



Chapter 2 - Section 7

Learning from Videos

Dr Yali Li

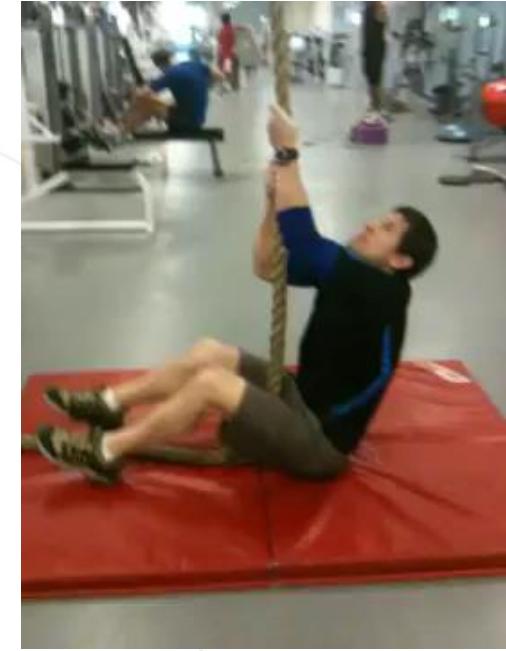
Friday, April 8, 2022

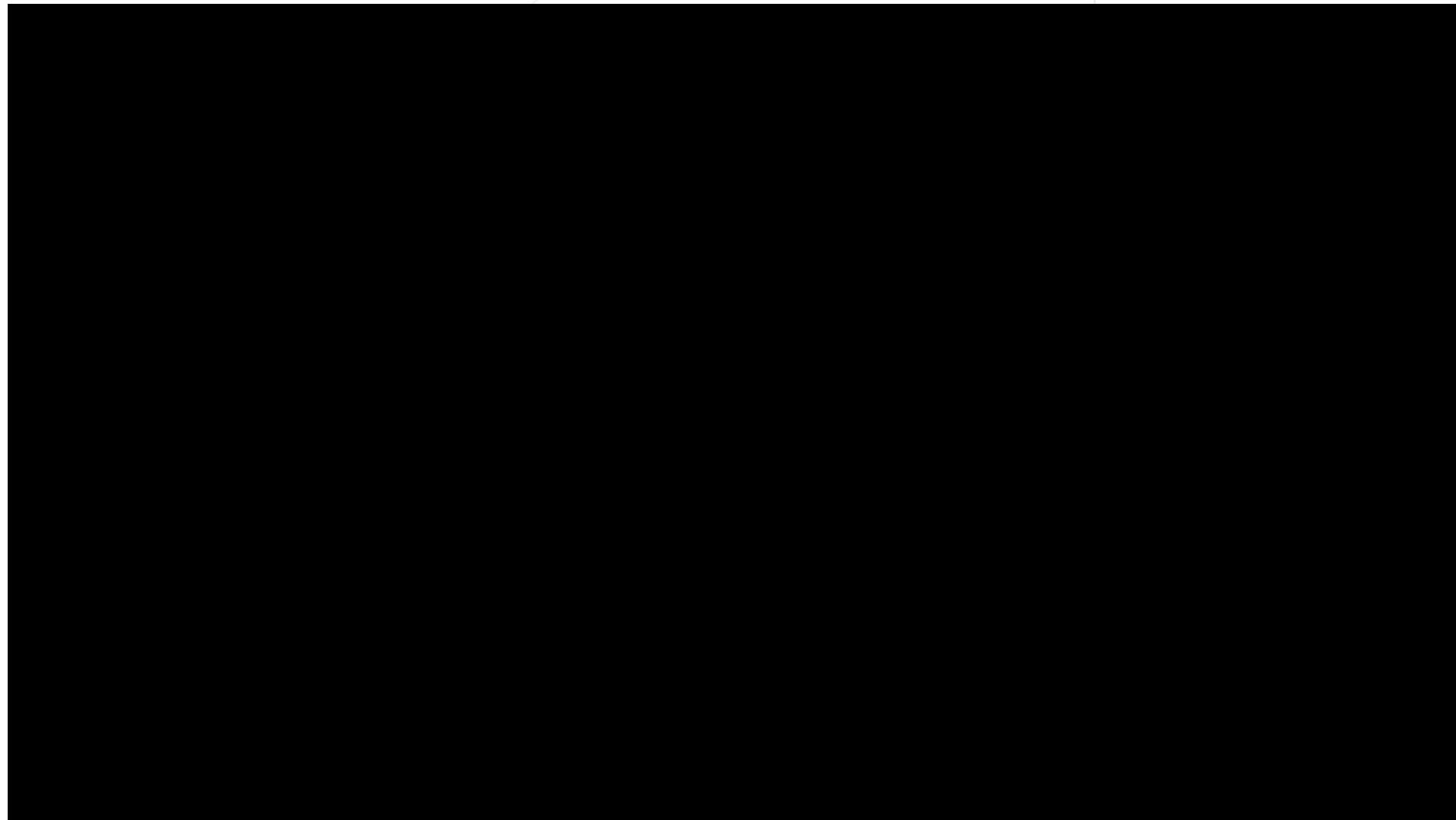
Acknowledge : Han Wang , Zhongdao Wang, Jian Han

Outline

Part 1	Video classification	P01-P08
Part 2	Datasets	P09-P18
Part 3	Solutions	P19-P98
Part 4	Summary & Future Topics	P99-P100

- We collect videos everywhere
 - Websites, mobiles, etc.







Sports Video Classification

https://cs.stanford.edu/people/karpathy/deepvideo/cnn_video_classify_demo.m4v

- Example tasks: video recognition vs image recognition

Images:



Recognize objects



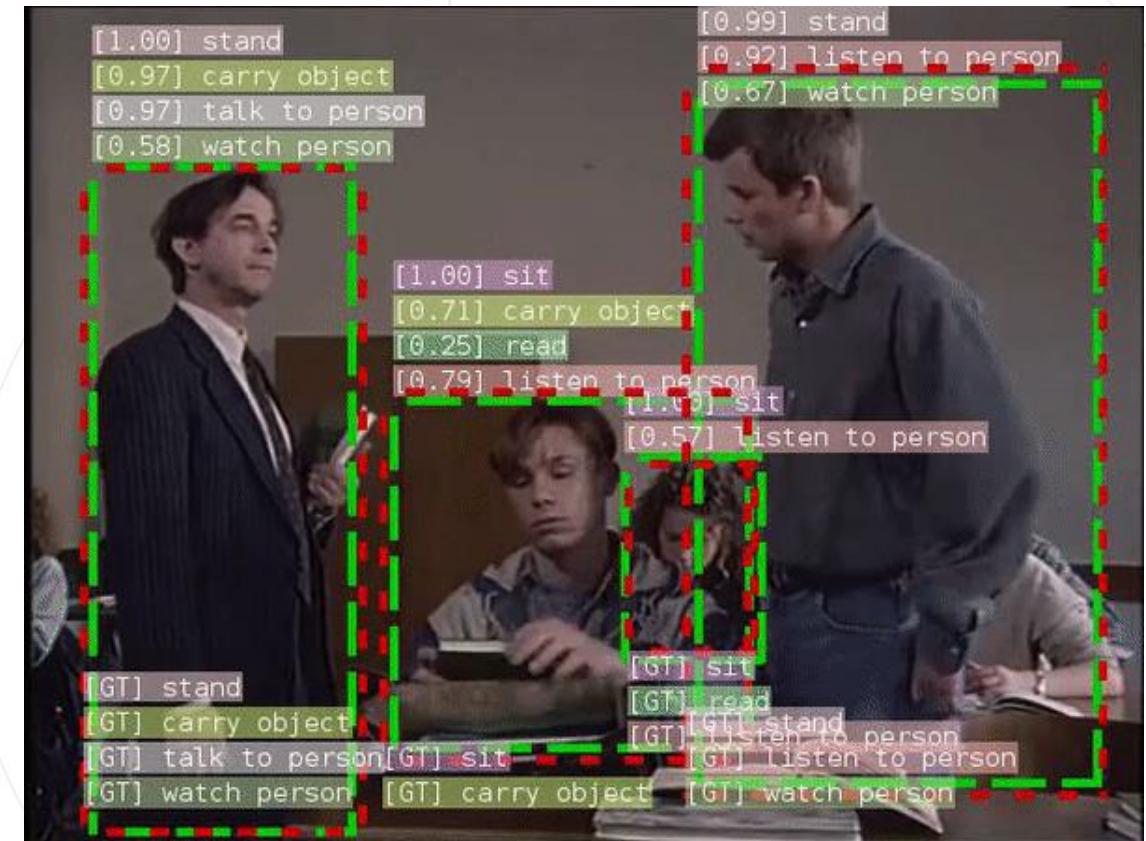
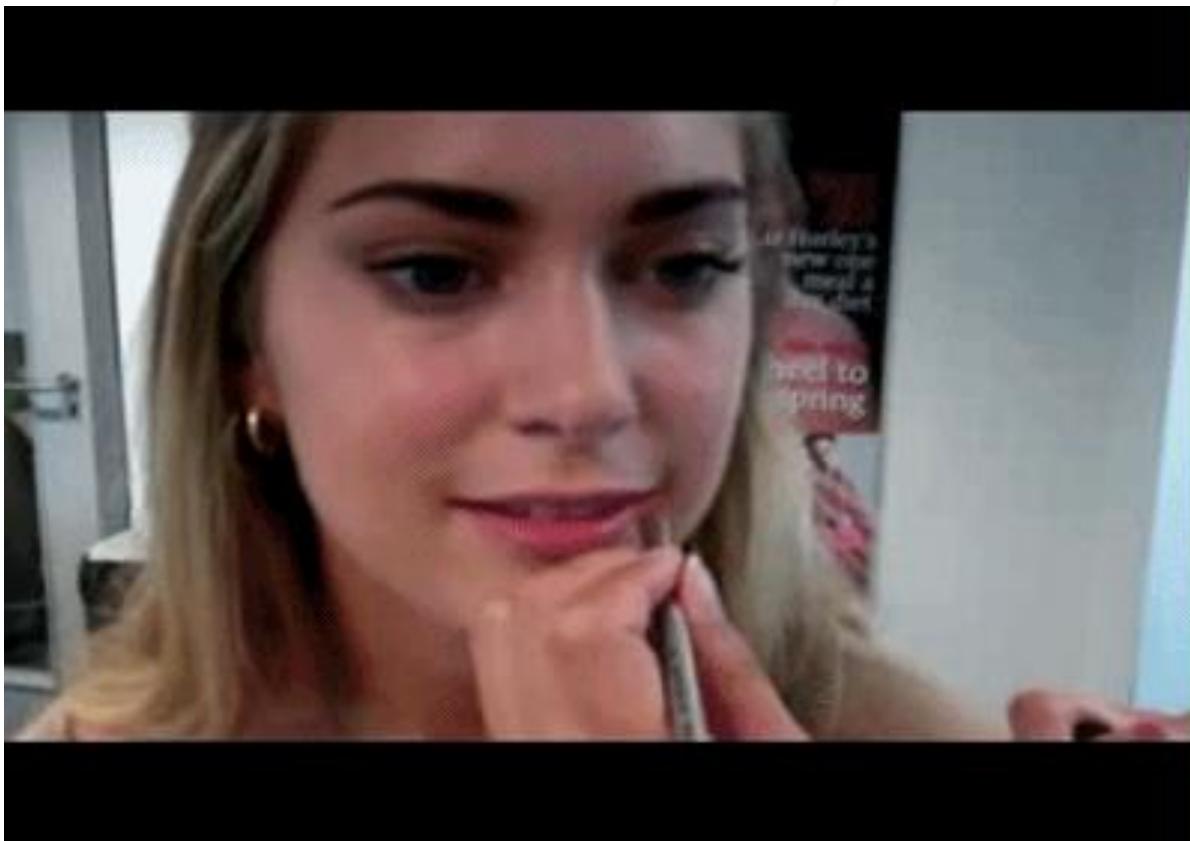
Videos:



Recognize actions



- Example tasks: video classification /action recognition



Outline

Part 1	Video classification	P01-P08
Part 2	Datasets	P09-P18
Part 3	Solutions	P19-P98
Part 4	Summary & Future Topics	P99-P100

Video Classification - Datasets

- UCF101
 - Youtube videos
 - 101 action classes, 13320 videos, ~27 hours of video data
 - Large variations in camera motion, object appearance and pose, viewpoint, background, illumination, etc.
 - HOG/HOF descriptors + SVM: yield an overall accuracy 43.9%.

K Soomro et al, UCF101: A Dataset of 101 Human Action Classes From Videos in The Wild, CRCV-TR-12-01, November, 2012.

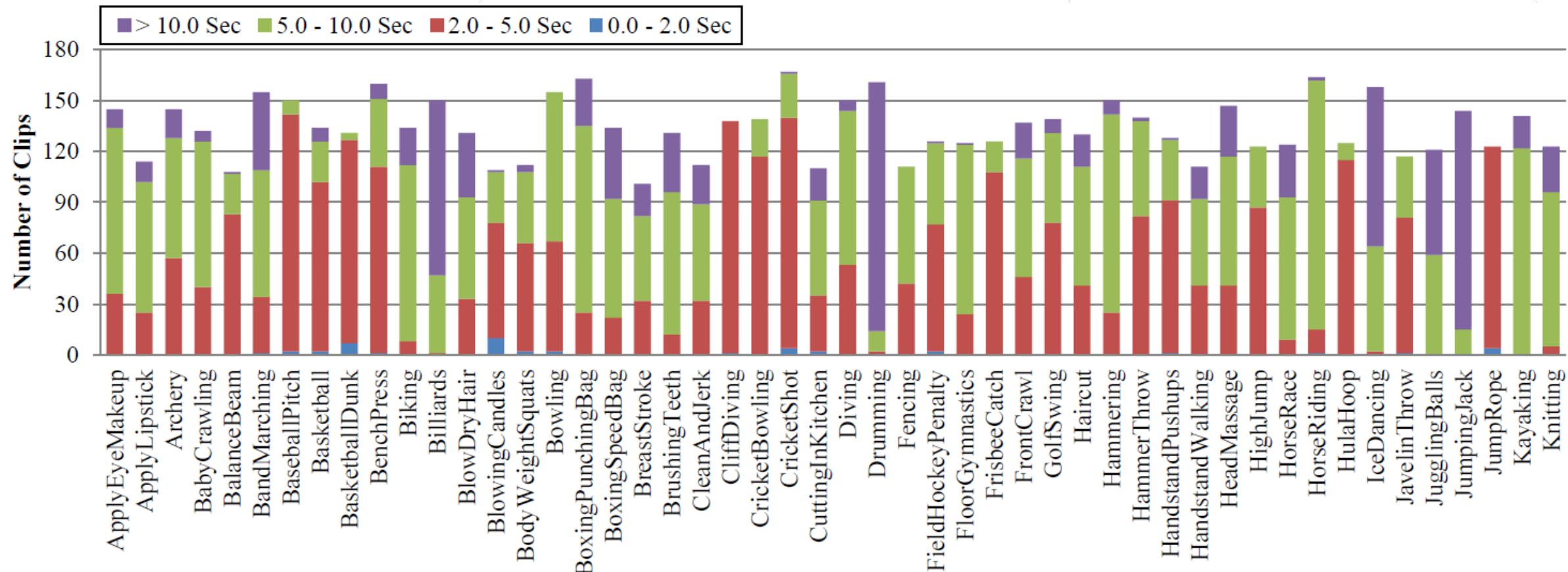
<https://www.crcv.ucf.edu/research/data-sets/ucf101/>



Video Classification - Datasets

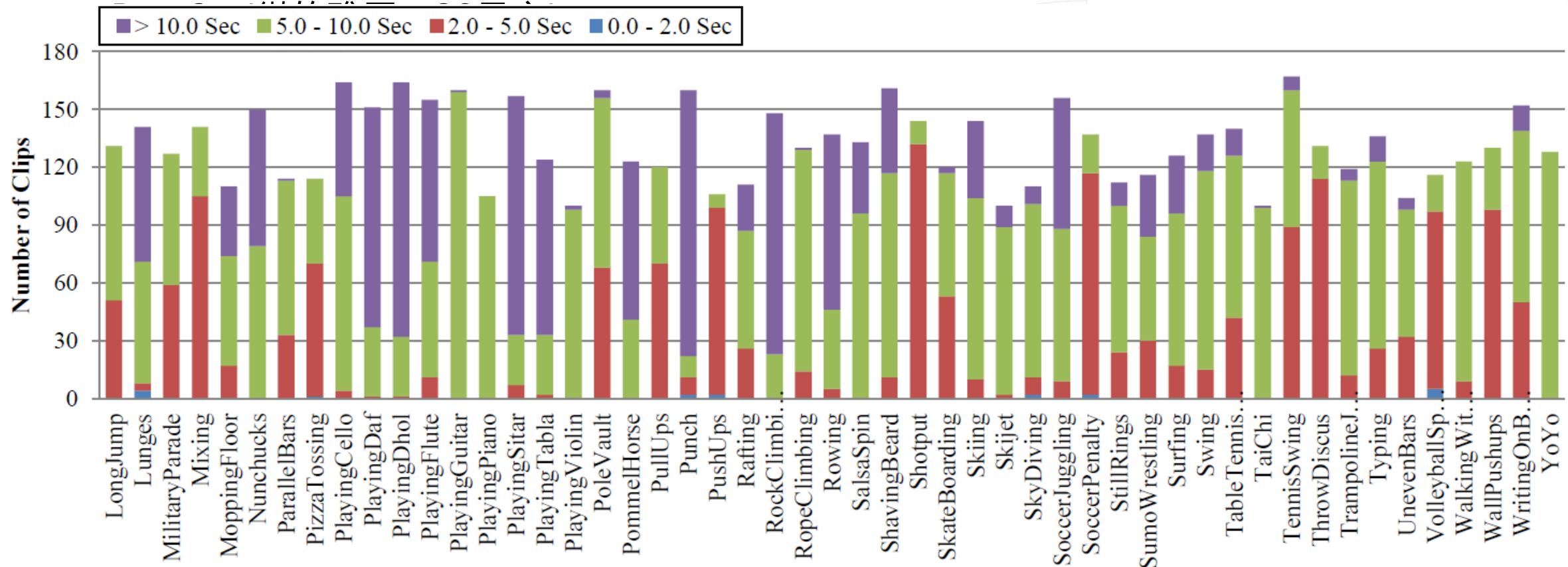
- UCF101

- 101 action classes, 13320 videos, ~27 hours of video data

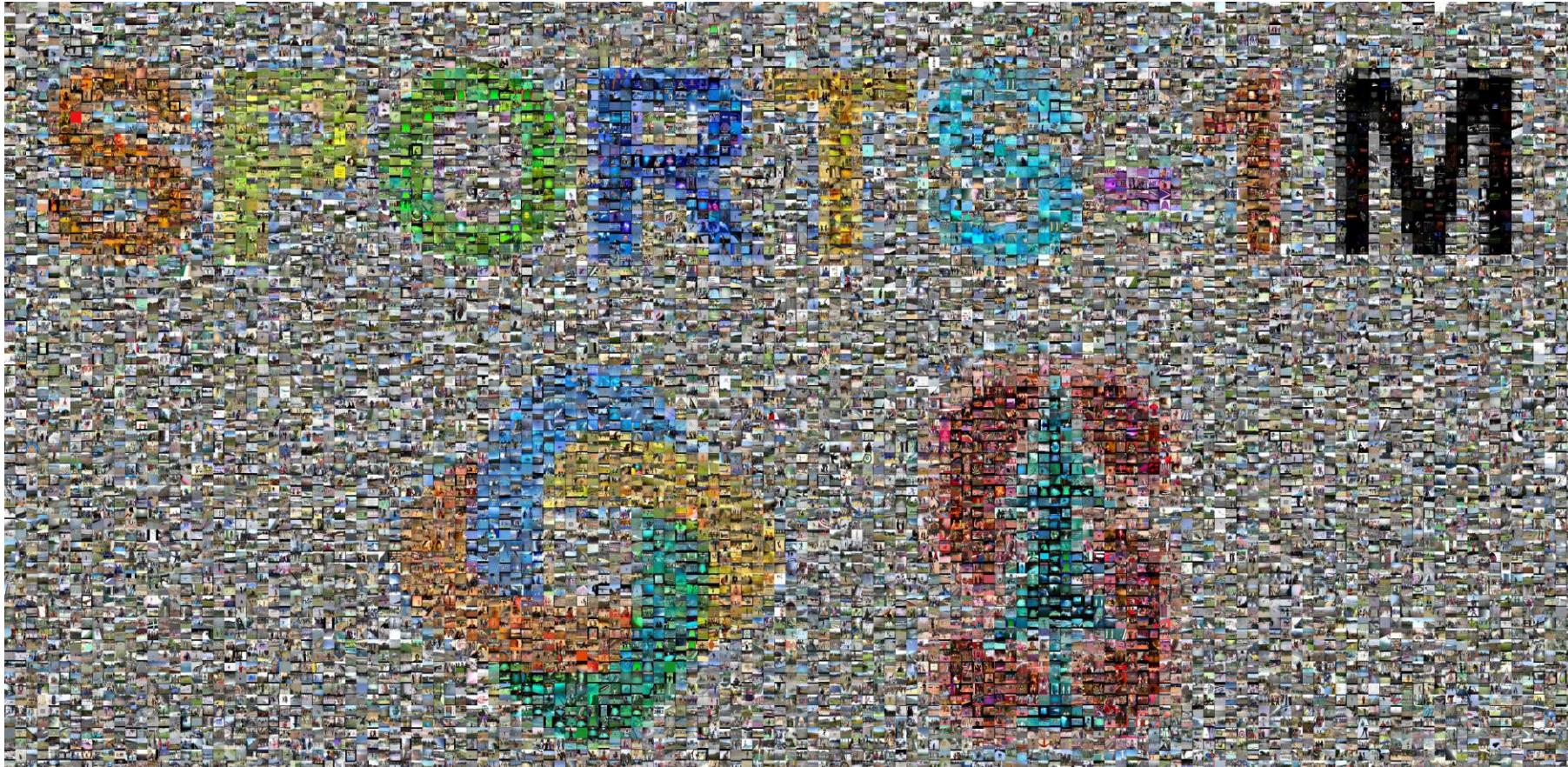


- UCF101

- 101 action classes, 13320 videos, ~27 hours of video data



- Sports-1M
 - YouTube videos, 1,133,157 videos, 487 sports labels

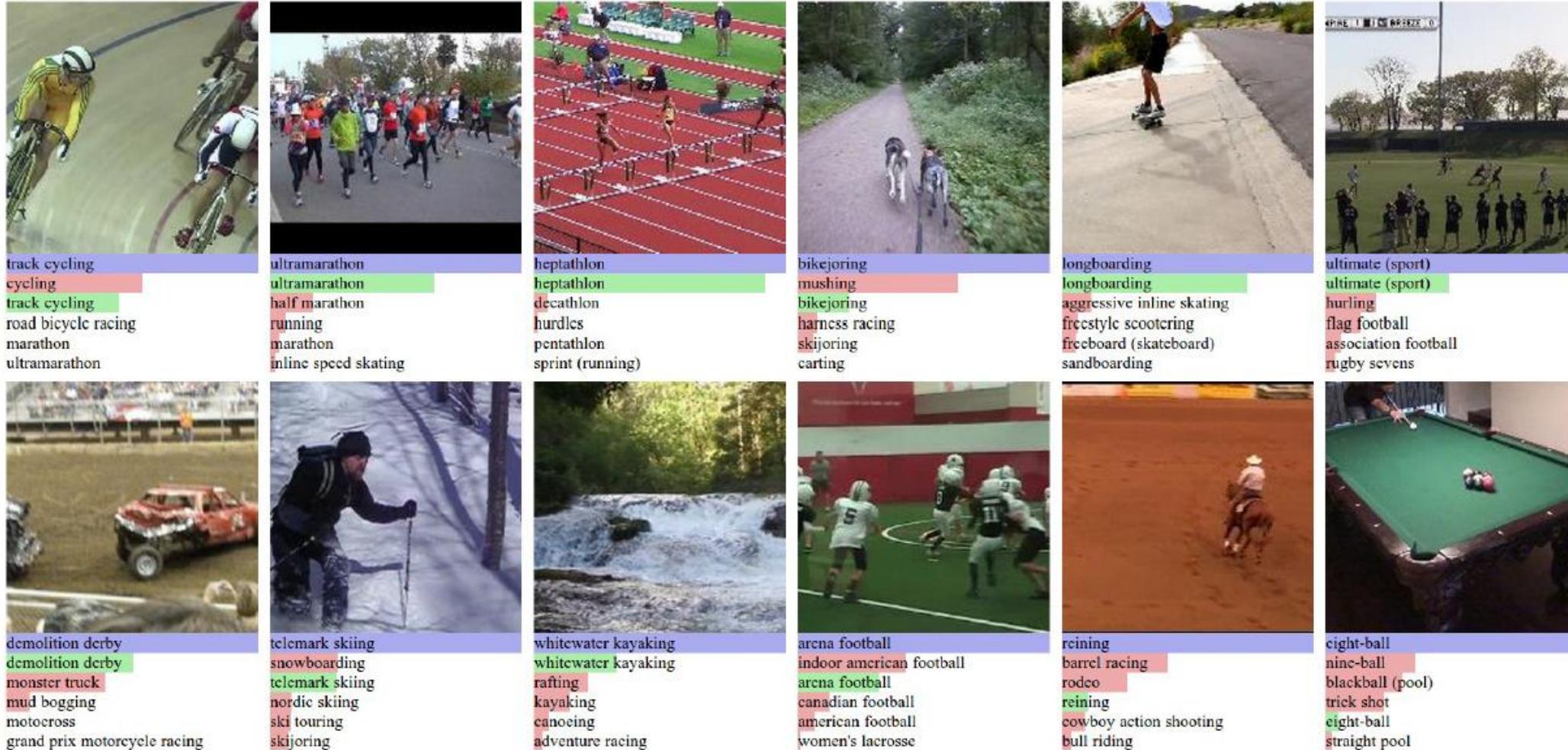


- Sports-1M
 - 1 million YouTube videos, belonging to 487 classes.
 - 1000-3000 videos per class
 - ~5% of the videos are annotated with more than one class.



Video Classification - Datasets

- Sports-1M



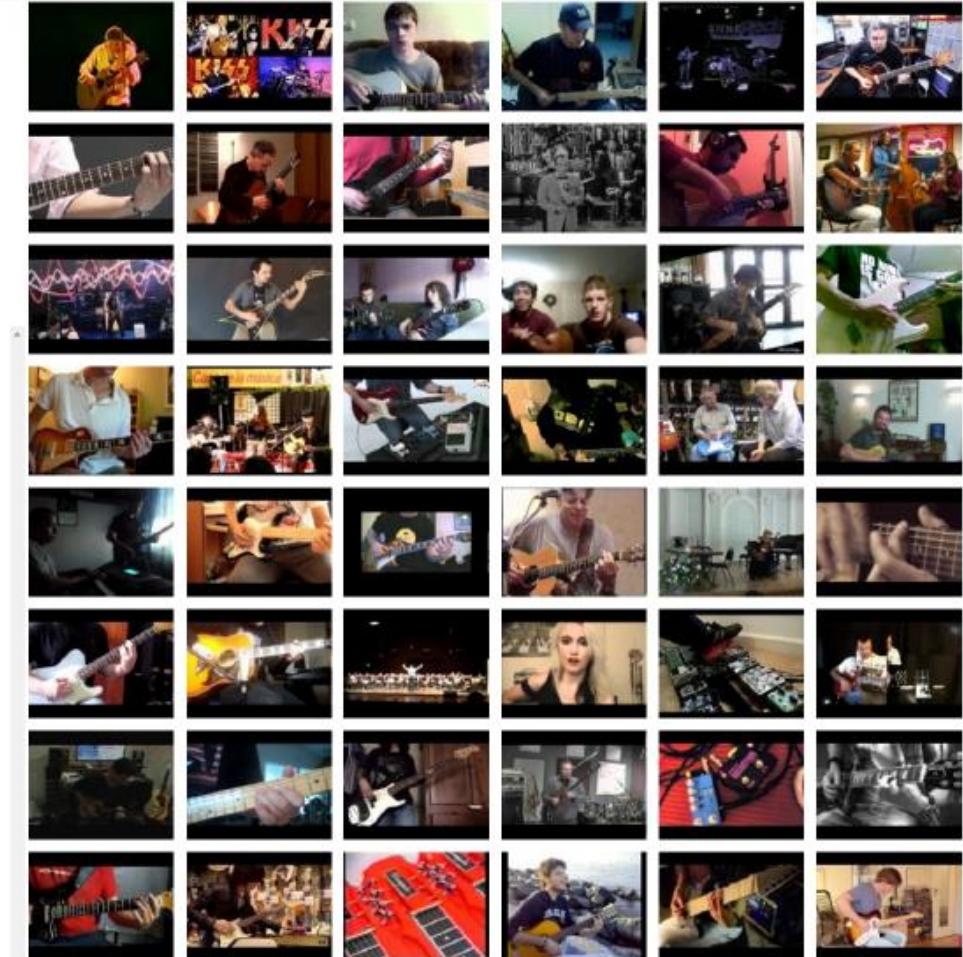
- YouTube-8M
 - millions of videos labeled with thousands of classes,
 - a large-scale benchmark dataset for general **multi-label** video classification.
 - 8 million videos
 - 500K hours of videos annotated with a vocabulary of 4800 visual entities.

Vertical
Arts & Entertainment ▾

Filter
Guitar

Entities

- Acoustic guitar
- Cort Guitars
- Electric guitar
- Flamenco guitar
- Guitar
- Guitar Center
- Guitar Hero
- Guitar Hero III: Legends of Rock
- Guitar amplifier
- Lead guitar
- PRS Guitars
- Pedal steel guitar
- Resonator guitar
- Steel guitar
- Steel-string acoustic guitar
- Twelve-string guitar
- Washburn Guitars



<http://research.google.com/youtube8m/>

- YouTube-8M



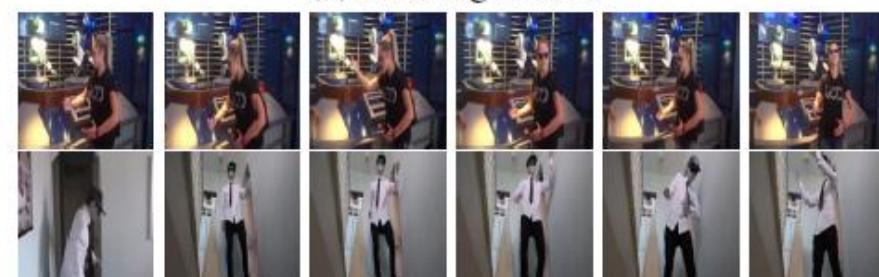
- Kinetics
 - Youtube videos
 - 650,000 video clips covering 400/600/700 human action classes
 - ≥ 600 video clips for each action class.
 - Each video clip lasts around 10 seconds and is labeled with a single action class



(a) headbanging



(c) shaking hands



(e) robot dancing

Outline

Part 1	Video classification	P01-P08
Part 2	Datasets	P09-P18
Part 3	Solutions	P19-P98
Part 4	Summary & Future Topics	P99-P100

Challenges:

- Computationally expensive
 - Lower quality:
Resolution, motion blur, occlusion
 - Requires lots of training data!
- Balance the efficiency and performance!

A video model need to:

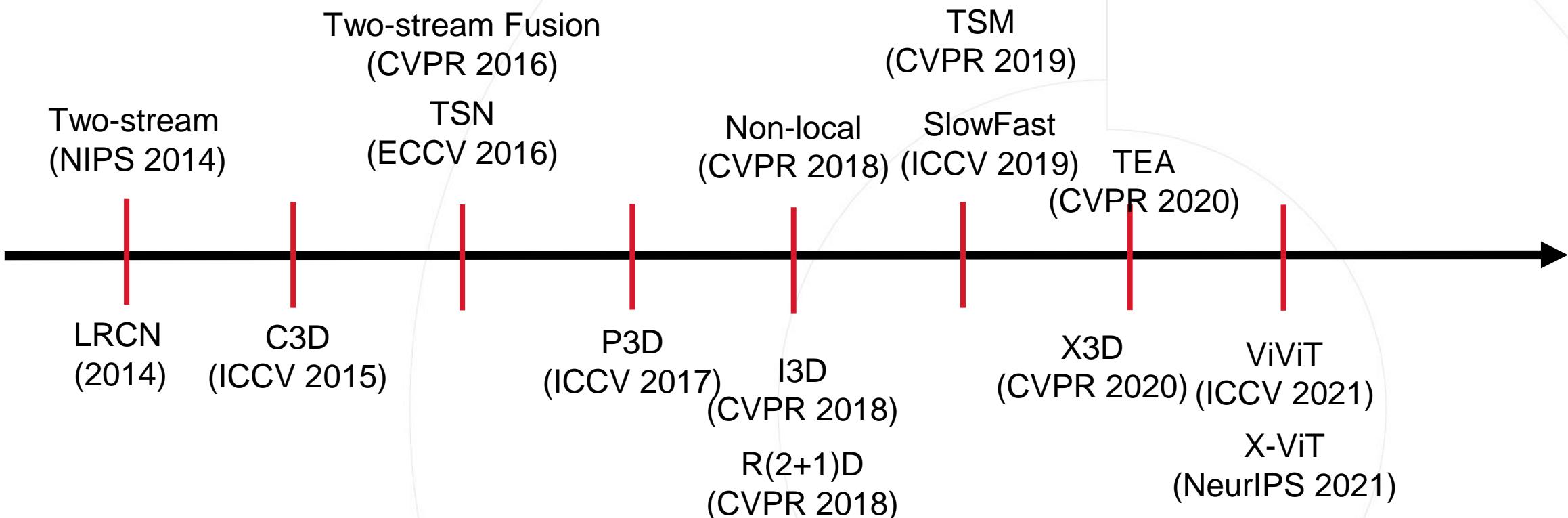
- Sequence modeling
- Temporal reasoning (receptive field)



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 (640x480): ~1.5GB per minute
HD (1920x1080): ~10GB per minute



Raw video: Long, with high FPS (frames per second)



Training: Train model to classify short clips with low FPS (frames per second)

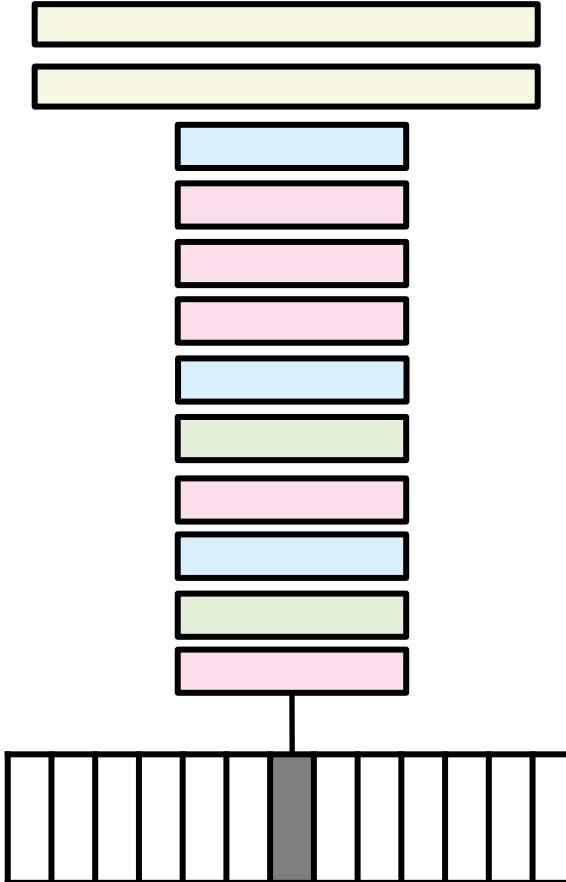


Testing: Run model on different clips, average predictions

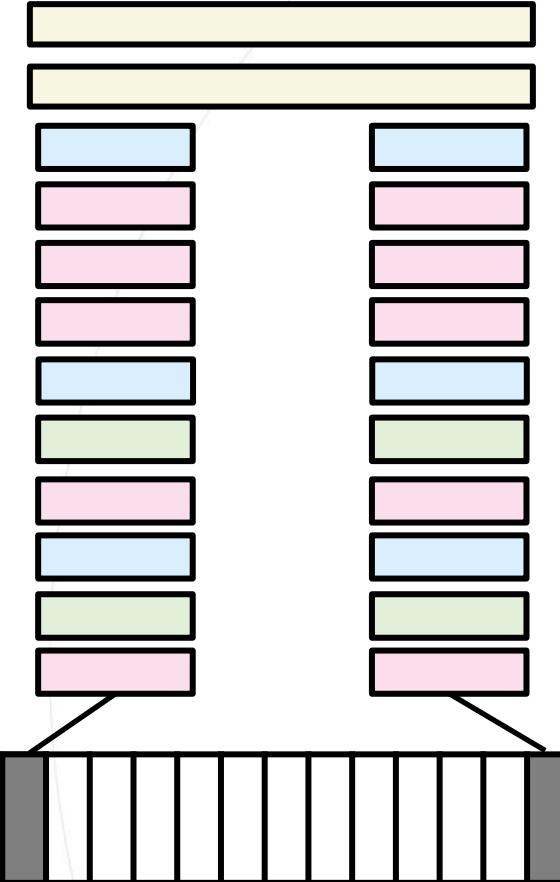


Video classification - pipelines

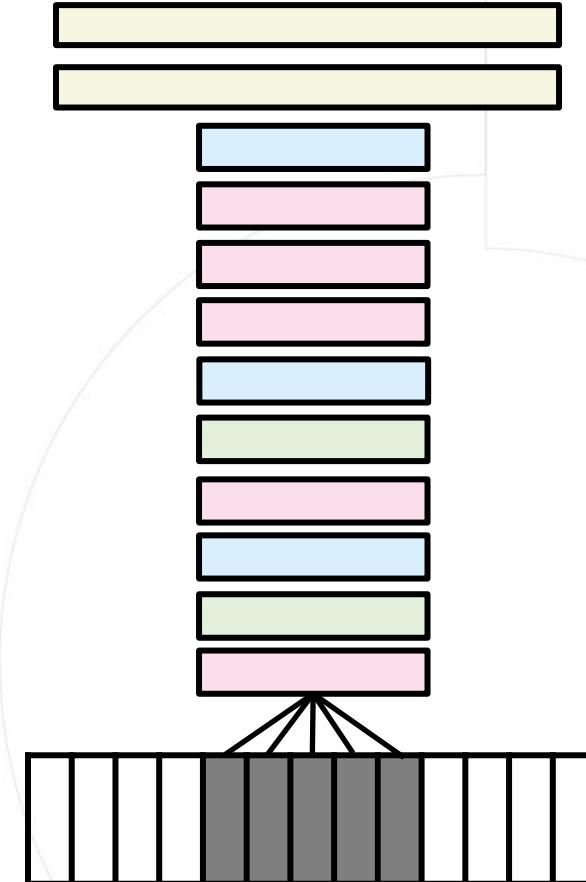
Single Frame



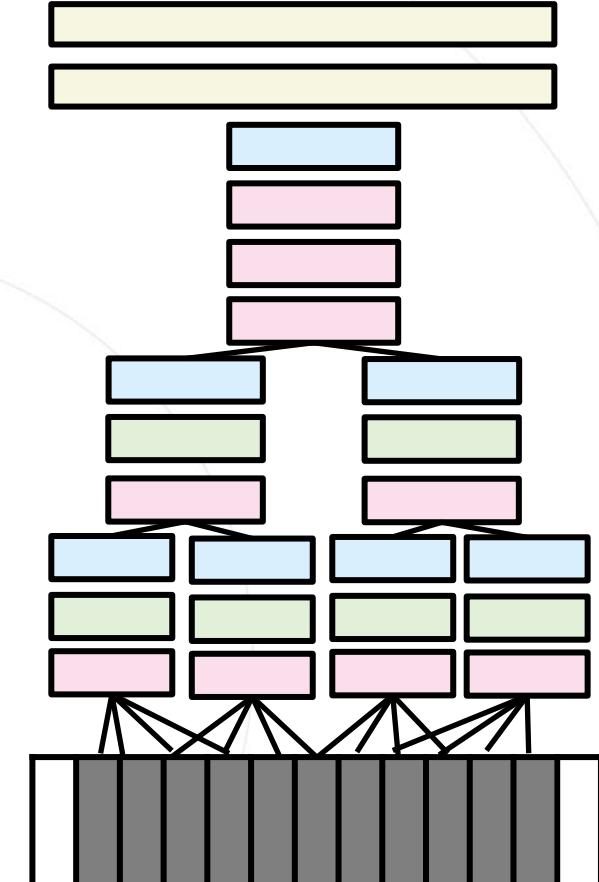
Late Fusion



Early Fusion



Slow Fusion



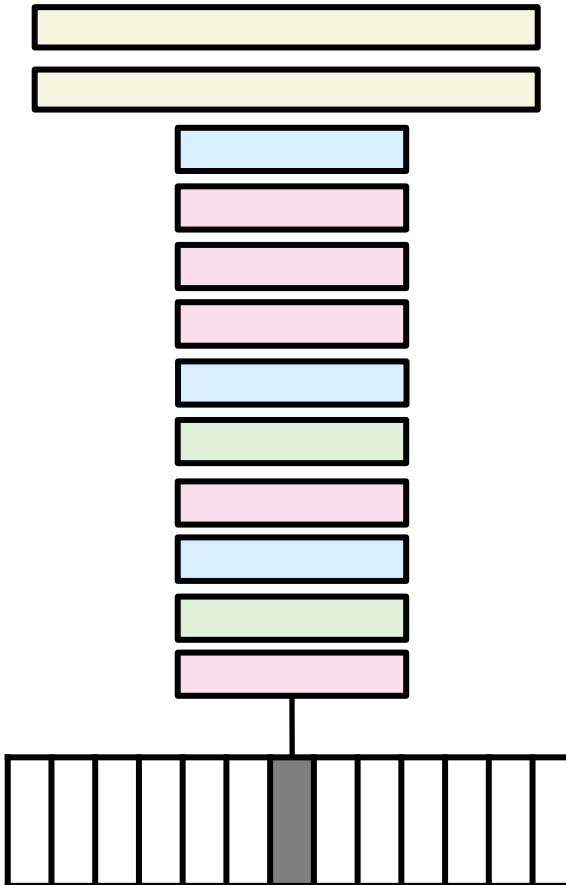
Convolutional layer

Normalization

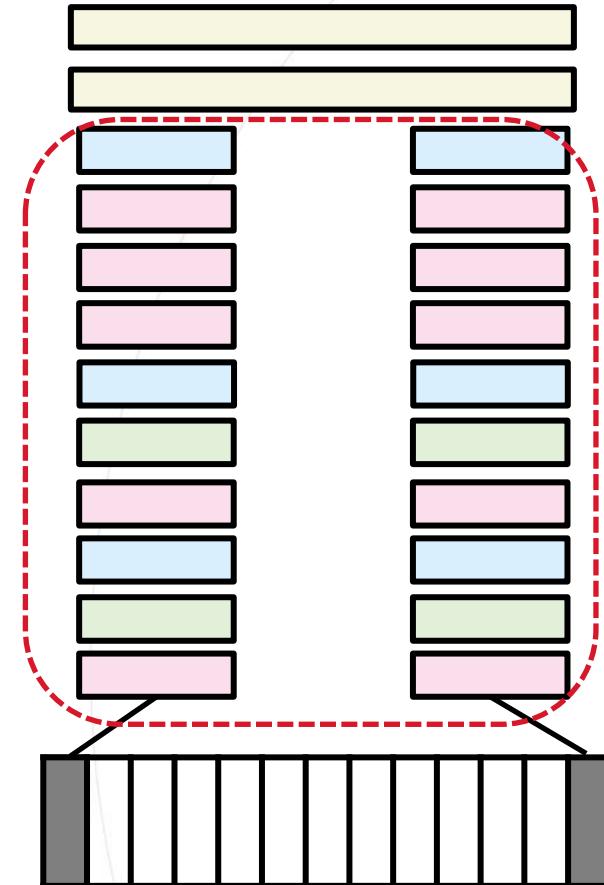
Pooling layer

Video classification - pipelines

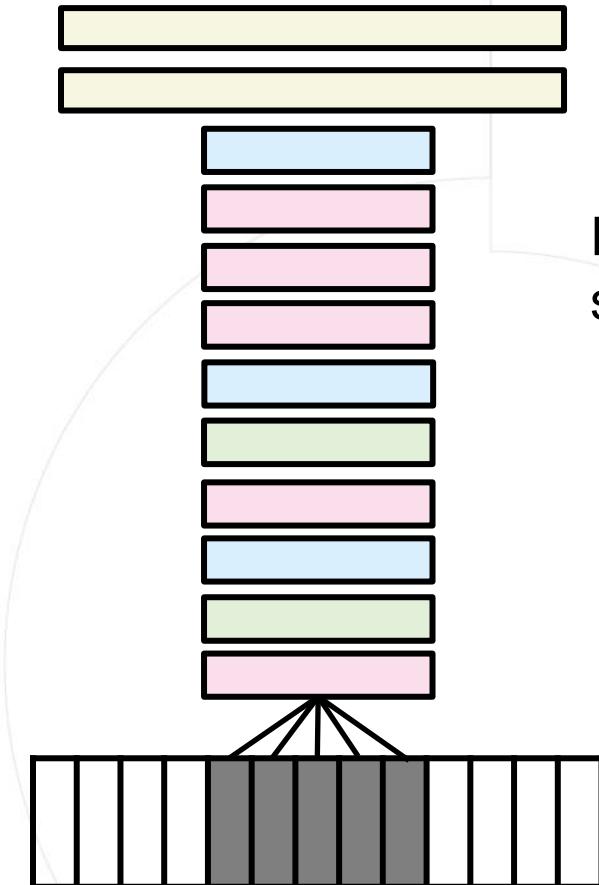
Single Frame



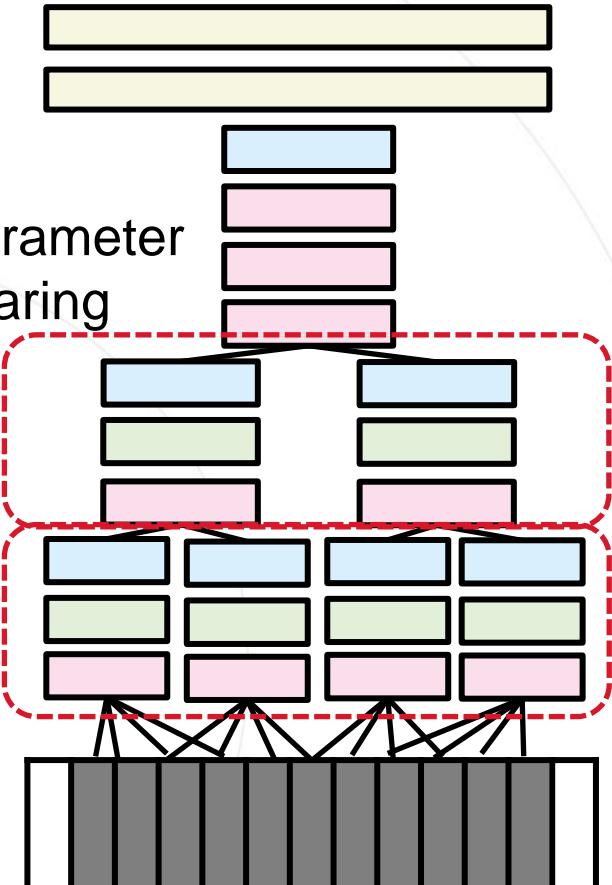
Late Fusion



Early Fusion



Slow Fusion



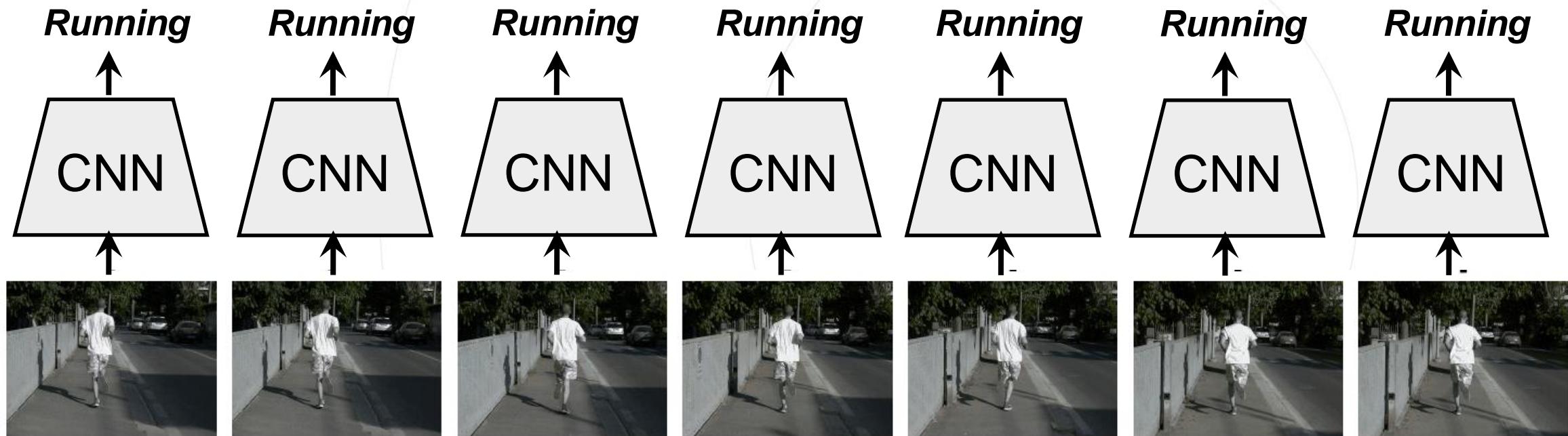
Convolutional layer

Normalization

Pooling layer

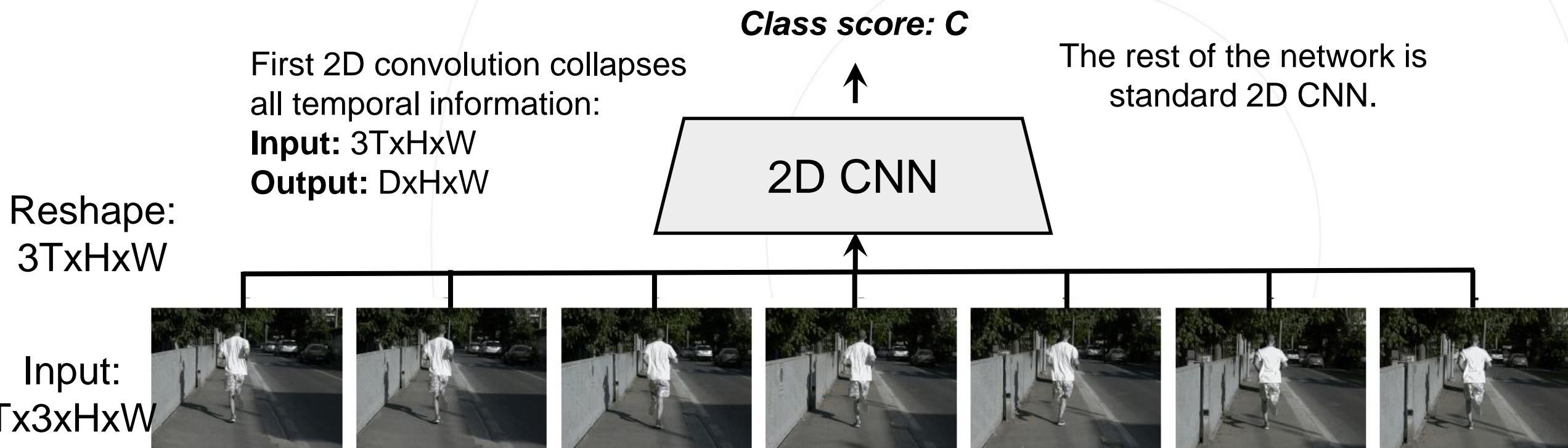
Video Classification: Single-Frame CNN

- Training: train normal 2D CNN to classify video frames independently
- Testing: average predicted probabilities
- Simple, yet strong baseline for video classification!



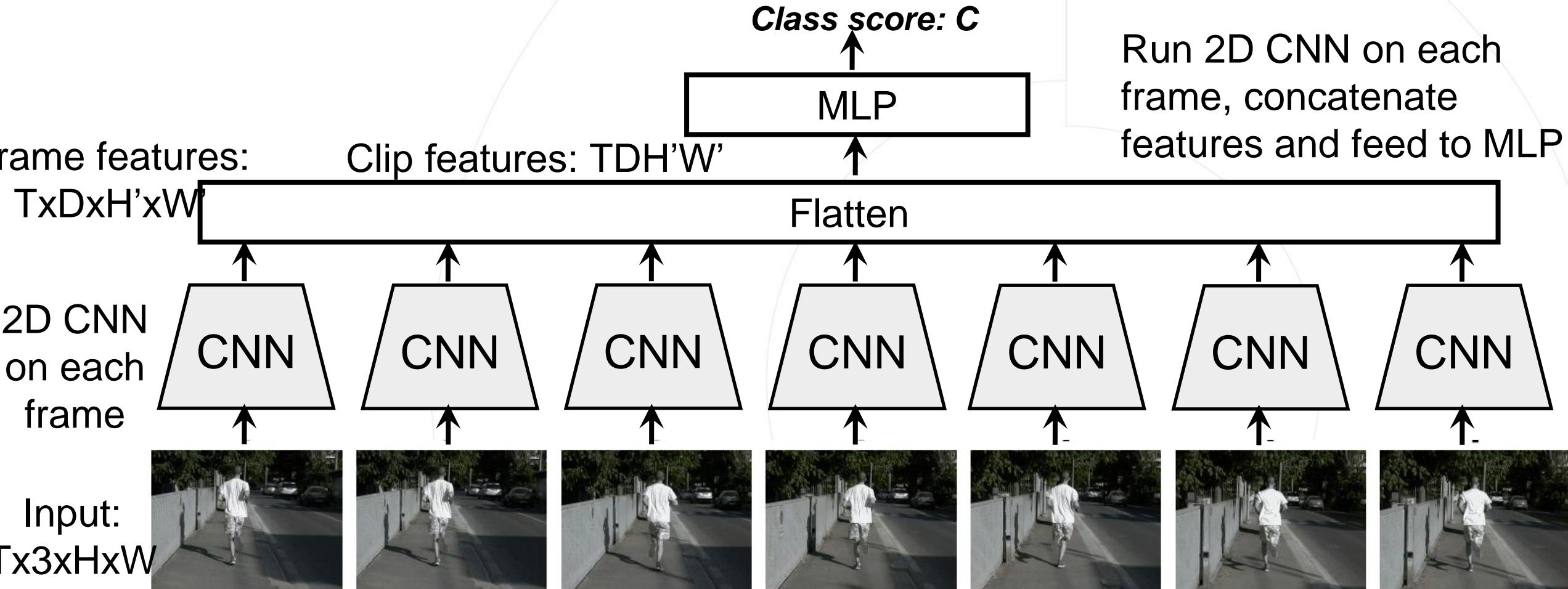
Video Classification: Early Fusion

- Compare frames with very first conv layer, after that normal 2D CNN
- Combine information across an entire time window immediately on the pixel level. The early and direct connectivity to pixel data allows the network to precisely detect local motion direction and speed.



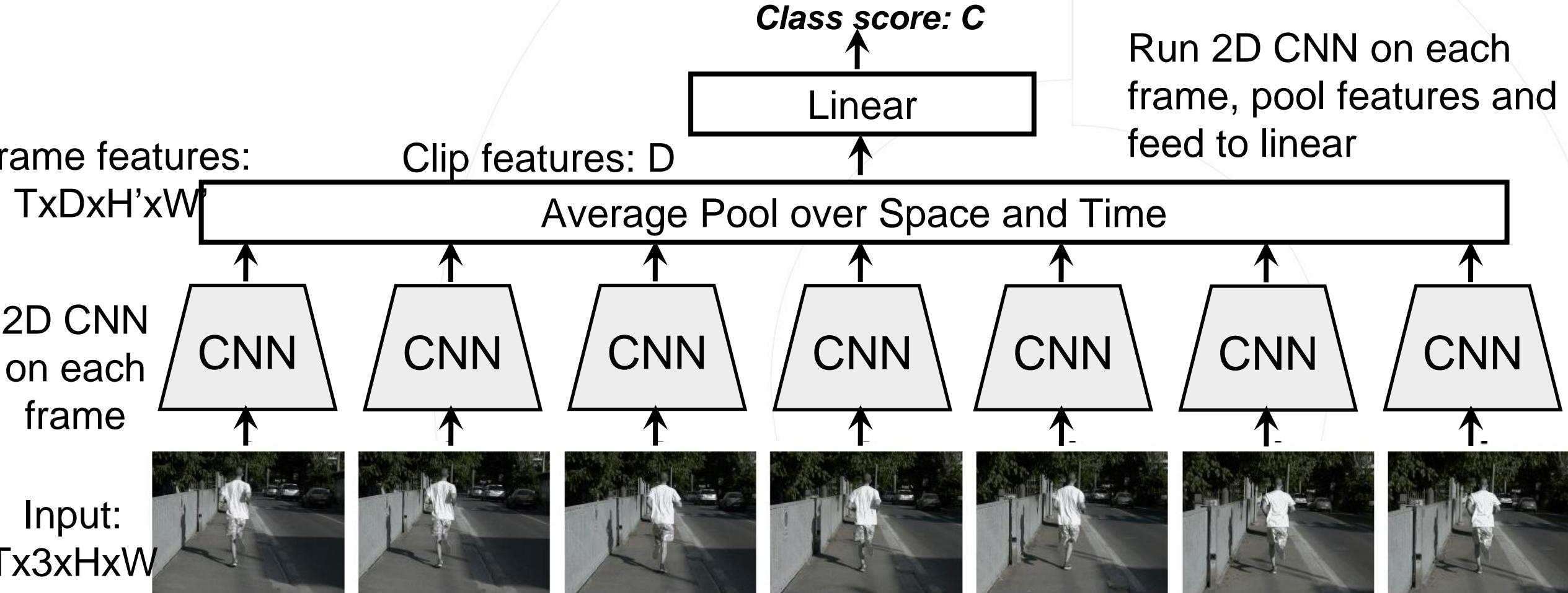
Video Classification: Late Fusion (with FC layers)

- Get high-level appearance of each frame, combine them together

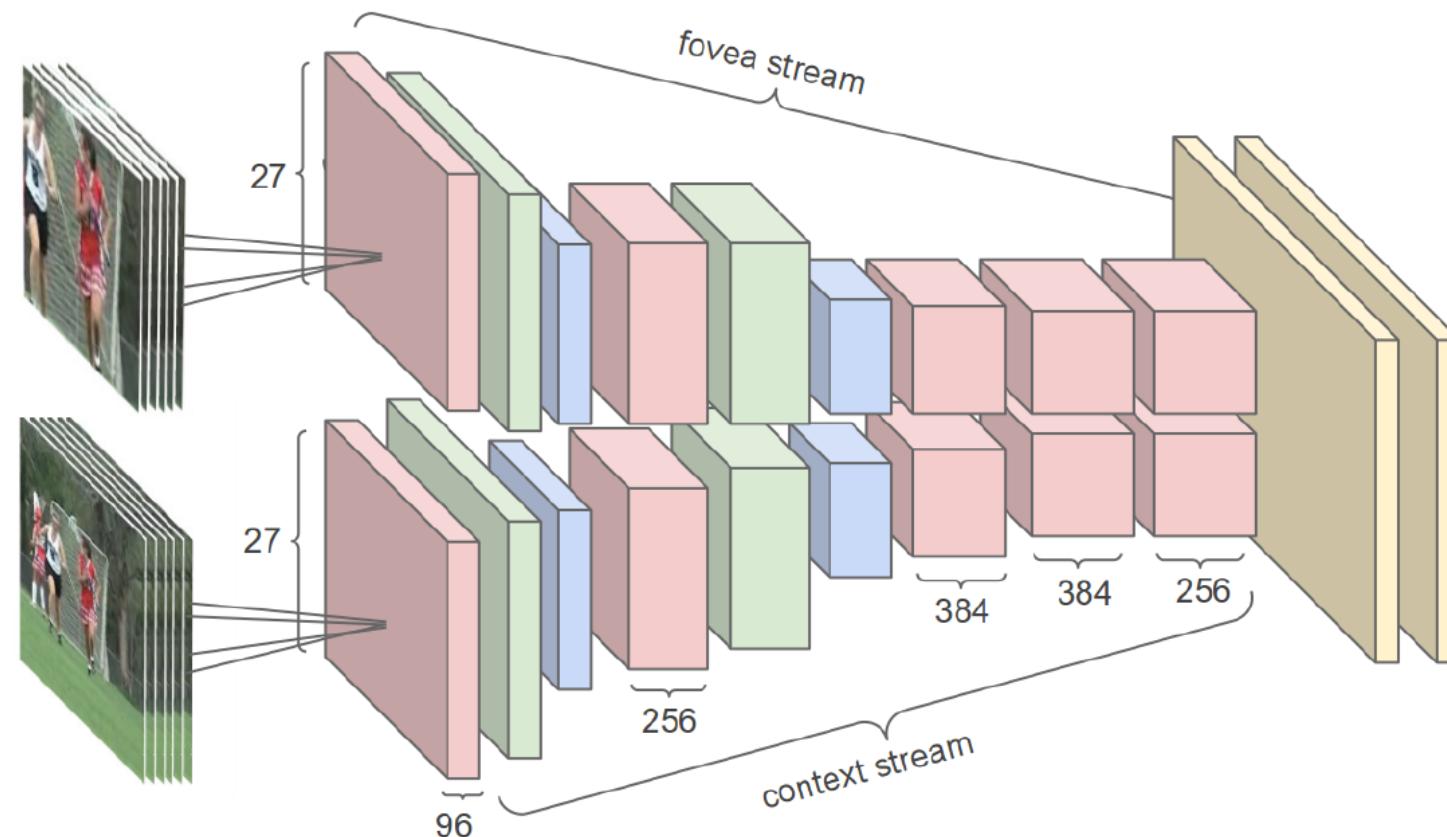


Video Classification: Late Fusion (with pooling layers)

- Get high-level appearance of each frame, combine them together



Fovea stream: receives the center 89×89 region at the original resolution.
Context stream: receives the down-sampled frames at half the original spatial resolution (89×89 pixels)



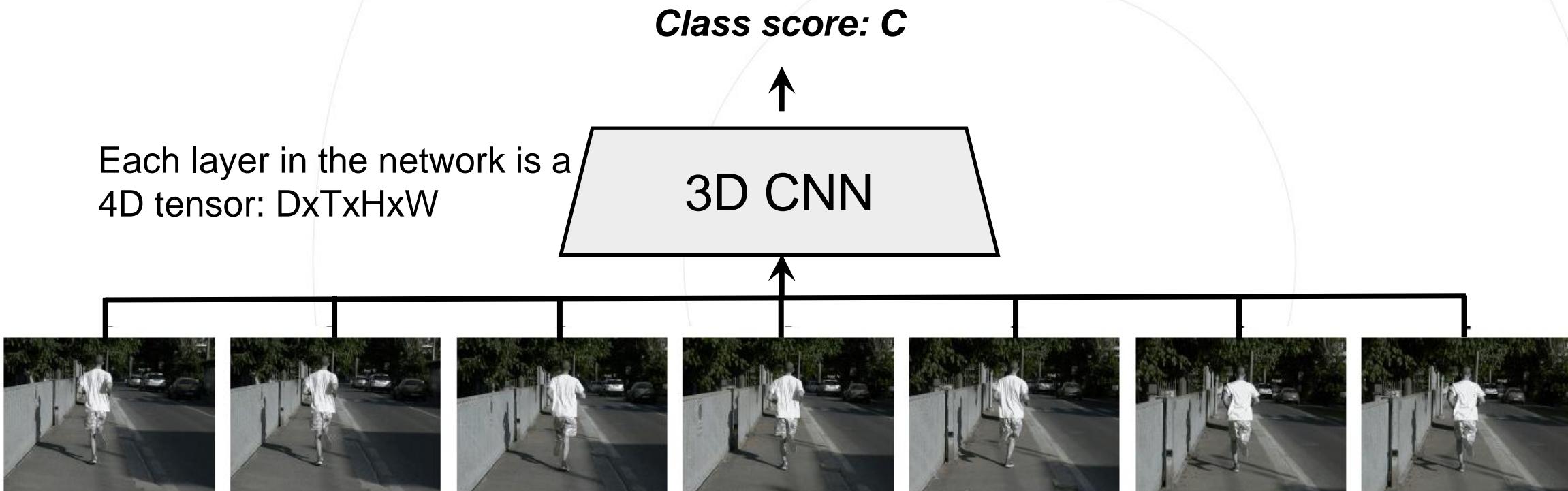
Video classification: Performance comparison

Model	Clip Hit@1	Video Hit@1	Video Hit@5
Feature Histograms + Neural Net	-	55.3	-
Single-Frame	41.1	59.3	77.7
Single-Frame + Multires	42.4	60.0	78.5
Single-Frame Fovea Only	30.0	49.9	72.8
Single-Frame Context Only	38.1	56.0	77.2
Early Fusion	38.9	57.7	76.8
Late Fusion	40.7	59.3	78.7
<u>Slow Fusion</u>	41.9	60.9	80.2
CNN Average (Single+Early+Late+Slow)	41.4	63.9	82.4

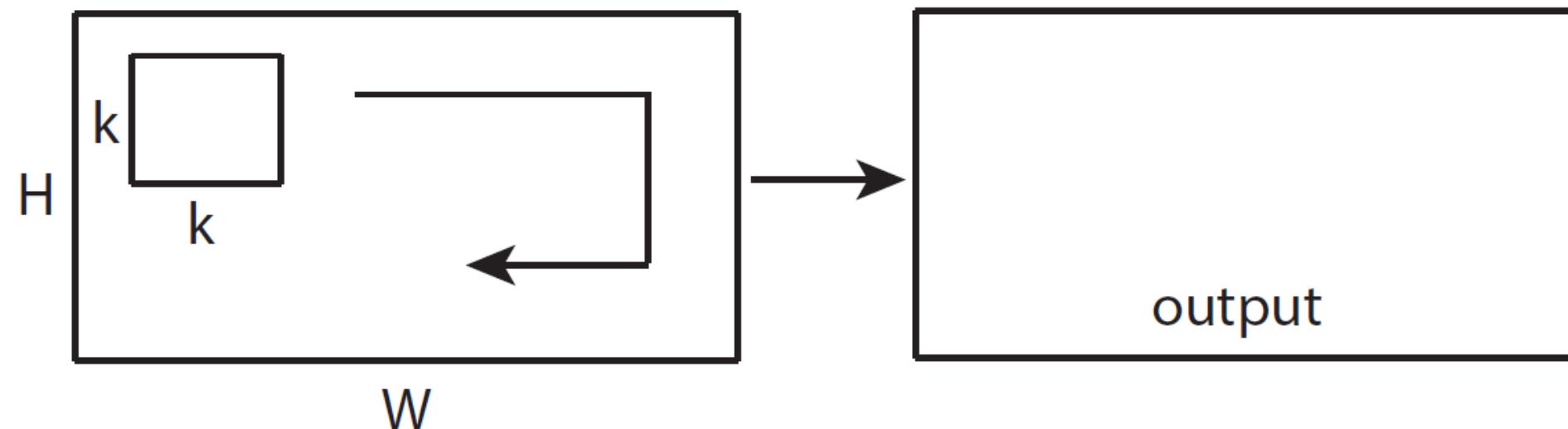
Table 1: Results on the 200,000 videos of the Sports-1M test set. Hit@k values indicate the fraction of test samples that contained at least one of the ground truth labels in the top k predictions.

Video Classification: 3D CNN

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



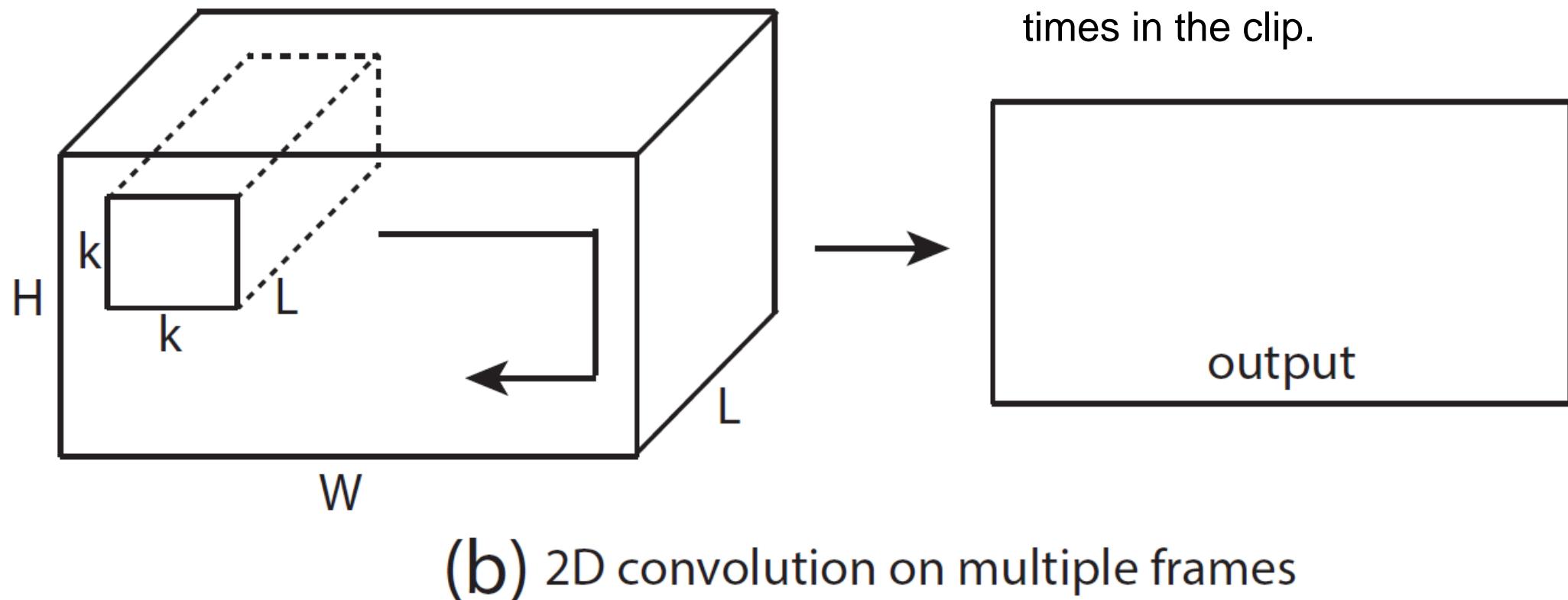
Video Classification: 3D CNN – C3D



(a) 2D convolution

Sliding over x, y

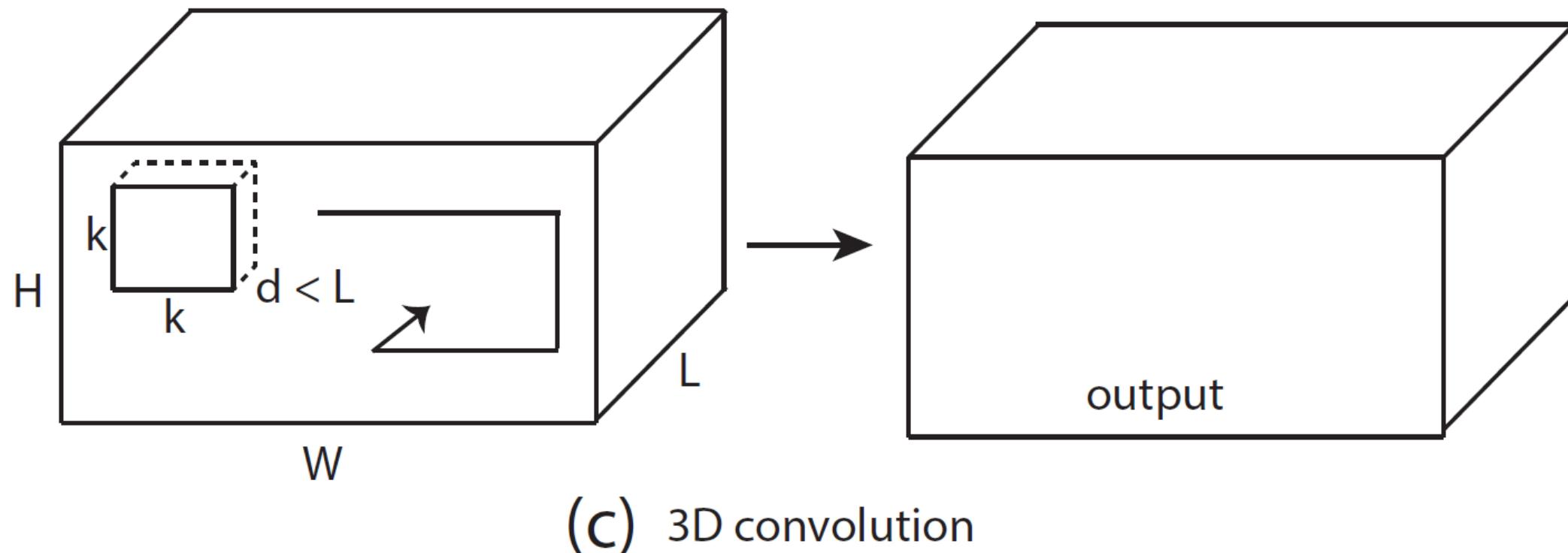
Video Classification: 3D CNN



No temporal shift-invariance! The convolutional filters are different for the same motion at different times in the clip.

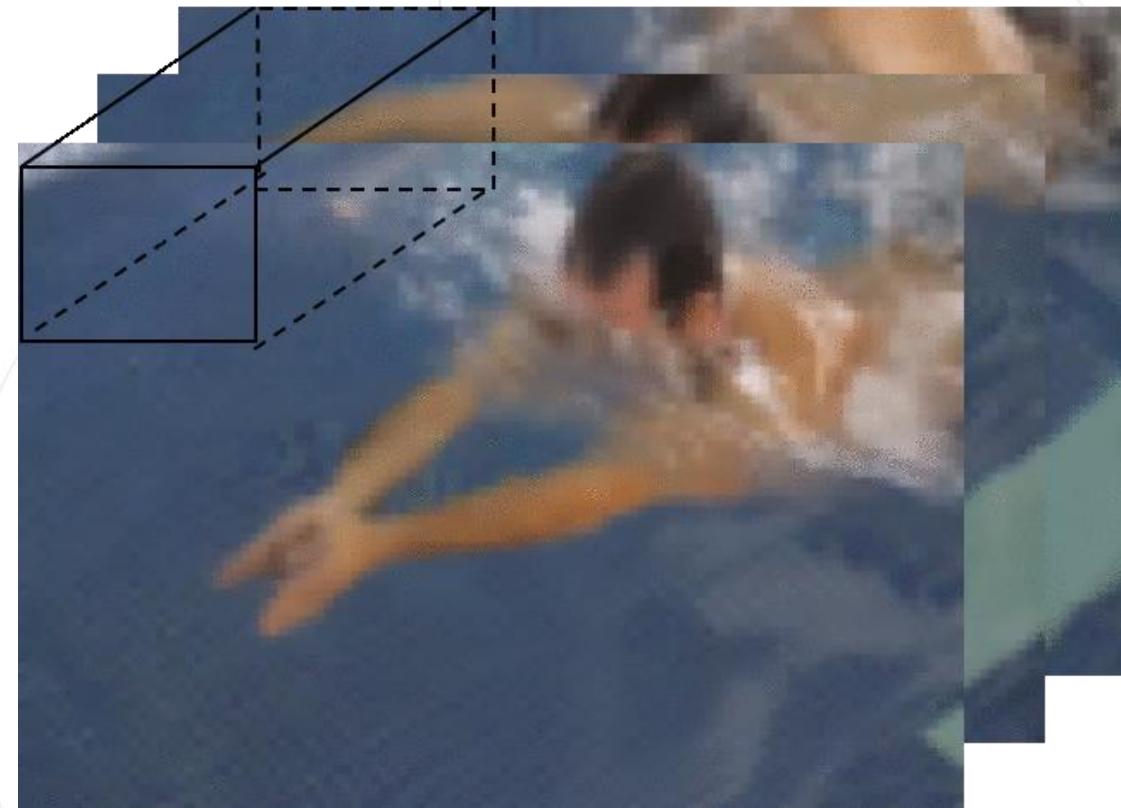
Video Classification: 3D CNN

Temporal shift-invariance.
Each convolutional filter slides
over time, x, y .



Video Classification: 3D CNN – C3D

- VGG of 3D CNNs
- Uses all $3 \times 3 \times 3$ conv and $2 \times 2 \times 2$ pooling (except Pool1)
- Release pre-trained model on Sports-1M: can be used as a video feature extractor



Du Tran, et al. Learning spatiotemporal features with 3D convolutional networks. In Proc. ICCV, 2015.

Video Classification: 3D CNN – C3D

- VGG of 3D CNNs
- Uses all $3 \times 3 \times 3$ conv and $2 \times 2 \times 2$ pooling (except Pool1)
- Release pre-trained model on Sports-1M: can be used as a video feature extractor
- Problem: 3D conv is very expensive!
- AlexNet: 0.7 GFLOP
- VGG-16: 13.6 GFLOP
- **C3D: 39.5 GFLOP**

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

Video Classification: 3D CNN – C3D



Figure 4. **Visualization of C3D model, using the method from [46].** Interestingly, C3D captures appearance for the first few frames but thereafter only attends to salient motion. Best viewed on a color screen.

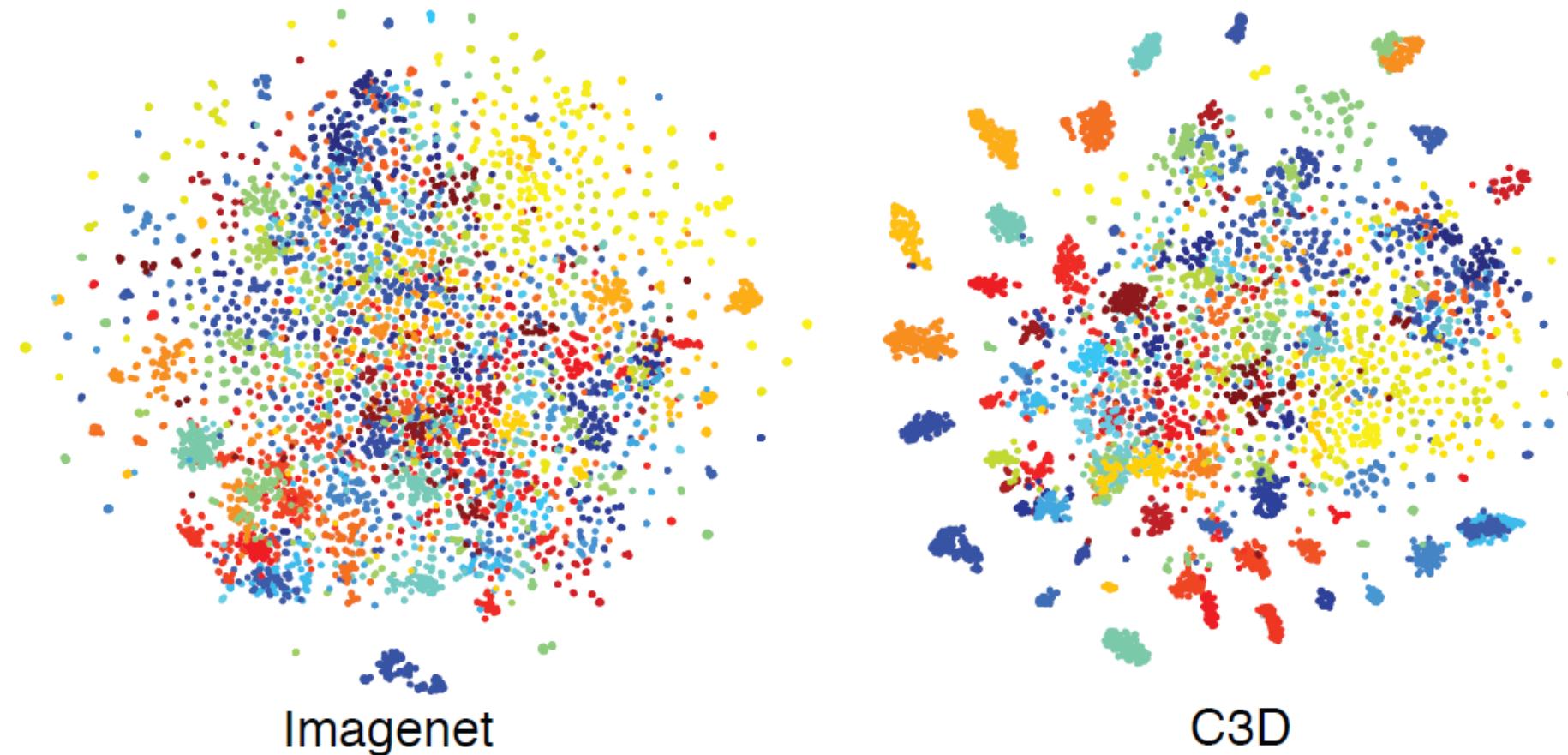
Video Classification: 3D CNN – C3D



<https://vlg.cs.dartmouth.edu/c3d/>

Du Tran, et al. Learning spatiotemporal features with 3D convolutional networks. In Proc. ICCV, 2015.

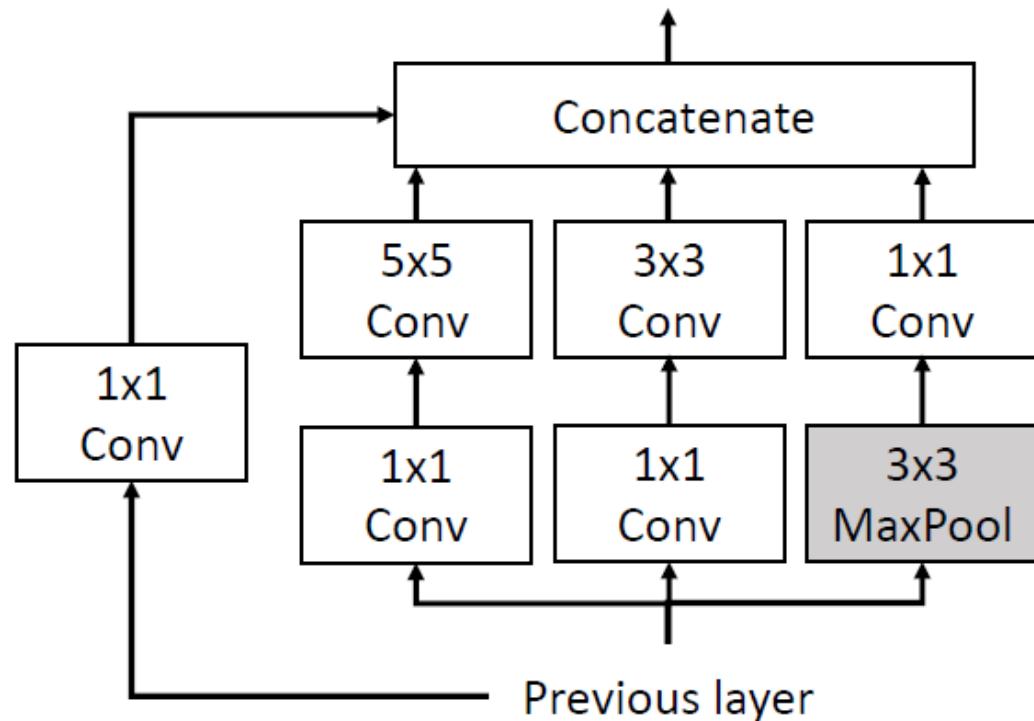
Video Classification: 3D CNN – C3D



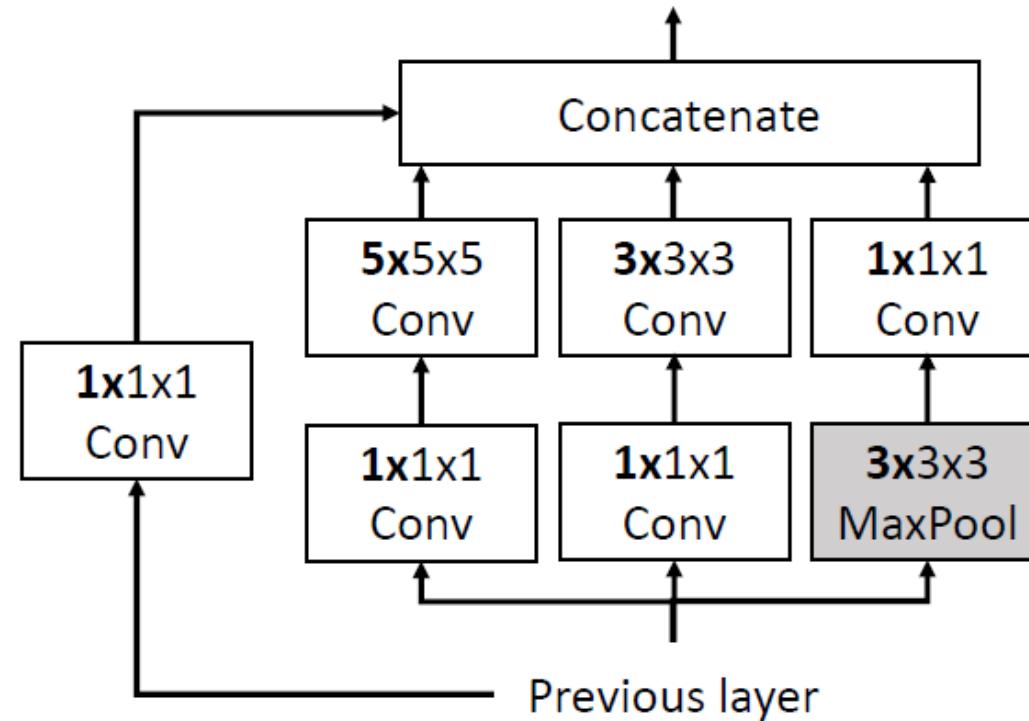
Du Tran, et al. Learning spatiotemporal features with 3D convolutional networks. In Proc. ICCV, 2015.

Video Classification: 3D CNN – I3D

Inception Block: Original



Inception Block: Inflated

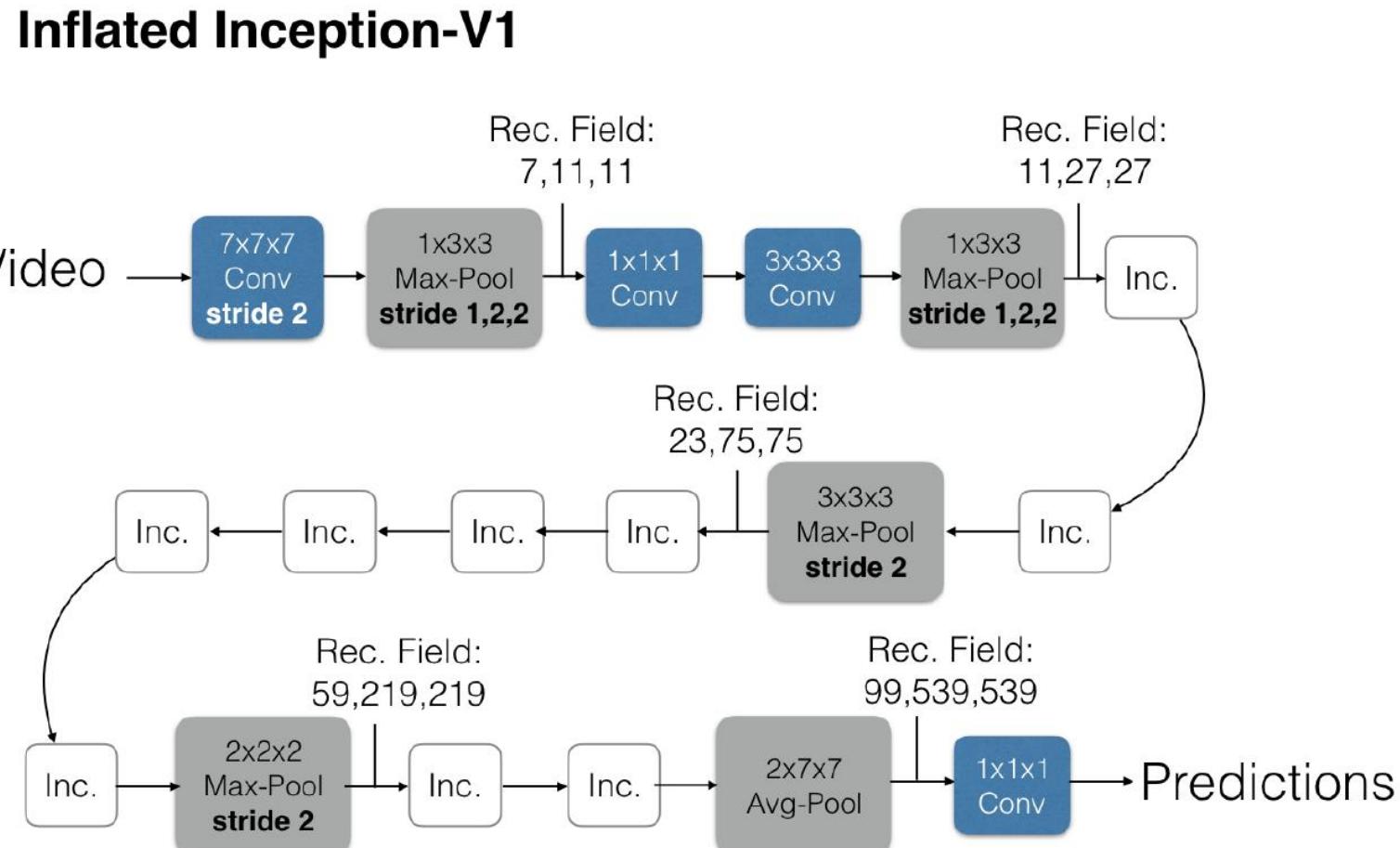


Video Classification: 3D CNN – I3D

- Inflating 2D ConvNets into 3D.**

Make square filters cubic – $N \times N$ filters become $N \times N \times N$.

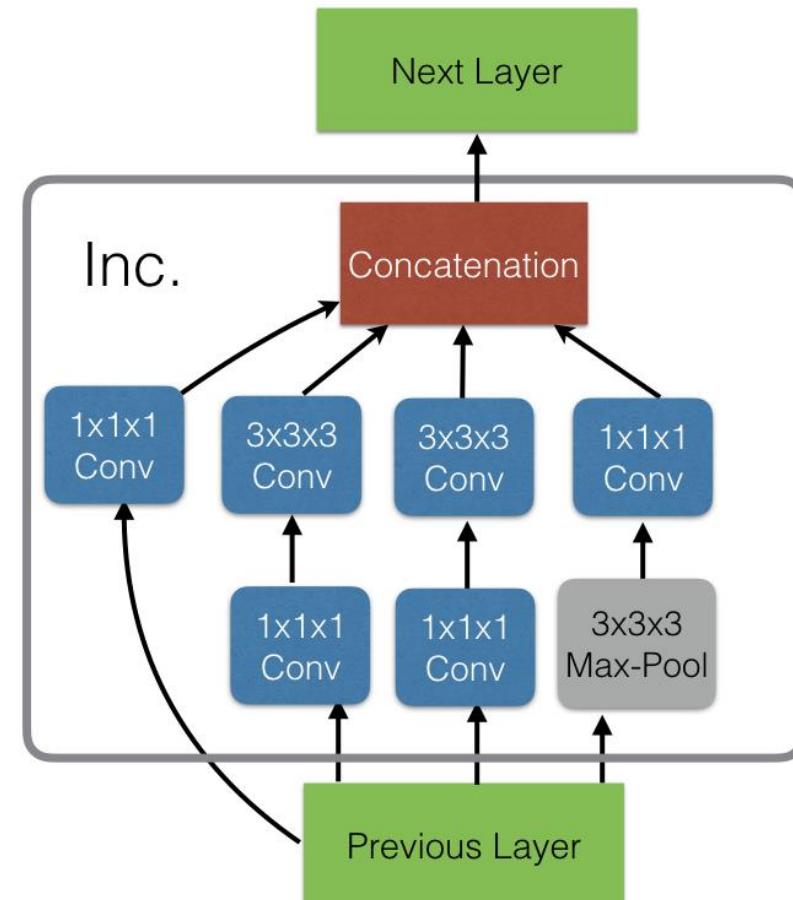
- Use ImageNet-pretrained Inception-V1 as base network.
Repeat the 2D pre-trained weights in the 3rd dimension



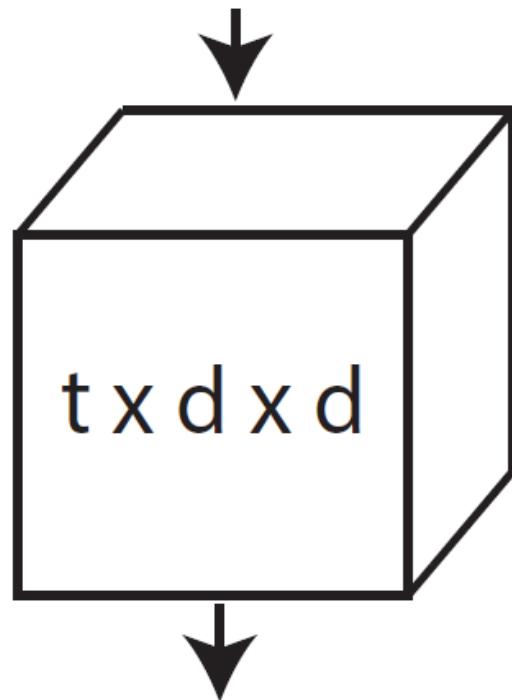
Video Classification: 3D CNN – I3D

- **Inflating 2D ConvNets into 3D.**
Make square filters cubic – $N \times N$ filters become $N \times N \times N$.
- Use ImageNet-pretrained Inception-V1 as base network.
Repeat the 2D pre-trained weights in the 3rd dimension

Inception Module (Inc.)



Video Classification: 3D CNN – R(2+1)D

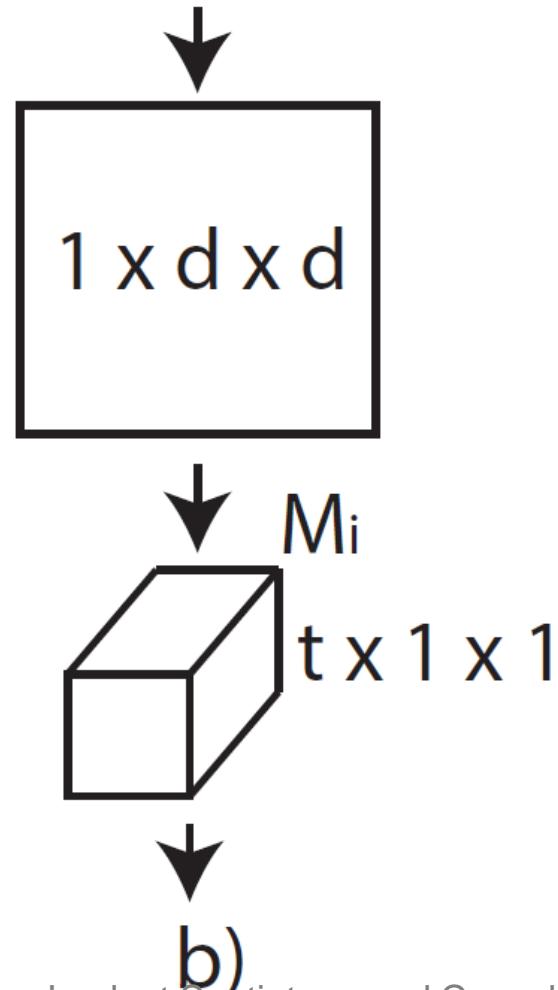


3D conv:

$$C_{out} \times C_{in} \times t \times d \times d$$

Full 3D convolution is carried out using a filter of size $t \times d \times d$ where t denotes the temporal extent and d is the spatial width and height.

Video Classification: 3D CNN – R(2+1)D

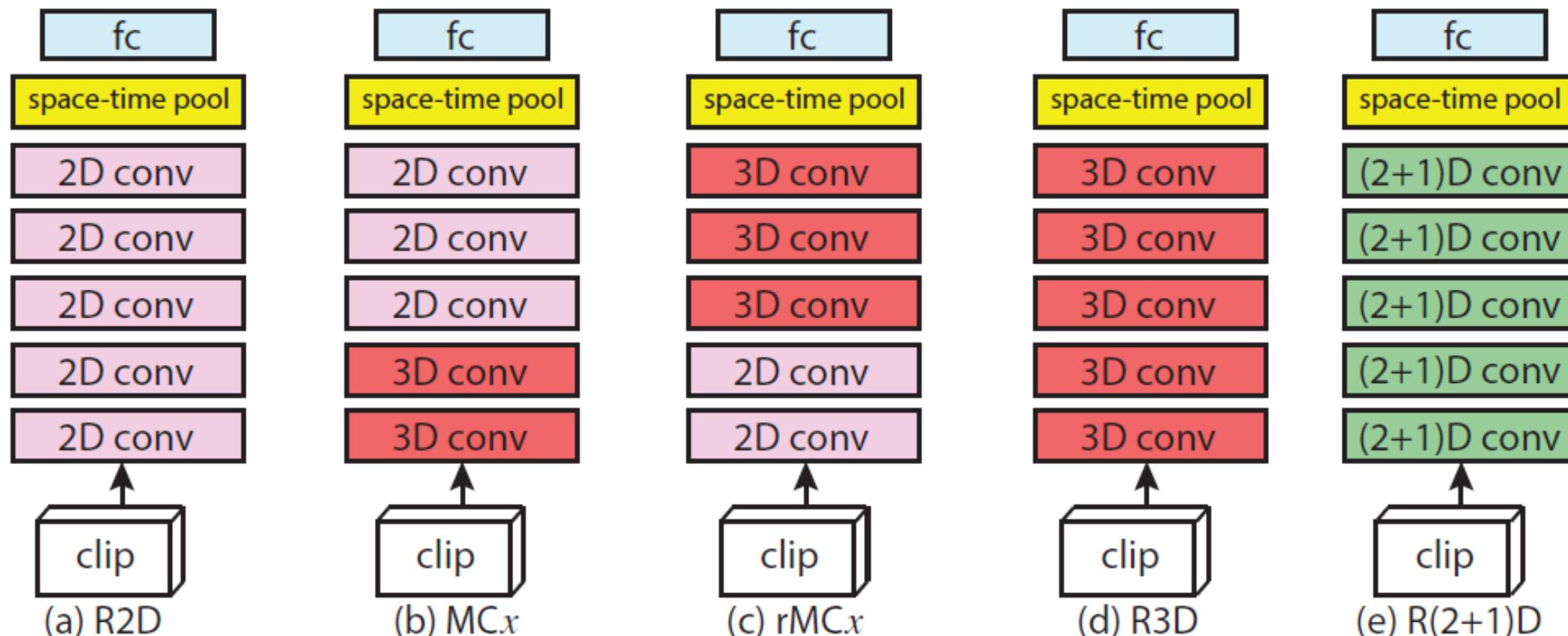


$(2+1)$ D conv:

$$C_{out} \times M \times 1 \times d \times d + M \times C_{in} \times t \times 1 \times 1$$

A $(2+1)$ D convolutional block splits the computation into a spatial 2D convolution followed by a temporal 1D convolution. The numbers of 2D filters (M) are hyper-parameters.

Video Classification: 3D CNN – R(2+1)D



Du Tran, et al. A Closer Look at Spatiotemporal Convolutions for Action Recognition. In Proc. CVPR, 2018.

Video classification: 3D CNN – performance comparison

method	Clip@1	Video@1	Video@5
DeepVideo [16]	41.9	60.9	80.2
C3D [36]	46.1	61.1	85.2
2D Resnet-152 [13]	46.5*	64.6*	86.4*
Conv pooling [42]	-	71.7	90.4
P3D [25]	47.9*	66.4*	87.4*
R3D-RGB-8frame	53.8	-	-
R(2+1)D-RGB-8frame	56.1	72.0	91.2
R(2+1)D-Flow-8frame	44.5	65.5	87.2
R(2+1)D-Two-Stream-8frame	-	72.2	91.4
R(2+1)D-RGB-32frame	57.0	73.0	91.5
R(2+1)D-Flow-32frame	46.4	68.4	88.7
R(2+1)D-Two-Stream-32frame	-	73.3	91.9

Table 4. Comparison with the state-of-the-art on Sports-1M. R(2+1)D outperforms C3D by 10.9%, and P3D by 9.1% and it achieves the best reported accuracy on this benchmark to date.

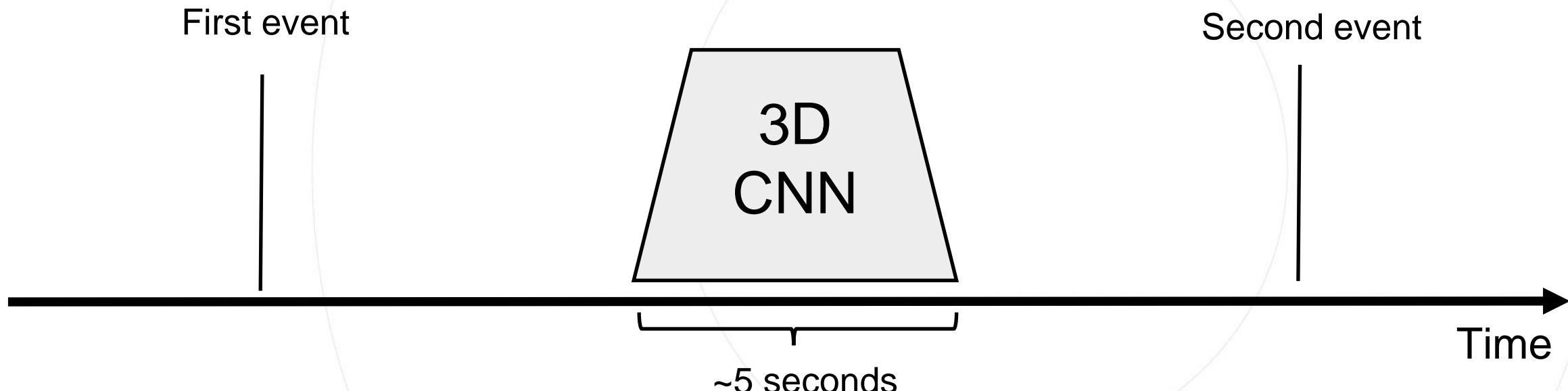
*These baseline numbers are taken from [25].

method	pretraining dataset	top1	top5
I3D-RGB [4]	none	67.5	87.2
I3D-RGB [4]	ImageNet	72.1	90.3
I3D-Flow [4]	ImageNet	65.3	86.2
I3D-Two-Stream [4]	ImageNet	75.7	92.0
R(2+1)D-RGB	none	72.0	90.0
R(2+1)D-Flow	none	67.5	87.2
R(2+1)D-Two-Stream	none	73.9	90.9
R(2+1)D-RGB	Sports-1M	74.3	91.4
R(2+1)D-Flow	Sports-1M	68.5	88.1
R(2+1)D-Two-Stream	Sports-1M	75.4	91.9

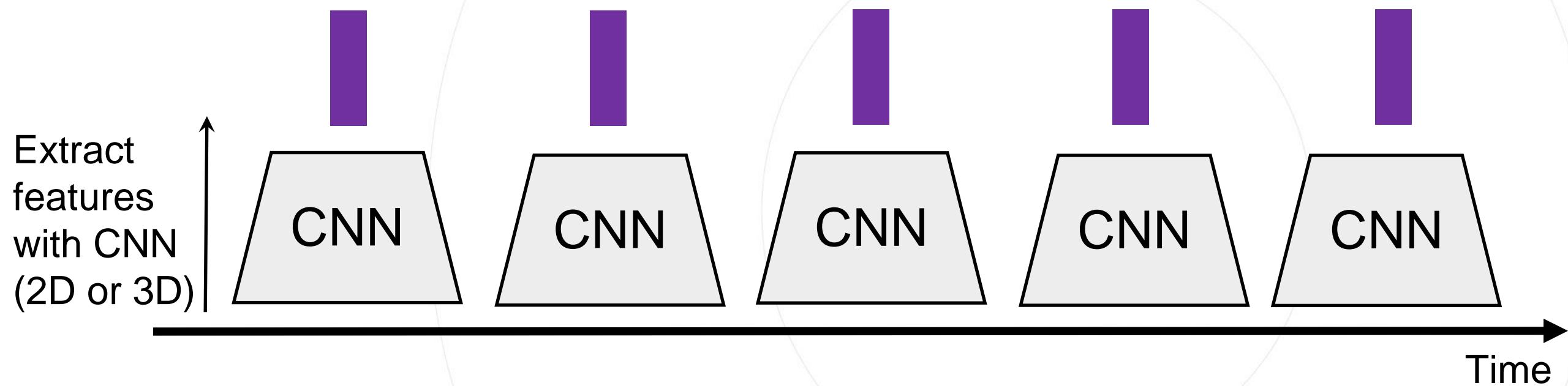
Table 5. Comparison with the state-of-the-art on Kinetics. R(2+1)D outperforms I3D by 4.5% when trained from scratch on RGB. R(2+1)D pretrained on Sports-1M outperforms I3D pretrained on ImageNet, for both RGB and optical flow. However, it is slightly worse than I3D (0.3%) when fusing the two streams.

Video classification: modeling long-term temporal structure

Temporal CNNs can only model local motion between frames in very short clips of ~2-5 seconds. **How to model long-term structure** is a remaining issue.

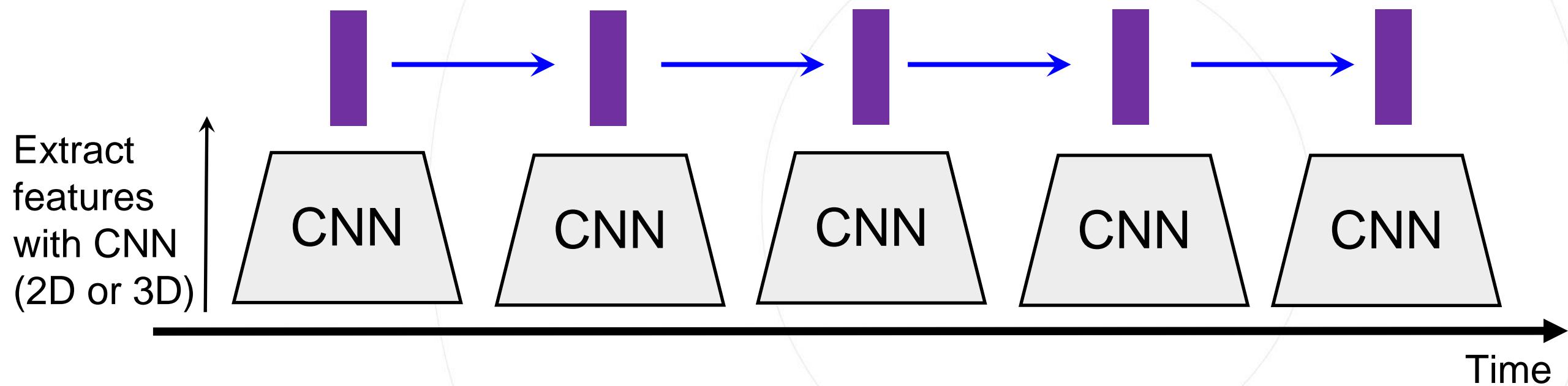


Video classification: modeling long-term temporal structure



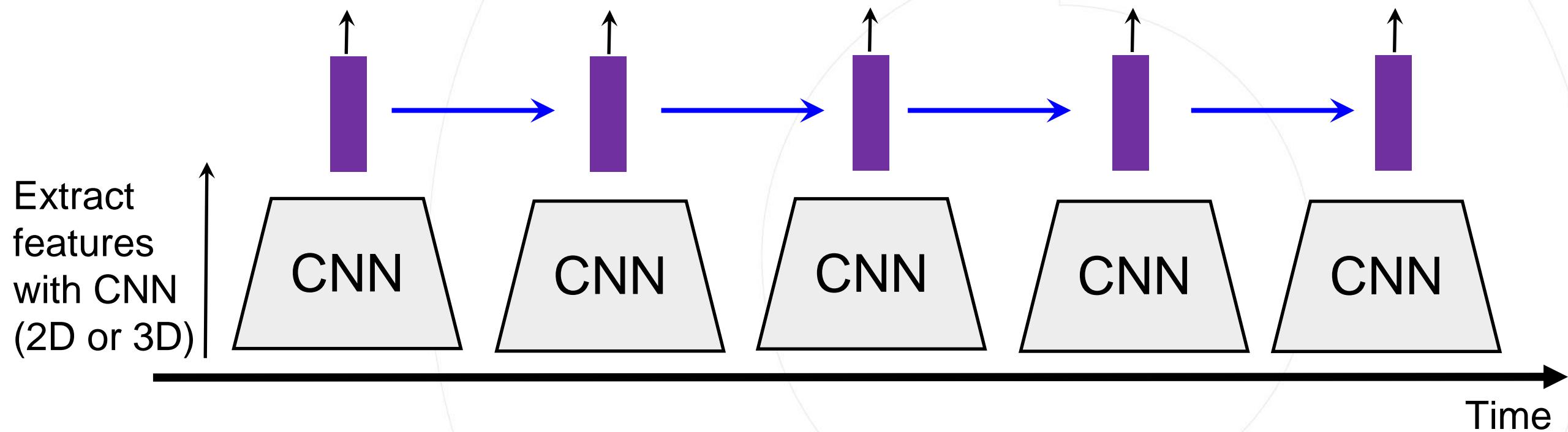
Video classification: modeling long-term temporal structure

Process local features with recurrent neural network (i.e., LSTM)



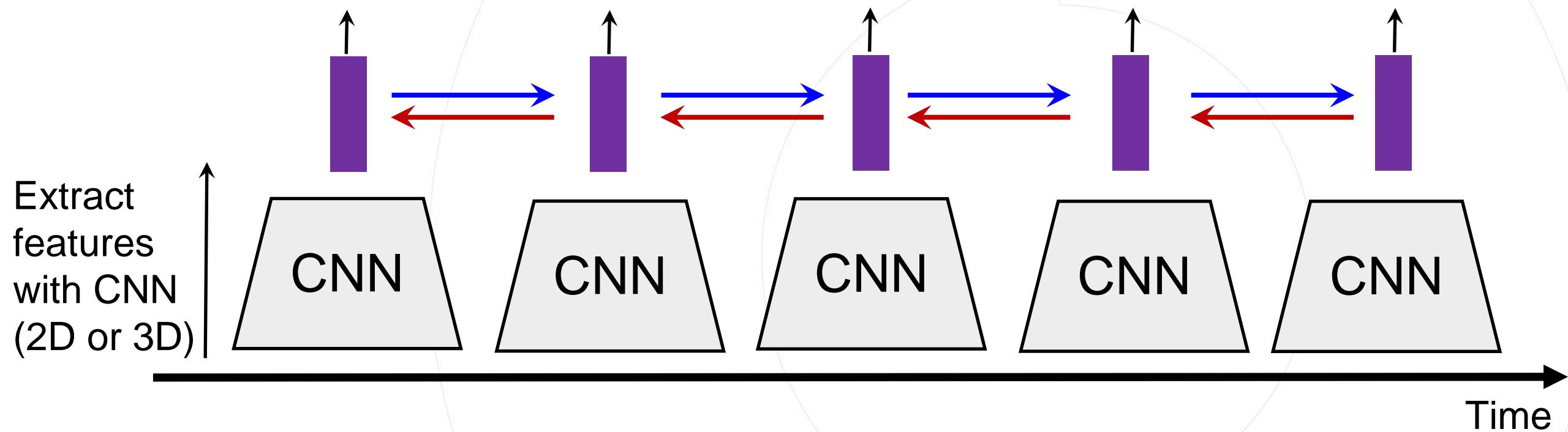
Video classification: modeling long-term temporal structure

Process local features with recurrent neural network (i.e., LSTM)
Many-to-many: one output per video frame

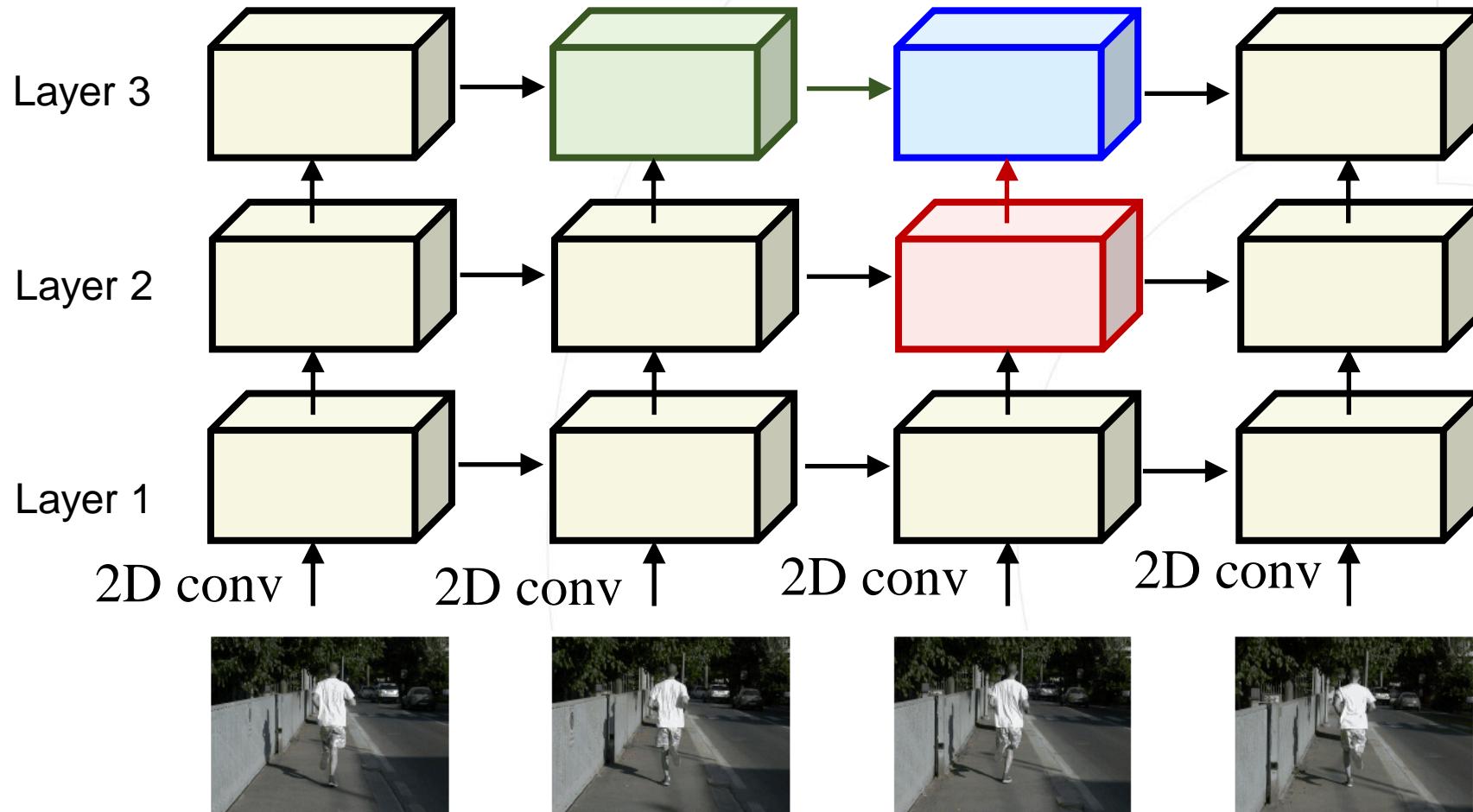


Video classification: modeling long-term temporal structure

Sometimes use the pre-trained CNN as the feature extractor,
Don't backpropagate to CNN to save memory.

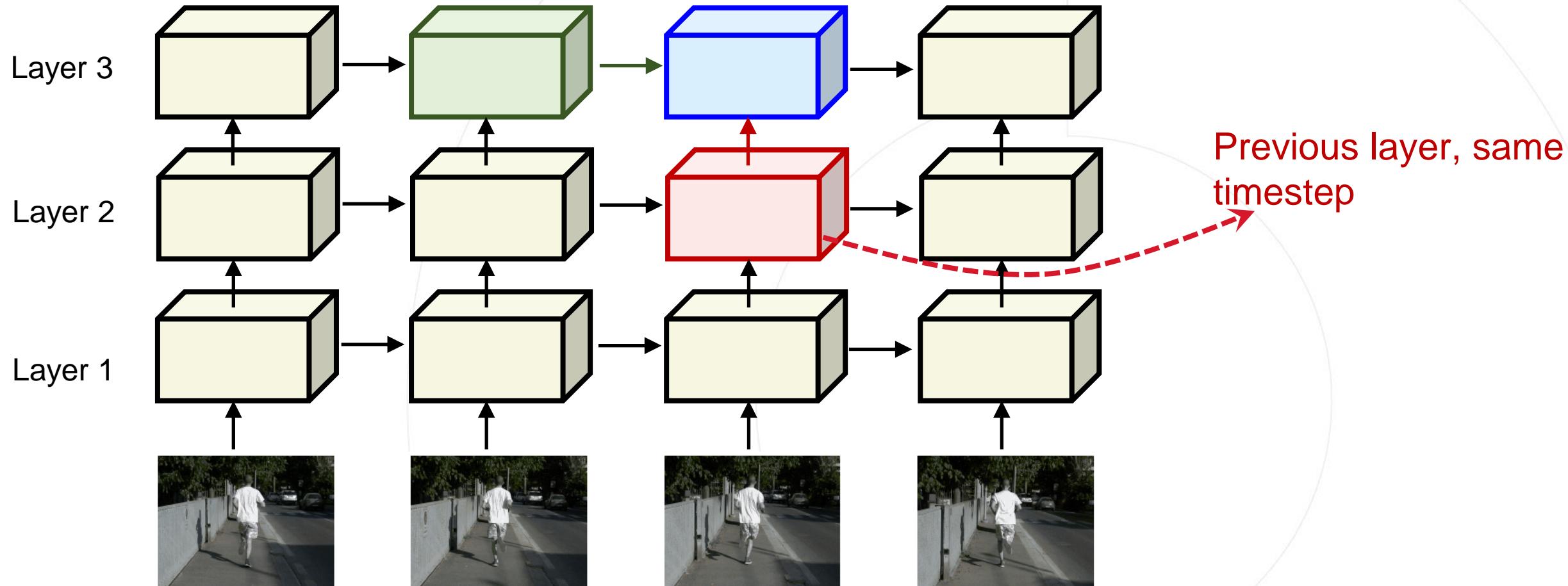


Recurrent Neural Networks

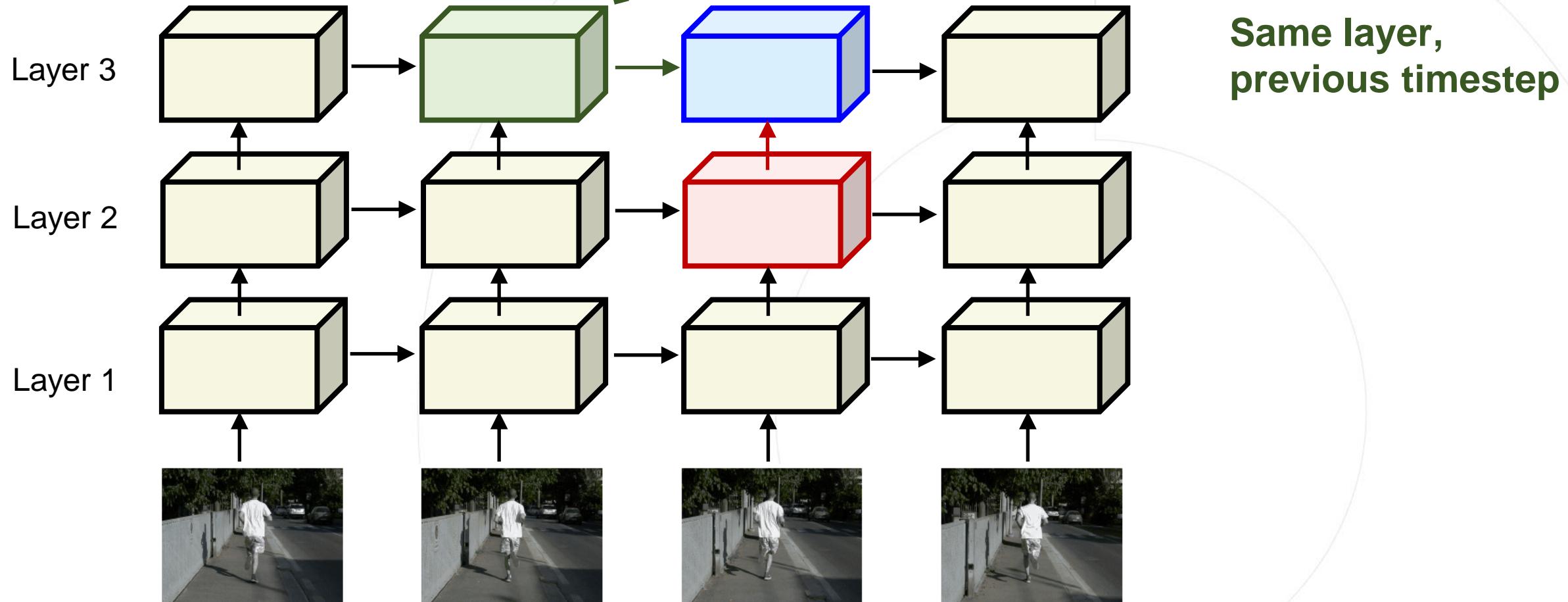


Entire network
uses 2D feature
maps: $C \times H \times W$

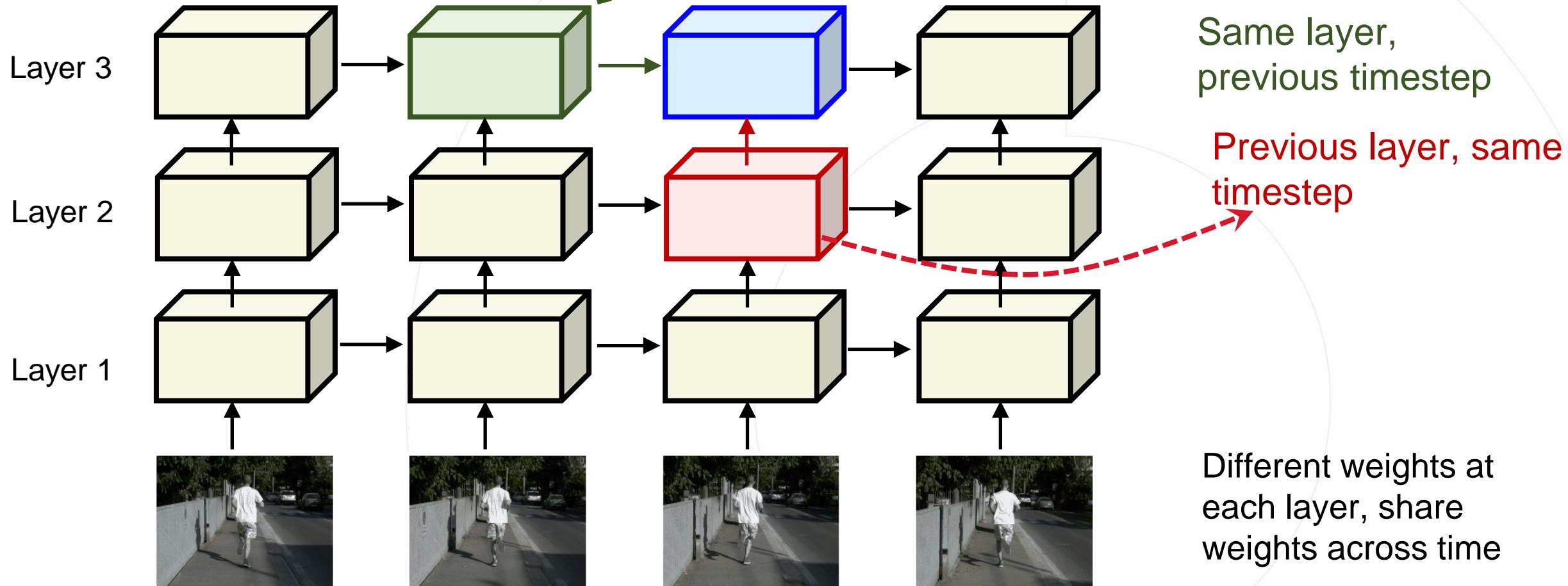
Recurrent Neural Networks



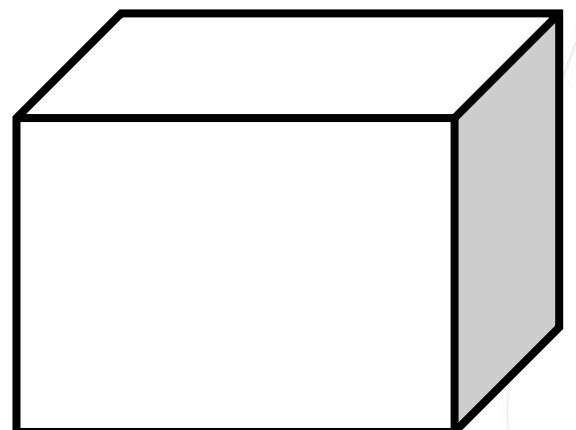
Recurrent Neural Networks



Recurrent Neural Networks

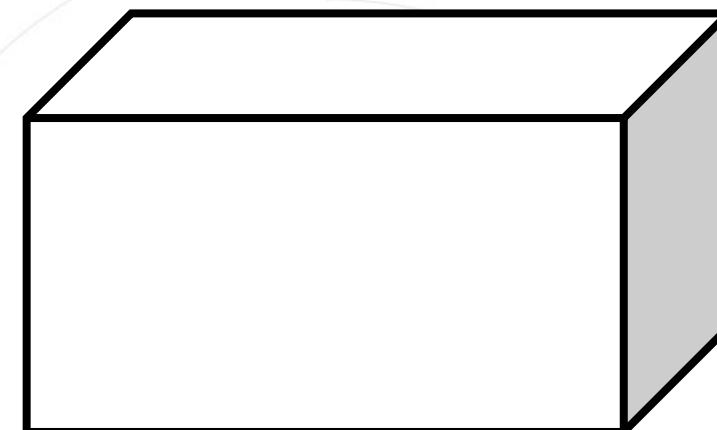


2D convolution



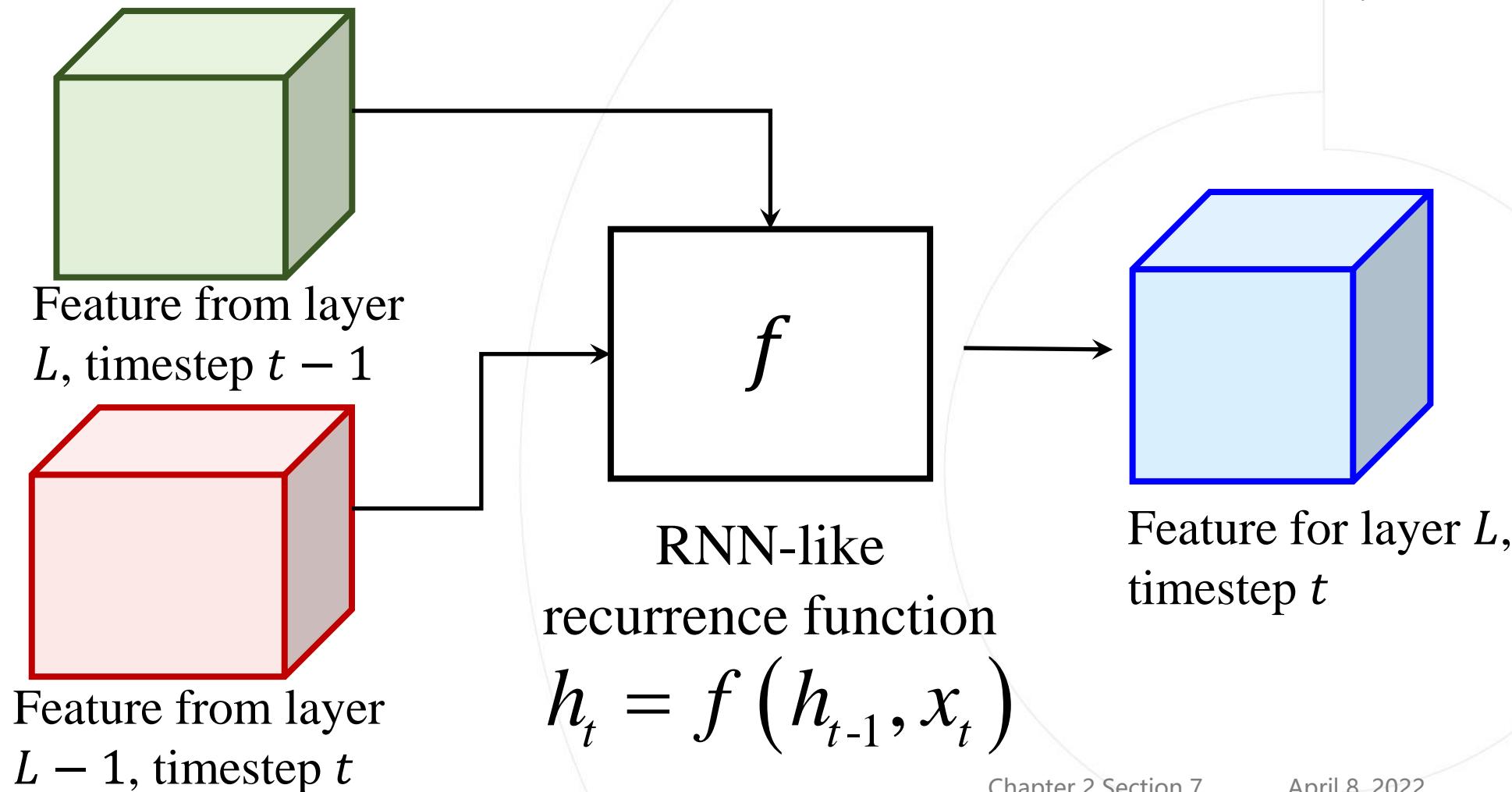
Input features:
 $C \times H \times W$

2D conv
→



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

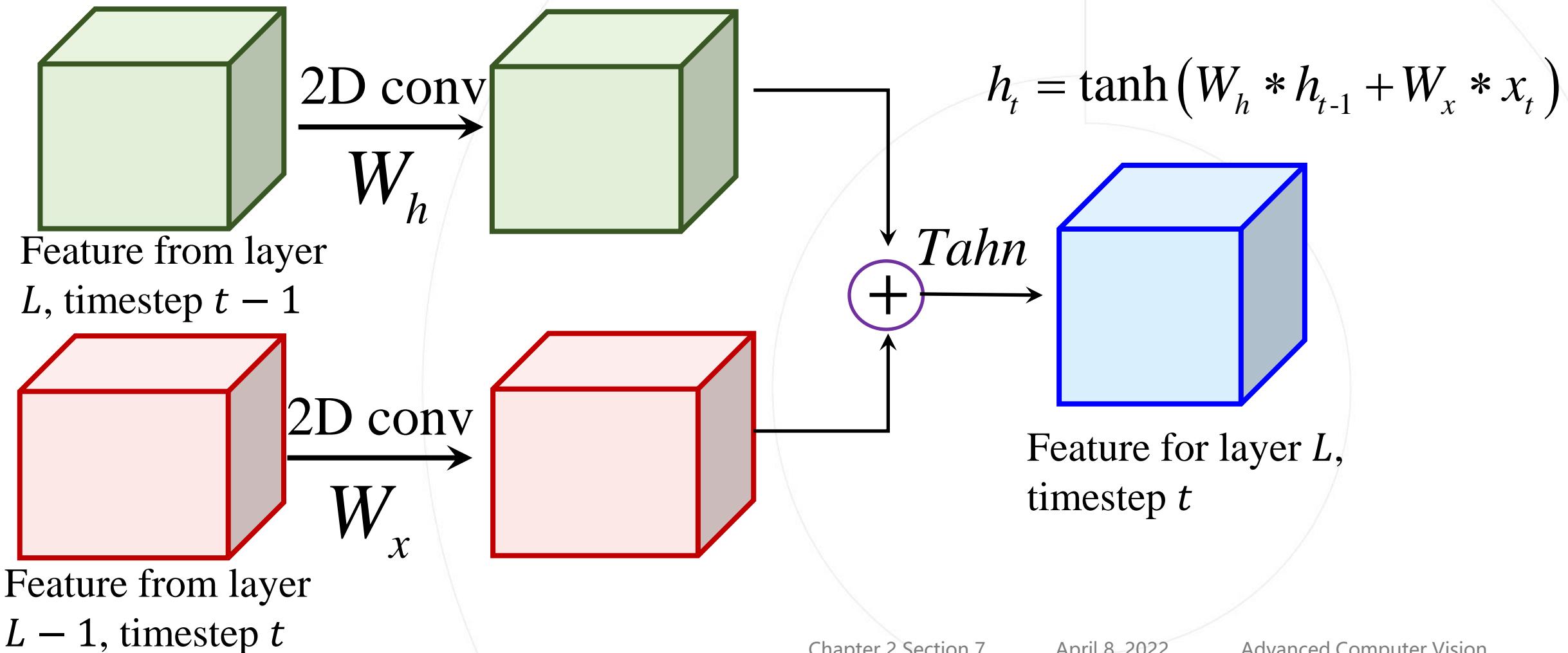
Recurrent Neural Networks



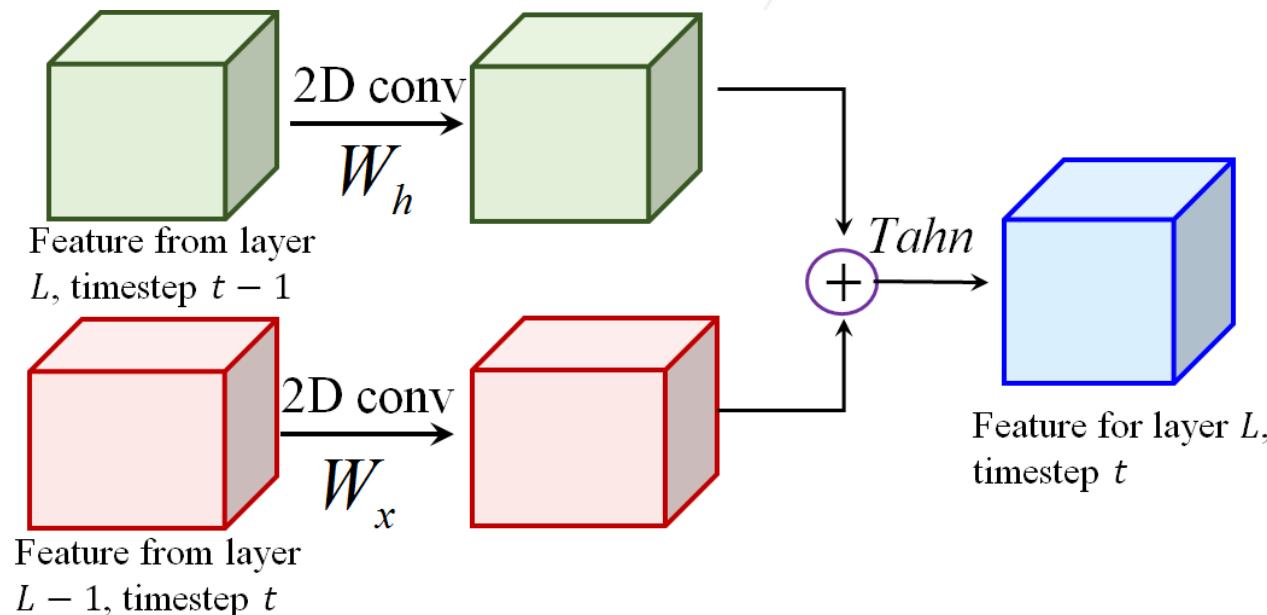
Vanilla RNN: $h_t = \tanh(W_h h_{t-1} + W_x x_t)$

Recurrent Convolutional Networks

Replace all matrix multiply in Vanilla RNN with convolution



Recurrent Convolutional Networks



GRU

$$r_t = \sigma(W_{xr} * x_t + W_{hr} * h_{t-1} + b_r)$$

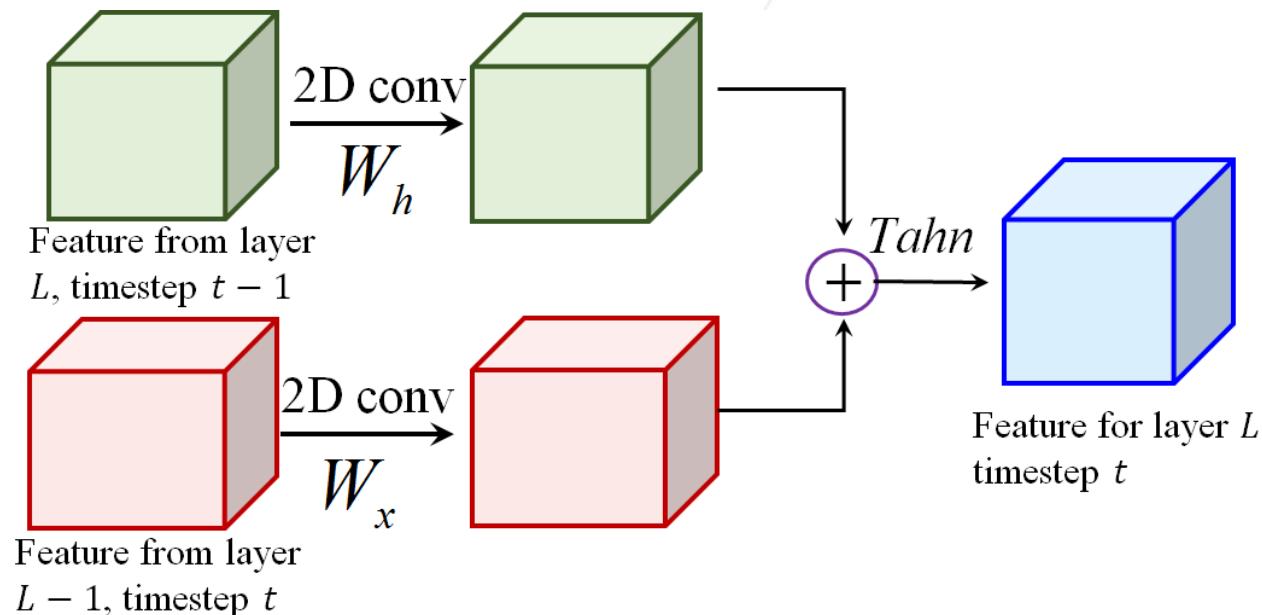
$$z_t = \sigma(W_{xz} * x_t + W_{hz} * h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W_{xh} * x_t + W_{hh} * (r_t \square h_{t-1}) + b_h)$$

$$h_t = z_t \square h_{t-1} + (1 - z_t) \square \tilde{h}_t$$

Do similar transform for other RNN variants

Recurrent Convolutional Networks



LSTM

$$g_t = \tanh(W_{xg} * x_t + W_{hg} * h_{t-1} + b_g)$$

$$i_t = \sigma(W_{xi} * x_t + W_{hi} * h_{t-1} + b_i)$$

$$f_t = \sigma(W_{xf} * x_t + W_{hf} * h_{t-1} + b_f)$$

$$o_t = \sigma(W_{xo} * x_t + W_{ho} * h_{t-1} + b_o)$$

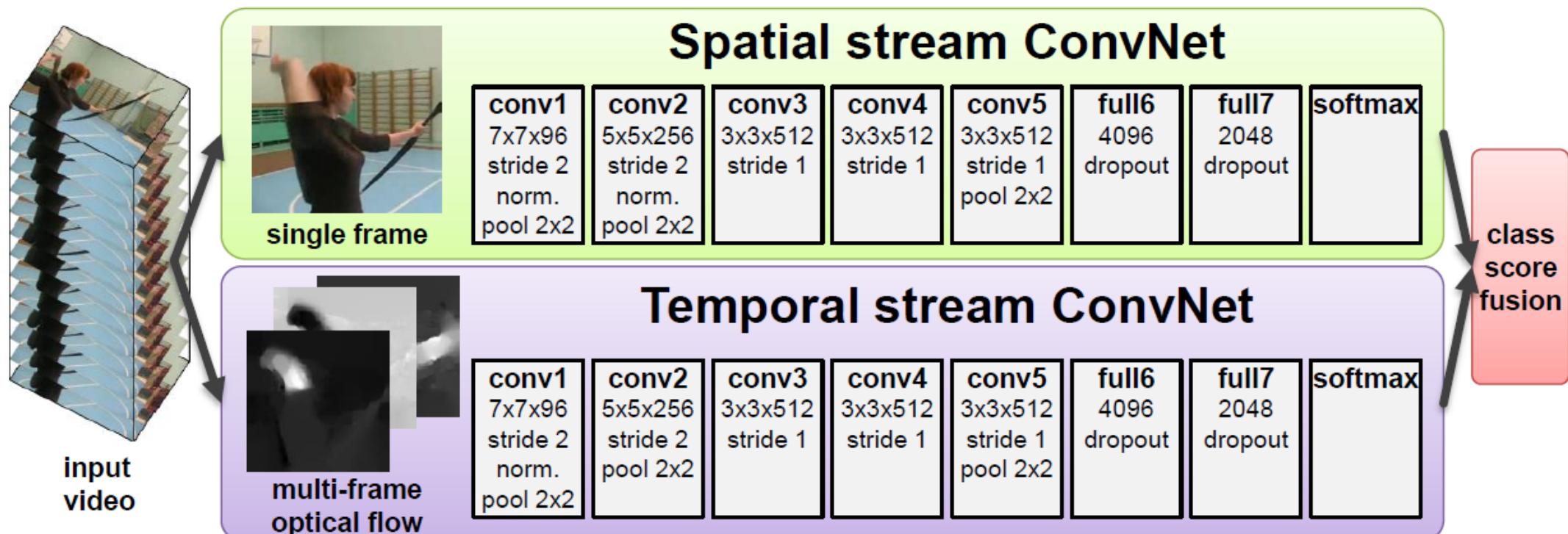
$$c_t = i_t \square g_t + f_t \square c_{t-1}$$

$$h_t = o_t \square \tanh(c_t)$$

Do similar transform for other RNN variants

Two-stream networks - separating motion & appearance

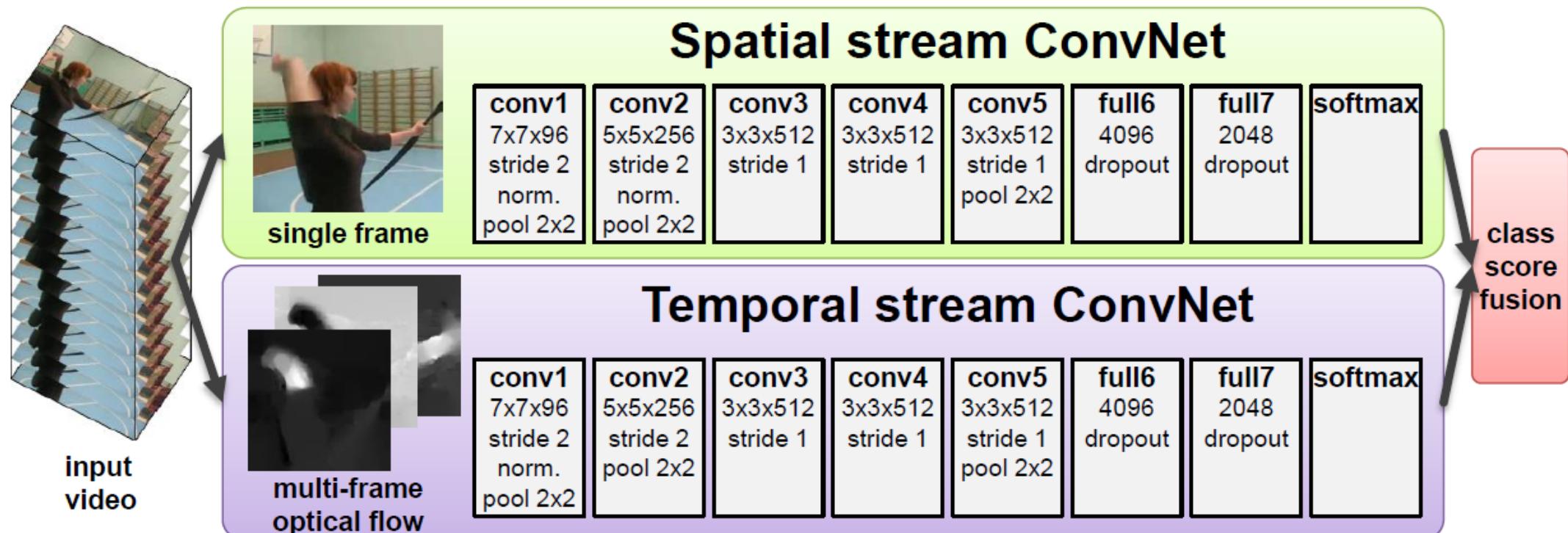
- **Spatial stream:** performs action recognition from still video frames
- **Temporal stream:** recognizes action from motion in the form of dense optical flow.



Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In NIPS, 2014.

Video classification: two-stream networks – separating motion & appearance

- **Input of spatial stream:** single image $3 \times H \times W$
- **Input of temporal stream:** stack of optical flow $2(T - 1) \times H \times W$, first 2D conv processes all flow images



Optical Flow: measuring motion

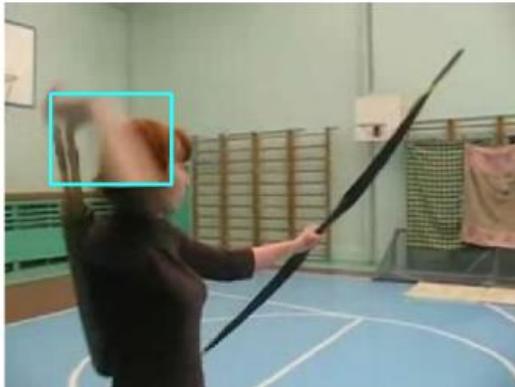


Image frame at time t

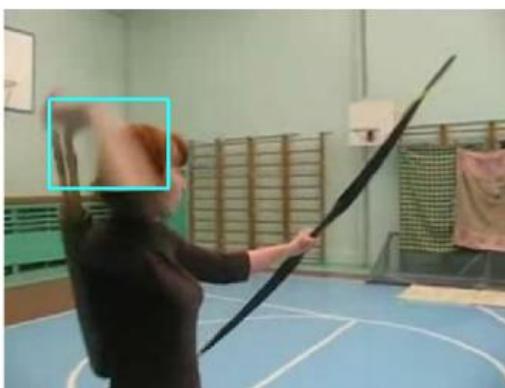
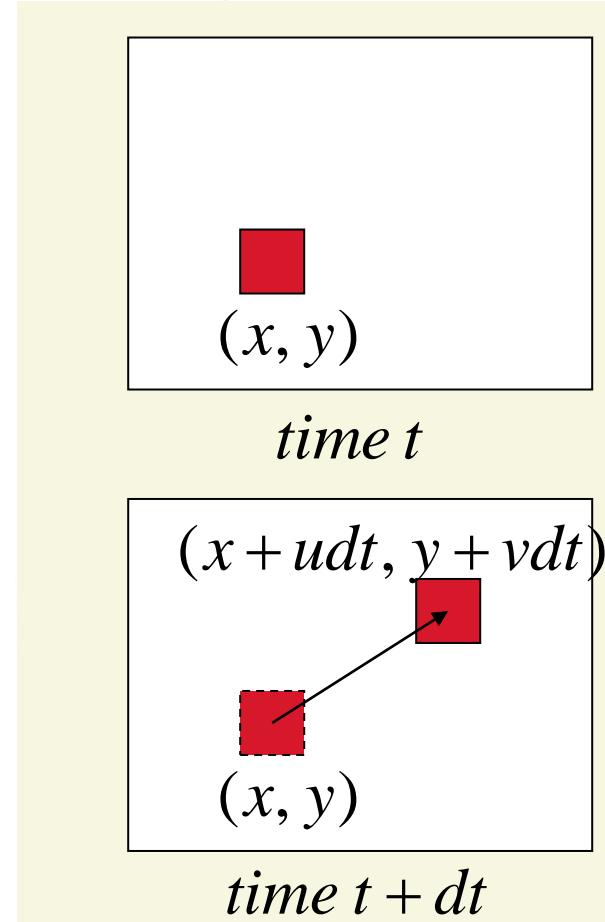


Image frame at time $t + 1$



Optical Flow: Velocities (u, v)

Displacement: $(dx, dy) = (udt, vdt)$

- **Brightness constancy:** assume brightness of patch remains same in both images:

$$I(x + dx, y + dy, t + dt) = I(x, y, t)$$

- **Small motion:** assume brightness of patch remains same in both images:

$$I(x + dx, y + dy, t + dt)$$

$$\square I(x, y, t) + \frac{\partial I}{\partial x} dx + \frac{\partial I}{\partial y} dy + \frac{\partial I}{\partial t} dt$$

Optical Flow: measuring motion

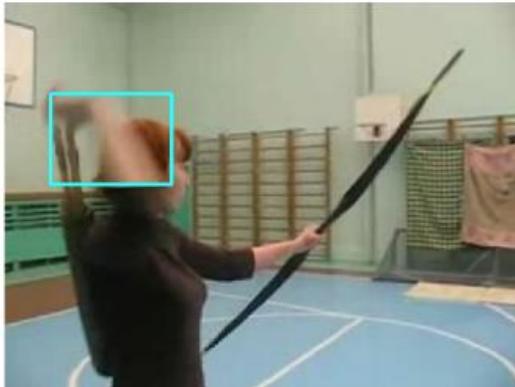


Image frame at time t

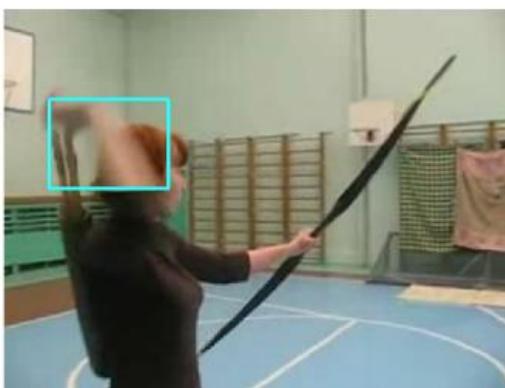
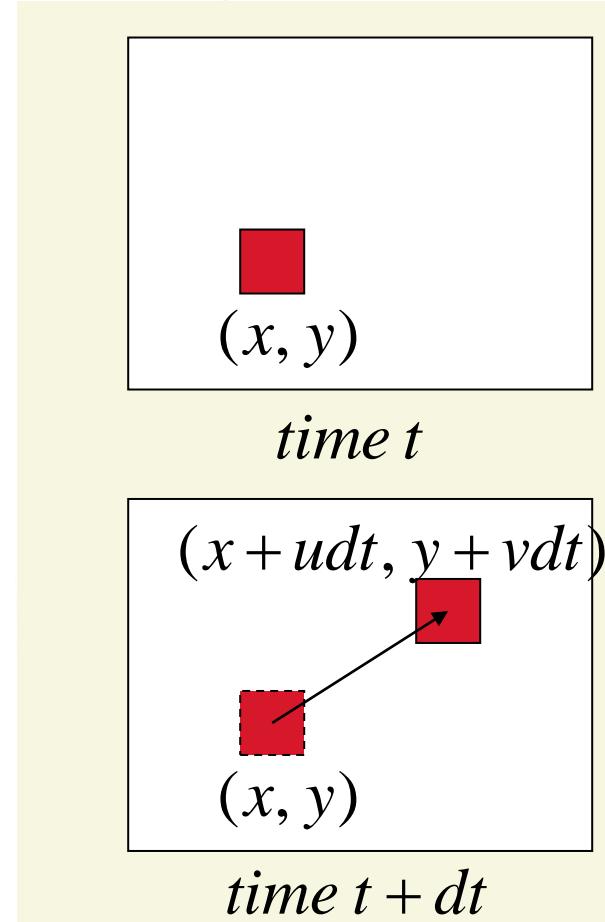


Image frame at time $t + 1$



Optical Flow: Velocities (u, v)

Displacement: $(dx, dy) = (udt, vdt)$

- **Brightness constancy:** assume brightness of patch remains same in both images

$$I(x + dx, y + dy, t + dt) = I(x, y, t)$$

- **Small motion:** assume brightness of patch remains same in both images

$$\cancel{I(x, y, t) + \frac{\partial I}{\partial x} dx + \frac{\partial I}{\partial y} dy + \frac{\partial I}{\partial t} dt = I(x, y, t)}$$

Optical Flow: measuring motion

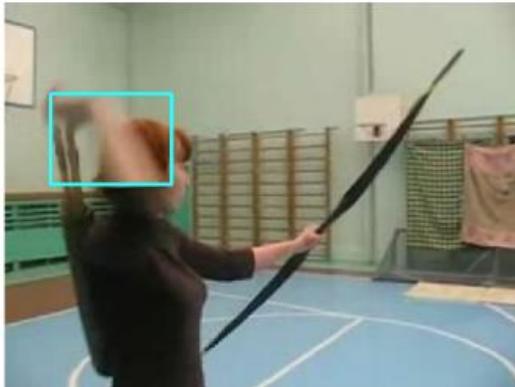


Image frame at time t

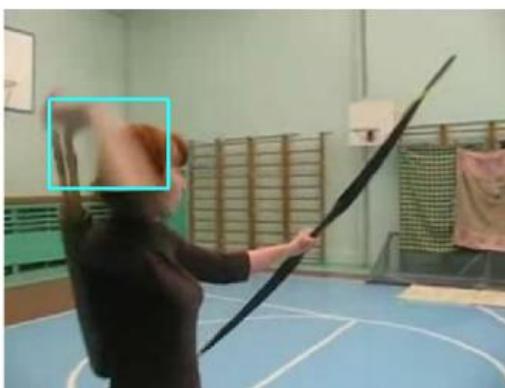
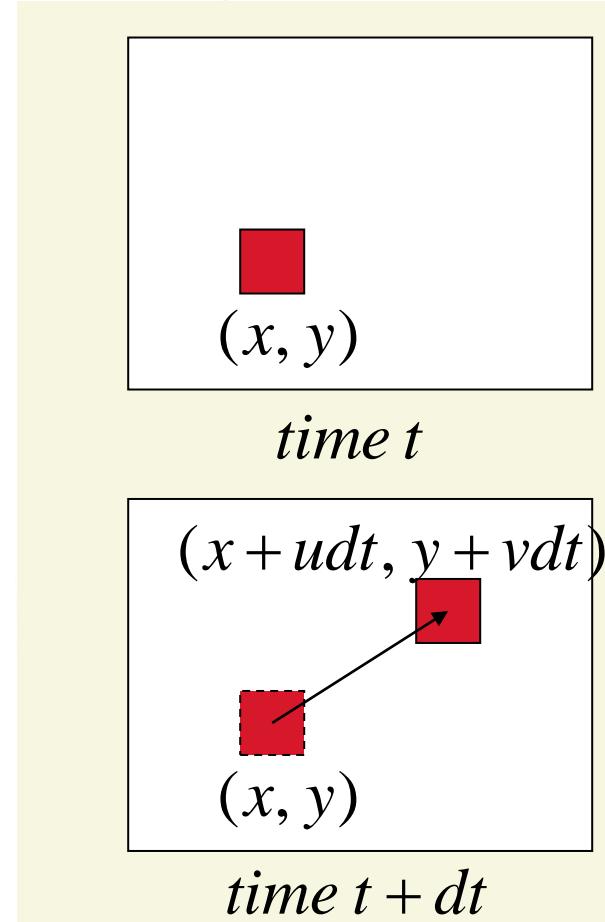


Image frame at time $t + 1$



Optical Flow: Velocities (u, v)

Displacement: $(dx, dy) = (udt, vdt)$

- **Brightness constancy:** assume brightness of patch remains same in both images

$$I(x + dx, y + dy, t + dt) = I(x, y, t)$$

- **Small motion:** assume brightness of patch remains same in both images

$$\begin{aligned} I_x dx + I_y dy + I_t dt &= 0 \Rightarrow I_x \frac{dx}{dt} + I_y \frac{dy}{dt} + I_t = 0 \\ \Rightarrow I_x u + I_y v + I_t &= 0 \end{aligned}$$

Optical Flow: measuring motion

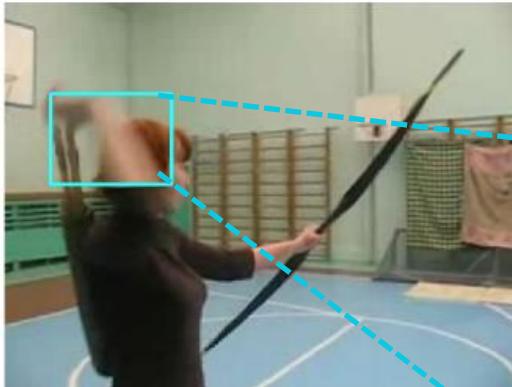


Image frame at time t

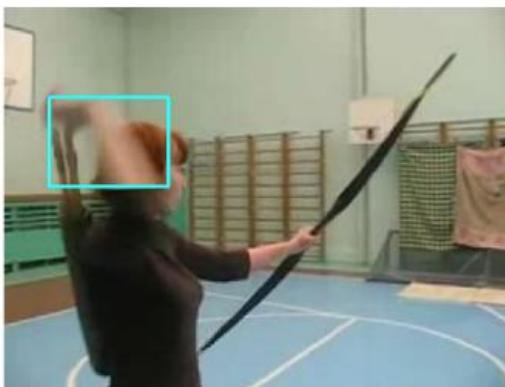
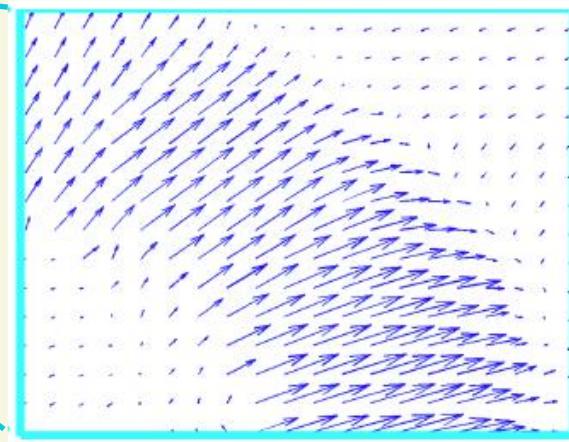


Image frame at time $t + 1$

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



Optical flow tells where each pixel will move in the next frames

$$F(x, y) = (u, v)$$

$$I_{t+1}(x+u, y+v) = I_t(x, y)$$

Optical Flow: measuring motion

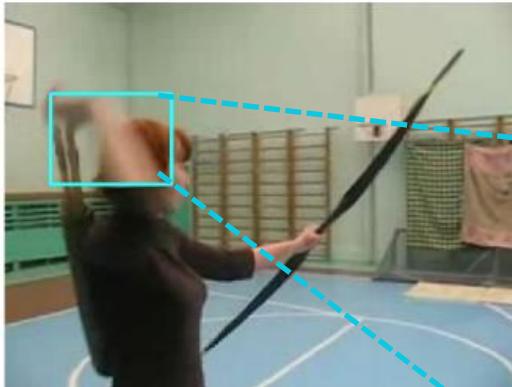


Image frame at time t

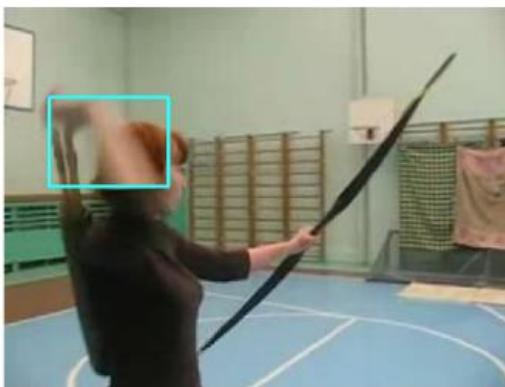
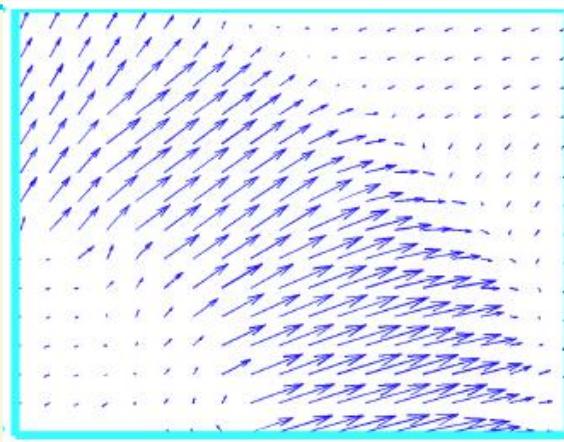


Image frame at time $t + 1$

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



Optical flow tells where each pixel will move in the next frames

$$F(x, y) = (u, v)$$

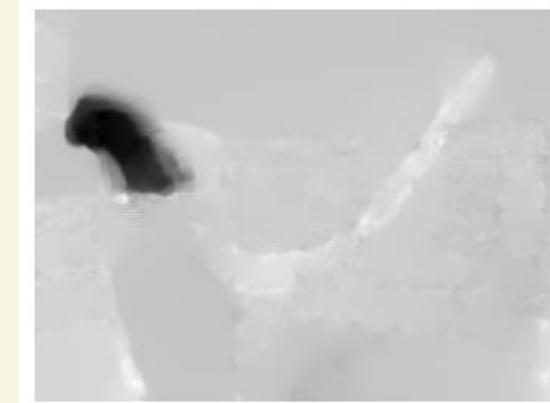
$$I_{t+1}(x+u, y+v) = I_t(x, y)$$

Optical flow highlights local motion.

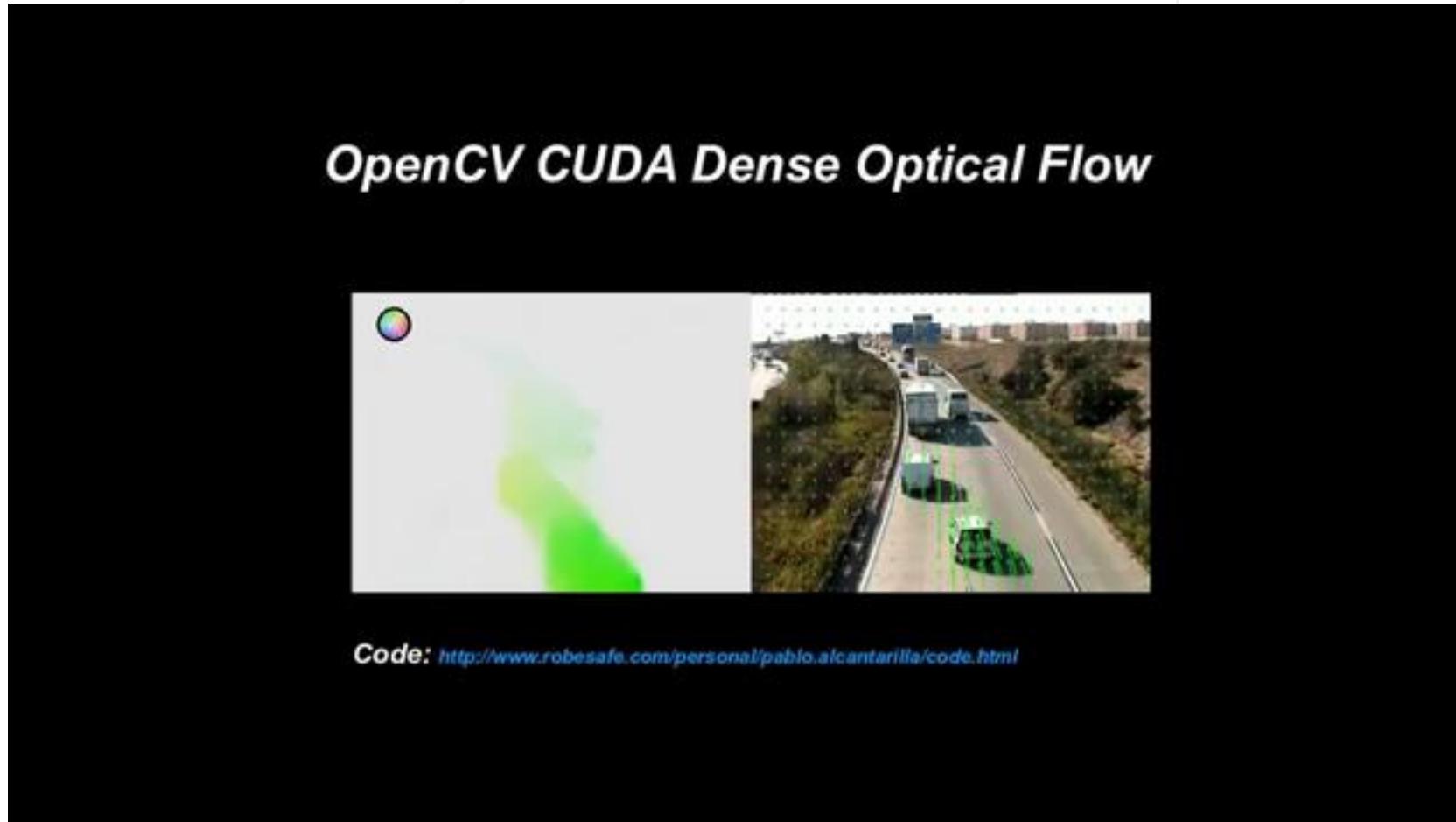
Horizontal flow u



Vertical flow v



Optical Flow: measuring motion



Optical Flow: measuring motion

FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks

Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, Thomas Brox

University of Freiburg, Germany

—— Supplementary Material ——

Optical Flow: measuring motion

- a) Optical flow stacking;
- c) Bi-directional optical flow;

- b) Trajectory stacking;
- d) Mean flow subtraction.

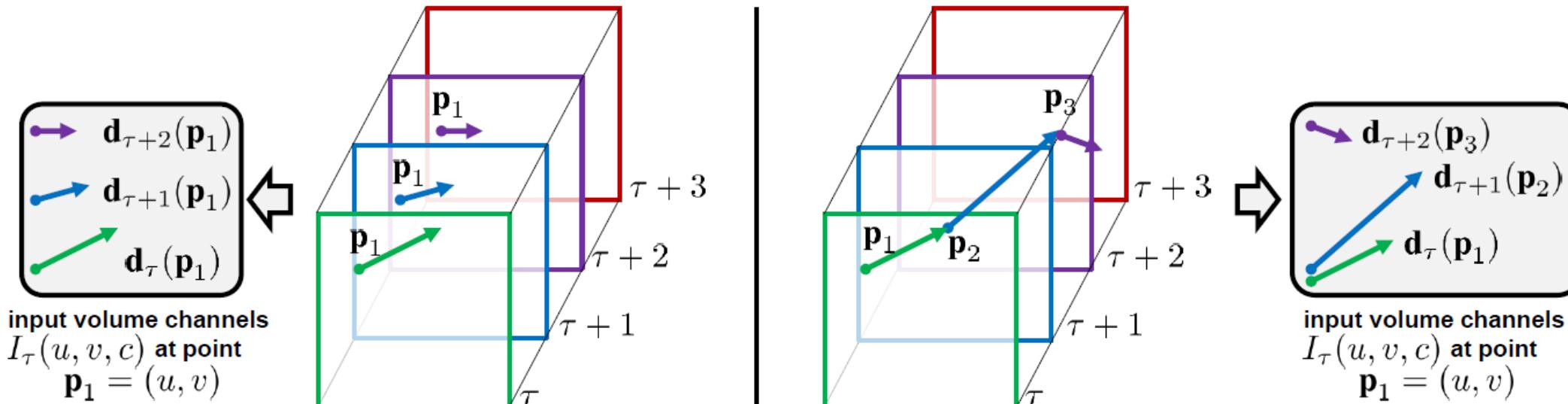
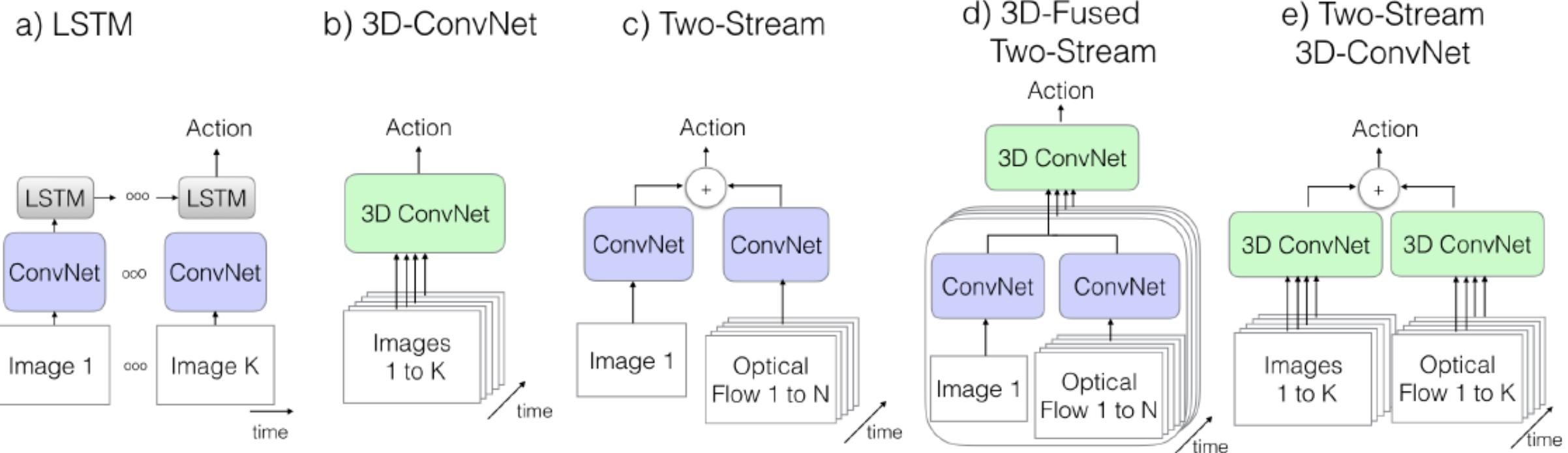
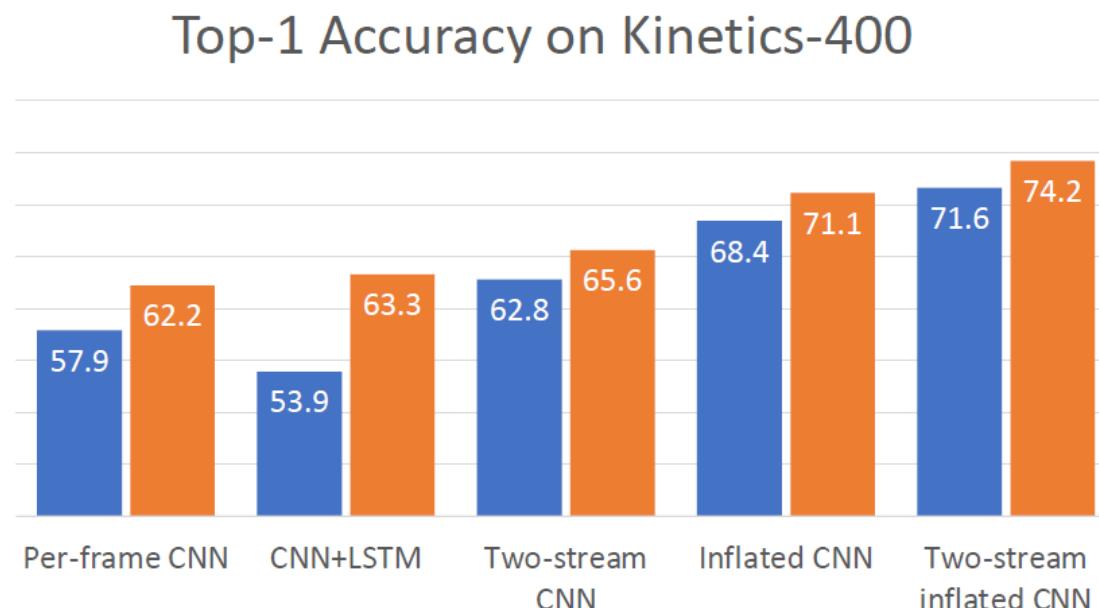


Figure 3: **ConvNet input derivation from the multi-frame optical flow.** *Left:* optical flow stacking (1) samples the displacement vectors \mathbf{d} at the same location in multiple frames. *Right:* trajectory stacking (2) samples the vectors along the trajectory. The frames and the corresponding displacement vectors are shown with the same colour.

Video classification: I3D (Two-stream)



Video classification: performance comparison

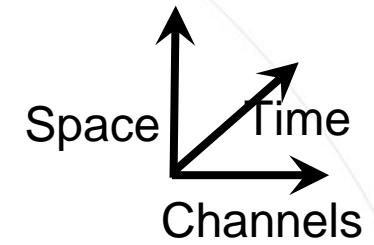
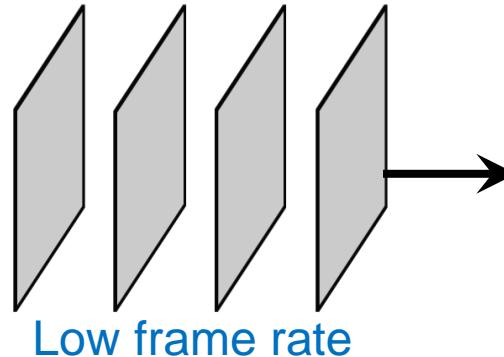


method	pretraining dataset	top1	top5
I3D-RGB [4]	none	67.5	87.2
I3D-RGB [4]	ImageNet	72.1	90.3
I3D-Flow [4]	ImageNet	65.3	86.2
I3D-Two-Stream [4]	ImageNet	75.7	92.0
R(2+1)D-RGB	none	72.0	90.0
R(2+1)D-Flow	none	67.5	87.2
R(2+1)D-Two-Stream	none	73.9	90.9
R(2+1)D-RGB	Sports-1M	74.3	91.4
R(2+1)D-Flow	Sports-1M	68.5	88.1
R(2+1)D-Two-Stream	Sports-1M	75.4	91.9

Table 5. **Comparison with the state-of-the-art on Kinetics.**
R(2+1)D outperforms I3D by 4.5% when trained from scratch on RGB. R(2+1)D pretrained on Sports-1M outperforms I3D pretrained on ImageNet, for both RGB and optical flow. However, it is slightly worse than I3D (0.3%) when fusing the two streams.

SlowFast Networks: treating time and space differently

- **Slow pathway.** Operating at low frame rate, to capture spatial semantics



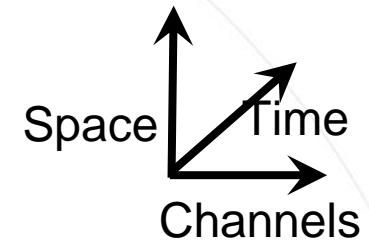
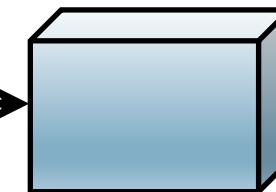
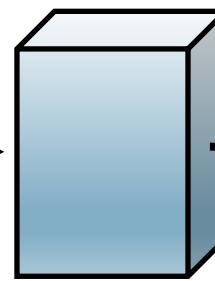
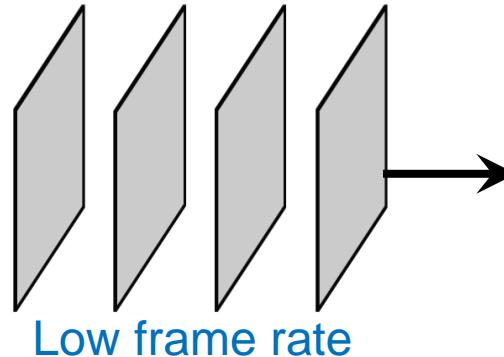
A *large* temporal stride τ on input frames.
Processes only one out of τ frames.
Typical value: $\tau = 16$ (roughly 2 frames sampled per second for 30-fps videos).

SlowFast Networks: treating time and space differently

- **Slow pathway.** Operating at low frame rate, to capture spatial semantics



Slow



Fast

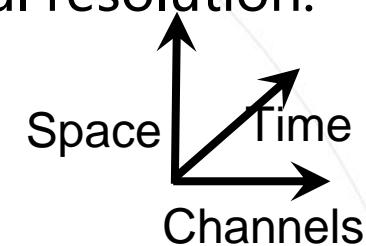
Denoting the number of frames sampled by the **Slow pathway** as T , the raw clip length is $T \times \tau$ frames.

SlowFast Networks: treating time and space differently

- **Fast pathway.** Operates at high frame rate, to capture motion at fine temporal resolution.



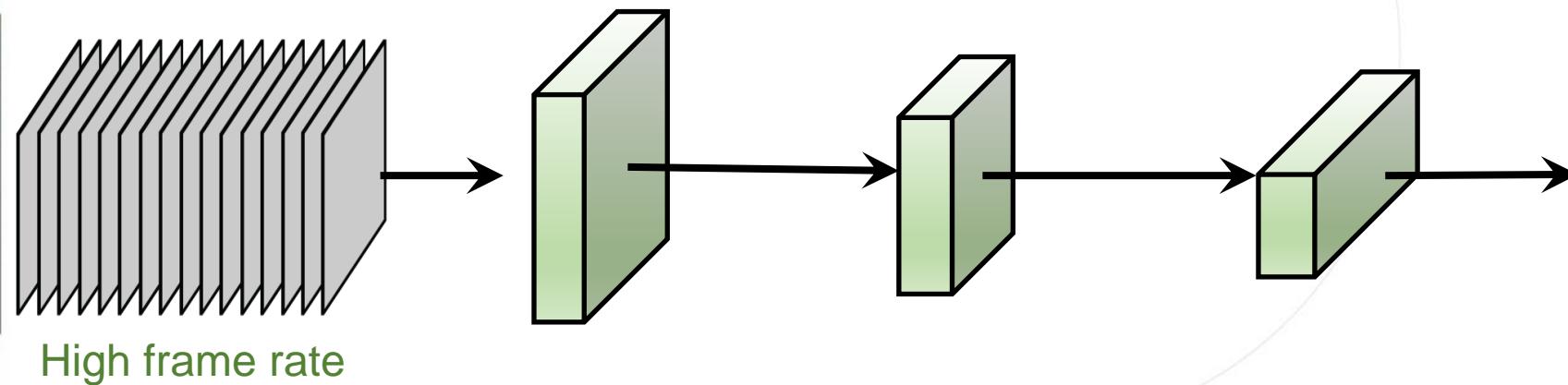
High frame rate. Works with a *small* temporal stride of τ/α , where $\alpha > 1$ is the frame rate ratio between the Fast and Slow pathways.



Slow



Fast



High frame rate

SlowFast Networks: treating time and space differently

- **Fast pathway.** Operates at high frame rate, to capture motion at fine temporal resolution.



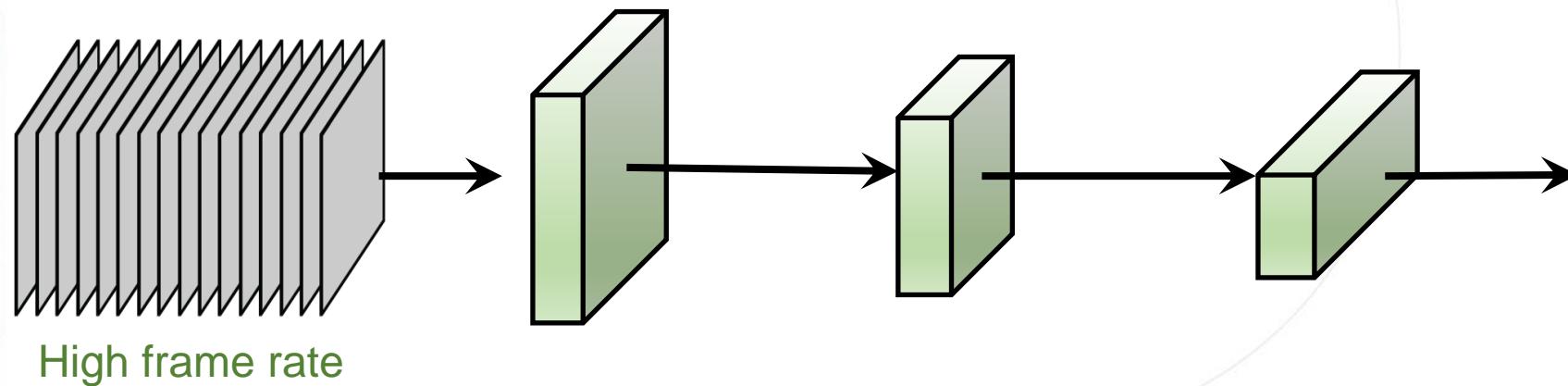
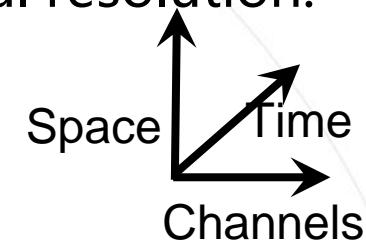
Slow



Fast

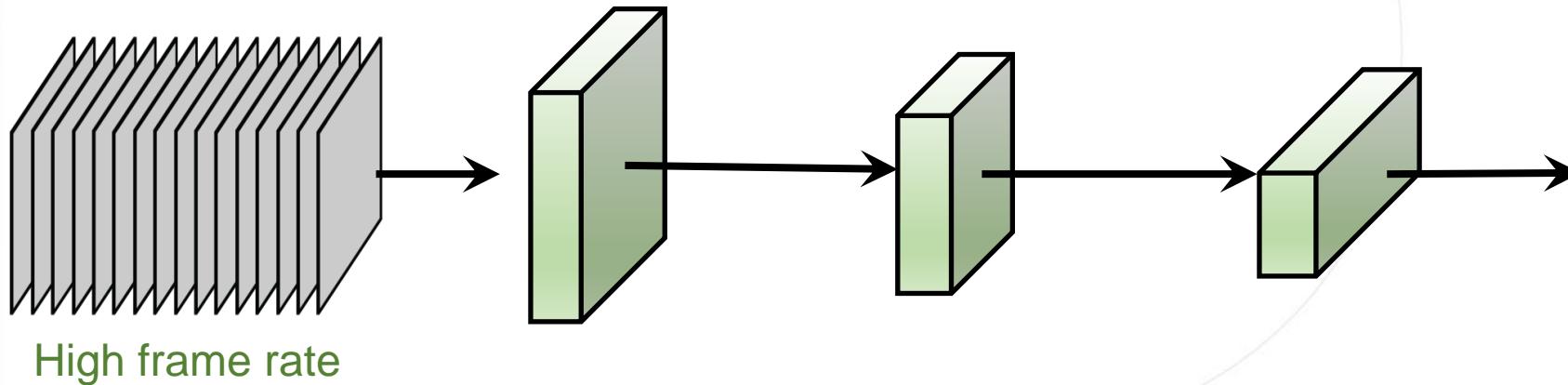
High temporal resolution features.

No temporal down-sampling layers
(neither temporal pooling nor time-strided convolutions).



SlowFast Networks: treating time and space differently

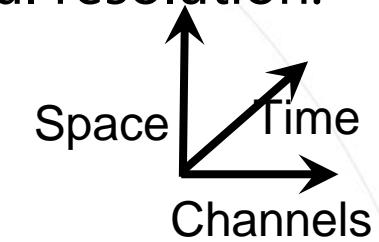
- **Fast pathway.** Operates at high frame rate, to capture motion at fine temporal resolution.



Low channel capacity. Analogous to the Slow pathway, but has a ratio of β ($\beta < 1$) channels of the Slow pathway.

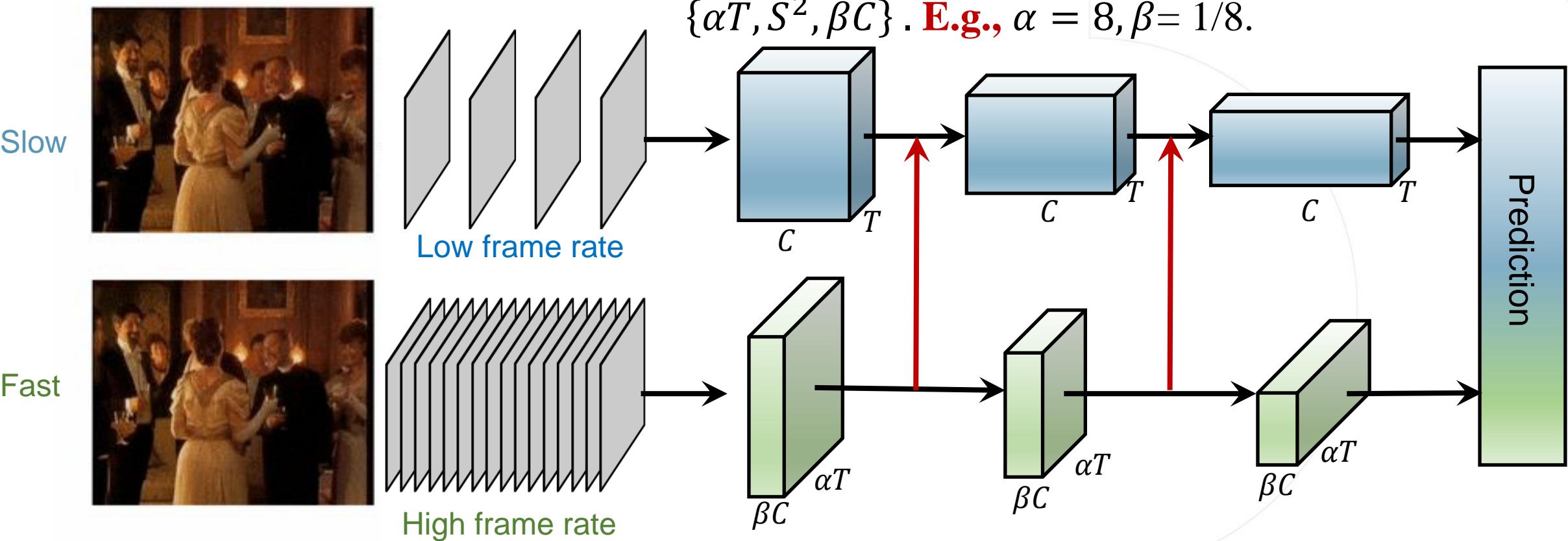
Typical value: $\beta = 1/8$.

~20% of the total computation



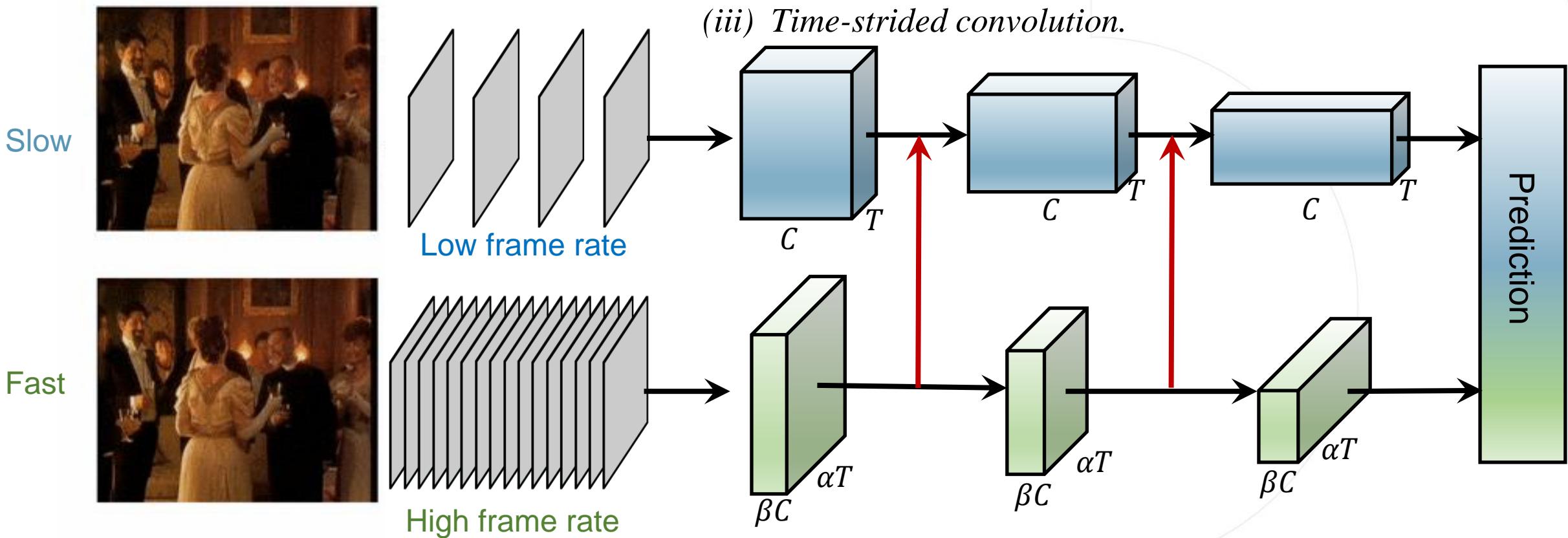
SlowFast Networks: treating time and space differently

Denoting the feature shape of *Slow pathway* as $\{T, S^2, C\}$, the feature shape of the *Fast pathway* is $\{\alpha T, S^2, \beta C\}$. E.g., $\alpha = 8, \beta = 1/8$.



SlowFast Networks: treating time and space differently **Lateral connections.**

- (i) Time-to-channel: Reshape and transpose.
- (ii) Time-strided sampling.
- (iii) Time-strided convolution.



SlowFast Networks

- Dimensions are $\{T \times S^2, C\}$
- Strides are {temporal, spatial}
- Residual blocks are shown by brackets
- Non-degenerate temporal filters are underlined
- 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	-	-	64×224^2
data layer	stride 16, 1^2	stride 2, 1^2	<i>Slow</i> : 4×224^2 <i>Fast</i> : 32×224^2
conv ₁	$1 \times 7^2, 64$ stride 1, 2^2	$5 \times 7^2, 8$ stride 1, 2^2	<i>Slow</i> : 4×112^2 <i>Fast</i> : 32×112^2
pool ₁	1×3^2 max stride 1, 2^2	1×3^2 max stride 1, 2^2	<i>Slow</i> : 4×56^2 <i>Fast</i> : 32×56^2
res ₂	$\left[\begin{array}{l} 1 \times 1^2, 64 \\ 1 \times 3^2, 64 \\ 1 \times 1^2, 256 \end{array} \right] \times 3$	$\left[\begin{array}{l} \frac{3 \times 1^2, 8}{1 \times 3^2, 8} \\ 1 \times 1^2, 32 \end{array} \right] \times 3$	<i>Slow</i> : 4×56^2 <i>Fast</i> : 32×56^2
res ₃	$\left[\begin{array}{l} 1 \times 1^2, 128 \\ 1 \times 3^2, 128 \\ 1 \times 1^2, 512 \end{array} \right] \times 4$	$\left[\begin{array}{l} \frac{3 \times 1^2, 16}{1 \times 3^2, 16} \\ 1 \times 1^2, 64 \end{array} \right] \times 4$	<i>Slow</i> : 4×28^2 <i>Fast</i> : 32×28^2
res ₄	$\left[\begin{array}{l} \frac{3 \times 1^2, 256}{1 \times 3^2, 256} \\ 1 \times 1^2, 1024 \end{array} \right] \times 6$	$\left[\begin{array}{l} \frac{3 \times 1^2, 32}{1 \times 3^2, 32} \\ 1 \times 1^2, 128 \end{array} \right] \times 6$	<i>Slow</i> : 4×14^2 <i>Fast</i> : 32×14^2
res ₅	$\left[\begin{array}{l} \frac{3 \times 1^2, 512}{1 \times 3^2, 512} \\ 1 \times 1^2, 2048 \end{array} \right] \times 3$	$\left[\begin{array}{l} \frac{3 \times 1^2, 64}{1 \times 3^2, 64} \\ 1 \times 1^2, 256 \end{array} \right] \times 3$	<i>Slow</i> : 4×7^2 <i>Fast</i> : 32×7^2
	global average pool, concat, fc		# classes

Biological studies in the primate visual system.

- **Parvocellular (P-cells):** 80%, provide fine spatial detail and color, but lower temporal resolution, responding slowly to stimuli.
- **Magnocellular (M-cells):** 15-20%, operate at *high temporal frequency*, responsive to fast temporal changes, but not sensitive to spatial detail or color.

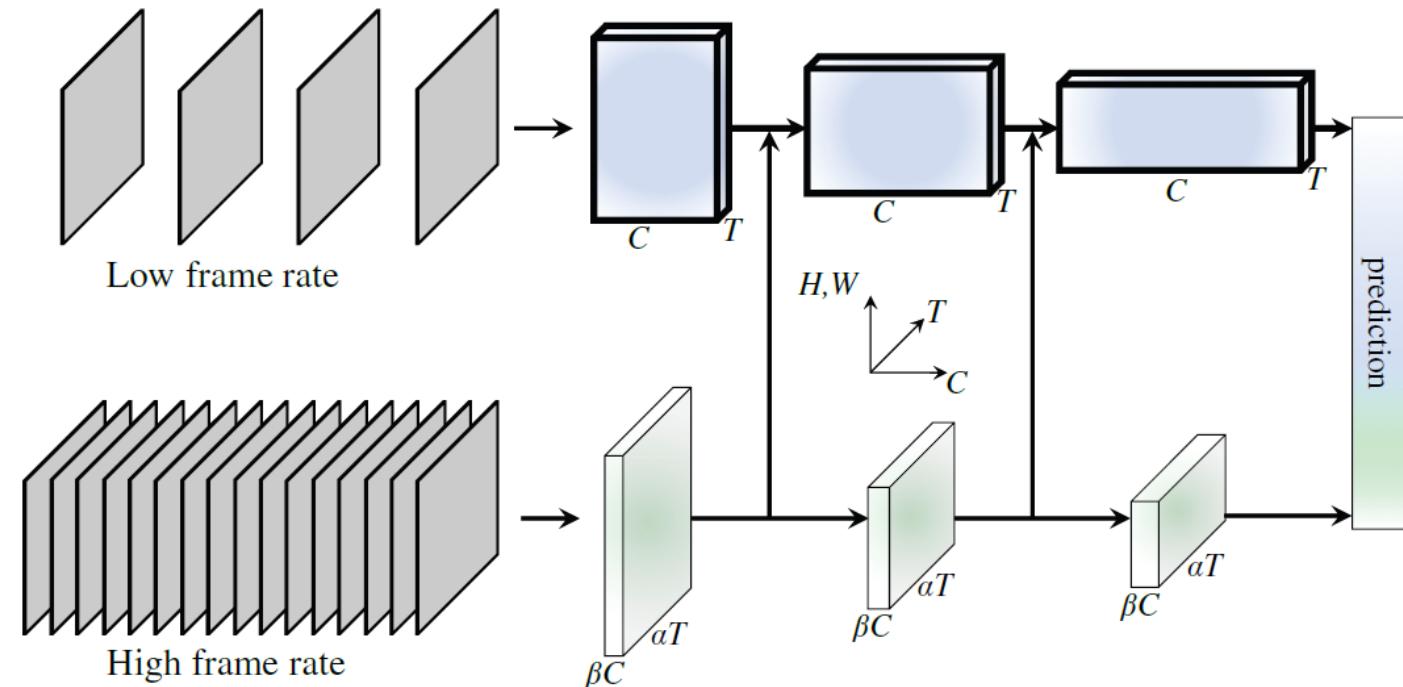


Figure 1. A **SlowFast network** has a low frame rate, low temporal resolution *Slow* pathway and a high frame rate, $\alpha \times$ higher temporal resolution *Fast* pathway. The Fast pathway is lightweight by using a fraction (β , e.g., 1/8) of channels. Lateral connections fuse them.

Video classification: performance comparison

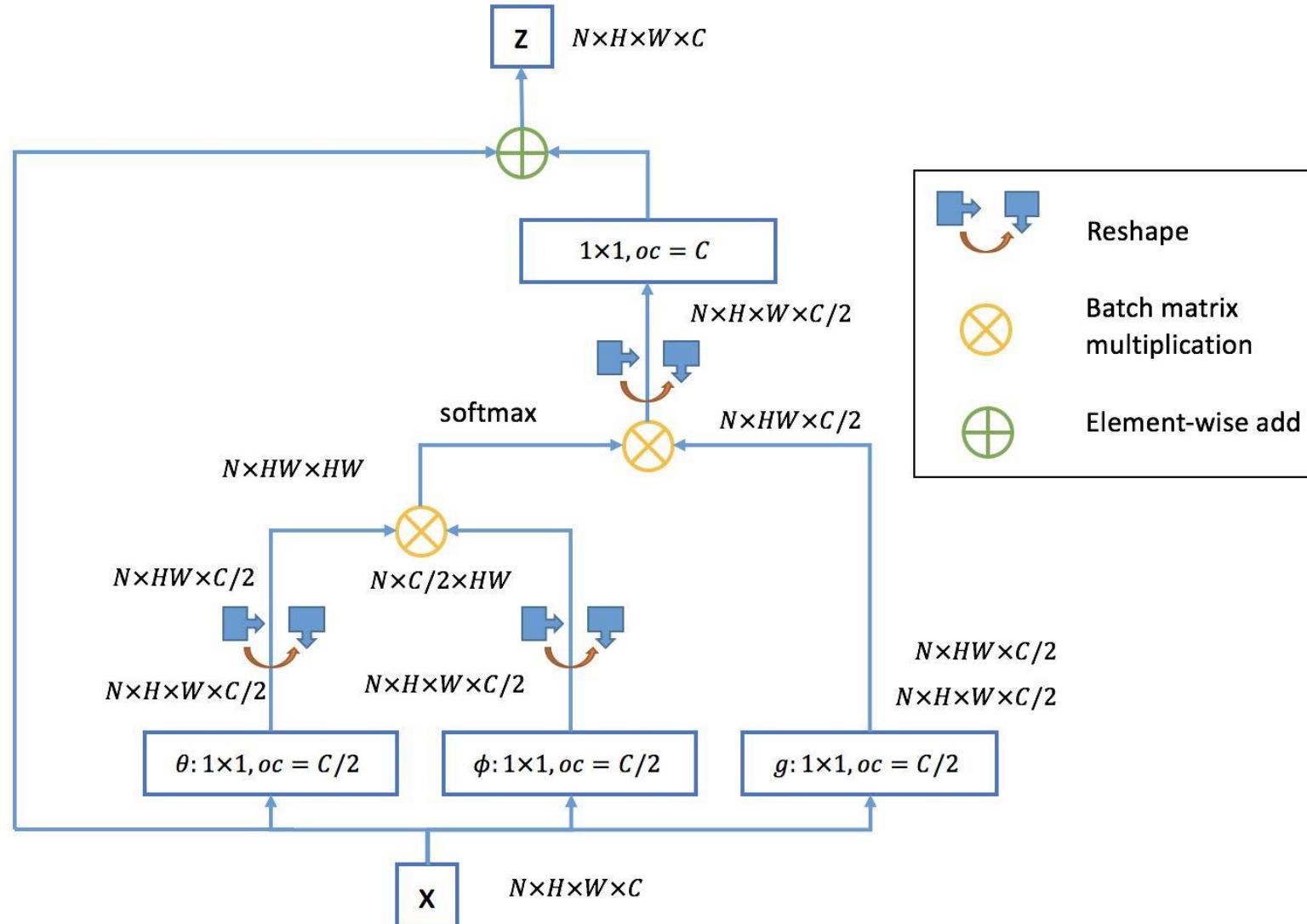
model	flow	pretrain	top-1	top-5	GFLOPs × views
I3D [5]		ImageNet	72.1	90.3	108 × N/A
Two-Stream I3D [5]	✓	ImageNet	75.7	92.0	216 × N/A
S3D-G [61]	✓	ImageNet	77.2	93.0	143 × N/A
Nonlocal R50 [56]		ImageNet	76.5	92.6	282 × 30
Nonlocal R101 [56]		ImageNet	77.7	93.3	359 × 30
R(2+1)D Flow [50]	✓	-	67.5	87.2	152 × 115
STC [9]		-	68.7	88.5	N/A × N/A
ARTNet [54]		-	69.2	88.3	23.5 × 250
S3D [61]		-	69.4	89.1	66.4 × N/A
ECO [63]		-	70.0	89.4	N/A × N/A
I3D [5]	✓	-	71.6	90.0	216 × N/A
R(2+1)D [50]	✓	-	72.0	90.0	152 × 115
R(2+1)D [50]	✓	-	73.9	90.9	304 × 115
SlowFast 4×16, R50		-	75.6	92.1	36.1 × 30
SlowFast 8×8, R50		-	77.0	92.6	65.7 × 30
SlowFast 8×8, R101		-	77.9	93.2	106 × 30
SlowFast 16×8, R101		-	78.9	93.5	213 × 30
SlowFast 16×8, R101+N		-	79.8	93.9	234 × 30

Comparison on Kinetics-400. In the last column, we report the inference cost with a single “view” (temporal clip with spatial crop) × the numbers of such views used. The SlowFast models are with different input sampling ($T \times$) and backbones (R-50, R-101, NL). “N/A” indicates the numbers are not available for us.

Spatio-Temporal Self-Attention (Nonlocal Block)



Spatio-Temporal Self-Attention (Nonlocal Block)

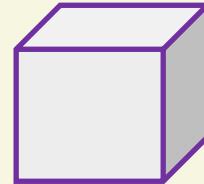


Video classification: Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



3D CNN



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

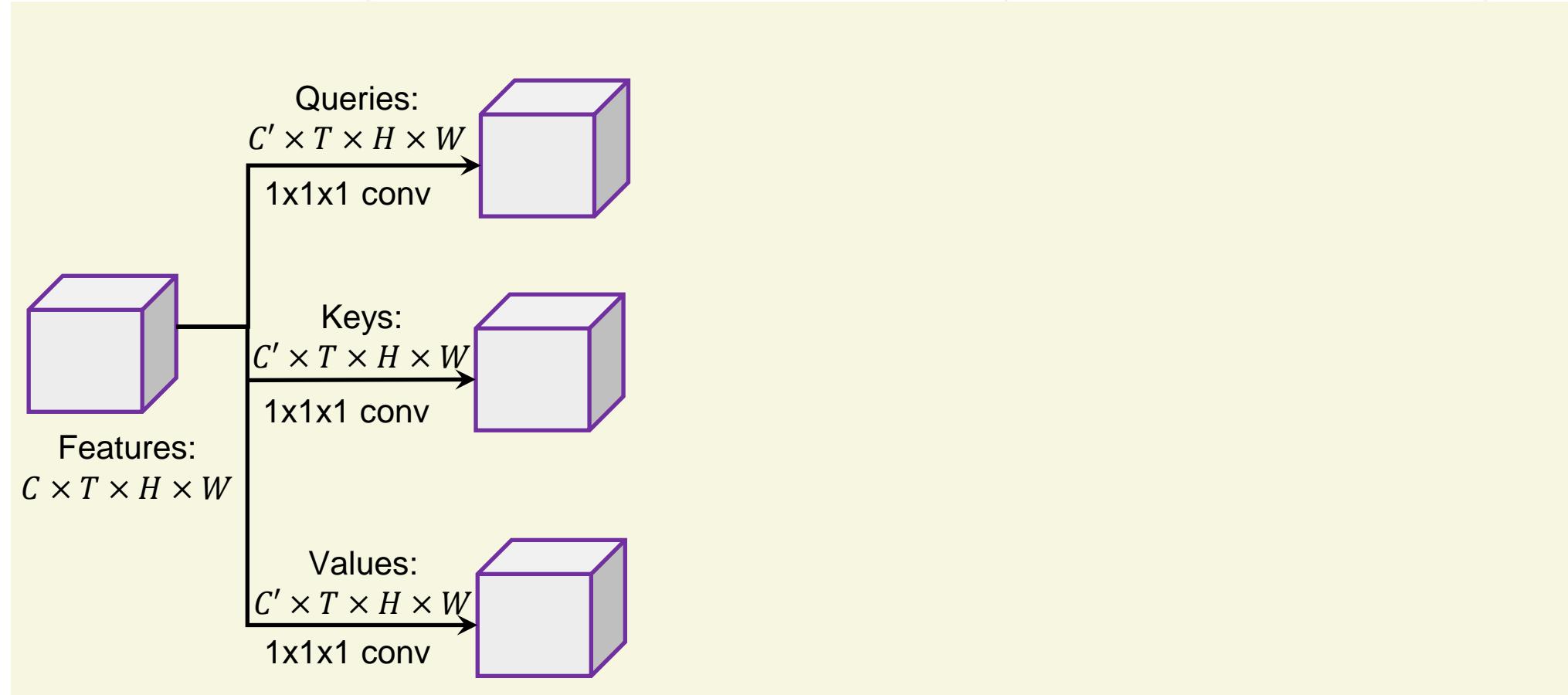
Nonlocal Block

Video classification: Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



3D CNN



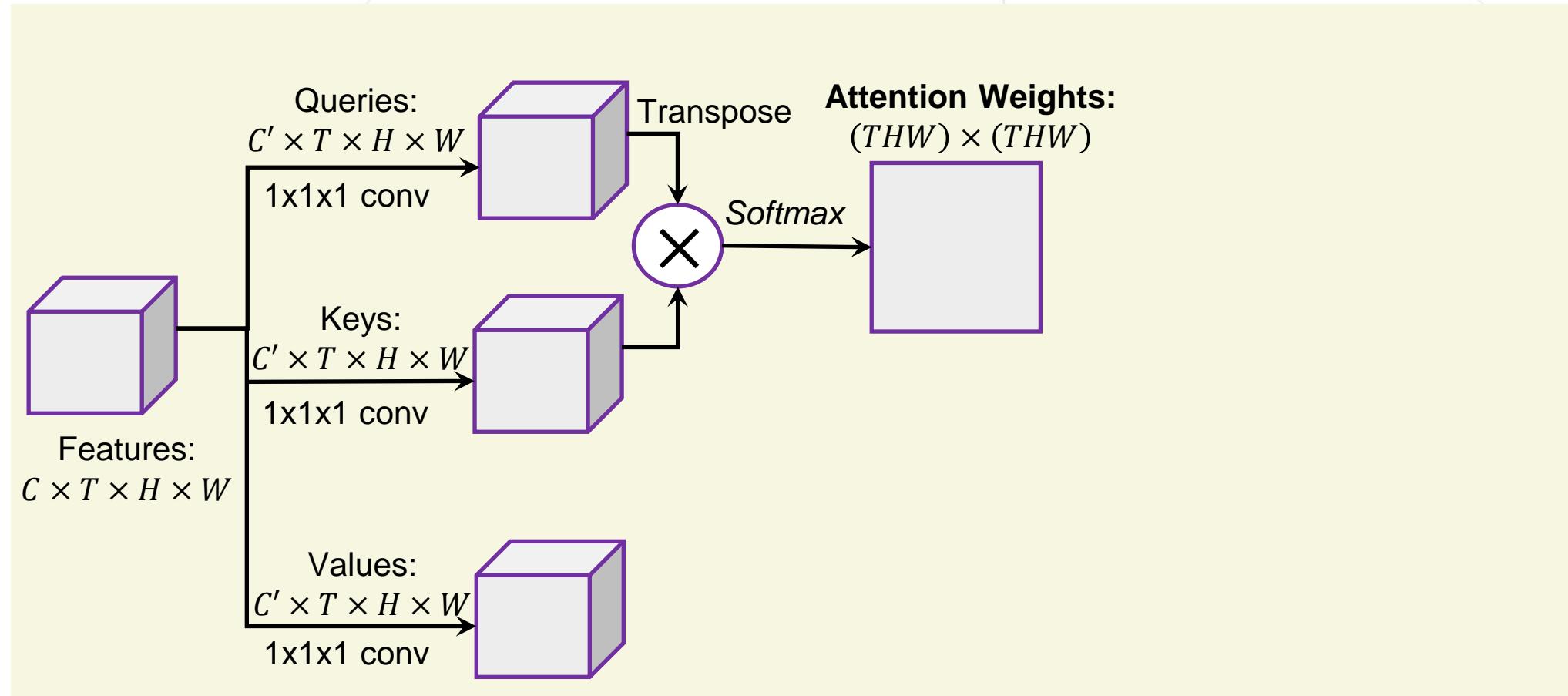
Nonlocal Block

Video classification: Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



3D CNN



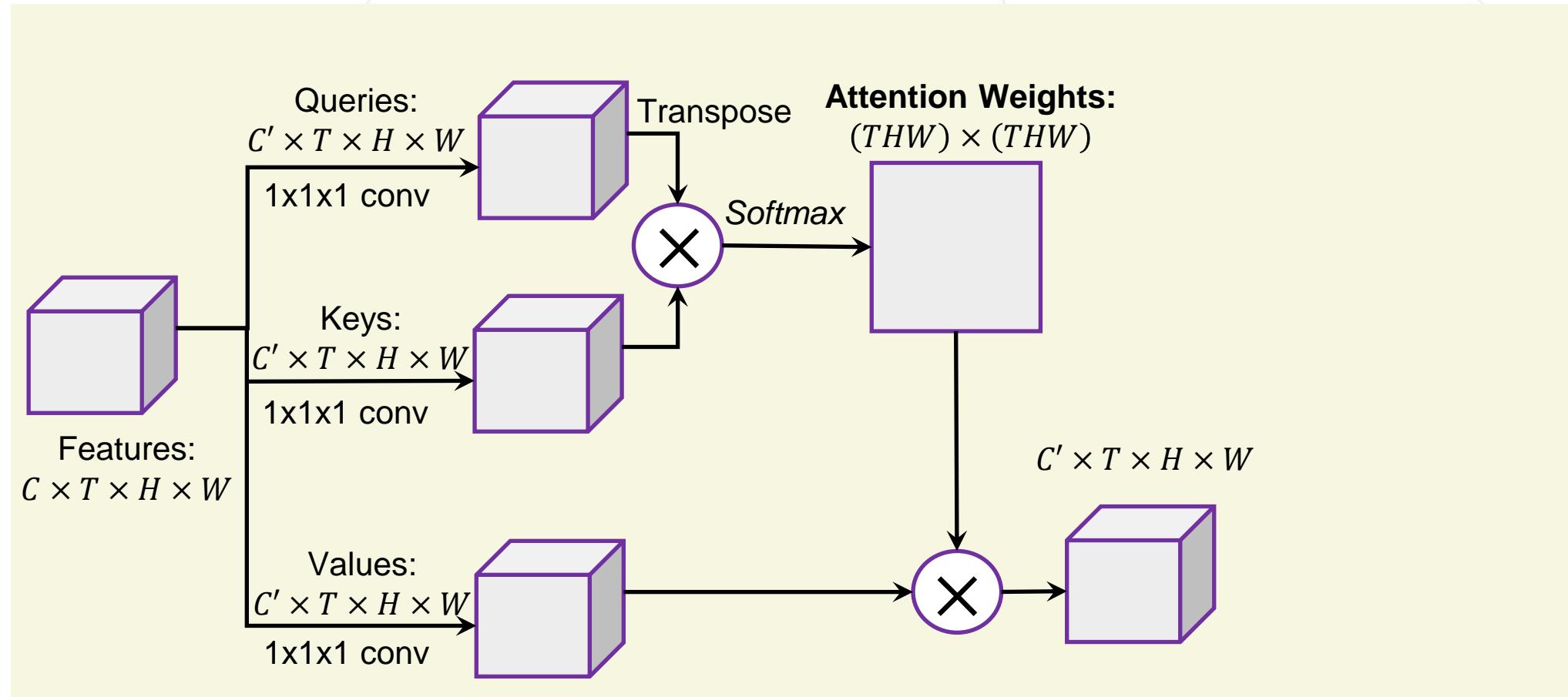
Nonlocal Block

Video classification: Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



3D CNN



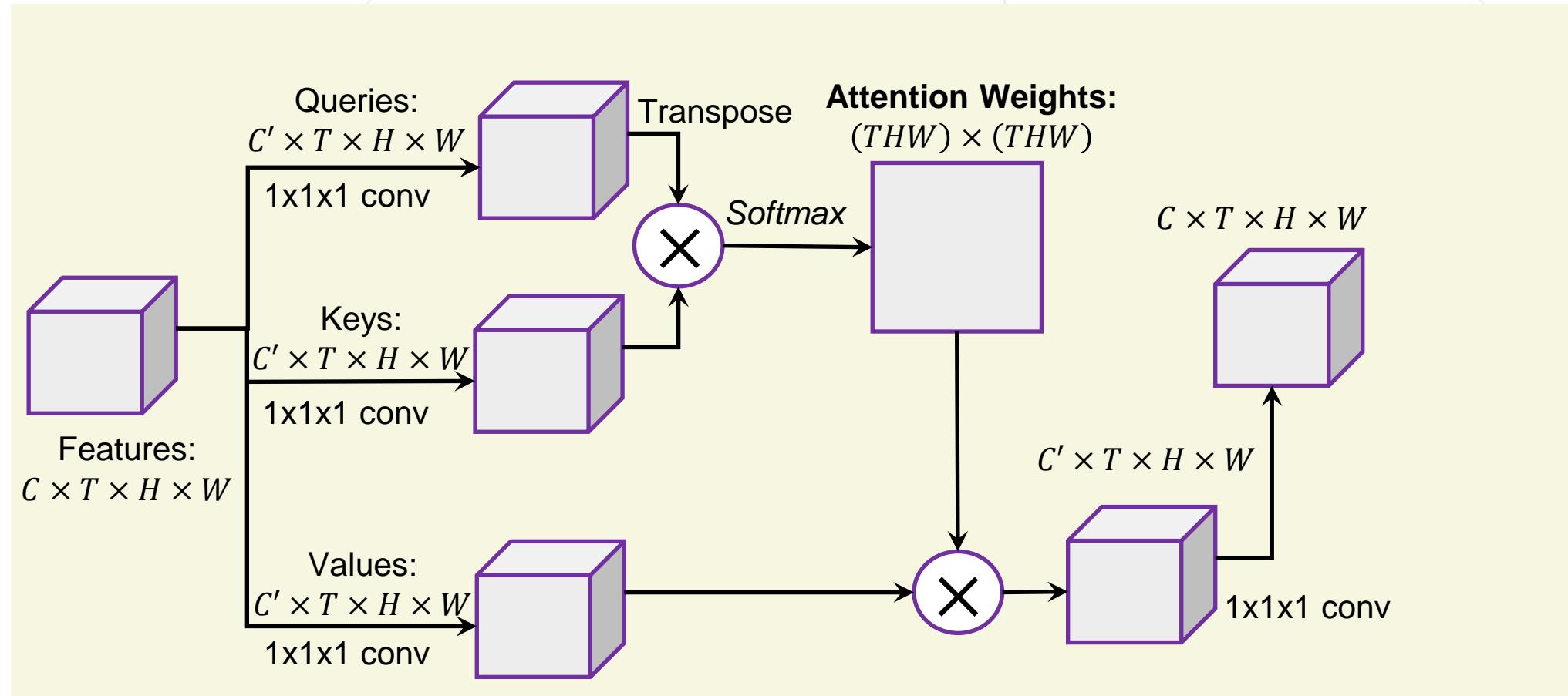
Nonlocal Block

Video classification: Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



3D CNN



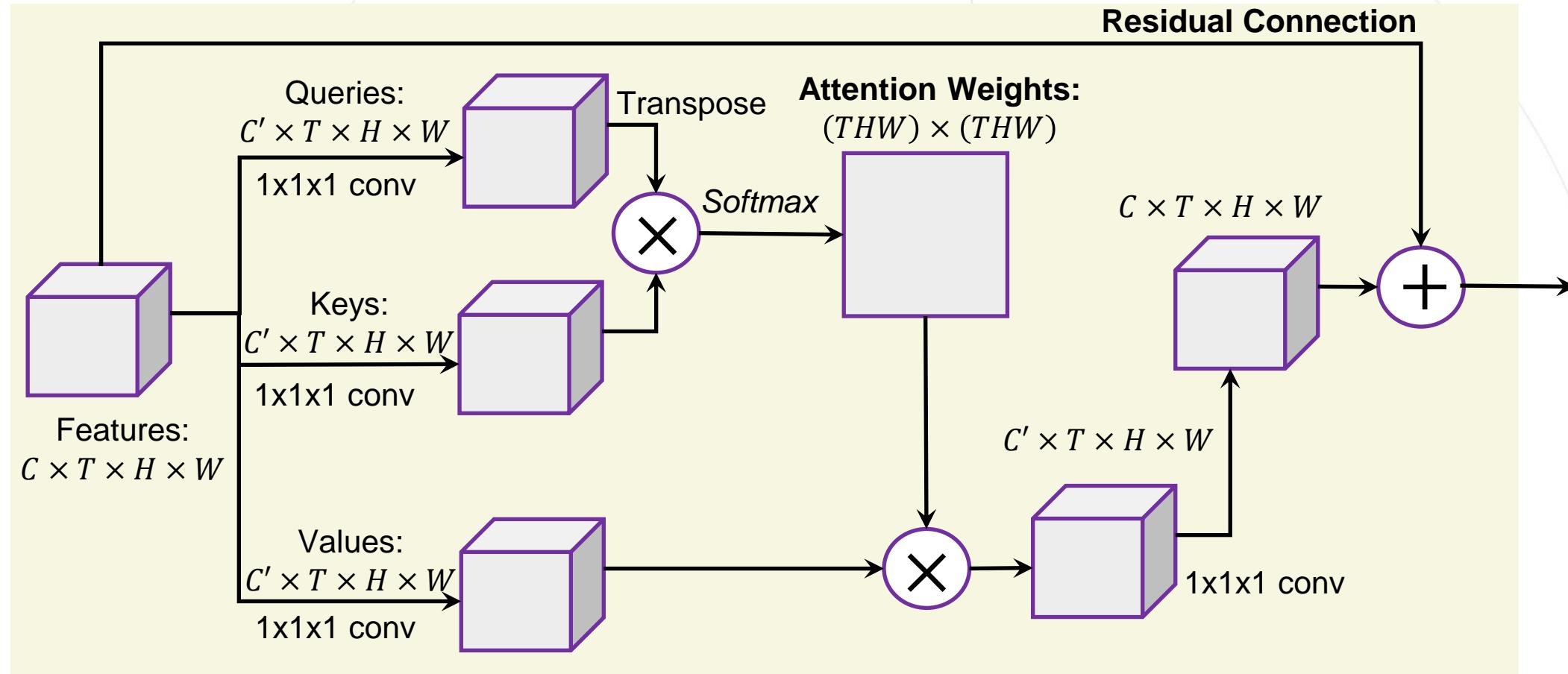
Nonlocal Block

Video classification: Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip



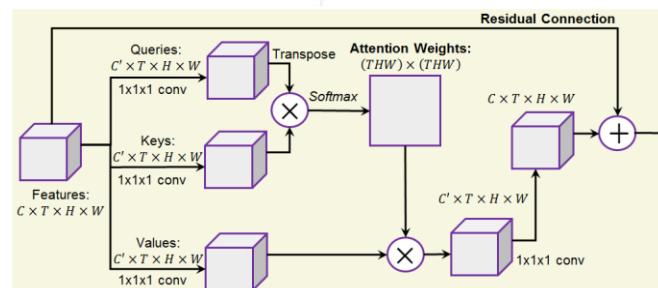
3D CNN



Nonlocal Block

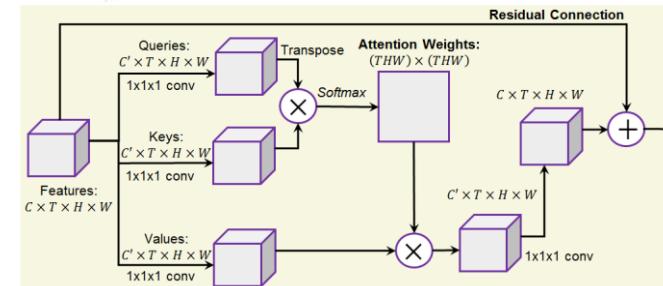
Video classification: Spatio-Temporal Self-Attention (Nonlocal Block)

Input clip

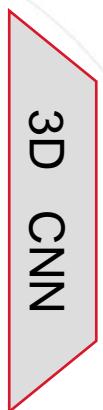


Nonlocal Block

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



Nonlocal Block

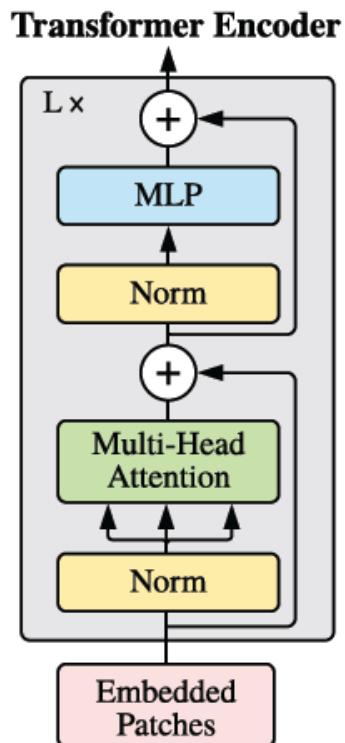
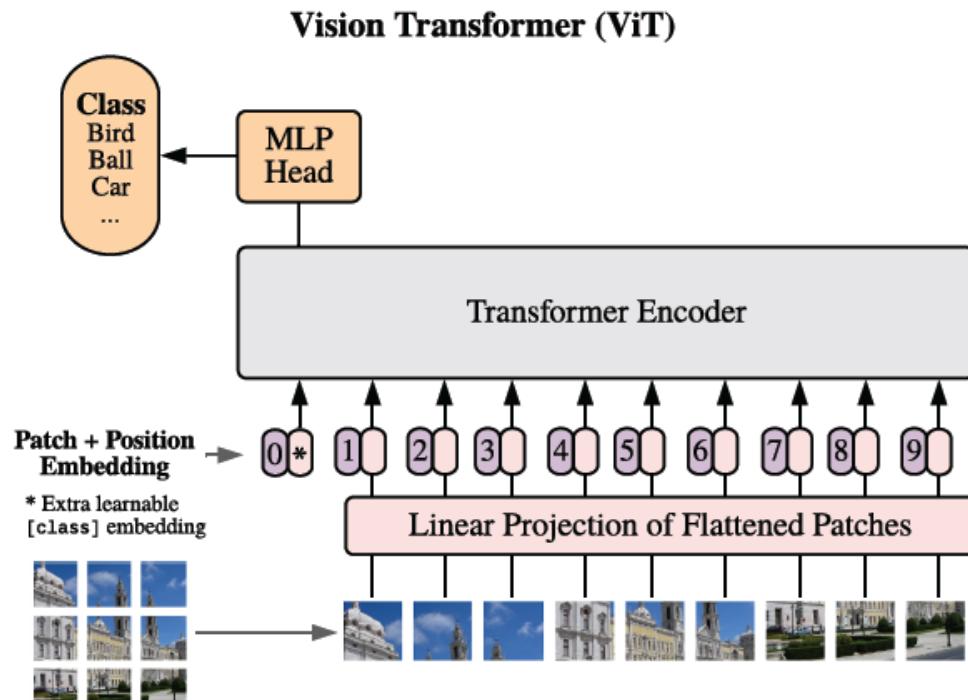
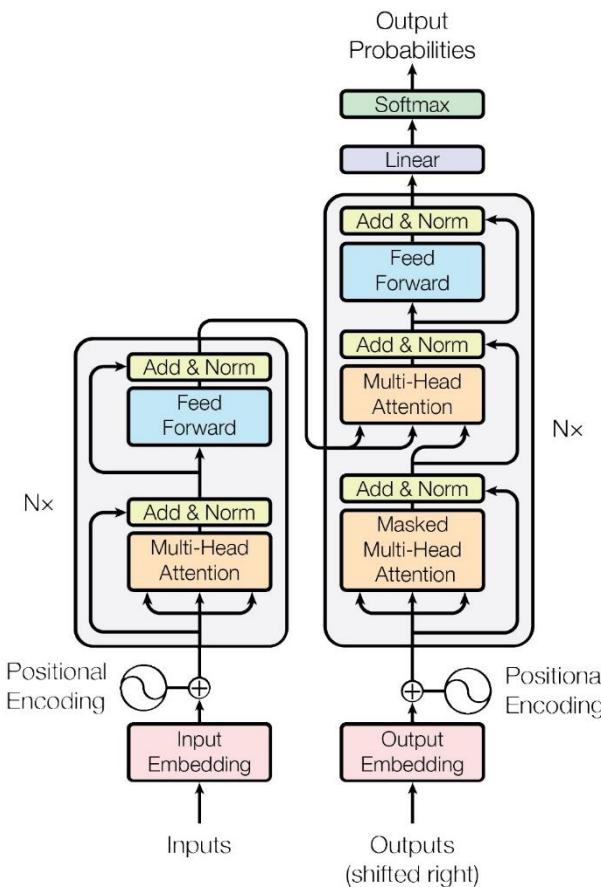


Spatio-Temporal Self-Attention (Nonlocal Block)

model	backbone	modality	top-1 val	top-5 val	top-1 test	top-5 test	avg test [†]
I3D in [7]	Inception	RGB	72.1	90.3	71.1	89.3	80.2
2-Stream I3D in [7]	Inception	RGB + flow	75.7	92.0	74.2	91.3	82.8
RGB baseline in [3]	Inception-ResNet-v2	RGB	73.0	90.9	-	-	-
3-stream late fusion [3]	Inception-ResNet-v2	RGB + flow + audio	74.9	91.6	-	-	-
3-stream LSTM [3]	Inception-ResNet-v2	RGB + flow + audio	77.1	93.2	-	-	-
3-stream SATT [3]	Inception-ResNet-v2	RGB + flow + audio	77.7	93.2	-	-	-
NL I3D [ours]	ResNet-50	RGB	76.5	92.6	-	-	-
	ResNet-101	RGB	77.7	93.3	-	-	83.8

Table 3. Comparisons with state-of-the-art results in **Kinetics**, reported on the val and test sets. We include the Kinetics 2017 competition winner's results [3], but their best results exploited audio signals (marked in gray) so were not vision-only solutions. [†]: “avg” is the average of top-1 and top-5 accuracy; individual top-1 or top-5 numbers are not available from the test server at the time of submitting this manuscript.

Vision Transformer : nlp → cv



Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In ICLR, 2021.

ViViT

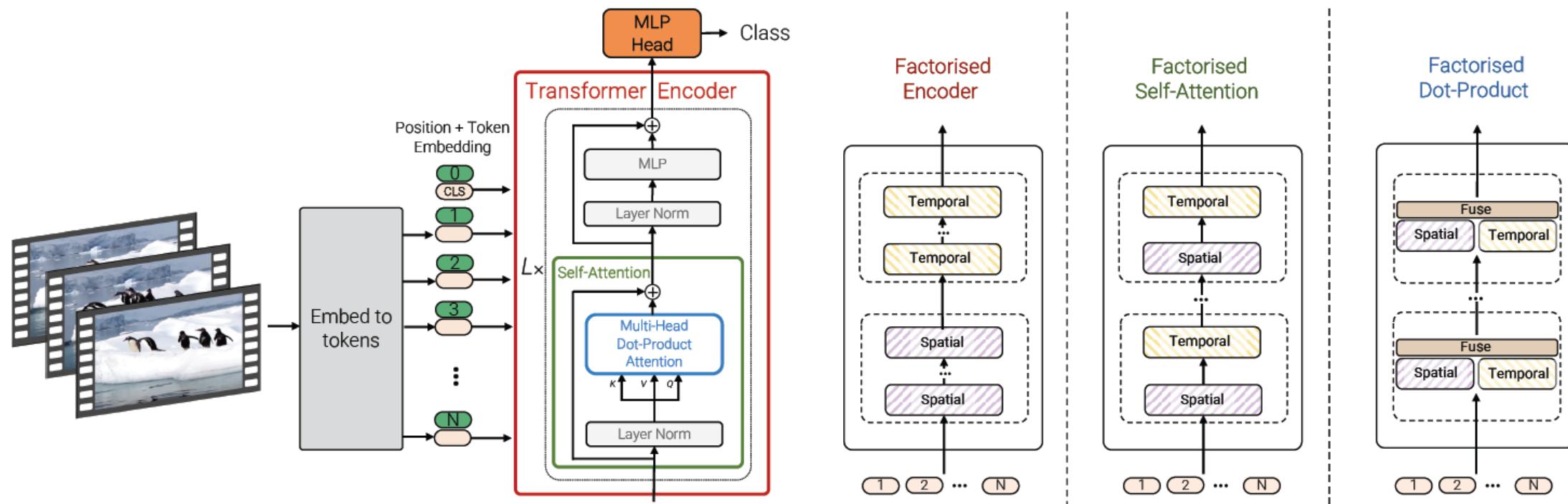


Figure 1: We propose a pure-transformer architecture for video classification, inspired by the recent success of such models for images [17]. To effectively process a large number of spatio-temporal tokens, we develop several model variants which factorise different components of the transformer encoder over the spatial- and temporal-dimensions. As shown on the right, these factorisations correspond to different attention patterns over space and time.

ViViT

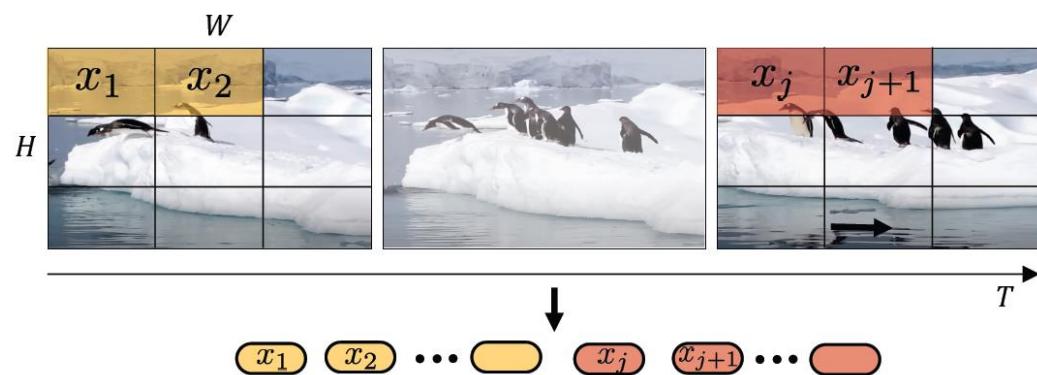


Figure 2: Uniform frame sampling: We simply sample n_t frames, and embed each 2D frame independently following ViT [17].

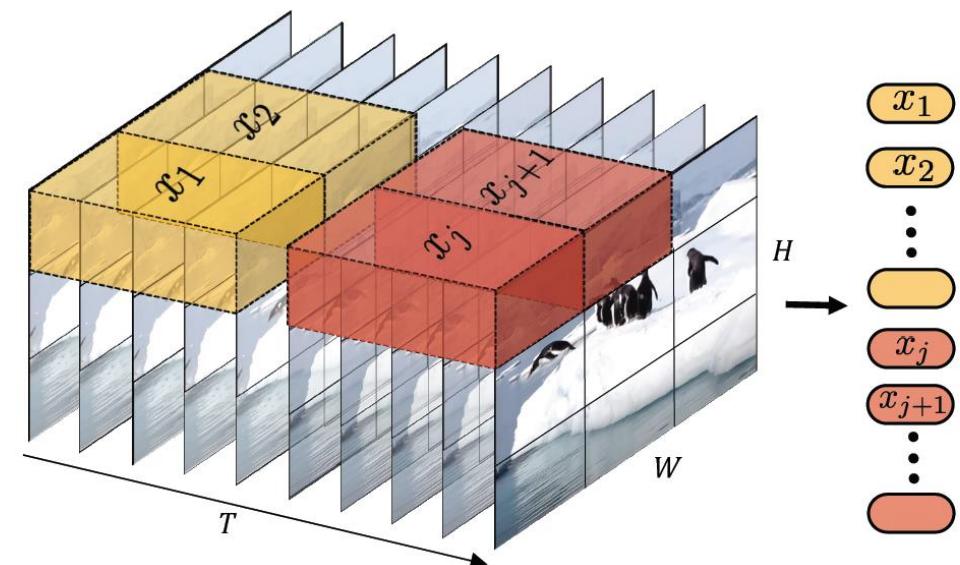


Figure 3: Tubelet embedding. We extract and linearly embed non-overlapping tubelets that span the spatio-temporal input volume.

(a) Kinetics 400

Method	Top 1	Top 5	Views	TFLOPs
blVNet [18]	73.5	91.2	–	–
STM [32]	73.7	91.6	–	–
TEA [41]	76.1	92.5	10 × 3	2.10
TSM-ResNeXt-101 [42]	76.3	–	–	–
I3D NL [74]	77.7	93.3	10 × 3	10.77
CorrNet-101 [69]	79.2	–	10 × 3	6.72
ip-CSN-152 [65]	79.2	93.8	10 × 3	3.27
LGD-3D R101 [50]	79.4	94.4	–	–
SlowFast R101-NL [20]	79.8	93.9	10 × 3	7.02
X3D-XXL [19]	80.4	94.6	10 × 3	5.82
TimeSformer-L [4]	80.7	94.7	1 × 3	7.14
ViViT-L/16x2 FE	80.6	92.7	1 × 1	3.98
ViViT-L/16x2 FE	81.7	93.8	1 × 3	11.94

Methods with large-scale pretraining

ip-CSN-152 [65] (IG [43])	82.5	95.3	10 × 3	3.27
ViViT-L/16x2 FE (JFT)	83.5	94.3	1 × 3	11.94
ViViT-H/16x2 (JFT)	84.9	95.8	4 × 3	47.77

(b) Kinetics 600

Method	Top 1	Top 5
AttentionNAS [75]	79.8	94.4
LGD-3D R101 [50]	81.5	95.6
SlowFast R101-NL [20]	81.8	95.1
X3D-XL [19]	81.9	95.5
TimeSformer-L [4]	82.2	95.6
ViViT-L/16x2 FE	82.9	94.6
ViViT-L/16x2 FE (JFT)	84.3	94.9
ViViT-H/16x2 (JFT)	85.8	96.5

X-ViT

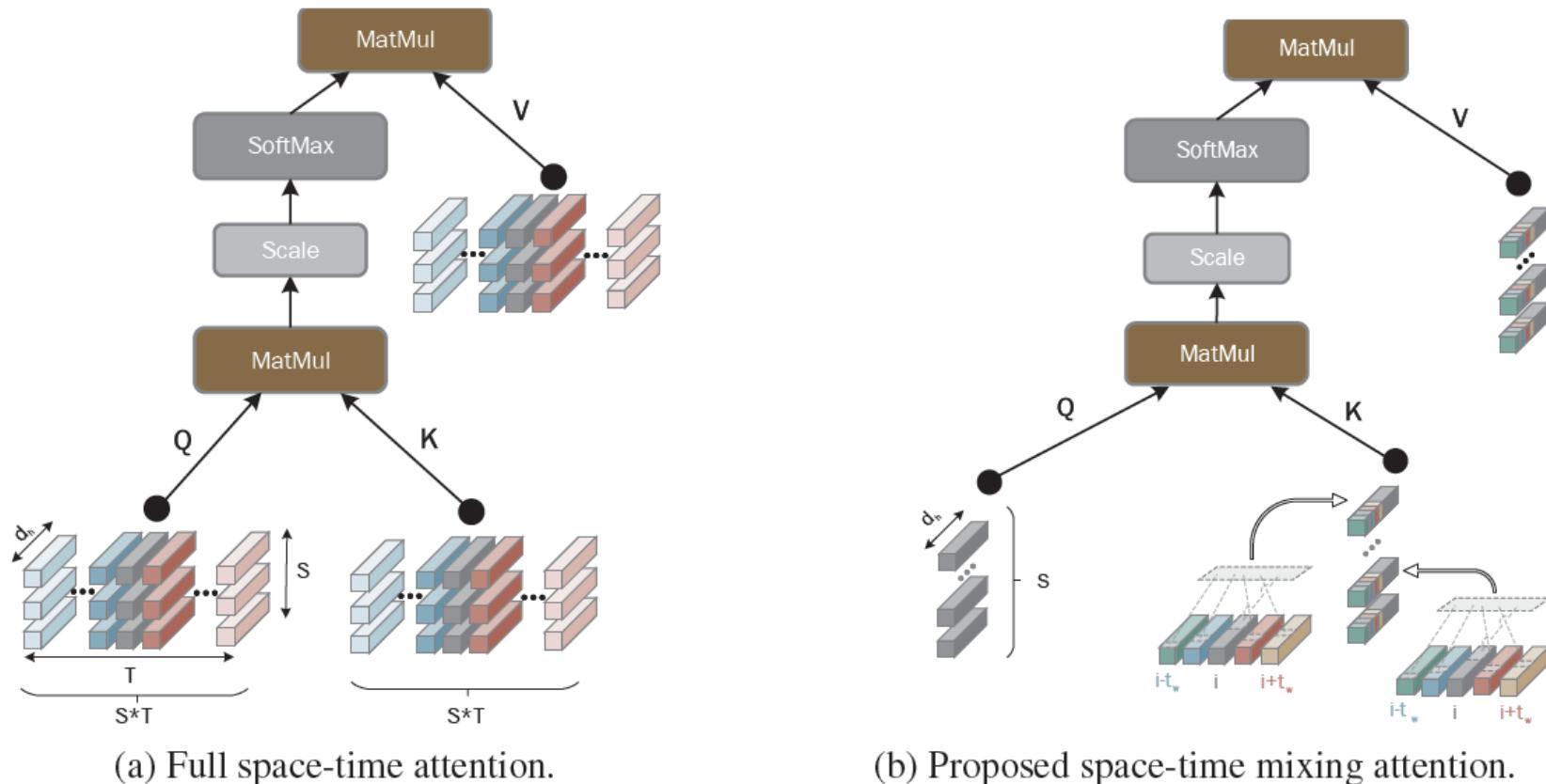


Figure 2: Detailed self-attention computation graph for (a) full space-time attention and (b) the proposed space-time mixing approximation. Notice that in our case only S tokens participate instead of ST . The temporal information is aggregated by indexing channels from adjacent frames. Tokens of identical colors share the same temporal index.

X-ViT

Table 5: Comparison with state-of-the-art on the Kinetics-400.

Method	Top-1	Top-5	# Frames	Views	Params	FLOPs ($\times 10^9$)
bLVNet [14]	73.5	91.2	24 × 2	3 × 3	25M	840
STM [19]	73.7	91.6	16	-	24M	-
TEA [25]	76.1	92.5	16	10 × 3	25.6M	2,100
TSM R50 [26]	74.7	-	16	10 × 3	25.6M	650
I3D NL [44]	77.7	93.3	128	10 × 3	-	10,800
CorrNet-101 [40]	79.2	-	32	10 × 3	-	6,700
ip-CSN-152 [38]	79.2	93.8	8	10 × 3	-	3,270
LGD-3D R101 [31]	79.4	94.4	16	-	-	-
SlowFast 8×8 R101+NL [16]	78.7	93.5	8	10 × 3	-	3,480
SlowFast 16×8 R101+NL [16]	79.8	93.9	16	10 × 3	-	7,020
X3D-XXL [15]	80.4	94.6	-	10 × 3	20.3M	5,823
TimeSformer-L [3]	80.7	94.7	96	1 × 3	121M	7,140
ViViT-L/16x2 [1]	80.6	94.7	32	4 × 3	312M	17,352
X-ViT (Ours)	78.5	93.7	8	1 × 3	92M	425
X-ViT (Ours)	79.4	93.9	8	2 × 3	92M	850
X-ViT (Ours)	80.2	94.7	16	1 × 3	92M	850
X-ViT (Ours)	80.7	94.7	16	2 × 3	92M	1700

Outline

Part 1	Video classification	P01-P08
Part 2	Datasets	P09-P18
Part 3	Solutions	P19-P98
Part 4	Summary & Future Topics	P99-P100

Video classification / action recognition

- Single-frame CNN
- Late fusion
- Early fusion
- 3D CNN
- CNN + RNN
- Two-stream networks
- Self attention