

# Action Recognition

Computer Vision (SJK02)

Universitat Jaume I

# Examples (demos)



[Pose estimation and action recognition](#)



[EPIC-Fusion: Audio-Visual Temporal Binding for Egocentric Action Recognition](#)

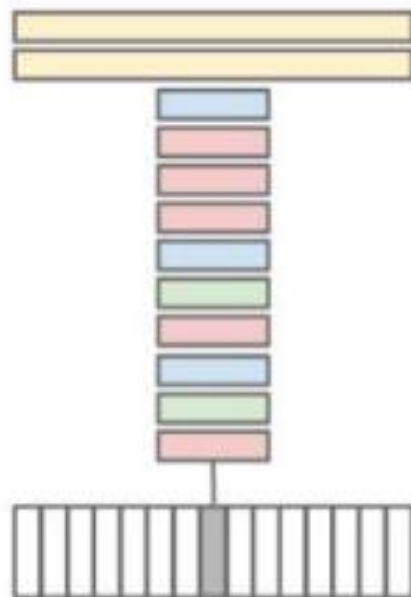
# Choice #1: CNNs

[Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset](#) (CVPR 2017)

[Large-scale video classification with convolutional neural networks](#) (CVPR 2014)

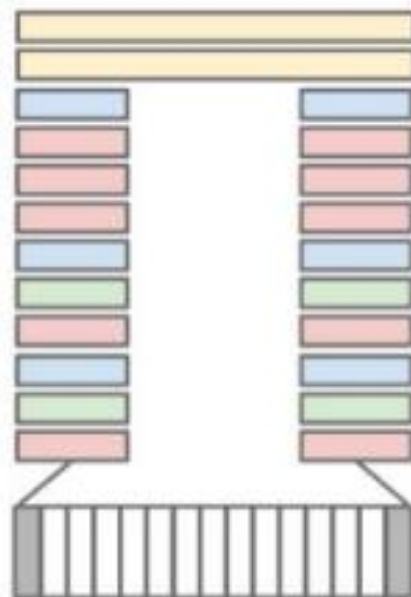
# What about the temporal information?

Single Frame



$W \times H \times 3$

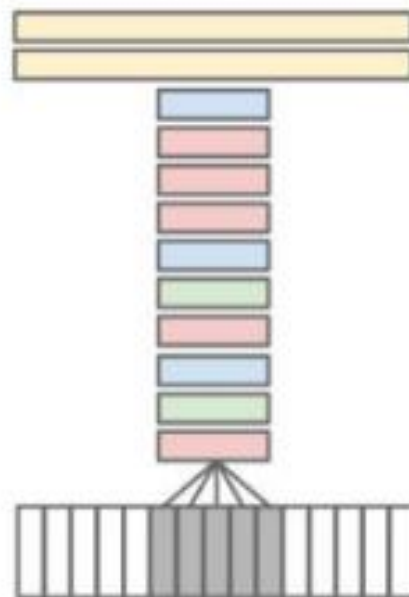
Late Fusion



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

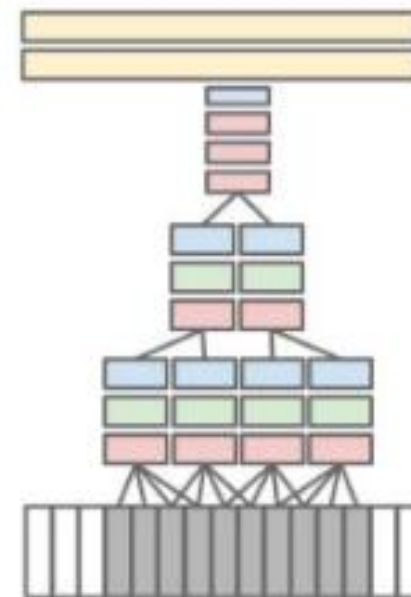
Shared params

Early Fusion

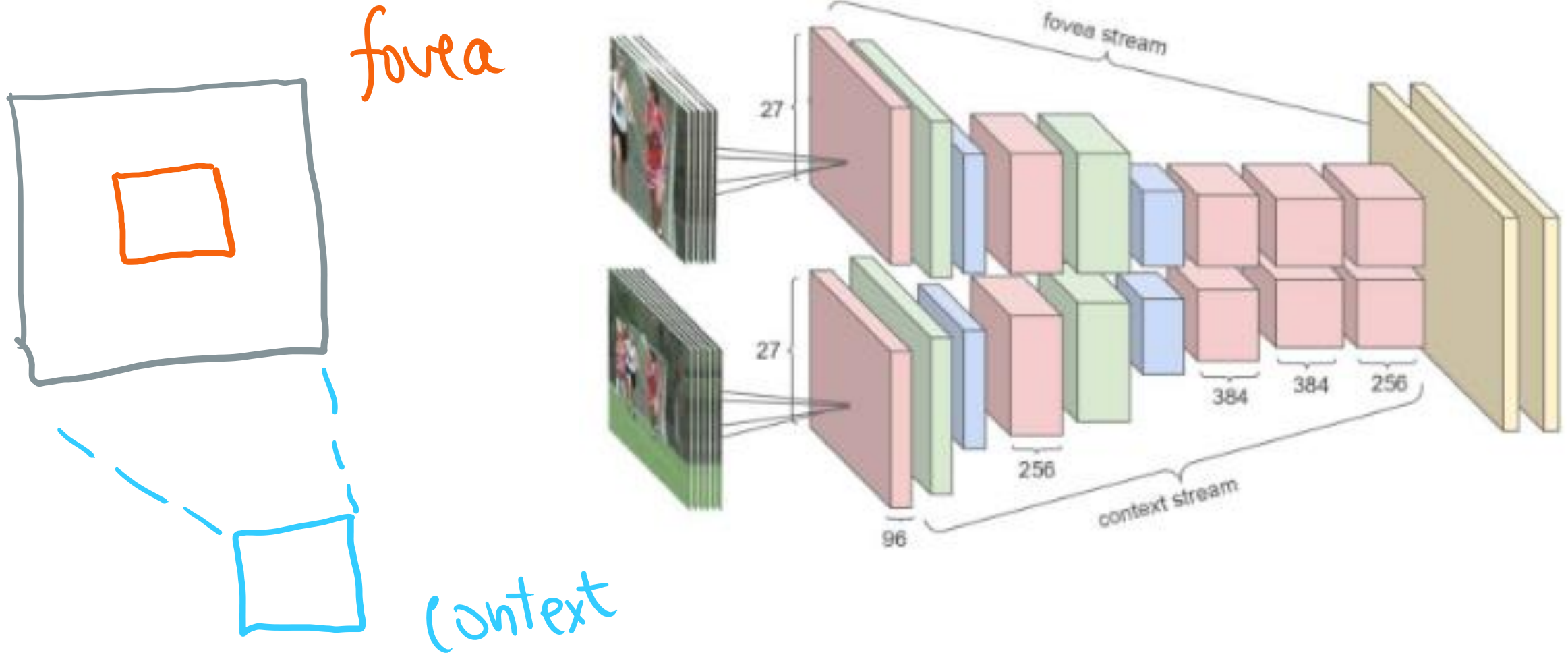


$W \times h \times 3 \times T$

Slow Fusion



# Multiresolution: fovea + context

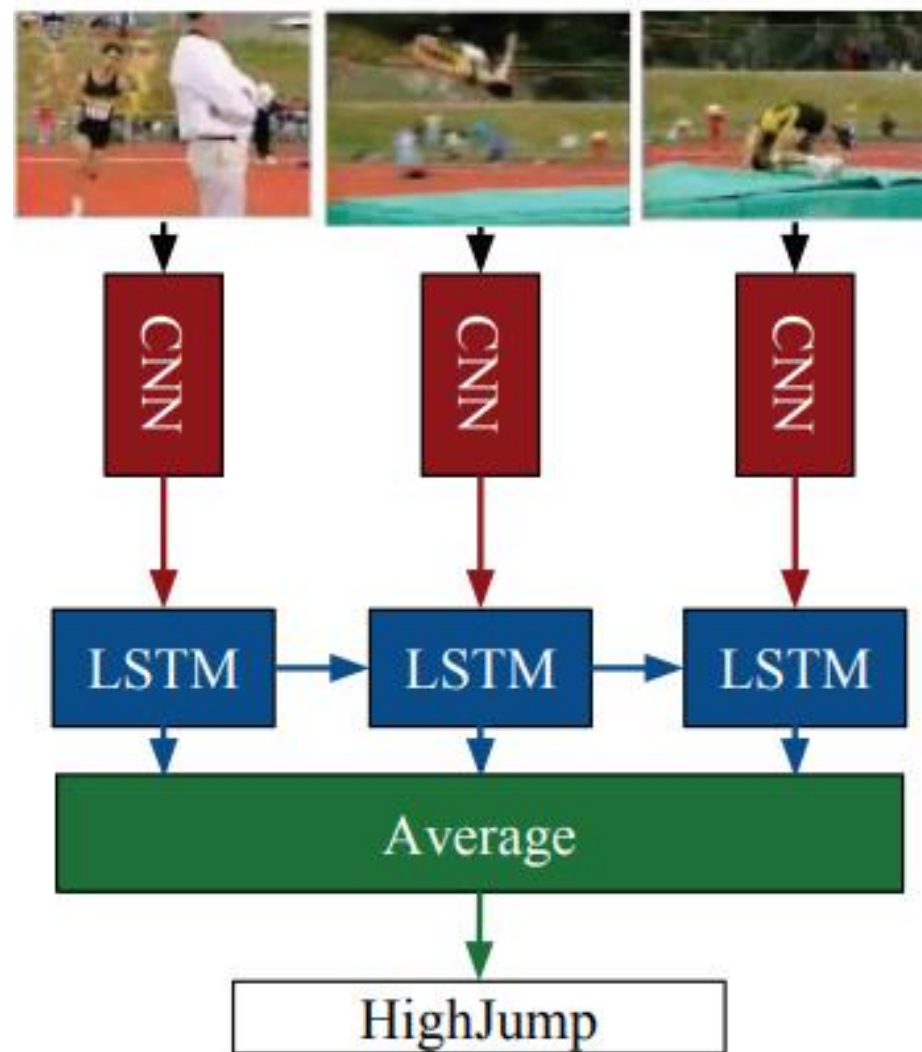
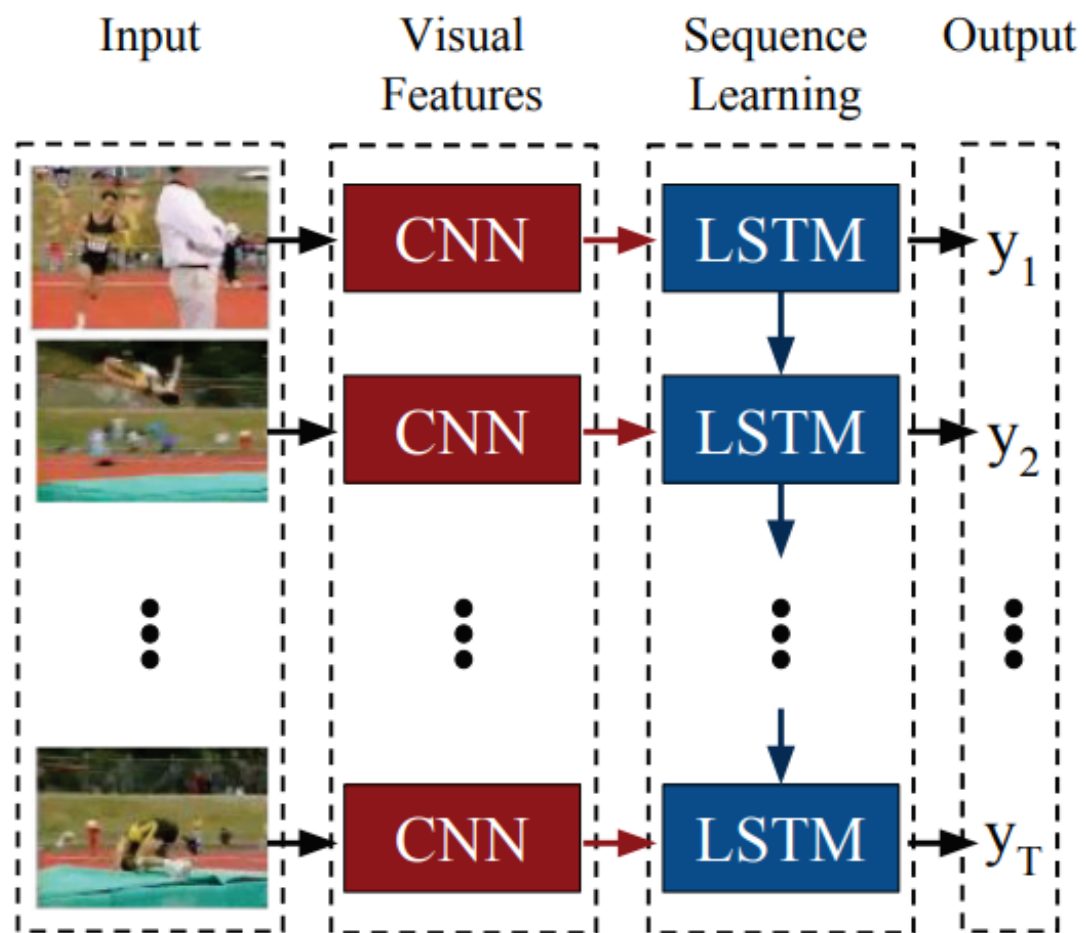


# Issue

Generally, temporal structure is ignored (or not well captured)

## **Choice #2: CNNs + LSTMs**

Long-term recurrent convolutional networks for visual recognition and description [[CVPR 2015](#), [PAMI 2017](#)]





# Benefit & Issue

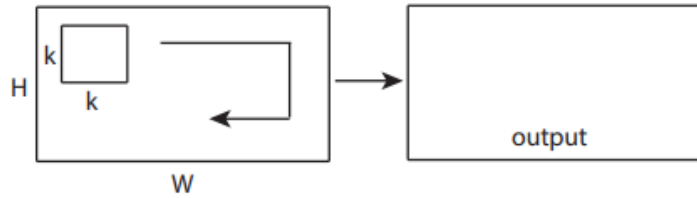
Can model order and long-range dependencies

May miss (potentially critical) low-level motions

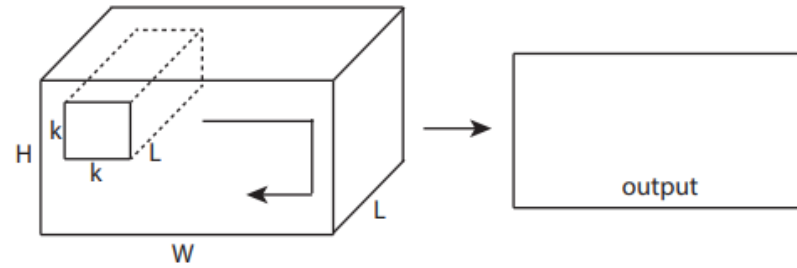
## Choice #3: 3D CNNs

[Learning spatiotemporal features with 3d convolutional networks](#) [ICCV 2015]  
[3D Convolutional Neural Networks for Human Action Recognition](#) [PAMI 2013]

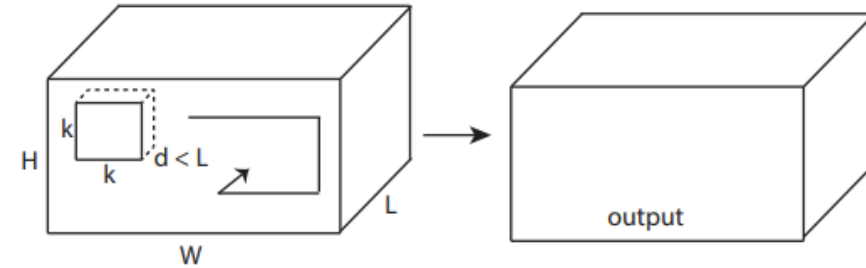
# 2D and 3D convolutions



(a) 2D convolution



(b) 2D convolution on multiple frames

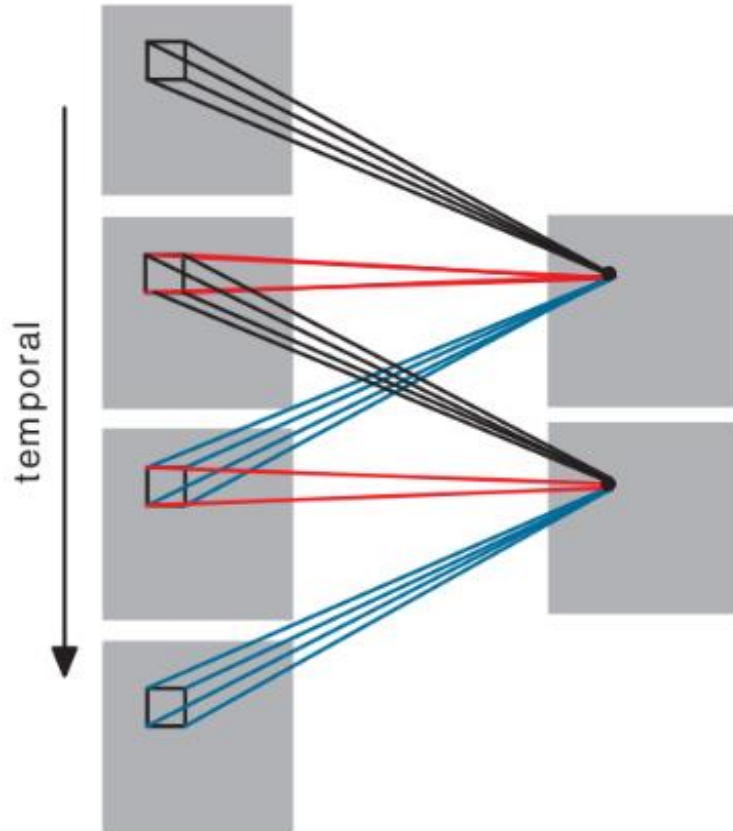


(c) 3D convolution

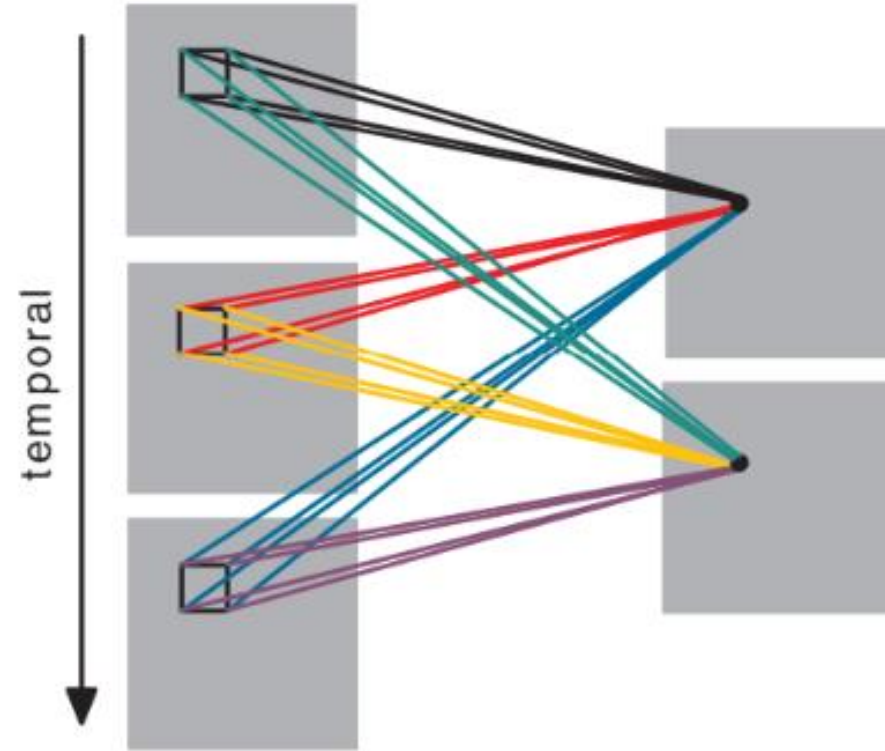
# 2D and 3D convolutions



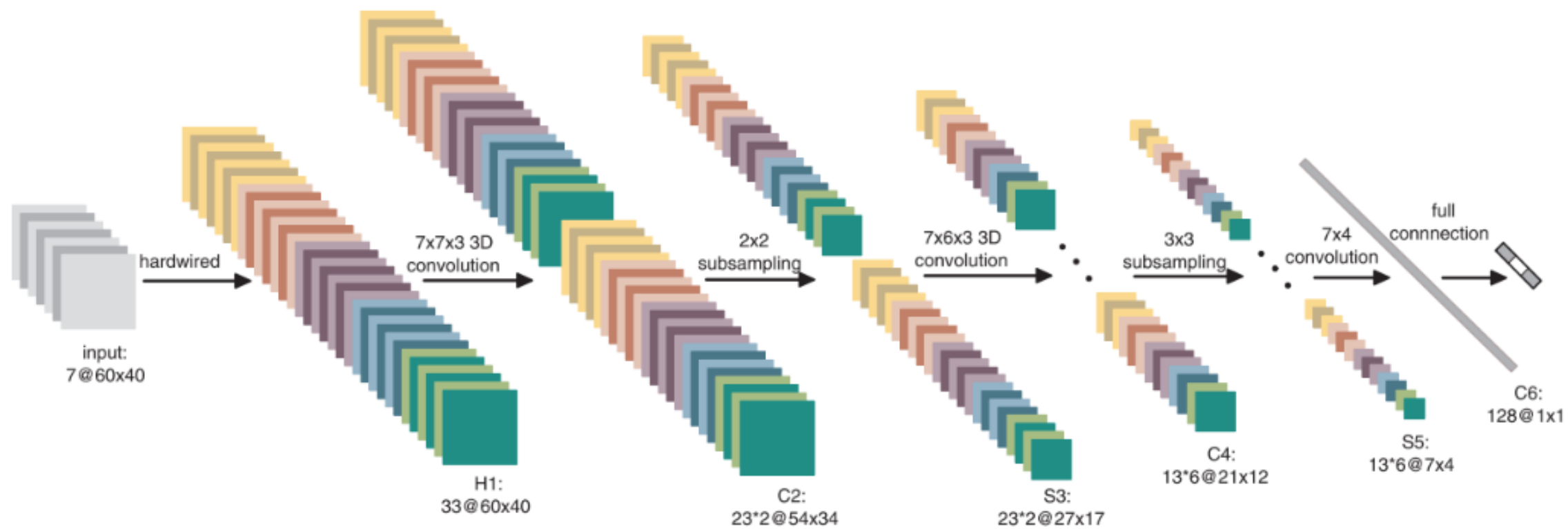
(a) 2D convolution



(b) 3D convolution



Multiple 3D convolutions  
for multiple features



# Benefit & Issue

Hierarchical representation of spatio-temporal information

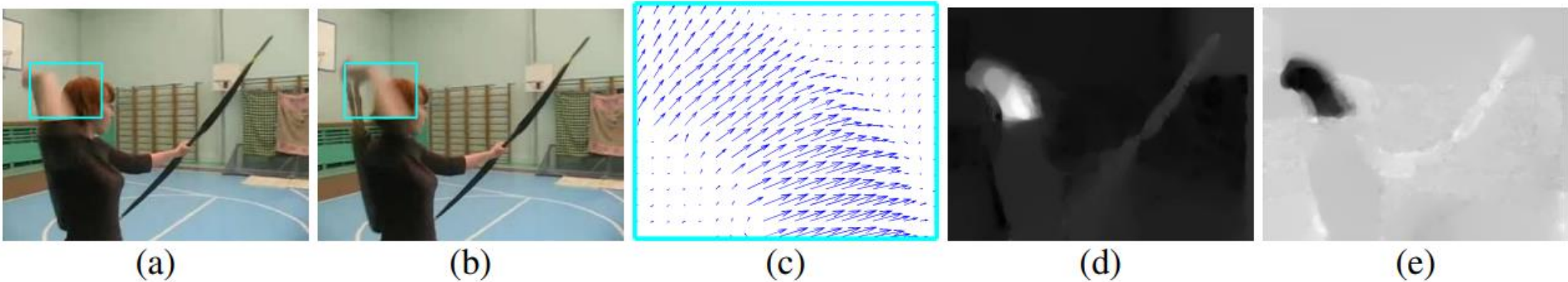
Many more parameters than 2D CNNs

## Choice #4: Two-stream nets

[Two-stream convolutional networks for action recognition in videos](#) (NIPS 2014)

[Convolutional Two-Stream Network Fusion for Video Action Recognition](#) (CVPR 2016)

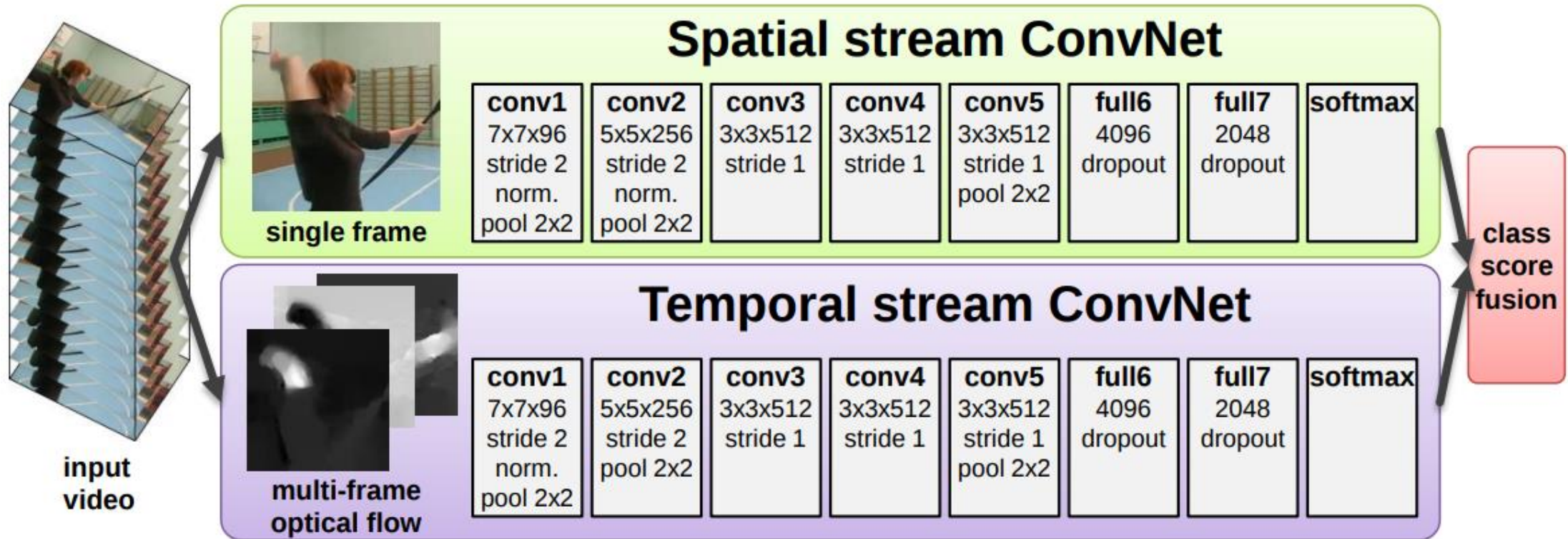
# Optic flow (precomputed, dense)



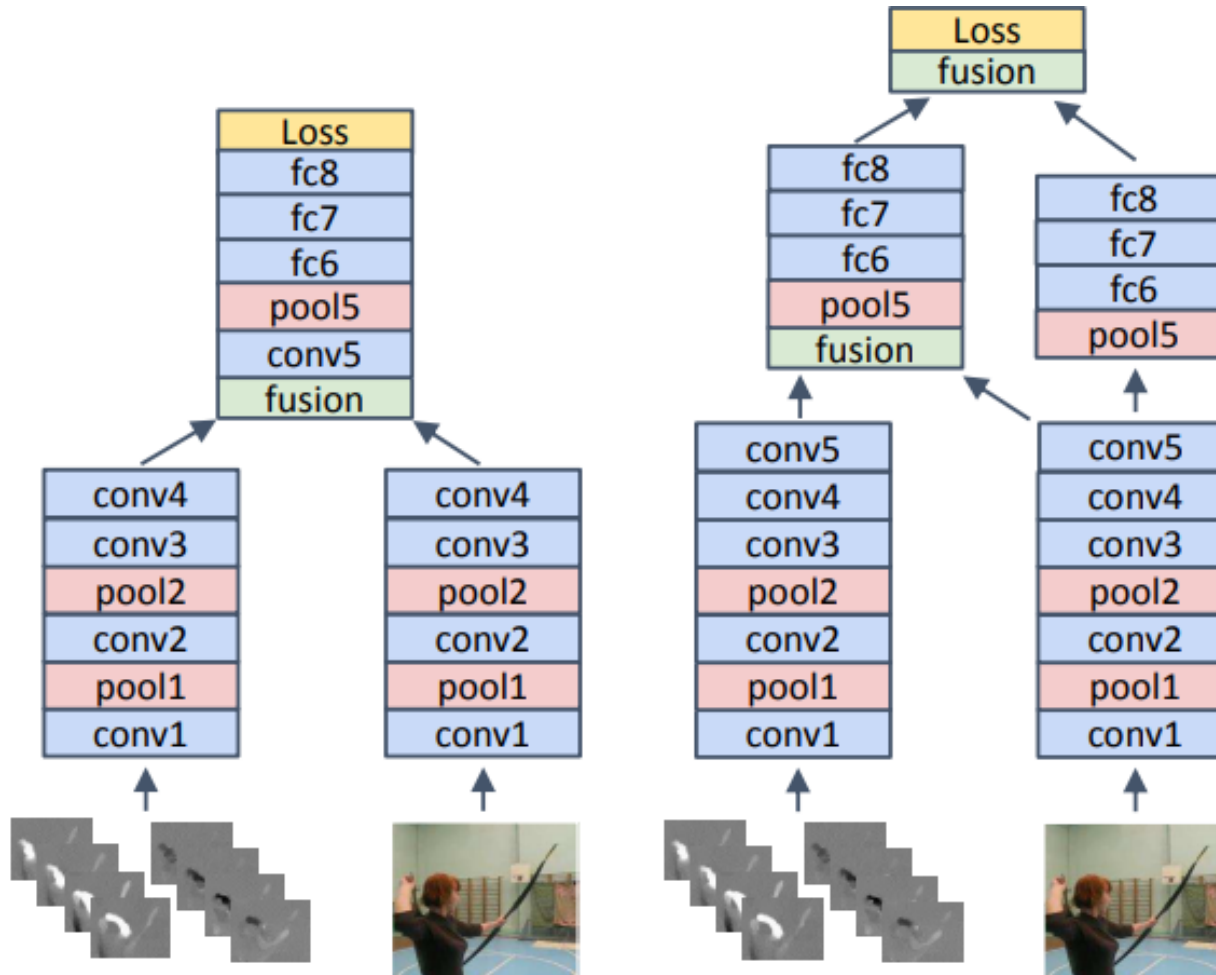
Not "end-to-end" (OF precomputed)



# Spatial + Temporal streams



# Next version: fusing the streams



# Benefit

Although 3D CNNs can capture temporal information, the two-stream architecture improves performance by leveraging on optic flow

# Choice #5: Two-stream inflated 3D CNNs

[Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset](#) (CVPR 2017)

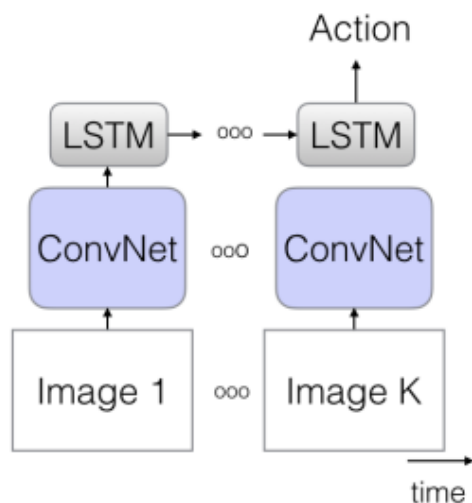
# Turn successful 2D CNNs into 3D CNNs

How? By turning  $N \times N$  filters into  $N \times N \times N$  ones

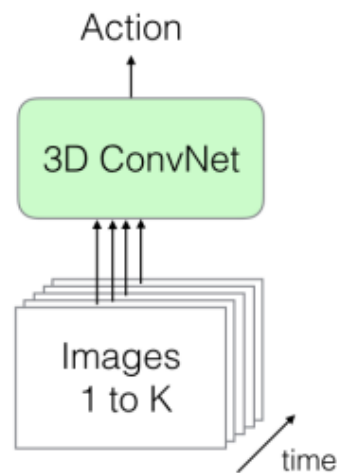
+ some other modifications

# Comparison study

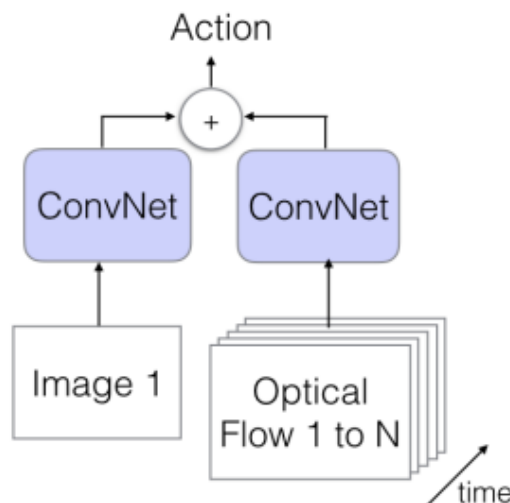
a) LSTM



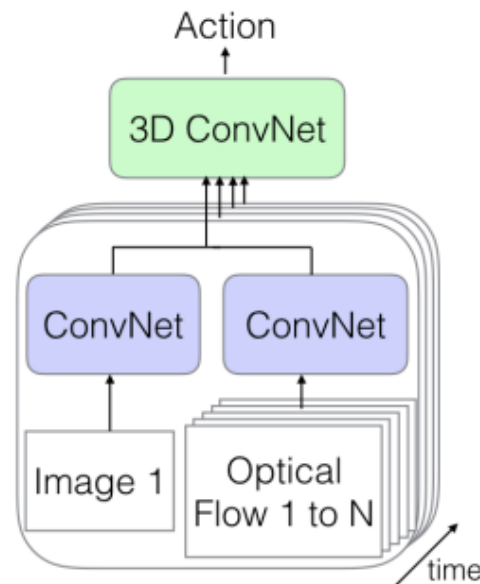
b) 3D-ConvNet



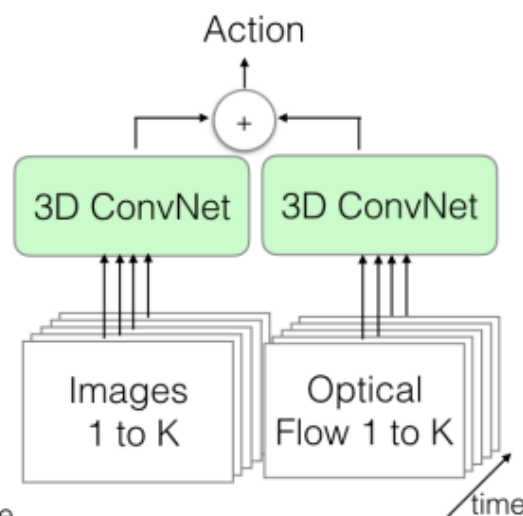
c) Two-Stream



d) 3D-Fused Two-Stream



e) Two-Stream 3D-ConvNet



Method	#Params	Training		Testing	
		# Input Frames	Temporal Footprint	# Input Frames	Temporal Footprint
ConvNet+LSTM	9M	25 rgb	5s	50 rgb	10s
3D-ConvNet	79M	16 rgb	0.64s	240 rgb	9.6s
Two-Stream	12M	1 rgb, 10 flow	0.4s	25 rgb, 250 flow	10s
3D-Fused	39M	5 rgb, 50 flow	2s	25 rgb, 250 flow	10s
Two-Stream I3D	25M	64 rgb, 64 flow	2.56s	250 rgb, 250 flow	10s

Architecture	UCF-101			HMDB-51			miniKinetics		
	RGB	Flow	RGB + Flow	RGB	Flow	RGB + Flow	RGB	Flow	RGB + Flow
(a) LSTM	81.0	–	–	36.0	–	–	69.9	–	–
(b) 3D-ConvNet	51.6	–	–	24.3	–	–	60.0	–	–
(c) Two-Stream	83.6	85.6	91.2	43.2	56.3	58.3	70.1	58.4	72.9
(d) 3D-Fused	83.2	85.8	89.3	49.2	55.5	56.8	71.4	61.0	74.0
(e) Two-Stream I3D	<b>84.5</b>	<b>90.6</b>	<b>93.4</b>	<b>49.8</b>	<b>61.9</b>	<b>66.4</b>	<b>74.1</b>	<b>69.6</b>	<b>78.7</b>