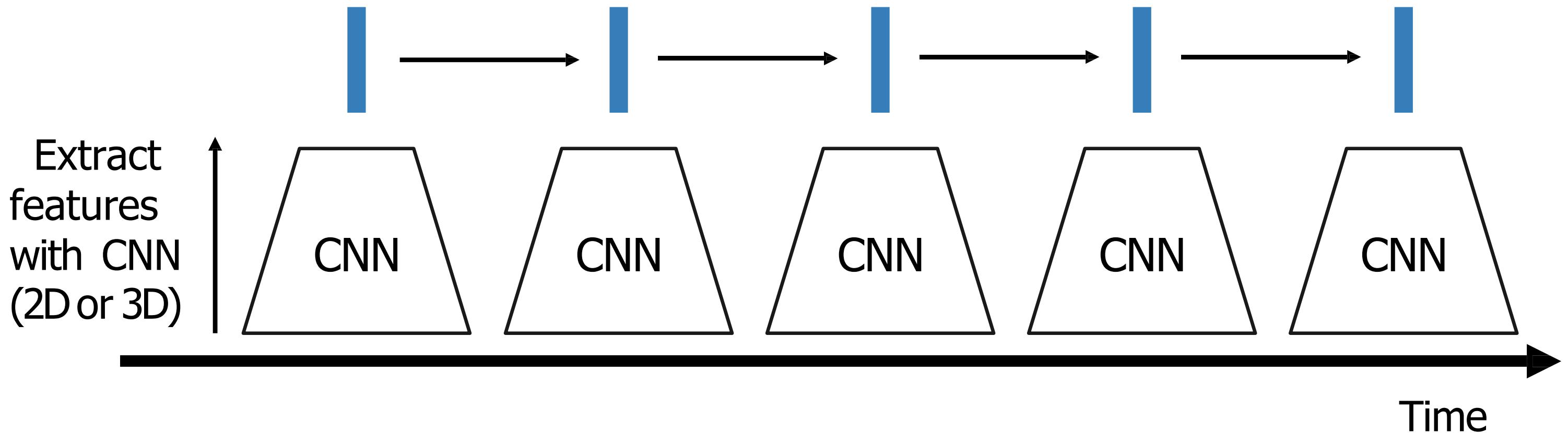


# Modeling long-term temporal structure



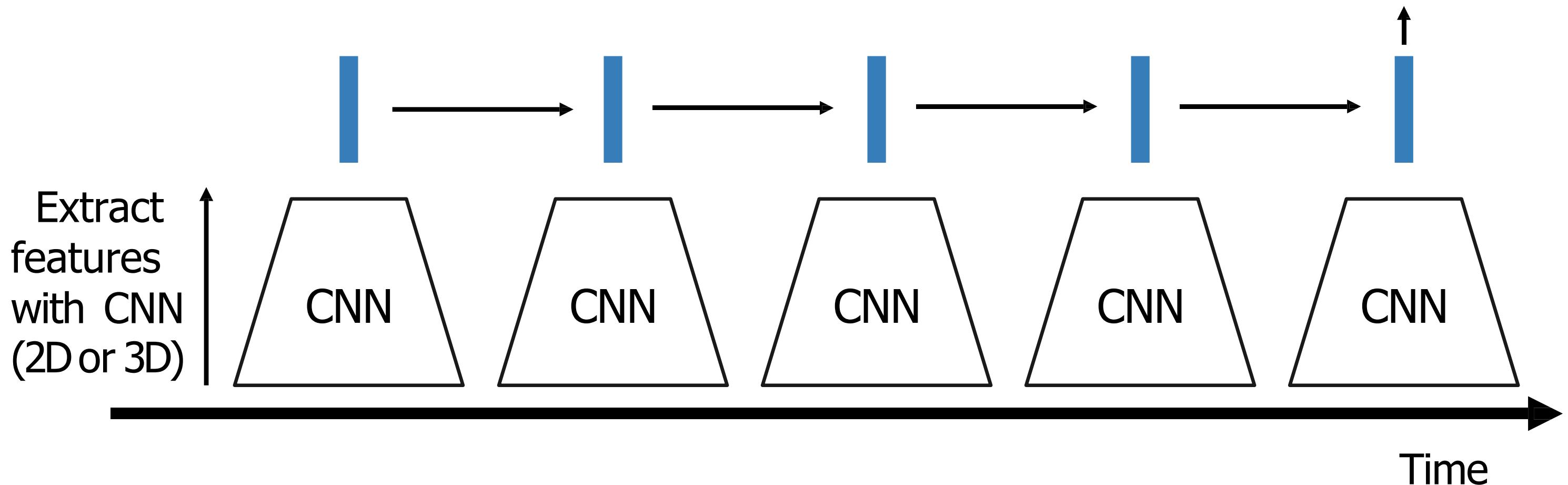
# Modeling long-term temporal structure

Process local features using recurrent network (e.g. LSTM)



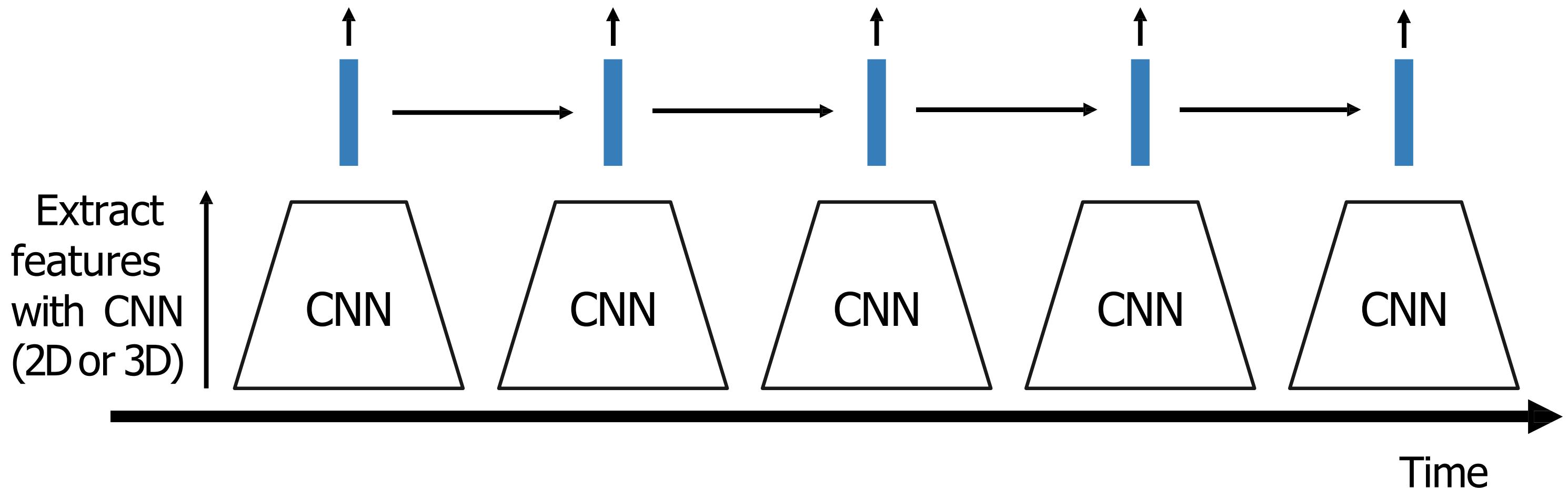
# Modeling long-term temporal structure

Process local features using recurrent network (e.g. LSTM)  
Many to one: One output at end of video



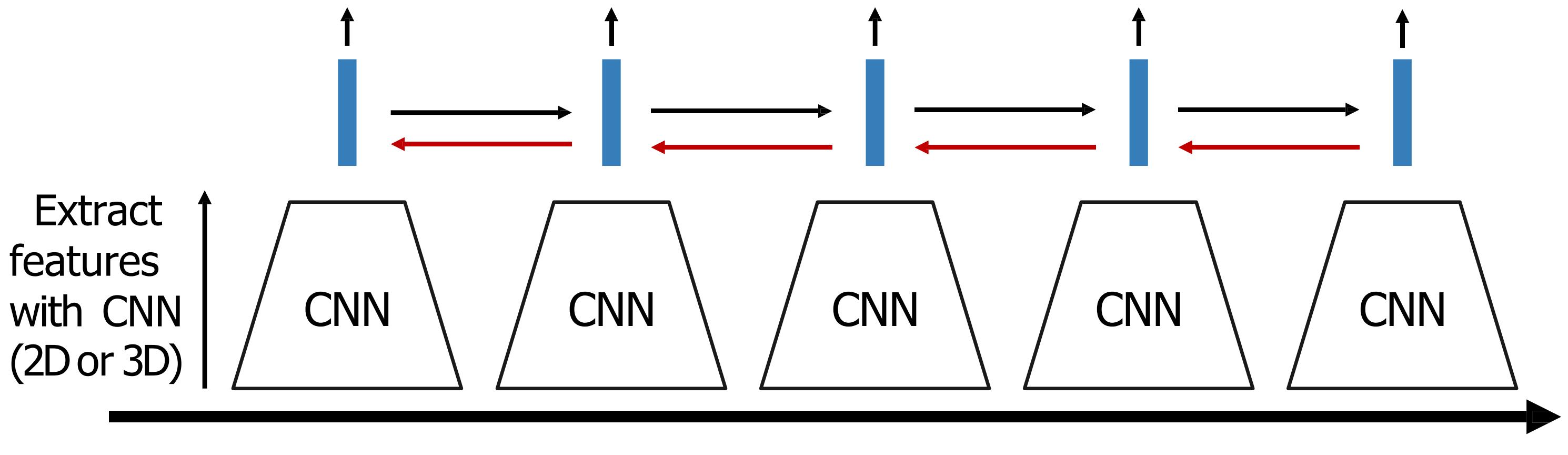
# Modeling long-term temporal structure

Process local features using recurrent network (e.g. LSTM)  
Many to many: one output per video frame



# Modeling long-term temporal structure

Sometimes don't backprop to CNN to save memory;  
pretrain and use it as a feature extractor



Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011

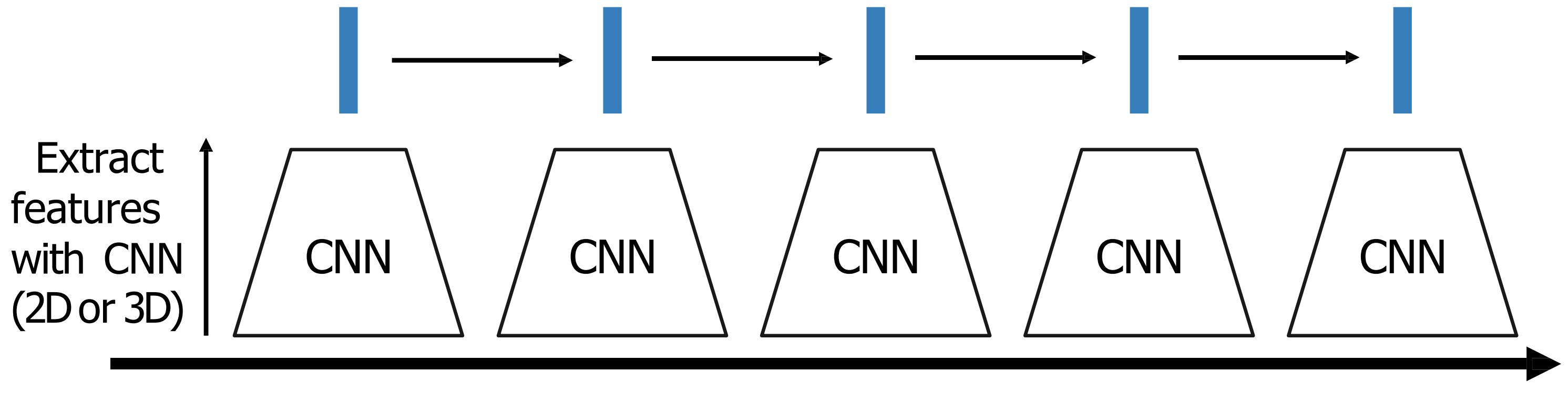
Donahue et al, "Long-term recurrent convolutional networks for visual recognition and description", CVPR 2015

# Modeling long-term temporal structure

Inside CNN: Each value is a function of a fixed temporal window (local temporal structure)

Inside RNN: Each vector is a function of all previous vectors (global temporal structure)

Can we merge both approaches?

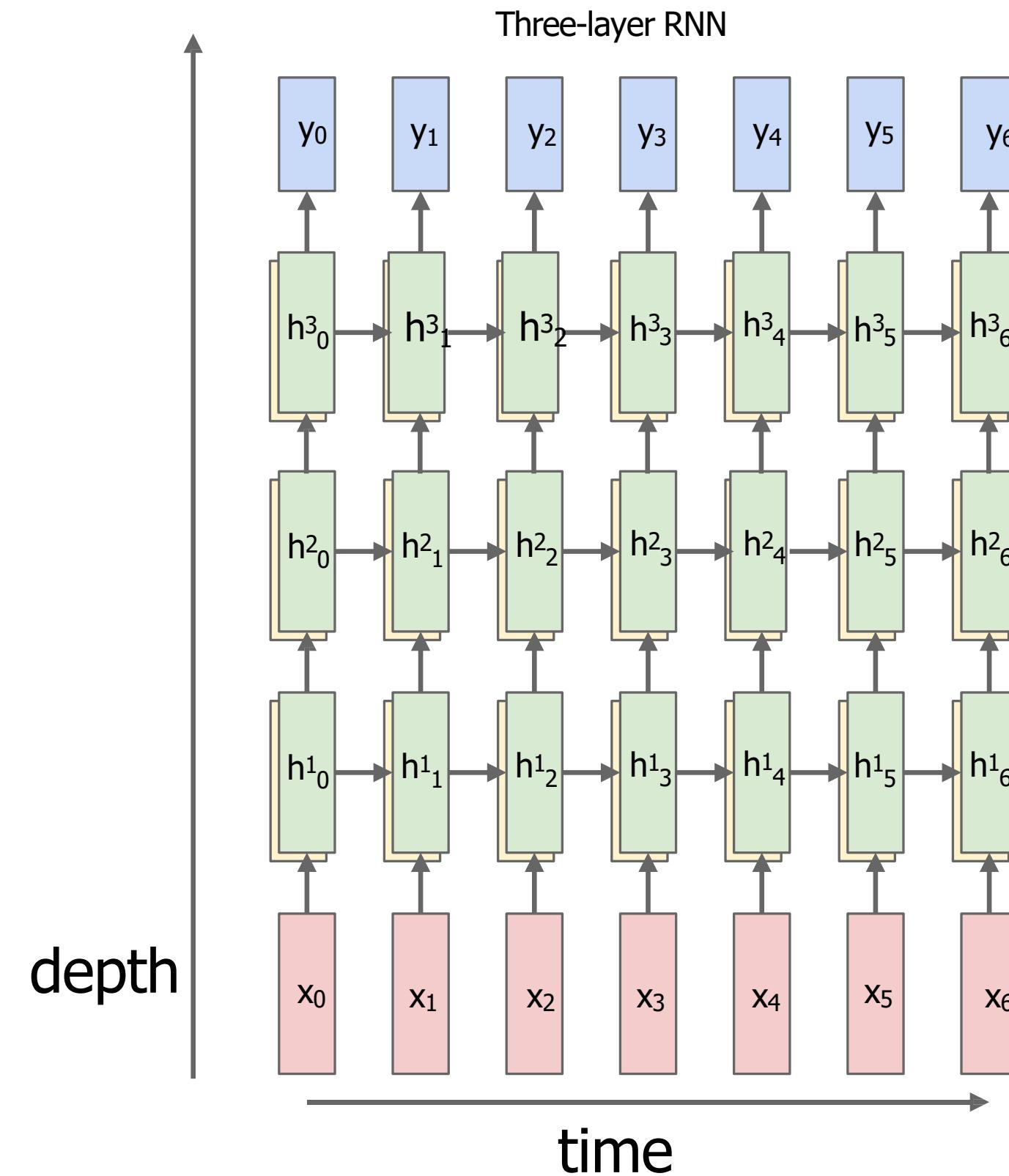


Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011

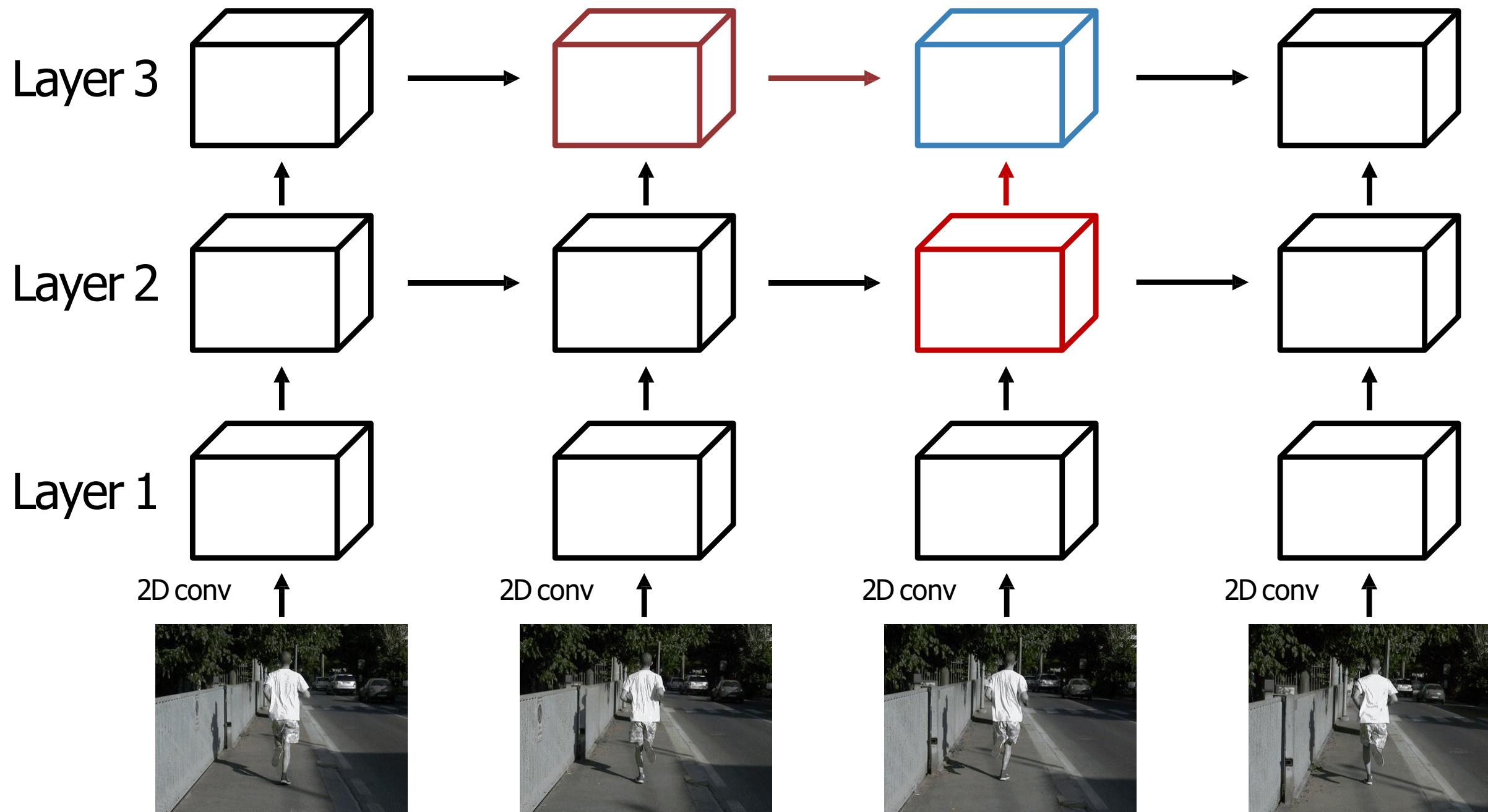
Donahue et al, "Long-term recurrent convolutional networks for visual recognition and description", CVPR 2015

# Recall: Multi-layer RNN

We can use a similar structure to process videos!



# Recurrent Convolutional Network



Entire network uses 2D feature maps:  $C \times H \times W$

Each depends on two inputs:

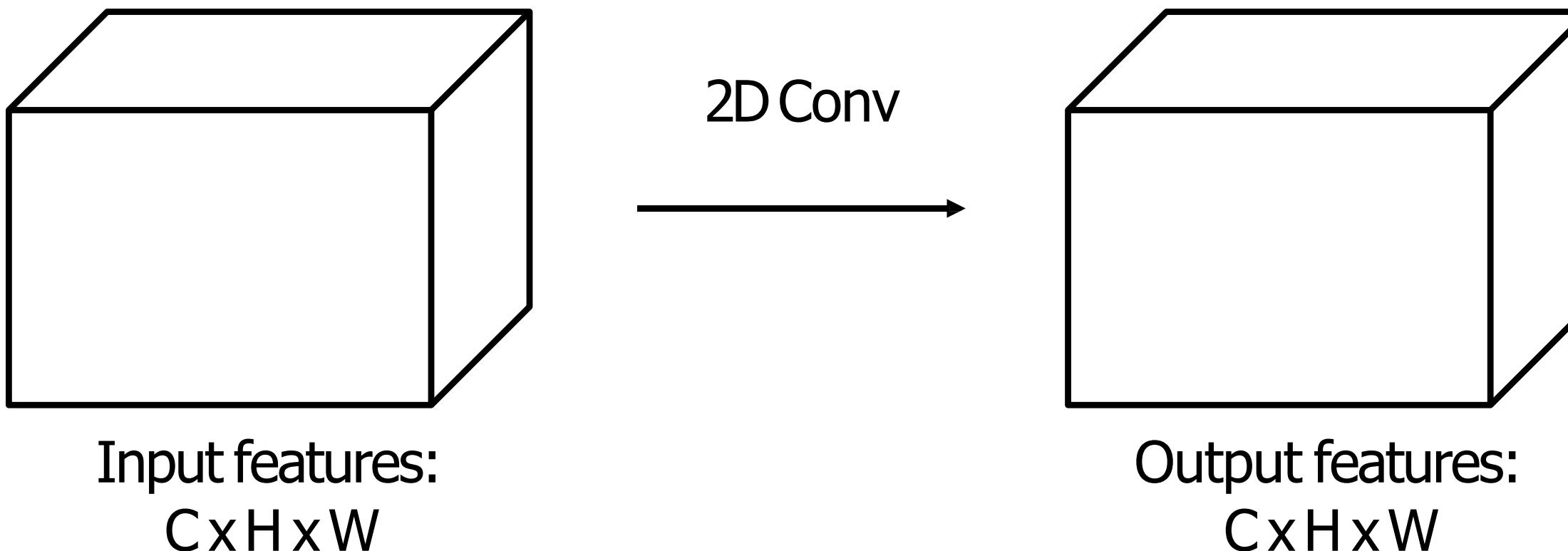
1. Same layer, previous timestep
2. Prev layer, same timestep

Use different weights at each layer, share weights across time

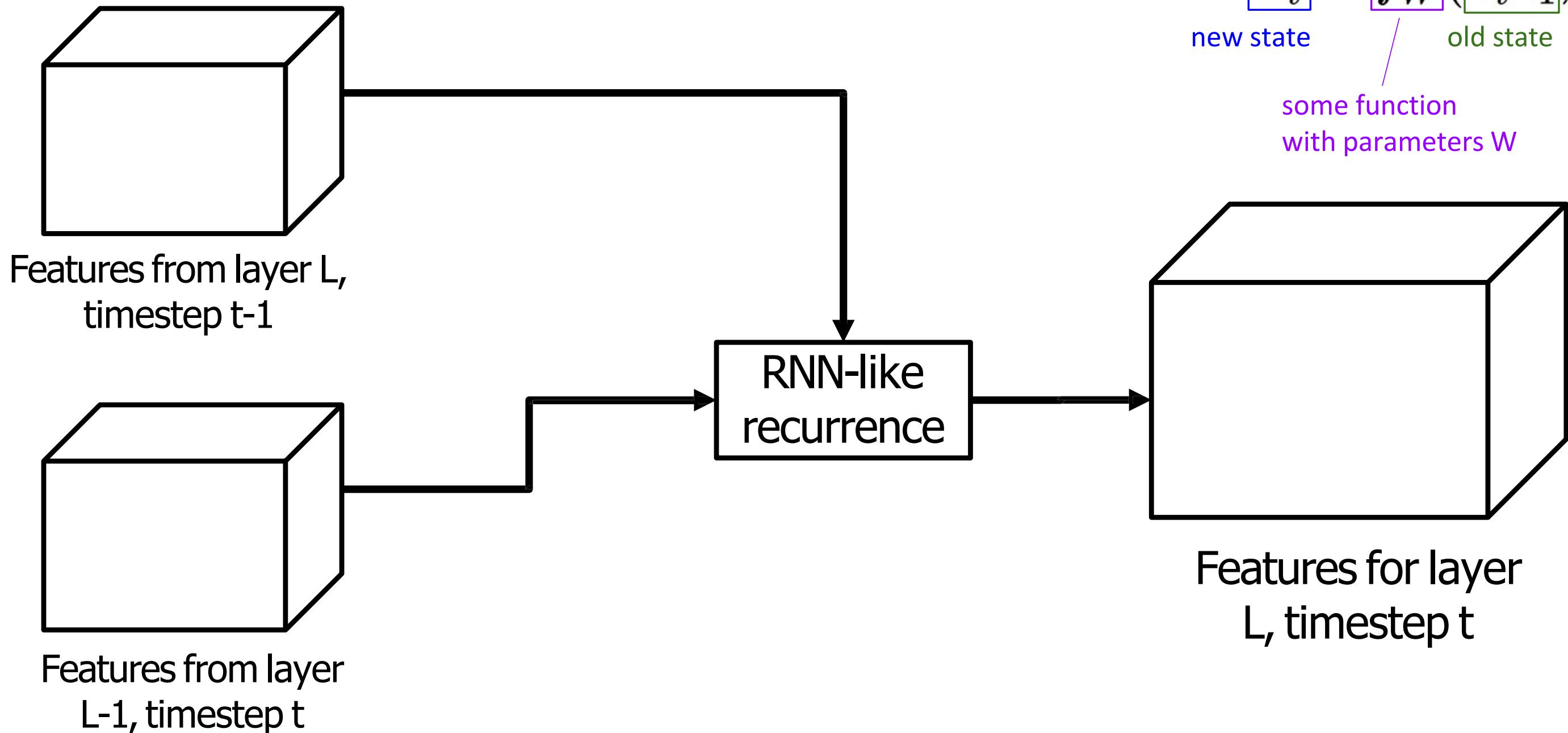
Ballas et al, "Delving Deeper into Convolutional Networks for Learning Video Representations", ICLR 2016

# Recurrent Convolutional Network

Normal 2D CNN:



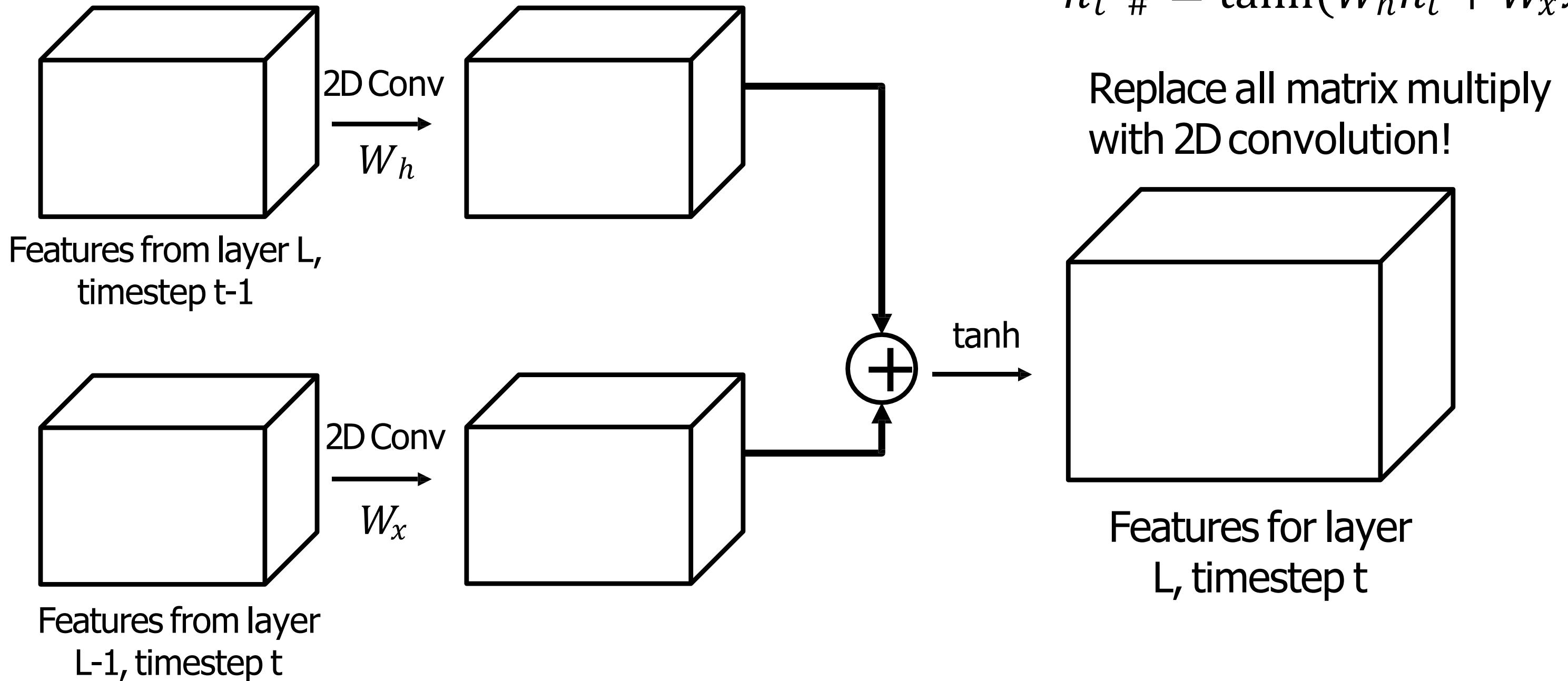
# Recurrent Convolutional Network



# Recurrent Convolutional Network

Recall: Vanilla RNN

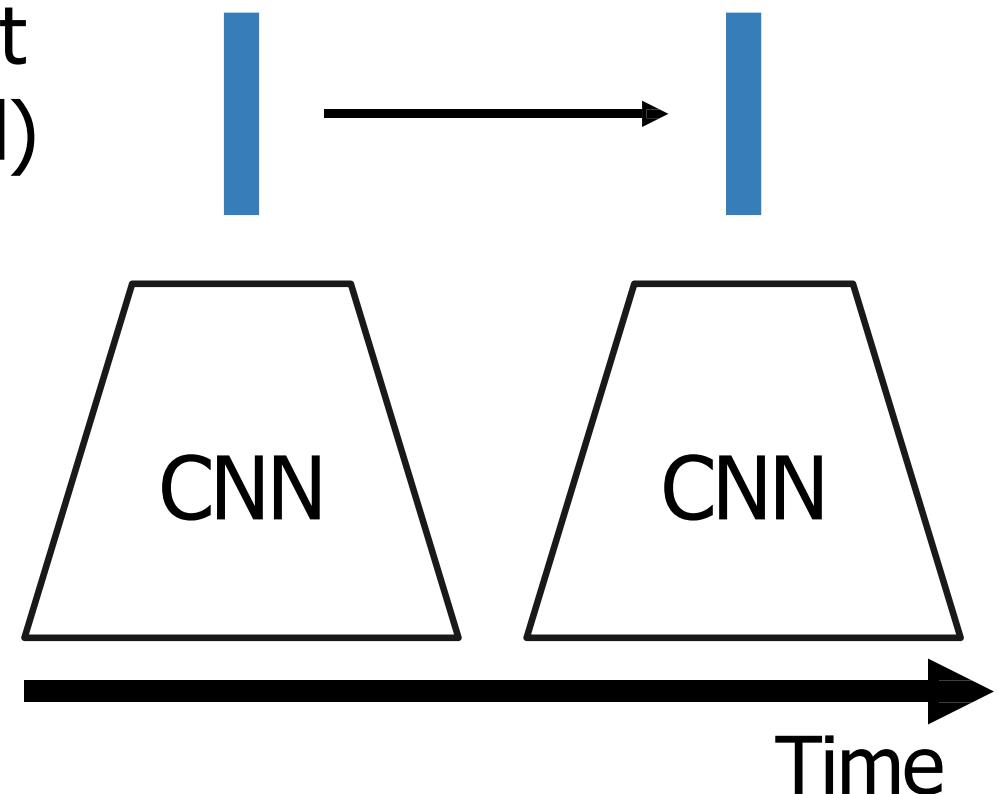
$$h_t'' \# = \tanh(W_h h_t + W_x x)$$



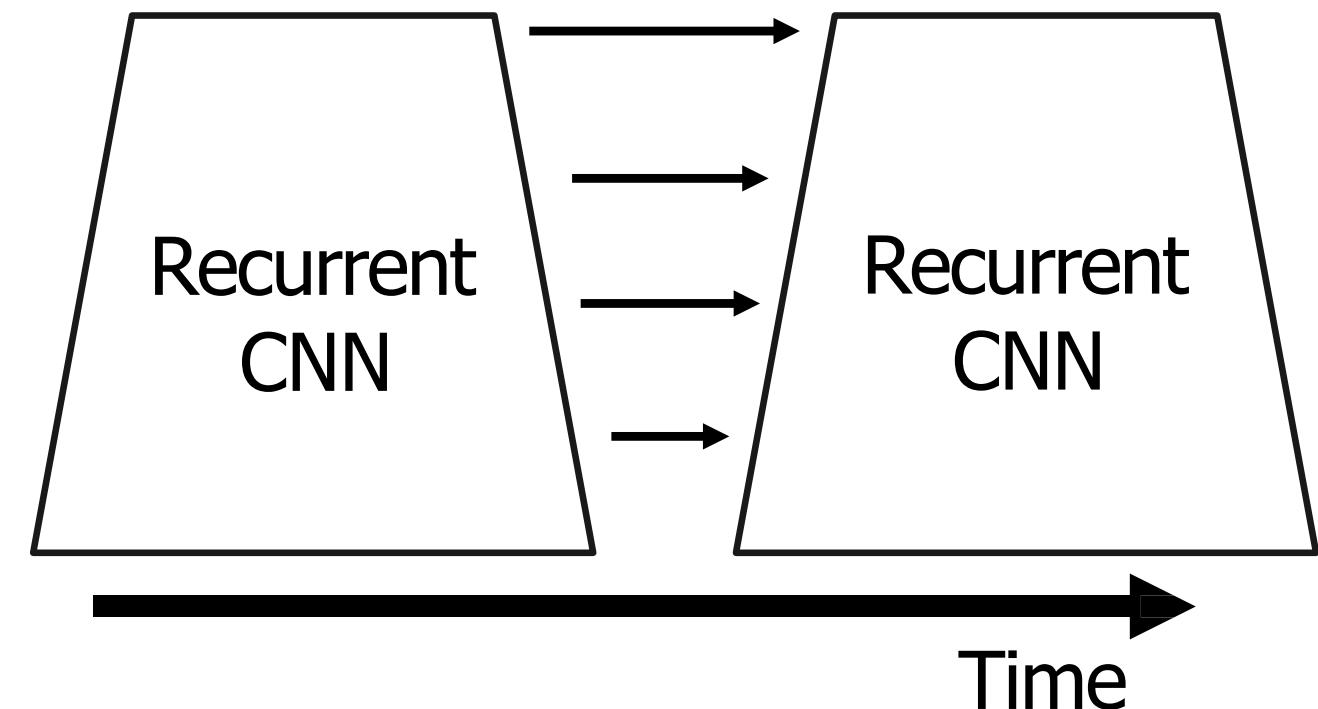
# Modeling long-term temporal structure

RNN: Infinite  
temporal extent  
(fully-connected)

CNN: finite  
temporal extent  
(convolutional)



Recurrent CNN: Infinite  
temporal extent  
(convolutional)



Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011

Donahue et al, "Long-term recurrent convolutional networks for visual recognition and description", CVPR 2015

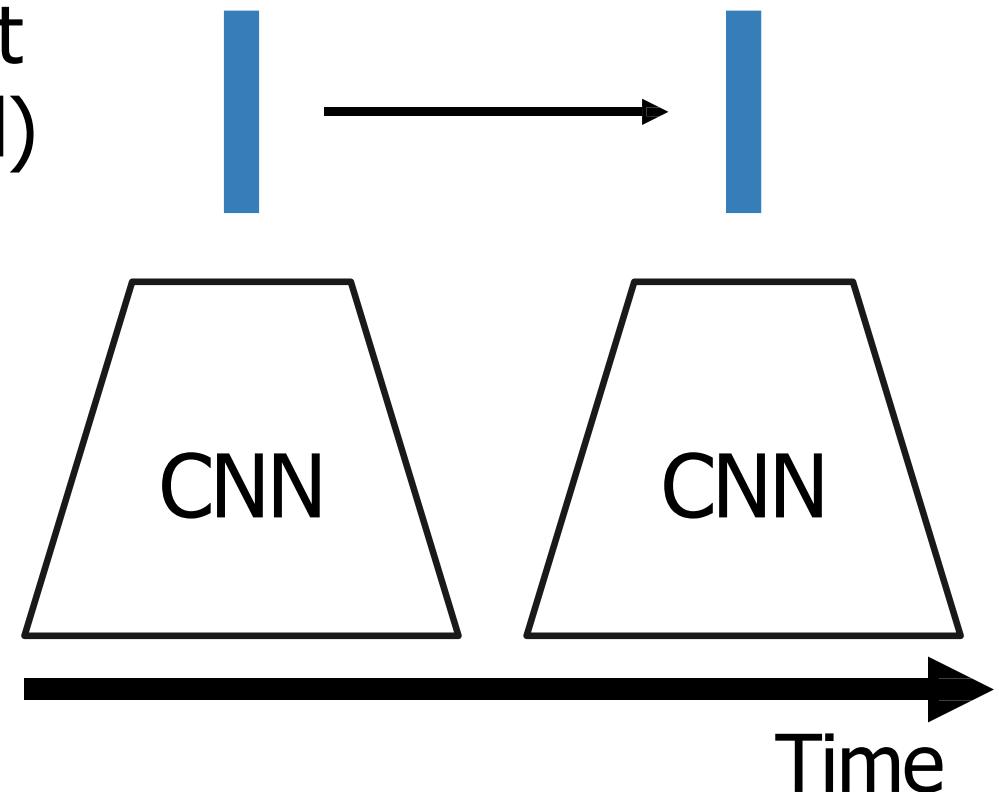
Ballas et al, "Delving Deeper into Convolutional Networks for Learning Video Representations", ICLR 2016

# Modeling long-term temporal structure

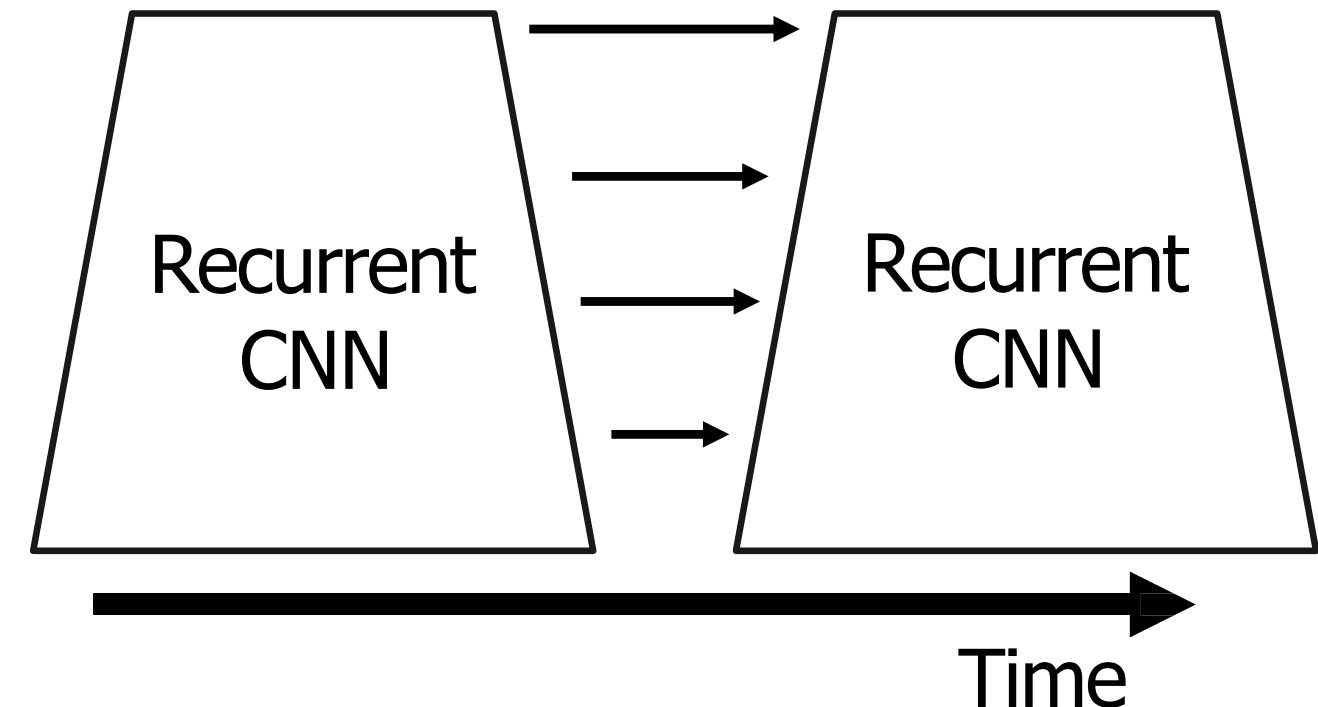
Problem: RNNs are slow for long sequences (can't be parallelized)

RNN: Infinite temporal extent (fully-connected)

CNN: finite temporal extent (convolutional)



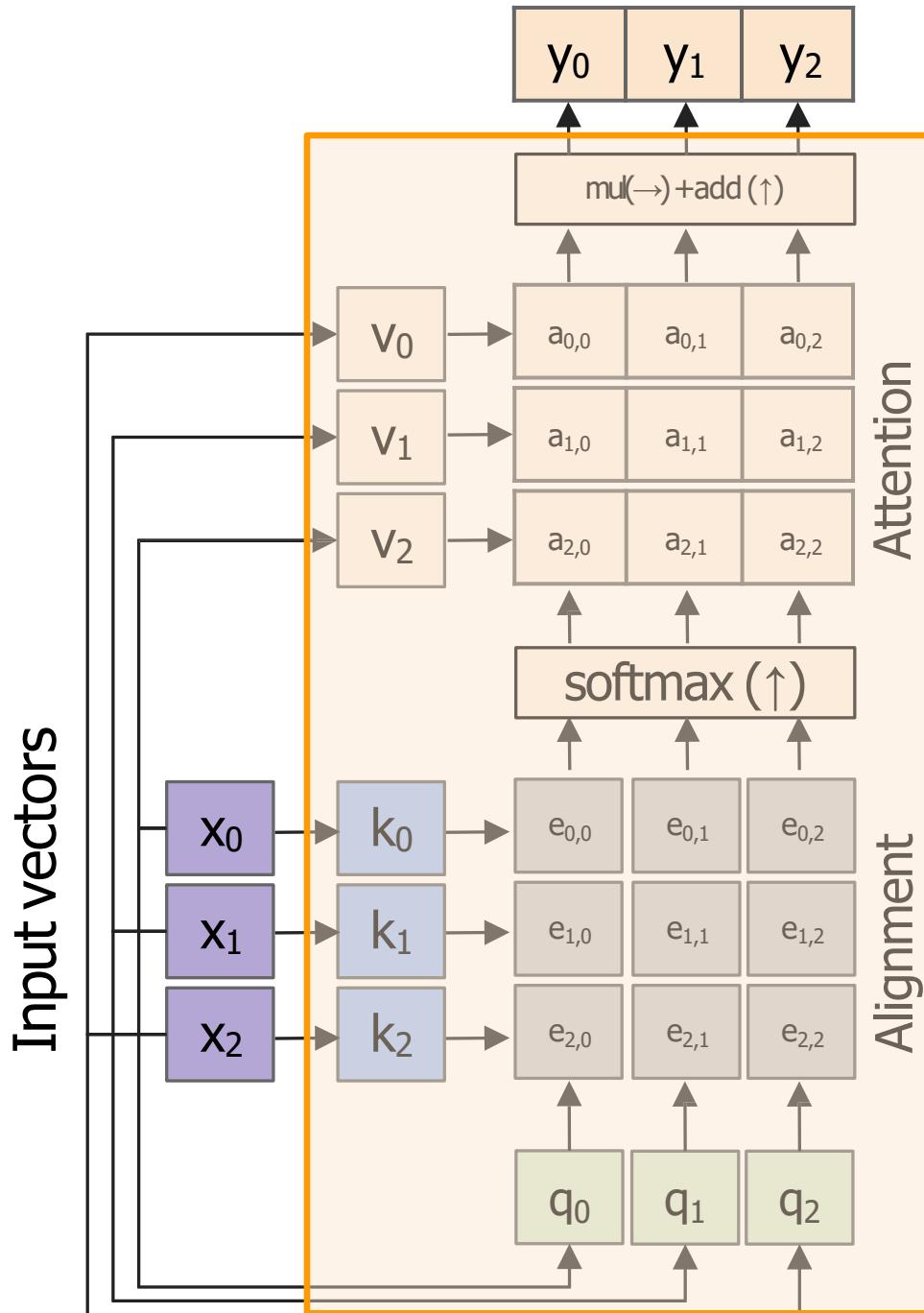
Recurrent CNN: Infinite temporal extent (convolutional)



Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011  
Donahue et al, "Long-term recurrent convolutional networks for visual recognition and description", CVPR 2015

Ballas et al, "Delving Deeper into Convolutional Networks for Learning Video Representations", ICLR 2016

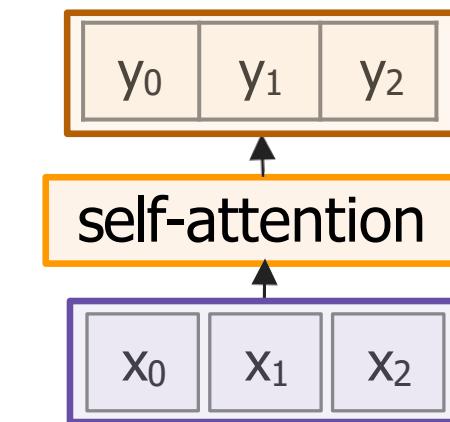
# Recall: Self-Attention



Outputs:  
context vectors:  $y$  (shape:  $D_v$ )

Operations:  
Key vectors:  $k = xW_k$   
Value vectors:  $v = xW_v$   
Query vectors:  $q = xW_q$   
Alignment:  $e_{i,j} = q_j \cdot k_i / \sqrt{D}$   
Attention:  $a = \text{softmax}(e)$   
Output:  $y_j = \sum a_{i,j} v_i$

Inputs:  
Input vectors:  $x$  (shape:  $N \times D$ )

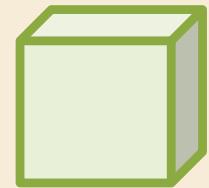


# Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



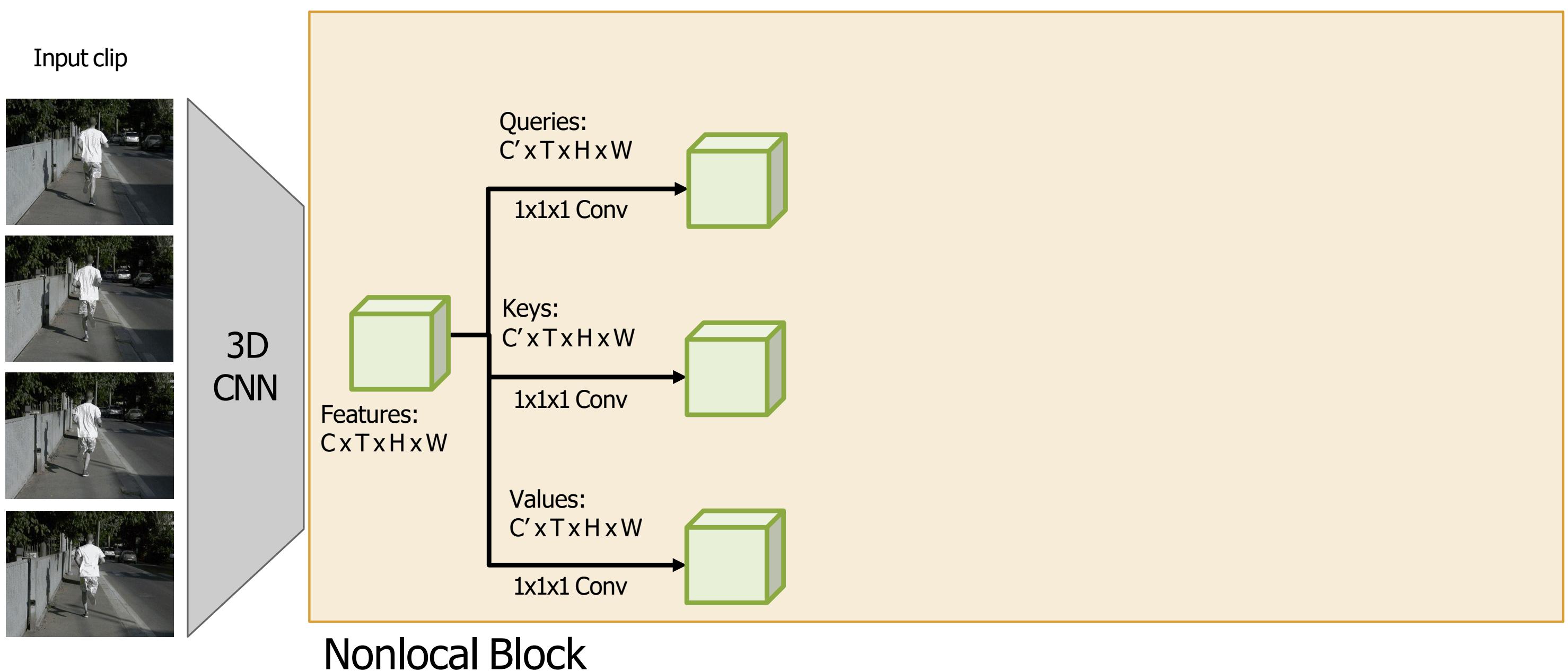
3D  
CNN



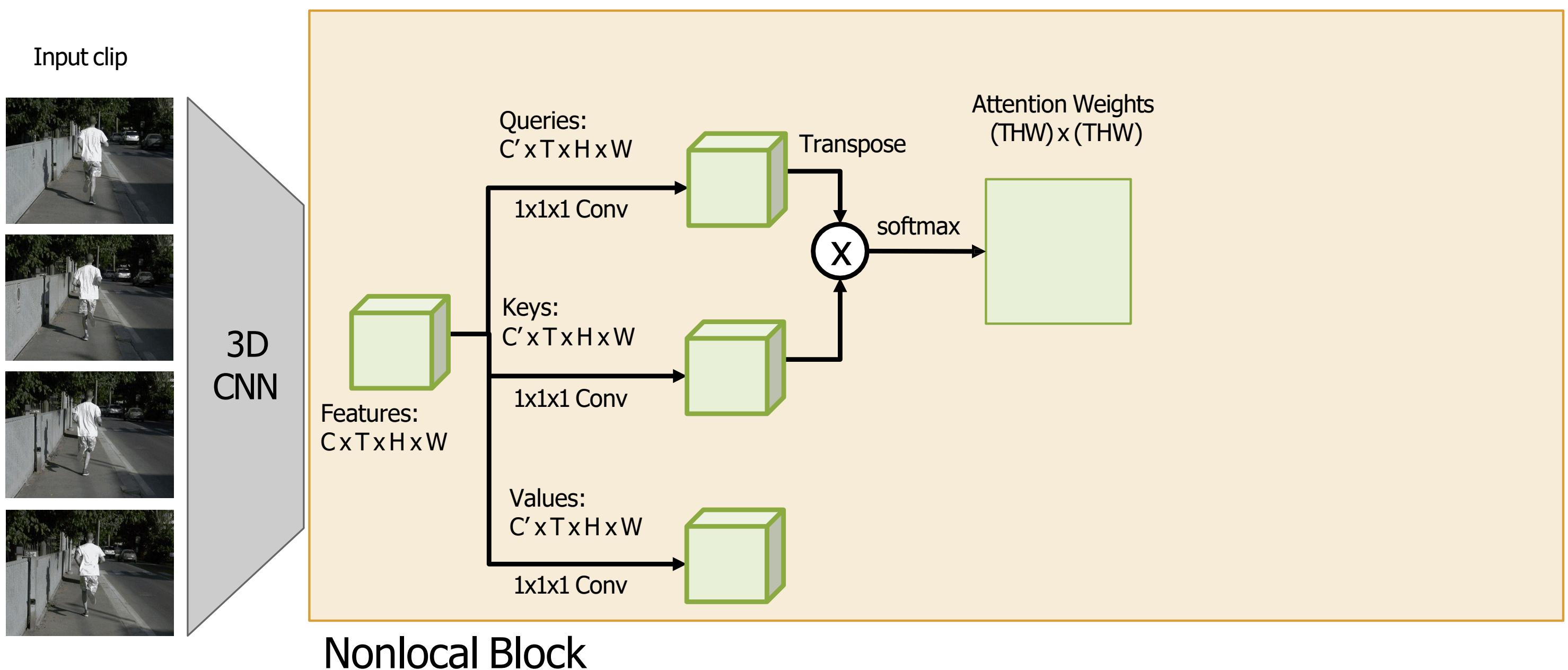
Features:  
 $C \times T \times H \times W$

Nonlocal Block

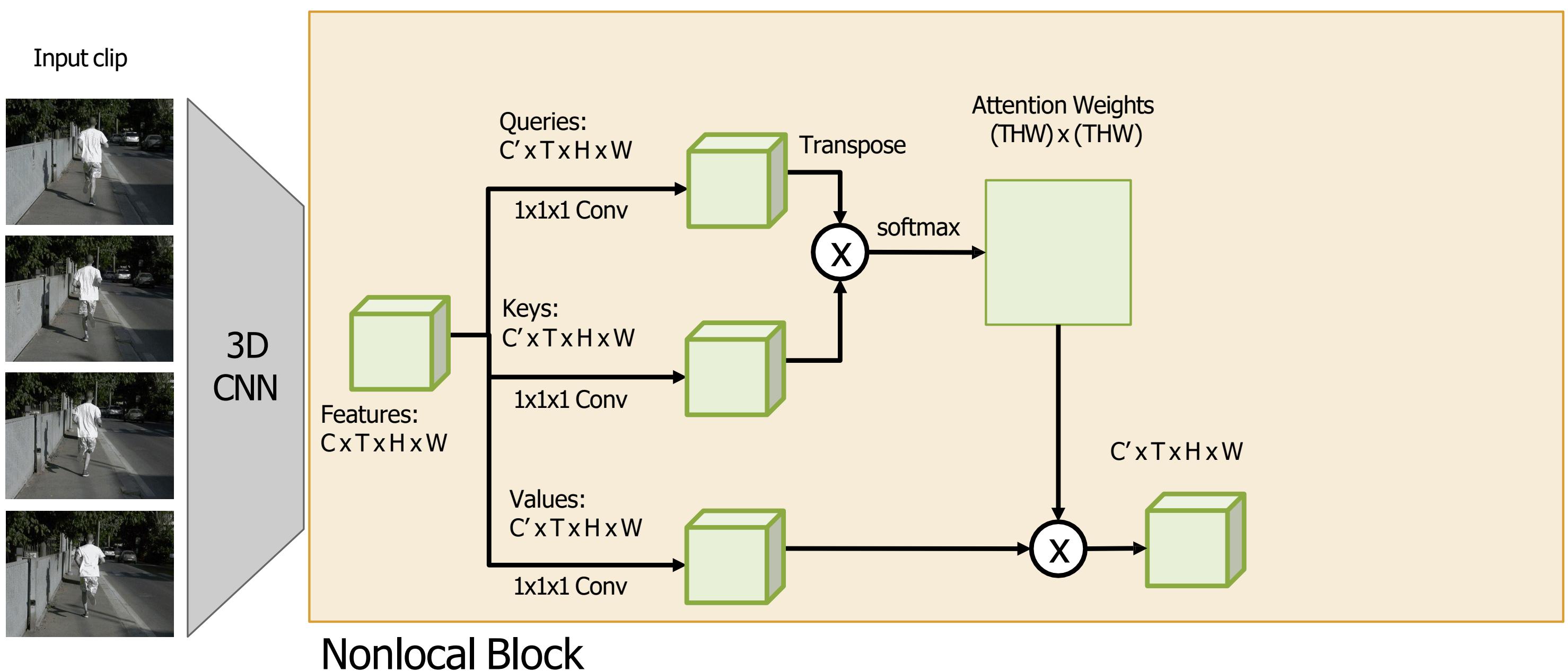
# Spatio-Temporal Self-Attention (Nonlocal Block)



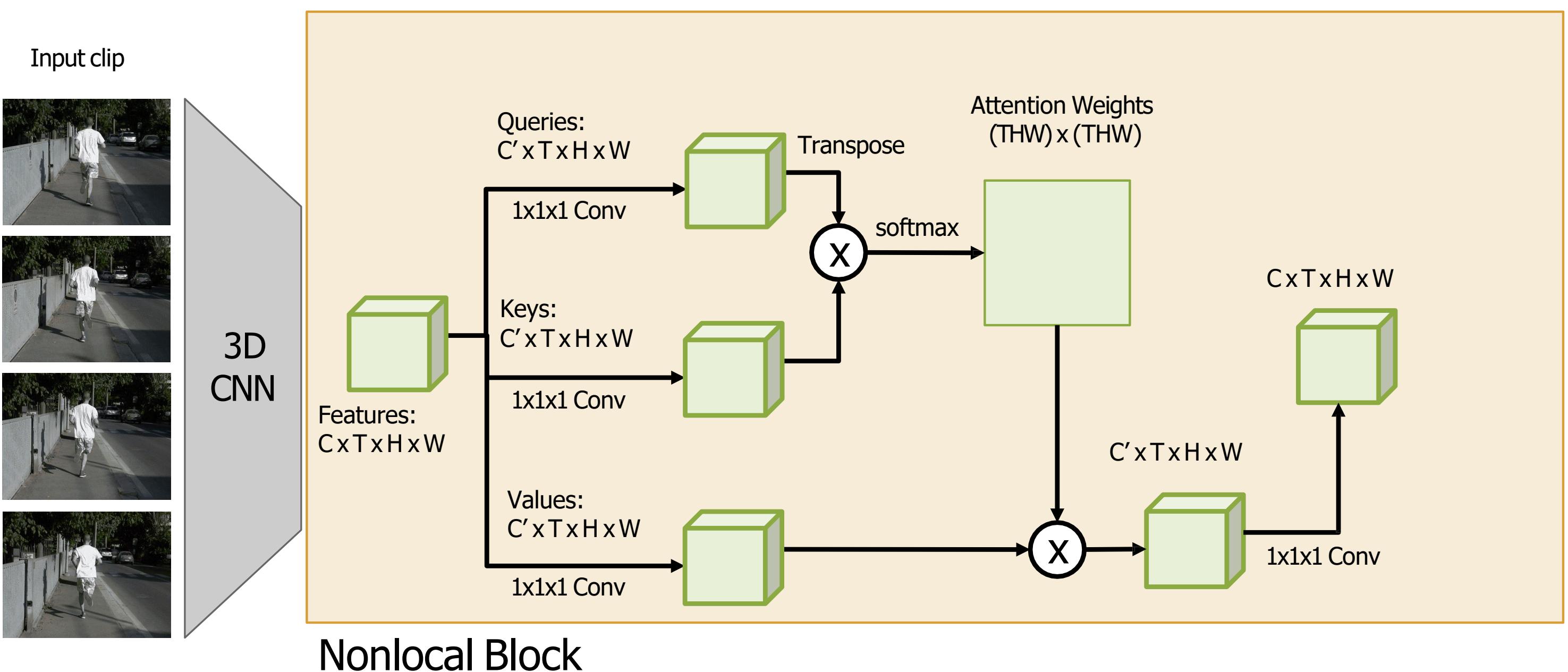
# Spatio-Temporal Self-Attention (Nonlocal Block)



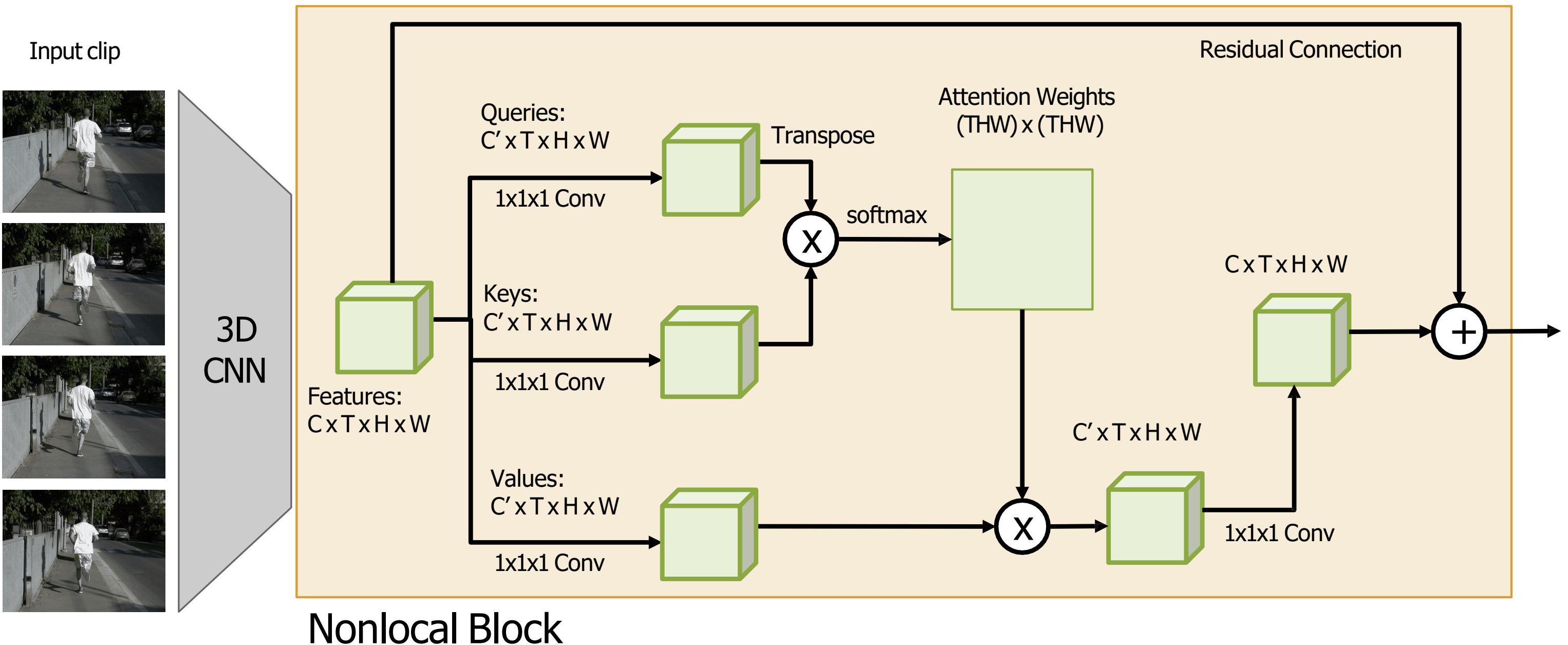
# Spatio-Temporal Self-Attention (Nonlocal Block)



# Spatio-Temporal Self-Attention (Nonlocal Block)

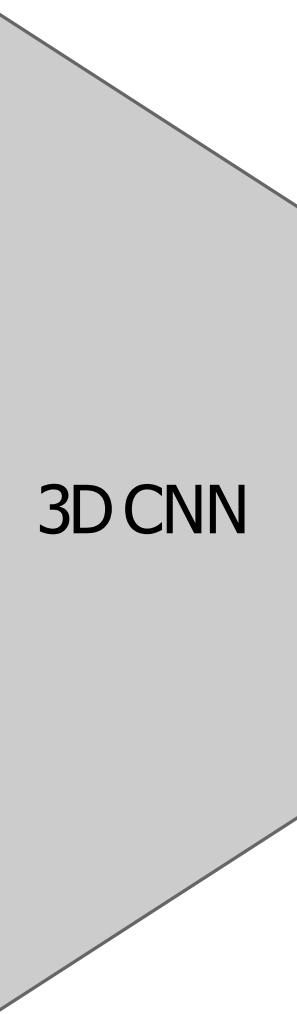


# Spatio-Temporal Self-Attention (Nonlocal Block)

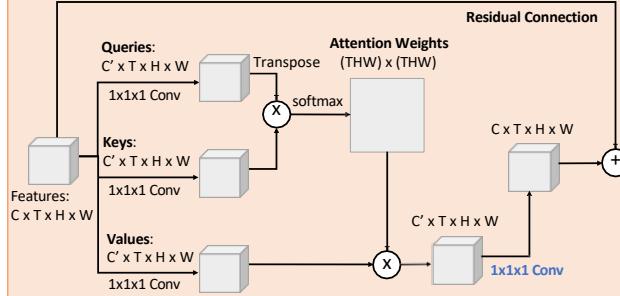


# Spatio-Temporal Self-Attention (Nonlocal Block)

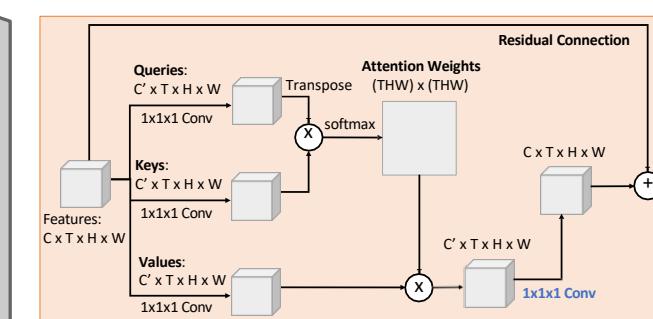
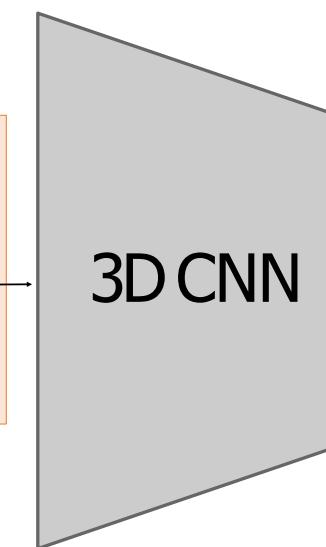
Input clip



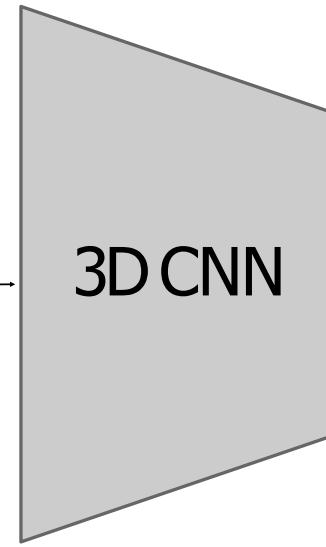
We can add nonlocal blocks into existing 3D CNN architectures.  
But what is the best 3D CNN architecture?



Nonlocal Block



Nonlocal Block



Running

# Inflating 2D Networks to 3D (I3D)

There has been a lot of work on architectures for images.  
Can we reuse image architectures for video?

Idea: take a 2D CNN architecture.

Replace each  $2D K_h \times K_w$  conv/pool  
layer with a  $3D K_t \times K_h \times K_w$  version

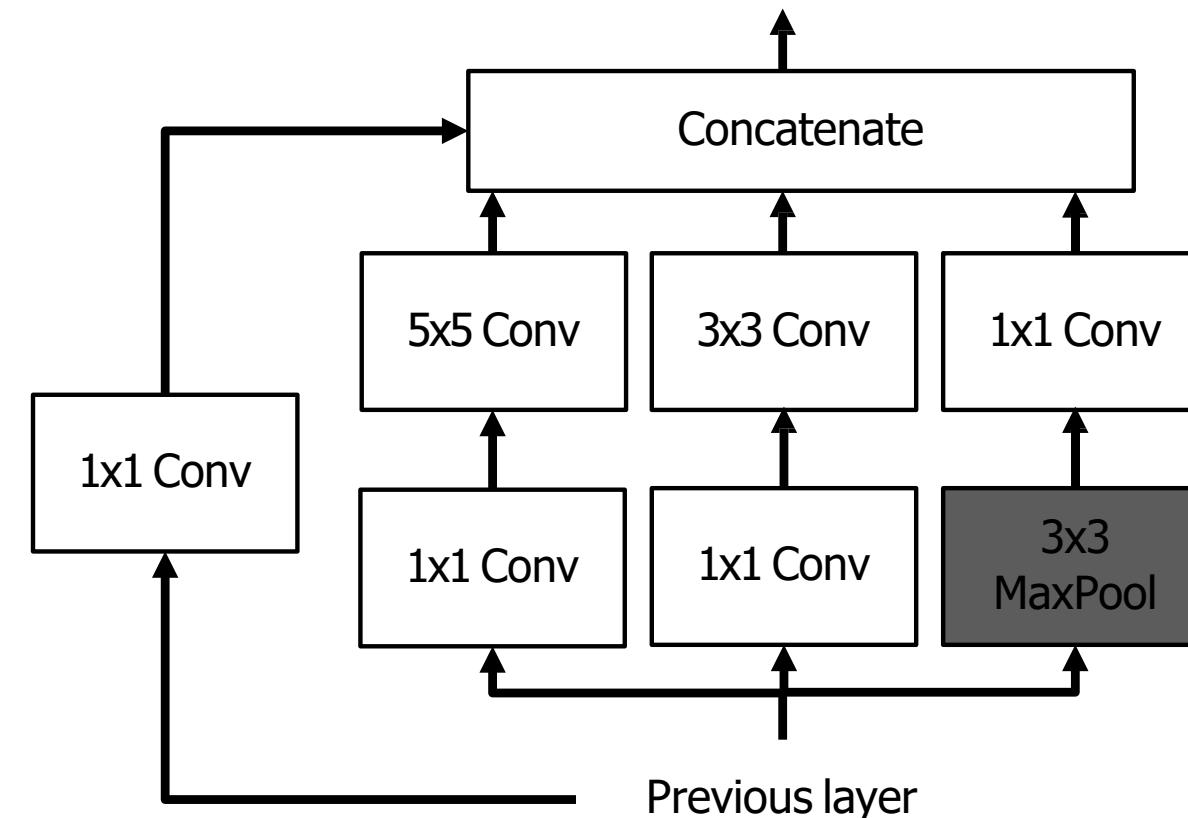
# Inflating 2D Networks to 3D (I3D)

There has been a lot of work on architectures for images.  
Can we reuse image architectures for video?

Idea: take a 2D CNN architecture.

Replace each  $2D K_h \times K_w$  conv/pool layer with a  $3D K_t \times K_h \times K_w$  version

Inception Block: Original



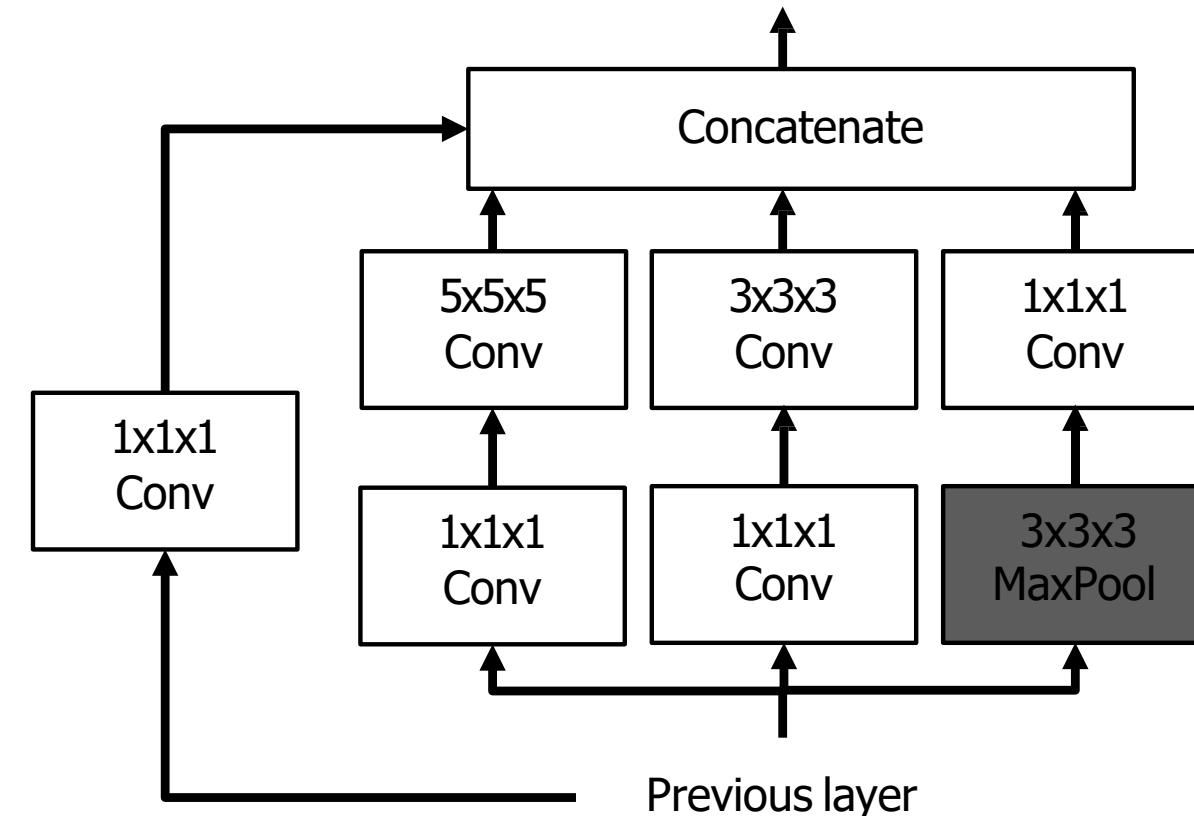
# Inflating 2D Networks to 3D (I3D)

There has been a lot of work on architectures for images.  
Can we reuse image architectures for video?

Idea: take a 2D CNN architecture.

Replace each  $2D K_h \times K_w$  conv/pool layer with a  $3D K_t \times K_h \times K_w$  version

Inception Block: Inflated



# Inflating 2D Networks to 3D (I3D)

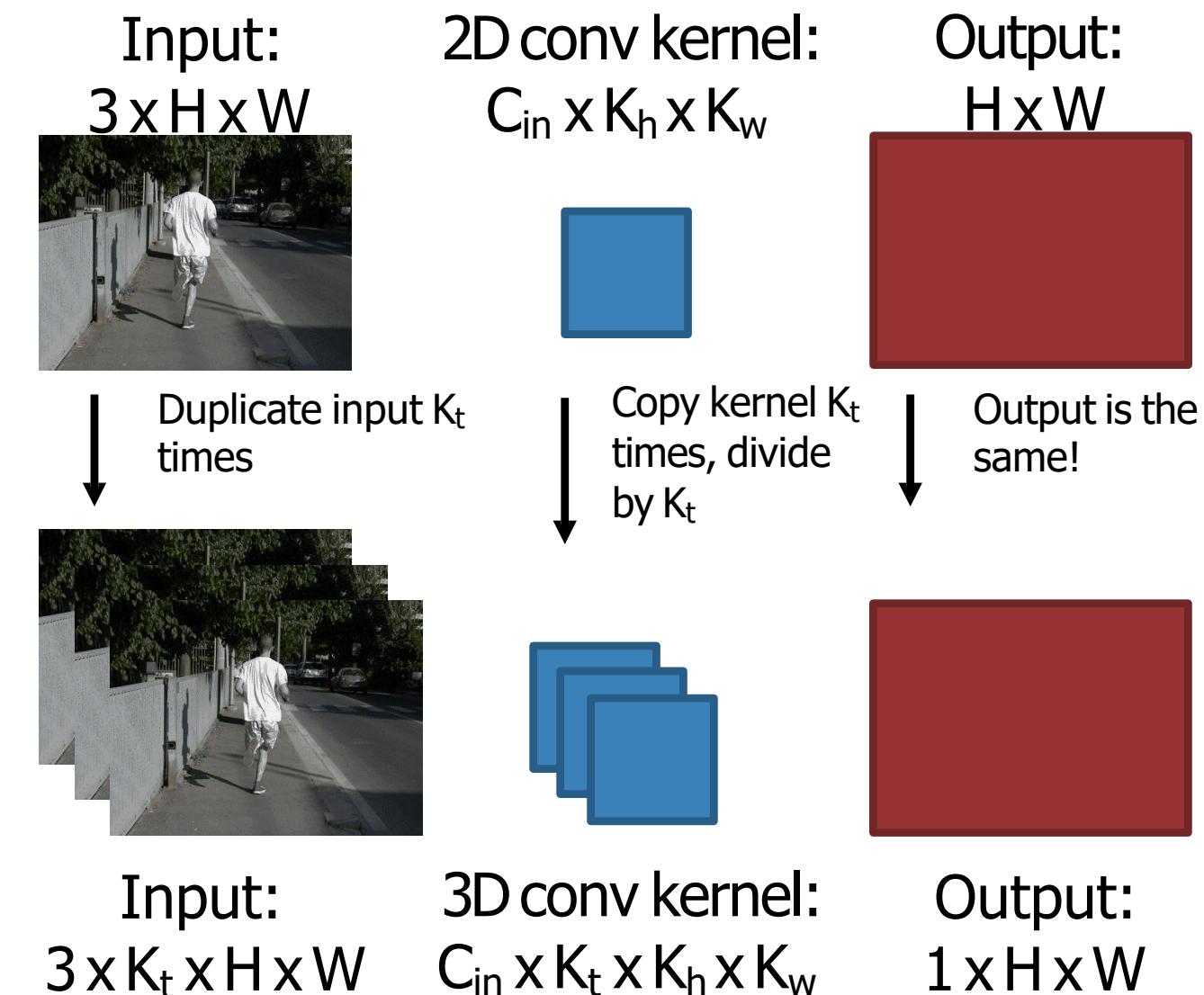
There has been a lot of work on architectures for images.  
Can we reuse image architectures for video?

Idea: take a 2D CNN architecture.

Replace each  $2D K_h \times K_w$  conv/pool layer with a  $3D K_t \times K_h \times K_w$  version

Can use weights of 2D conv to initialize 3D conv: copy  $K_t$  times in space and divide by  $K_t$

This gives the same result as 2D conv given “constant” video input



# Inflating 2D Networks to 3D (I3D)

There has been a lot of work on architectures for images.  
Can we reuse image architectures for video?

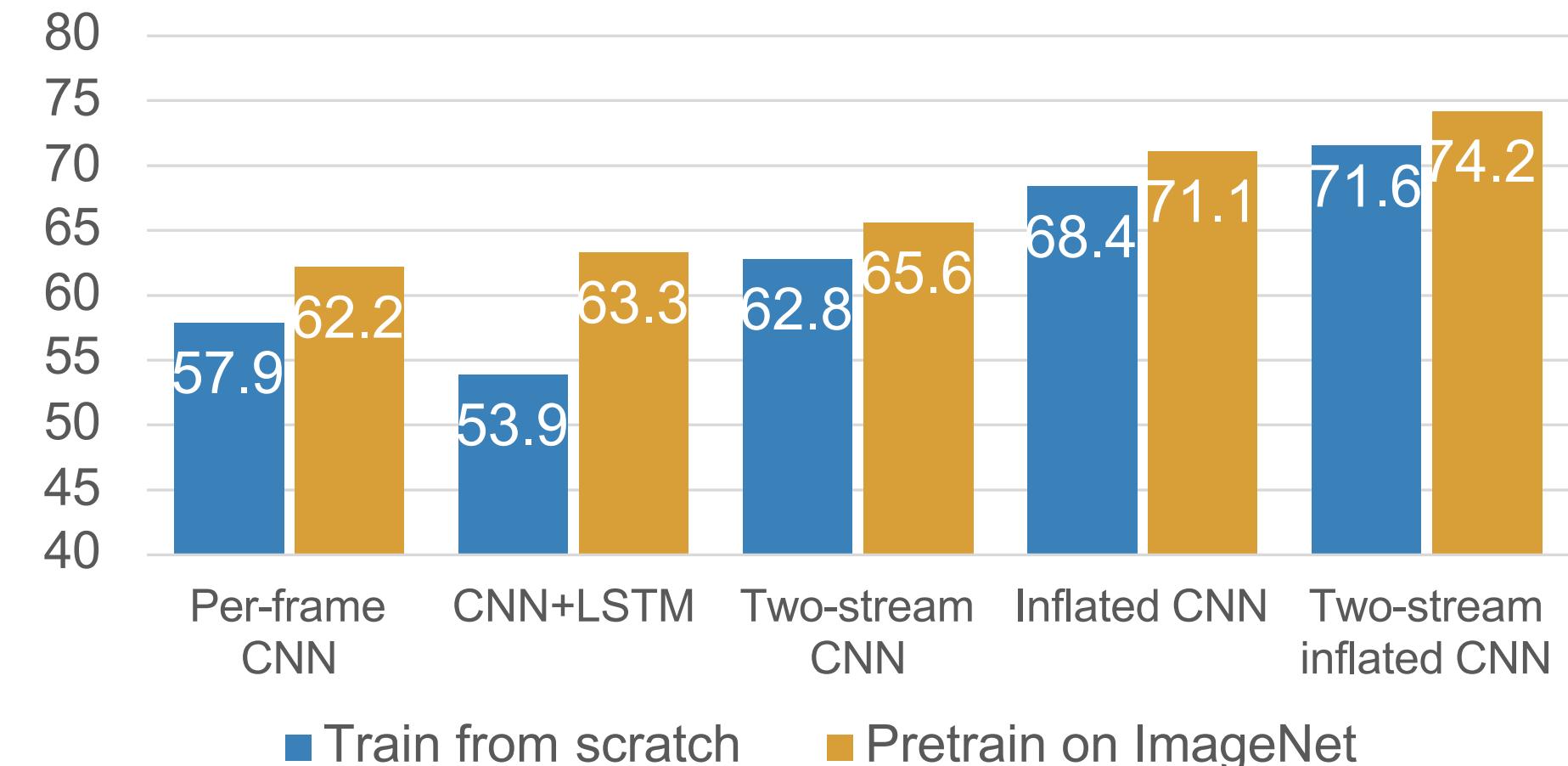
Idea: take a 2D CNN architecture.

Replace each  $2D K_h \times K_w$  conv/pool layer with a  $3D K_t \times K_h \times K_w$  version

Can use weights of 2D conv to initialize 3D conv: copy  $K_t$  times in space and divide by  $K_t$

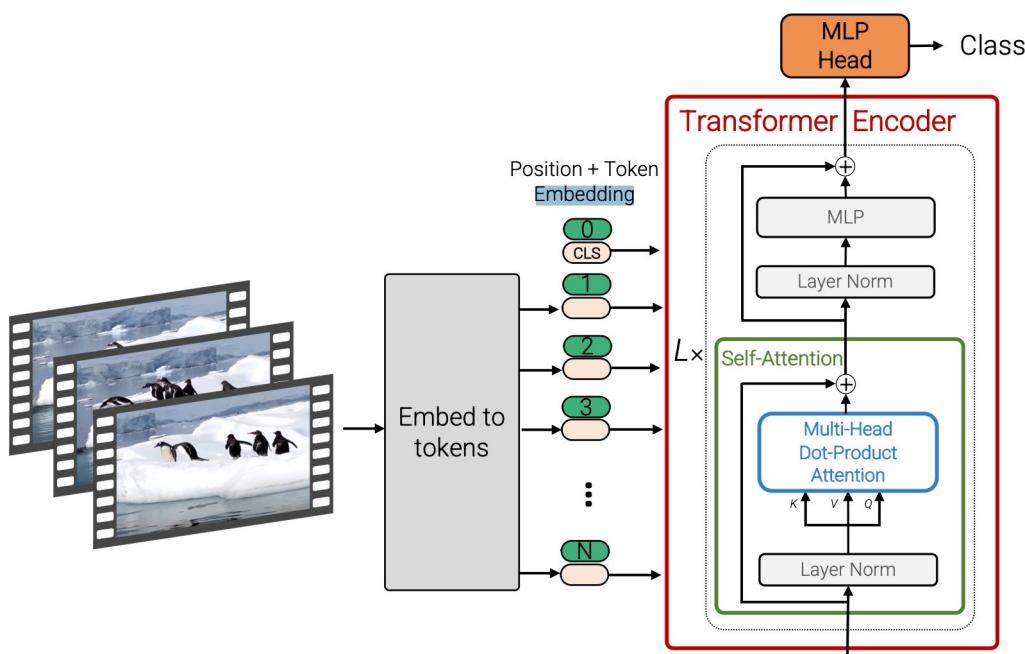
This gives the same result as 2D conv given “constant” video input

Top-1 Accuracy on Kinetics-400

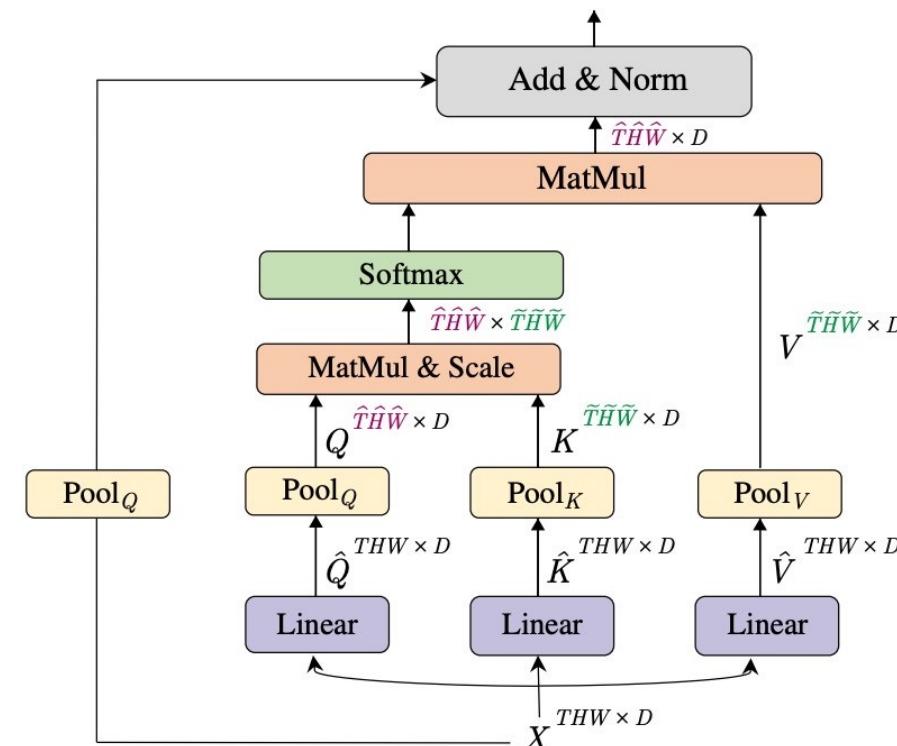


# Vision Transformers for Video

Factorized attention:  
Attend over space / time



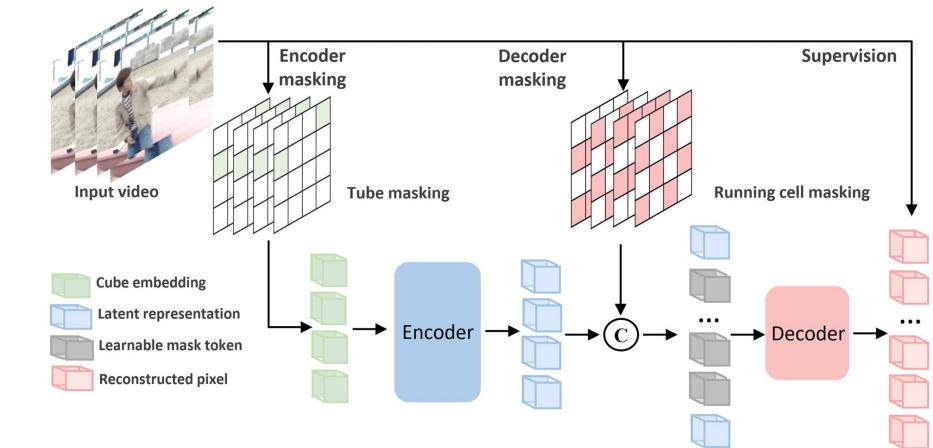
Pooling module:  
Reduce number of tokens



Bertasius et al, "Is Space-Time Attention All You Need for Video Understanding?", ICML 2021  
Arnab et al, "ViViT: A Video Vision Transformer", ICCV 2021  
 Neimark et al, "Video Transformer Network", ICCV 2021

Fan et al, "Multiscale Vision Transformers", ICCV 2021  
 Li et al, "MViTv2: Improved Multiscale Vision Transformers for Classification and Detection", CVPR 2022

Video masked autoencoders:  
Efficient scalable pretraining

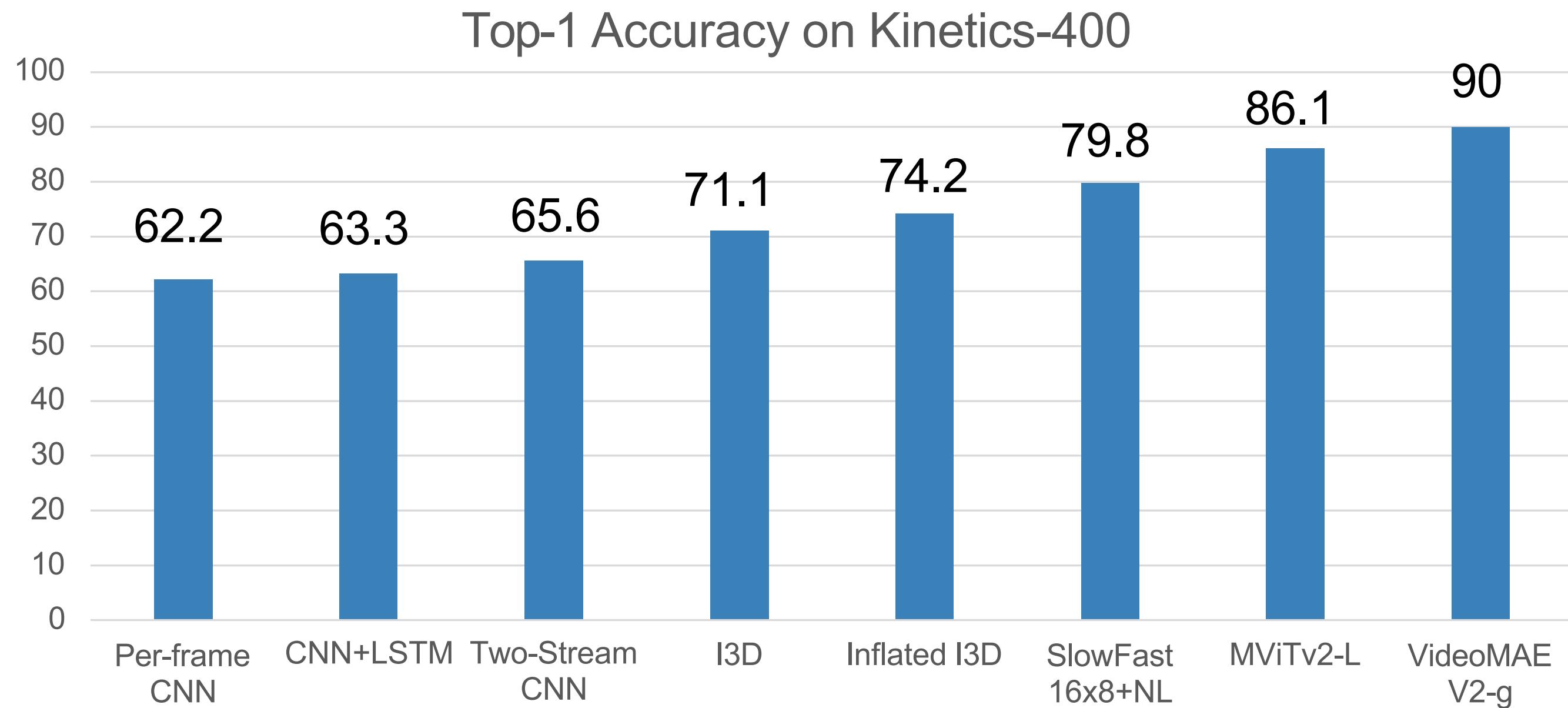


Wang et al. VideoMAE V2: Scaling Video Masked Autoencoders with Dual Making. CVPR 2023.

Tong et al. Video MAE: Masked Autoencoders are Data-Efficient Learners for Self-Supervised Video Pre-Training. NeurIPS 2022.

Feichtenhofer et al. Masked autoencoders as spatiotemporal learners. NeurIPS 2022.

# Vision Transformers for Video



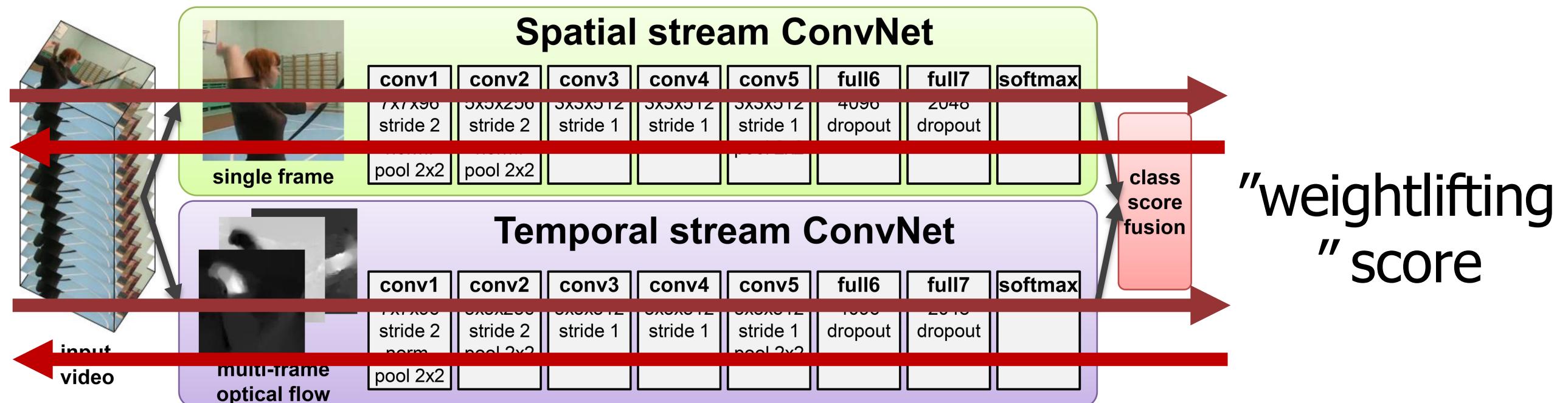
# Visualizing Video Models

Image



Flow

Forward: Compute class score



Backward: Compute gradient

Add a term to encourage spatially smooth flow; tune penalty to pick out “slow” vs “fast” motion

Figure credit: Simonyan and Zisserman, “Two-stream convolutional networks for action recognition in videos”, NeurIPS 2014  
Feichtenhofer et al, “What have we learned from deep representations for action recognition?”, CVPR 2018  
Feichtenhofer et al, “Deep insights into convolutional networks for video recognition?”, IJCV 2019.

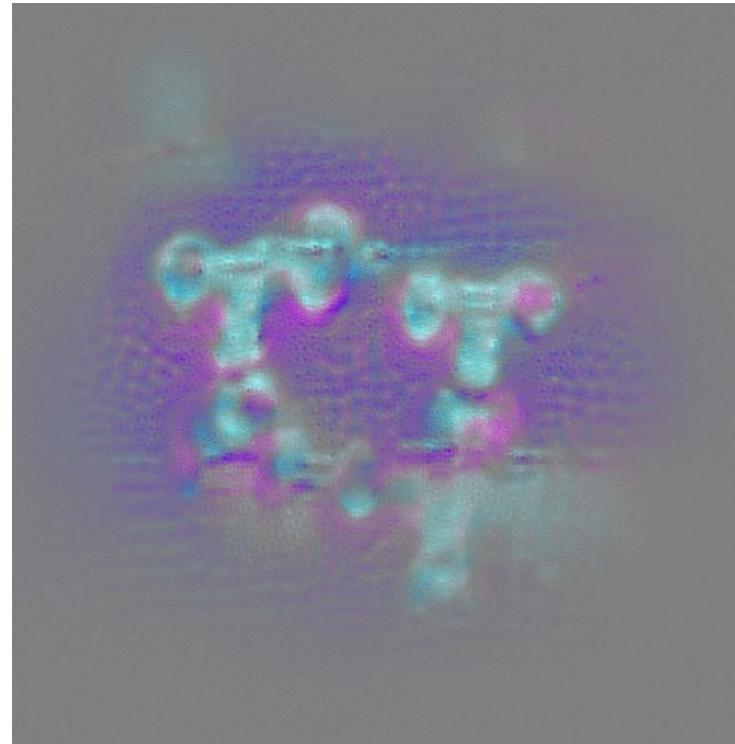
# Can you guess the action?

Feichtenhofer et al, "What have we learned from deep representations for action recognition?", CVPR 2018  
Feichtenhofer et al, "Deep insights into convolutional networks for video recognition?", IJCV 2019.  
Slide credit: Christoph Feichtenhofers

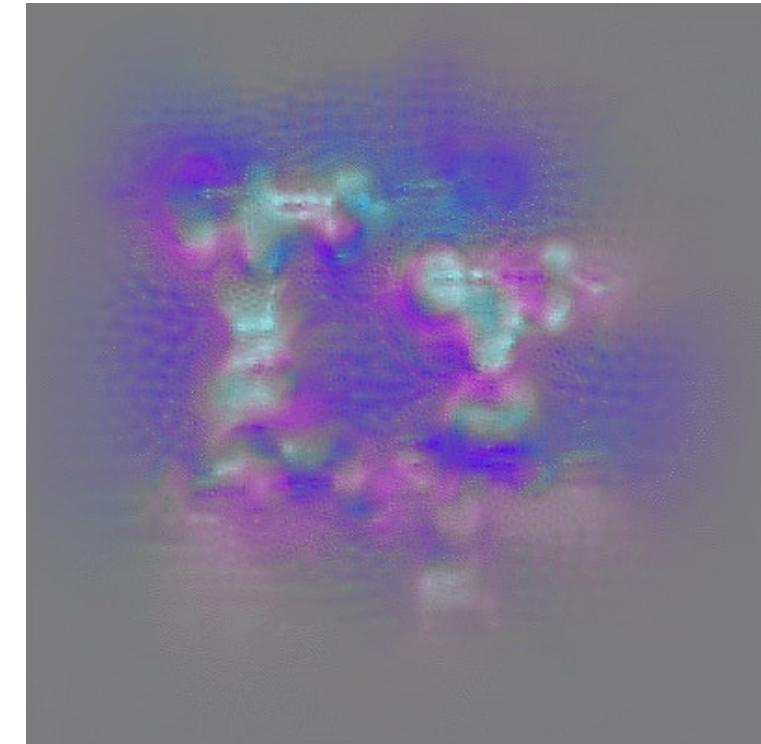
Appearance



"Slow" motion



"Fast" motion

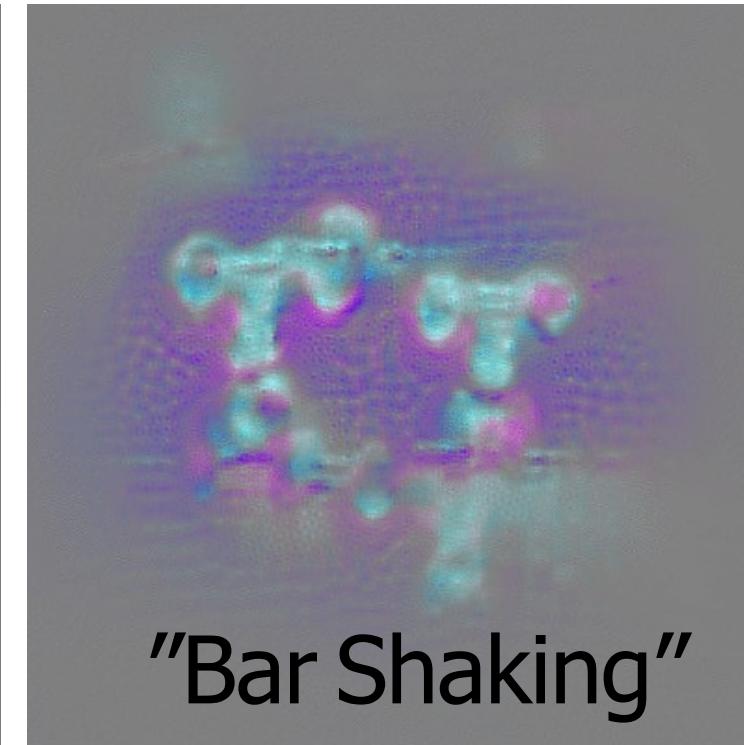


# Can you guess the action? Weightlifting

Appearance

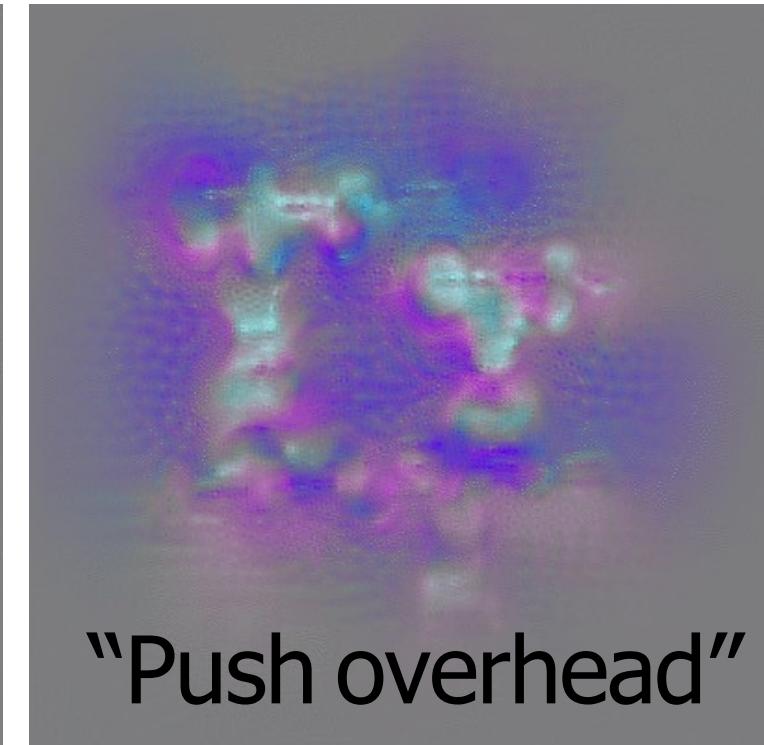


“Slow” motion



“Bar Shaking”

“Fast” motion



“Push overhead”



# So far: Classify short clips



Videos: Recognize actions



Swimming  
Running  
Jumping  
Eating  
Standing

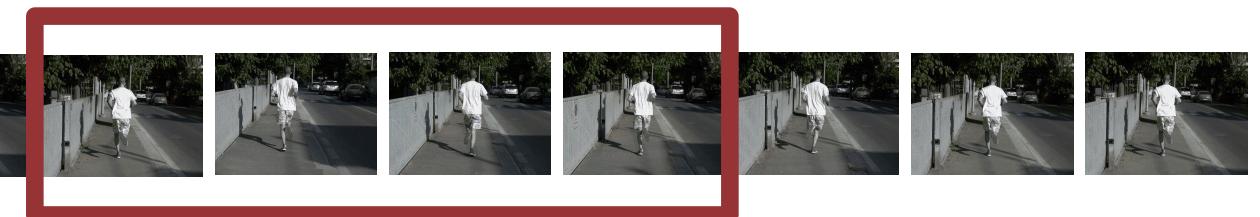
# Temporal Action Localization

Given a long untrimmed video sequence, identify frames corresponding to different actions

Running



Jumping

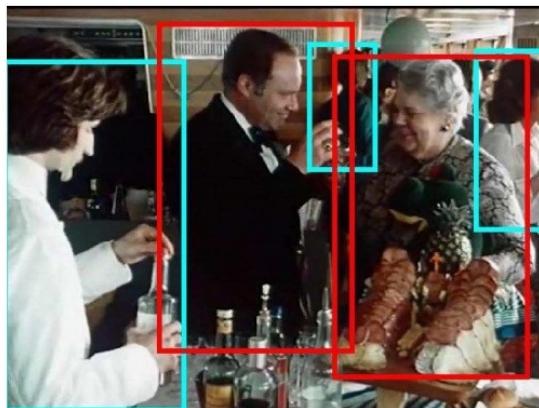


Can use architecture similar to Faster R-CNN: first generate temporal proposals then classify

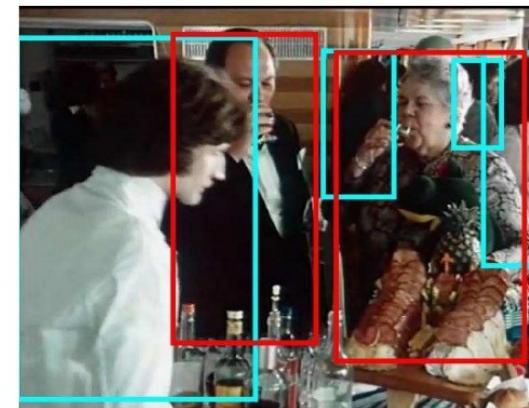
Chao et al, "Rethinking the Faster R-CNN Architecture for Temporal Action Localization", CVPR 2018

# Spatio-Temporal Detection

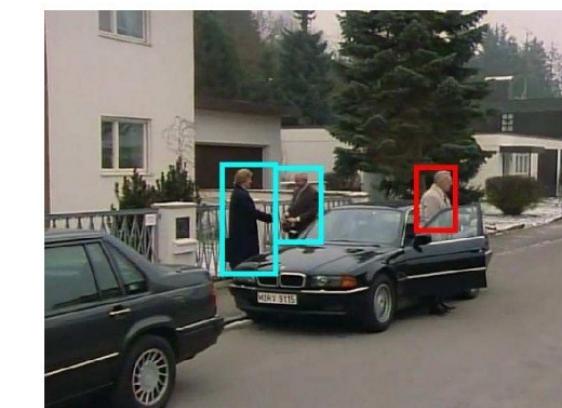
Given a long untrimmed video, detect all the people in both space and time and classify the activities they are performing.  
Some examples from AVA Dataset:



clink glass → drink



open → close

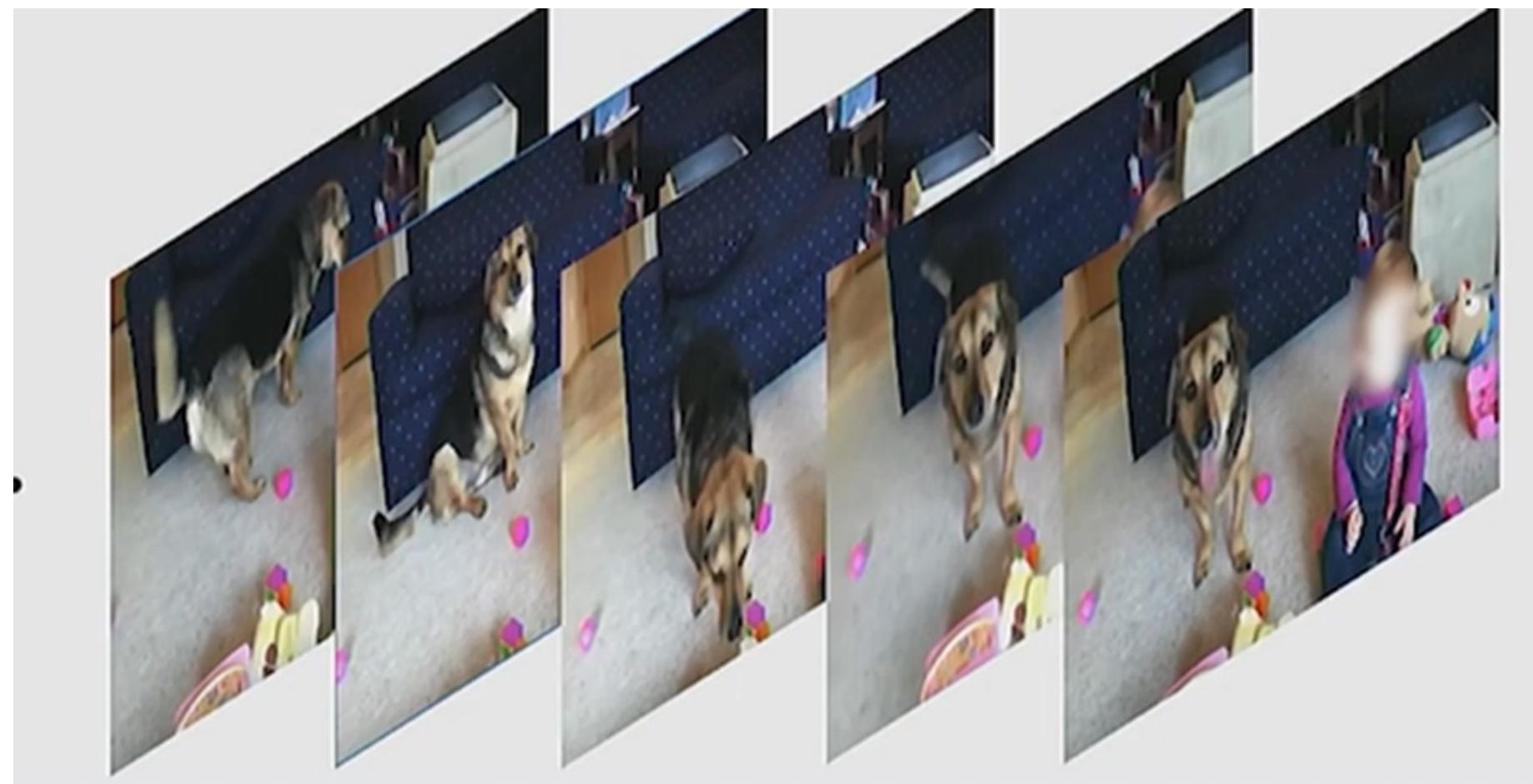


grab (a person) → hug



look at phone → answer phone

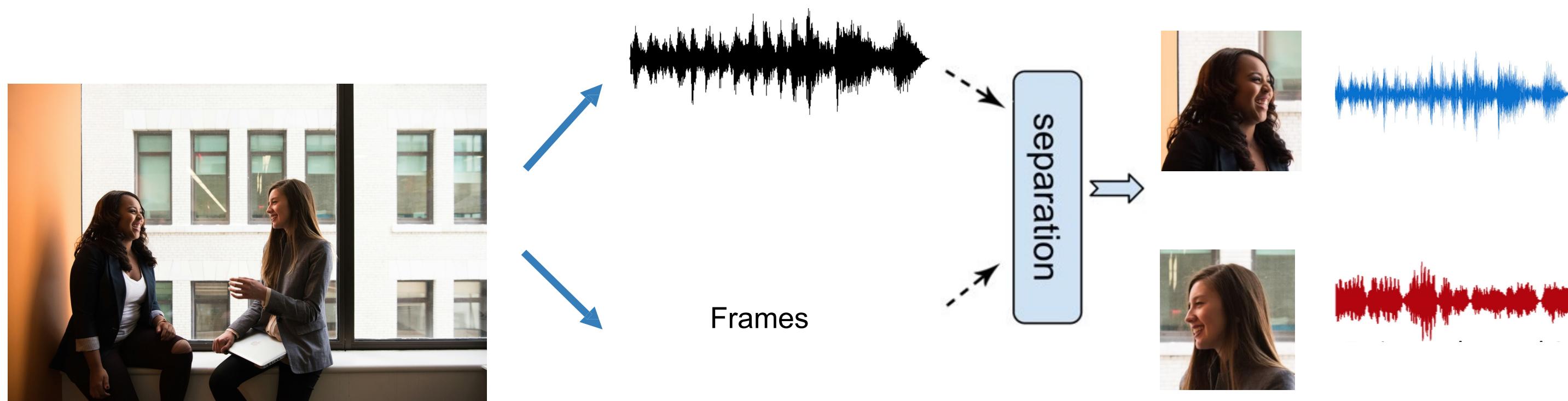
# Today: Temporal Stream



3D CNN, Two-Stream Neural Network, Spatial-Temporal Self-Attention.....



# Visually-guided audio source separation



[Gao et al. ECCV 2018, Afouras et al. Interspeech'18, Gabby et al. Interspeech'18, Owens & Efros ECCV'18, Ephrat et al. SIGGRAPH'18, Zhao et al. ECCV 2018, Gao & Grauman ICCV 2019, Zhao et al. ICCV 2019, Xu et al. ICCV 2019, Gan et al. CVPR 2020, Gao et al. CVPR 2021, Tzinis et al. ECCV 2022, Chen et al. CVPR 2023] Images: Public Domain- Unsplash

# Speech Mixture Video:

<https://www.youtube.com/watch?v=Kxb9AUmSrIA>

# Musical instruments source separation

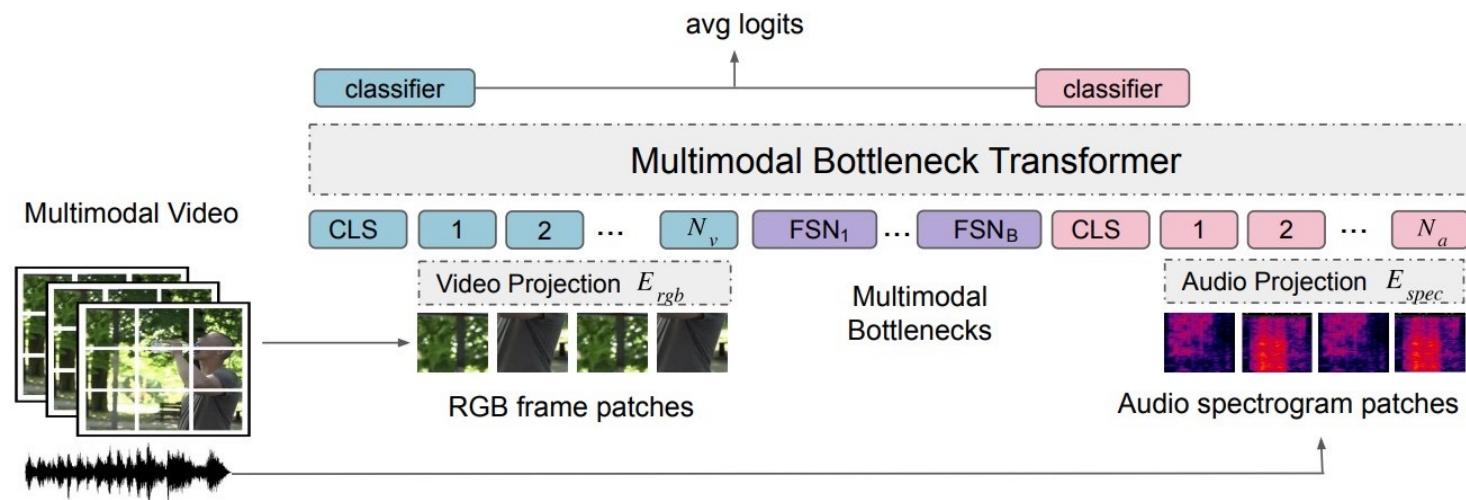
Train on 100,000 unlabeled multi-source video clips,  
then separate audio for novel video.



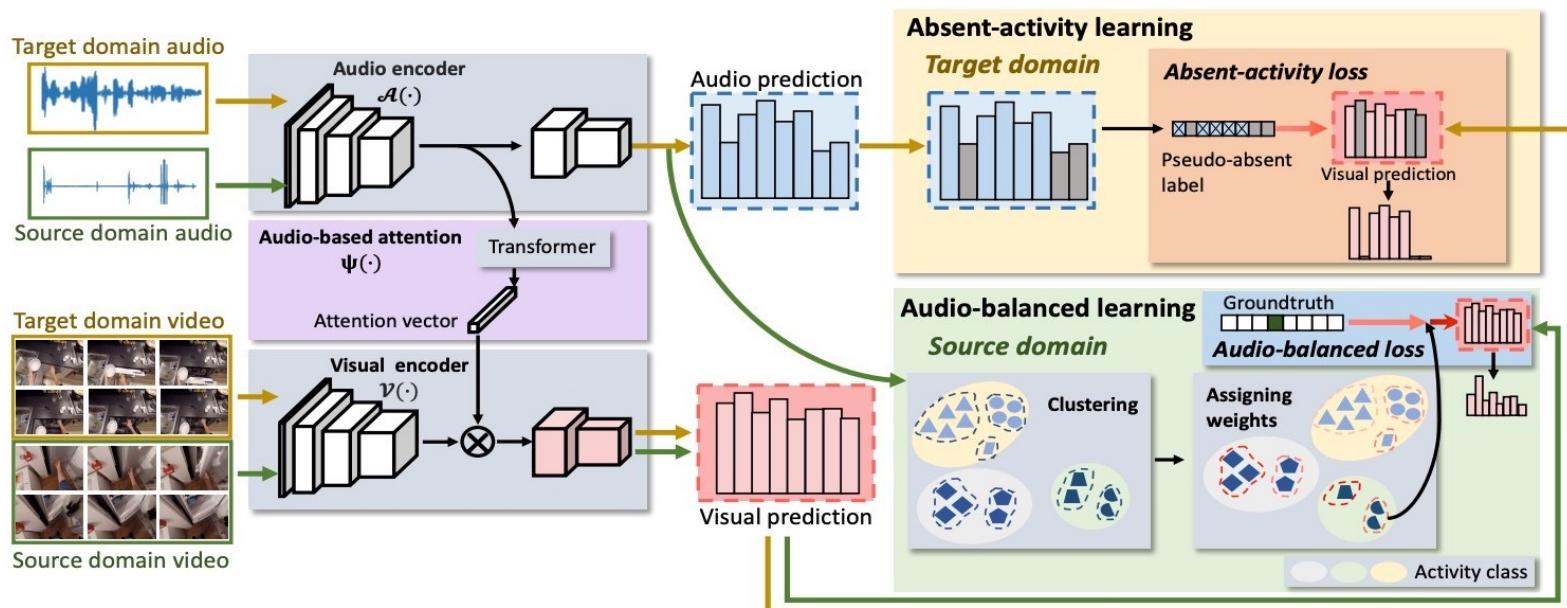
original video  
(before separation)

object detections:  
violin & flute

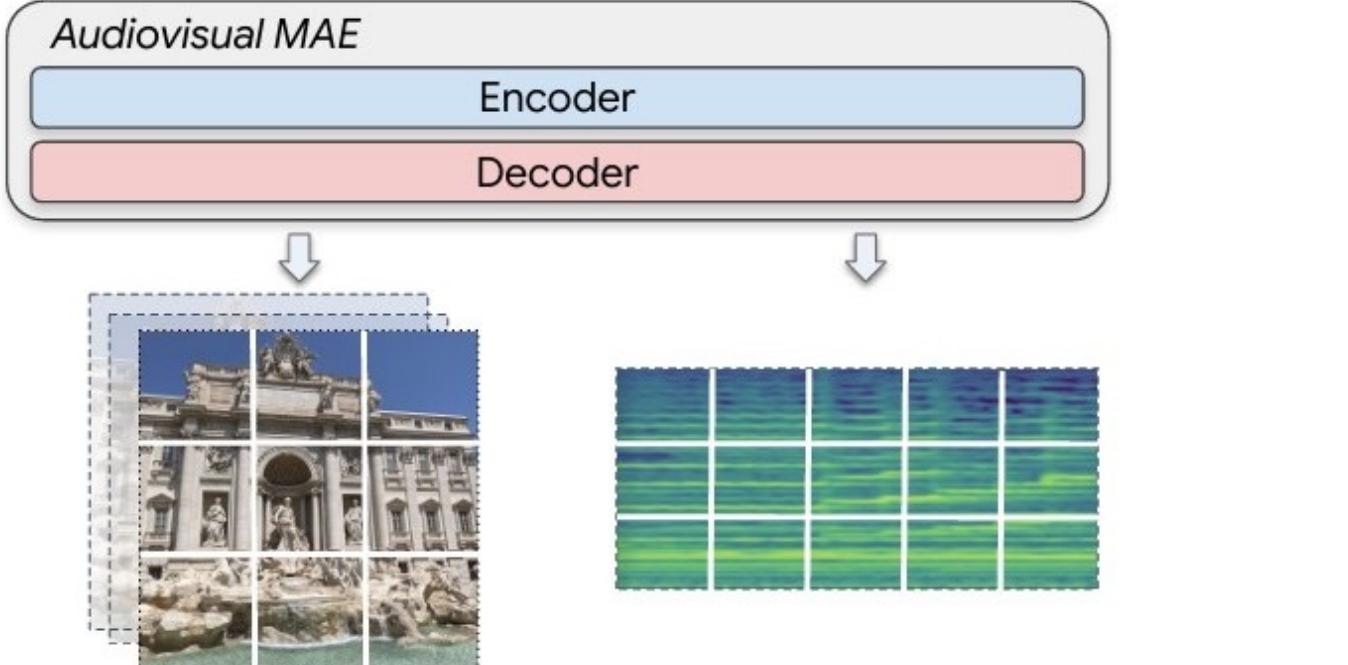
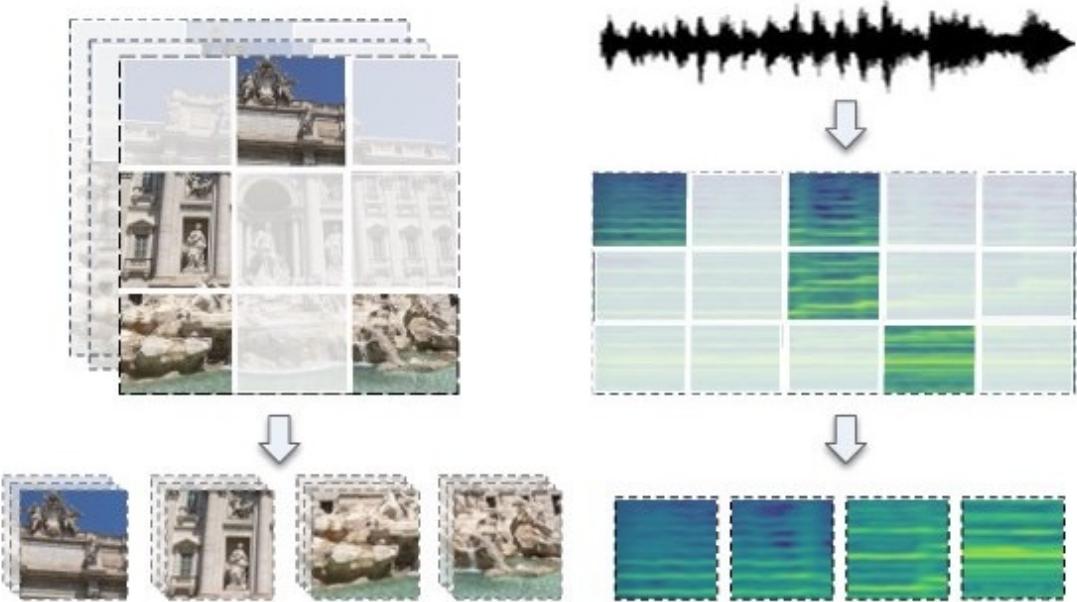
# Audio-Visual Video Understanding



Attention Bottlenecks for Multimodal Fusion, Nagrani et al. NeurIPS 2021



Audio-Adaptive Activity Recognition Across Video Domains,  
Zhang et al. CVPR 2022



Audio-Visual Masked Autoencoders. Georgescu et al. ICCV 2023.

# Efficient Video Understanding

Action recognition in long videos

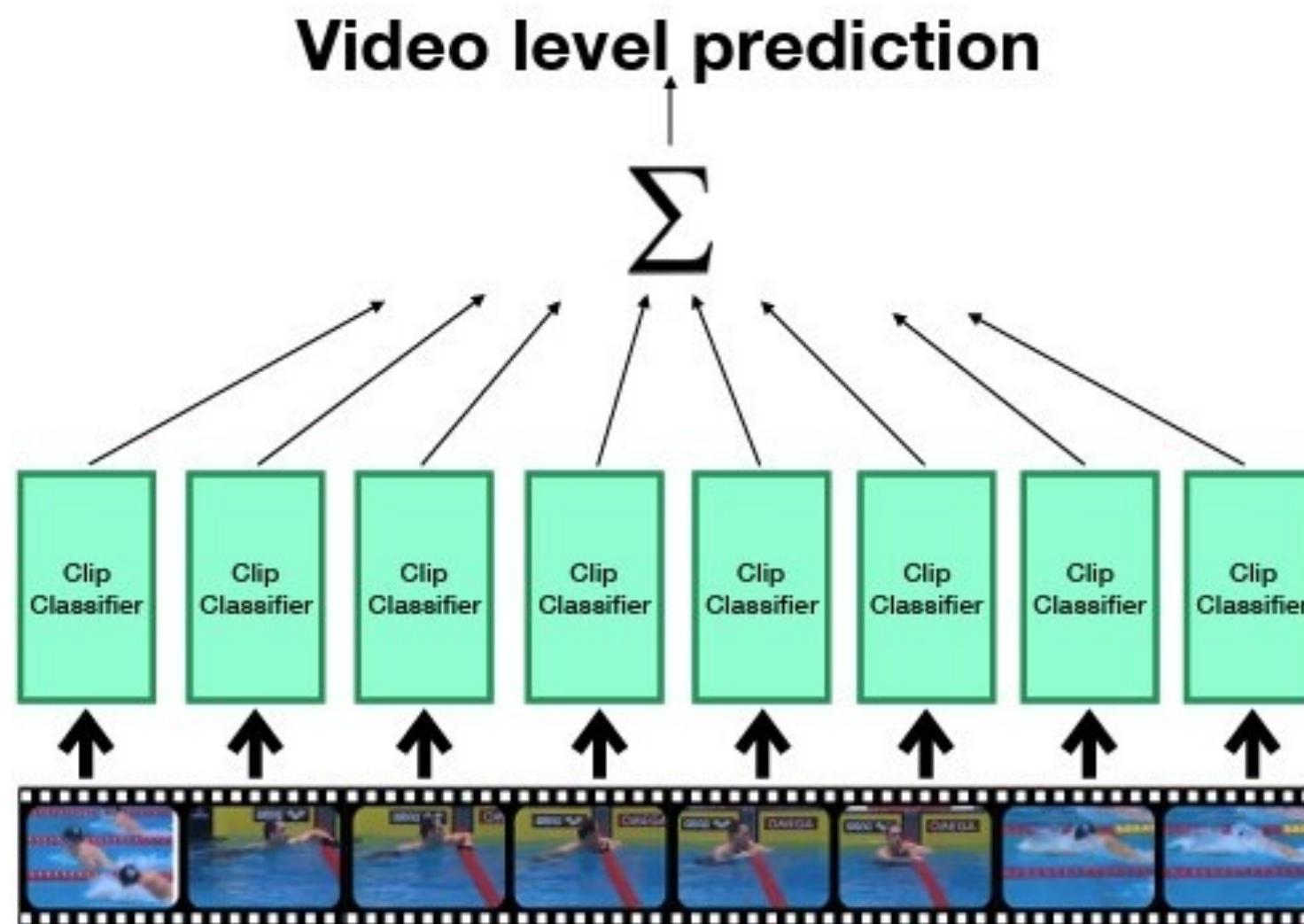
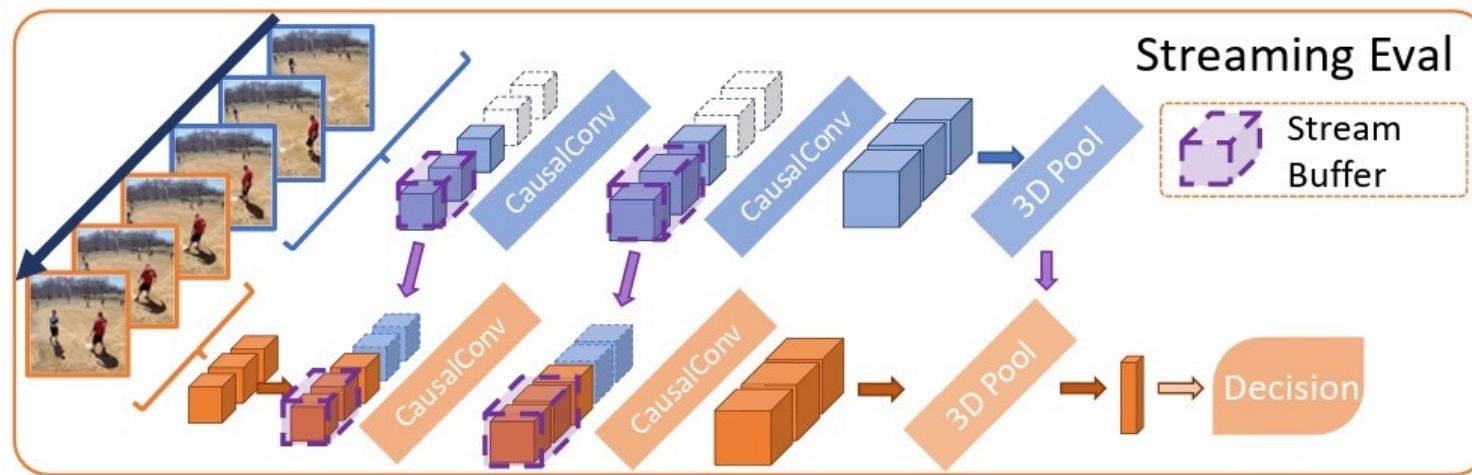


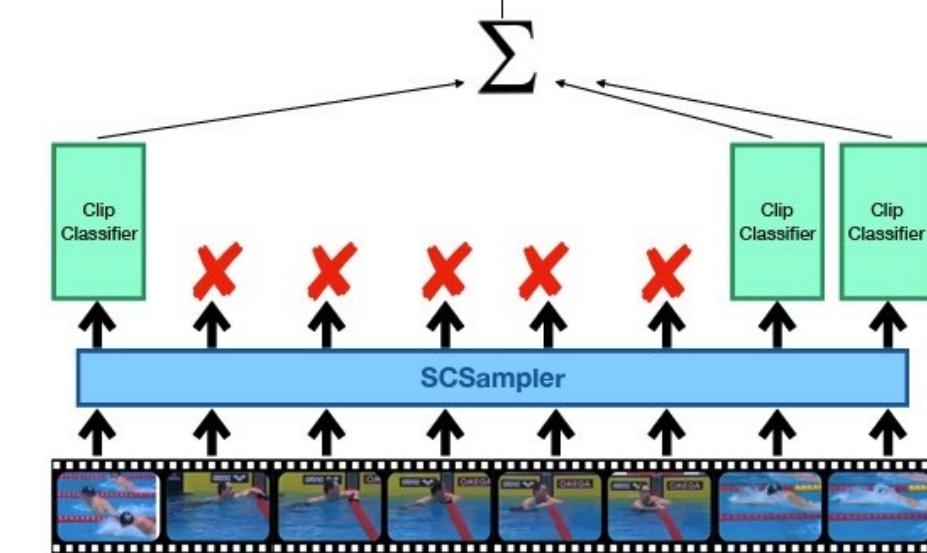
Image Credit: Korbar et al.

# Efficient Video Understanding

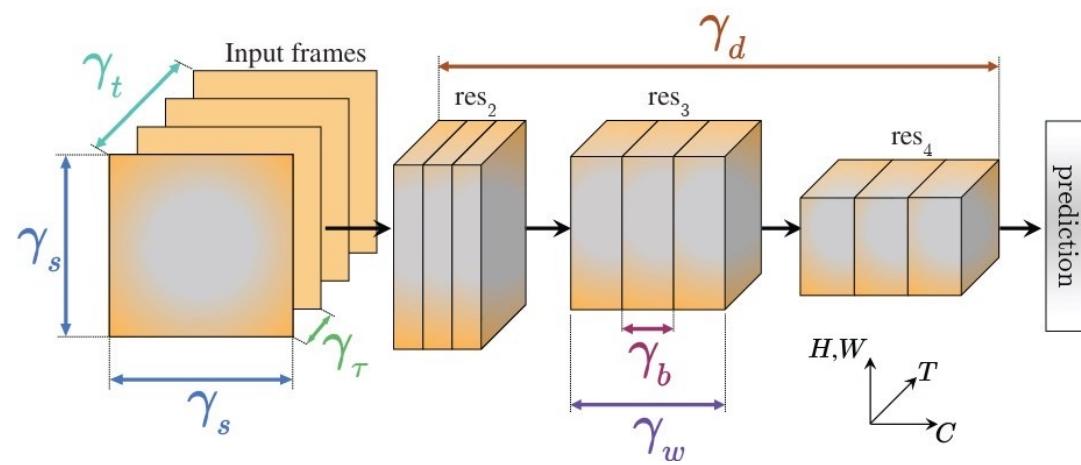


MoViNets: Mobile Video Networks for Efficient Video Recognition.  
Kondratyuk et al. CVPR 2021

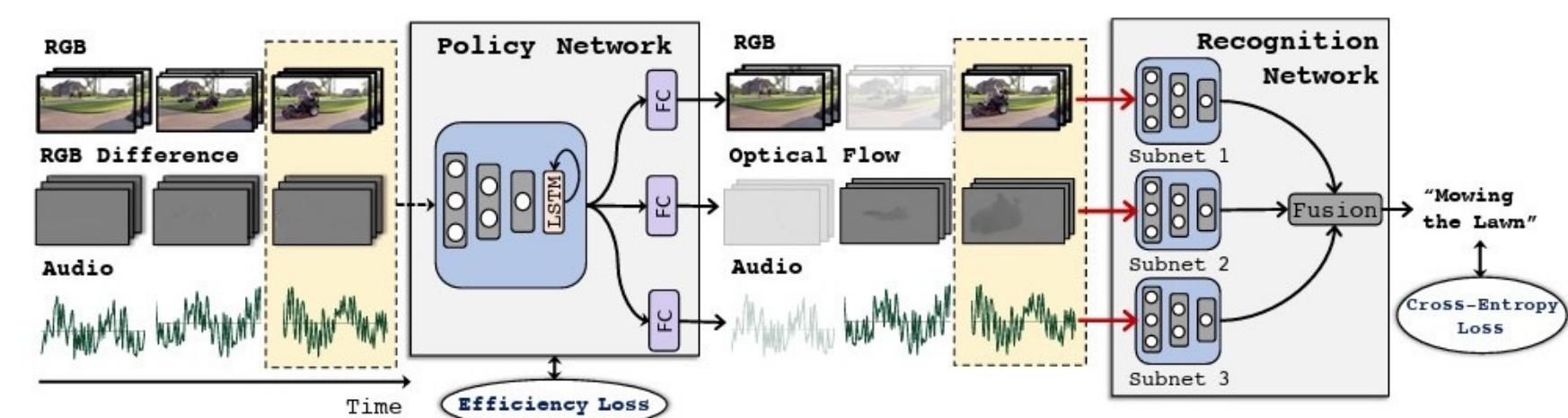
## Video level prediction



SCSampler: Sampling Salient Clips from Video for Efficient Action Recognition. Korbar et al. ICCV 2019

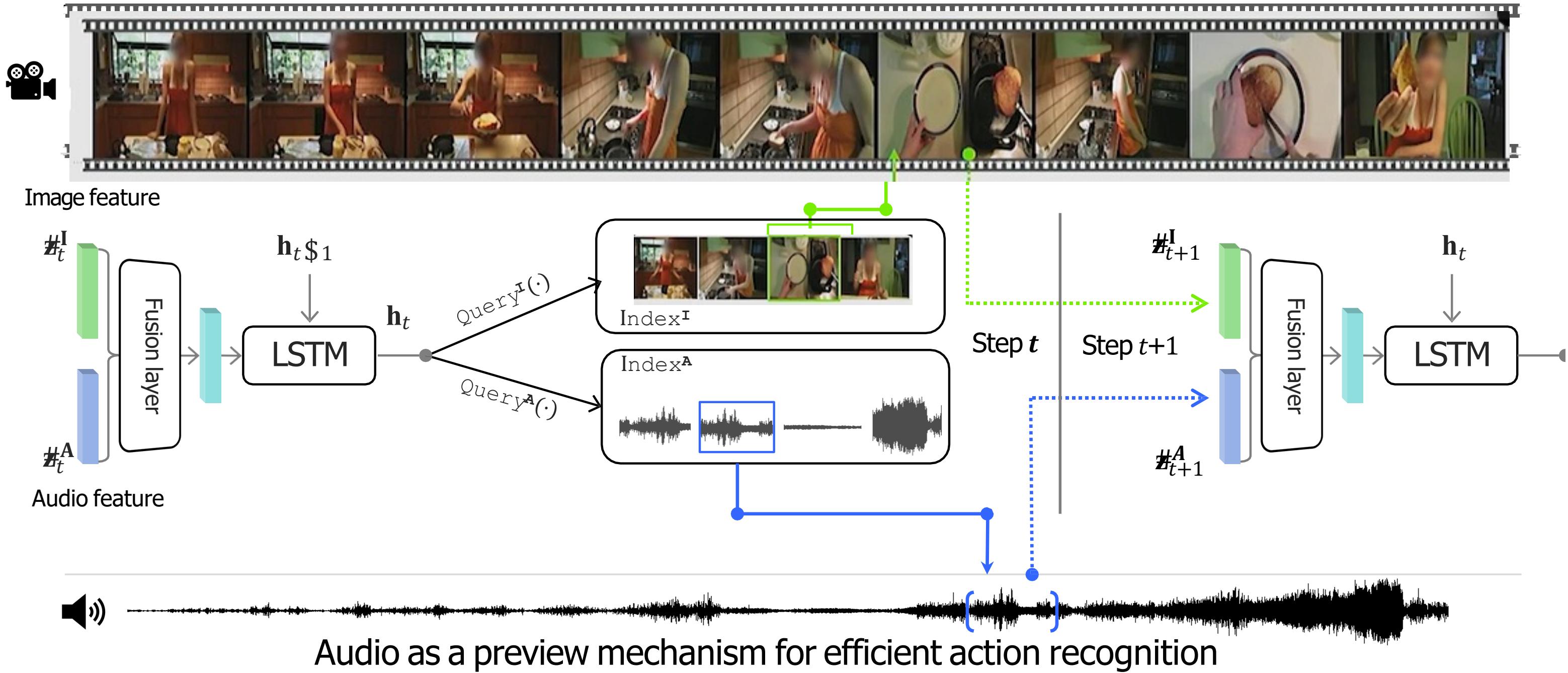


X3D: Expanding Architecture for Efficient Video Recognition.  
Christoph Feichtenhofer. CVPR 2020.



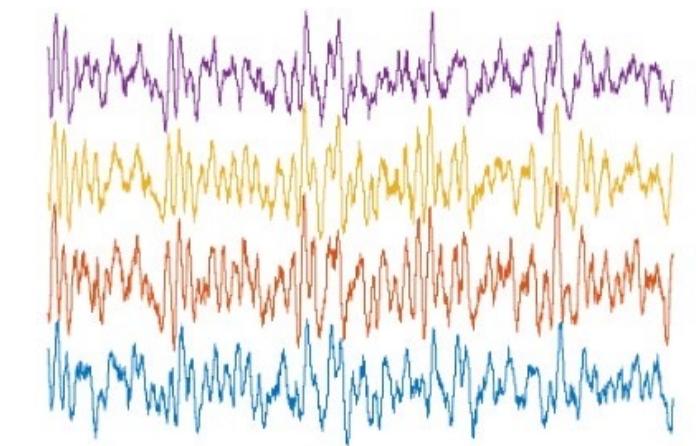
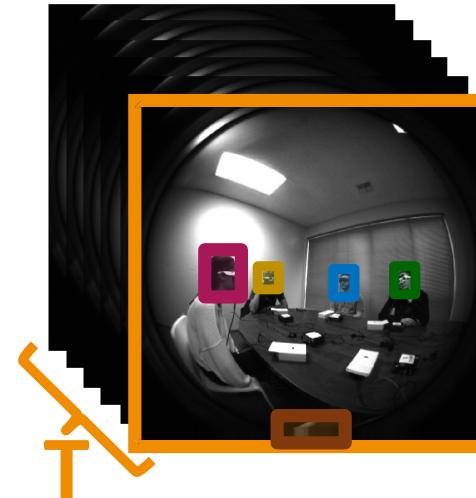
AdaMML: Adaptive Multi-Modal Learning for Efficient Video Recognition. Pandal et al. ICCV 2021

# Efficient Video Understanding

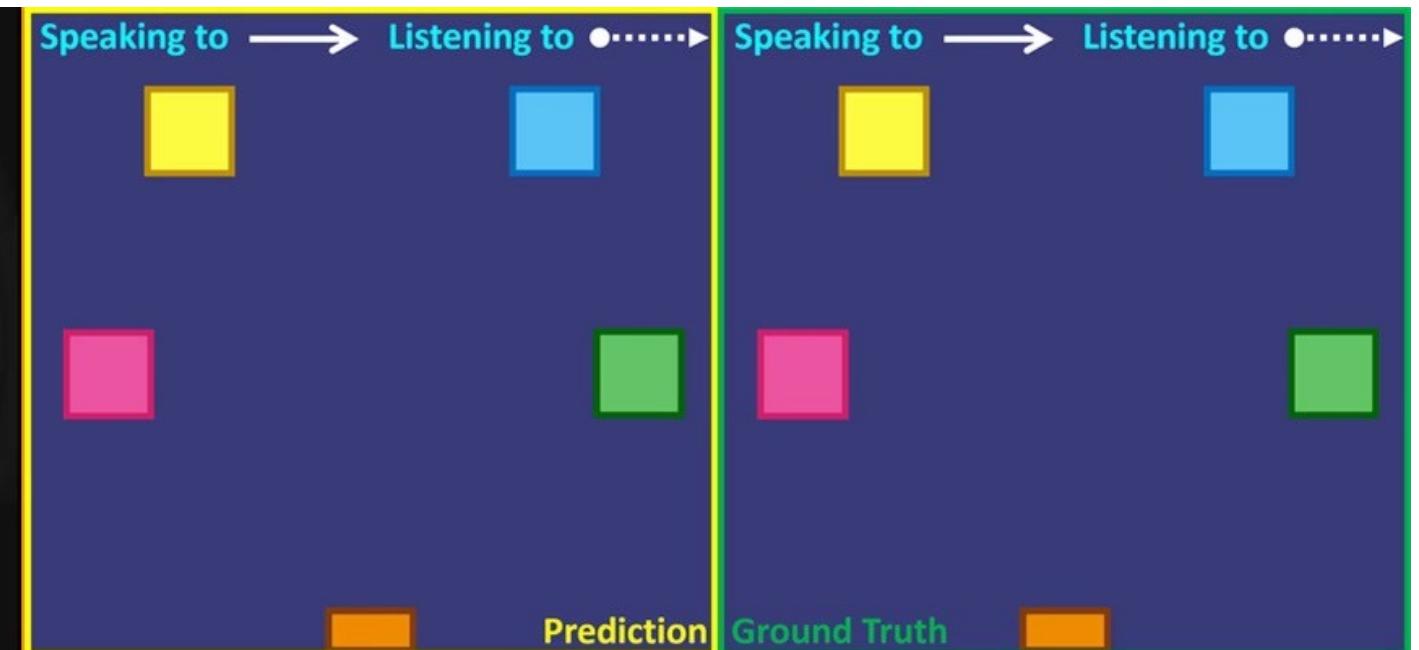


Gao et al., Listen to Look: Action Recognition by Previewing Audio, CVPR 2020

# Multimodal Egocentric Video Understanding

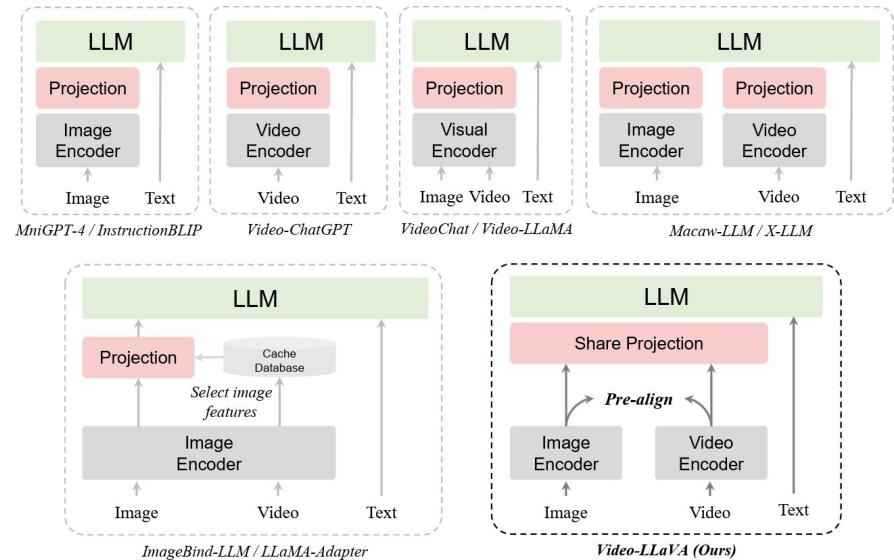


## Ego-Exo Conversational Graph Prediction

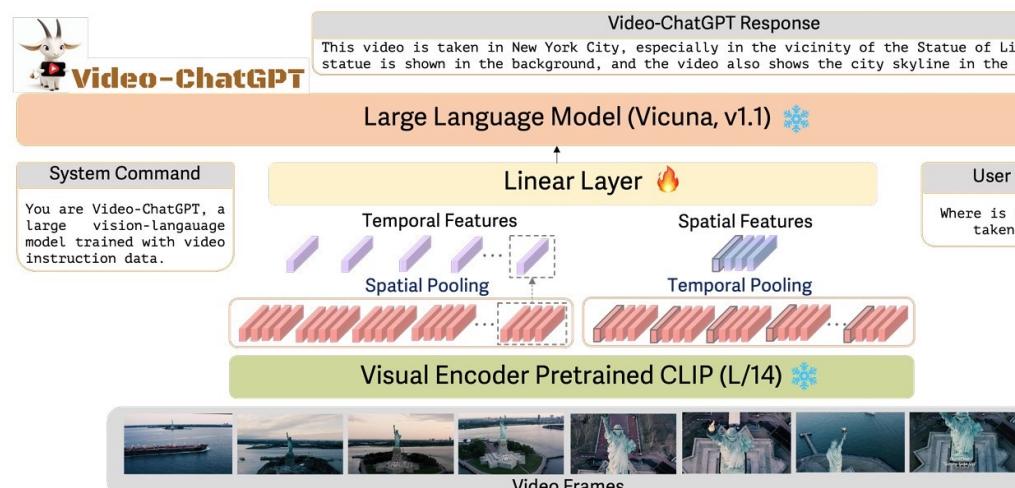


The Audio-Visual Conversational Graph: From and Egocentric-Exocentric Perspective. Jia et al. CVPR 2024

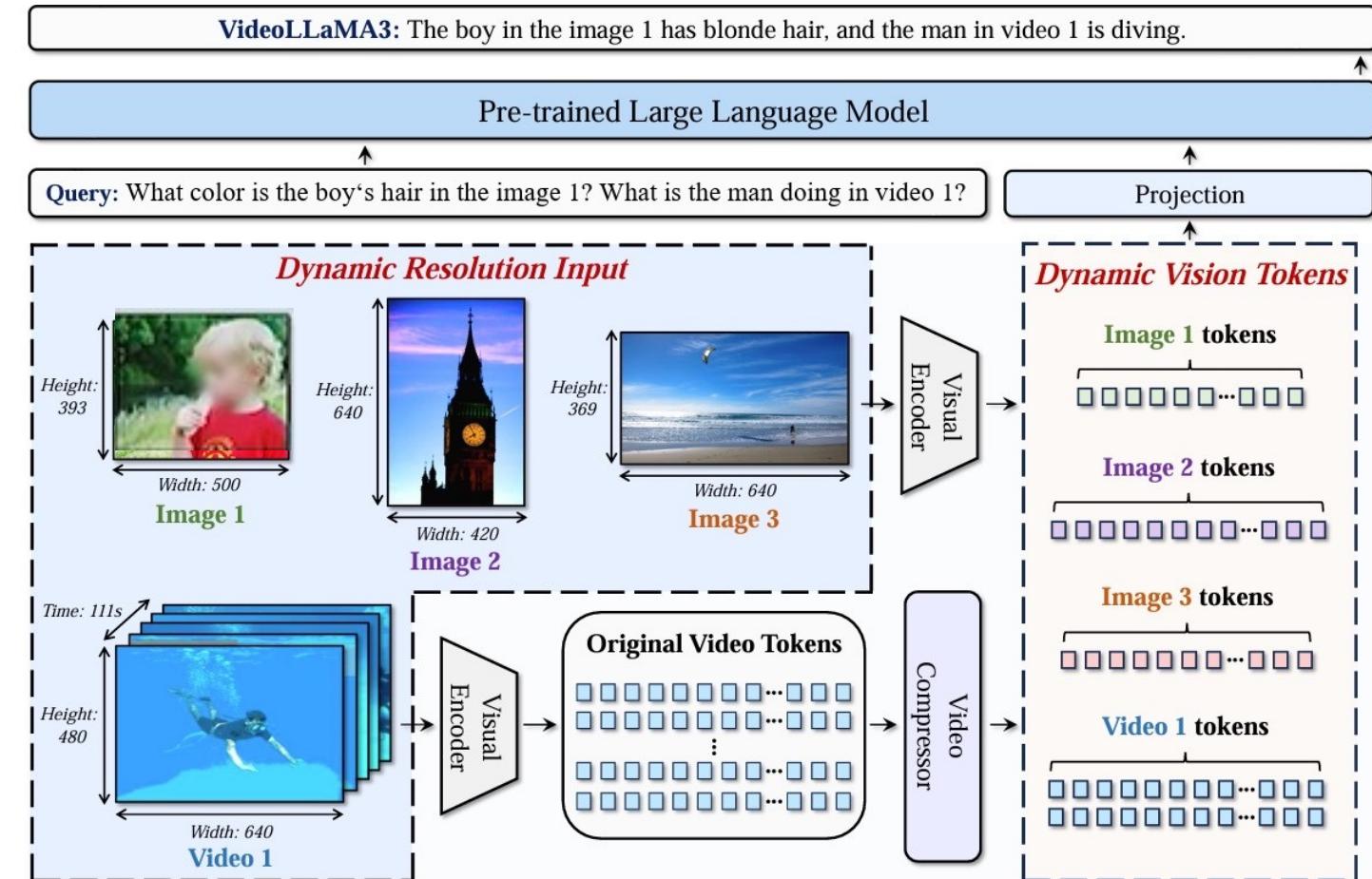
# Video Understanding +LLMs



**Video-LLaVA: Learning United Visual Representations by Alignment Before Projection.** Lin et al. EMNLP 2024



**Video-ChatGPT: Towards Detailed Video Understanding via Large Vision and Language Models.** Maaz et al. ACL 2024.



**VideoLLaMA 3: Frontier Multimodal Foundation Models for Video Understanding.** Zhang et al. arXiv 2025

# Next time: Large Scale Distributed Training

