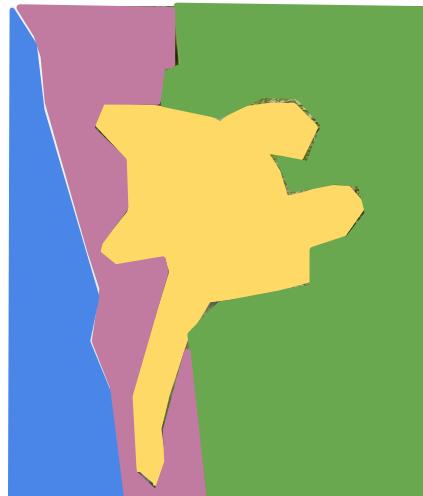


# Lecture 18: Videos

# Computer Vision Tasks: 2D Recognition

## Semantic Segmentation



GRASS, CAT, TREE,  
SKY

No spatial extent  
No objects, just pixels

## Object Detection



DOG, DOG, CAT

Multiple Objects

## Instance Segmentation



DOG, DOG, CAT

This image is CC0 public domain

Justin Johnson

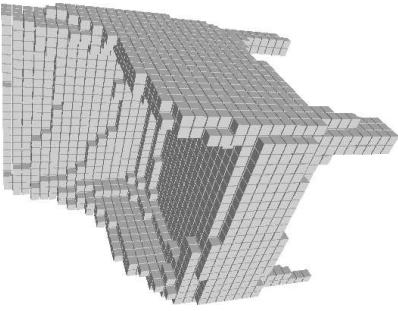
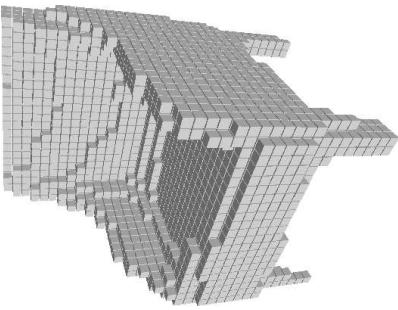
Lecture 18 - 2

November 18, 2019

# Last Time: 3D Shapes

Predicting 3D Shapes  
from single image

Processing 3D  
input data



Input Image

3D Shape

3D Shape



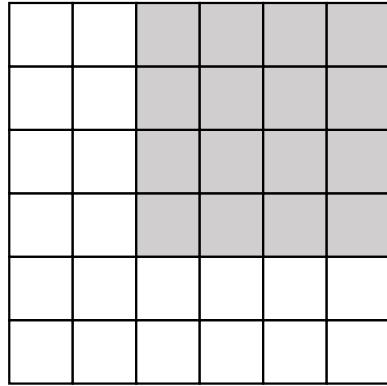
Chair

# Last Time: 3D Shape Representations

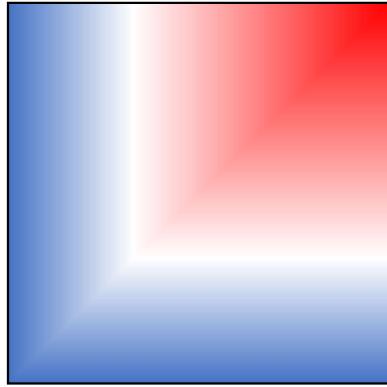
Depth Map



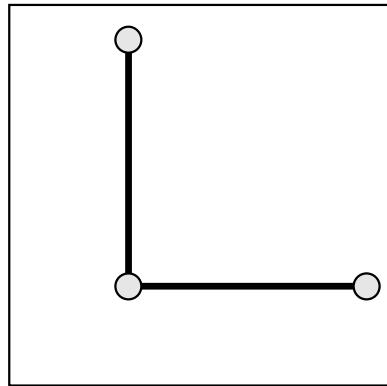
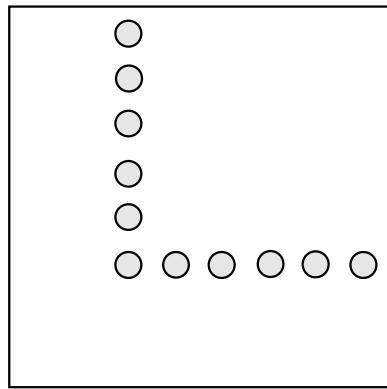
Voxel Grid



Implicit Surface



Pointcloud  
Mesh



Today: Video = 2D + Time

A video is a **sequence** of images  
4D tensor:  $T \times 3 \times H \times W$   
(or  $3 \times T \times H \times W$ )



# Example task: Video Classification

Swimming  
Running  
Jumping  
Eating  
Standing



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

[Running video](#) is in the [public domain](#)

Justin Johnson

Lecture 18 - 6

November 18, 2019

# Example task: Video Classification



Images: Recognize objects



Dog  
**Cat**  
Fish  
Truck

Videos: Recognize actions



Swimming  
**Running**  
Jumping  
Eating  
Standing



[Running video](#) is in the [public domain](#)

Justin Johnson

Lecture 18 - 7

November 18, 2019

**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)



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

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



Justin Johnson

Lecture 18 - 10

November 18, 2019

# Training on Clips

**Raw video:** Long, high FPS



**Training:** Train model to classify short clips with low FPS



# Training on Clips

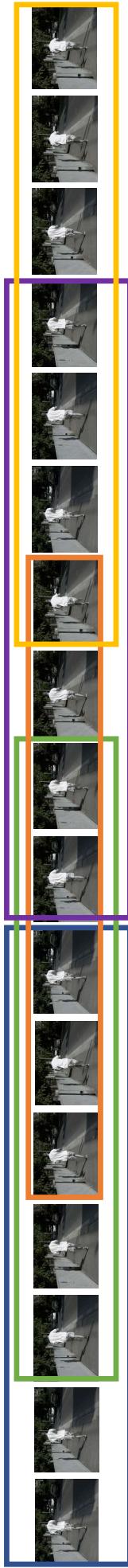
**Raw video:** Long, high FPS



**Training:** Train model to classify short clips with low FPS

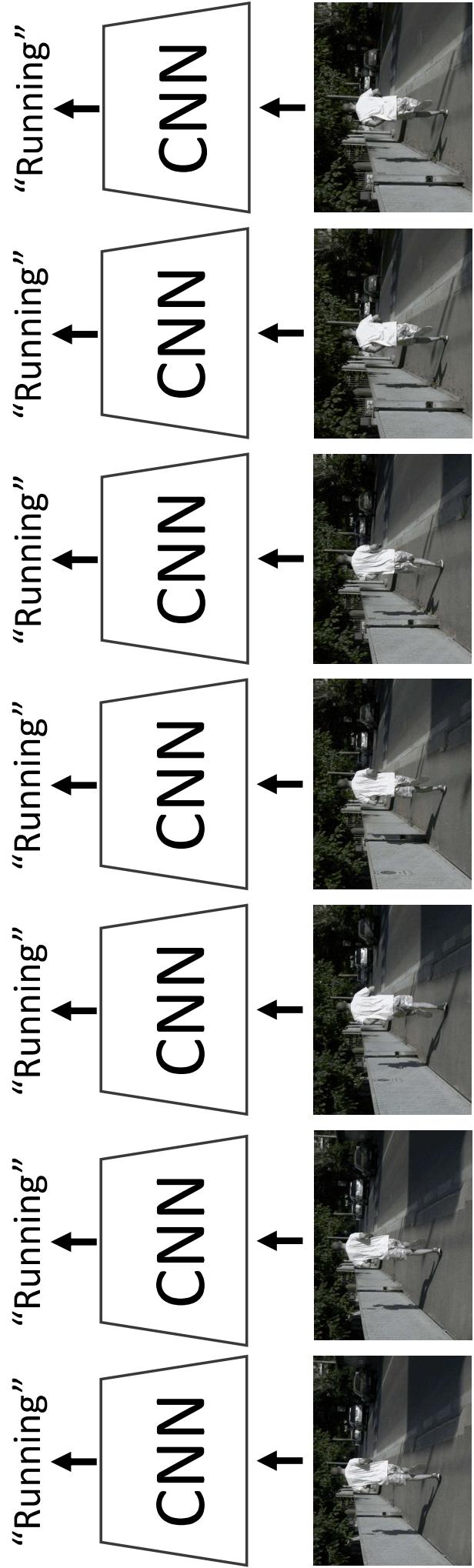


**Testing:** Run model on different clips, average predictions



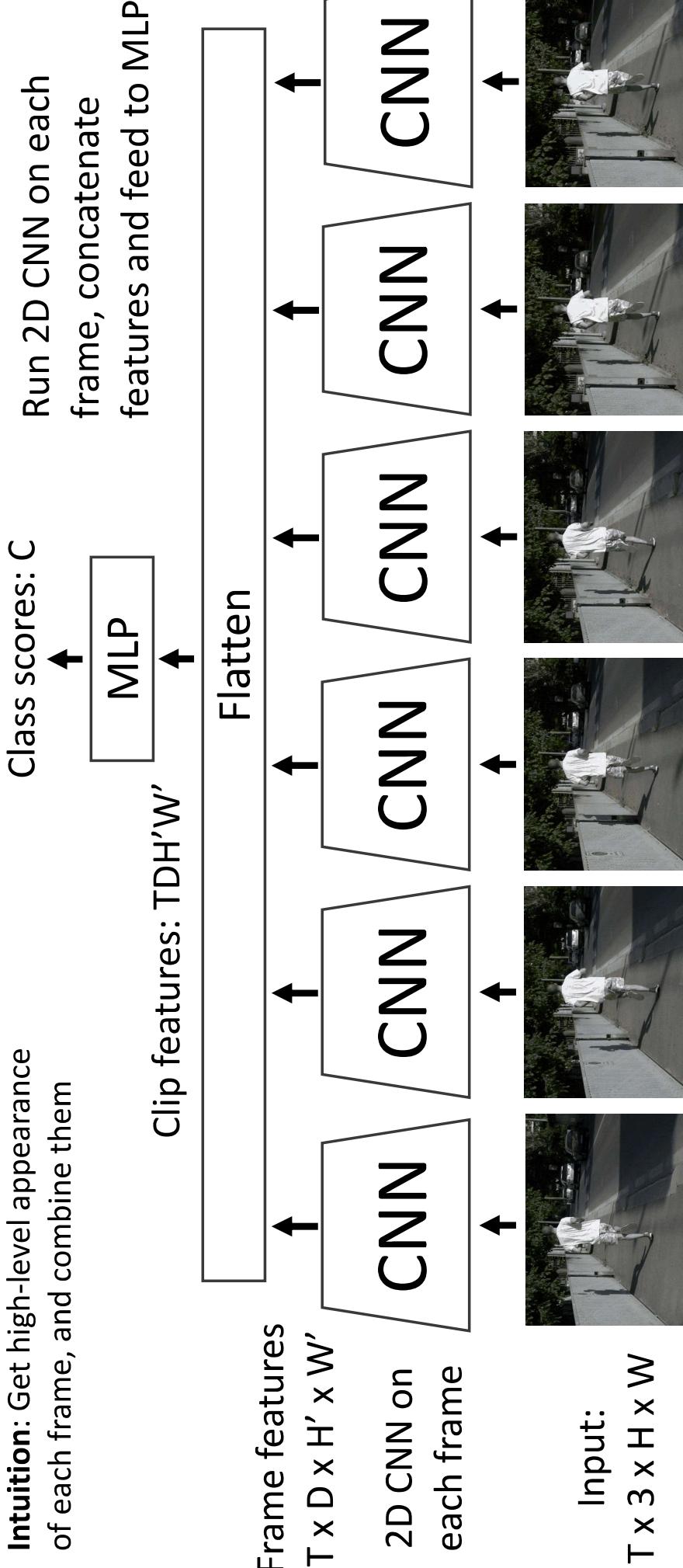
# Video Classification: Single-Frame CNN

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

Justin Johnson

Lecture 18 - 14

November 18, 2019

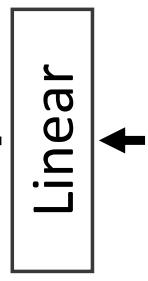
# Video Classification: Late Fusion (with pooling)

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

Class scores: C

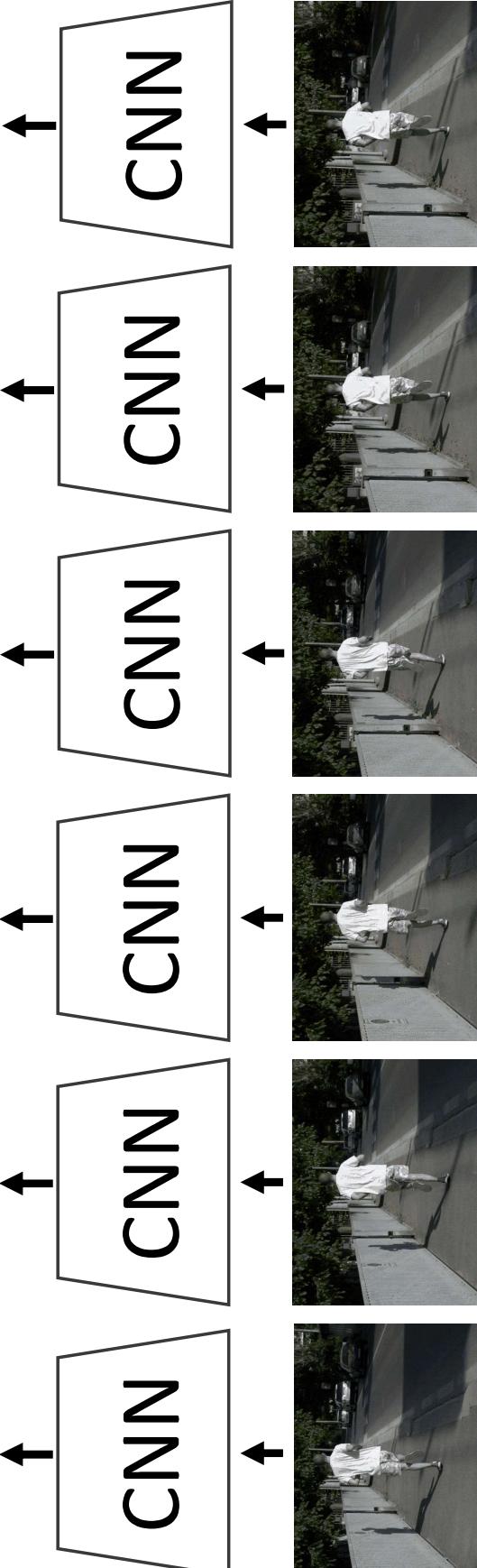
Run 2D CNN on each frame, pool features and feed to Linear

Clip features: D



Frame features  
 $T \times D \times H' \times W'$   
2D CNN on each frame

Average Pool over space and time



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

# 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

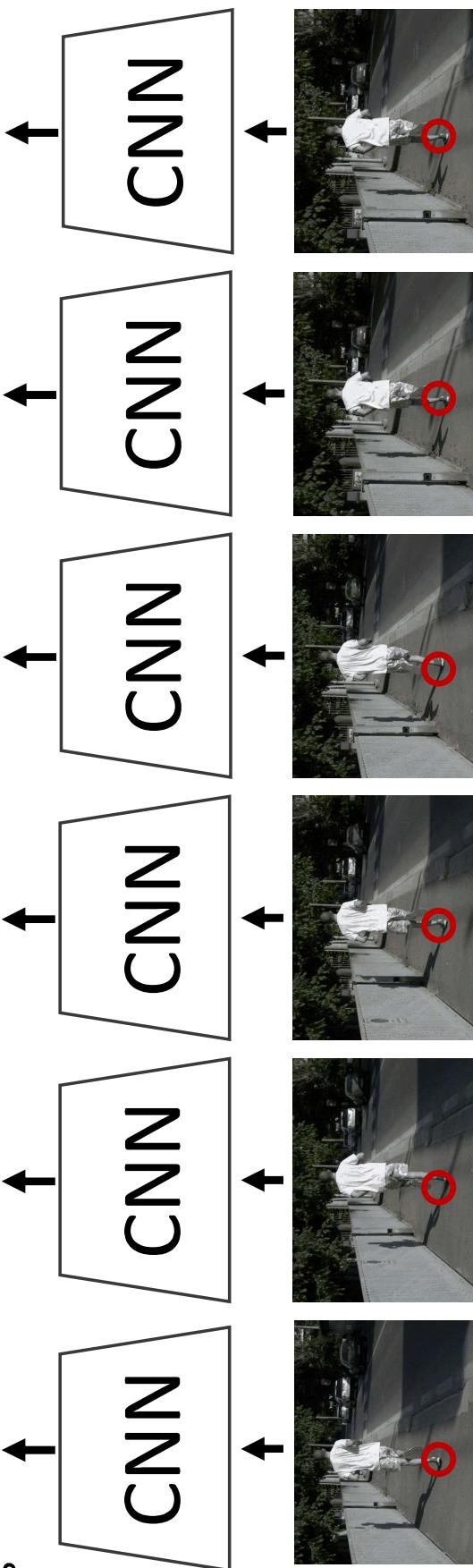
Class scores: C

Run 2D CNN on each frame, pool features and feed to Linear

Clip features: D

Frame features  
 $T \times D \times H' \times W'$   
2D CNN on each frame

Average Pool over space and time



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

# Video Classification: Early Fusion

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

Class scores:  $C$

**2D CNN**

Rest of the network  
is standard 2D CNN

First 2D convolution collapses  
all temporal information:

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

**Output:**  $D \times H \times W$

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

Justin Johnson

Lecture 18 - 17

November 18, 2019

# 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!

Class scores:  $C$

2D CNN

Rest of the network  
is standard 2D CNN

First 2D convolution collapses  
all temporal information:

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

**Output:**  $D \times H \times W$

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

Justin Johnson

Lecture 18 - 18

November 18, 2019

# Video Classification: 3D CNN

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

Class scores:  $C$

3D CNN

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

# Early Fusion vs Late Fusion $\times$ 3D CNN

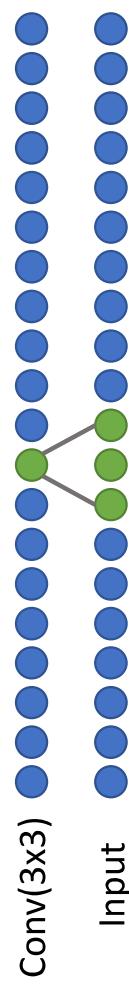
Layer	Size (C x T x H x W)	Receptive Field (T x H x W)
Input	3 x 20 x 64 x 64	
Conv2D(3x3, 3->12)	12 x 20 x 64 x 64	1 x 3 x 3

Late  
Fusion

# Early Fusion vs Late Fusion vs 3D CNN

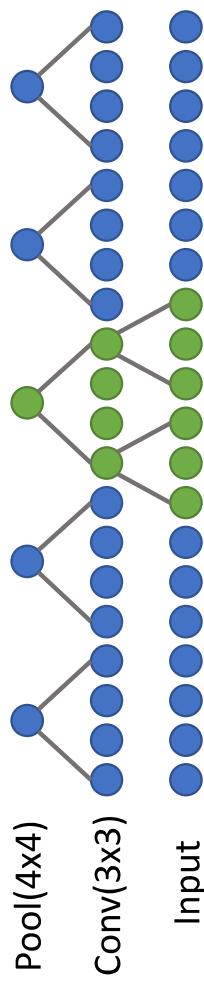
Layer	Size (C x T x H x W)	Receptive Field (T x H x W)
Input	$3 \times 20 \times 64 \times 64$	
Conv2D(3x3, 3->12)	$12 \times 20 \times 64 \times 64$	$1 \times 3 \times 3$

Late  
Fusion



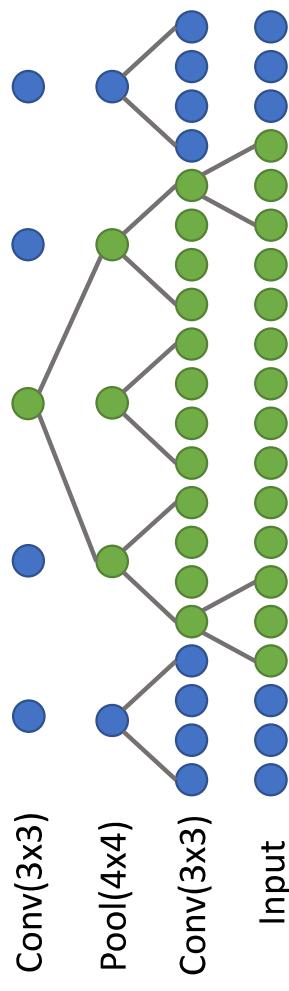
# Early Fusion vs Late Fusion vs 3D CNN

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



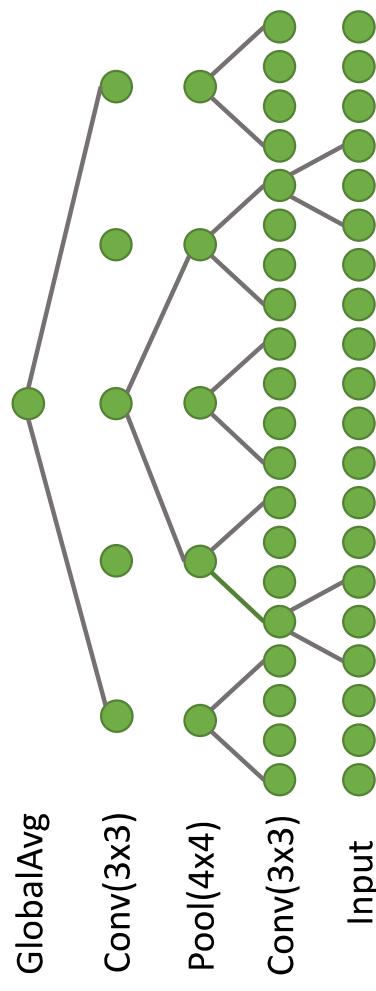
# Early Fusion vs Late Fusion $\times$ 3D CNN

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



# Early Fusion vs Late Fusion $\text{vs}$ 3D CNN

	Layer	Size (C x T x H x W)	Receptive Field (T x H x W)
	Input	$3 \times 20 \times 64 \times 64$	
Late Fusion	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	$24 \times 1 \times 1 \times 1$	$20 \times 64 \times 64$



# Early Fusion vs Late Fusion $\times$ 3D CNN

	Layer	Size (C x T x H x W)	Receptive Field (T x H x W)
Late Fusion	Input	3 x 20 x 64 x 64	
	Conv2D(3x3, 3->12)	12 x 20 x 64 x 64	1 x 3 x 3
	Pool2D(4x4)	12 x 20 x 16 x 16	1 x 6 x 6
	Conv2D(3x3, 12->24)	24 x 20 x 16 x 16	1 x 14 x 14
	GlobalAvgPool	24 x 1 x 1 x 1	20 x 64 x 64
Early Fusion	Input	3 x 20 x 64 x 64	
	Conv2D(3x3, 3*10->12)	12 x 64 x 64	20 x 3 x 3
	Pool2D(4x4)	12 x 16 x 16	20 x 6 x 6
	Conv2D(3x3, 12->24)	24 x 16 x 16	20 x 14 x 14
	GlobalAvgPool	24 x 1 x 1	20 x 64 x 64

Early Fusion

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

# Early Fusion vs Late Fusion VS 3D CNN

(Small example architectures,  
in practice much bigger)

	Layer	Size (C x T x H x W)	Receptive Field (T x H x W)
Late Fusion	Input	3 x 20 x 64 x 64	
	Conv2D(3x3, 3->12)	12 x 20 x 64 x 64	1 x 3 x 3
	Pool2D(4x4)	12 x 20 x 16 x 16	1 x 6 x 6
	Conv2D(3x3, 12->24)	24 x 20 x 16 x 16	1 x 14 x 14
Early Fusion	GlobalAvgPool	24 x 1 x 1 x 1	20 x 64 x 64
	Input	3 x 20 x 64 x 64	
	Conv2D(3x3, 3*10->12)	12 x 64 x 64	20 x 3 x 3
	Pool2D(4x4)	12 x 16 x 16	20 x 6 x 6
3D CNN	Conv2D(3x3, 12->24)	24 x 16 x 16	20 x 14 x 14
	GlobalAvgPool	24 x 1 x 1	20 x 64 x 64
	Input	3 x 20 x 64 x 64	
	Conv3D(3x3x3, 3->12)	12 x 20 x 64 x 64	3 x 3 x 3
	Pool3D(4x4x4)	12 x 5 x 16 x 16	6 x 6 x 6
	Conv3D(3x3x3, 12->24)	24 x 5 x 16 x 16	14 x 14 x 14
	GlobalAvgPool	24 x 1 x 1	20 x 64 x 64

# Early Fusion vs Late Fusion $\text{vs}$ 3D CNN

What is the difference?

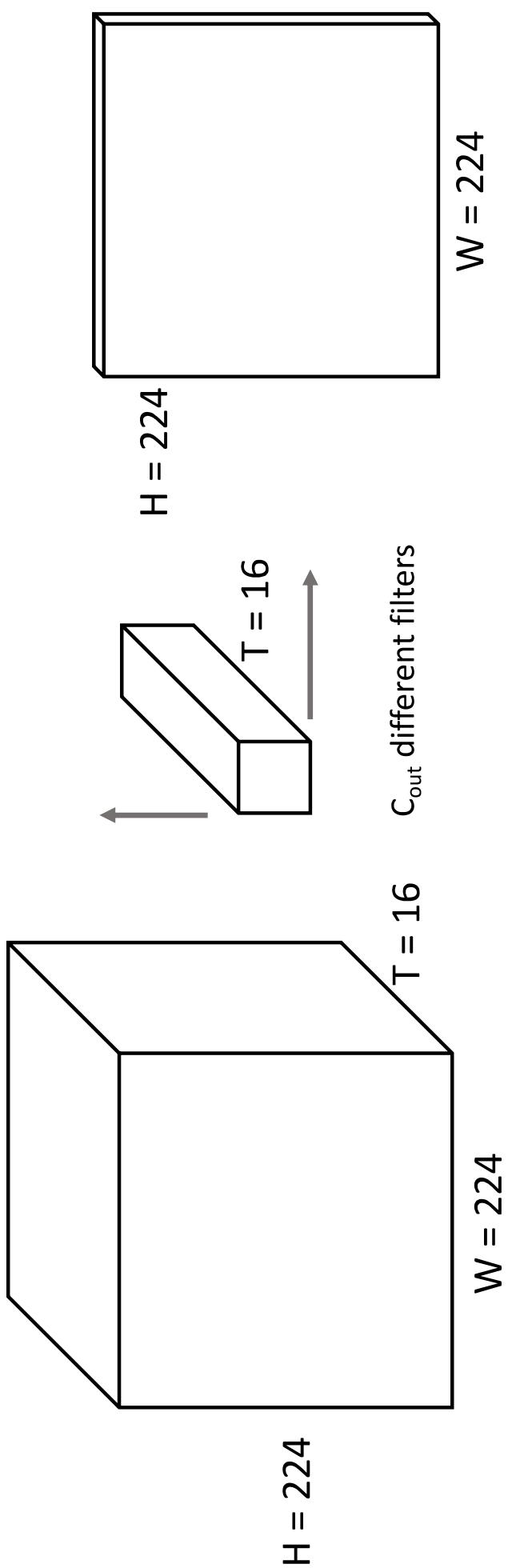
	Layer	Size (C $\times$ T $\times$ H $\times$ W)	Receptive Field (T $\times$ H $\times$ W)	
Late Fusion	Input	3 $\times$ 20 $\times$ 64 $\times$ 64		
	Conv2D(3x3, 3->12)	12 $\times$ 20 $\times$ 64 $\times$ 64	1 $\times$ 3 $\times$ 3	Build slowly in space, All-at-once in time at end
	Pool2D(4x4)	12 $\times$ 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	
Early Fusion	GlobalAvgPool	24 $\times$ 1 $\times$ 1 $\times$ 1	20 $\times$ 64 $\times$ 64	
	Input	3 $\times$ 20 $\times$ 64 $\times$ 64		
	Conv2D(3x3, 3*10->12)	12 $\times$ 64 $\times$ 64	20 $\times$ 3 $\times$ 3	Build slowly in space, All-at-once in time at start
	Pool2D(4x4)	12 $\times$ 16 $\times$ 16	20 $\times$ 6 $\times$ 6	
3D CNN	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	
	Input	3 $\times$ 20 $\times$ 64 $\times$ 64		
	Conv3D(3x3x3, 3->12)	12 $\times$ 20 $\times$ 64 $\times$ 64	3 $\times$ 3 $\times$ 3	Build slowly in space, Build slowly in time "Slow Fusion"
	Pool3D(4x4x4)	12 $\times$ 5 $\times$ 16 $\times$ 16	6 $\times$ 6 $\times$ 6	
	Conv3D(3x3x3, 12->24)	24 $\times$ 5 $\times$ 16 $\times$ 16	14 $\times$ 14 $\times$ 14	
	GlobalAvgPool	24 $\times$ 1 $\times$ 1	20 $\times$ 64 $\times$ 64	

## 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



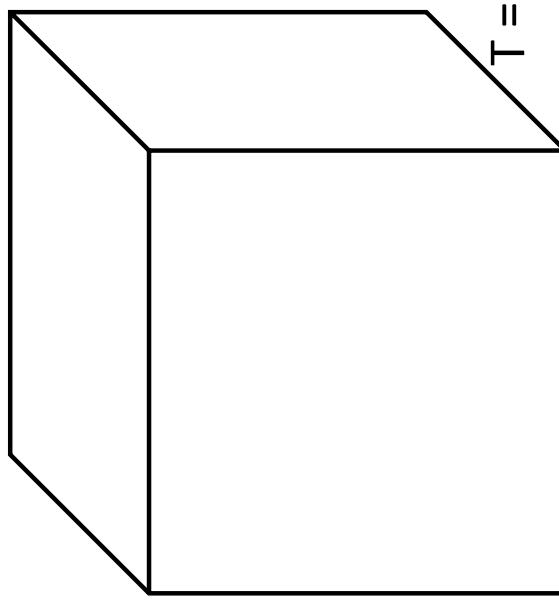
## 2D Conv (Early Fusion) vs 3D Conv (3D CNN)

**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

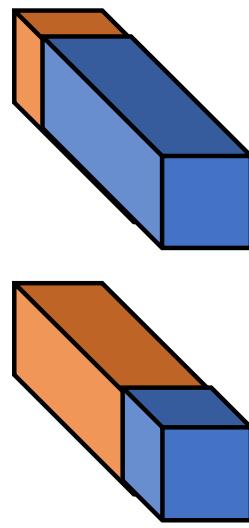
**Output:**  
 $C_{out} \times H \times W$   
2D grid with  $C_{out}$ -dim  
feat at each point



$H = 224$   
 $T = 16$   
 $C_{out}$  different filters



$W = 224$



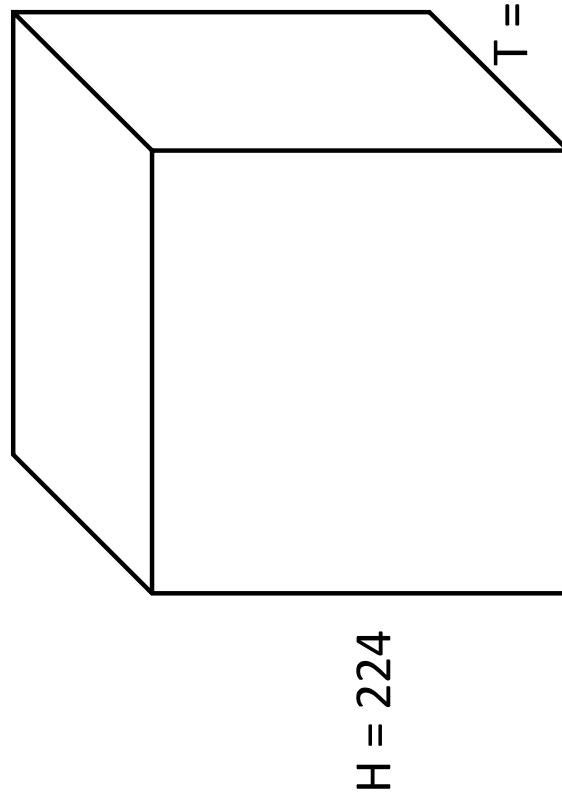
No temporal shift-invariance! Needs  
to learn separate filters for the same  
motion at different times in the clip

## 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



$$W = 224$$

$$H = 224$$

$$T = 3$$

$$C_{out}$$
 different filters

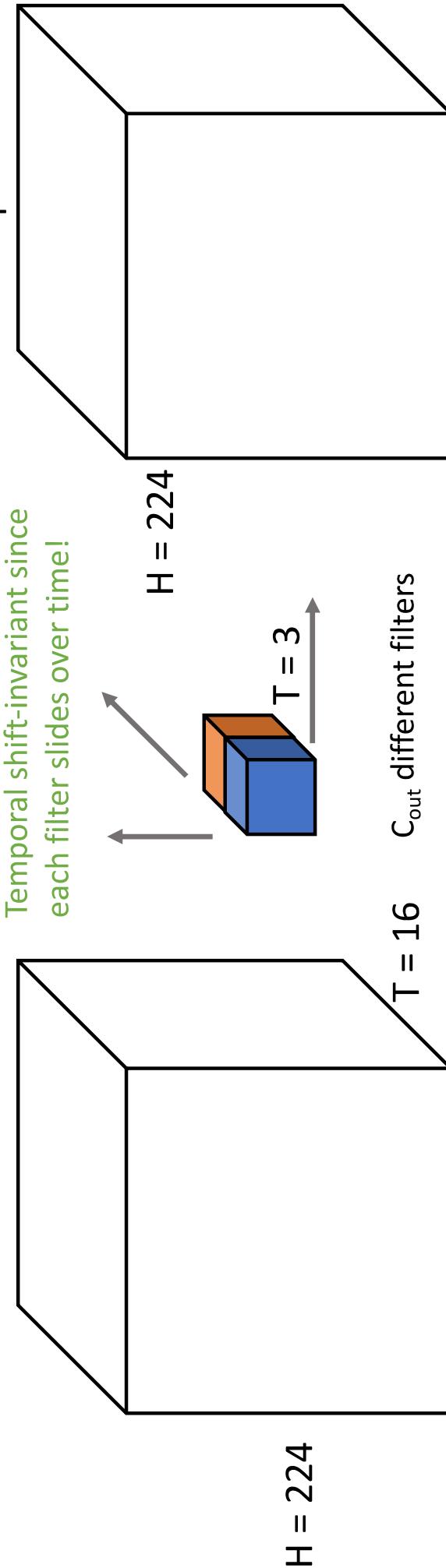
## 2D Conv (Early Fusion) vs 3D Conv (3D CNN)

**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

Temporal shift-invariant since  
each filter slides over time!



$$W = 224$$

$$T = 16$$

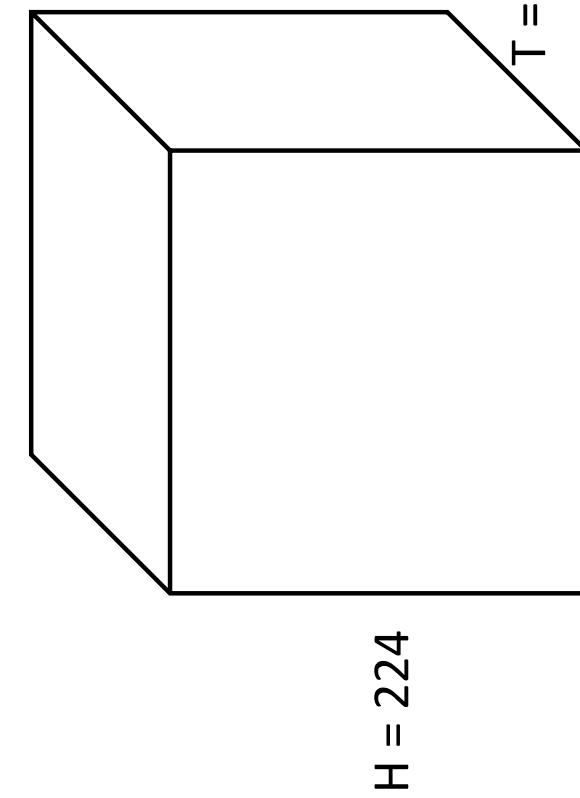
$$H = 224$$

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

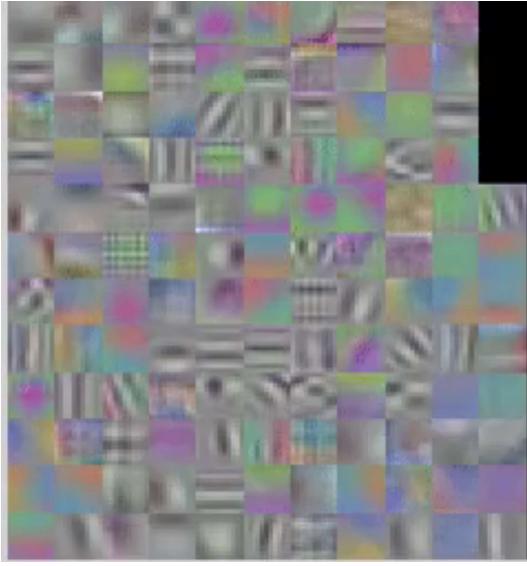
**Input:**  $C_{in} \times T \times H \times W$   
(3D grid with  $C_{in}$ -dim  
feat at each point)

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

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



Temporal shift-invariant since  
each filter slides over time!



$W = 224$

$T = 16$   $C_{out}$  different filters

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

# Example Video Dataset: Sports-1M



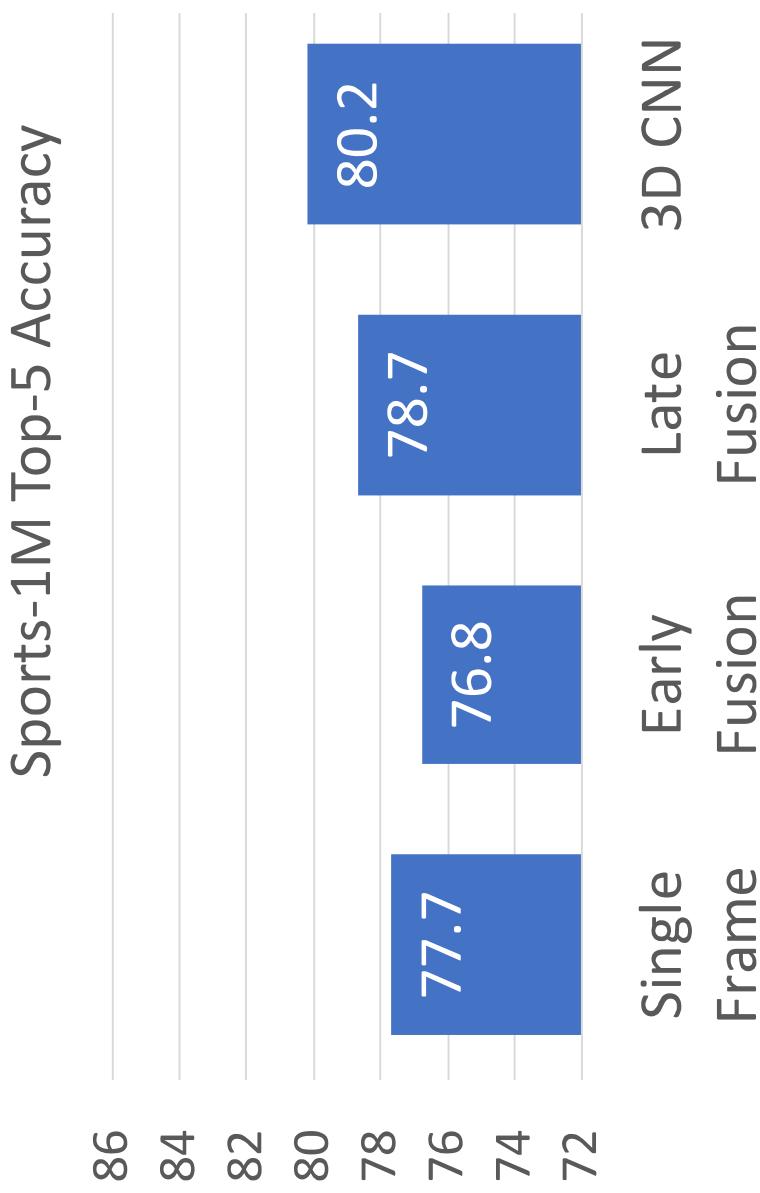
1 million YouTube videos  
annotated with labels for  
487 different types of sports

Ground Truth

Correct prediction

Incorrect prediction

# Early Fusion vs Late Fusion vs 3D CNN



Single Frame model  
works well – always  
try this first!

3D CNNs have  
improved a lot  
since 2014!

## 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

Layer	Size
Input	$3 \times 16 \times 112 \times 112$
Conv1 ( $3 \times 3 \times 3$ )	$64 \times 16 \times 112 \times 112$
Pool1 ( $1 \times 2 \times 2$ )	$64 \times 16 \times 56 \times 56$
Conv2 ( $3 \times 3 \times 3$ )	$128 \times 16 \times 56 \times 56$
Pool2 ( $2 \times 2 \times 2$ )	$128 \times 8 \times 28 \times 28$
Conv3a ( $3 \times 3 \times 3$ )	$256 \times 8 \times 28 \times 28$
Conv3b ( $3 \times 3 \times 3$ )	$256 \times 8 \times 28 \times 28$
Pool3 ( $2 \times 2 \times 2$ )	$256 \times 4 \times 14 \times 14$
Conv4a ( $3 \times 3 \times 3$ )	$512 \times 4 \times 14 \times 14$
Conv4b ( $3 \times 3 \times 3$ )	$512 \times 4 \times 14 \times 14$
Pool4 ( $2 \times 2 \times 2$ )	$512 \times 2 \times 7 \times 7$
Conv5a ( $3 \times 3 \times 3$ )	$512 \times 2 \times 7 \times 7$
Conv5b ( $3 \times 3 \times 3$ )	$512 \times 2 \times 7 \times 7$
Pool5	$512 \times 1 \times 3 \times 3$
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

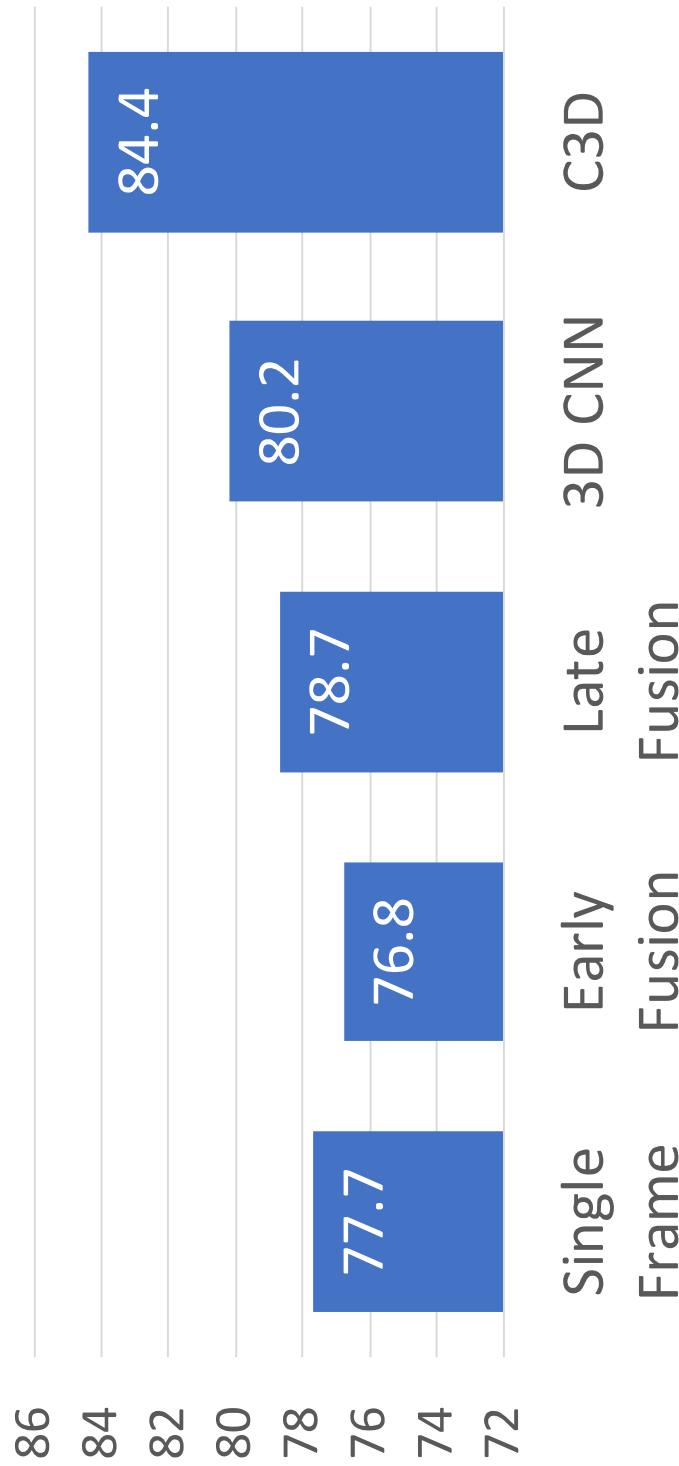
VGG-16: 13.6 GFLOP

**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

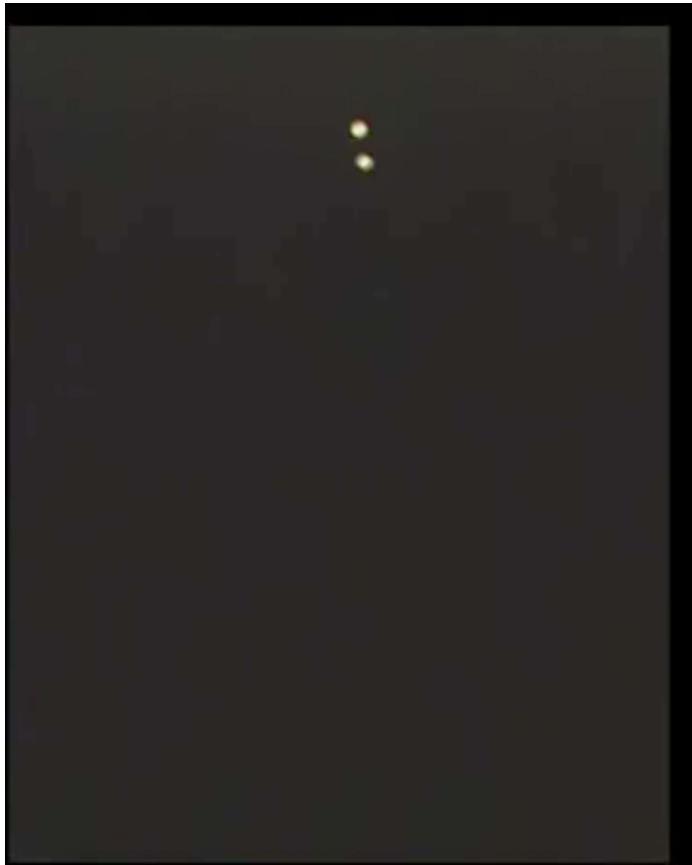
## Sports-1M Top-5 Accuracy



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

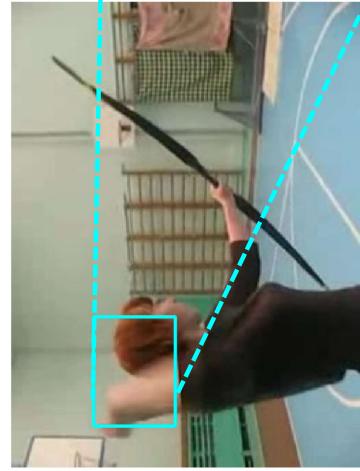
Justin Johnson

Lecture 18 - 39

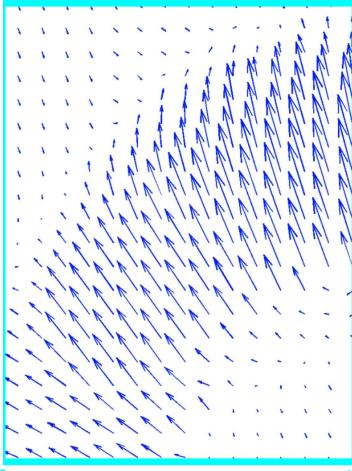
November 18, 2019

# Measuring Motion: Optical Flow

Image at frame  $t$



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:

$$F(x, y) = (dx, dy)$$
$$I_{t+1}(x+dx, y+dy) = I_t(x, y)$$

Image at frame  $t+1$



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

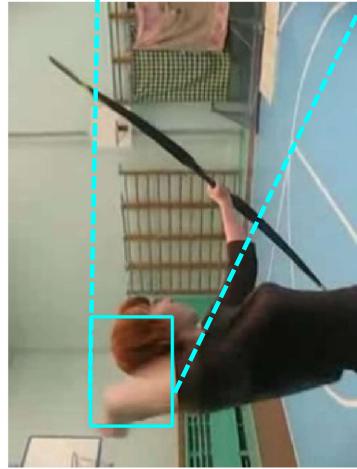
Justin Johnson

Lecture 18 - 40

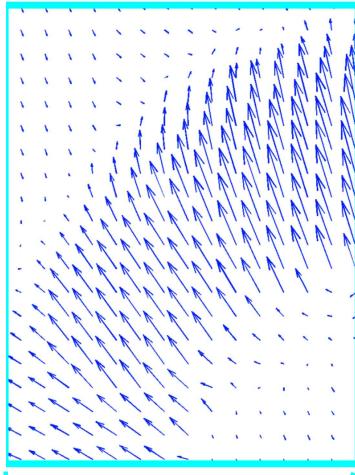
November 18, 2019

# Measuring Motion: Optical Flow

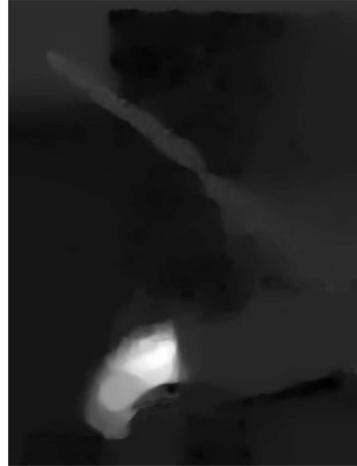
Image at frame t



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



Horizontal flow  $dx$



Vertical flow  $dy$

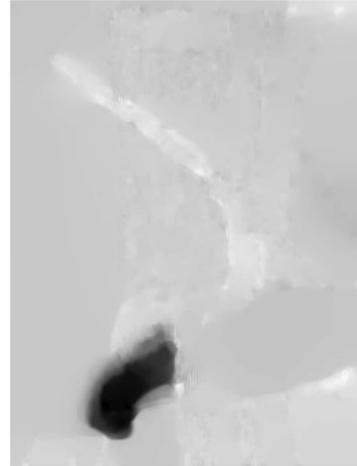


Tells where each pixel will move in the next frame:  
 $F(x, y) = (dx, dy)$   
 $I_{t+1}(x+dx, y+dy) = I_t(x, y)$

Image at frame  $t+1$

Vertical Flow  $dy$

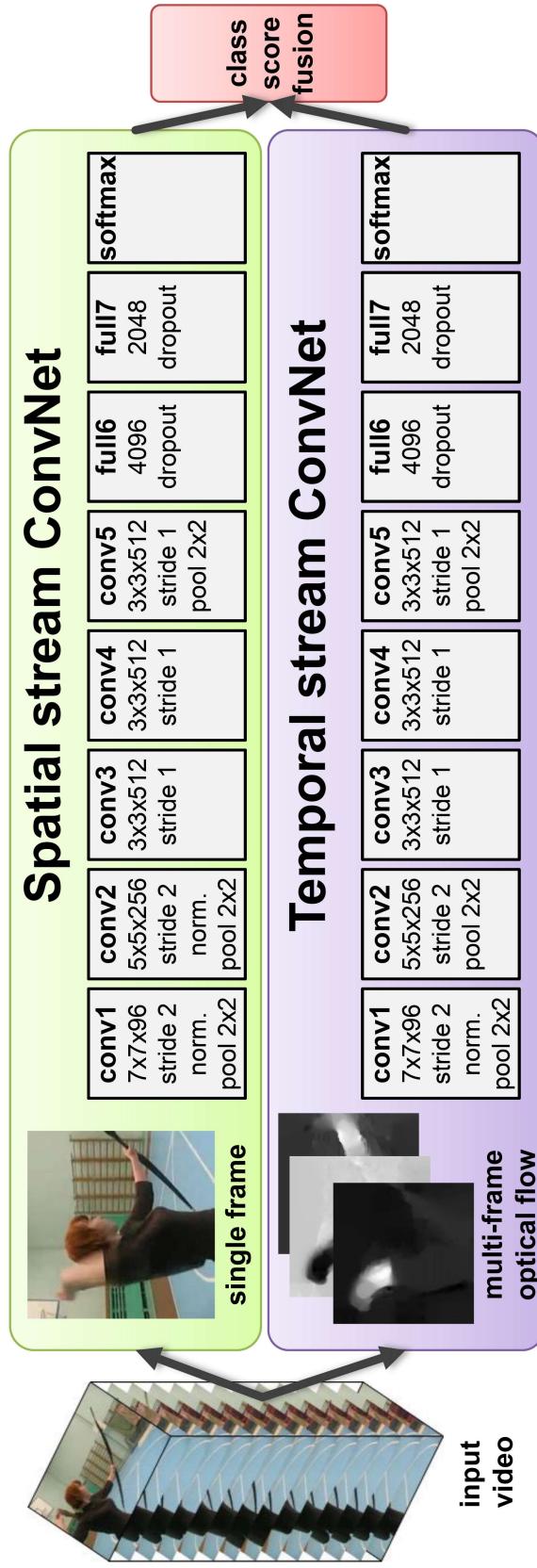
Horizontal flow  $dx$



Optical Flow highlights local motion

# Separating Motion and Appearance: Two-Stream Networks

**Input:** Single Image  
 $3 \times H \times W$



**Inpout:** 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

Accuracy on UCF-101



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

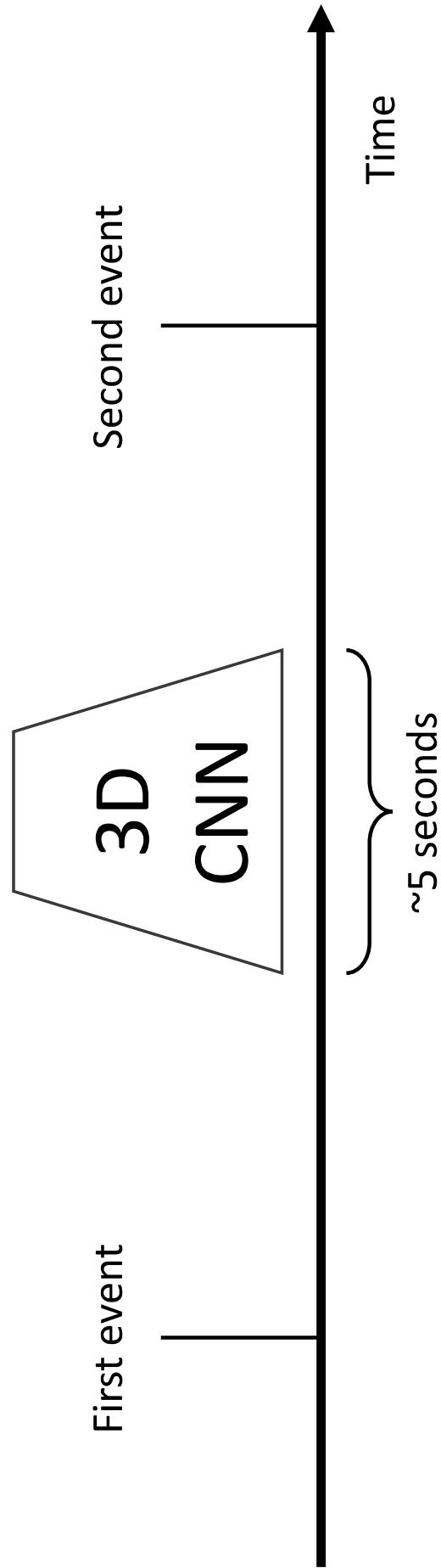
Justin Johnson

Lecture 18 - 43

November 18, 2019

# Modeling long-term temporal structure

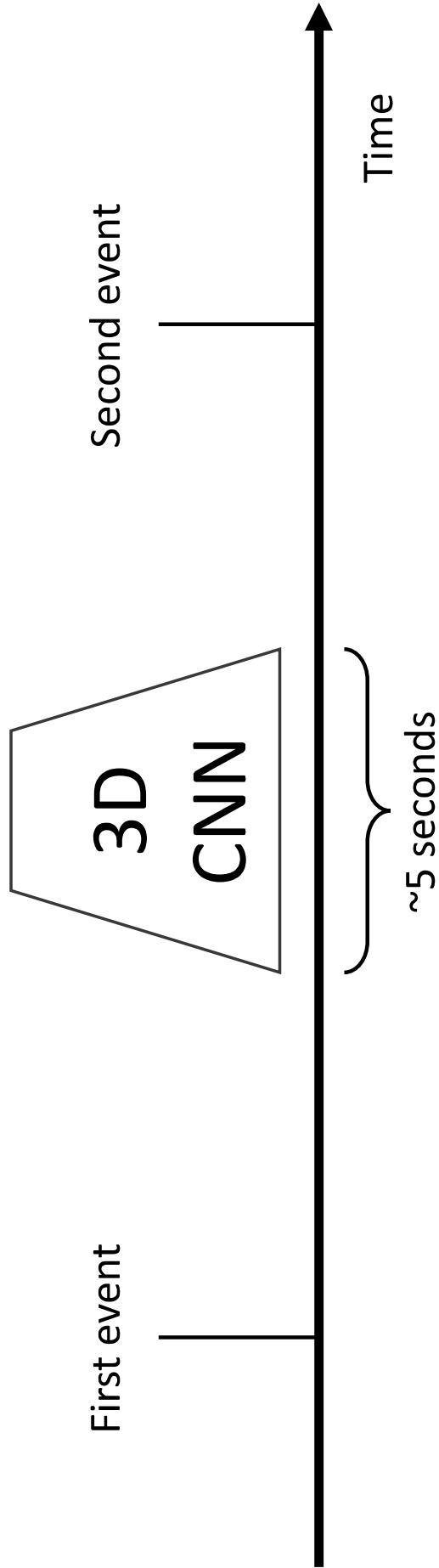
So far all our temporal CNNs only model local motion between frames in very short clips of ~2-5 seconds. What about long-term structure?



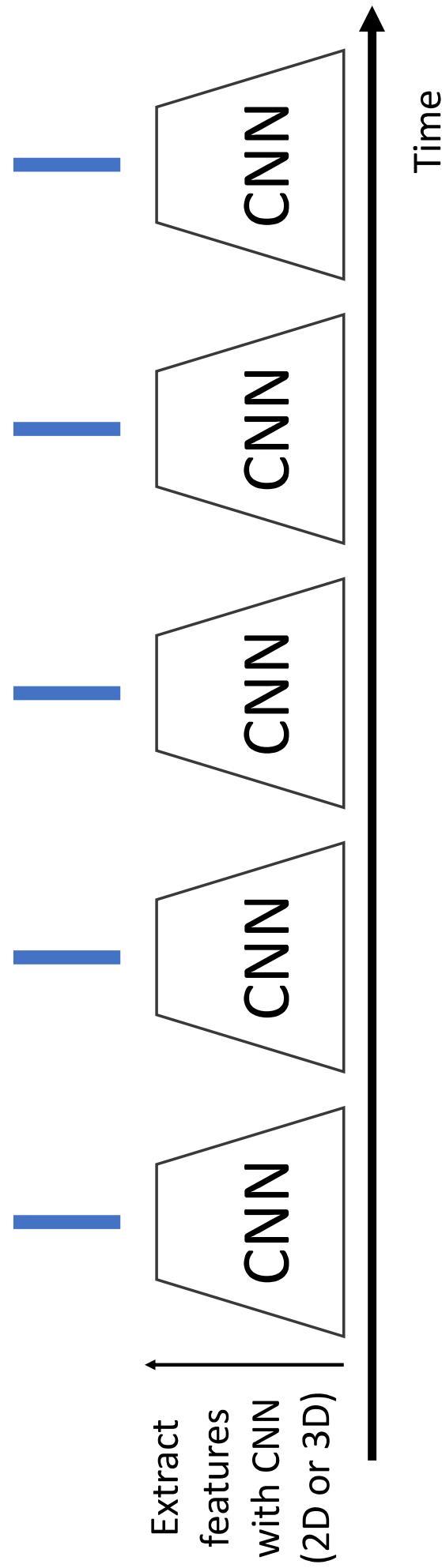
# Modeling long-term temporal structure

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?

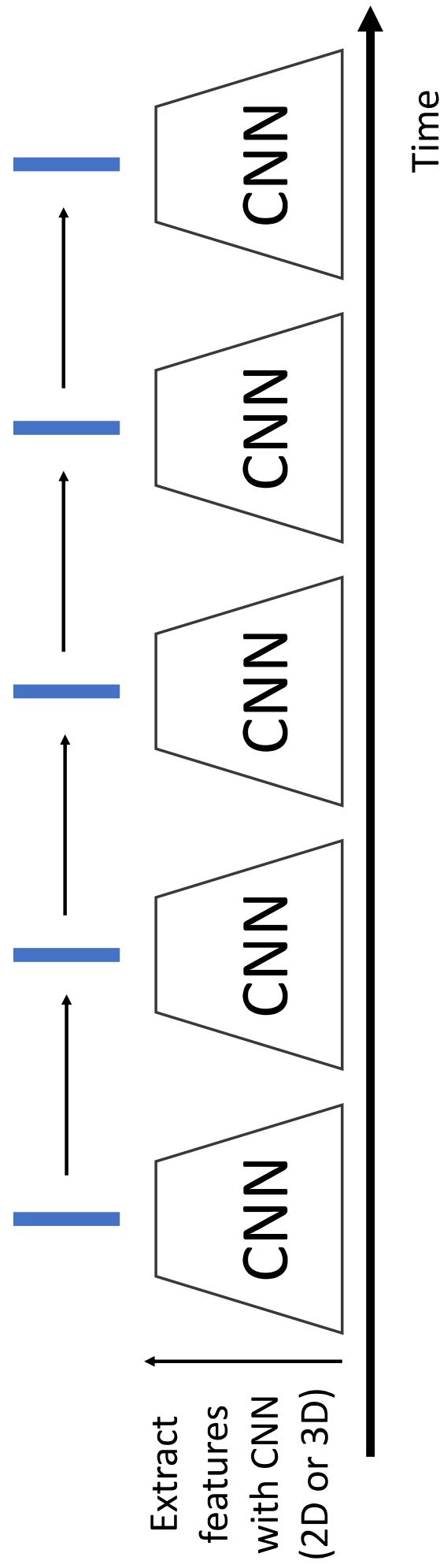


# Modeling long-term temporal structure



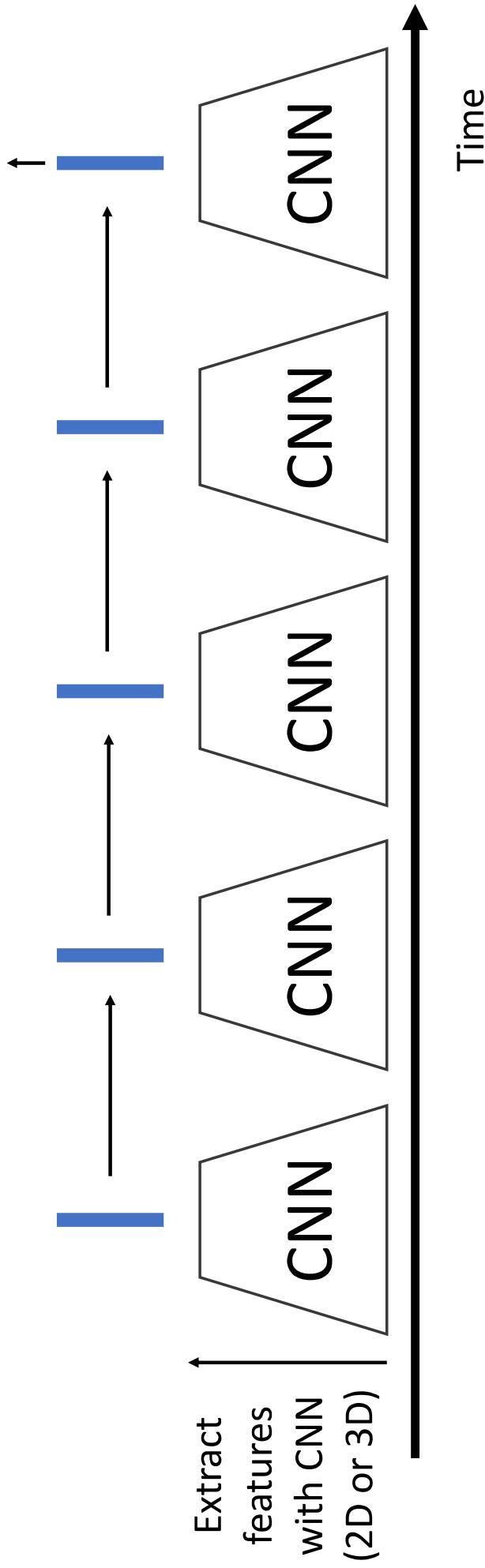
# Modeling long-term temporal structure

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



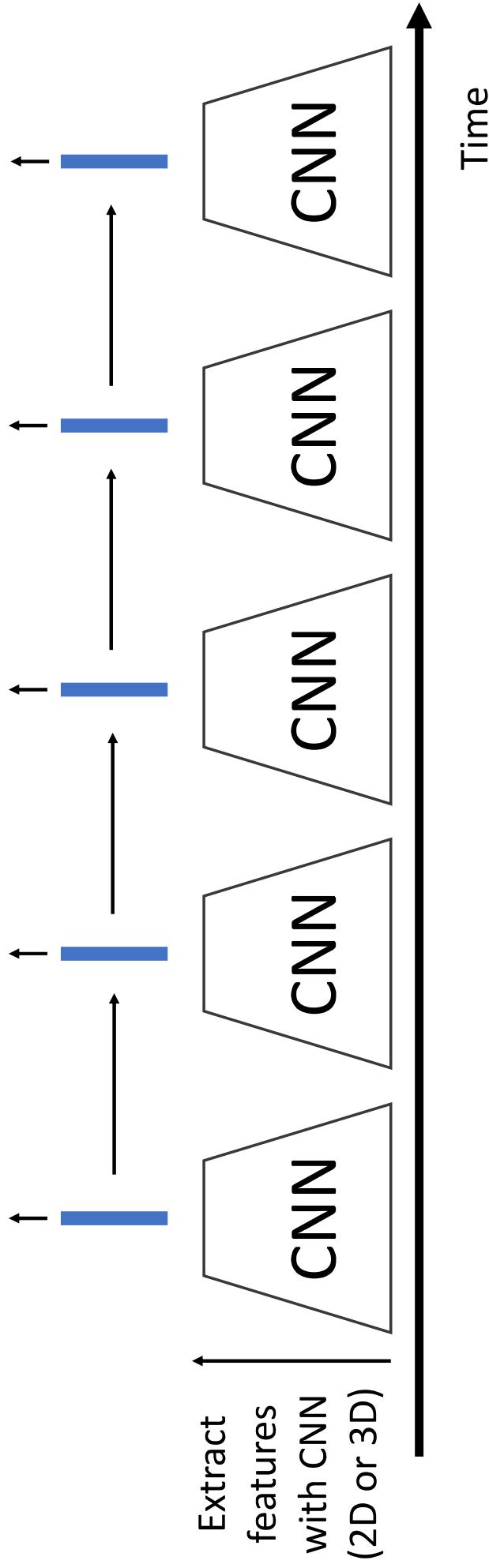
# Modeling long-term temporal structure

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



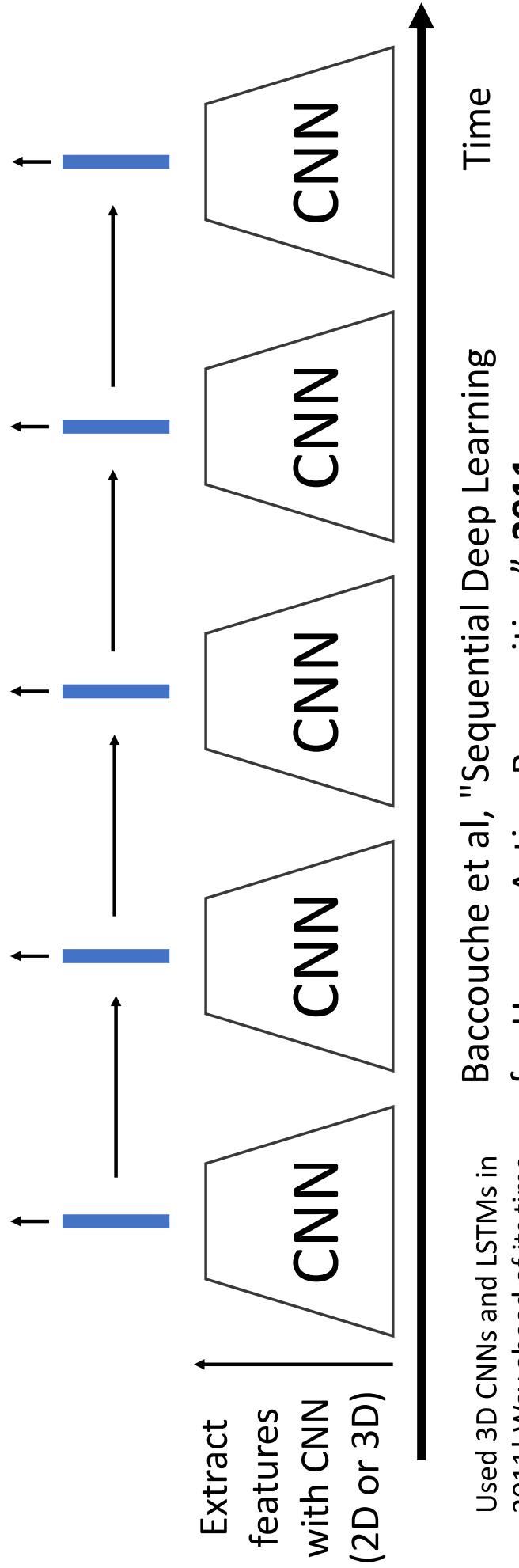
# Modeling long-term temporal structure

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



# Modeling long-term temporal structure

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

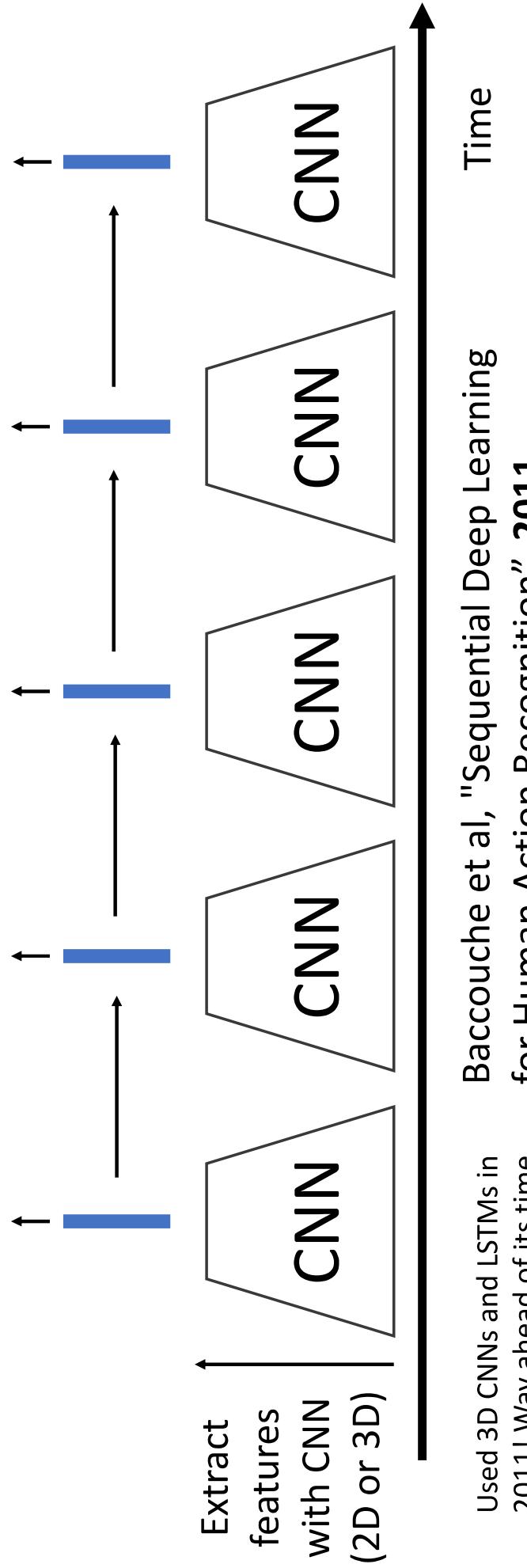


Used 3D CNNs and LSTMs in  
2011! Way ahead of its time

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

# Modeling long-term temporal structure

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

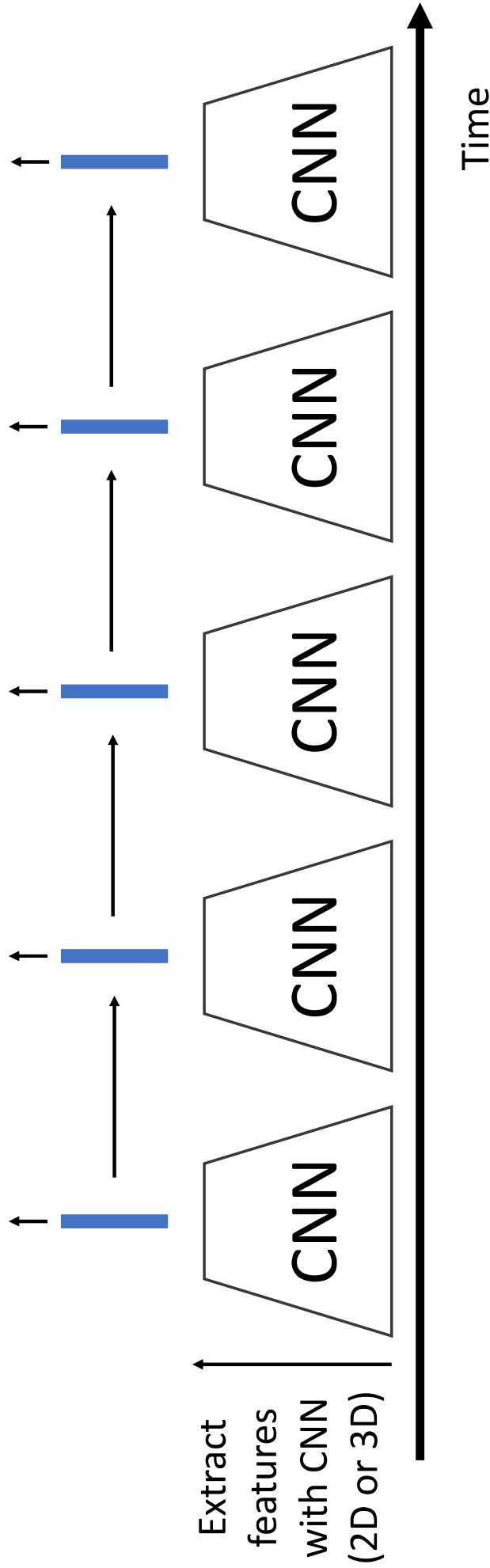


Used 3D CNNs and LSTMs in  
2011! Way ahead of its time

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

# Modeling long-term temporal structure

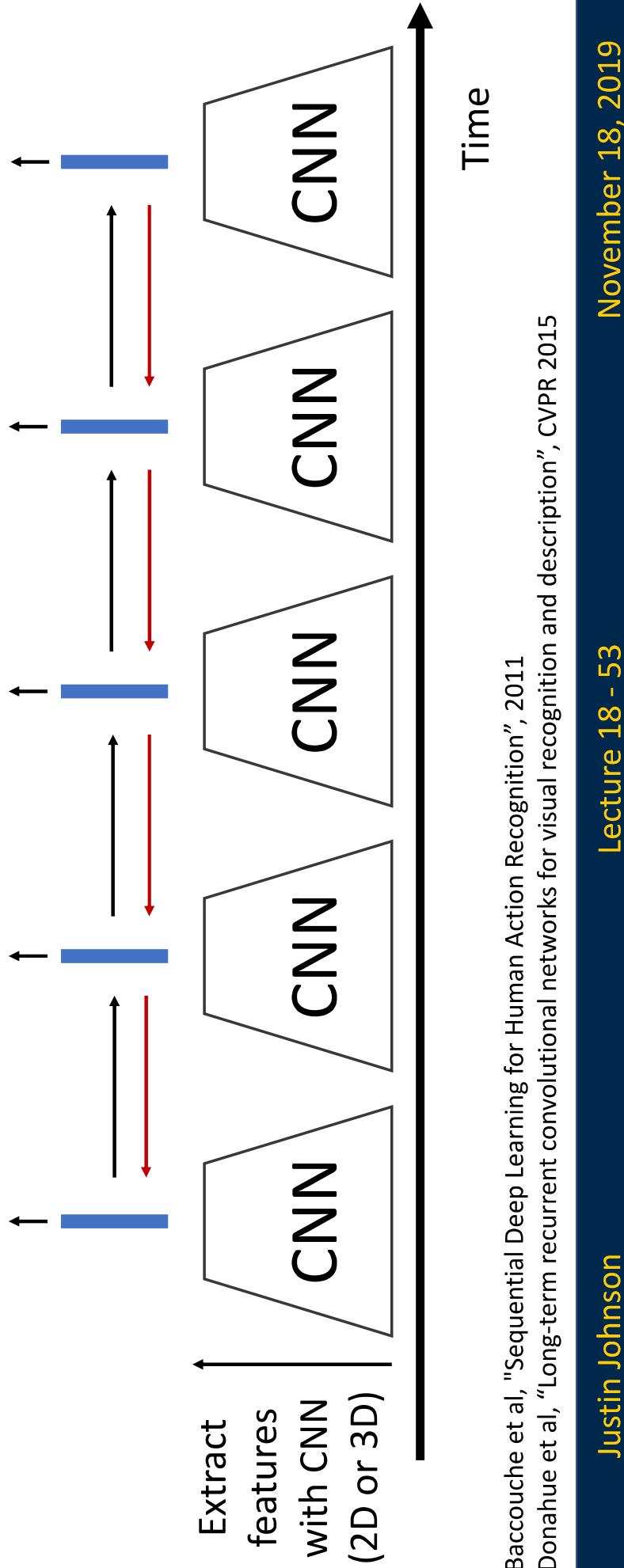
Process local features using recurrent network (e.g. LSTM)  
Many to many: one output per video frame



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

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

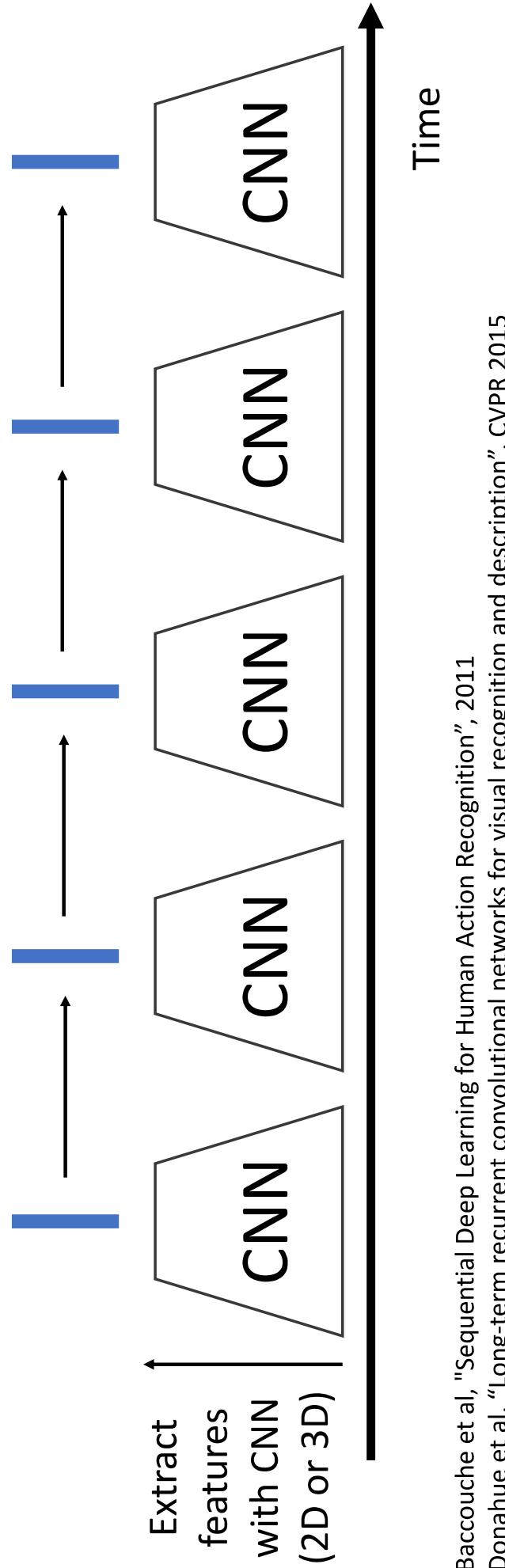


Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011  
Donahue et al, "Long-term recurrent convolutional networks for visual recognition and description", CVPR 2015

# Modeling long-term temporal structure

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

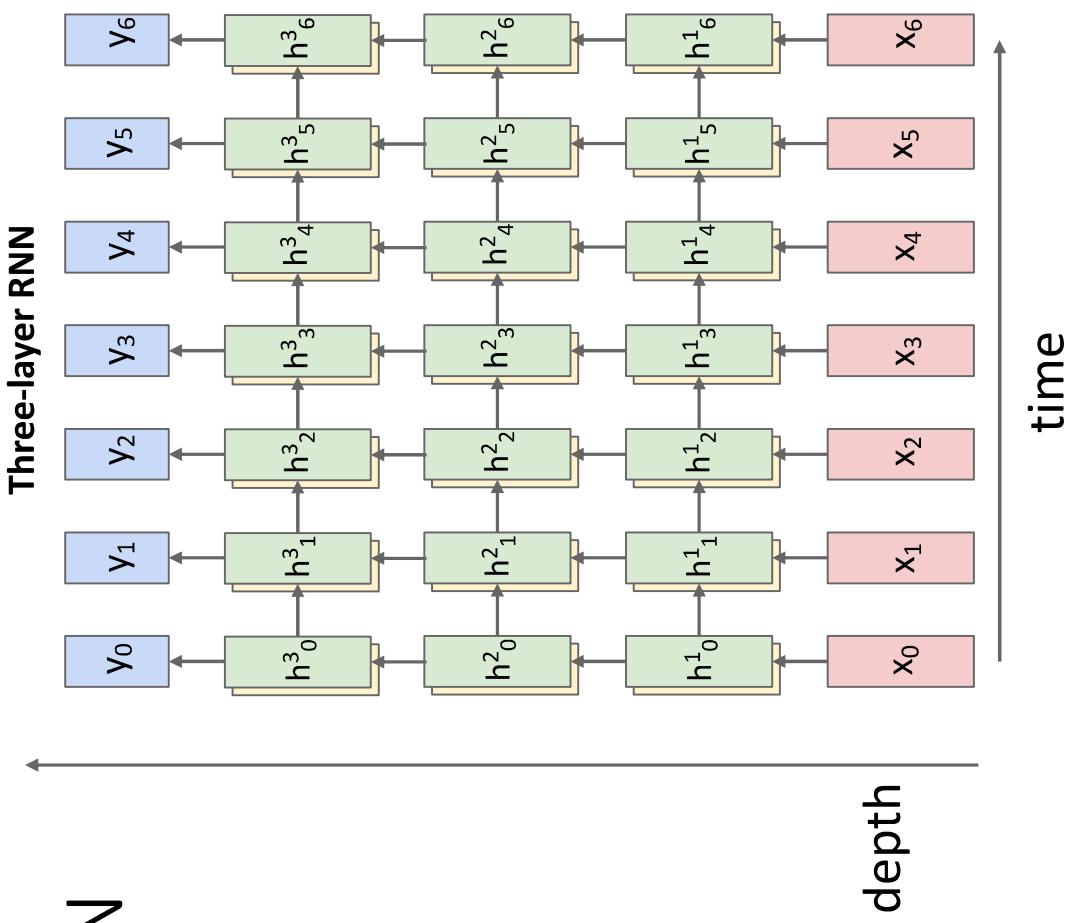
Inside RNN: Each vector is a function of all previous vectors (global temporal structure)  
Can we merge both approaches?



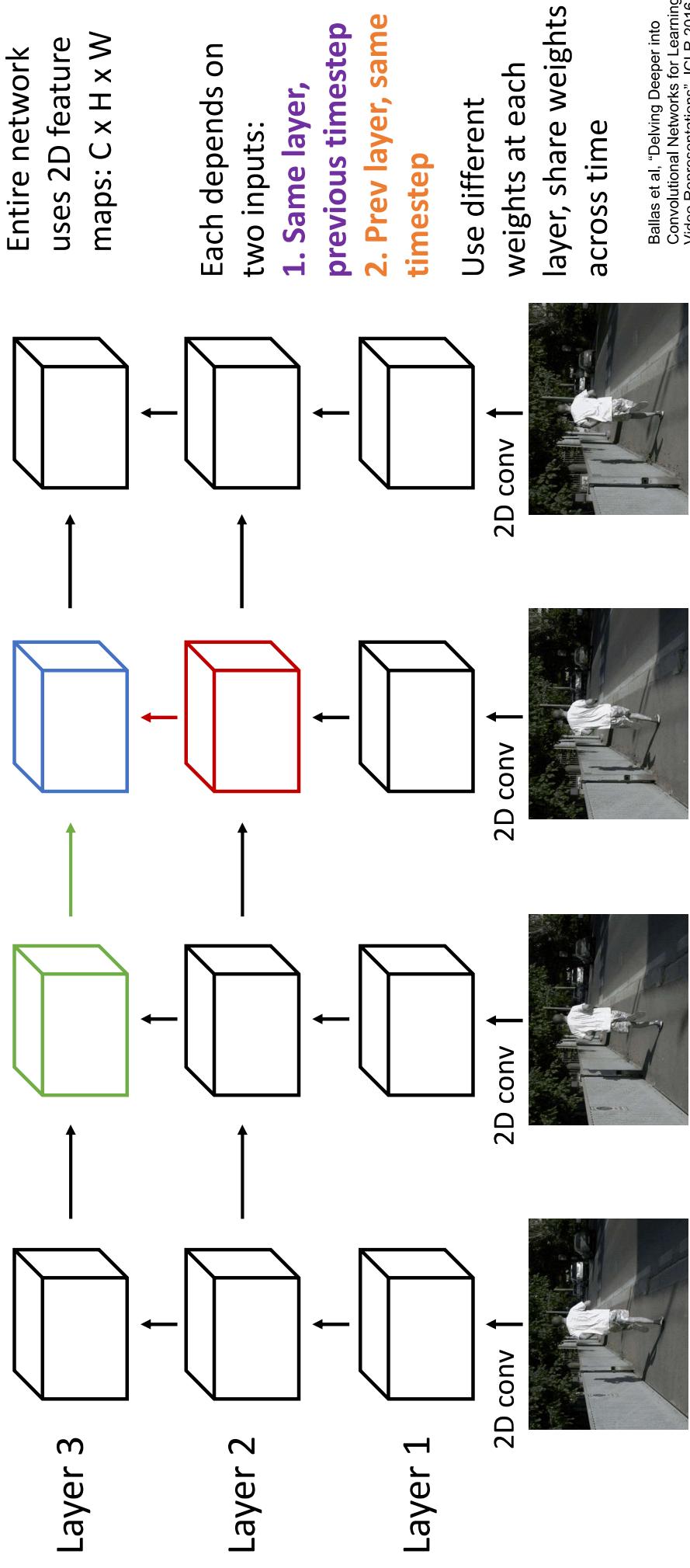
Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011  
Donahue et al, "Long-term recurrent convolutional networks for visual recognition and description", CVPR 2015

## Recall: Multi-layer RNN

We can use a  
similar structure to  
process videos!



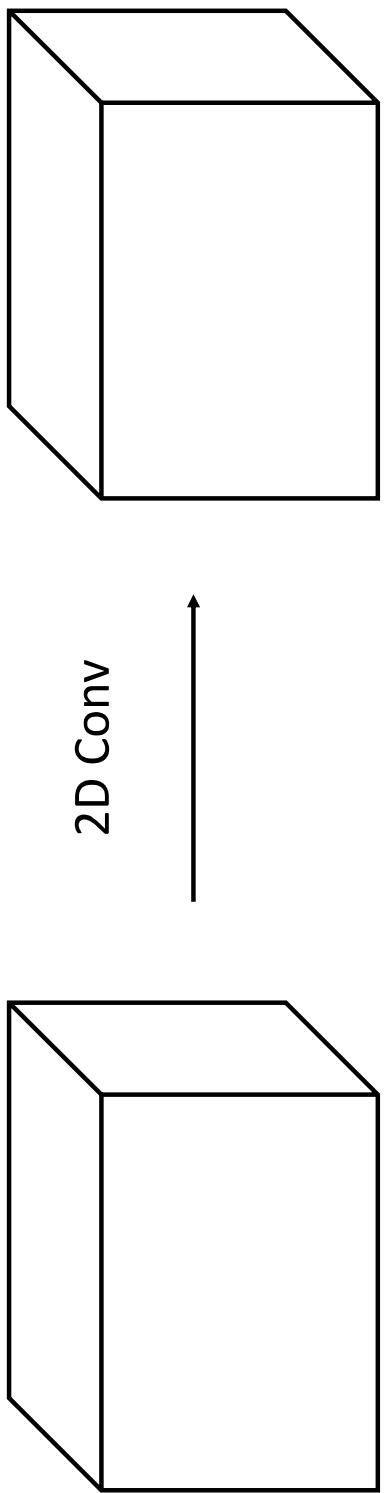
# Recurrent Convolutional Network



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

# Recurrent Convolutional Network

Normal 2D CNN:



Input features:  
 $C \times H \times W$

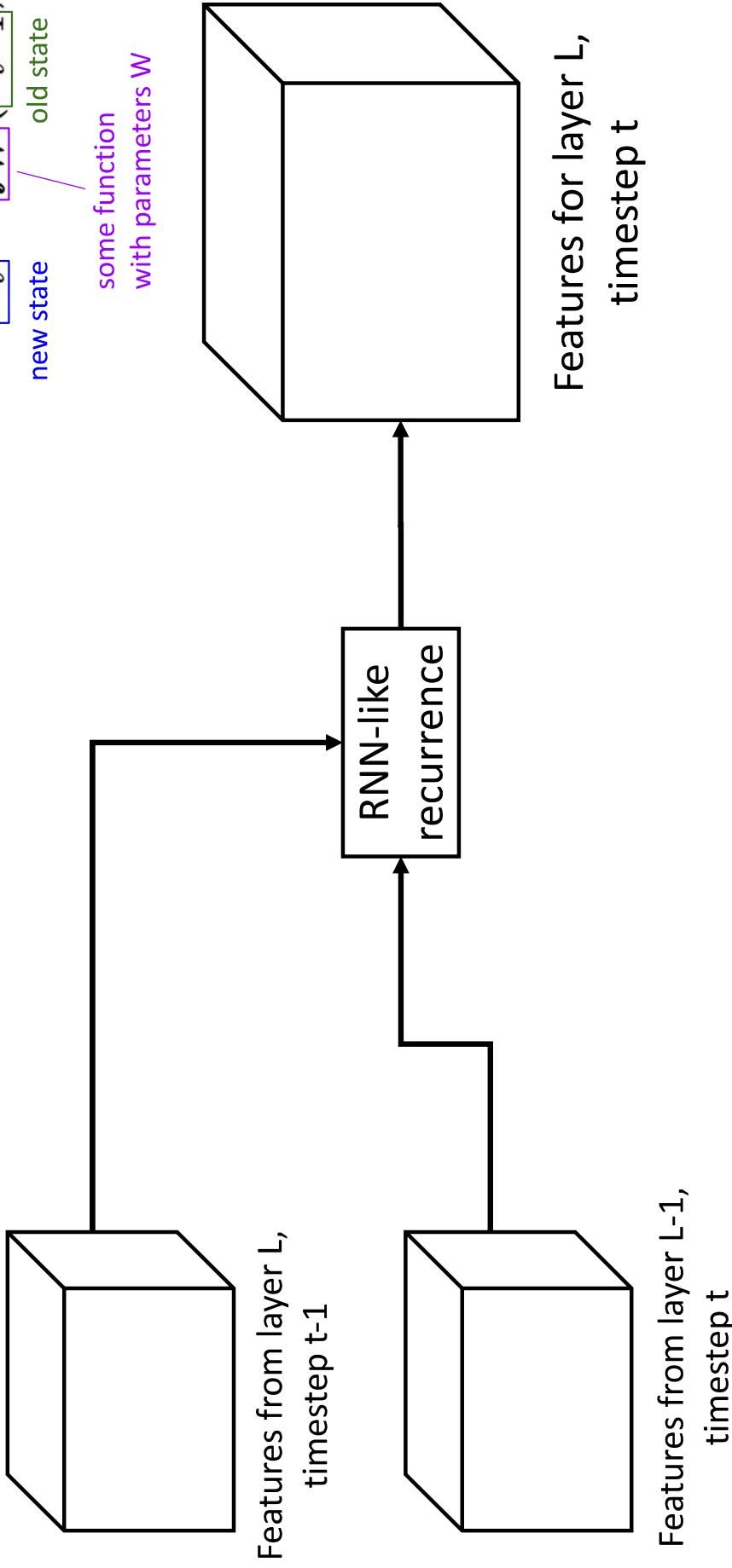
Output features:  
 $C \times H \times W$

# Recurrent Convolutional Network

Recall: Recurrent Network

$$h_t = f_W(h_{t-1}, x_t)$$

new state  
old state  
some function  
with parameters W



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

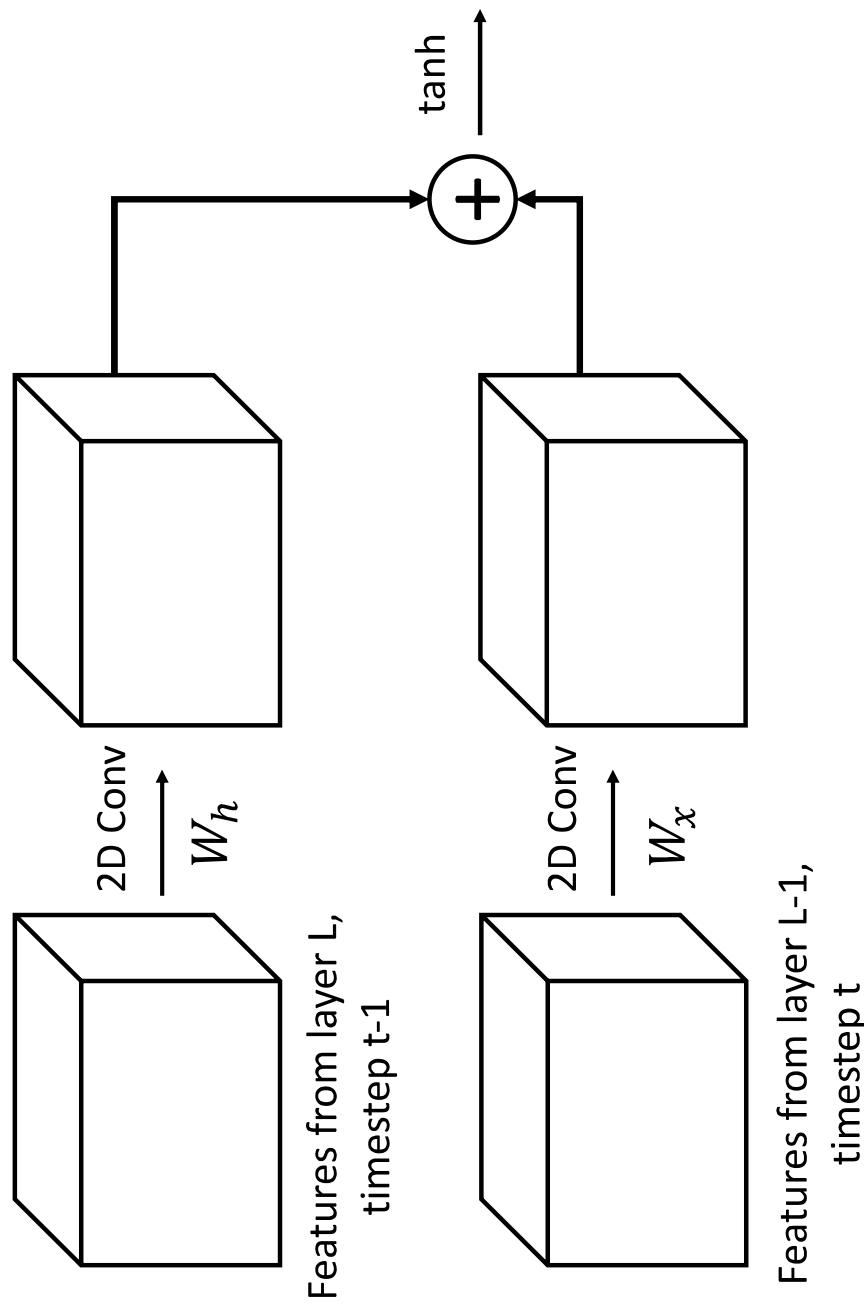
Justin Johnson

Lecture 18 - 58

November 18, 2019

## Recurrent Convolutional Network

Recall: Vanilla RNN  
$$h_{t+1} = \tanh(W_h h_t + W_x x)$$
  
Replace all matrix multiply  
with 2D convolution!



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

Justin Johnson

Lecture 18 - 59

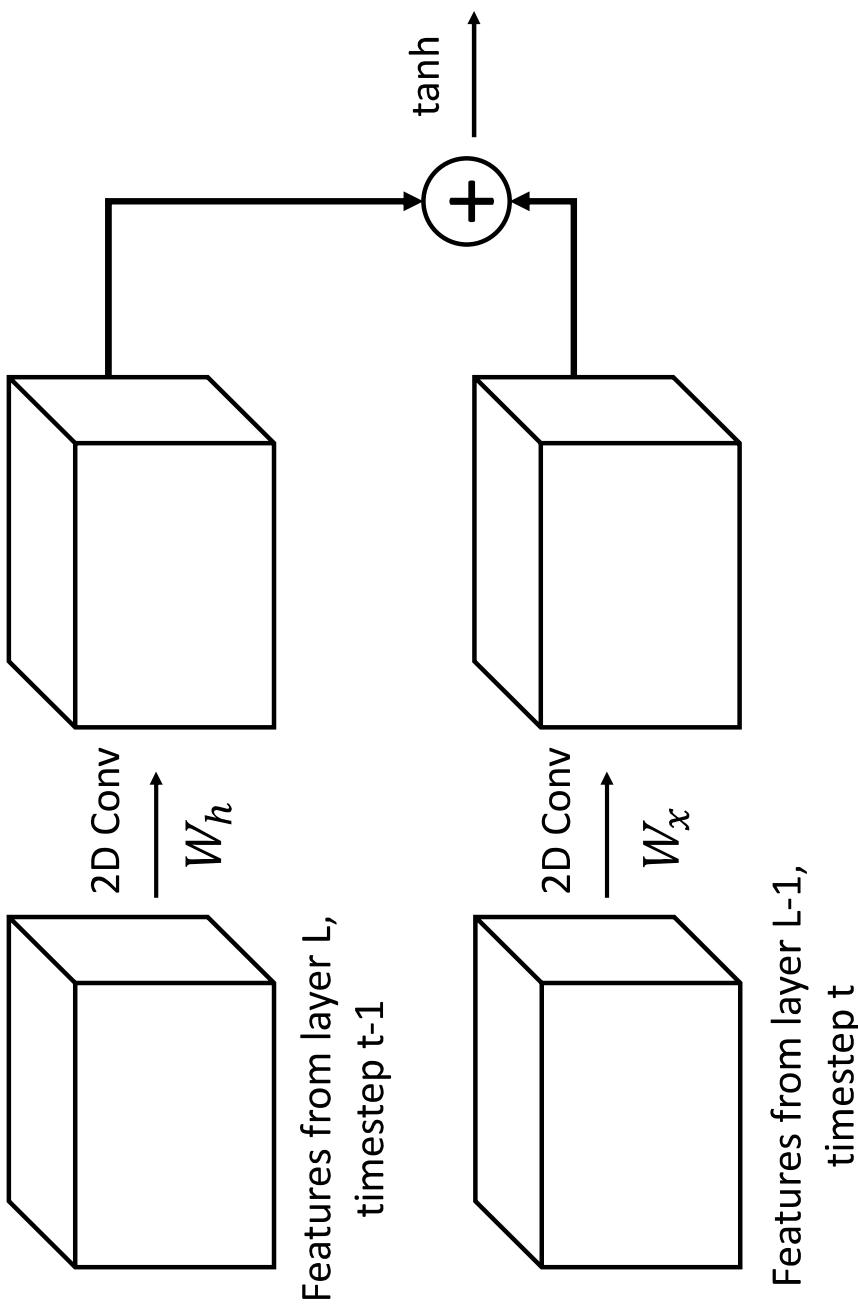
November 18, 2019

# Recurrent Convolutional Network

## Recall: GRU

$$\begin{aligned}r_t &= \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r) \\z_t &= \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z) \\ \tilde{h}_t &= \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h) \\ h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t\end{aligned}$$

Can do similar transform for other RNN variants (GRU, LSTM)



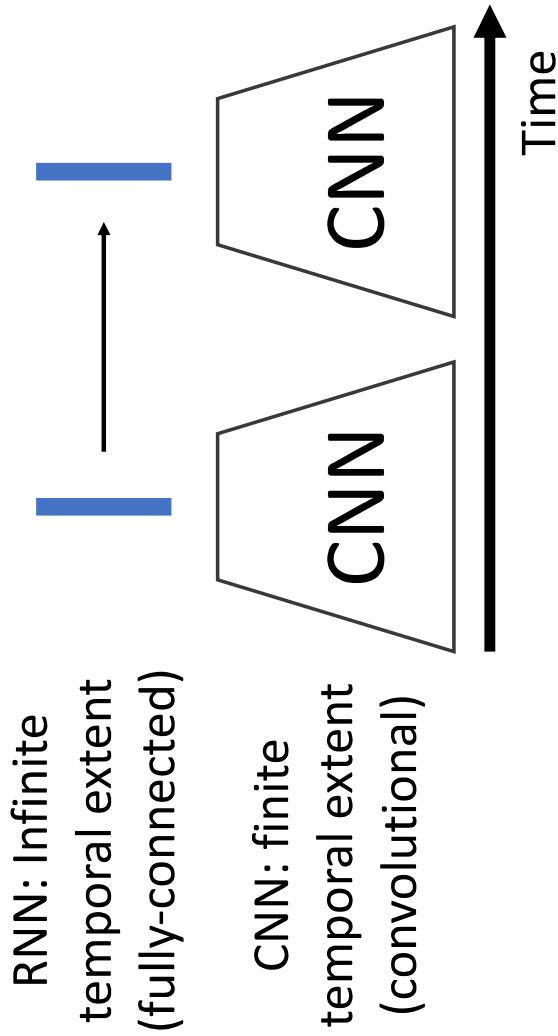
Ballas et al, "Delving Deeper into Convolutional Networks for Learning Video Representations", ICLR 2016

Justin Johnson

Lecture 18 - 60

November 18, 2019

# Modeling long-term temporal structure



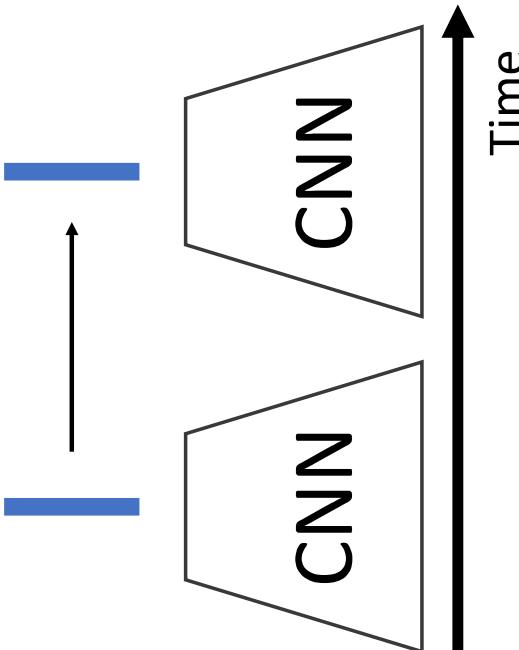
Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011  
Donahue et al, "Long-term recurrent convolutional networks for visual recognition and description", CVPR 2015

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

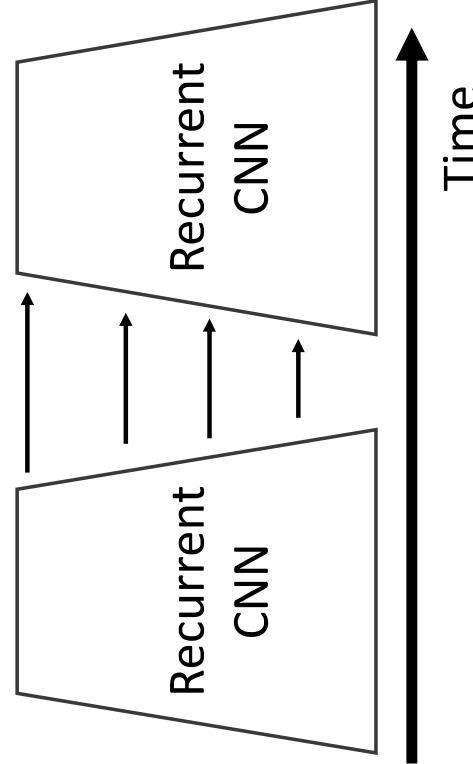
# Modeling long-term temporal structure

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

RNN: Infinite temporal extent (fully-connected)



Recurrent CNN: Infinite temporal extent (convolutional)

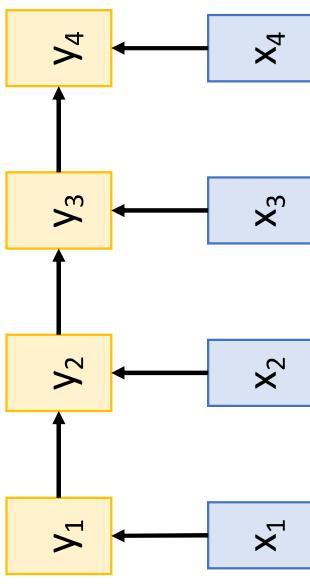


Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011  
Donahue et al, "Long-term recurrent convolutional networks for visual recognition and description", CVPR 2015

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

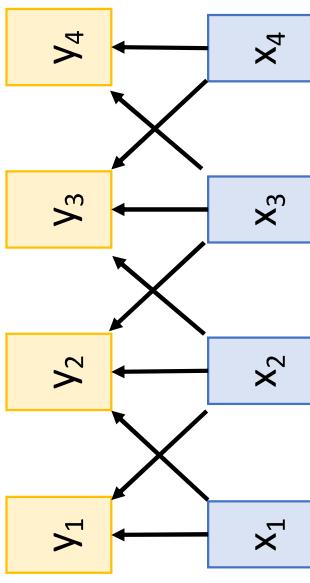
# Recall: Different ways of processing sequences

## Recurrent Neural Network



Works on **Ordered Sequences**  
(+)**Good at long sequences:** After one RNN layer,  $h_T$  "sees" the whole sequence  
(-)**Not parallelizable:** need to compute hidden states sequentially  
In video: CNN+RNN, or recurrent CNN

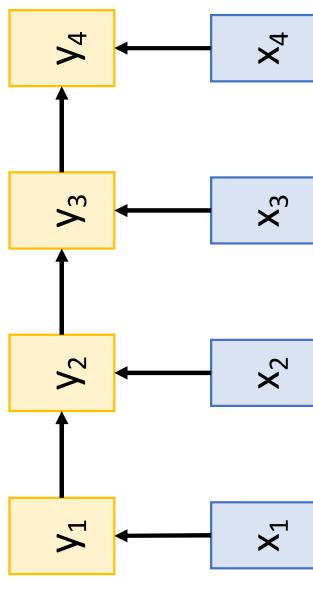
## 1D Convolution



Works on **Multidimensional Grids**  
(-)**Bad at long sequences:** Need to stack many conv layers for outputs to "see" the whole sequence  
(+)**Highly parallel:** Each output can be computed in parallel  
In video: 3D convolution

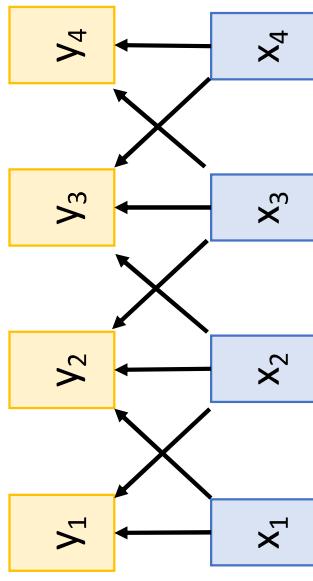
# Recall: Different ways of processing sequences

## Recurrent Neural Network



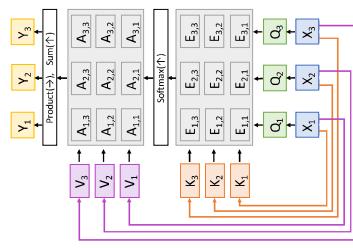
Works on **Ordered Sequences**  
(+ ) **Good at long sequences:** After one RNN layer,  $h_T$  "sees" the whole sequence  
(- ) **Not parallelizable:** need to compute **hidden states sequentially**  
In video: CNN+RNN, or recurrent CNN

## 1D Convolution



Works on **Multidimensional Grids**  
(- ) **Bad at long sequences:** Need to stack many conv layers for outputs to "see" the whole sequence  
(+ ) **Highly parallel:** Each output can be computed in parallel  
In video: 3D convolution

## Self-Attention



Works on **Sets of Vectors**  
(- ) **Good at long sequences:** after one self-attention layer, each output "sees" all inputs!  
(+) Highly parallel: Each output can be computed in parallel  
(-) **Very memory intensive**  
In video: ????

## Recall: Self-Attention

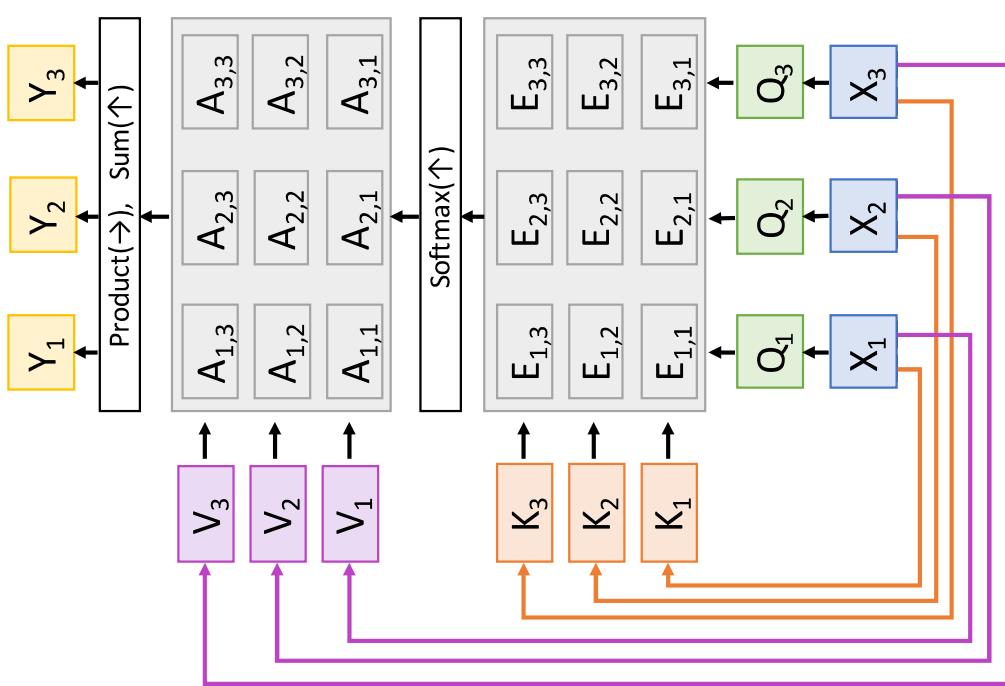
**Input:** Set of vectors  $X_1, \dots, X_N$

**Keys, Queries, Values:** Project each  $x$  to a key, query, and value using linear layer

**Affinity matrix:** Compare each pair of  $x$ , (using scaled dot-product between keys and values) and normalize using softmax

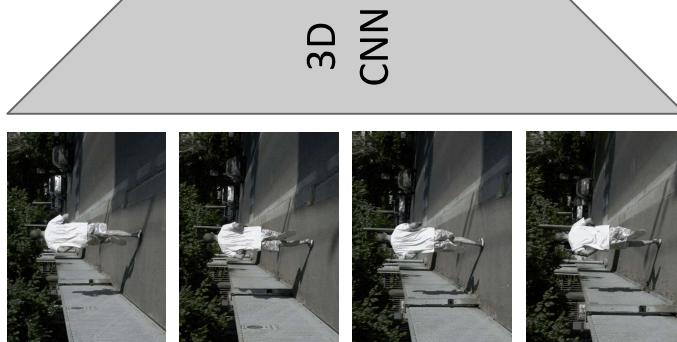
**Output:** Weighted sum of values, with weights given by affinity matrix

Features in 3D CNN:  $C \times T \times H \times W$   
Interpret as a set of THW vectors of dim  $C$



# Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



3D  
CNN

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

Nonlocal Block

Wang et al, "Non-local neural networks", CVPR 2018

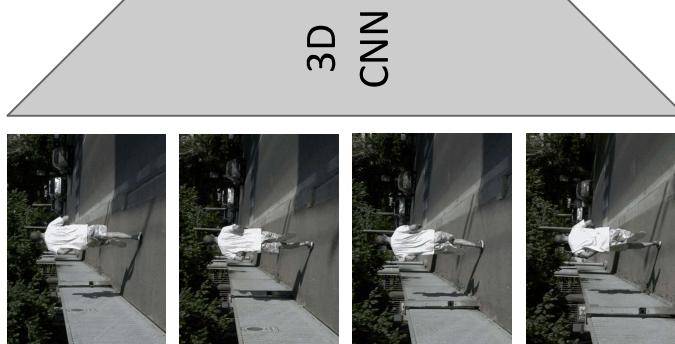
Justin Johnson

Lecture 18 - 66

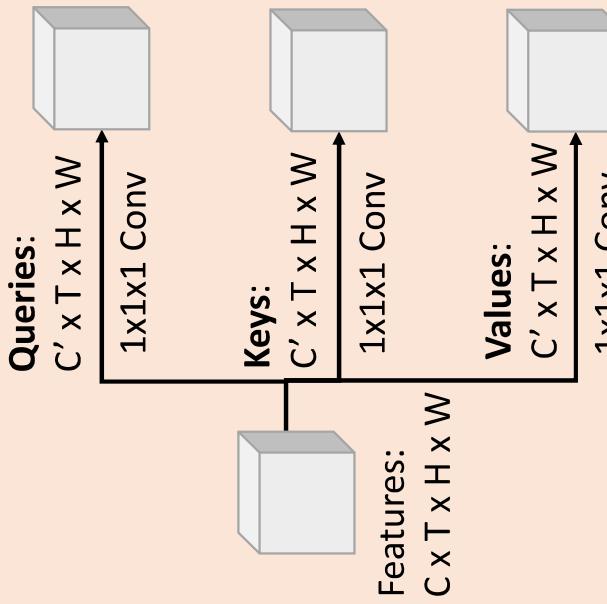
November 18, 2019

# Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



3D CNN



Nonlocal Block

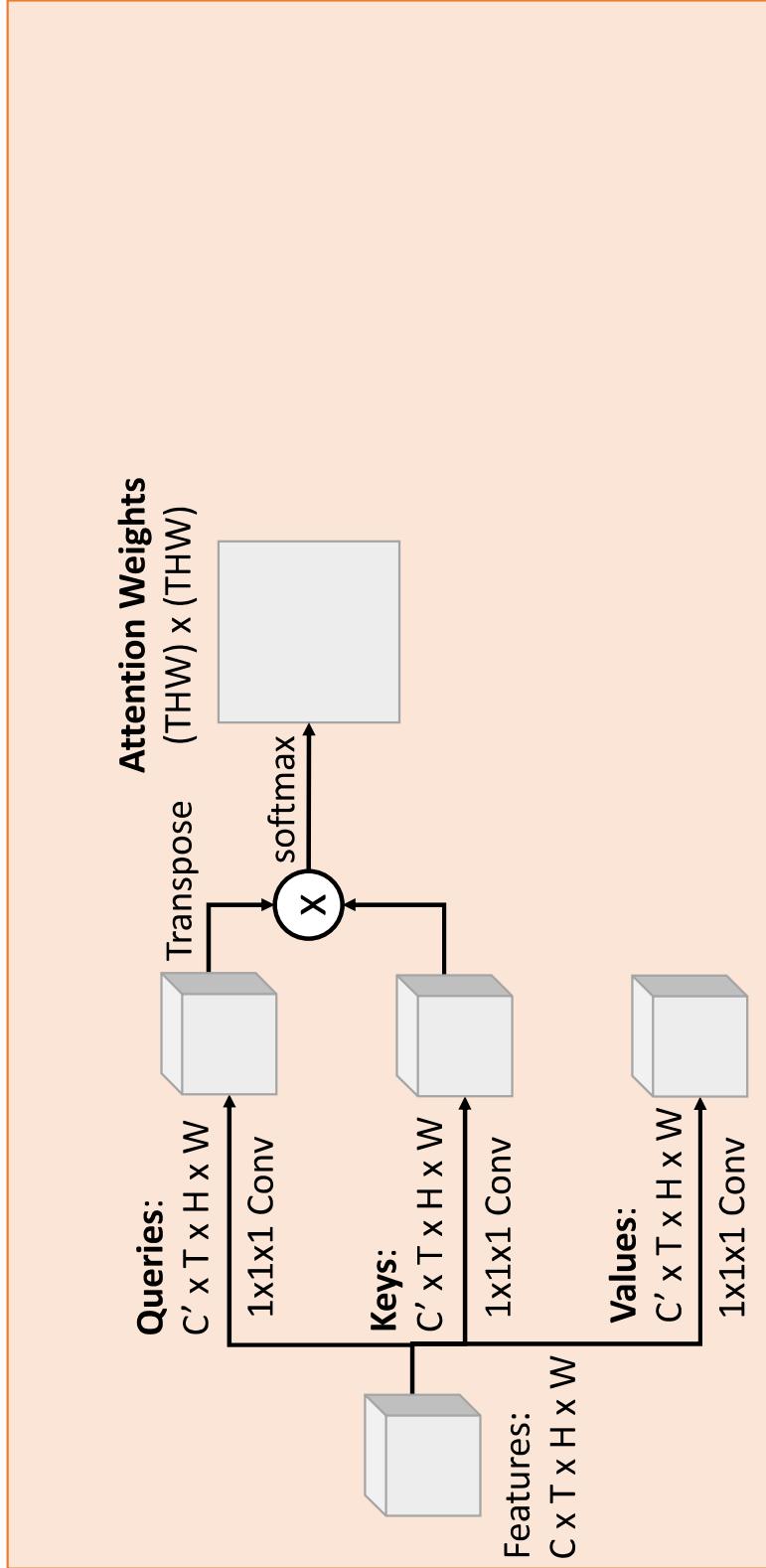
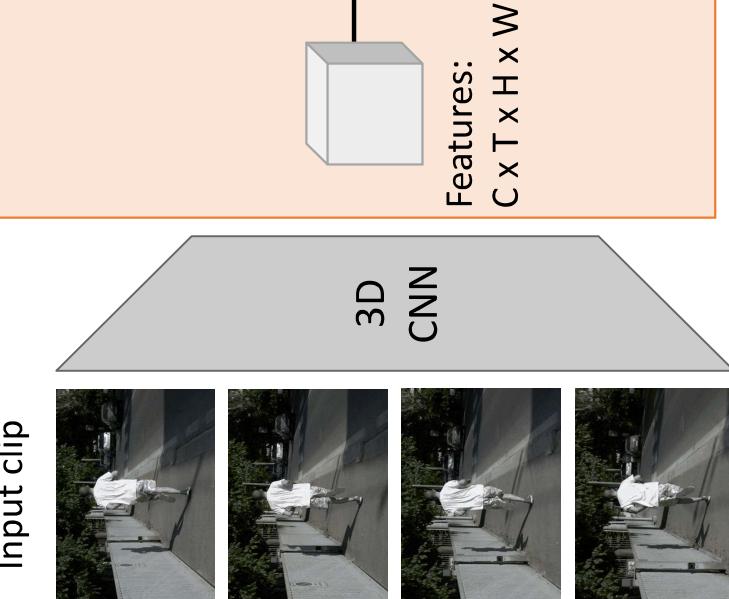
Wang et al., "Non-local neural networks", CVPR 2018

Justin Johnson

Lecture 18 - 67

November 18, 2019

# Spatio-Temporal Self-Attention (Nonlocal Block)



Nonlocal Block

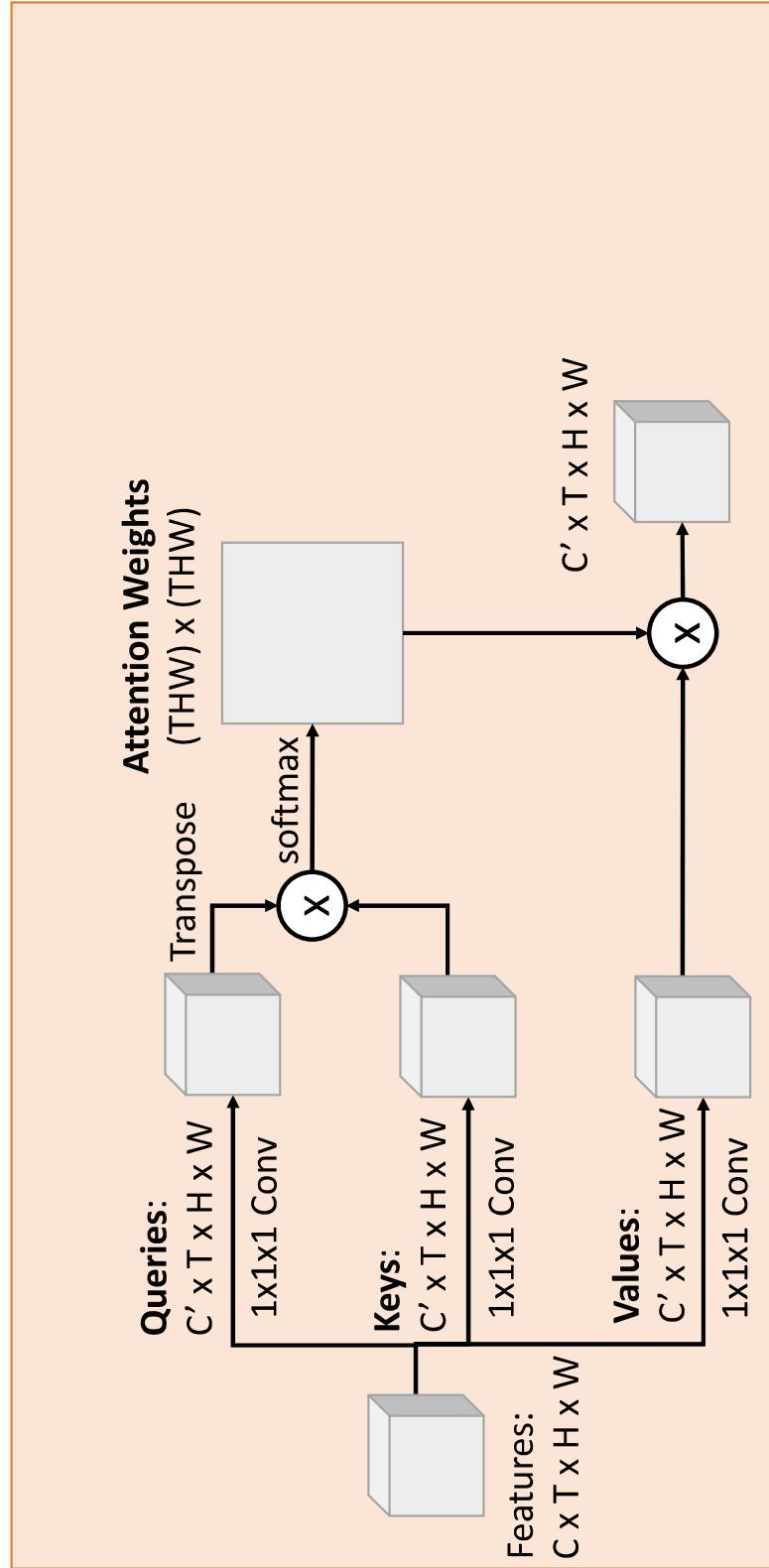
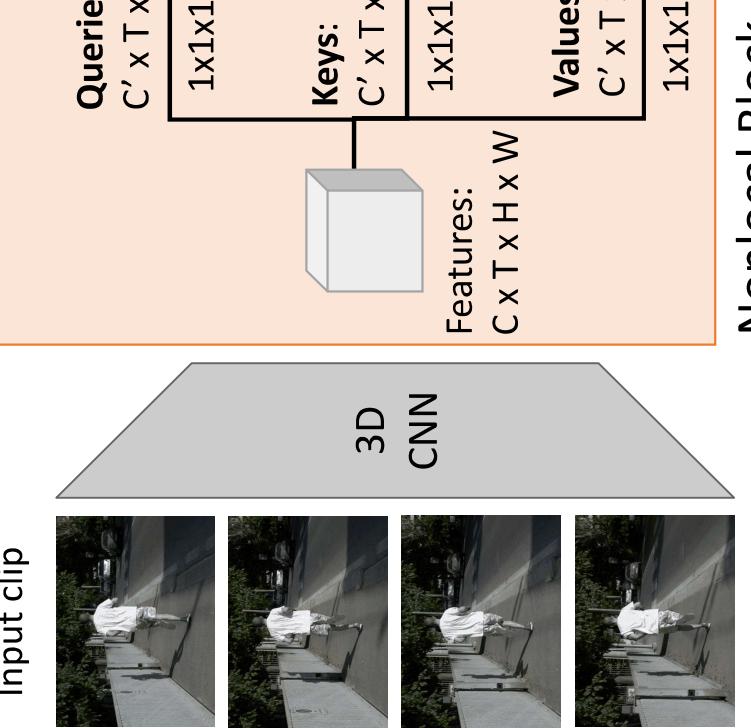
Wang et al., "Non-local neural networks", CVPR 2018

Justin Johnson

Lecture 18 - 68

November 18, 2019

# Spatio-Temporal Self-Attention (Nonlocal Block)



Nonlocal Block

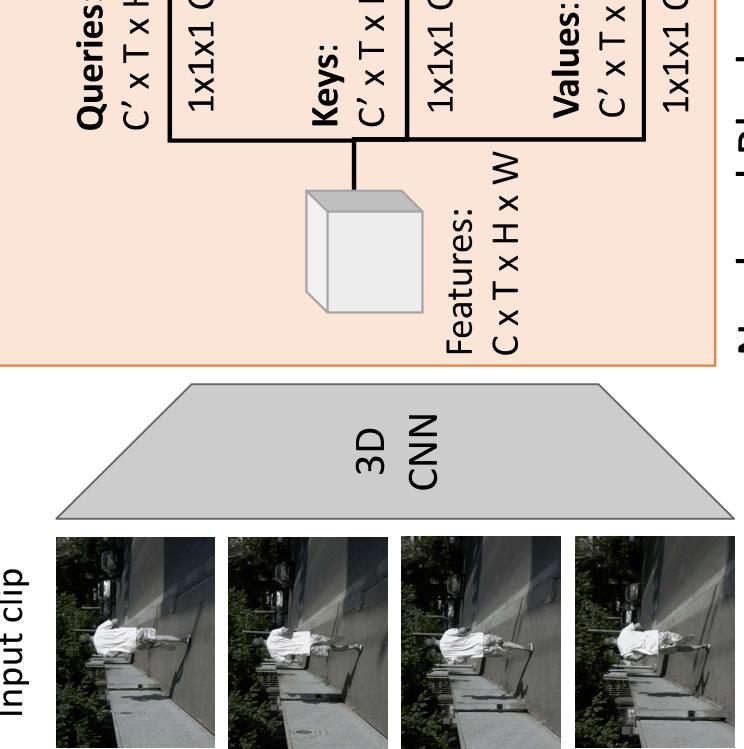
Wang et al., "Non-local neural networks", CVPR 2018

Justin Johnson

Lecture 18 - 69

November 18, 2019

# Spatio-Temporal Self-Attention (Nonlocal Block)



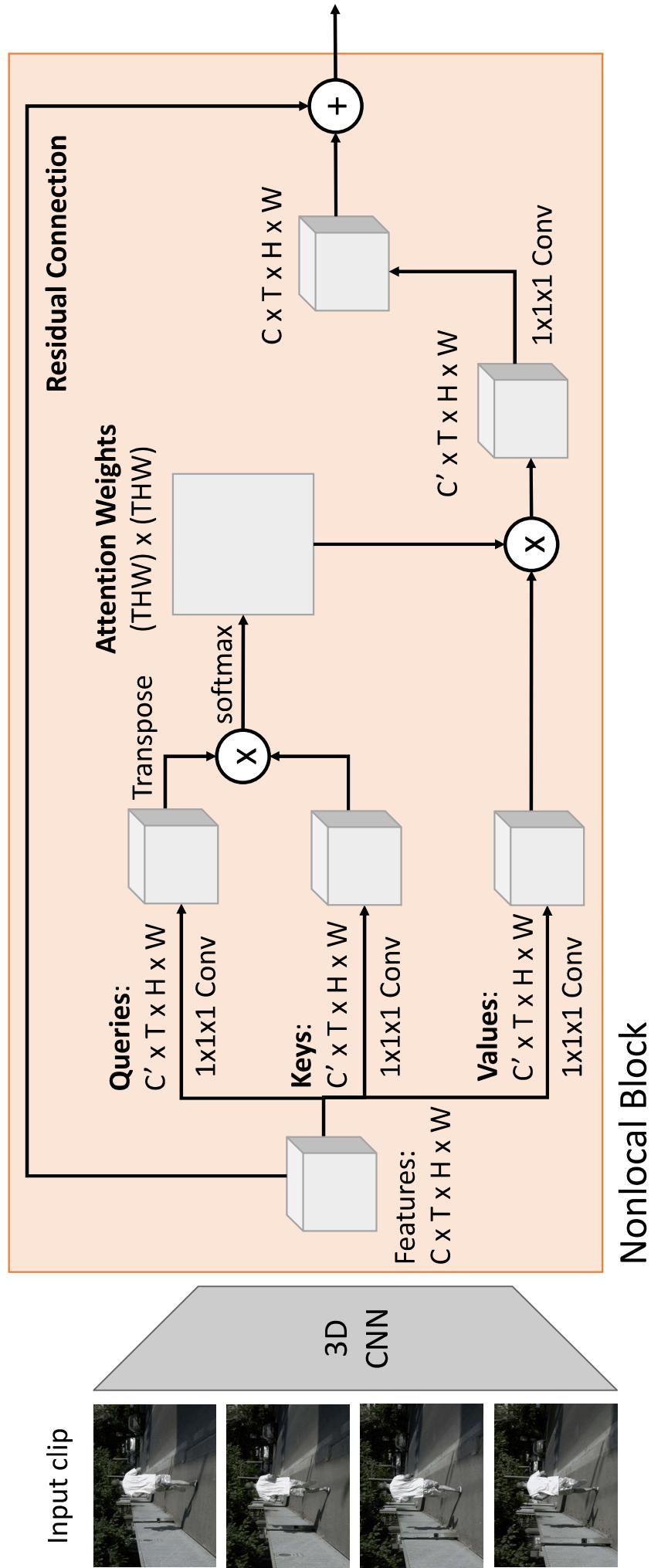
Wang et al., "Non-local neural networks", CVPR 2018

Justin Johnson

Lecture 18 - 70

November 18, 2019

# Spatio-Temporal Self-Attention (Nonlocal Block)



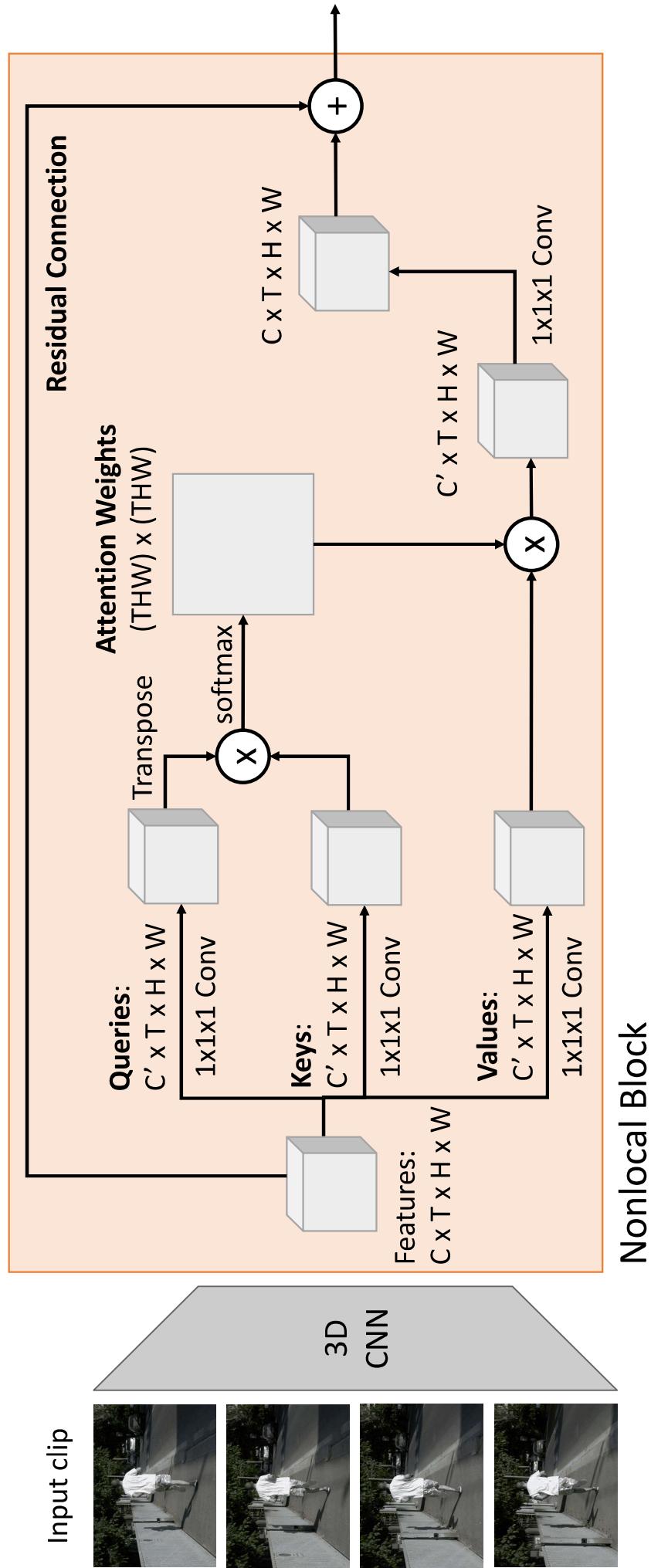
Wang et al., "Non-local neural networks", CVPR 2018

Justin Johnson

Lecture 18 - 71

November 18, 2019

# Spatio-Temporal Self-Attention (Nonlocal Block)



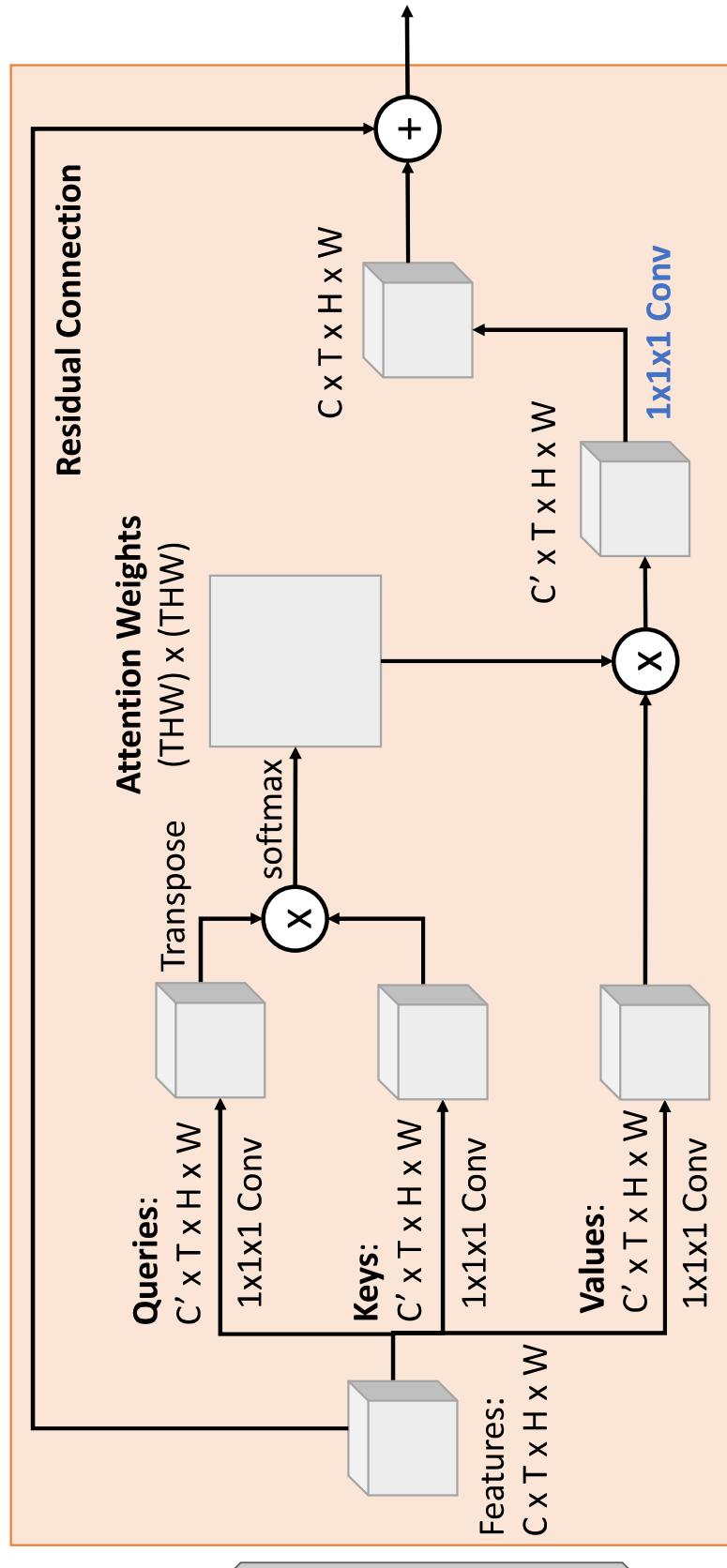
Wang et al., "Non-local neural networks", CVPR 2018

Justin Johnson

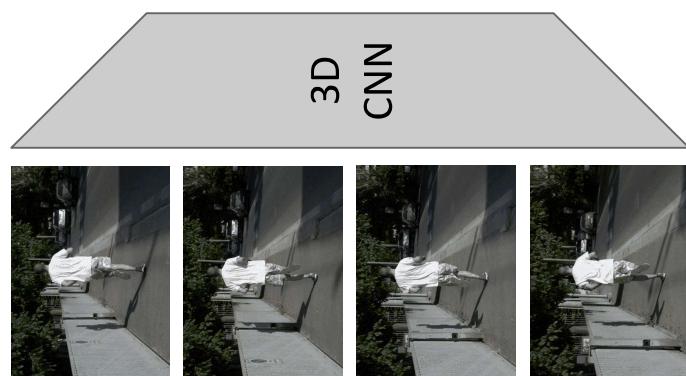
Lecture 18 - 72

November 18, 2019

# Spatio-Temporal Self-Attention (Nonlocal Block)



Input clip



3D  
CNN

**Nonlocal Block Trick:** Initialize **last conv** to 0, then entire block computes identity. Can insert into existing 3D CNNs

In practice, actually insert BatchNorm layer after final conv, and initialize scale parameter of BN layer to 0 rather than setting conv weight to 0

Wang et al, "Non-local neural networks", CVPR 2018

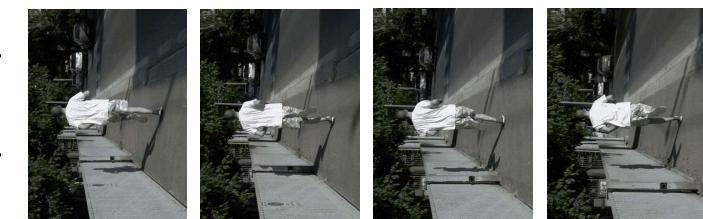
Justin Johnson

Lecture 18 - 73

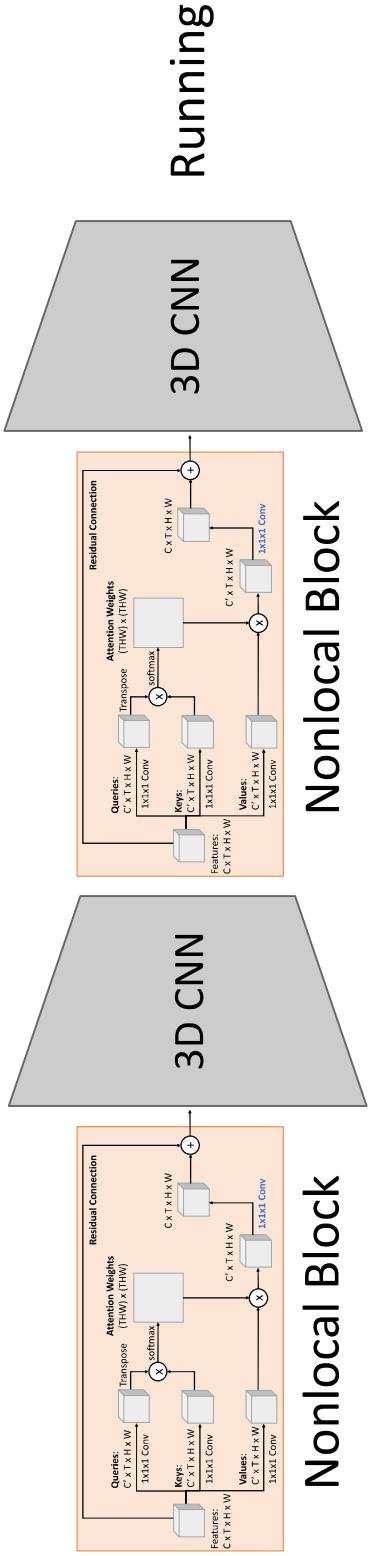
November 18, 2019

# 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?



# 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.

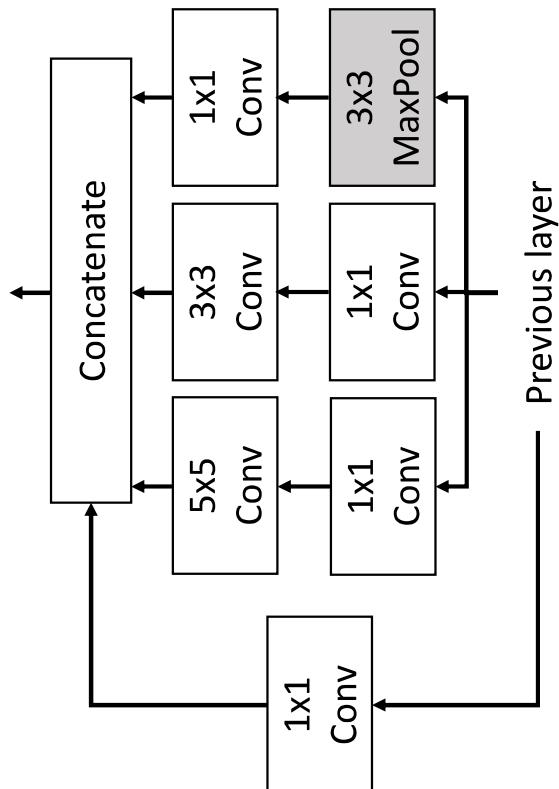
Replace each 2D  $K_h \times K_w$  conv/pool  
layer with a 3D  $K_t \times K_h \times K_w$  version

# 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.

Inception Block: Original



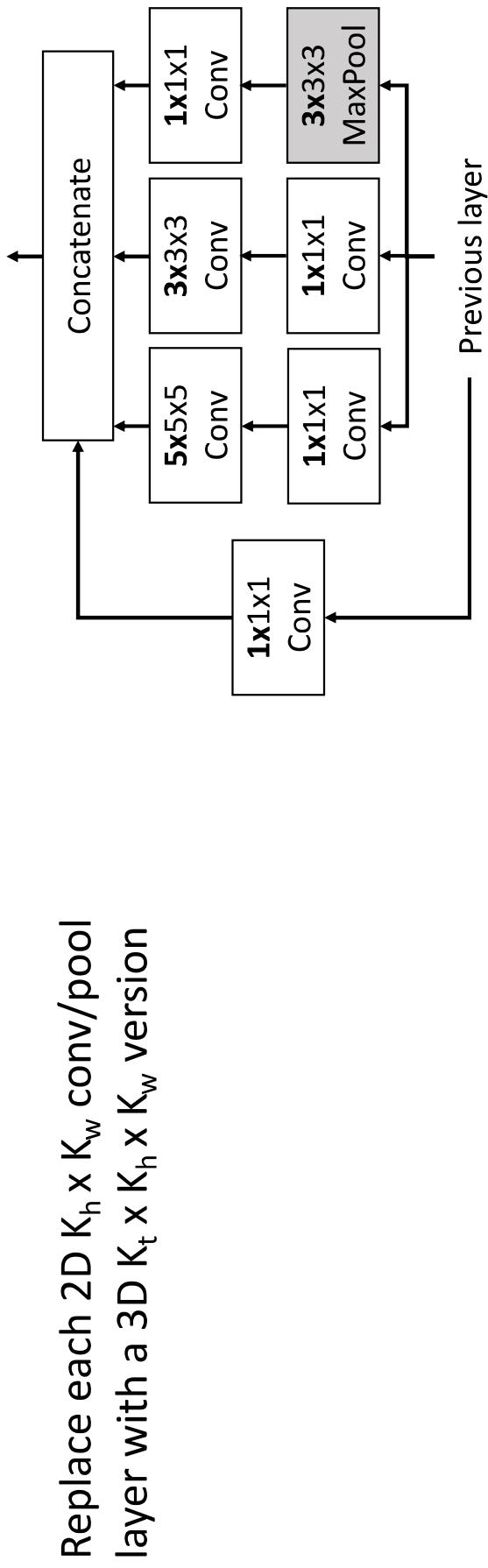
Replace each  $2D K_h \times K_w$  conv/pool  
layer with a  $3D K_t \times K_h \times K_w$  version

# 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.

Inception Block: Inflated



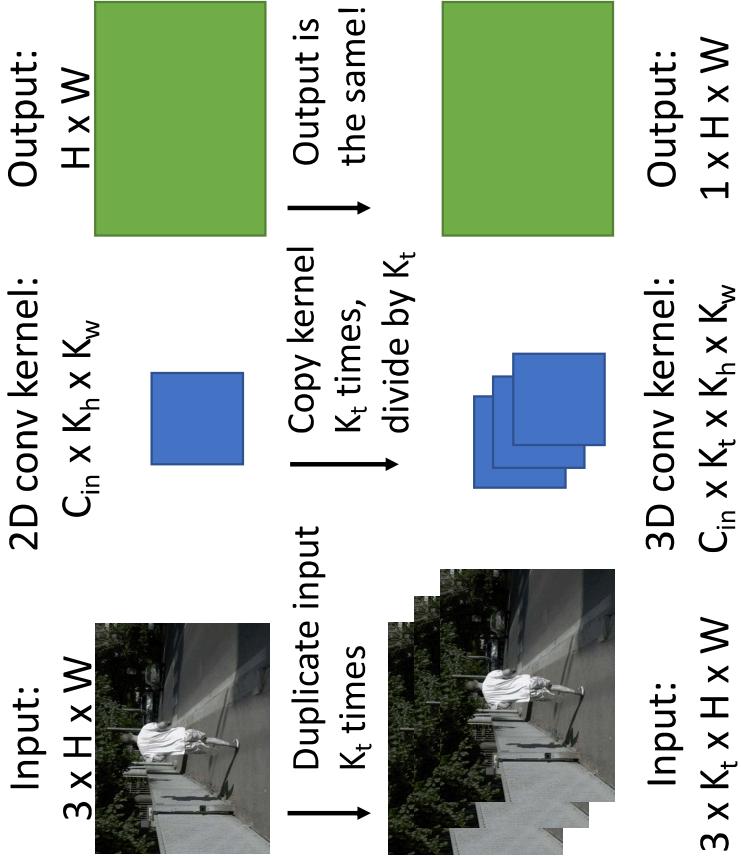
Replace each  $K_h \times K_w$  conv/pool  
layer with a  $3D K_t \times K_h \times K_w$  version

# 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  $2D \times K_h \times K_w$  conv/pool layer with a  $3D \times K_t \times K_h \times K_w$  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

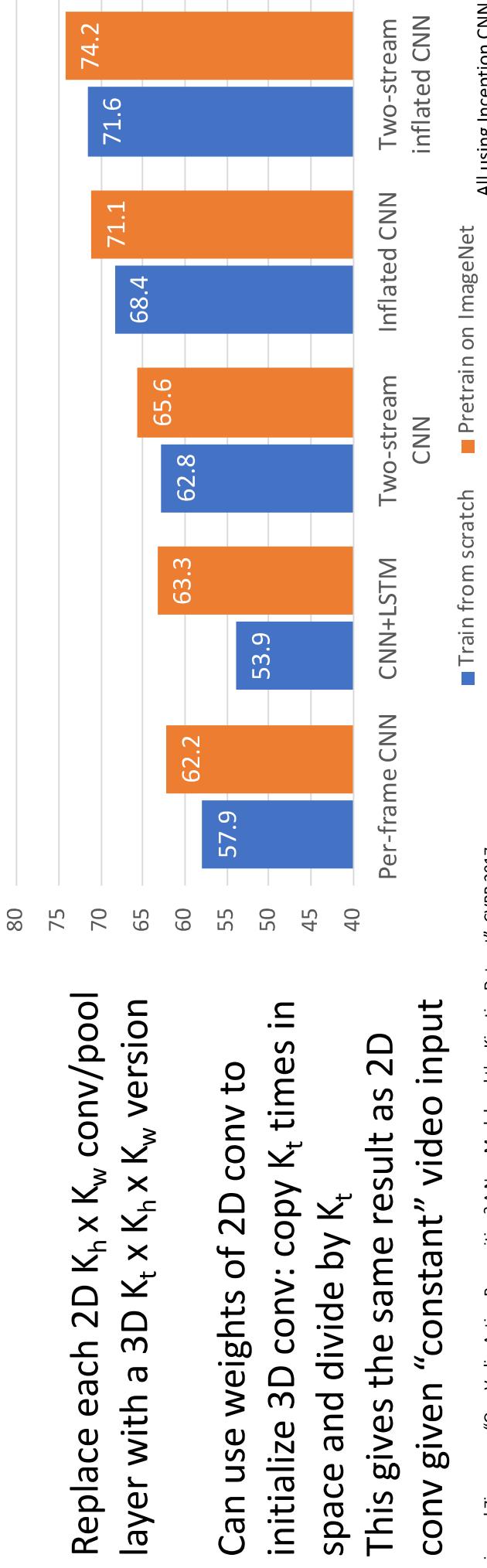
Carreira and Zisserman, “Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset”, CVPR 2017

# 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.

Top-1 Accuracy on Kinetics-400



Carreira and Zisserman, "Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset", CVPR 2017

Justin Johnson

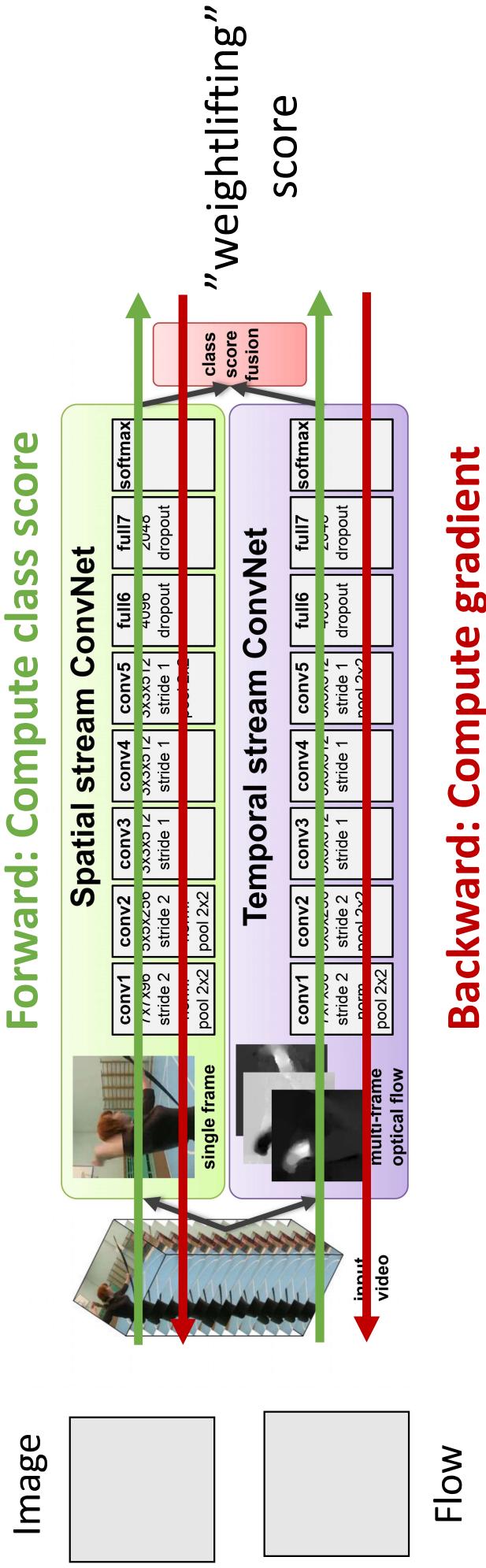
Lecture 18 - 79

November 18, 2019

All using Inception CNN

# Visualizing Video Models

Image



Flow

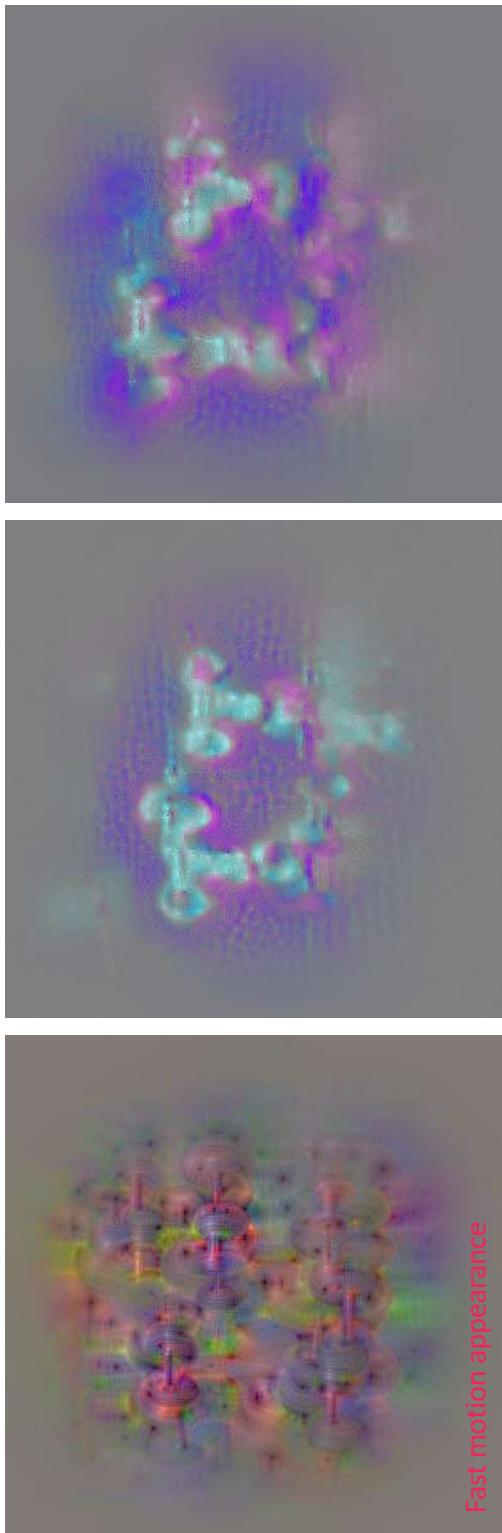
**Backward: Compute gradient**

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?

Appearance      “Slow” motion      “Fast” motion

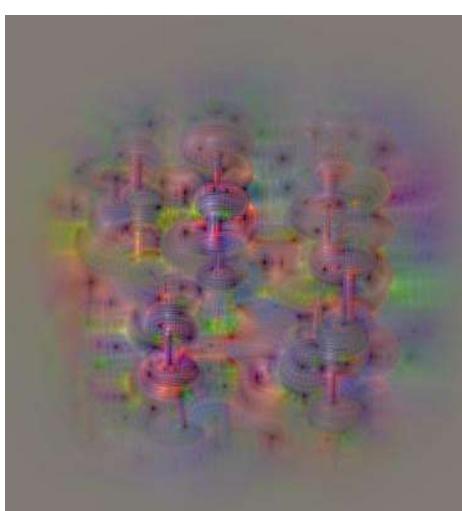


Feichtenhofer et al., “What have we learned from deep representations for action recognition?”, CVPR 2018  
Feichtenhofer et al., “Deep insights into convolutional networks for video recognition?”, IJCV 2019.  
Slide credit: Christoph Feichtenhofer

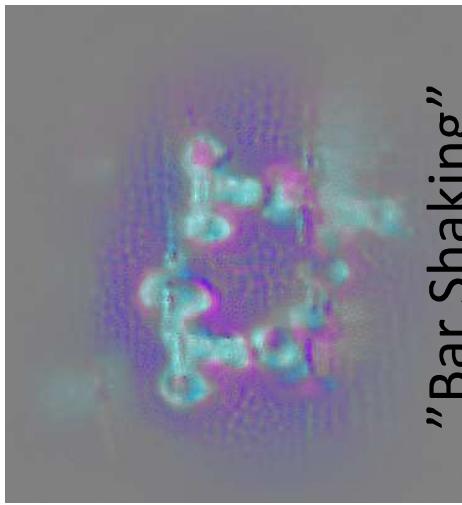
# Can you guess the action? Weightlifting

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 Feichtenhofer

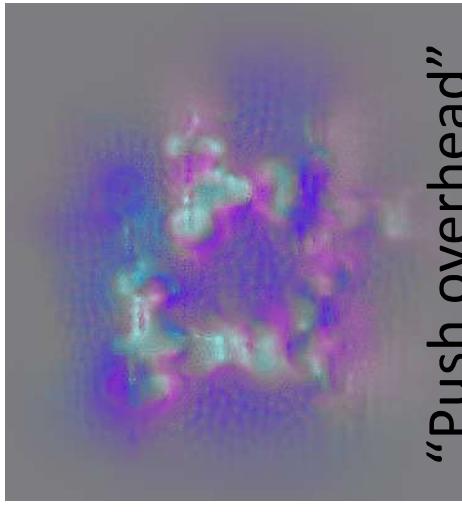
Appearance



"Slow" motion



"Fast" motion



Fast motion appearance

"Bar Shaking"



"Push overhead"

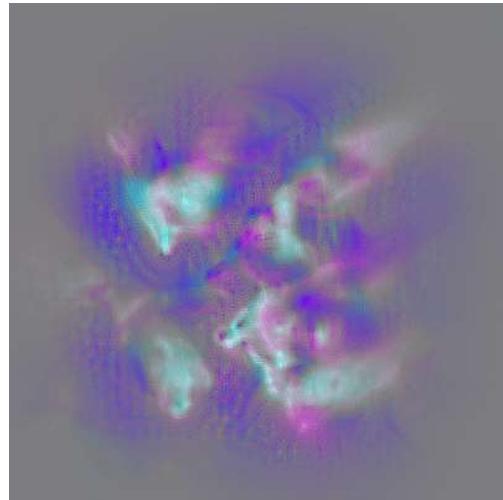


Can you guess the action?

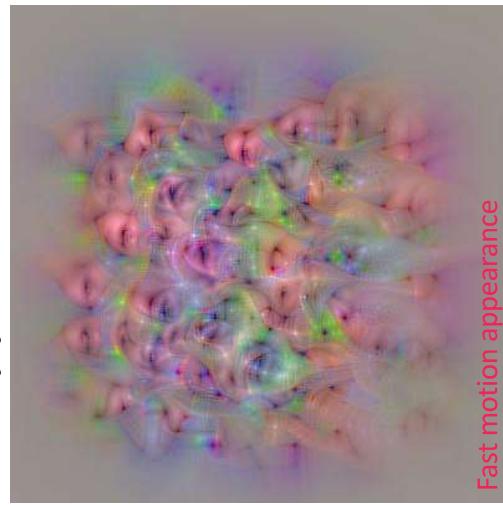
Appearance      “Slow” motion



“Fast” motion

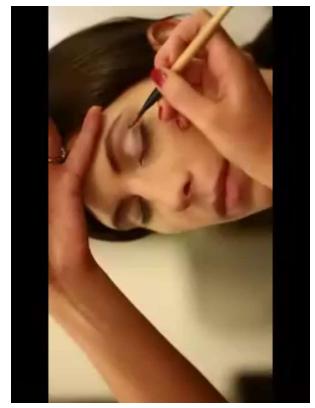
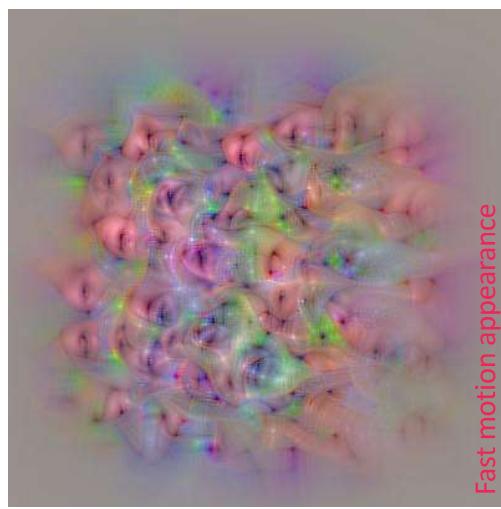
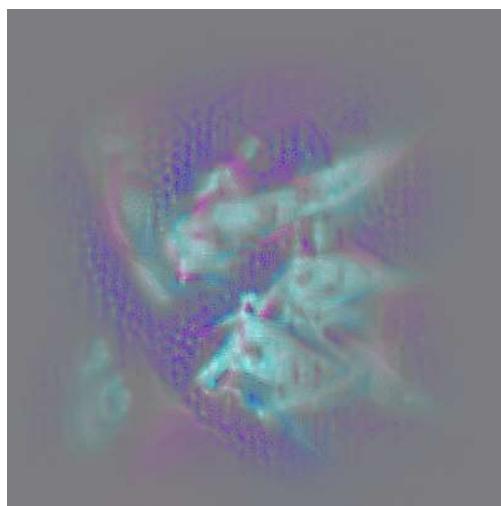
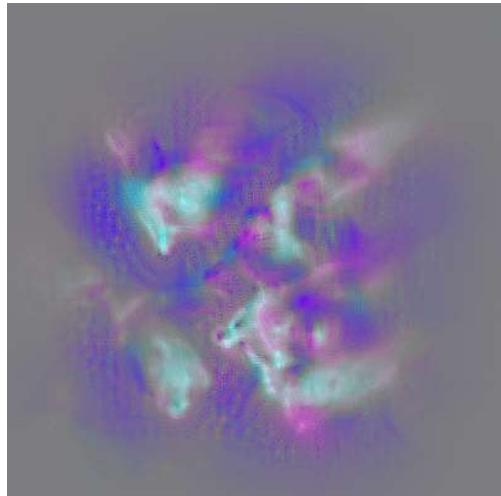


Fast motion appearance

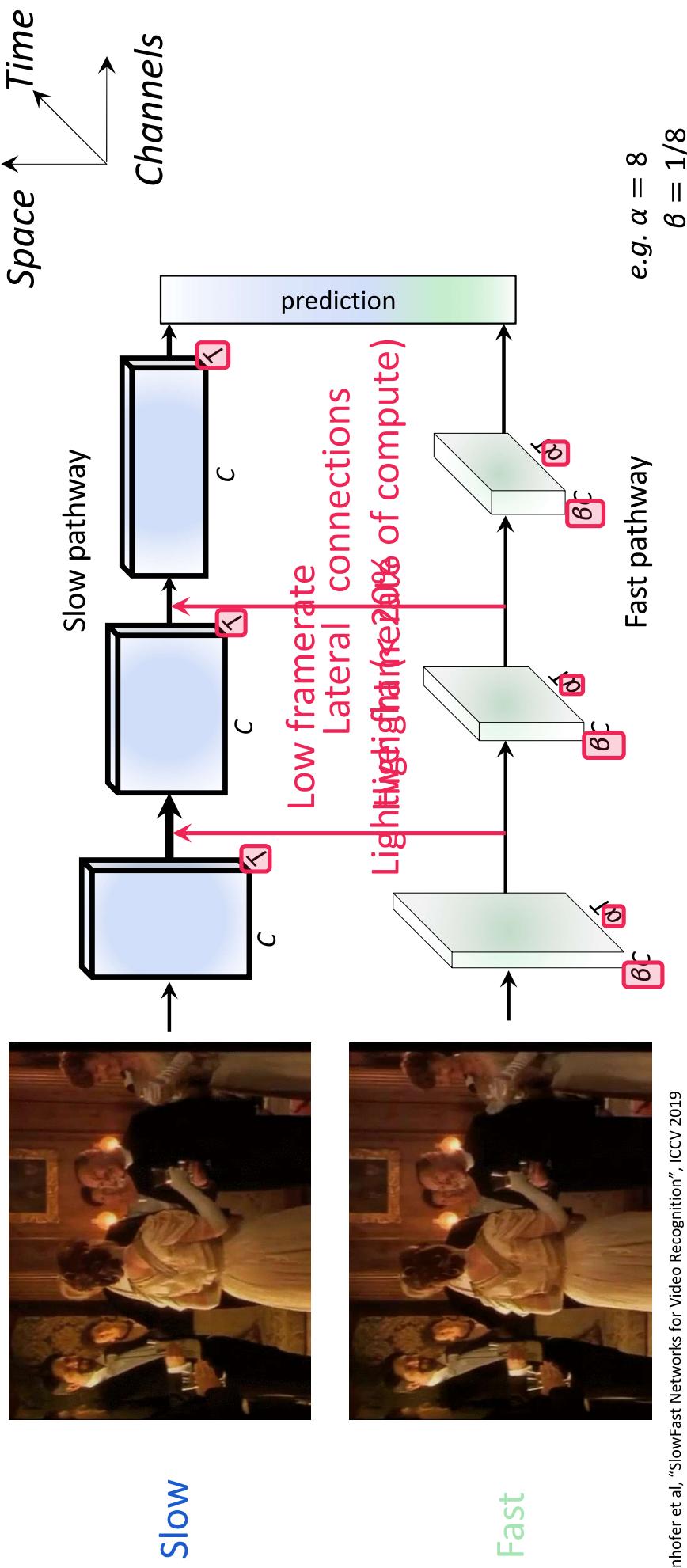


Can you guess the action? Apply Eye Makeup

Appearance      “Slow” motion      “Fast” motion



# Treating time and space differently: SlowFast Networks



Feichtenhofer et al., "SlowFast Networks for Video Recognition", ICCV 2019  
Slide credit: Christoph Feichtenhofer

Justin Johnson

Lecture 18 - 85

November 18, 2019

$$\begin{aligned} \text{e.g. } \alpha &= 8 \\ \theta &= 1/8 \end{aligned}$$

# Treating time and space differently: SlowFast Networks

- Dimensions are  $\{T \times S^2, C\}$
- Strides are {temporal, spatial $^2\}$
- The backbone is ResNet-50
- Residual blocks are shown by brackets
- Non-degenerate temporal filters are underlined
- Here the speed ratio is  $\alpha = 8$  and the channel ratio is  $\theta = 1/8$
- Orange numbers mark fewer channels, for the Fast pathway
- Green numbers mark higher temporal resolution of the Fast pathway
- No temporal pooling is performed throughout the hierarchy

stage	Slow pathway	Fast pathway	output sizes $T \times S^2$
raw clip	-	-	$64 \times 224^2$
data layer	stride $16, 1^2$	stride <b>2</b> , $1^2$	$Slow : 4 \times 224^2$ $Fast : \underline{32} \times 224^2$
conv <sub>1</sub>	$\underline{1 \times 7^2}, 64$ stride $1, 2^2$	$\underline{5 \times 7^2}, \underline{8}$ stride $1, 2^2$	$Slow : 4 \times 112^2$ $Fast : \underline{32} \times 112^2$
pool <sub>1</sub>	$1 \times 3^2$ max stride $1, 2^2$	$1 \times 3^2$ max stride $1, 2^2$	$Slow : 4 \times 56^2$ $Fast : \underline{32} \times 56^2$
res <sub>2</sub>	$\left[ \begin{array}{c} 1 \times 1^2, 64 \\ 1 \times 3^2, 64 \\ 1 \times 1^2, 256 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 3 \times 1^2, \underline{8} \\ 1 \times 3^2, \underline{8} \\ 1 \times 1^2, \underline{32} \end{array} \right] \times 3$	$Slow : 4 \times 56^2$ $Fast : \underline{32} \times 56^2$
res <sub>3</sub>	$\left[ \begin{array}{c} 1 \times 1^2, 128 \\ 1 \times 3^2, 128 \\ 1 \times 1^2, 512 \end{array} \right] \times 4$	$\left[ \begin{array}{c} 3 \times 1^2, \underline{16} \\ 1 \times 3^2, \underline{16} \\ 1 \times 1^2, \underline{64} \end{array} \right] \times 4$	$Slow : 4 \times 28^2$ $Fast : \underline{32} \times 28^2$
res <sub>4</sub>	$\left[ \begin{array}{c} 3 \times 1^2, 256 \\ 1 \times 3^2, 256 \\ 1 \times 1^2, 1024 \end{array} \right] \times 6$	$\left[ \begin{array}{c} 3 \times 1^2, \underline{32} \\ 1 \times 3^2, \underline{32} \\ 1 \times 1^2, \underline{128} \end{array} \right] \times 6$	$Slow : 4 \times 14^2$ $Fast : \underline{32} \times 14^2$
res <sub>5</sub>	$\left[ \begin{array}{c} 3 \times 1^2, 512 \\ 1 \times 3^2, 512 \\ 1 \times 1^2, 2048 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 3 \times 1^2, \underline{64} \\ 1 \times 3^2, \underline{64} \\ 1 \times 1^2, \underline{256} \end{array} \right] \times 3$	$Slow : 4 \times 7^2$ $Fast : \underline{32} \times 7^2$
	global average pool, concat, fc		# classes

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

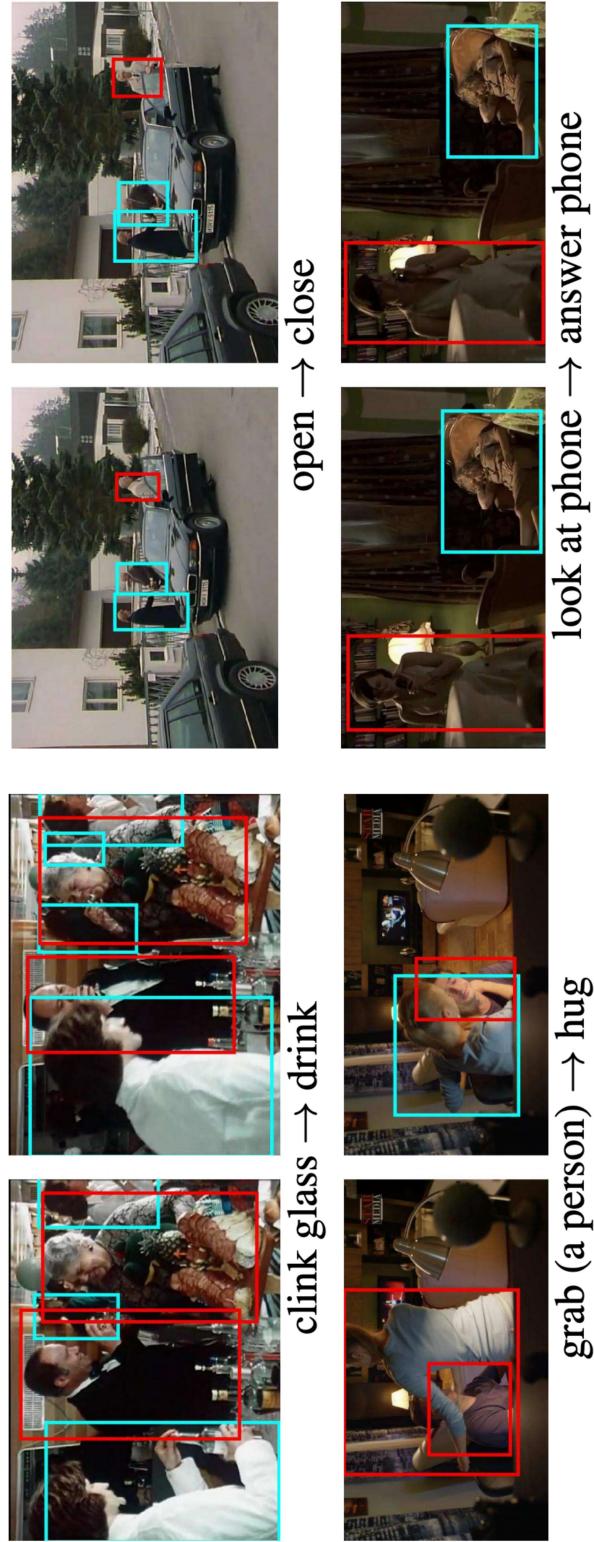
Justin Johnson

Lecture 18 - 89

November 18, 2019

# Spatio-Temporal Detection

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



Gu et al, "AVA: A Video Dataset of Spatio-temporally Localized Atomic Visual Actions", CVPR 2018

Justin Johnson

Lecture 18 - 90

November 18, 2019

# Recap: Video Models

## **Many video models:**

Single-frame CNN (Try this first!)

Late fusion

Early fusion

3D CNN / C3D

Two-stream networks

CNN + RNN

Convolutional RNN

Spatio-temporal self-attention

SlowFast networks (current SoTA)

Next time:  
Generative Models, part 1  
Generative Adversarial Networks