

Indian Institute of Information Technology, Allahabad



Video Recognition and Understanding

By

Dr. Satish Kumar Singh & Dr. Shiv Ram Dubey

Computer Vision and Biometrics Lab

Department of Information Technology

Indian Institute of Information Technology, Allahabad



TEAM

Computer Vision and Biometrics Lab (CVBL)

Department of Information Technology

Indian Institute of Information Technology Allahabad

Course Instructors

Dr. Satish Kumar Singh, Associate Professor, IIIT Allahabad (Email: sk.singh@iiita.ac.in)

Dr. Shiv Ram Dubey, Assistant Professor, IIIT Allahabad (Email: srdubey@iiita.ac.in)



DISCLAIMER

The content (text, image, and graphics) used in this slide are adopted from many sources for Academic purposes. Broadly, the sources have been given due credit appropriately. However, there is a chance of missing out some original primary sources. The authors of this material do not claim any copyright of such material.

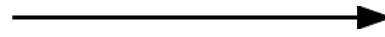


Recall: (2D) Image Classification



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

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



cat



Recall: (2D) Detection And Segmentation

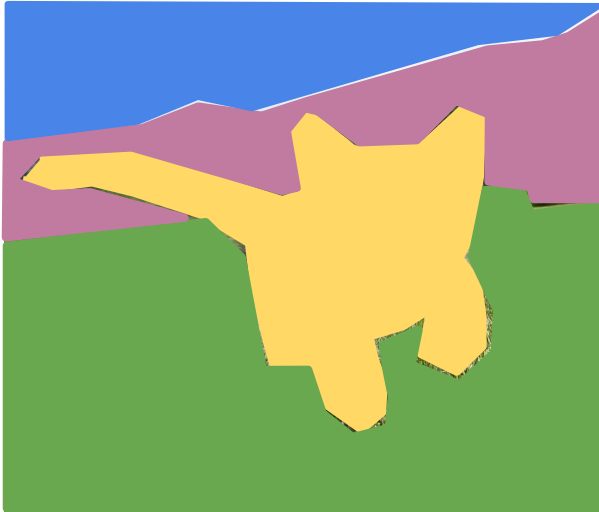
Classification



CAT

No spatial extent

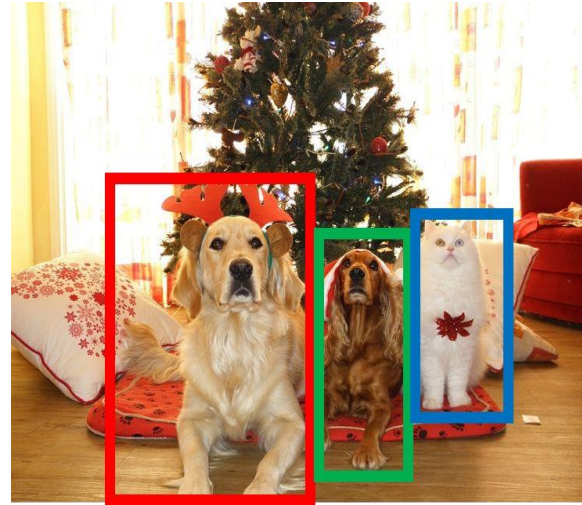
Semantic Segmentation



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](#)



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

Running video is in the public domain

Slide credit: Justin Johnson



Problem: Videos are big!



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

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**



Problem: Videos are big!



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

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**

Solution:

Train on short **clips**: low fps and low spatial resolution

e.g. $T = 16$ (3.2 seconds at 5 fps), $H=W=112$

~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



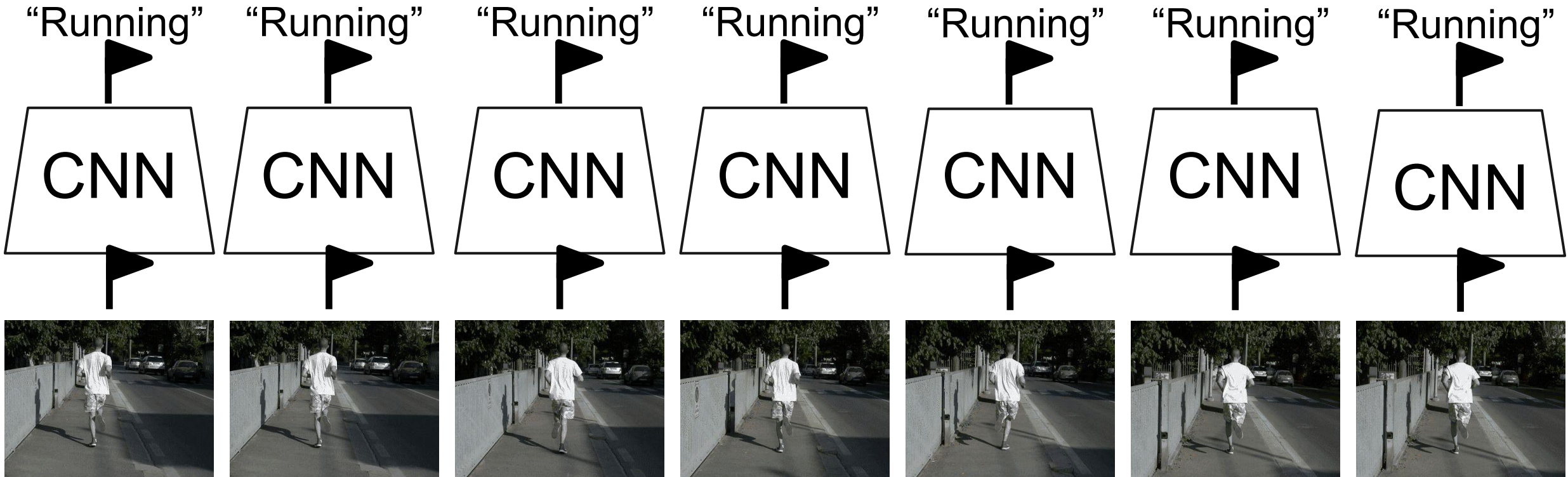
Testing: Run model on different clips, average predictions



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)

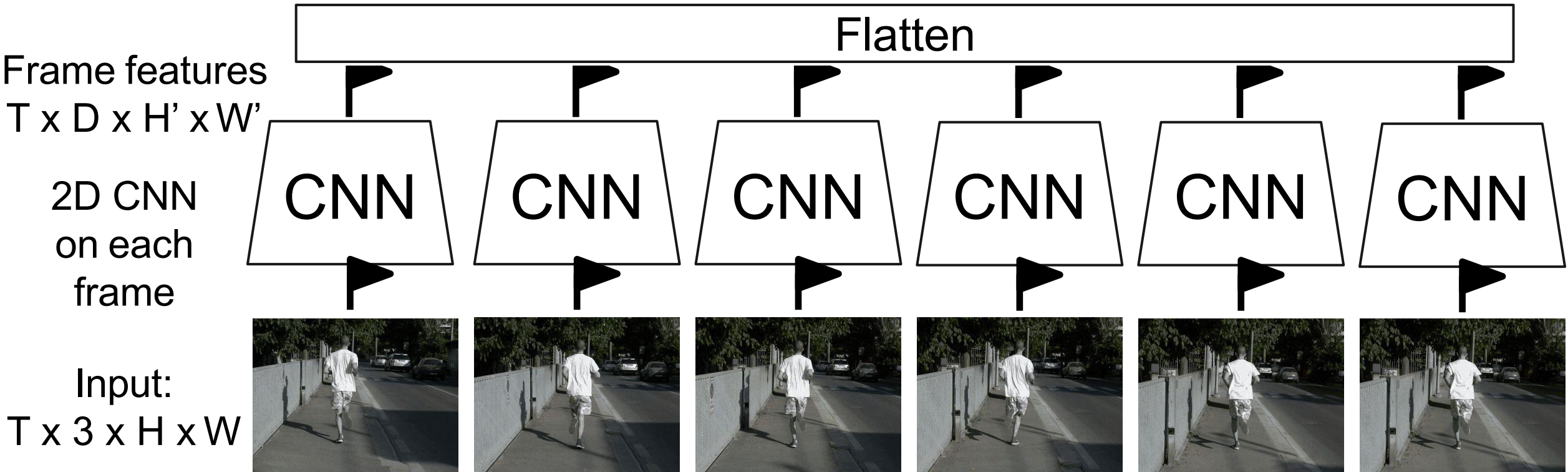
Intuition: Get high-level appearance of each frame, and combine them

Class scores: C

Clip features:
TDH'W'

MLP

Run 2D CNN on each frame, concatenate features and feed to MLP



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

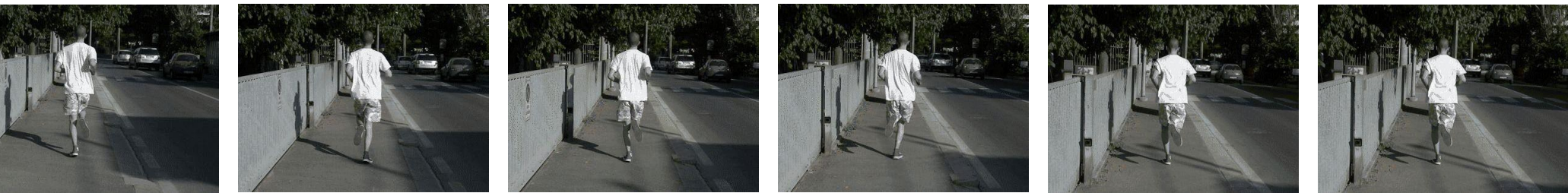
Clip features: D

Linear

Average Pool over space and time

Frame features
 $T \times D \times H' \times W'$

2D CNN
on each
frame



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

Class scores: C

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

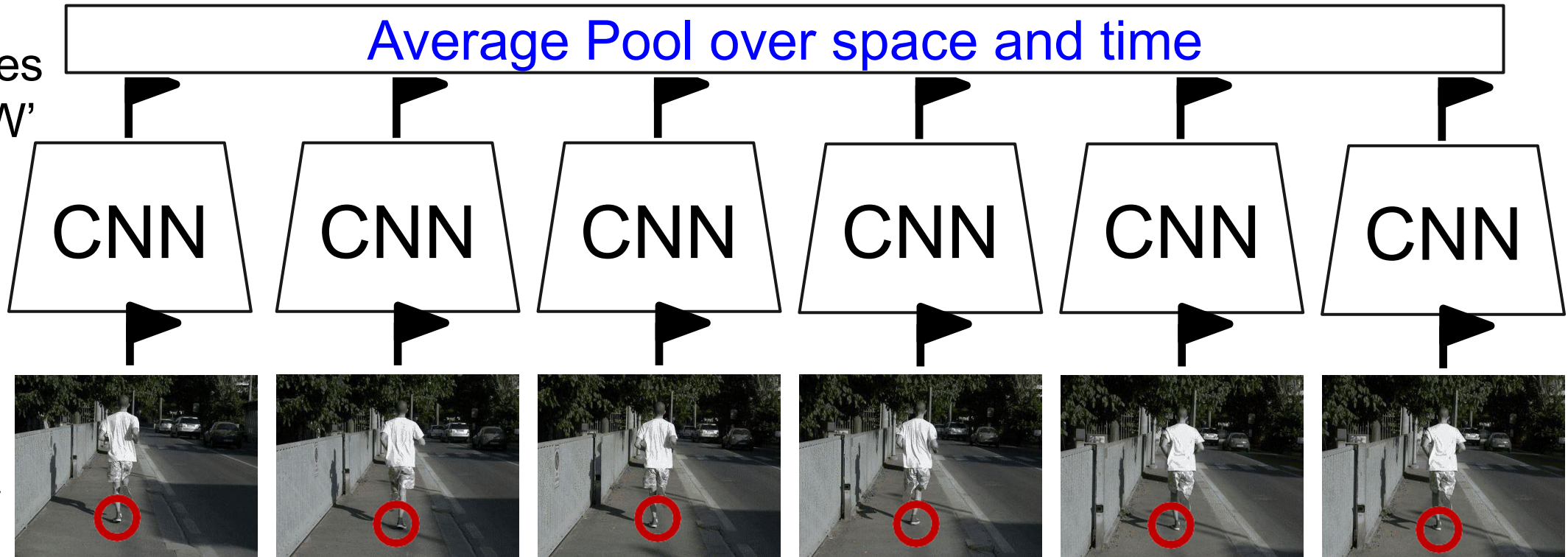
Clip features: D

Linear

Average Pool over space and time

Frame features
 $T \times D \times H' \times W'$

2D CNN
on each
frame



Input:

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

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

First 2D convolution collapses all temporal information:

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

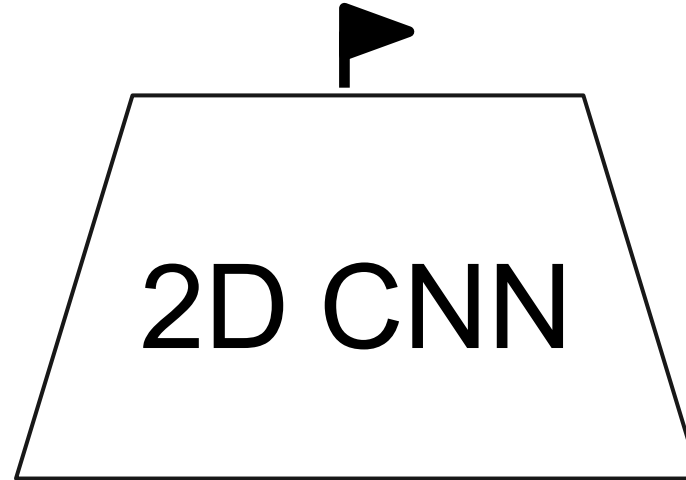
Output: $D \times H \times W$

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

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



Class scores: C



Rest of the network is standard 2D CNN



Video Classification: Early Fusion

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

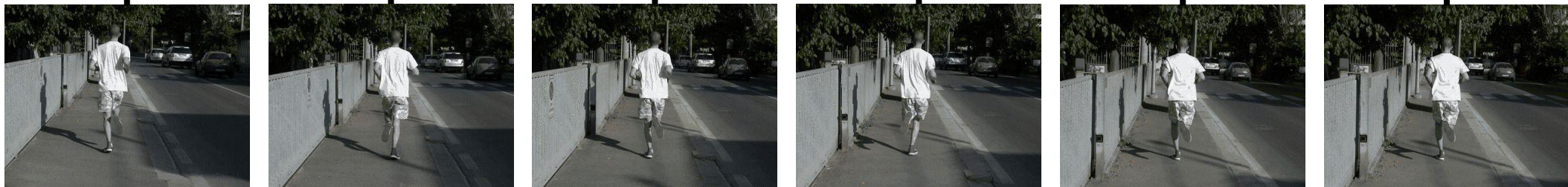
First 2D convolution collapses all temporal information:

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

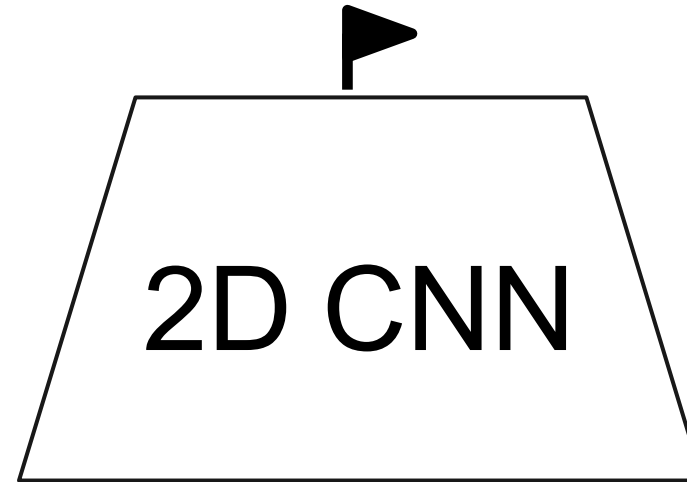
Output: $D \times H \times W$

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

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



Class scores: C



Rest of the network is standard 2D CNN

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



Video Classification: 3D CNN

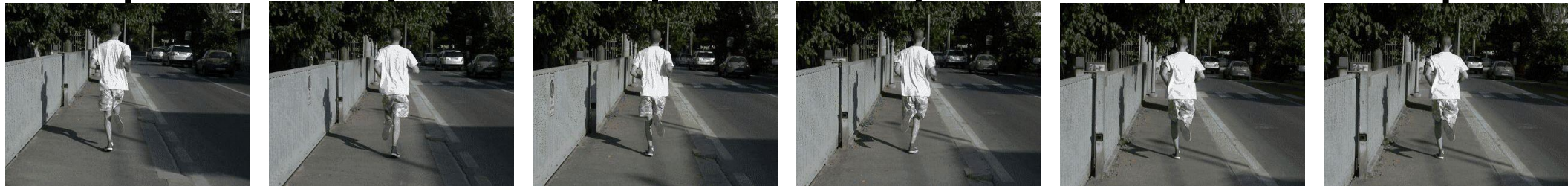
Intuition: Use 3D versions of convolution and pooling to slowly fuse temporal information over the course of the network

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

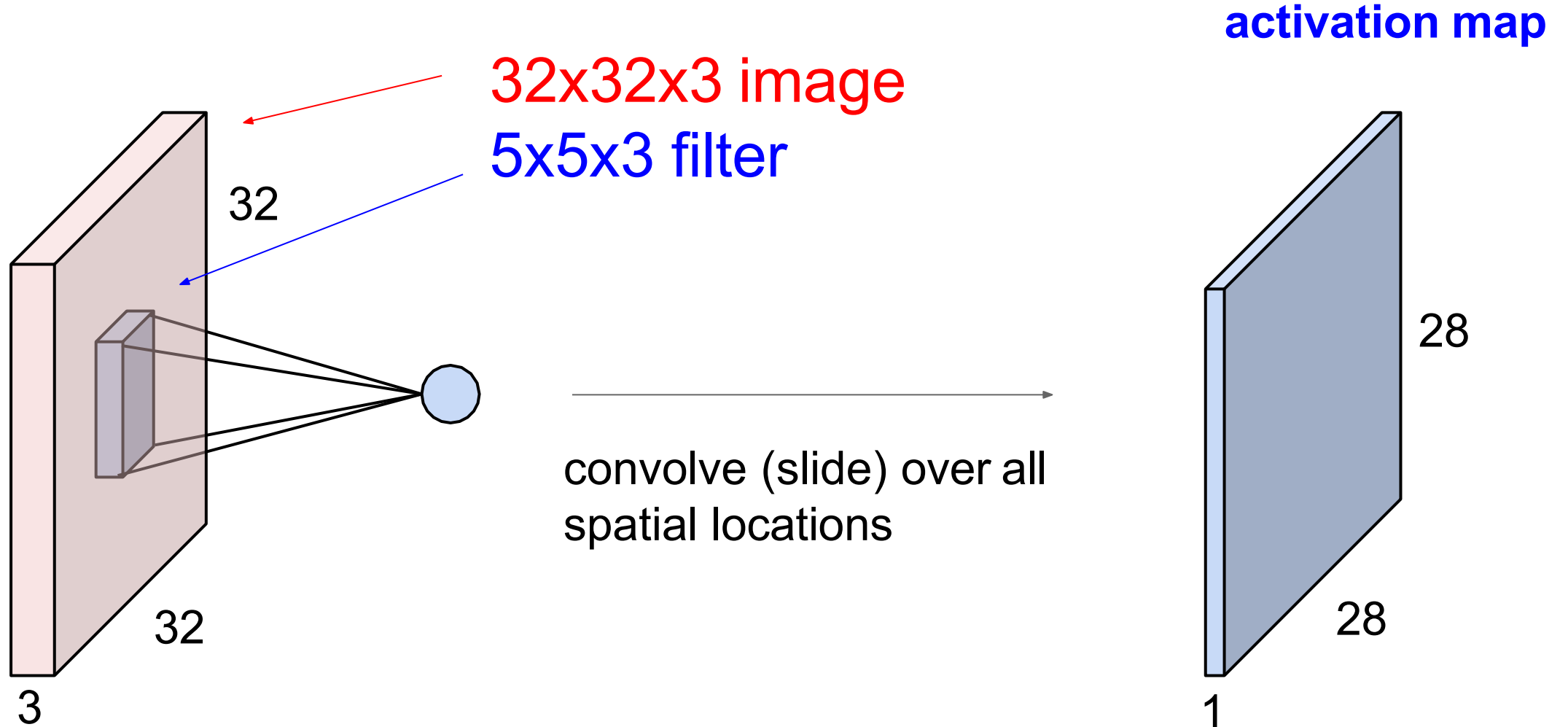
Class scores: C

3D CNN

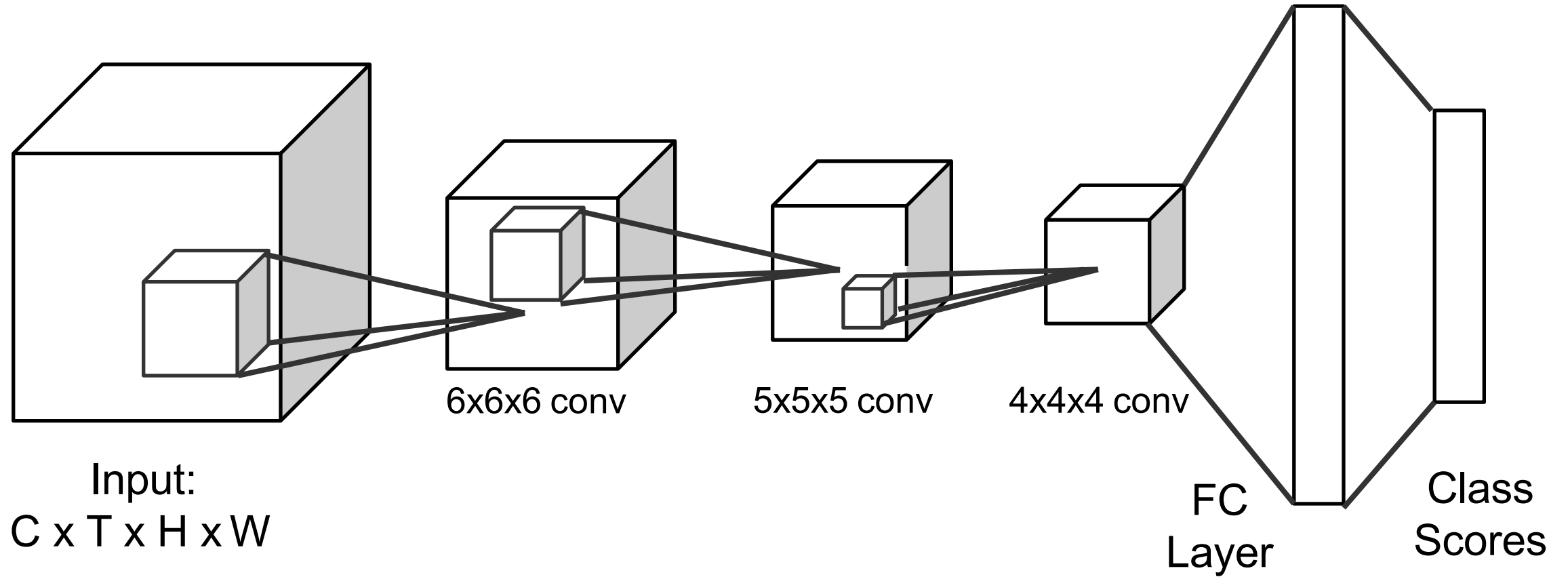
Input:
 $3 \times T \times H \times W$



Convolution Layer



3D Convolution



Early Fusion vs Late Fusion vs 3D CNN

Late
Fusion

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

(Small example architectures, in practice much bigger)



Early Fusion vs Late Fusion vs 3D CNN

Late Fusion

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

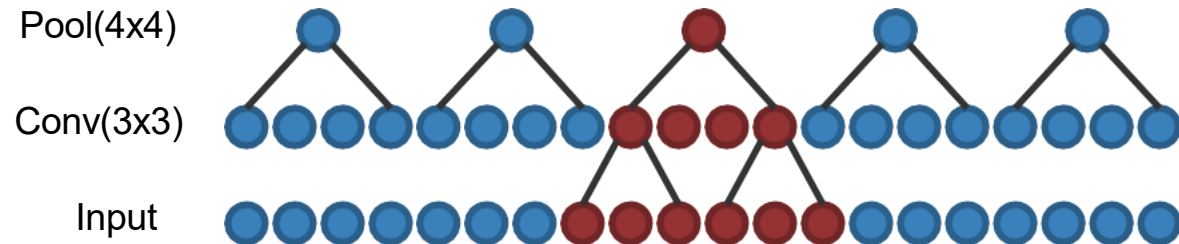


(Small example architectures, in practice much bigger)



Early Fusion vs Late Fusion vs 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



(Small example architectures, in practice much bigger)

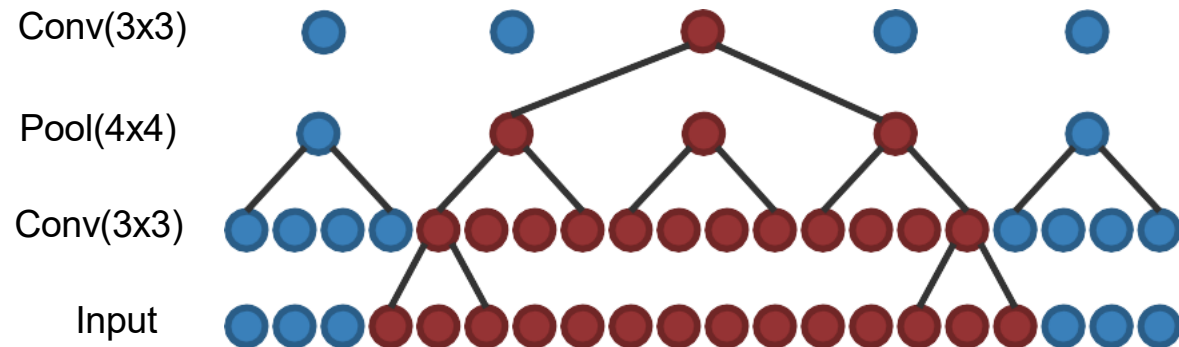


Early Fusion vs Late Fusion vs 3D CNN

Late Fusion

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

Build slowly in space



(Small example architectures, in practice much bigger)

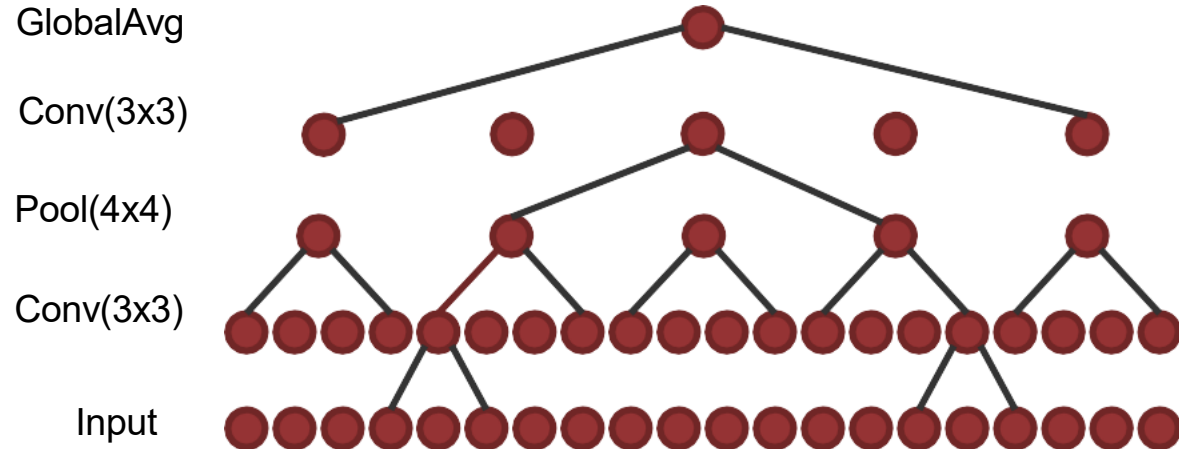


Early Fusion vs Late Fusion vs 3D CNN

Late Fusion

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
GlobalAvgPool	24 x 1 x 1 x 1	20 x 64 x 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

	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		Build slowly in space, All-at-once in time at end
	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		Build slowly in space, All-at-once in time at start
	Conv2D(3x3, 3*20->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	

(Small example architectures, in practice much bigger)



Early Fusion vs Late Fusion vs 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		Build slowly in space, All-at-once in time at end	
	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		Build slowly in space, All-at-once in time at start	
	Conv2D(3x3, 3*20->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		
3D CNN	Input	3 x 20 x 64 x 64		Build slowly in space, Build slowly in time "Slow Fusion"	(Small example architectures, in practice much bigger)
	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 vs 3D CNN

What is the difference?

Late Fusion

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
GlobalAvgPool	24 x 1 x 1 x 1	20 x 64 x 64

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

Early Fusion

Input	3 x 20 x 64 x 64	
Conv2D(3x3, 3*20->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

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

3D CNN

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

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

(Small example architectures, in practice much bigger)

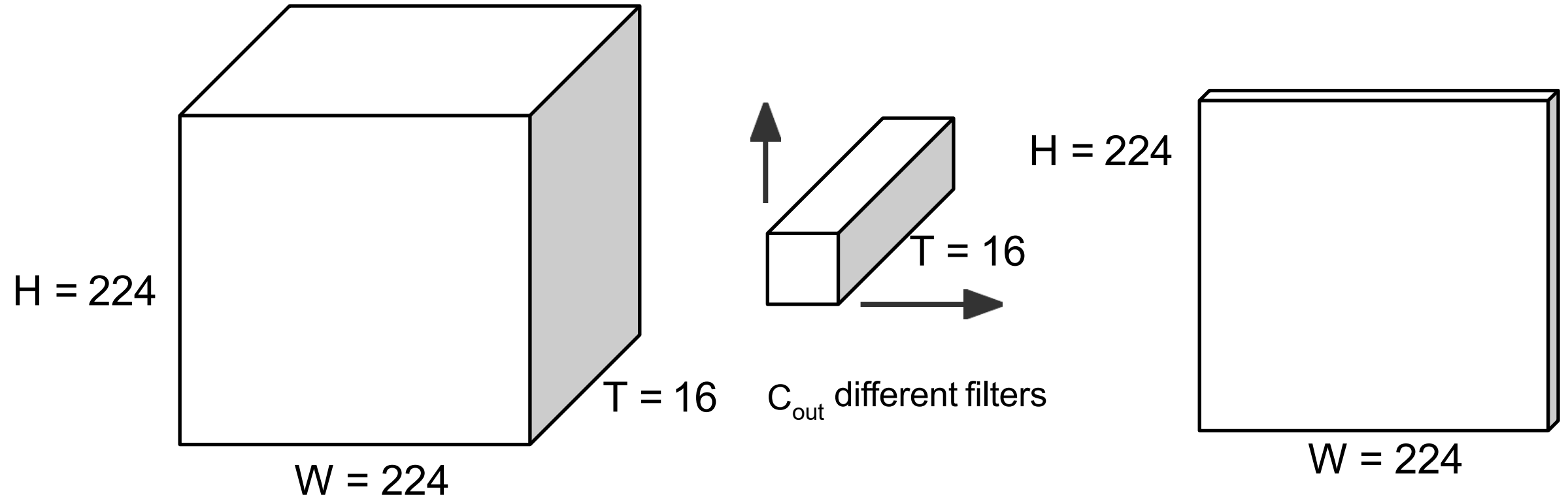


2D Conv (Early Fusion) vs 3D Conv (3D CNN)

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



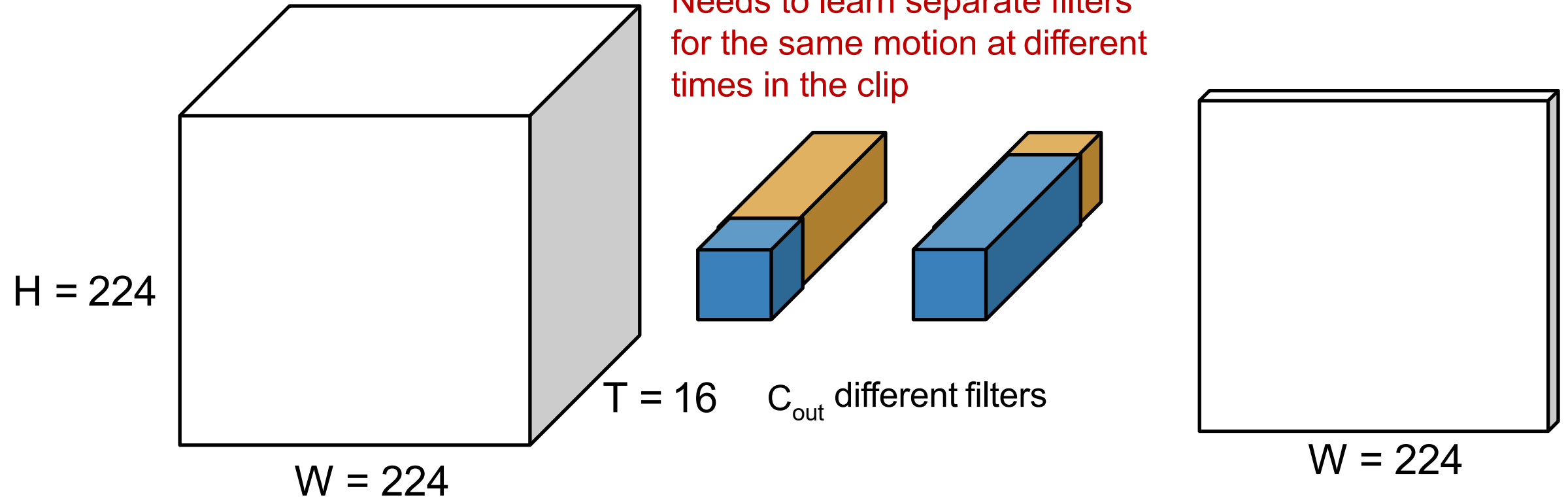
2D Conv (Early Fusion) vs 3D Conv (3D CNN)

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

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

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



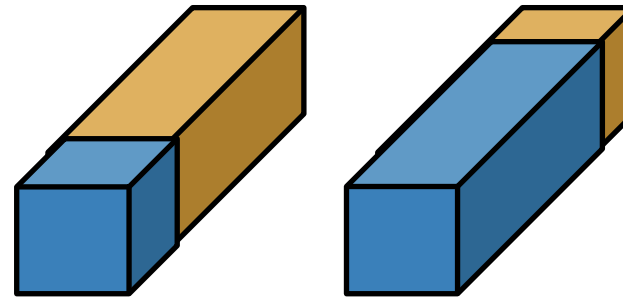
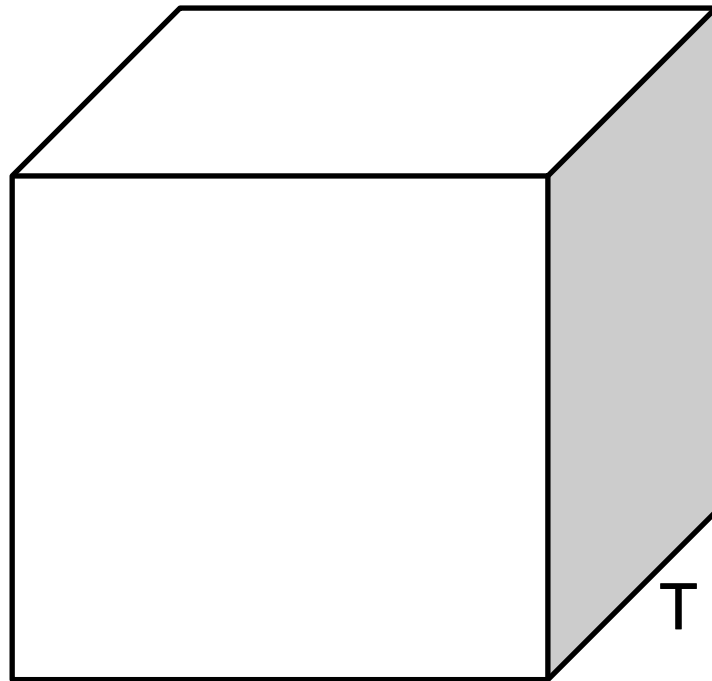
2D Conv (Early Fusion) vs 3D Conv (3D CNN)

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

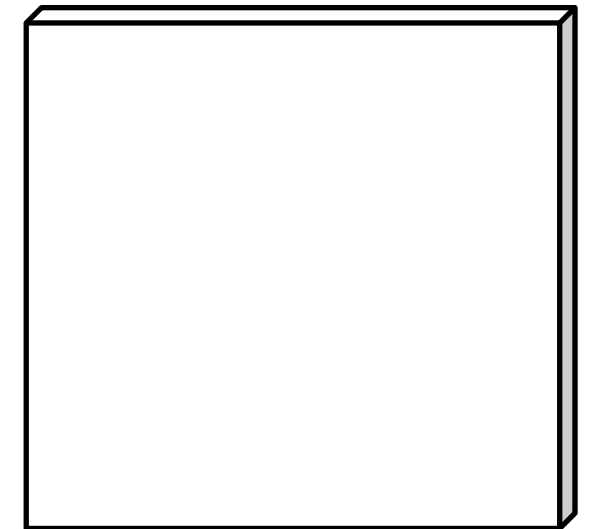
Weight:
 $C_{out} \times C_{in} \times T \times 3 \times 3$
Slide over x and y

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

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



$T = 16$ C_{out} different filters



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

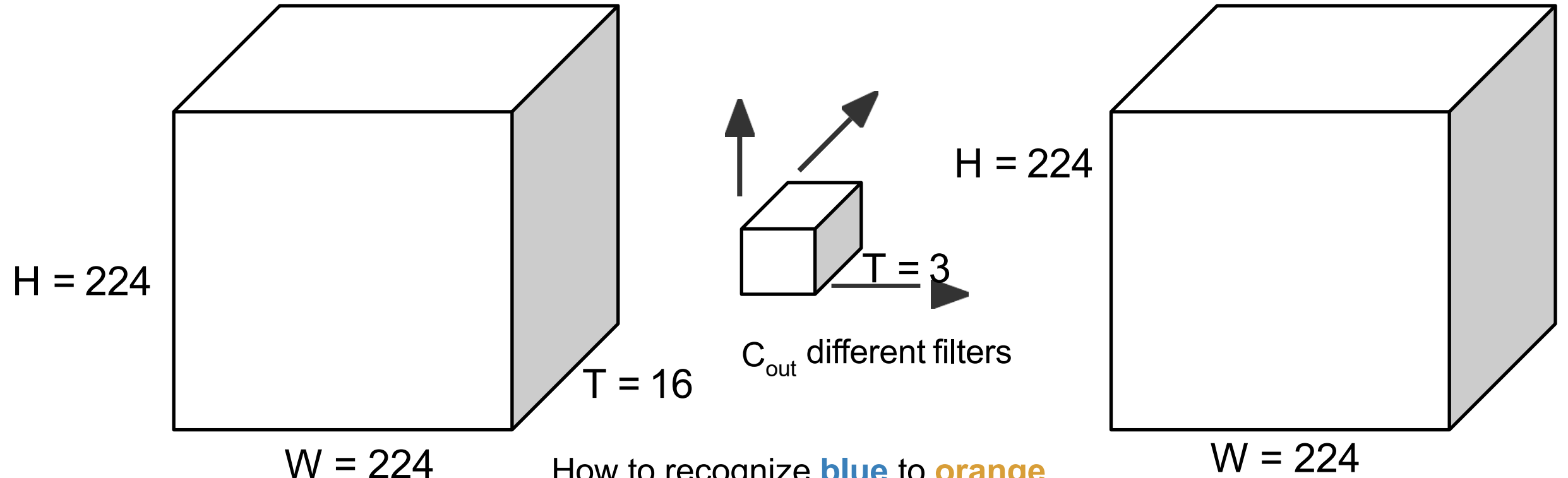


2D Conv (Early Fusion) vs 3D Conv (3D CNN)

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

Weight:
 $C_{out} \times C_{in} \times 3 \times 3 \times 3$

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



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



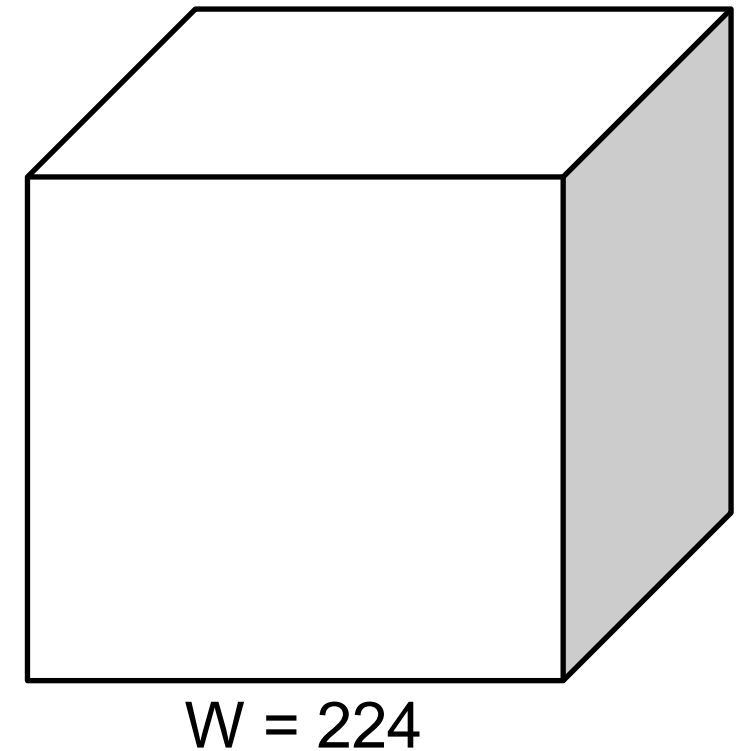
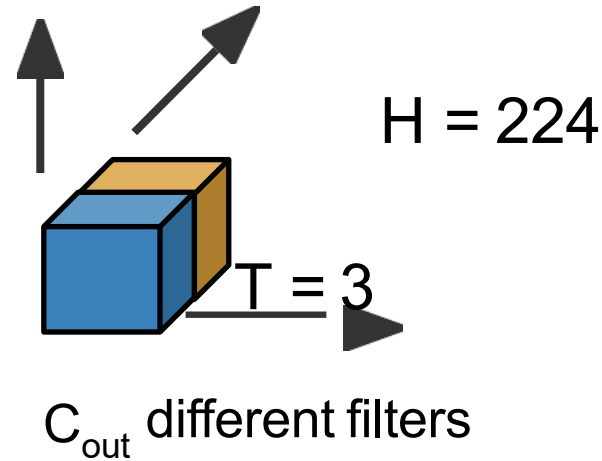
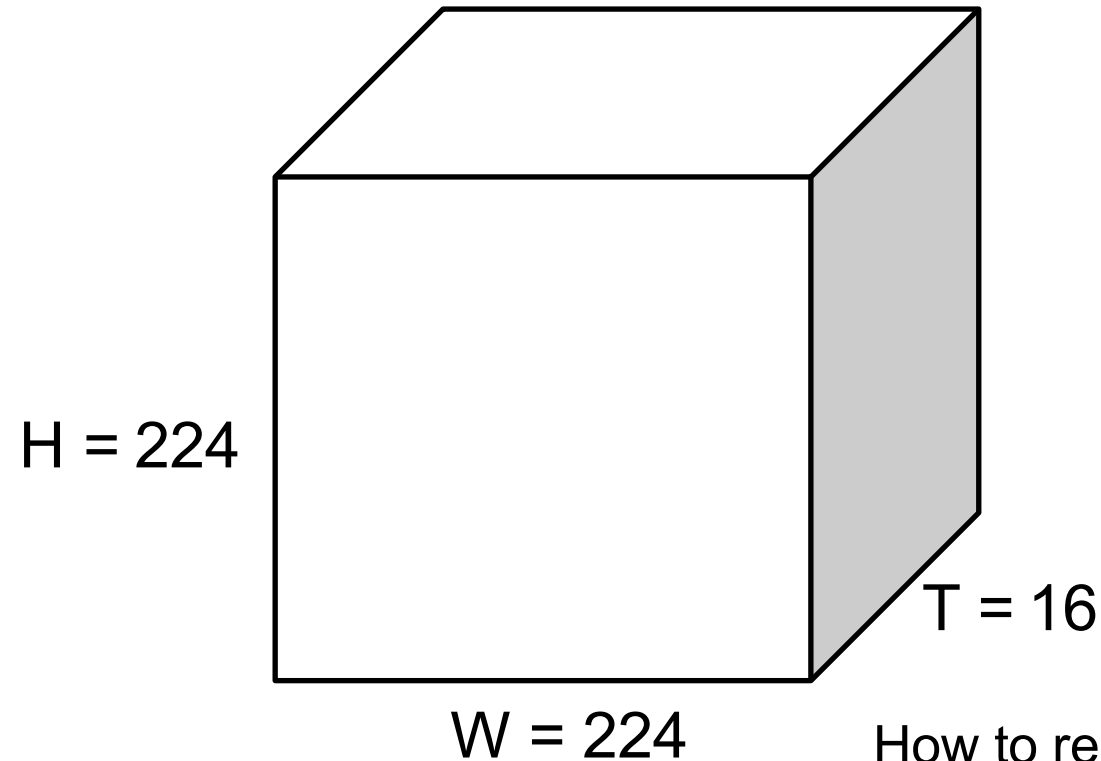
2D Conv (Early Fusion) vs 3D Conv (3D CNN)

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, y and t

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!



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



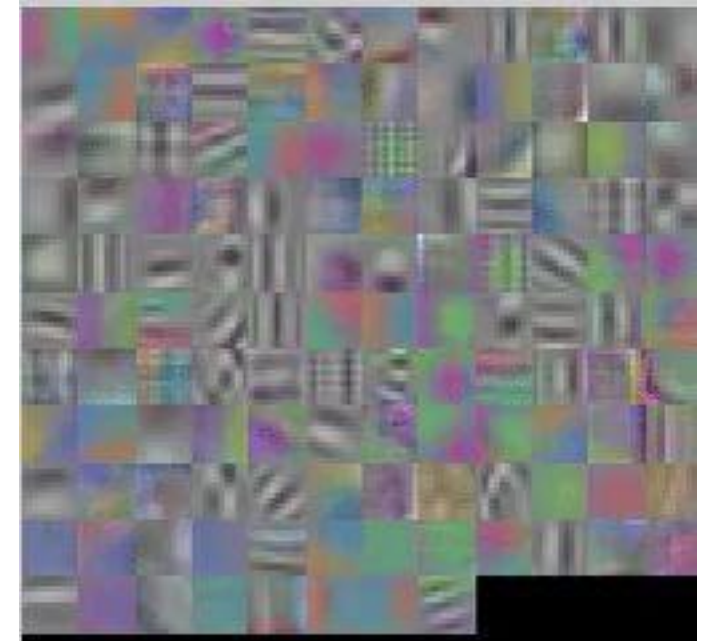
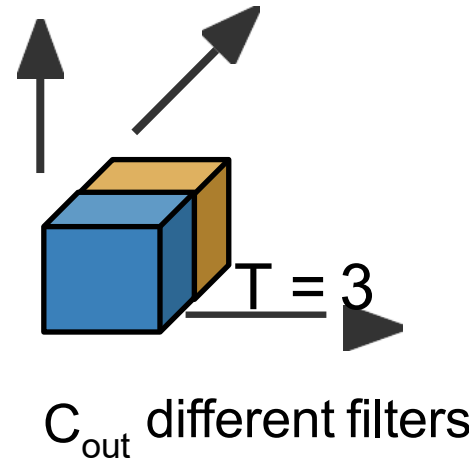
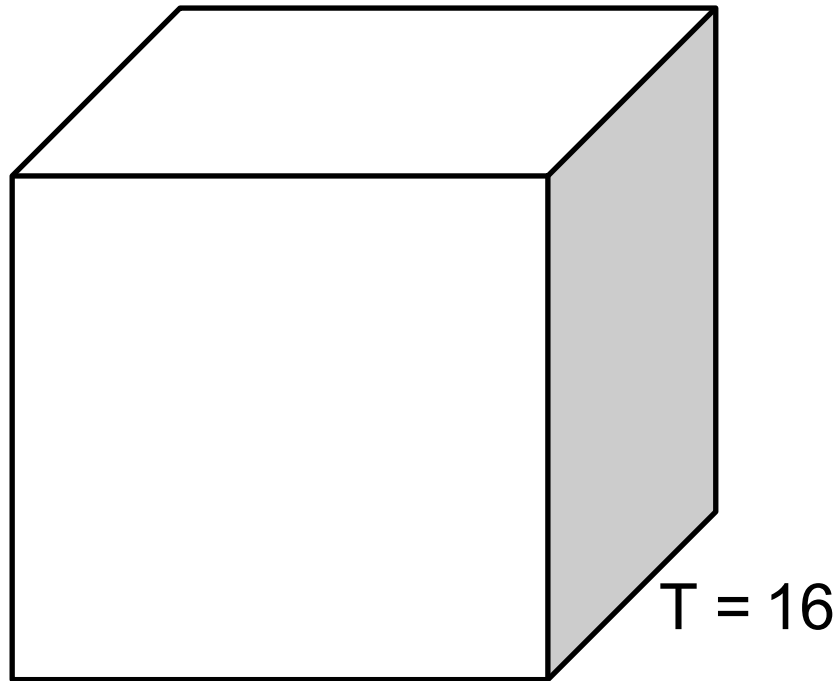
2D Conv (Early Fusion) vs 3D Conv (3D CNN)

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, y and t

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

Temporal shift-invariant
since each filter slides
over time!



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



Example Video Dataset: Sports-1M



track cycling
cycling
track cycling
road bicycle racing
marathon
ultramarathon



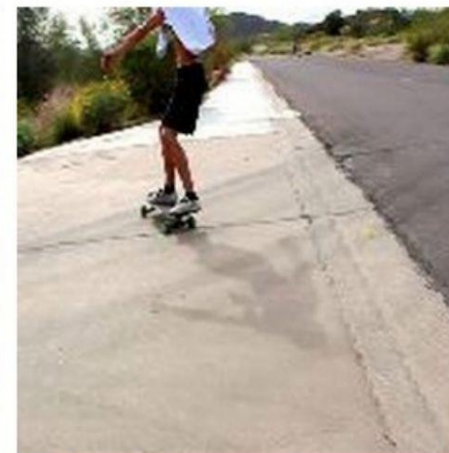
ultramarathon
ultramarathon
half marathon
running
marathon
inline speed skating



heptathlon
heptathlon
decathlon
hurdles
pentathlon
sprint (running)



bikejoring
mushing
bikejoring
harness racing
skijoring
carting



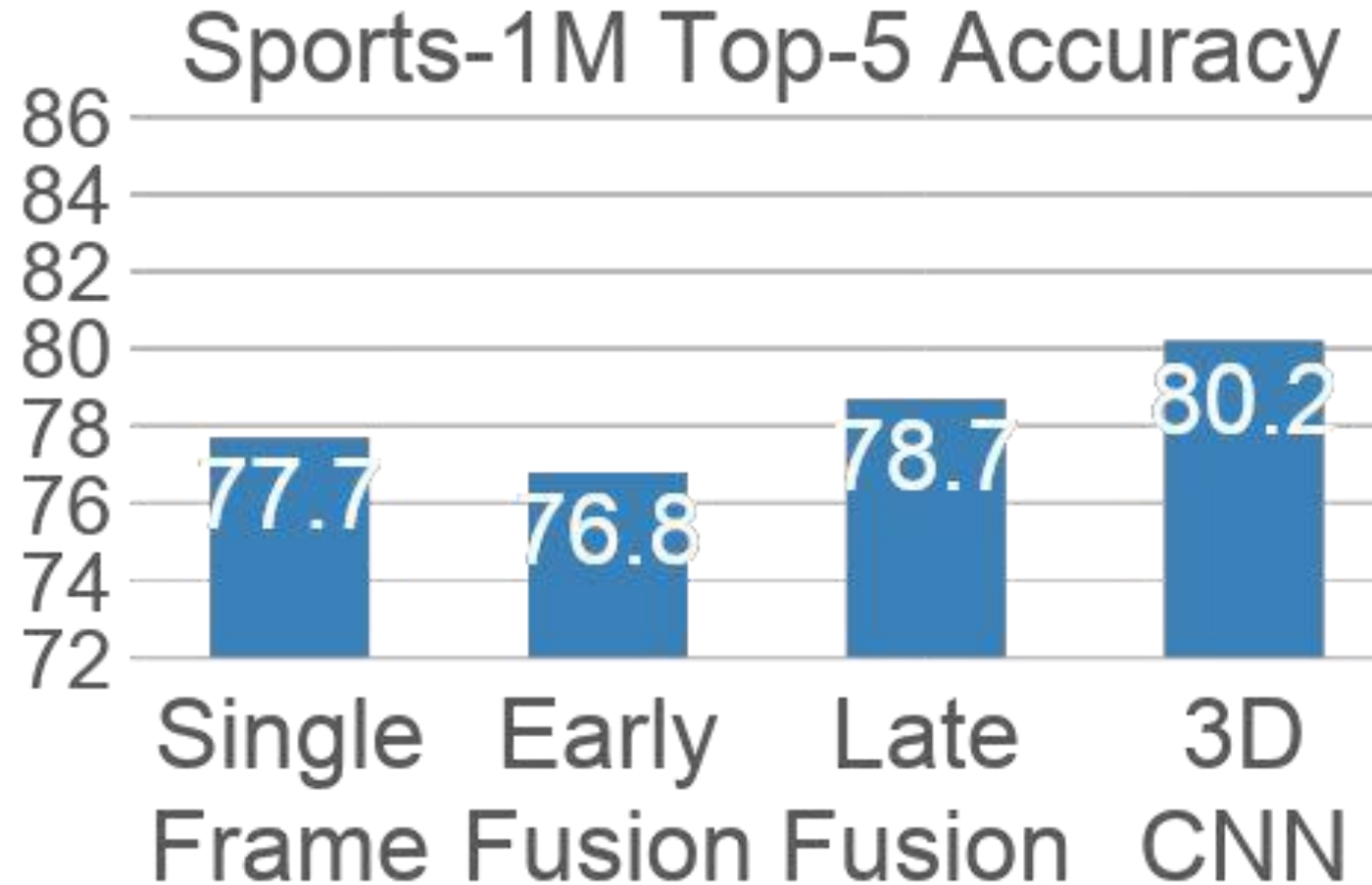
longboarding
longboarding
aggressive inline skating
freestyle scootering
freeboard (skateboard)
sandboarding

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 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

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 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

Problem: 3x3x3 conv is very expensive!

AlexNet: 0.7 GFLOP

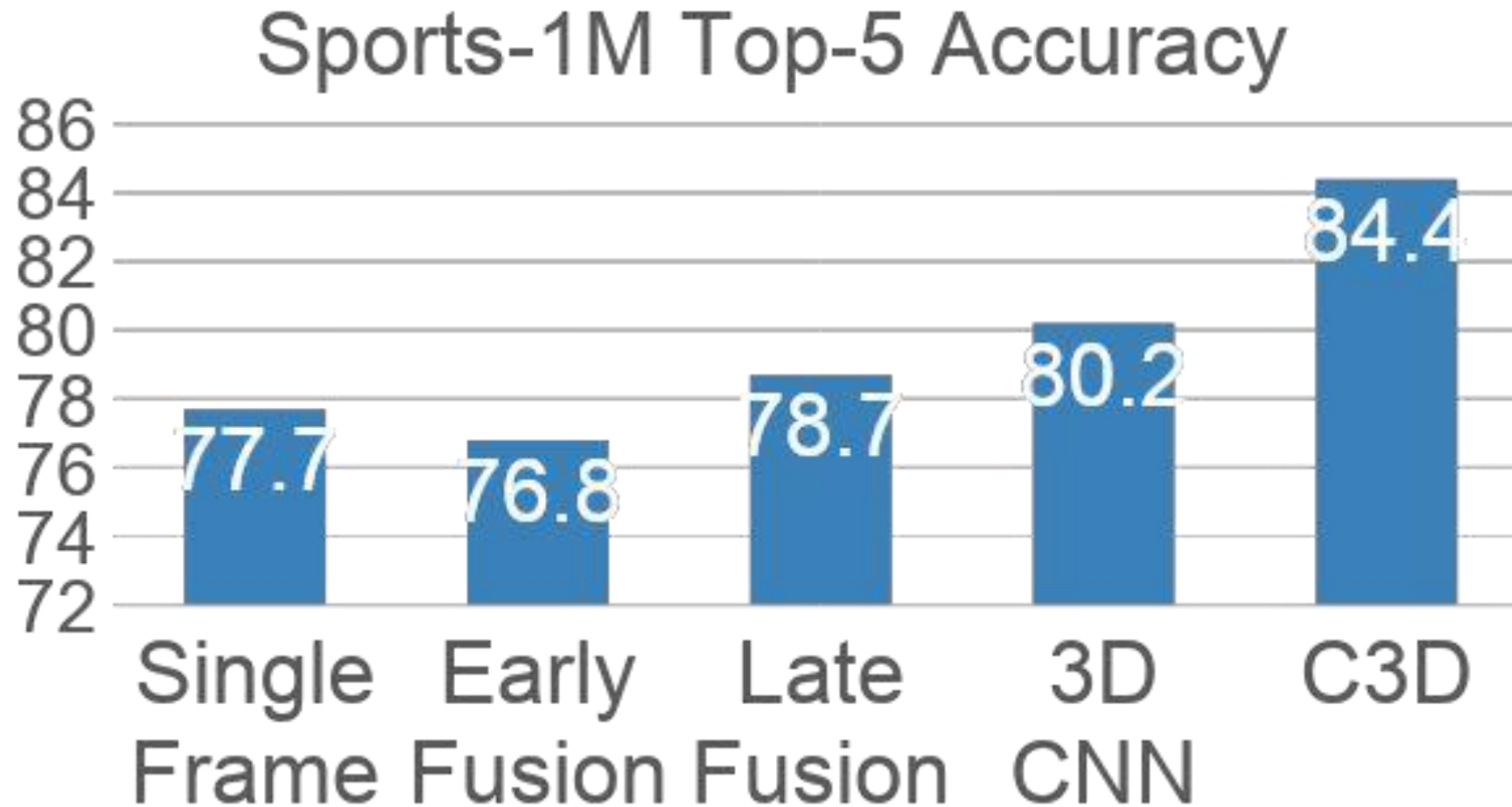
VGG-16: 13.6 GFLOP

C3D: **39.5 GFLOP (2.9x VGG!)**

Layer	Size	MFLOPs
Input	3 x 16 x 112 x 112	
Conv1 (3x3x3)	64 x 16 x 112 x 112	1.04
Pool1 (1x2x2)	64 x 16 x 56 x 56	
Conv2 (3x3x3)	128 x 16 x 56 x 56	11.10
Pool2 (2x2x2)	128 x 8 x 28 x 28	
Conv3a (3x3x3)	256 x 8 x 28 x 28	5.55
Conv3b (3x3x3)	256 x 8 x 28 x 28	11.10
Pool3 (2x2x2)	256 x 4 x 14 x 14	
Conv4a (3x3x3)	512 x 4 x 14 x 14	2.77
Conv4b (3x3x3)	512 x 4 x 14 x 14	5.55
Pool4 (2x2x2)	512 x 2 x 7 x 7	
Conv5a (3x3x3)	512 x 2 x 7 x 7	0.69
Conv5b (3x3x3)	512 x 2 x 7 x 7	0.69
Pool5	512 x 1 x 3 x 3	
FC6	4096	0.51
FC7	4096	0.45
FC8	C	0.05



Early Fusion vs Late Fusion vs 3D CNN



Recognizing Actions from Motion

We can easily recognize actions using only **motion information**



Measuring Motion: Optical Flow

Image at frame t

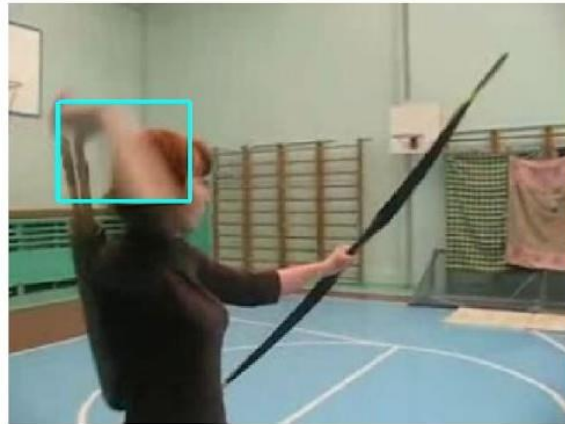
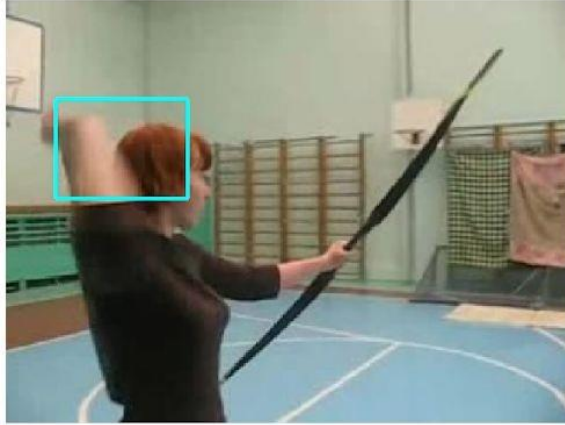


Image at frame $t+1$



Measuring Motion: Optical Flow

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

Image at frame t

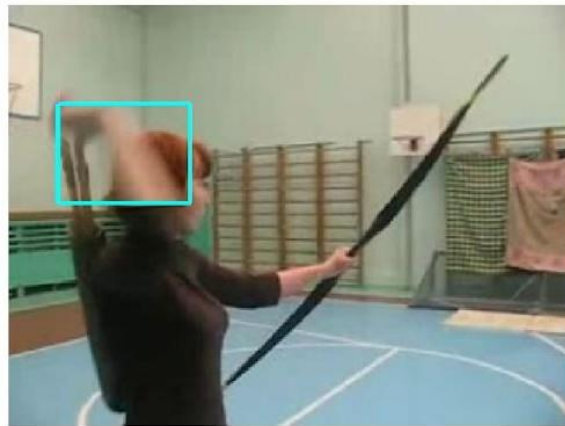
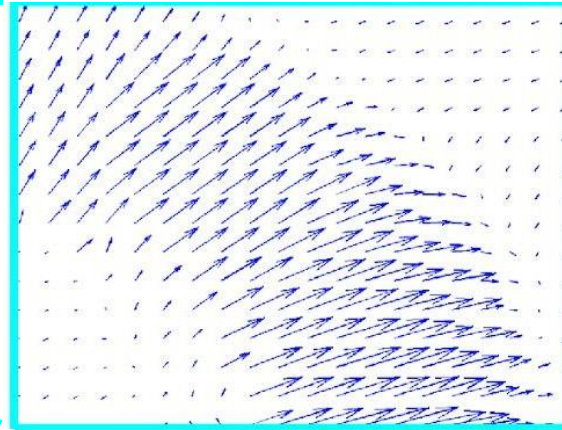
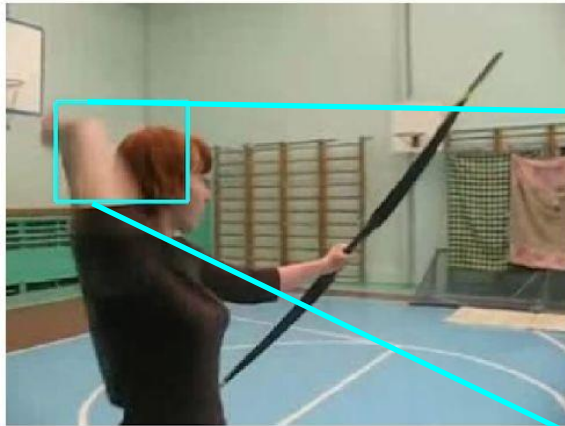


Image at frame $t+1$

Tells where each pixel will move in the next frame:

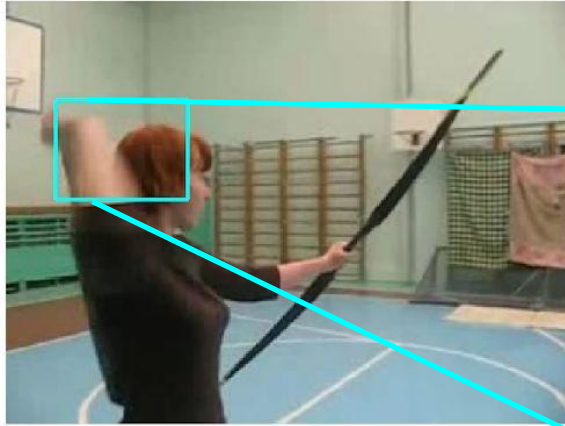
$$F(x, y) = (dx, dy)$$

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



Measuring Motion: Optical Flow

Image at frame t



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

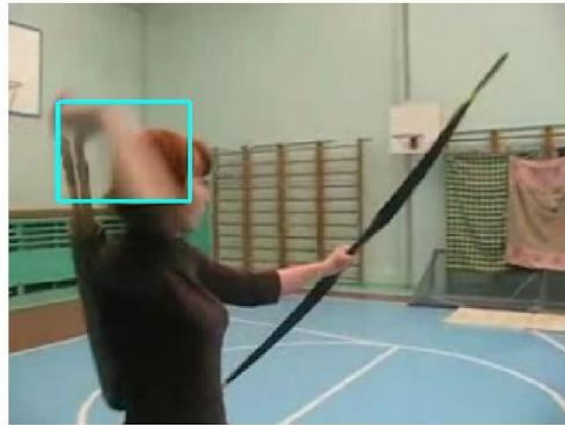
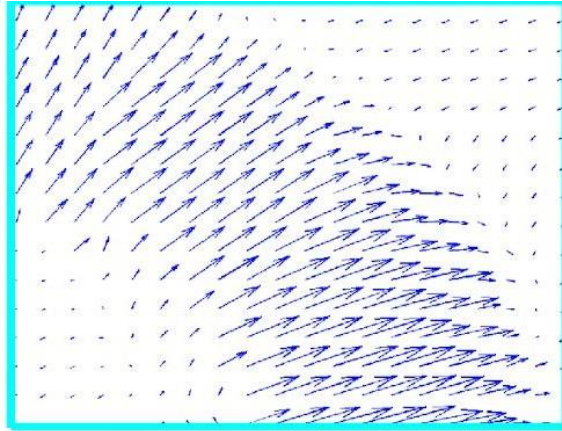


Image at frame t+1

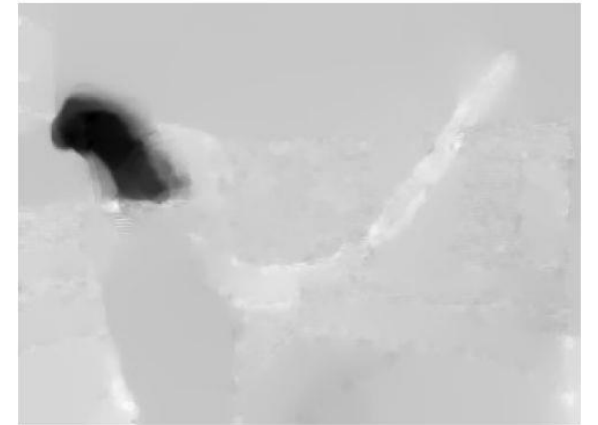
Tells where each pixel will move in the next frame:

$$F(x, y) = (dx, dy)$$

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

Optical Flow highlights **local motion**

Horizontal flow dx



Vertical Flow dy



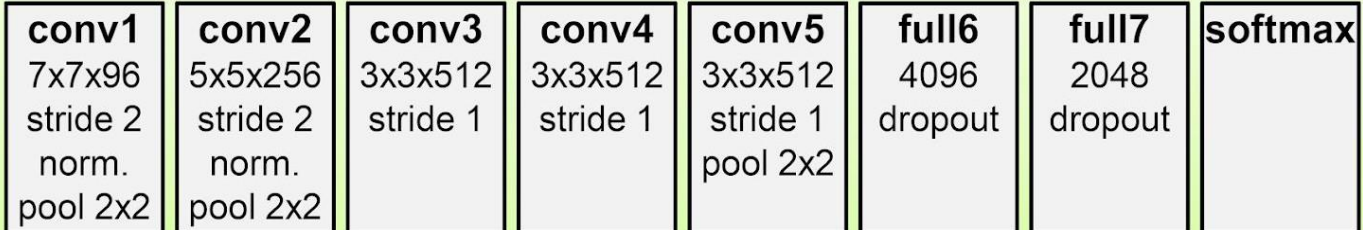
Separating Motion and Appearance: Two-Stream Networks

Input: Single Image
3 x H x W

Spatial stream ConvNet



single frame



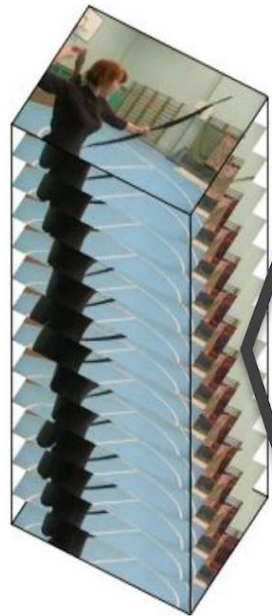
Temporal stream ConvNet



multi-frame
optical flow



class
score
fusion



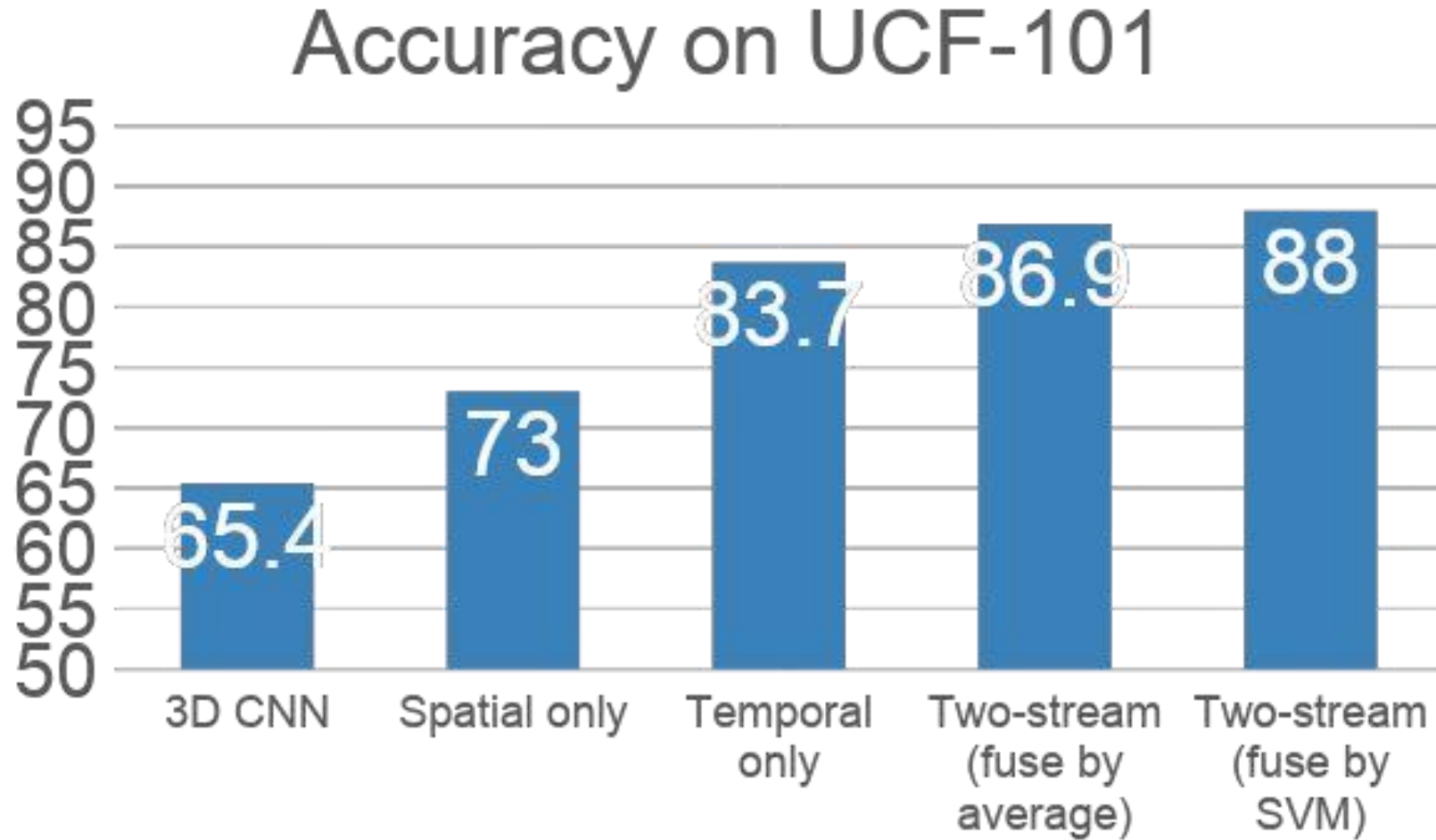
input
video

Input: Stack of optical flow:
 $[2*(T-1)] \times H \times W$

Early fusion: First 2D conv
processes all flow images

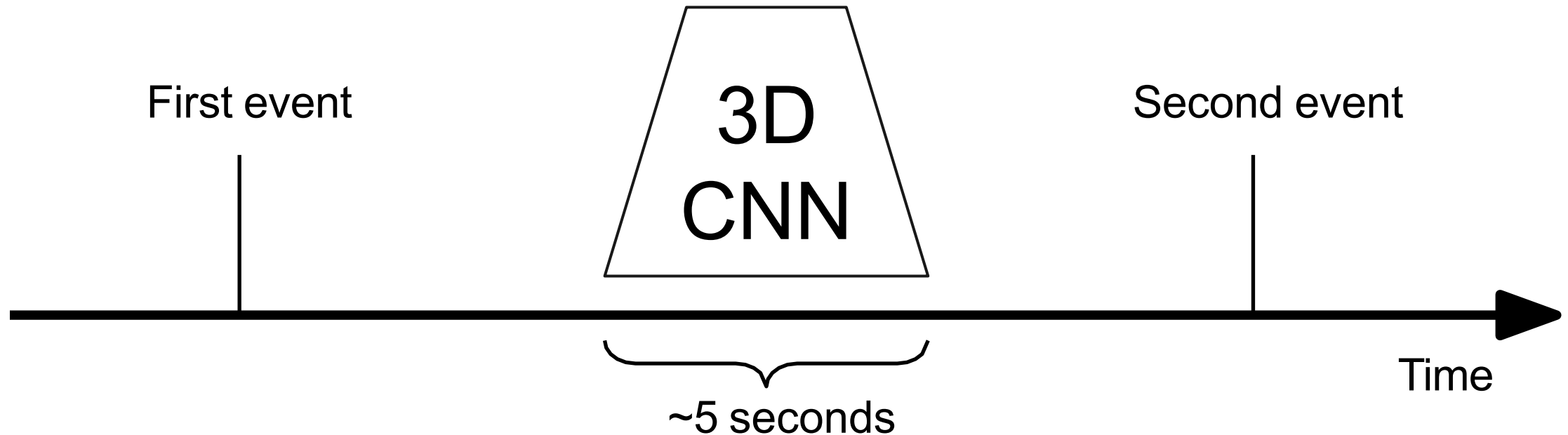


Separating Motion and Appearance: Two-Stream Networks



Modeling Long-term Temporal Structure

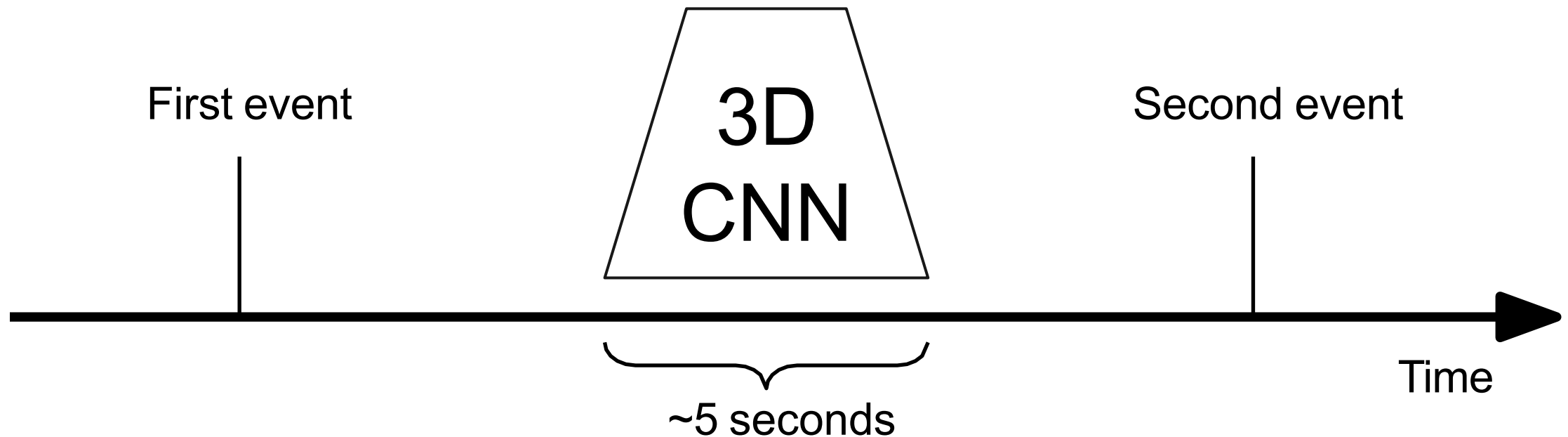
So far all our temporal CNNs only model local motion between frames in very short clips of ~2-5 seconds. What about long-term structure?



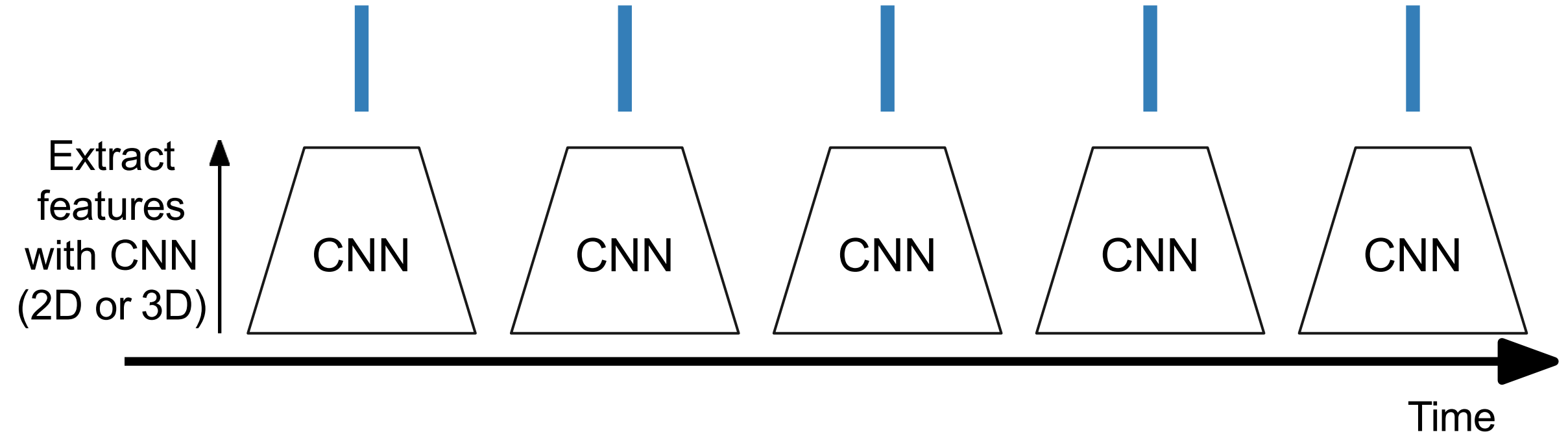
Modeling Long-term Temporal Structure

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?

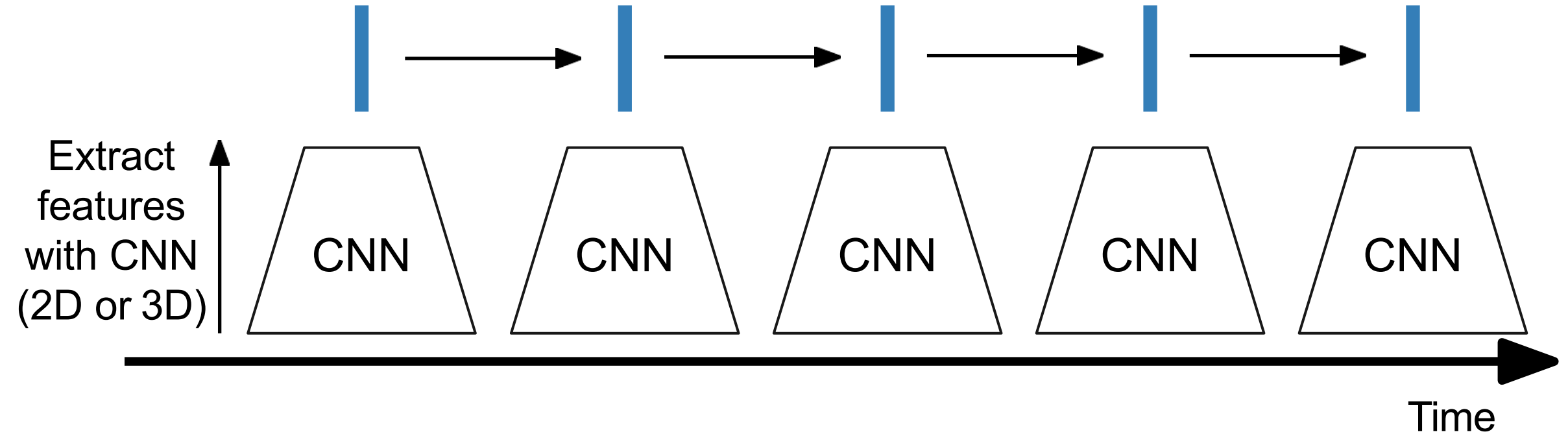


Modeling Long-term Temporal Structure



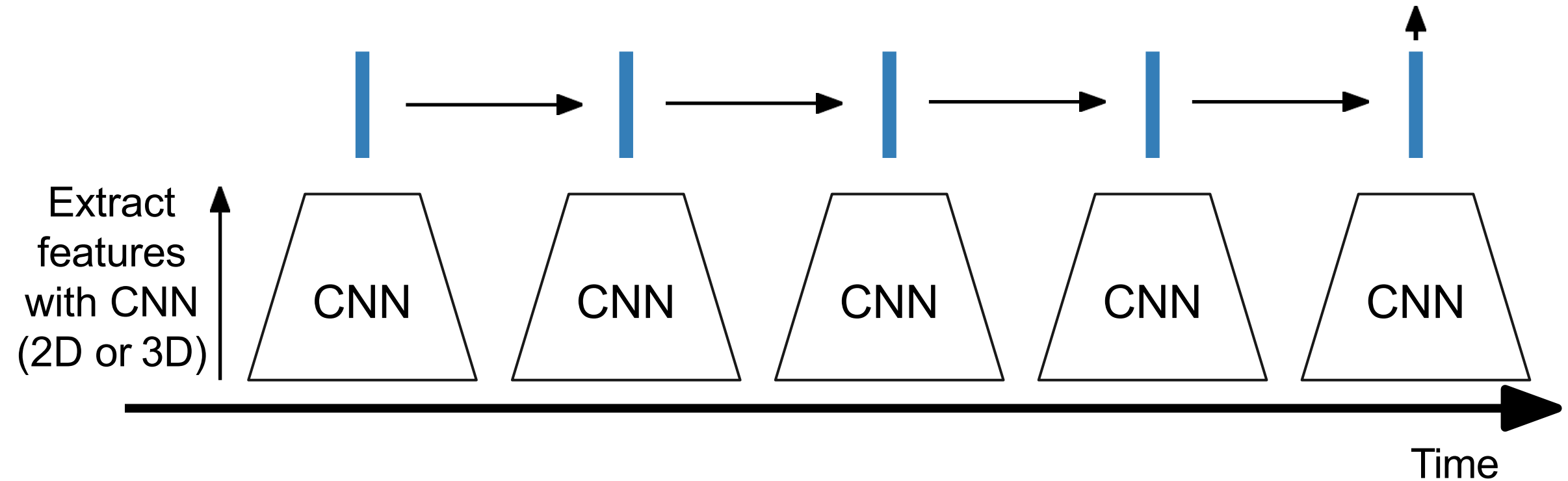
Modeling Long-term Temporal Structure

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



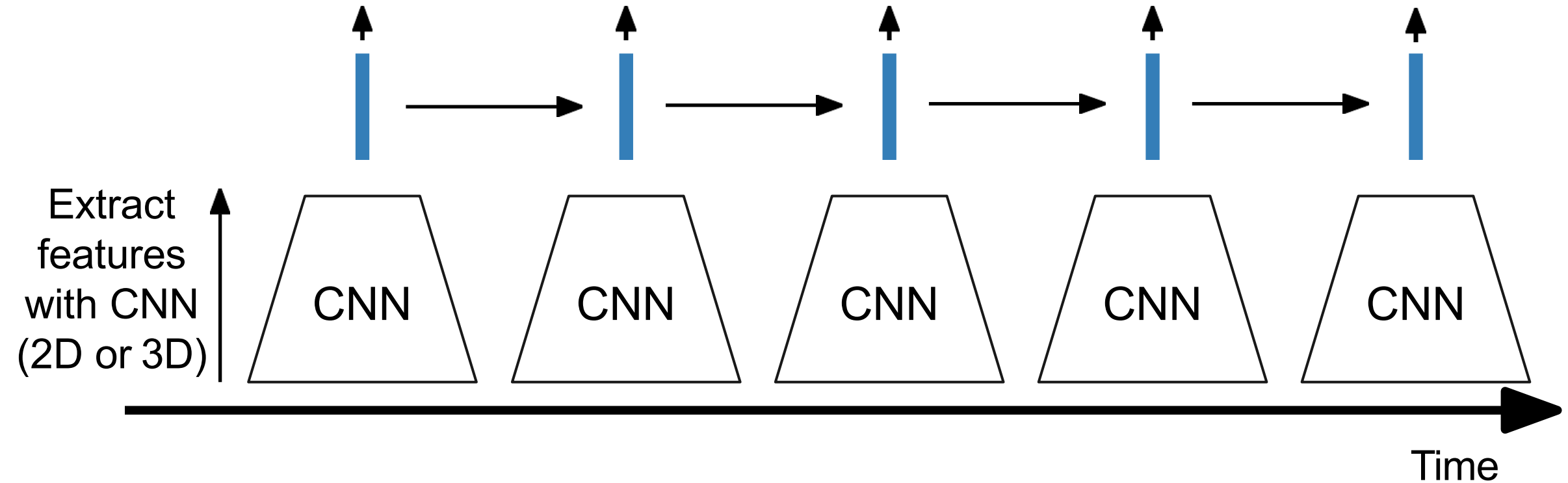
Modeling Long-term Temporal Structure

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



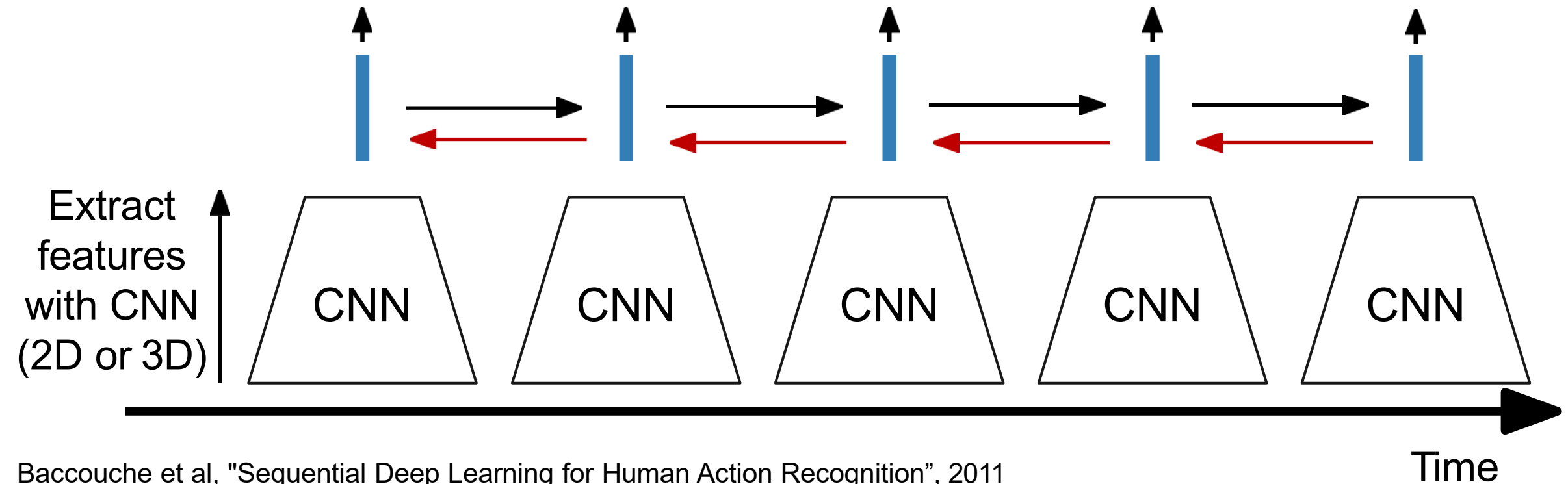
Modeling Long-term Temporal Structure

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

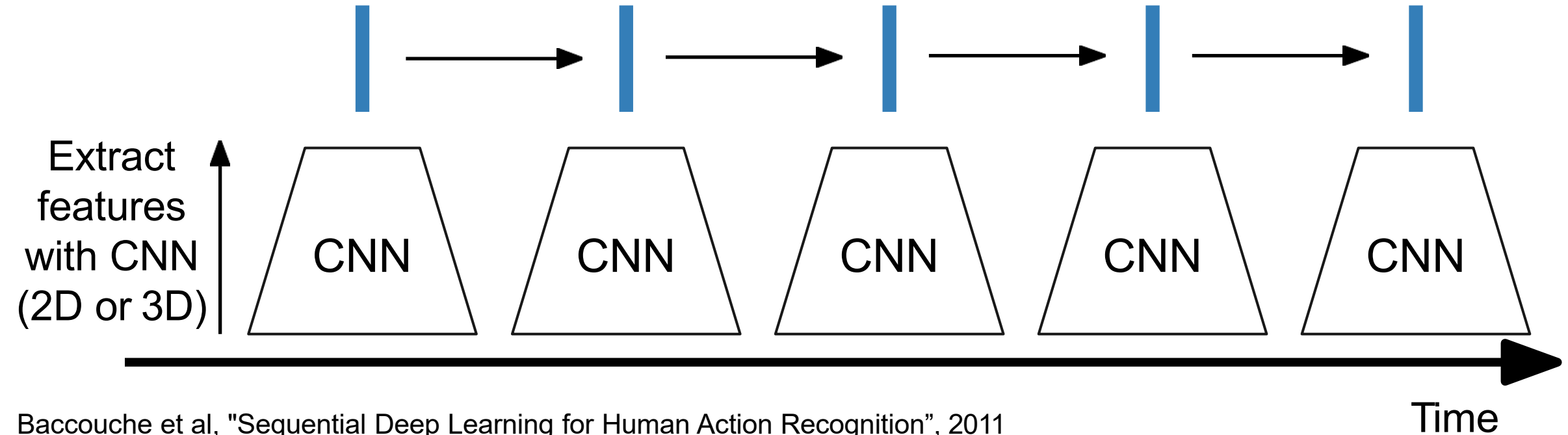


Modeling Long-term Temporal Structure

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?



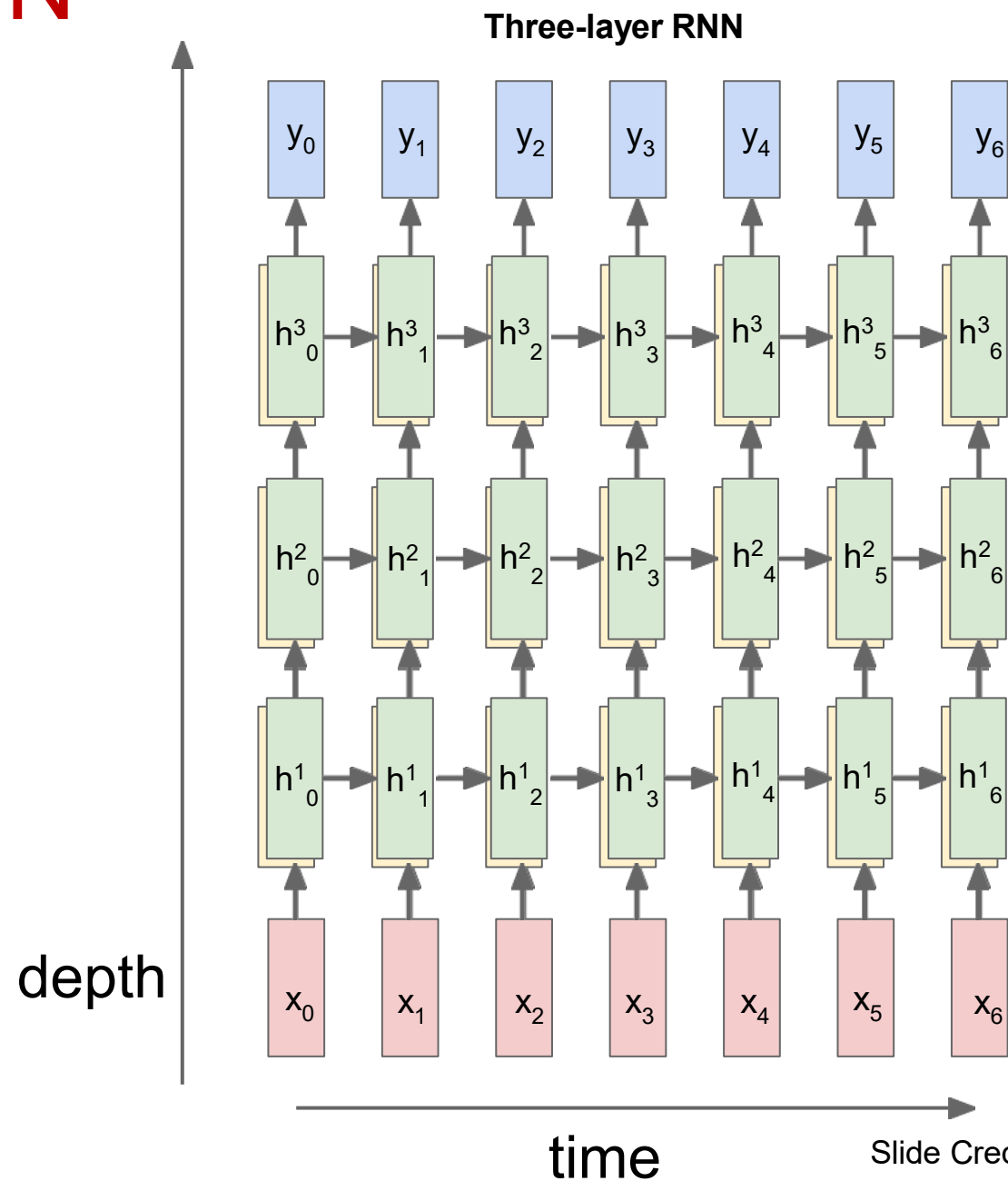
Baccouche et al, "Sequential Deep Learning for Human Action Recognition", 2011

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

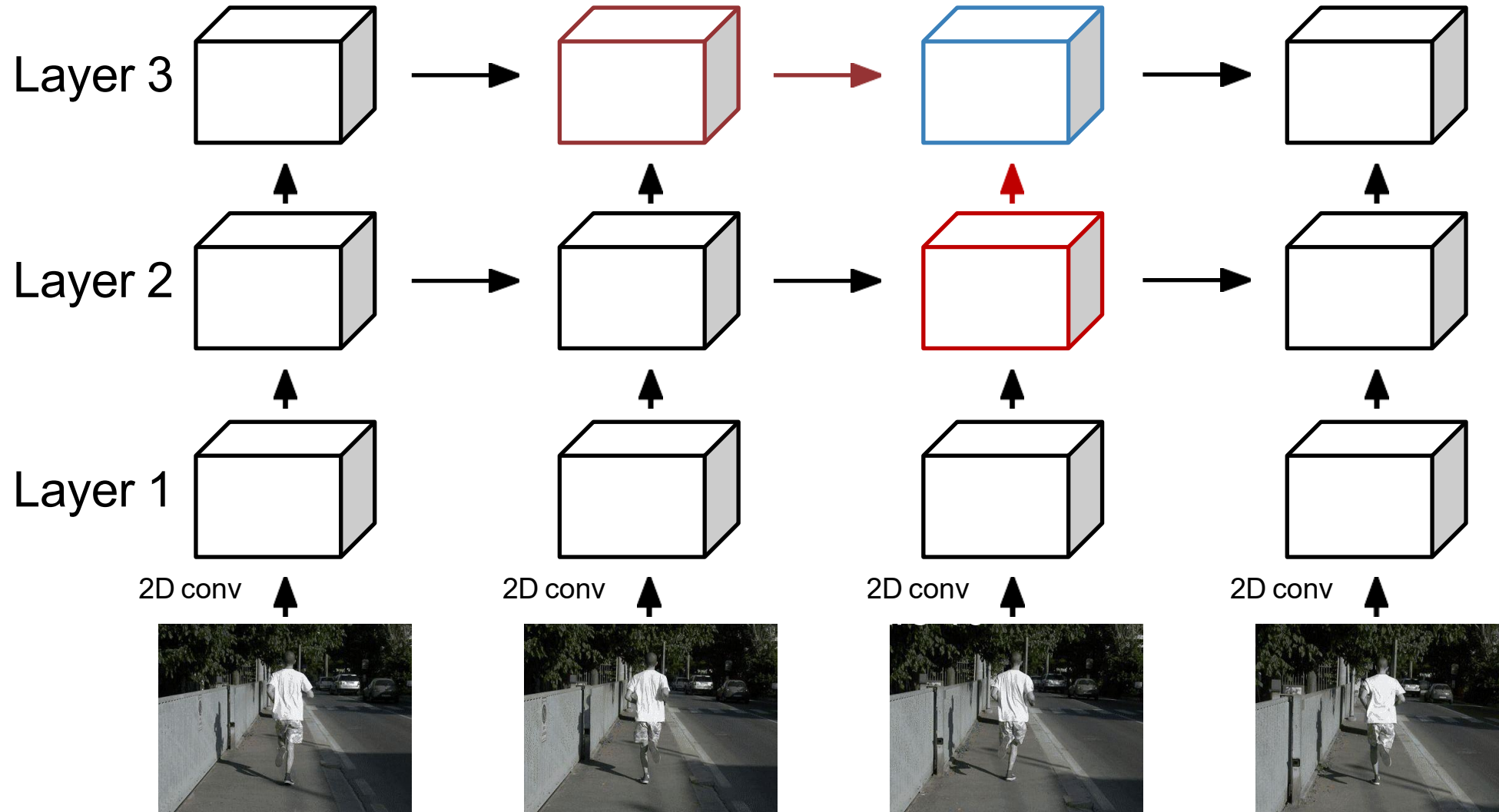


Recall: Multi-layer RNN

We can use a similar structure to process videos!



Recurrent Convolutional Network



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

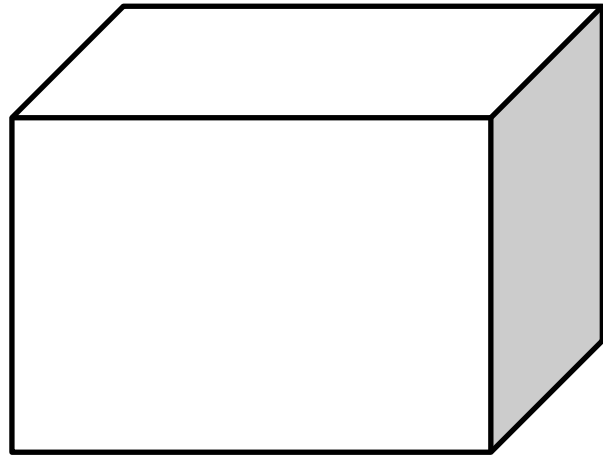
Each depends on two inputs:
1. Same layer, previous timestep
2. Prev layer, same timestep

Use different weights at each layer, share weights across time



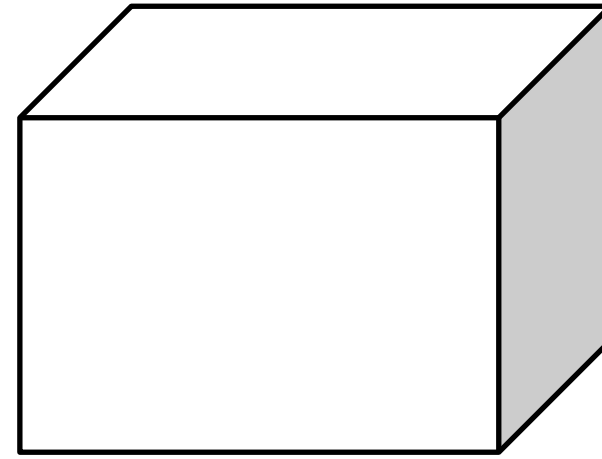
Recurrent Convolutional Network

Normal 2D CNN:



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

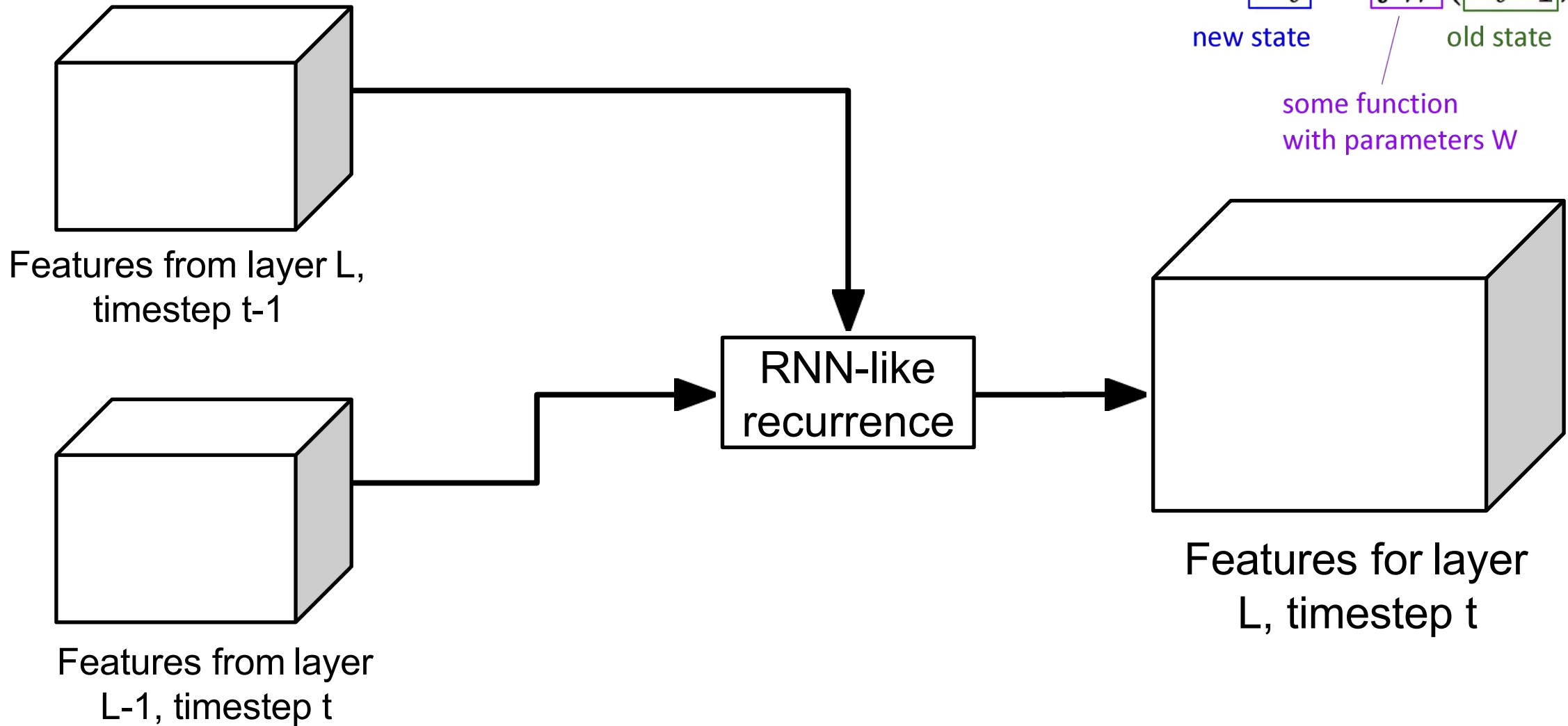
2D Conv
→



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



Recurrent Convolutional Network



Recall: Recurrent Network

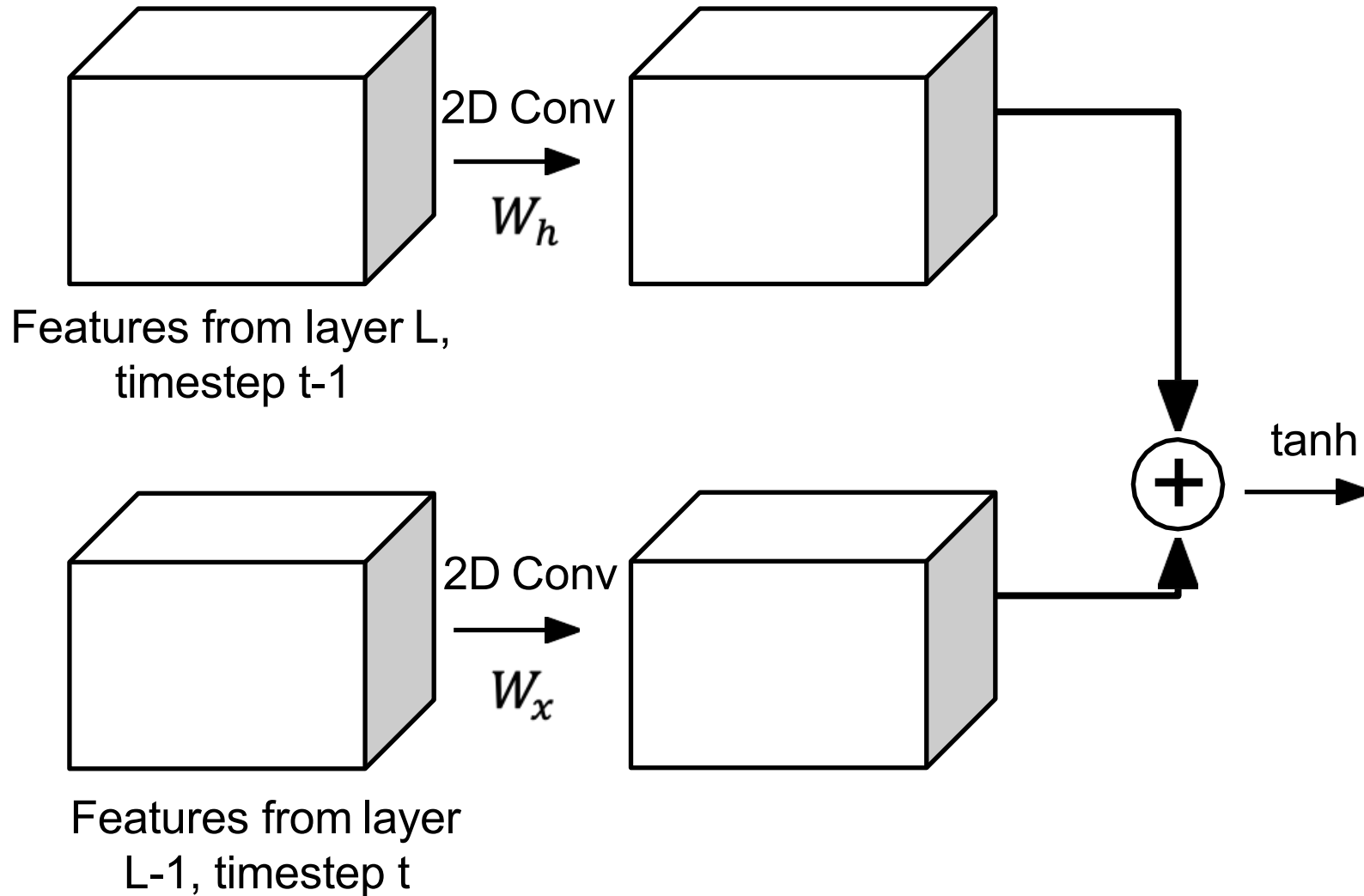
$$h_t = f_W(h_{t-1}, x_t)$$

new state / old state

some function with parameters W



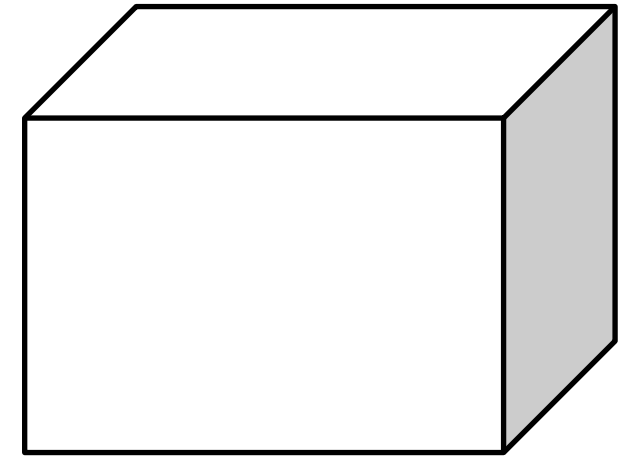
Recurrent Convolutional Network



Recall: Vanilla RNN

$$h_{t+1} = \tanh(W_h h_t + W_x x)$$

Replace all matrix multiply with 2D convolution!



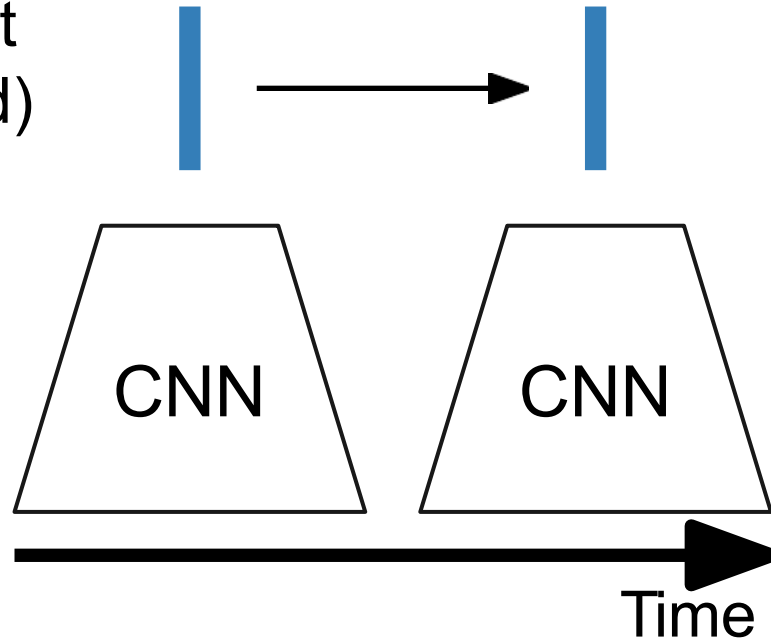
Features for layer L, timestep t



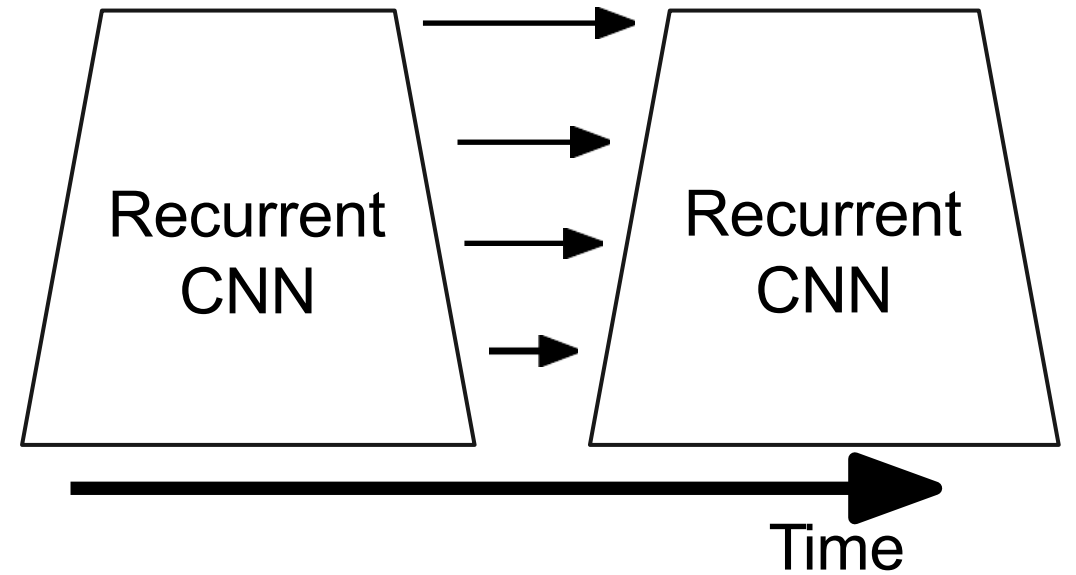
Modeling Long-term Temporal Structure

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

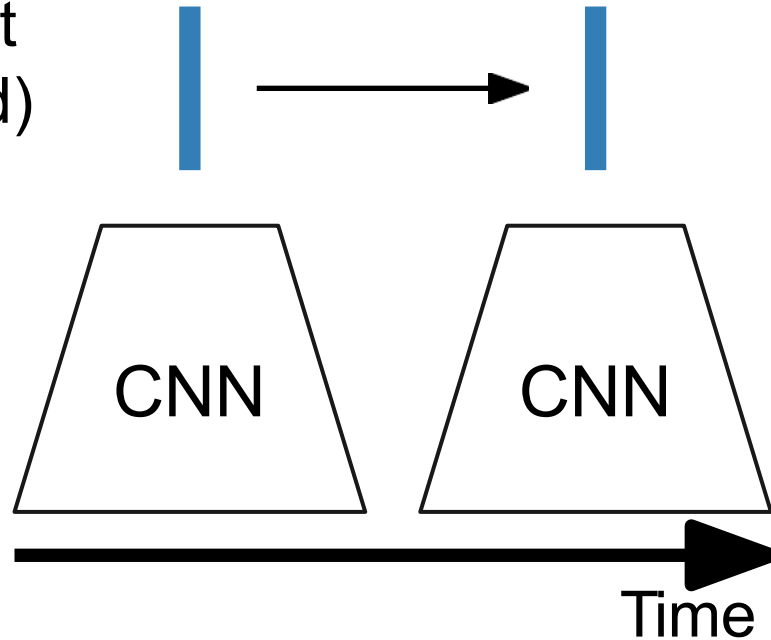
Ballas et al, "Delving Deeper into Convolutional Networks for Learning Video Representations", ICLR 2016



Modeling Long-term Temporal Structure

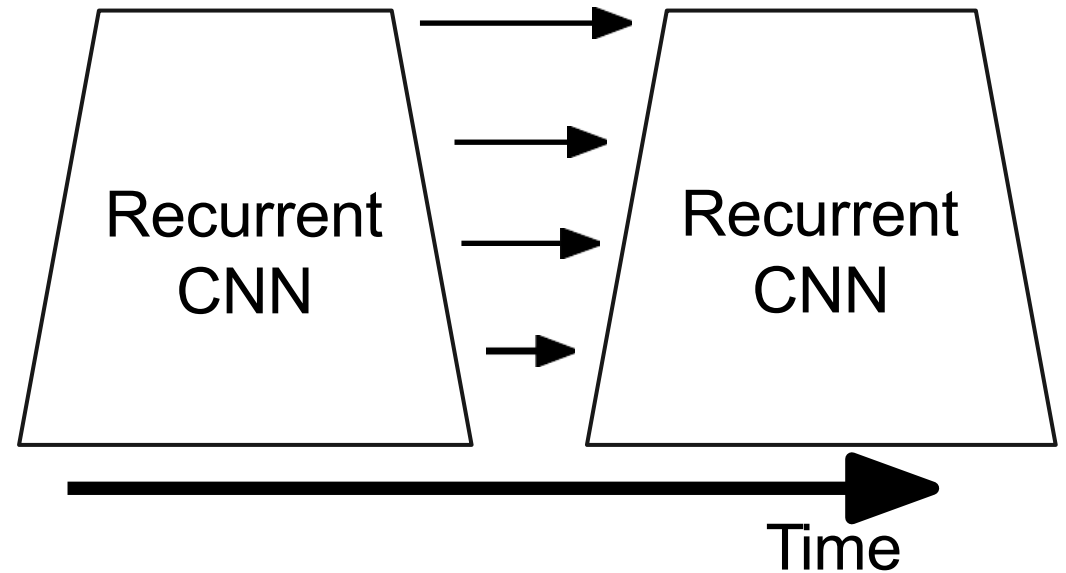
Problem: RNNs are slow for long sequences (can't be parallelized)

RNN: Infinite temporal extent (fully-connected)

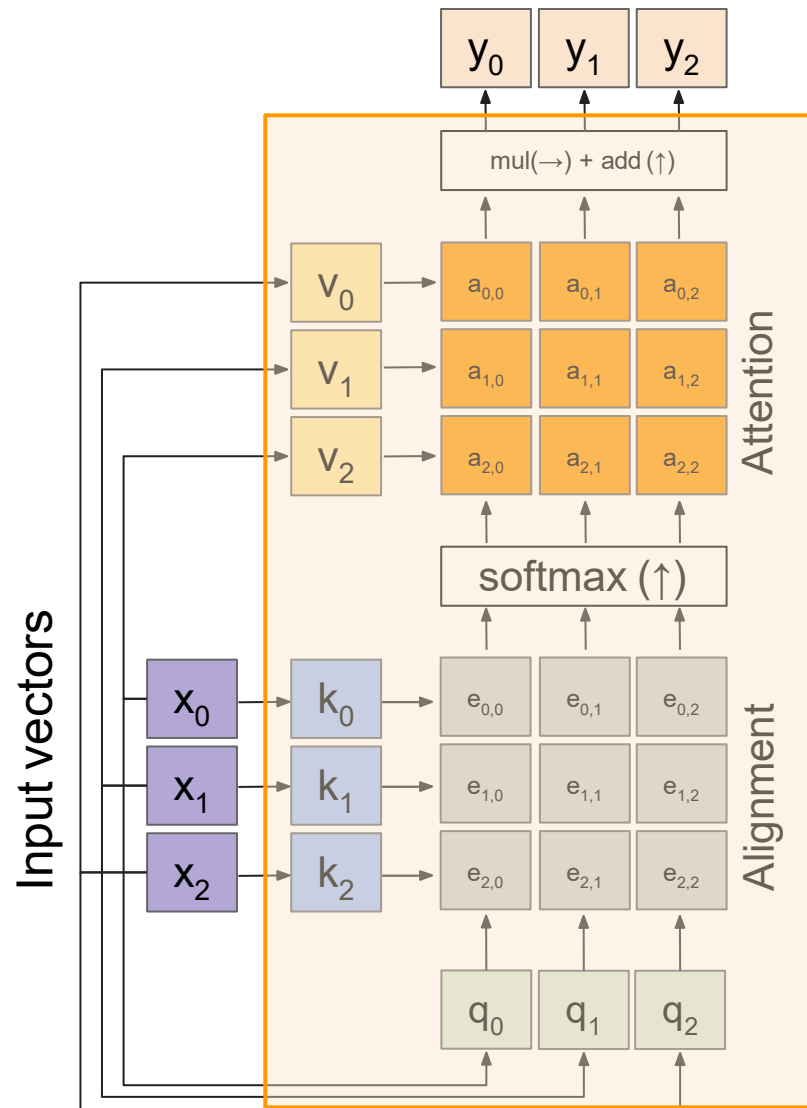


CNN: finite temporal extent (convolutional)

Recurrent CNN: Infinite temporal extent (convolutional)



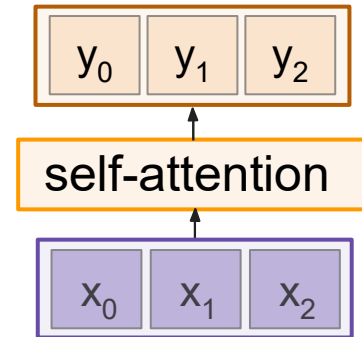
Recall: Self-Attention



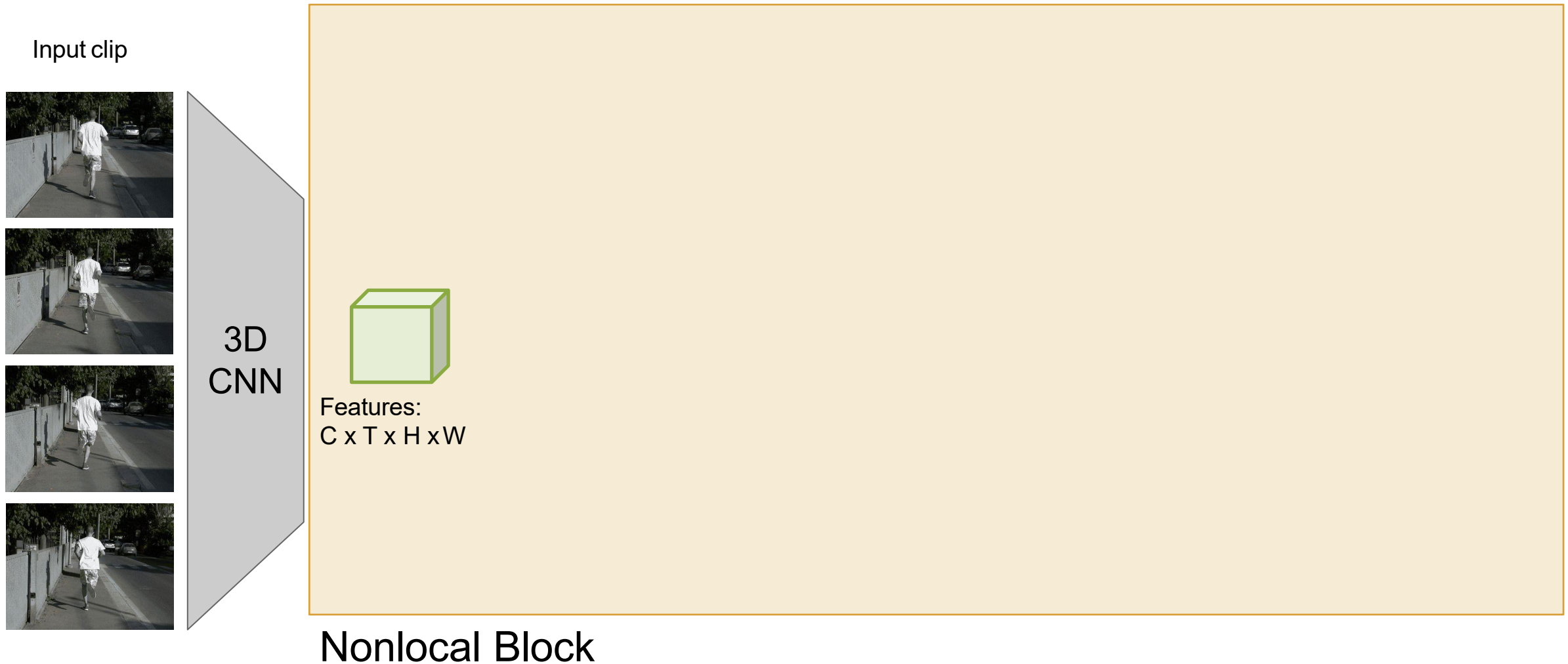
Outputs:
context vectors: \mathbf{y} (shape: D_v)

Operations:
 Key vectors: $\mathbf{k} = \mathbf{x}\mathbf{W}_k$
 Value vectors: $\mathbf{v} = \mathbf{x}\mathbf{W}_v$
 Query vectors: $\mathbf{q} = \mathbf{x}\mathbf{W}_q$
 Alignment: $e_{i,j} = \mathbf{q}_j \cdot \mathbf{k}_i / \sqrt{D}$
 Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$
 Output: $y_j = \sum_i a_{i,j} v_i$

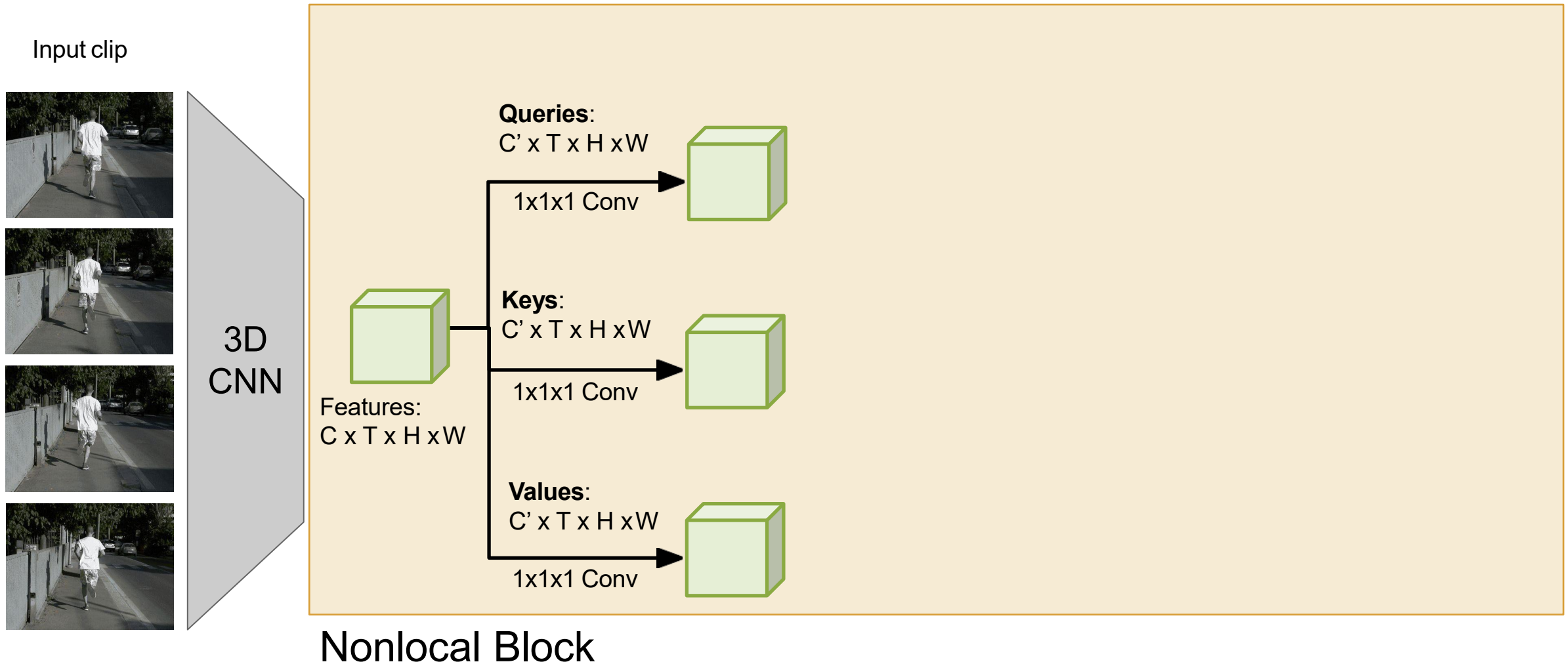
Inputs:
Input vectors: \mathbf{x} (shape: $N \times D$)



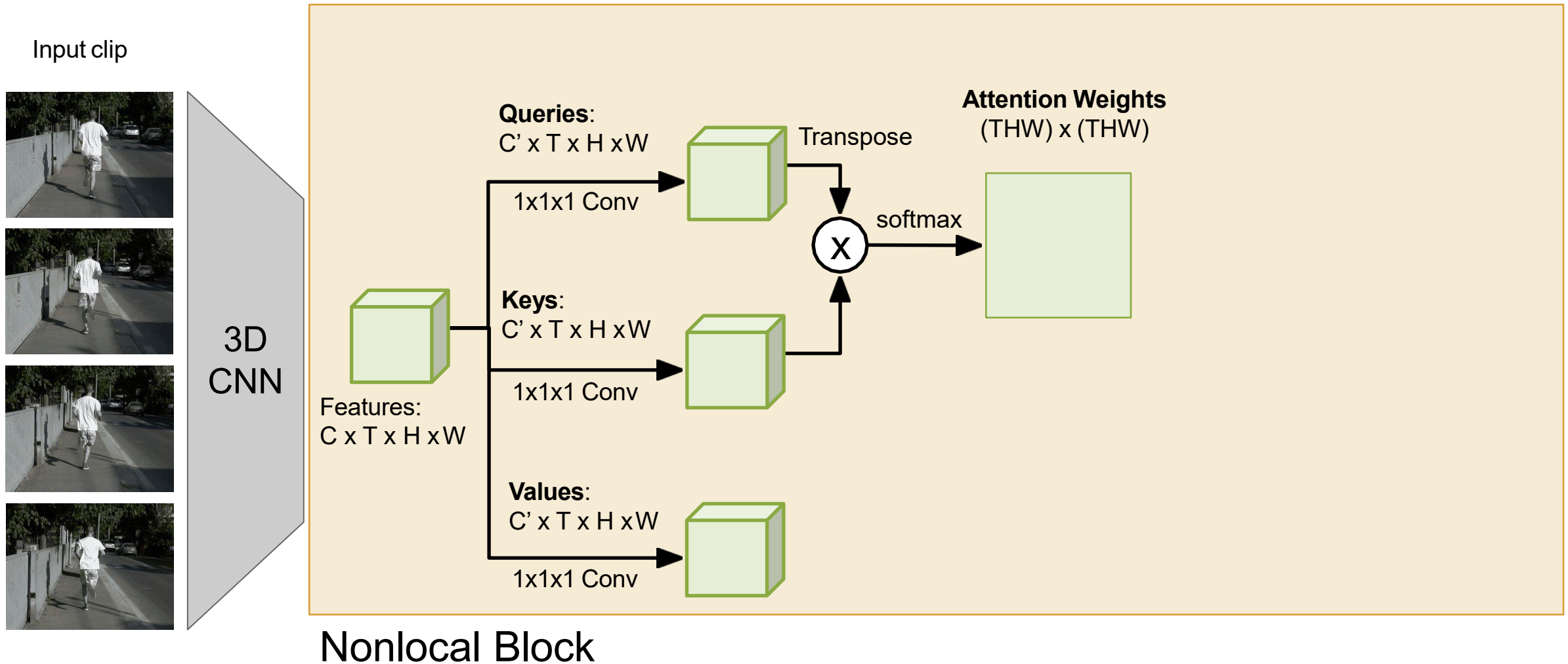
Spatio-Temporal Self-Attention (Nonlocal Block)



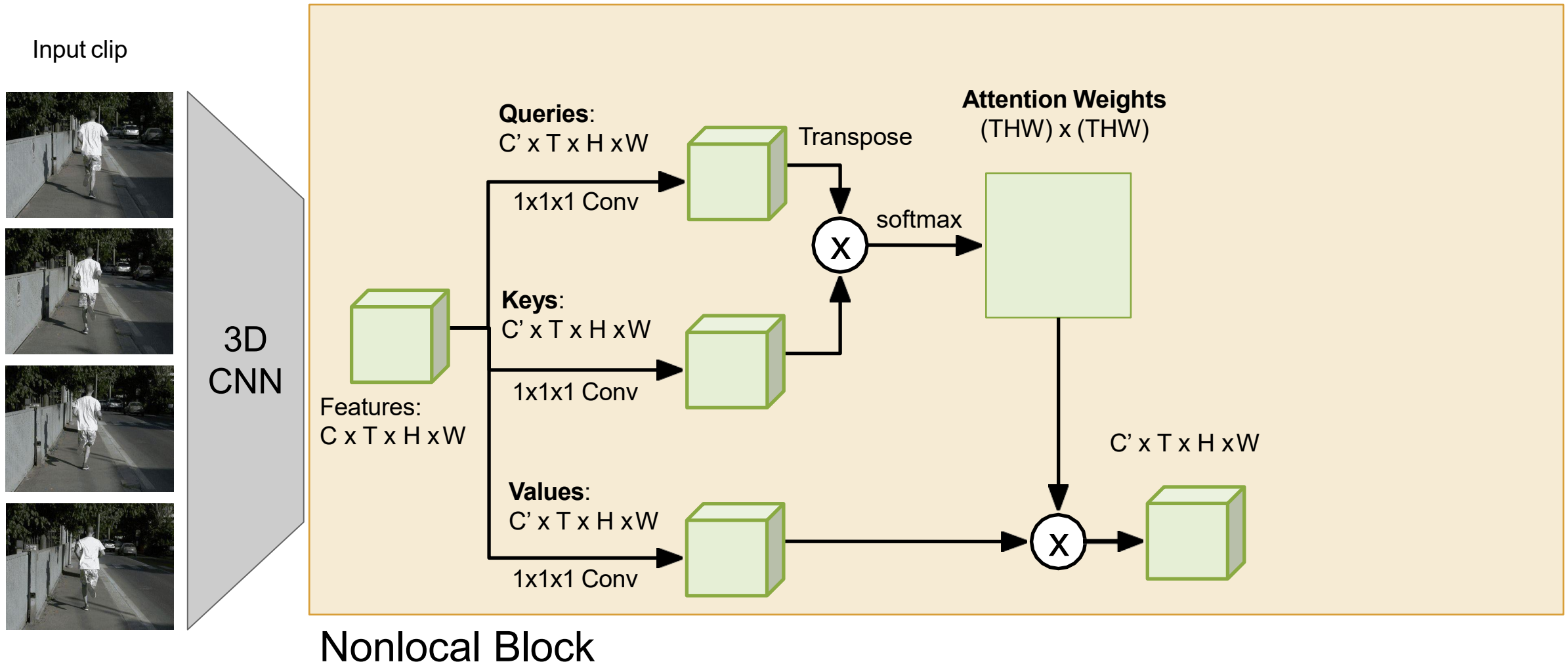
Spatio-Temporal Self-Attention (Nonlocal Block)



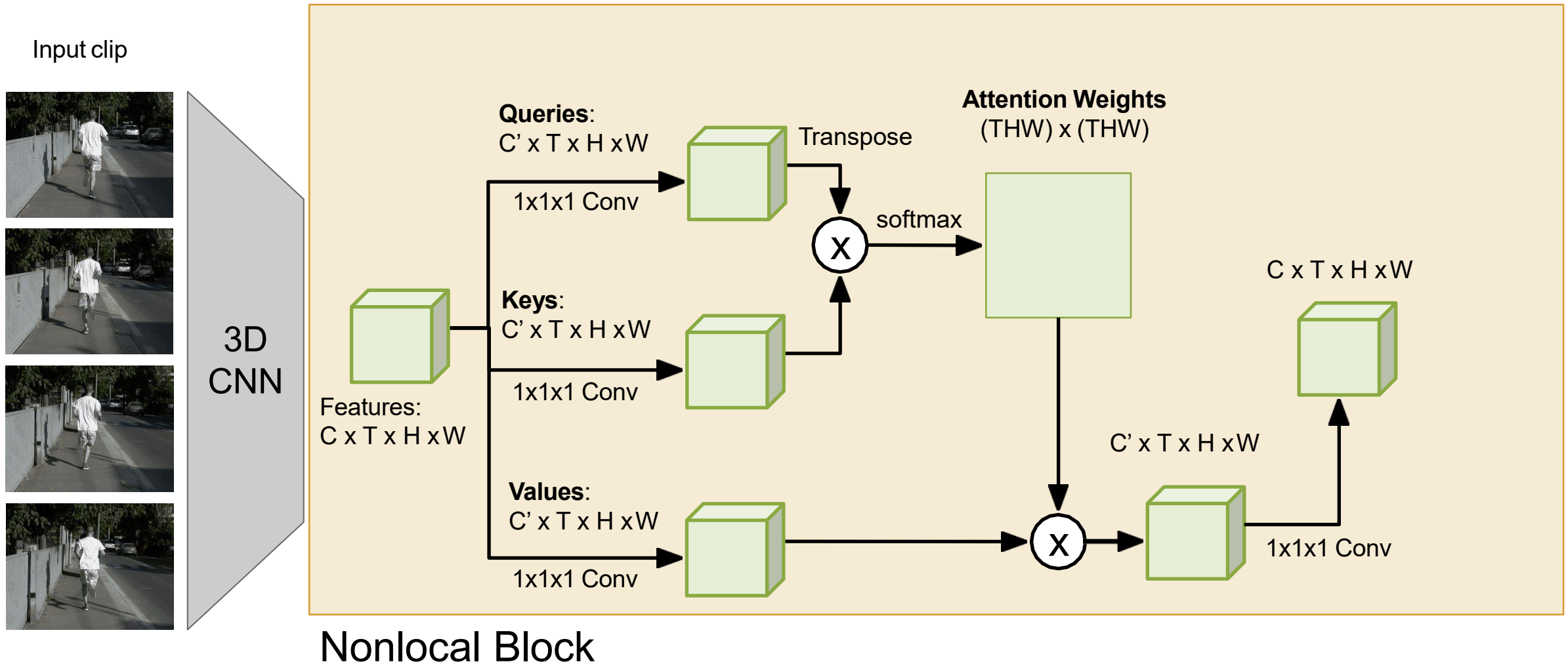
Spatio-Temporal Self-Attention (Nonlocal Block)



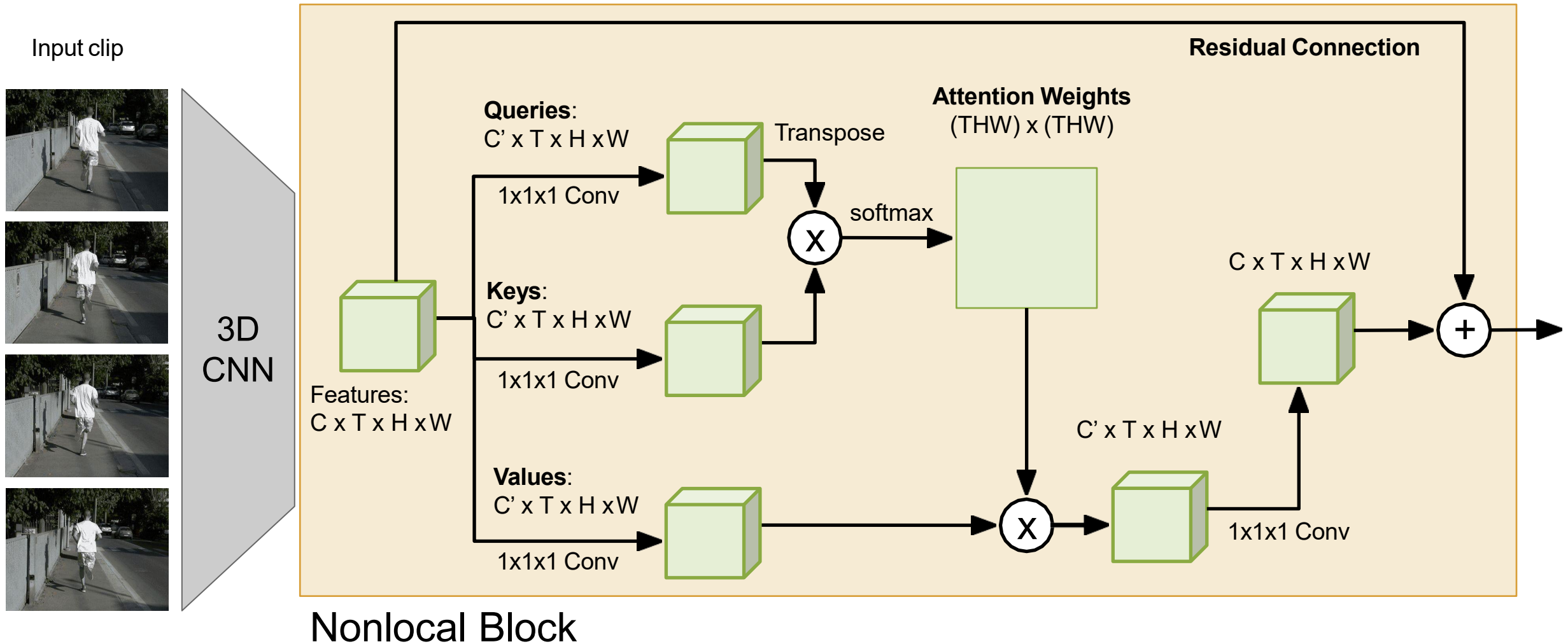
Spatio-Temporal Self-Attention (Nonlocal Block)



Spatio-Temporal Self-Attention (Nonlocal Block)



Spatio-Temporal Self-Attention (Nonlocal Block)

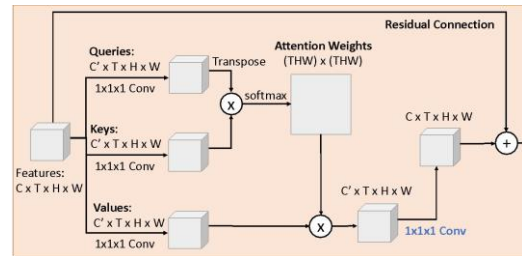
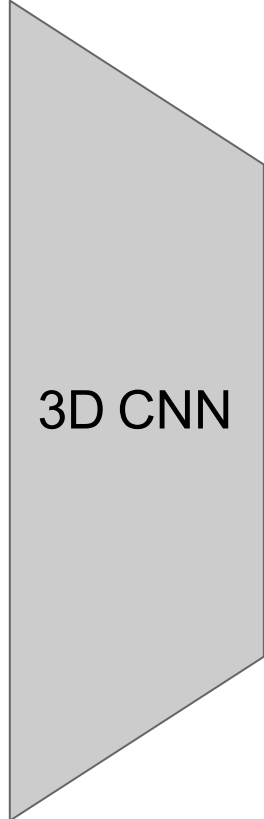


Spatio-Temporal Self-Attention (Nonlocal Block)

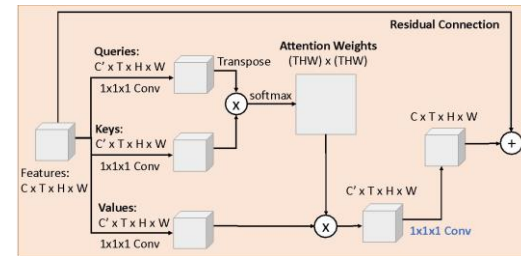
We can add nonlocal blocks into existing 3D CNN architectures.

But what is the best 3D CNN architecture?

Input clip



3D CNN



3D CNN

Running



Inflating 2D Networks to 3D (I3D)

We can add nonlocal blocks into existing 3D CNN architectures.

But what is the best 3D CNN architecture?

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



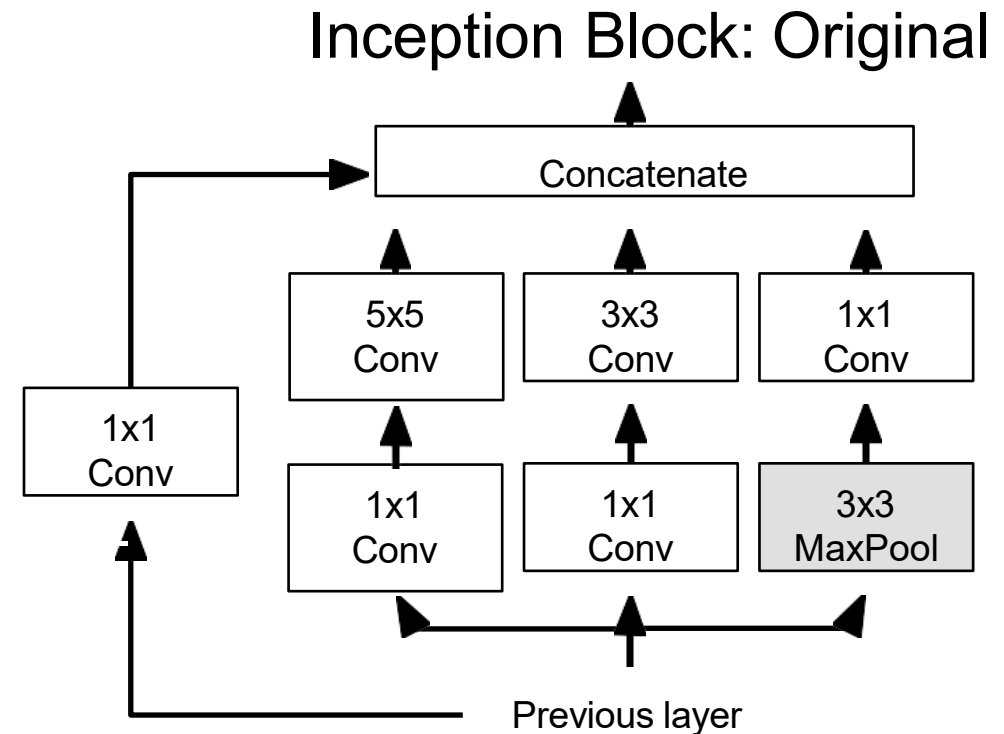
Inflating 2D Networks to 3D (I3D)

We can add nonlocal blocks into existing 3D CNN architectures.

But what is the best 3D CNN architecture?

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



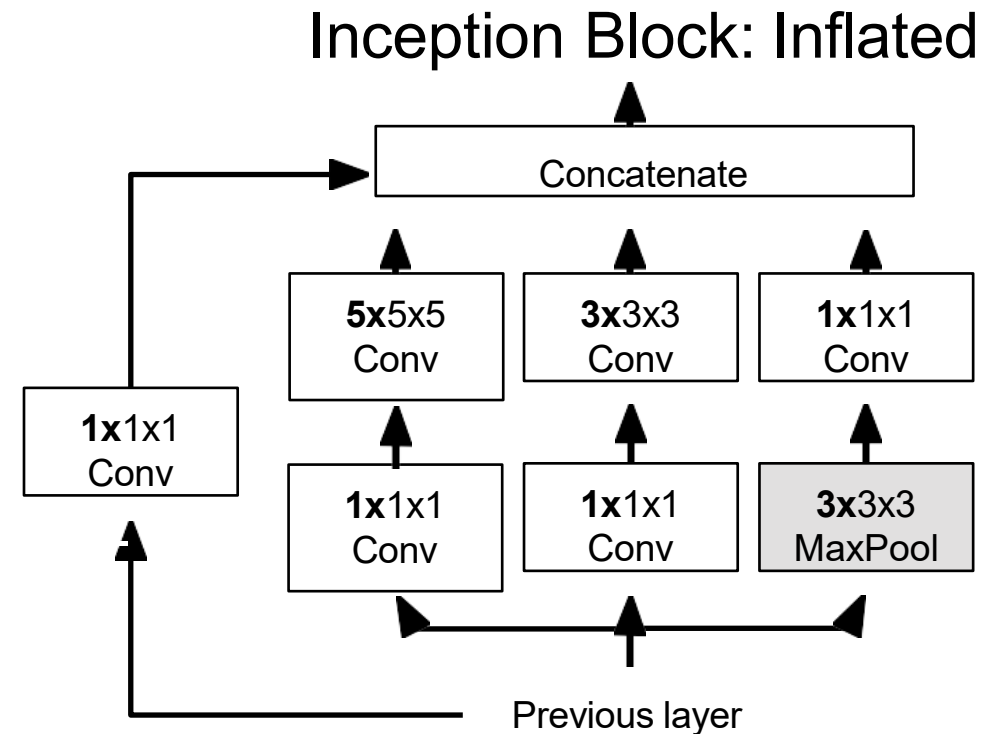
Inflating 2D Networks to 3D (I3D)

We can add nonlocal blocks into existing 3D CNN architectures.

But what is the best 3D CNN architecture?

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)

We can add nonlocal blocks into existing 3D CNN architectures.

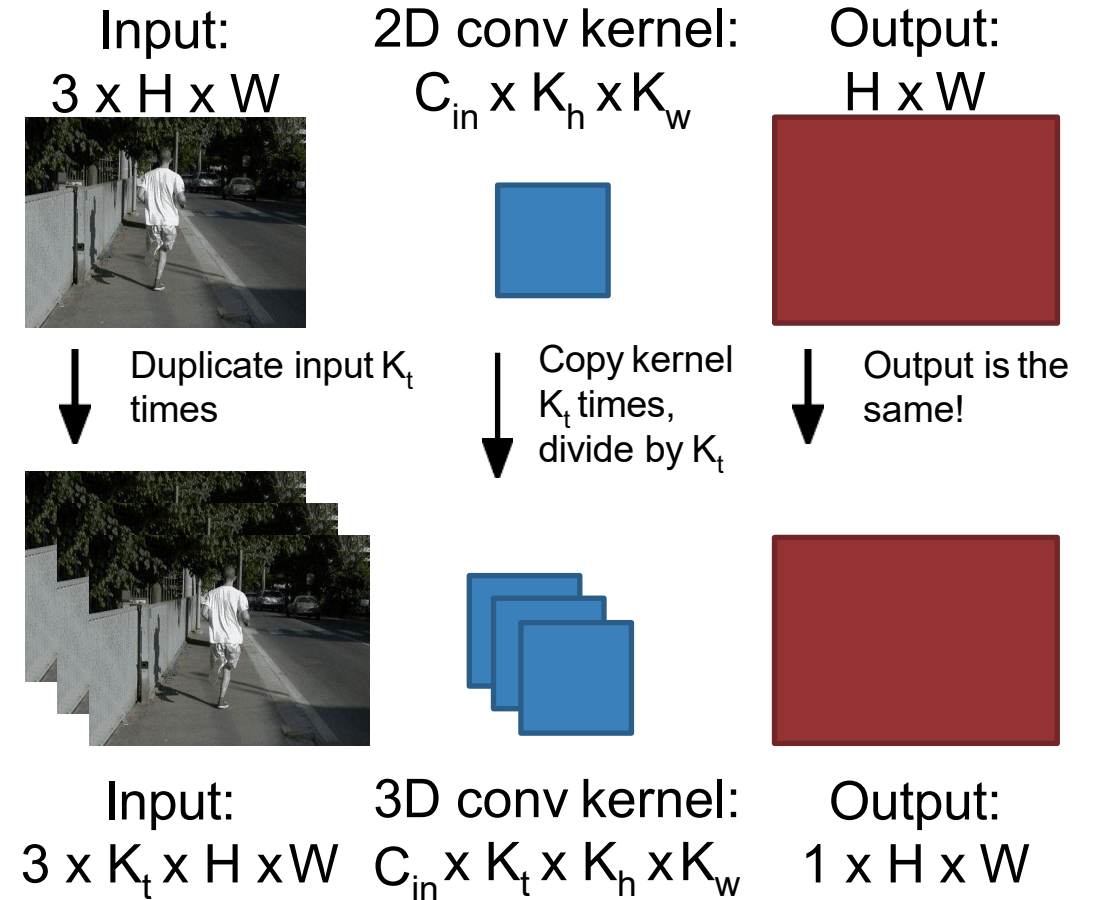
But what is the best 3D CNN architecture?

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)

We can add nonlocal blocks into existing 3D CNN architectures.

But what is the best 3D CNN architecture?

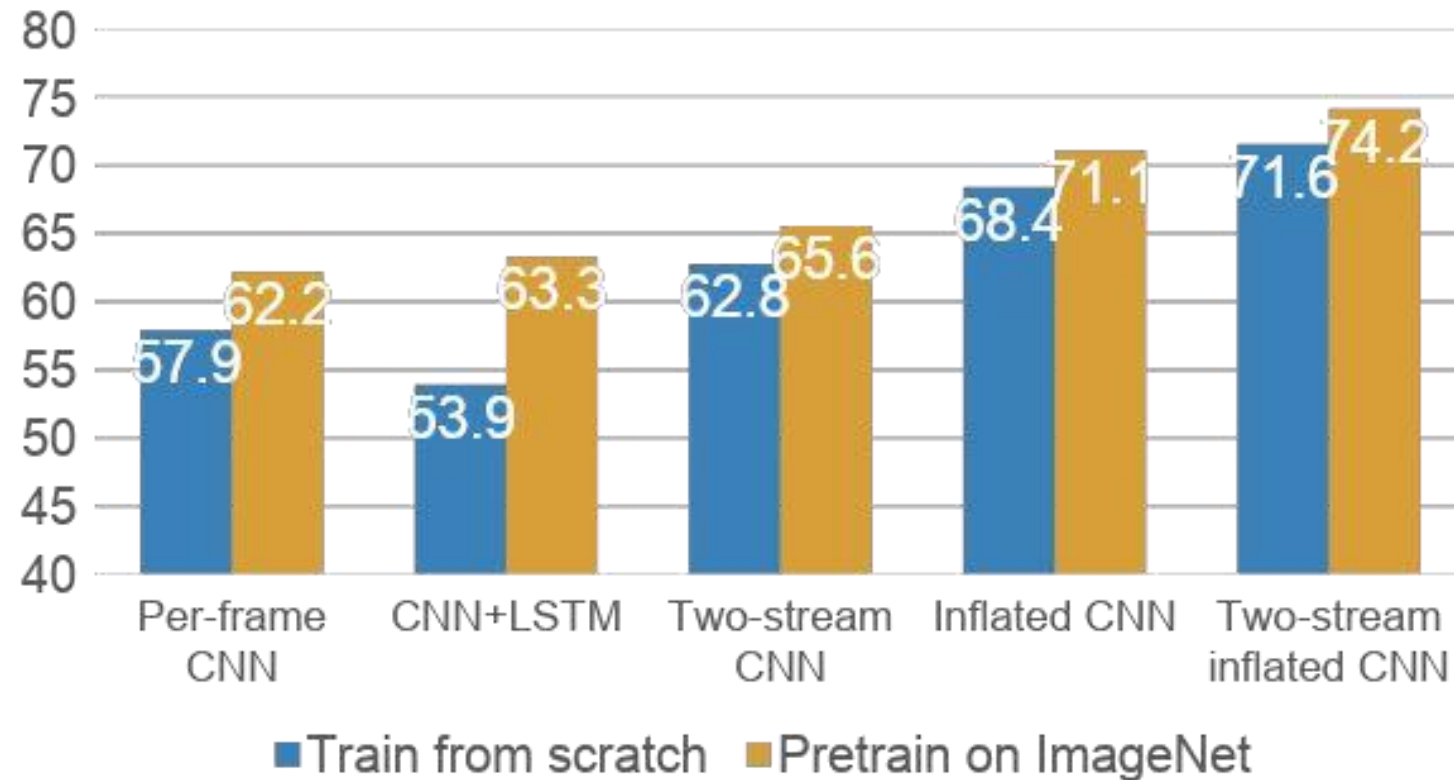
Idea: take a 2D CNN architecture.

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



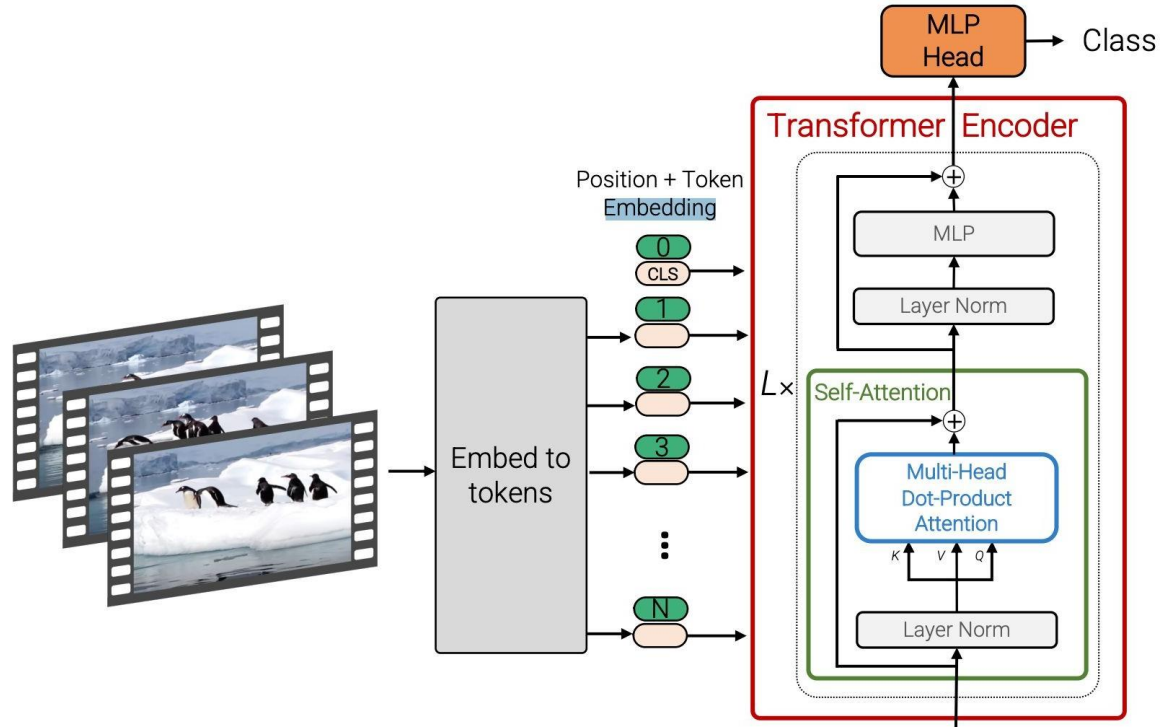
All using Inception CNN

Slide Credit: cs231n



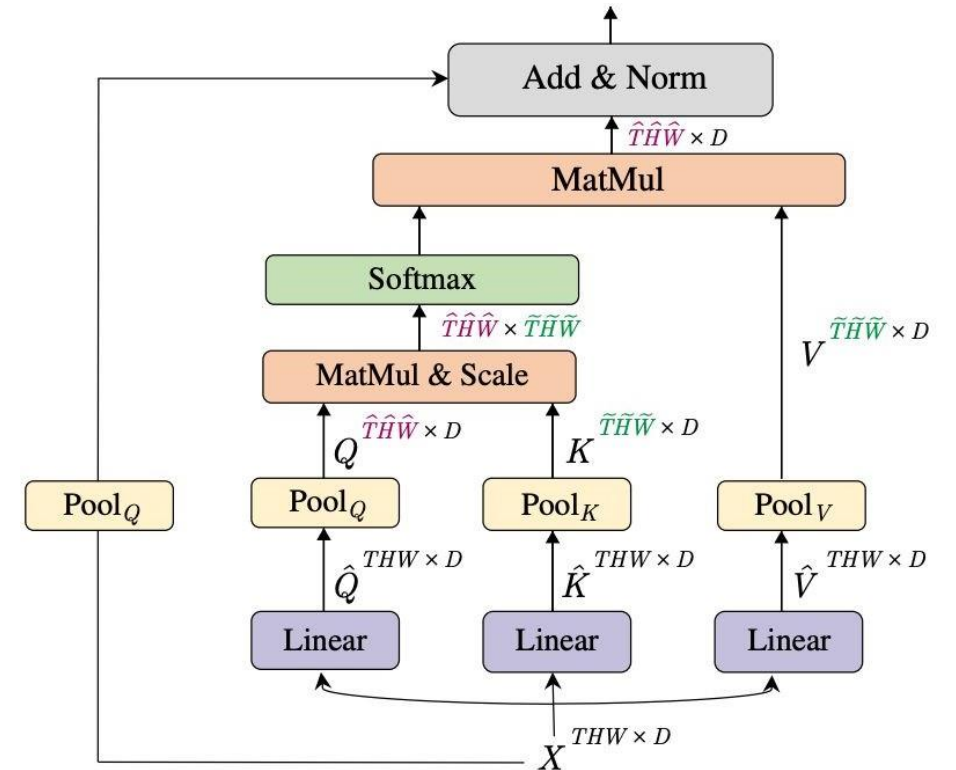
Vision Transformers for Video

Factorized attention: Attend over space / time



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

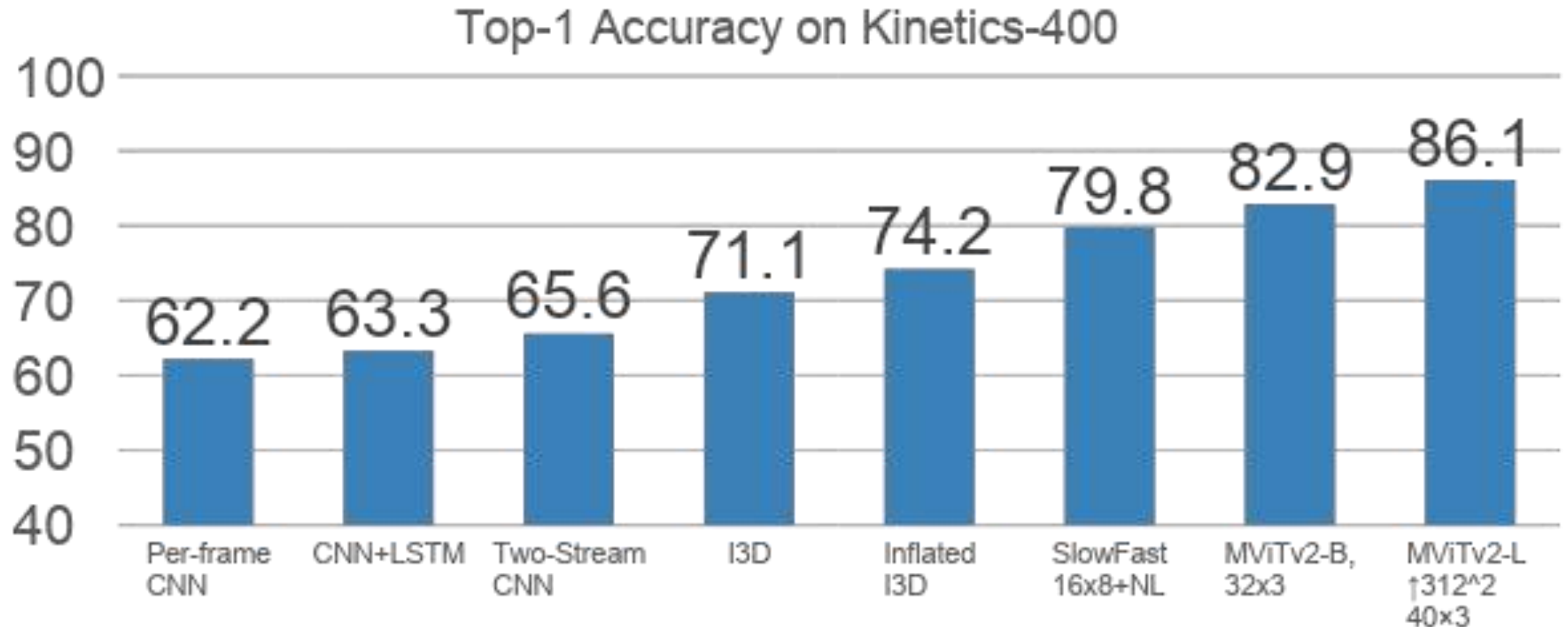
Pooling module: Reduce number of tokens



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



Vision Transformers for Video



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



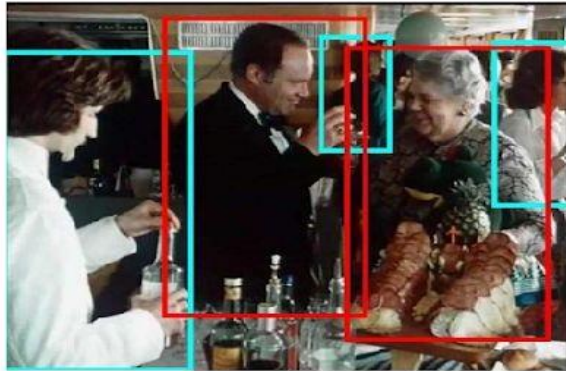
Can use architecture similar to Faster R-CNN:
first generate **temporal proposals** then **classify**



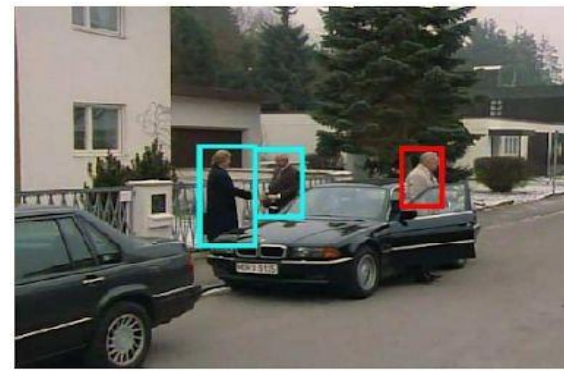
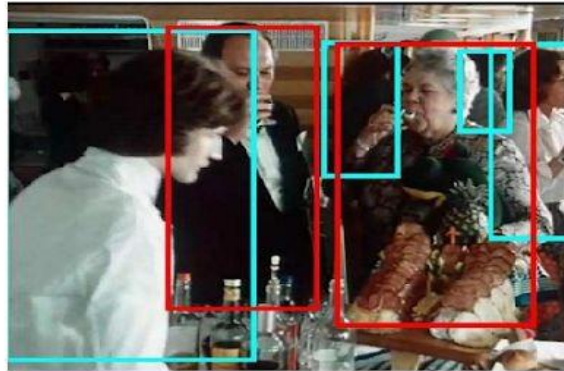
Spatio-Temporal Detection

Given a long untrimmed video, detect all the people in both space and time and classify the activities they are performing.

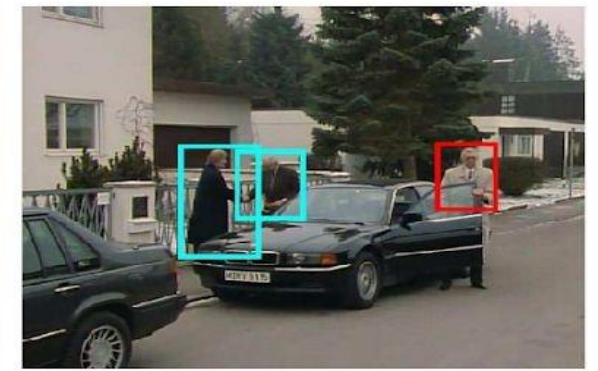
Some examples from AVADataset:



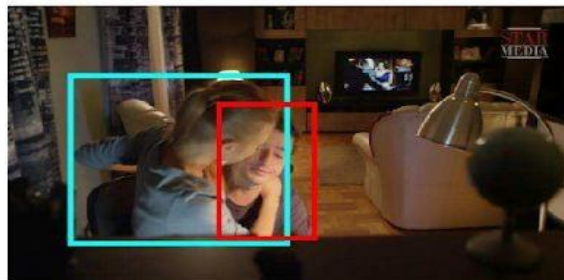
clink glass → drink



open → close



grab (a person) → hug



look at phone → answer phone



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]



Musical Instruments Source Separation

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

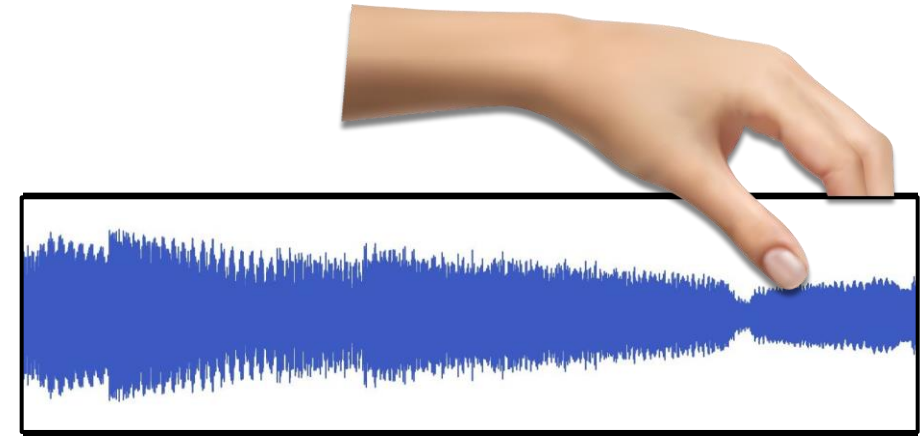


original video
(before separation)

object detections:
violin & flute

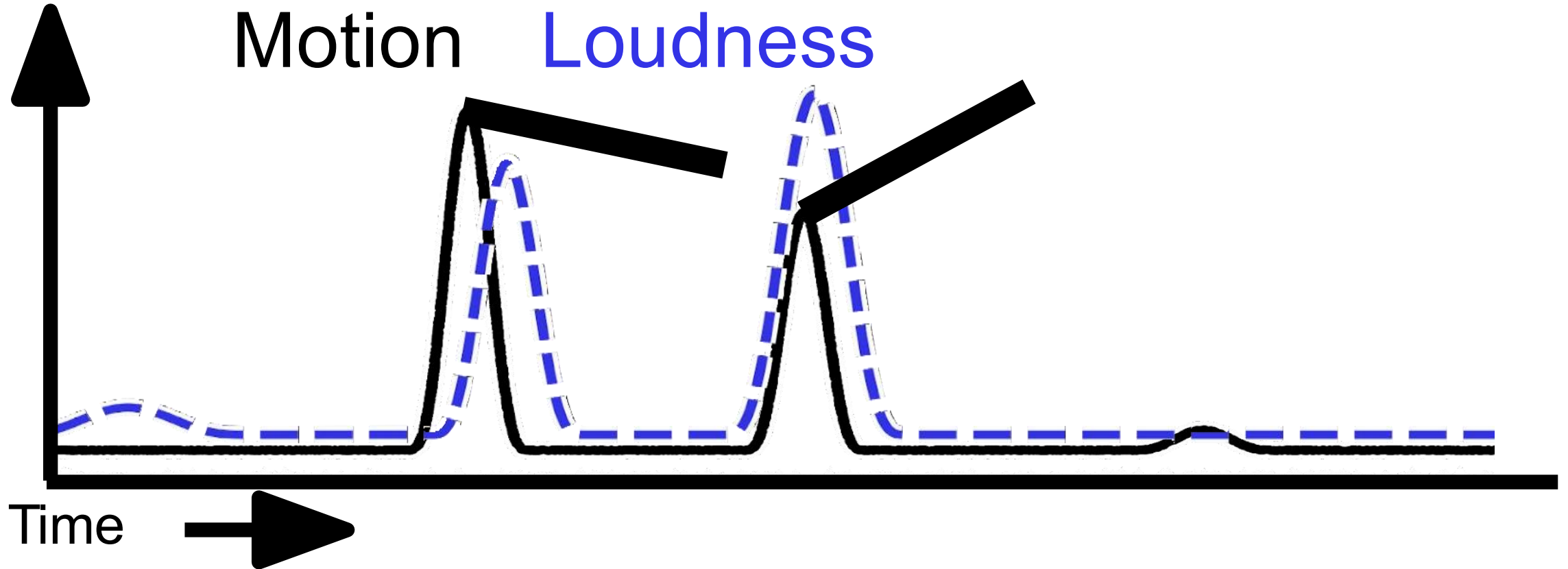


Learning Audio-Visual Synchronization



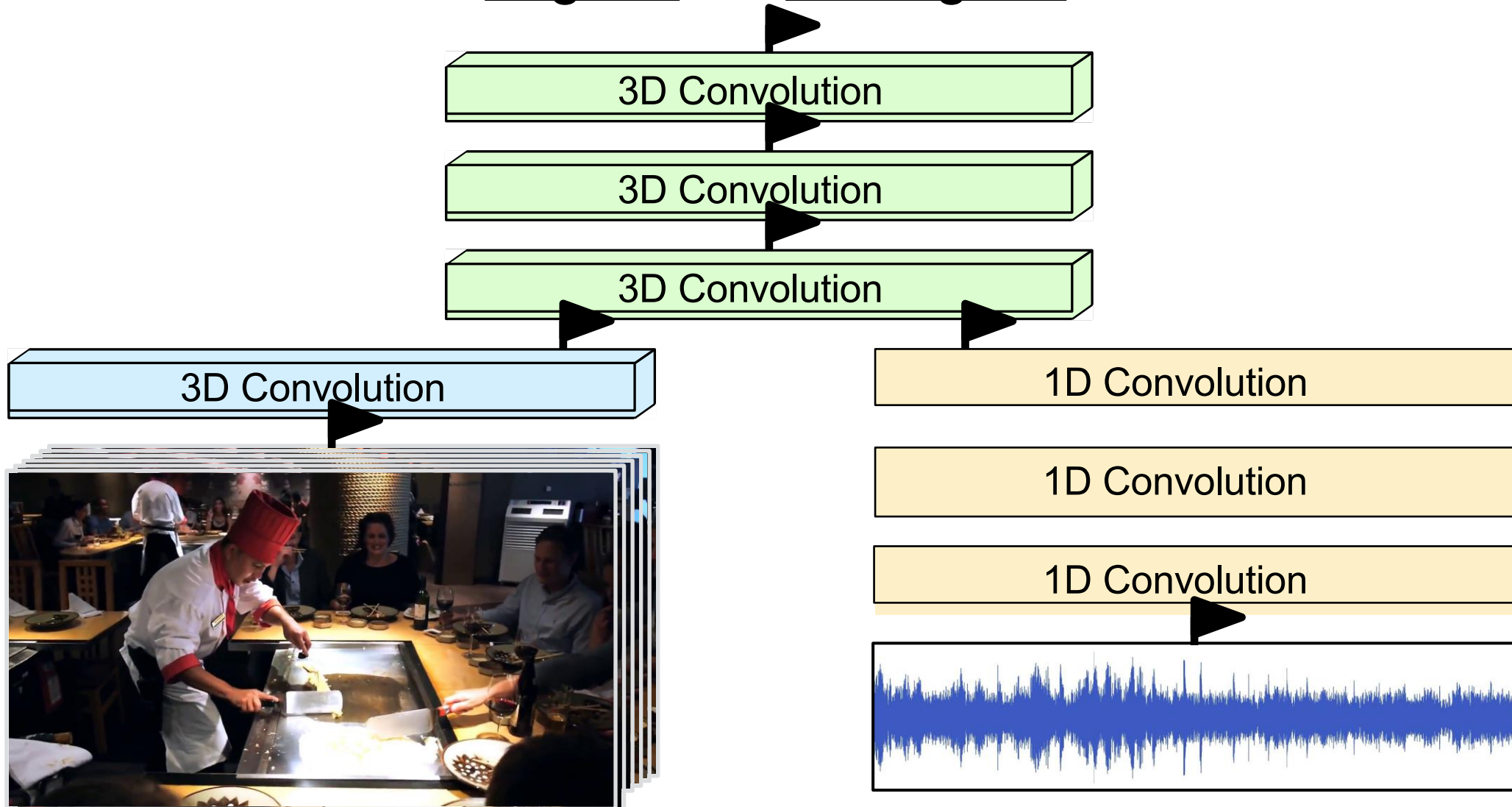
Owens & Efros, Audio-visual scene analysis with self-supervised multisensory features, ECCV 2018
Korbar et al., Co-training of audio and video representations from self-supervised temporal synchronization, NeurIPS 2018

Learning Audio-Visual Synchronization

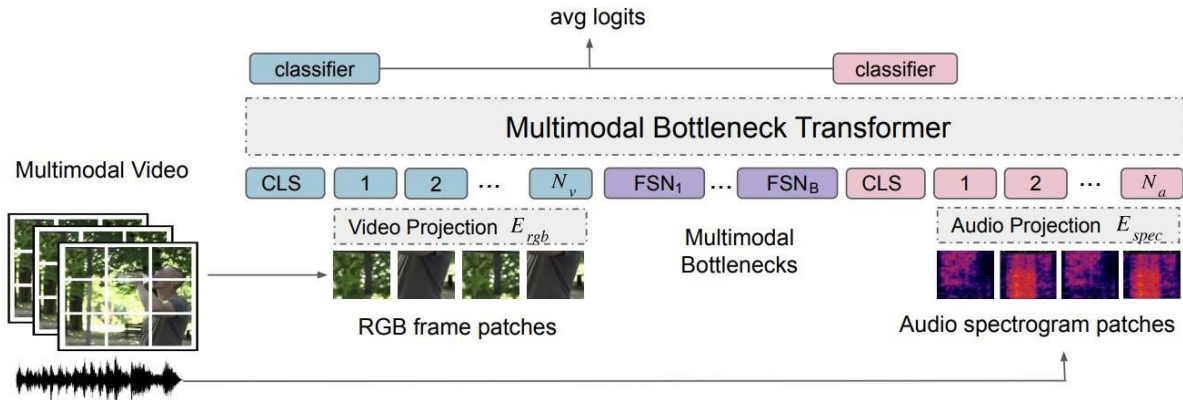


Learning Audio-Visual Synchronization

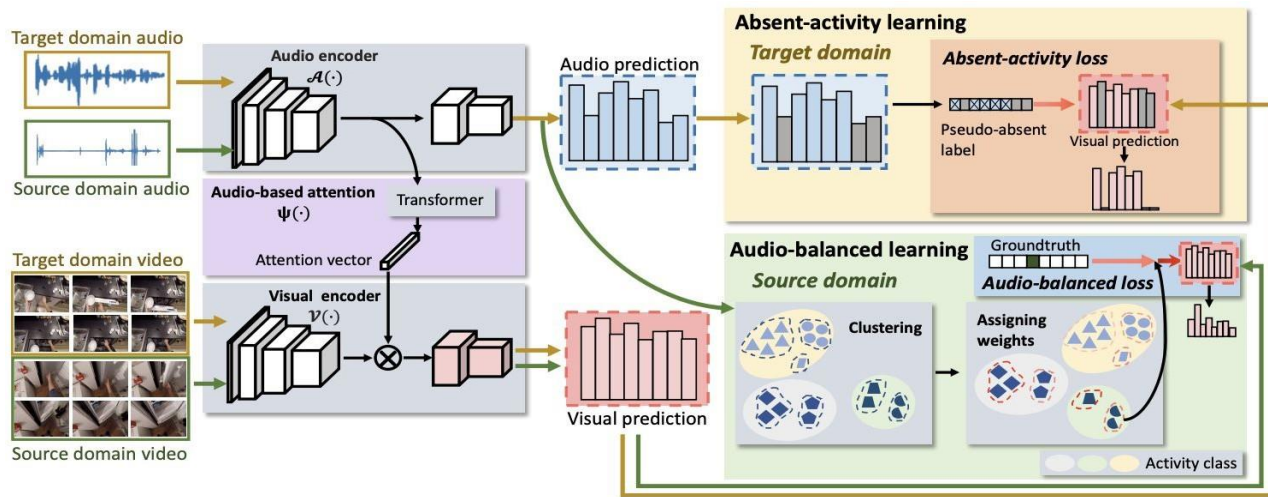
Aligned vs. misaligned



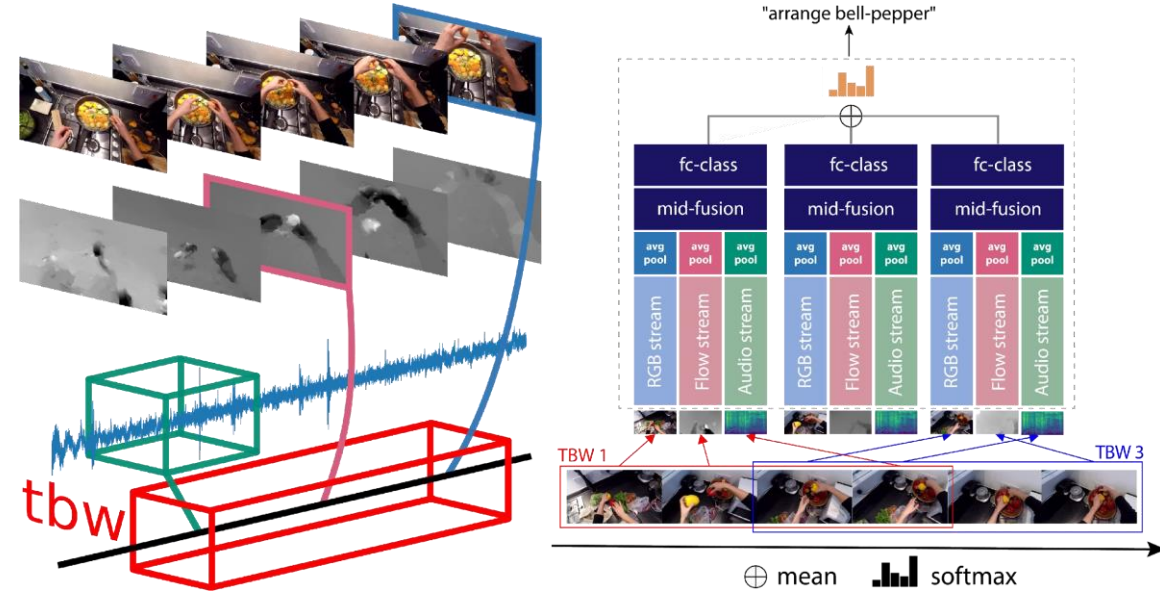
Multimodal Video Understanding



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



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



EPIC-Fusion: Audio-Visual Temporal Binding for Egocentric Action Recognition, Kazakos et al., ICCV 2019



ACKNOWLEDGEMENT

Thanks to the following courses and corresponding researchers for making their teaching/research material online

- CS231n: Deep Learning for Computer Vision, Stanford University
- Convolutional Neural Networks for Visual Recognition, Stanford University
- Deep Learning, Stanford University
- Introduction to Deep Learning, University of Illinois at Urbana-Champaign
- Introduction to Deep Learning, Carnegie Mellon University
- Natural Language Processing with Deep Learning, Stanford University
- And Many More Publicly Available Resources



Questions?

