

Lecture 11

Video Recognition, Optical Flow

Slides from: Du Tran, Rick Szeliski, Steve Seitz,
Christoph Feichtenhofer

Overview

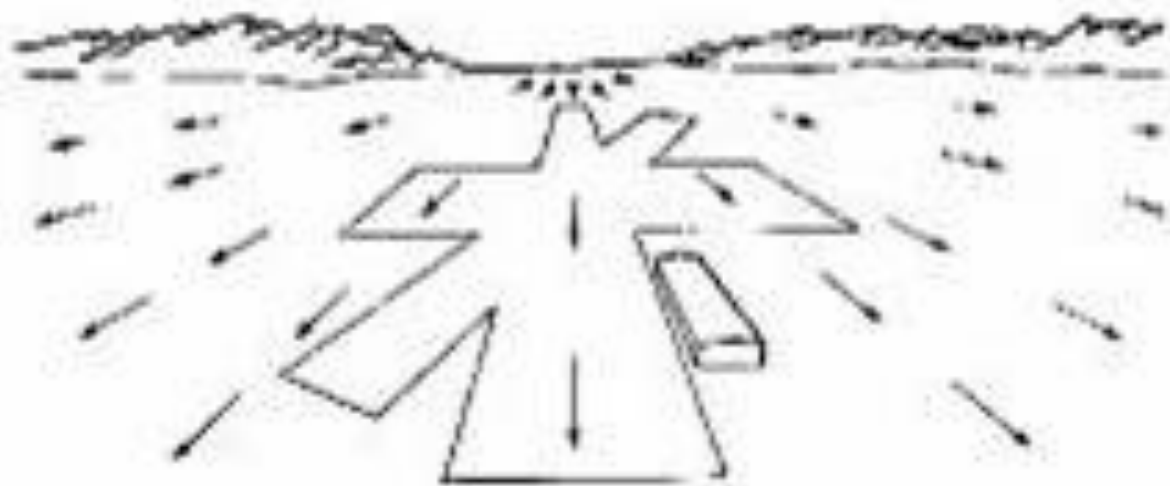
- Optical Flow
- ConvNets for Video
- Video Generation

Overview

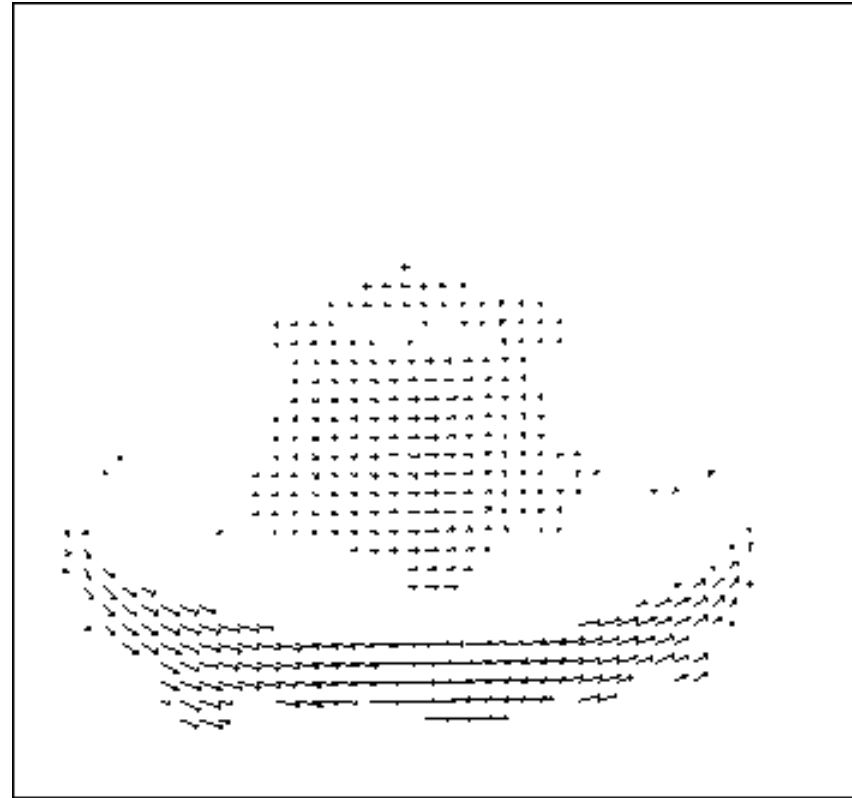
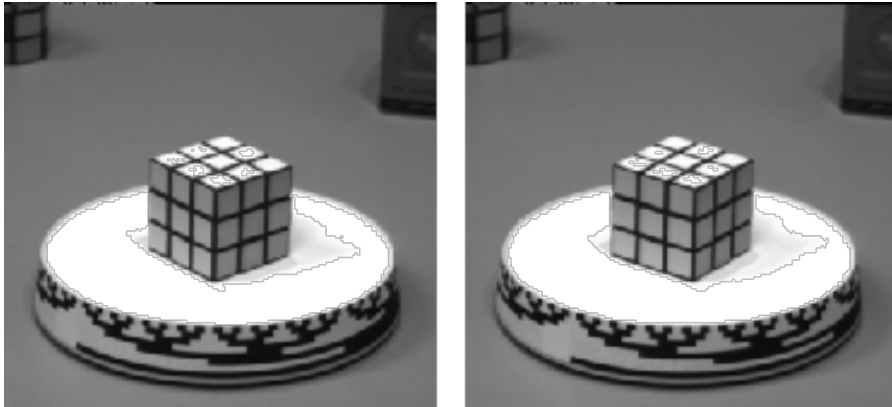
- Optical Flow
- ConvNets for Video
- Video Generation

Optical flow

Combination of slides from Rick Szeliski, Steve Seitz, Alyosha Efros and Bill Freeman and Fredo Durand



Motion estimation: Optical flow



Will start by estimating motion of each pixel separately
Then will consider motion of entire image

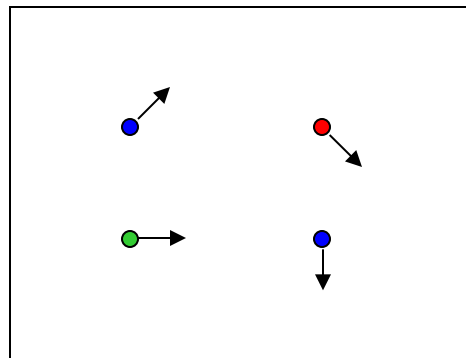
Why estimate motion?

Lots of uses

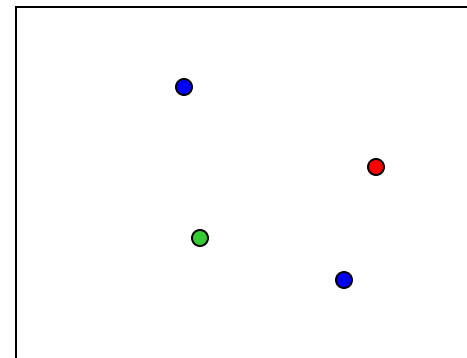
- Feature representation for DeepNets [coming up]
- Track object behavior
- Correct for camera jitter (stabilization)
- Align images (mosaics)
- 3D shape reconstruction
- Special effects



Problem definition: optical flow



$H(x, y)$



$I(x, y)$

How to estimate pixel motion from image H to image I?

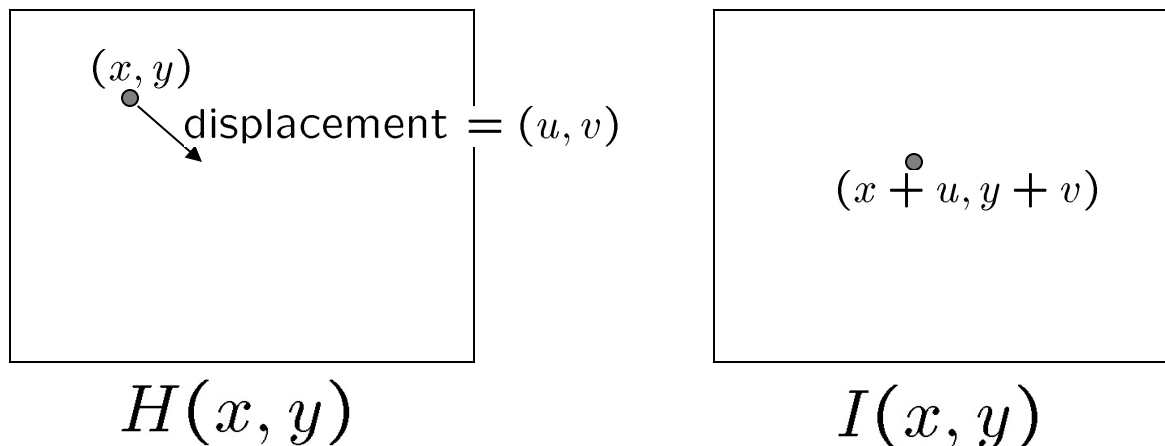
- Solve pixel correspondence problem
 - given a pixel in H, look for **nearby** pixels of the **same color** in I

Key assumptions

- **color constancy**: a point in H looks the same in I
 - For grayscale images, this is brightness constancy
- **small motion**: points do not move very far

This is called the optical flow problem

Optical flow constraints (grayscale images)



Let's look at these constraints more closely

- brightness constancy: Q: what's the equation?

$$H(x, y) = I(x + u, y + v)$$

- small motion: (u and v are less than 1 pixel)
 - suppose we take the Taylor series expansion of I :

$$\begin{aligned} I(x + u, y + v) &= I(x, y) + \frac{\partial I}{\partial x} u + \frac{\partial I}{\partial y} v + \text{higher order terms} \\ &\approx I(x, y) + \frac{\partial I}{\partial x} u + \frac{\partial I}{\partial y} v \end{aligned}$$

Optical flow equation

Combining these two equations

$$\begin{aligned}0 &= I(x + u, y + v) - H(x, y) && \text{shorthand: } I_x = \frac{\partial I}{\partial x} \\ &\approx I(x, y) + I_x u + I_y v - H(x, y) \\ &\approx (I(x, y) - H(x, y)) + I_x u + I_y v \\ &\approx I_t + I_x u + I_y v \\ &\approx I_t + \nabla I \cdot [u \ v]\end{aligned}$$

In the limit as u and v go to zero, this becomes exact

$$0 = I_t + \nabla I \cdot \left[\frac{\partial x}{\partial t} \ \frac{\partial y}{\partial t} \right]$$

Optical flow equation

$$0 = I_t + \nabla I \cdot [u \ v]$$

Q: how many unknowns and equations per pixel?

2 unknowns, one equation

Intuitively, what does this constraint mean?

- The component of the flow in the gradient direction is determined
- The component of the flow parallel to an edge is unknown

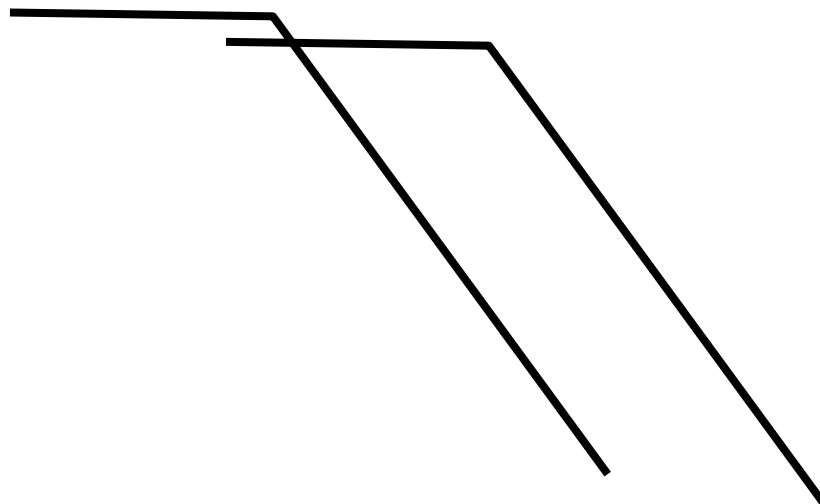
This explains the Barber Pole illusion

http://www.sandlotscience.com/Ambiguous/Barberpole_Illusion.htm

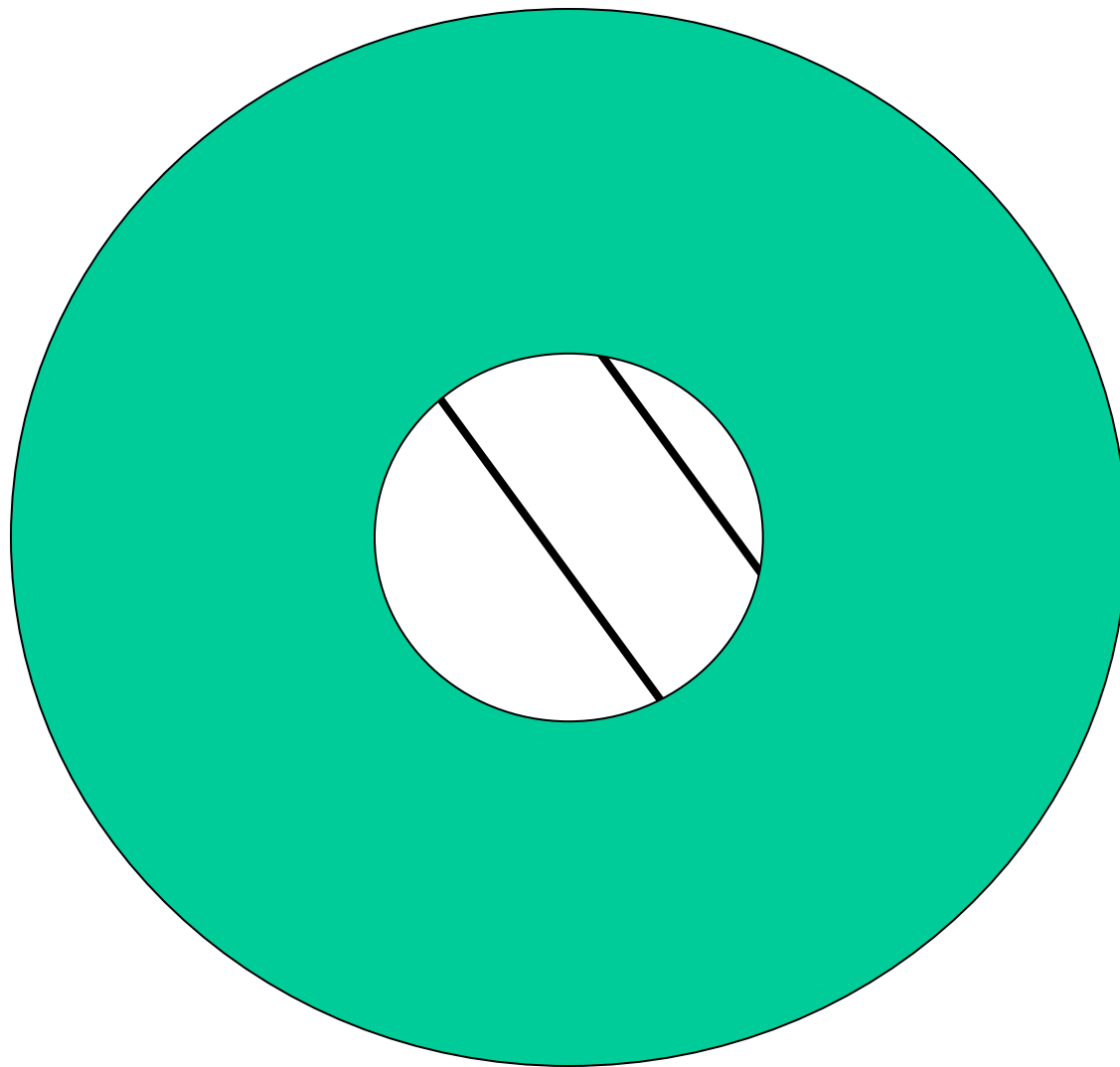
<http://www.liv.ac.uk/~marcob/Trieste/barberpole.html>



Aperture problem



Aperture problem



Solving the aperture problem

How to get more equations for a pixel?

- Basic idea: impose additional constraints
 - most common is to assume that the flow field is smooth locally
 - one method: pretend the pixel's neighbors have the same (u,v)
 - » If we use a 5x5 window, that gives us 25 equations per pixel!

$$0 = I_t(\mathbf{p}_i) + \nabla I(\mathbf{p}_i) \cdot [u \ v]$$

$$\begin{array}{c} \left[\begin{array}{cc} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ I_x(\mathbf{p}_2) & I_y(\mathbf{p}_2) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{25}) & I_y(\mathbf{p}_{25}) \end{array} \right] \begin{array}{c} \left[\begin{array}{c} u \\ v \end{array} \right] \\ \\ \\ \end{array} = - \begin{array}{c} \left[\begin{array}{c} I_t(\mathbf{p}_1) \\ I_t(\mathbf{p}_2) \\ \vdots \\ I_t(\mathbf{p}_{25}) \end{array} \right] \end{array} \\ \\ \begin{array}{ccc} A & d & b \\ 25 \times 2 & 2 \times 1 & 25 \times 1 \end{array} \end{array}$$

RGB version

How to get more equations for a pixel?

- Basic idea: impose additional constraints
 - most common is to assume that the flow field is smooth locally
 - one method: pretend the pixel's neighbors have the same (u,v)
 - » If we use a 5x5 window, that gives us 25*3 equations per pixel!

$$0 = I_t(\mathbf{p}_i)[0, 1, 2] + \nabla I(\mathbf{p}_i)[0, 1, 2] \cdot [u \ v]$$
$$\begin{bmatrix} I_x(\mathbf{p}_1)[0] & I_y(\mathbf{p}_1)[0] \\ I_x(\mathbf{p}_1)[1] & I_y(\mathbf{p}_1)[1] \\ I_x(\mathbf{p}_1)[2] & I_y(\mathbf{p}_1)[2] \\ \vdots & \vdots \\ I_x(\mathbf{p}_{25})[0] & I_y(\mathbf{p}_{25})[0] \\ I_x(\mathbf{p}_{25})[1] & I_y(\mathbf{p}_{25})[1] \\ I_x(\mathbf{p}_{25})[2] & I_y(\mathbf{p}_{25})[2] \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p}_1)[0] \\ I_t(\mathbf{p}_1)[1] \\ I_t(\mathbf{p}_1)[2] \\ \vdots \\ I_t(\mathbf{p}_{25})[0] \\ I_t(\mathbf{p}_{25})[1] \\ I_t(\mathbf{p}_{25})[2] \end{bmatrix}$$

A
75x2

d
2x1

b
75x1

Note that RGB is not enough to disambiguate
because R, G & B are correlated

Just provides better gradient

Lukas-Kanade flow

Prob: we have more equations than unknowns

$$\begin{array}{ccc} A & d = b & \longrightarrow \text{minimize } \|Ad - b\|^2 \\ 25 \times 2 & 2 \times 1 \quad 25 \times 1 & \end{array}$$

Solution: solve least squares problem

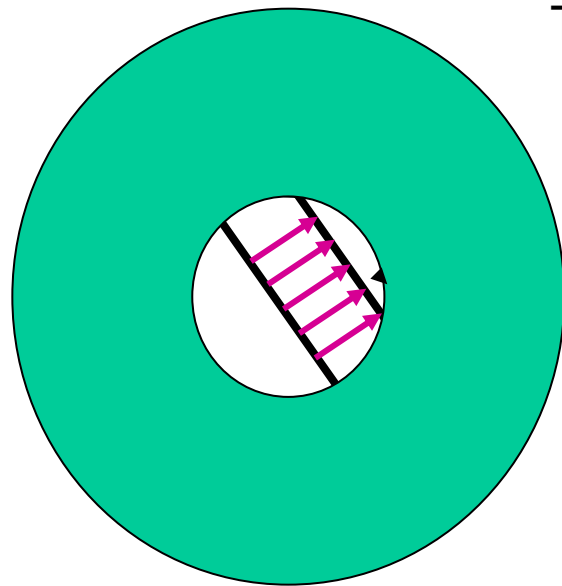
- minimum least squares solution given by solution (in d) of:

$$\begin{array}{ccc} (A^T A) & d = & A^T b \\ 2 \times 2 & 2 \times 1 & 2 \times 1 \end{array}$$

$$\begin{array}{ccc} \left[\begin{array}{cc} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{array} \right] & \left[\begin{array}{c} u \\ v \end{array} \right] & = - \left[\begin{array}{c} \sum I_x I_t \\ \sum I_y I_t \end{array} \right] \\ A^T A & & A^T b \end{array}$$

- The summations are over all pixels in the K x K window
- This technique was first proposed by Lukas & Kanade (1981)

Aperture Problem and Normal Flow



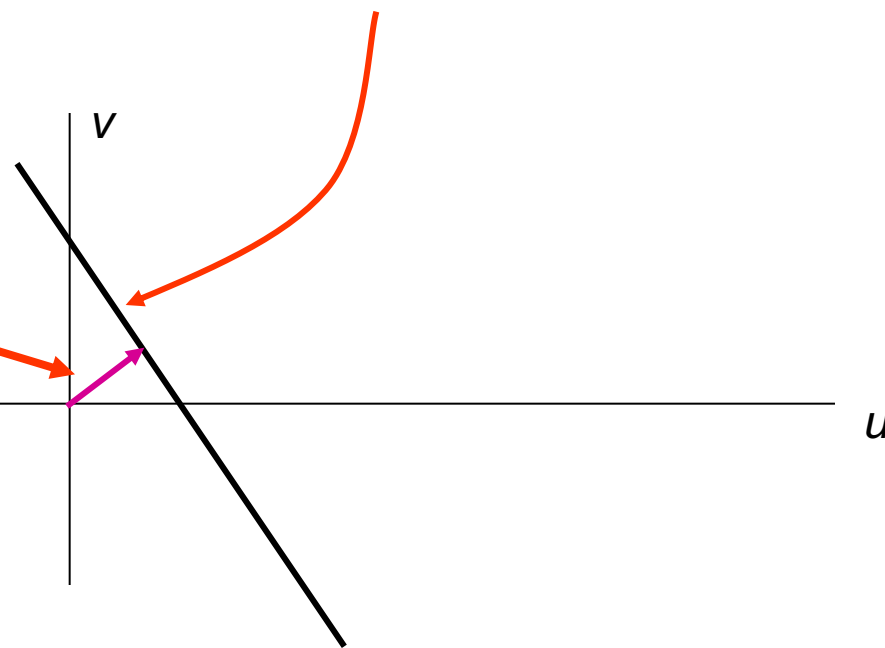
The gradient constraint:

$$I_x u + I_y v + I_t = 0$$
$$\nabla I \bullet \vec{U} = 0$$

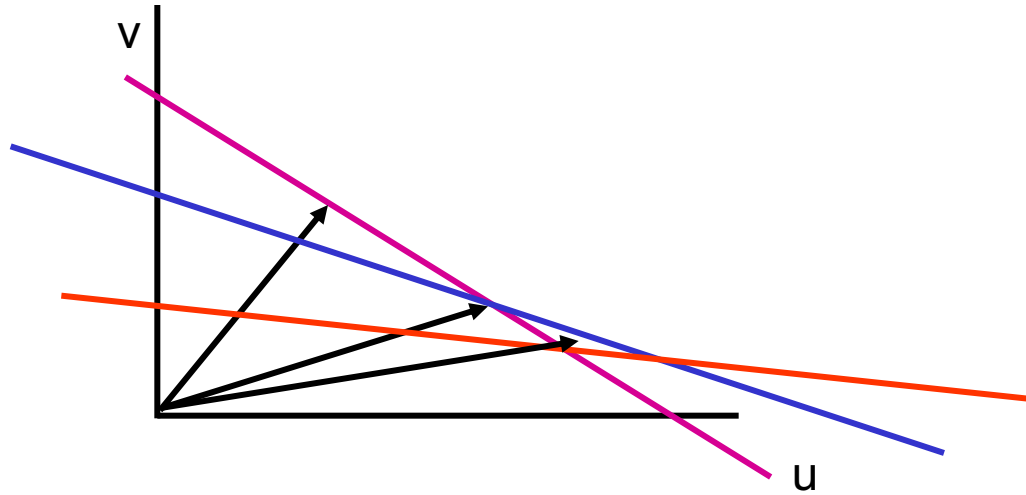
Defines a line in the (u, v) space

Normal Flow:

$$u_{\perp} = -\frac{I_t}{|\nabla I|} \frac{\nabla I}{|\nabla I|}$$



Combining Local Constraints



$$\nabla I^1 \bullet U = -I_t^1$$

$$\nabla I^2 \bullet U = -I_t^2$$

$$\nabla I^3 \bullet U = -I_t^3$$

etc.

Conditions for solvability

- Optimal (u, v) satisfies Lucas-Kanade equation

$$\begin{matrix} \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} & \begin{bmatrix} u \\ v \end{bmatrix} & = & - & \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix} \\ & A^T A & & & A^T b \end{matrix}$$

When is This Solvable?

- $A^T A$ should be invertible
- $A^T A$ should not be too small due to noise
 - eigenvalues λ_1 and λ_2 of $A^T A$ should not be too small
- $A^T A$ should be well-conditioned
 - λ_1 / λ_2 should not be too large ($\lambda_1 =$ larger eigenvalue)

$A^T A$ is solvable when there is no aperture problem

$$A^T A = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = \sum \begin{bmatrix} I_x \\ I_y \end{bmatrix} [I_x \ I_y] = \sum \nabla I (\nabla I)^T$$

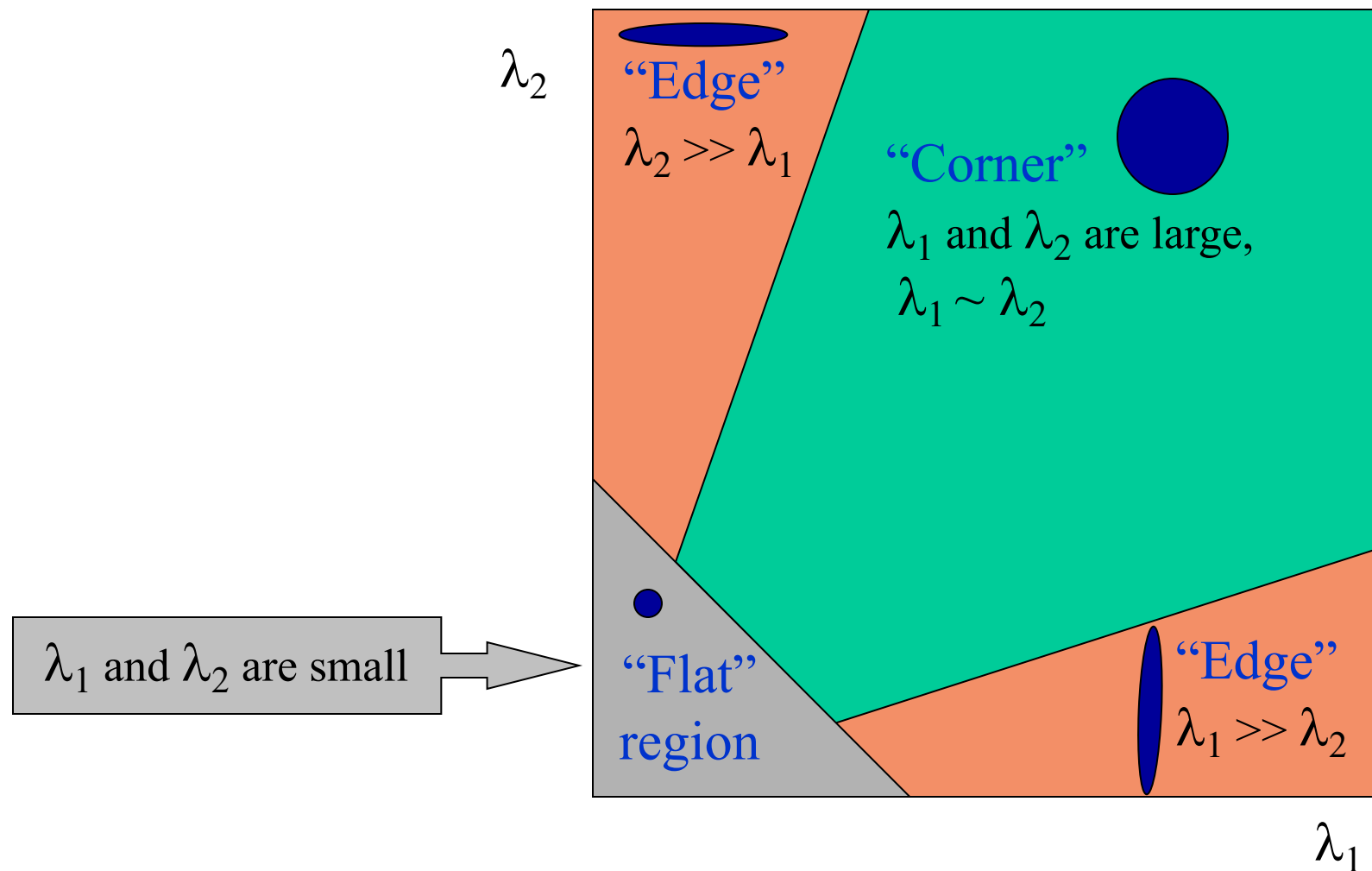
Eigenvectors of $A^T A$

$$A^T A = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = \sum \begin{bmatrix} I_x \\ I_y \end{bmatrix} [I_x \ I_y] = \sum \nabla I (\nabla I)^T$$

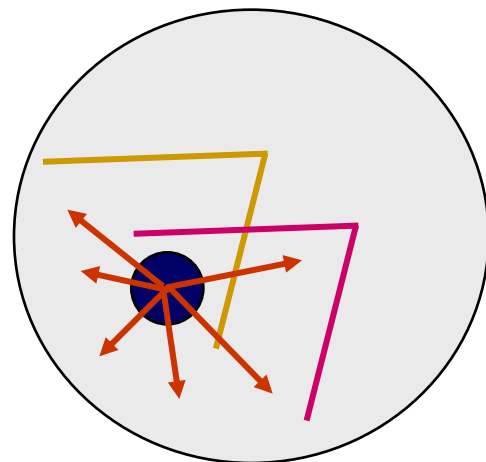
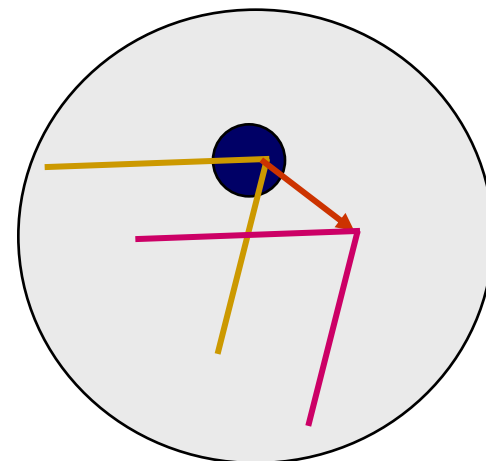
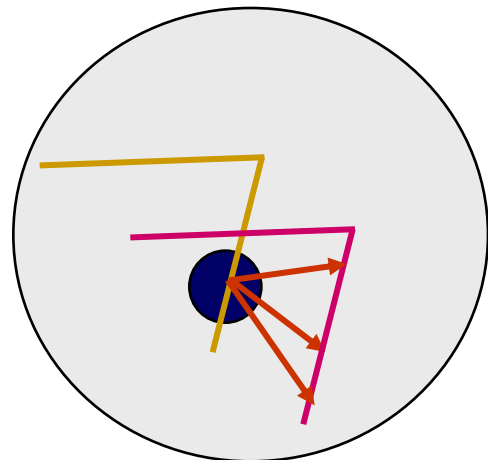
- $M = A^T A$ is the *second moment matrix*
- The eigenvectors and eigenvalues of M relate to edge direction and magnitude
 - The eigenvector associated with the larger eigenvalue points in the direction of fastest intensity change
 - The other eigenvector is orthogonal to it

Interpreting the eigenvalues

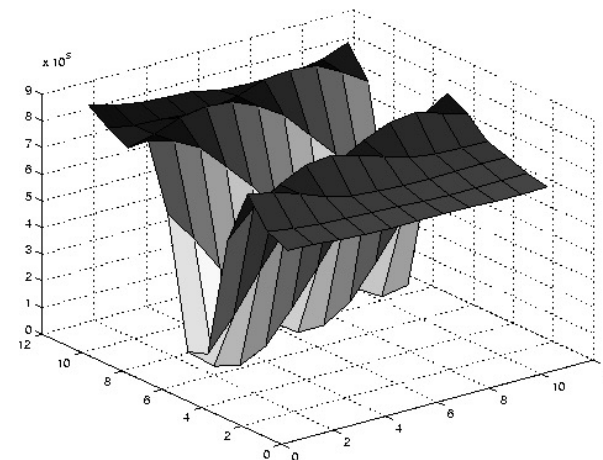
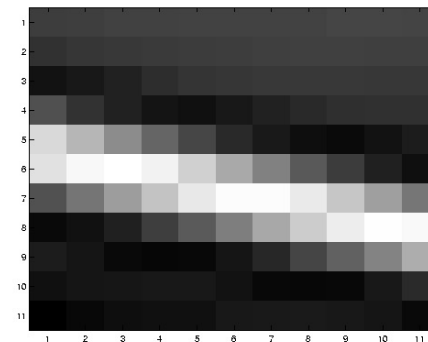
Classification of image points using eigenvalues of the second moment matrix:



Local Patch Analysis



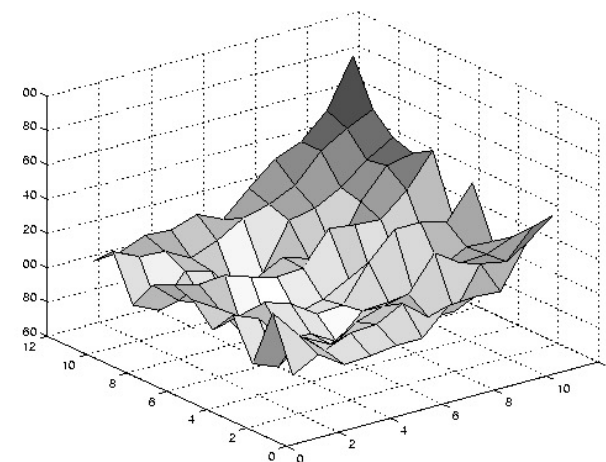
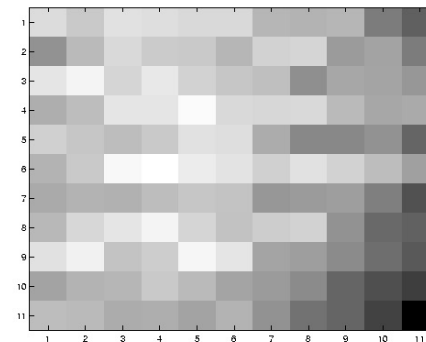
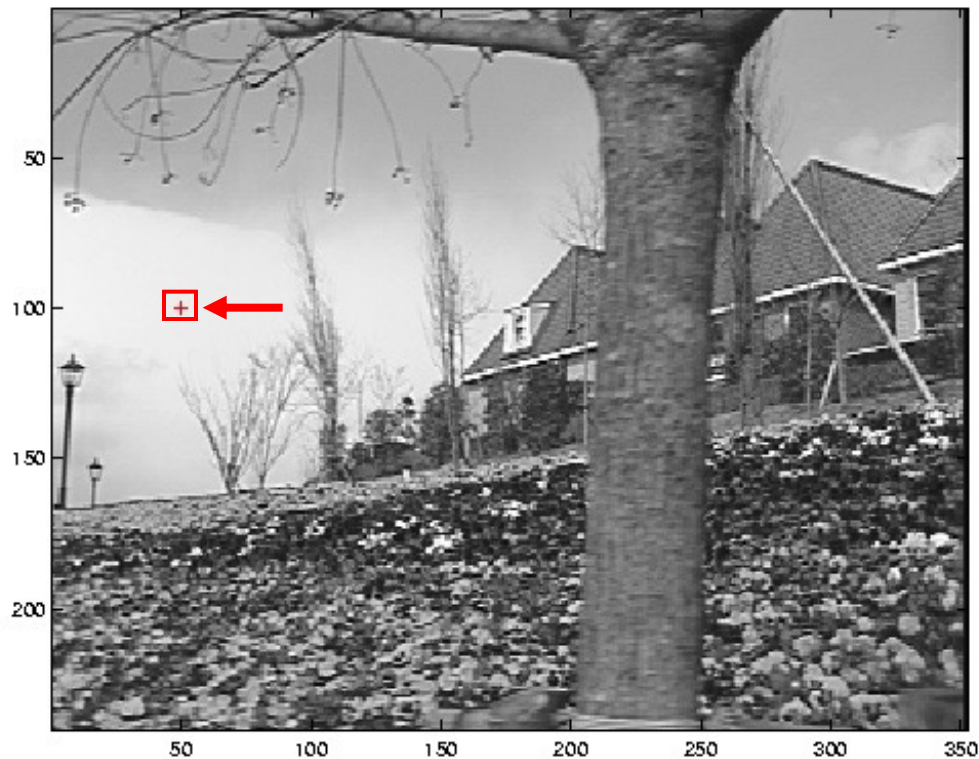
Edge



$$\sum \nabla I (\nabla I)^T$$

- large gradients, all the same
- large λ_1 , small λ_2

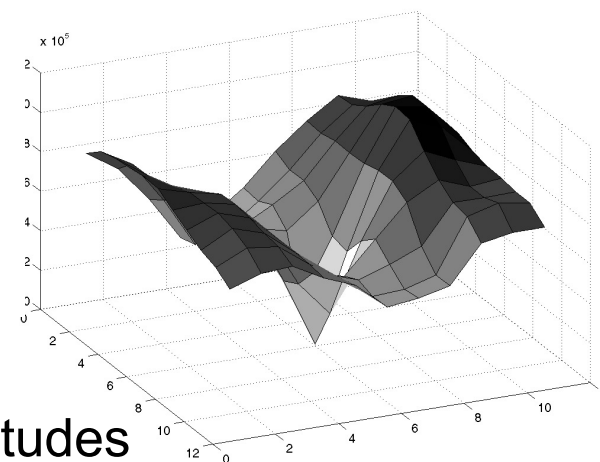
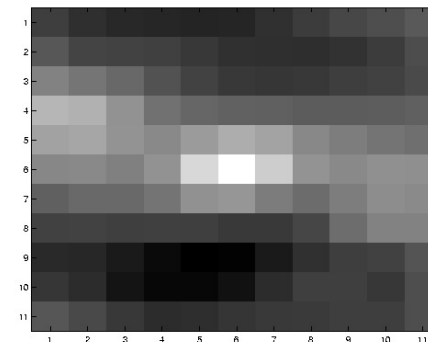
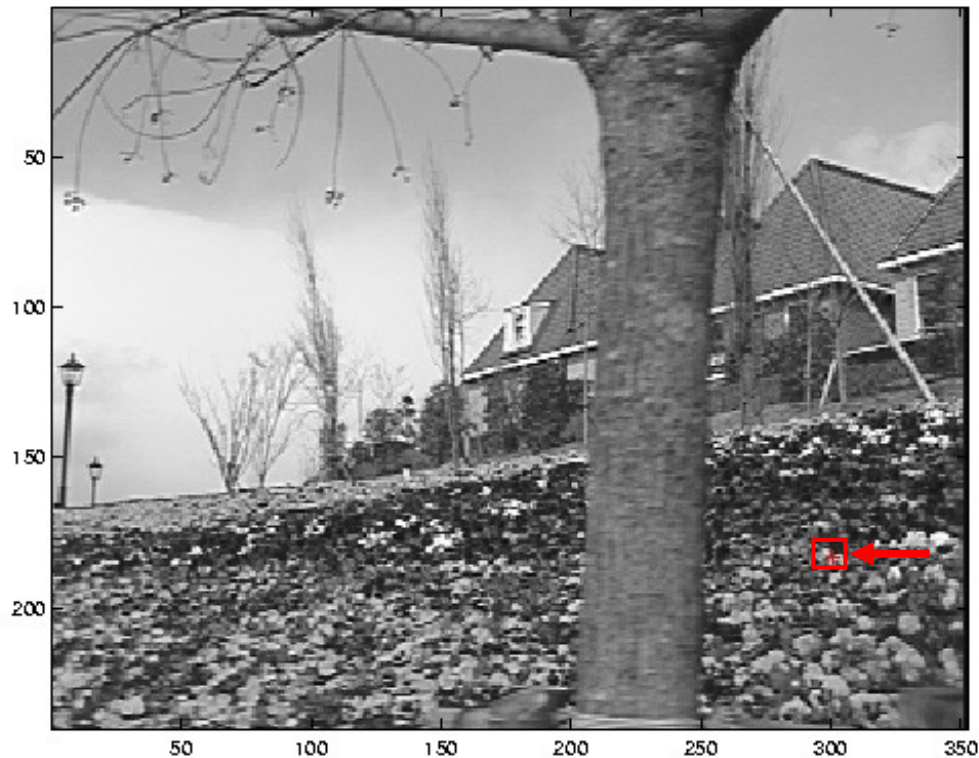
Low texture region



$$\sum \nabla I (\nabla I)^T$$

- gradients have small magnitude
- small λ_1 , small λ_2

High textured region



$$\sum \nabla I (\nabla I)^T$$

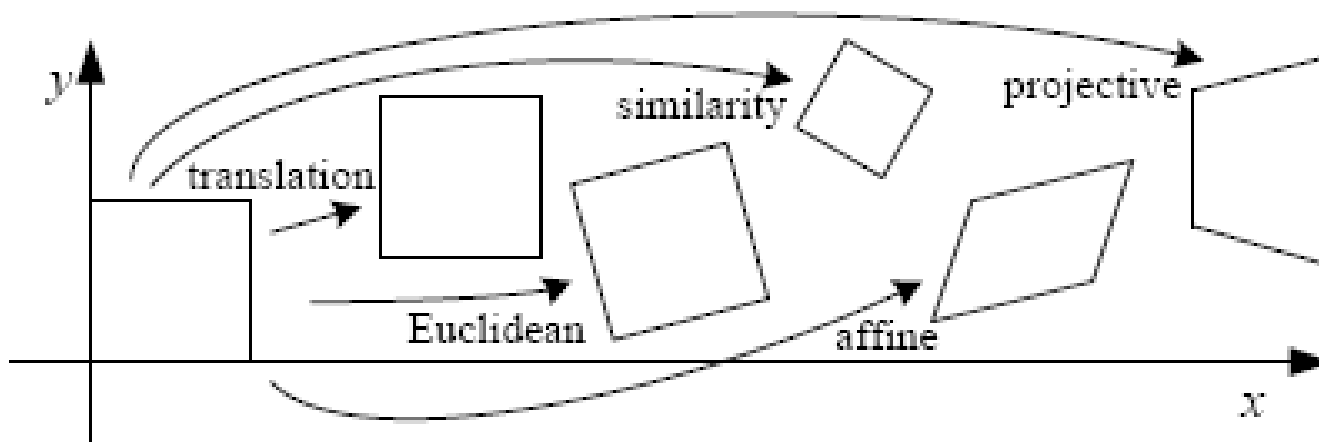
- gradients are different, large magnitudes
- large λ_1 , large λ_2

Observation

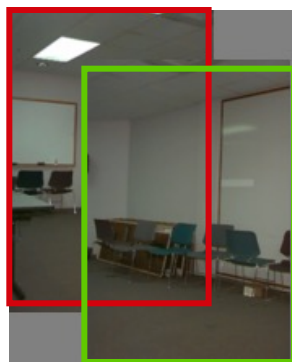
This is a two image problem BUT

- Can measure sensitivity by just looking at one of the images!
- This tells us which pixels are easy to track, which are hard
 - very useful later on when we do feature tracking...

Motion models

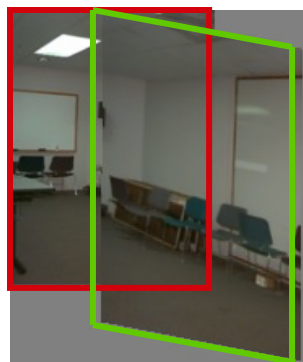


Translation



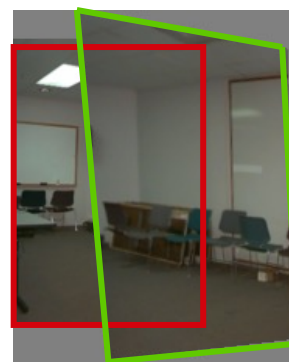
2 unknowns

Affine



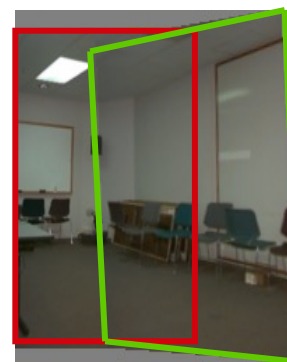
6 unknowns

Perspective



8 unknowns

3D rotation



3 unknowns

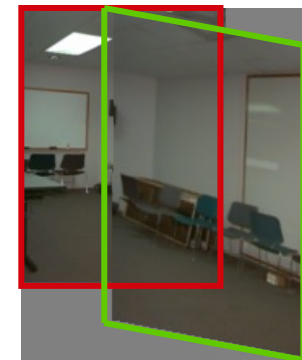
Affine motion

$$u(x, y) = a_1 + a_2x + a_3y$$

$$v(x, y) = a_4 + a_5x + a_6y$$

- Substituting into the brightness constancy equation:

$$I_x \cdot u + I_y \cdot v + I_t \approx 0$$

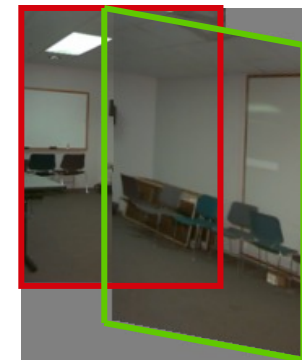


Affine motion

$$u(x, y) = a_1 + a_2x + a_3y$$

$$v(x, y) = a_4 + a_5x + a_6y$$

- Substituting into the brightness constancy equation:



$$I_x(a_1 + a_2x + a_3y) + I_y(a_4 + a_5x + a_6y) + I_t \approx 0$$

- Each pixel provides 1 linear constraint in 6 unknowns
- Least squares minimization:

$$Err(\vec{a}) = \sum \left[I_x(a_1 + a_2x + a_3y) + I_y(a_4 + a_5x + a_6y) + I_t \right]^2$$

Errors in Lukas-Kanade

What are the potential causes of errors in this procedure?

- Suppose $A^T A$ is easily invertible
- Suppose there is not much noise in the image

When our assumptions are violated

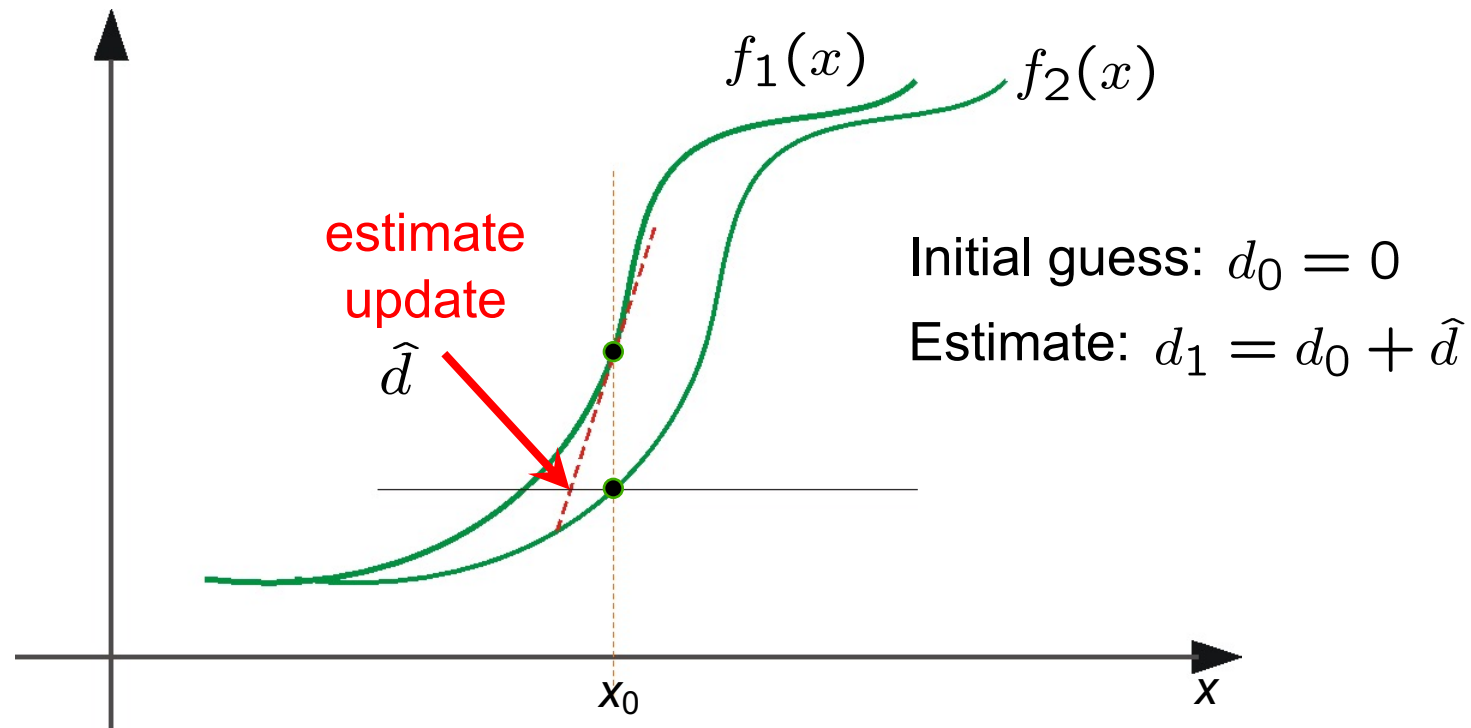
- Brightness constancy is not satisfied
- The motion is not small
- A point does not move like its neighbors
 - window size is too large
 - what is the ideal window size?

Iterative Refinement

Iterative Lukas-Kanade Algorithm

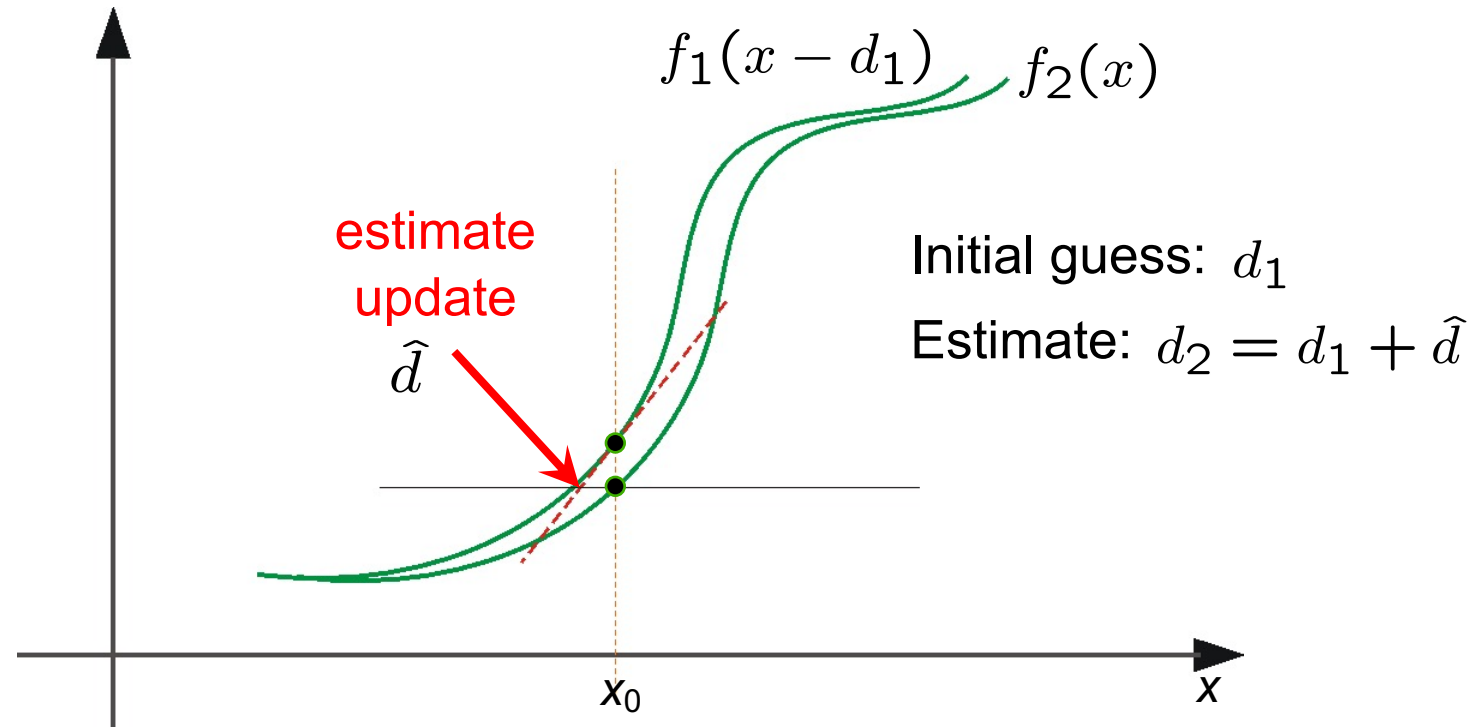
1. Estimate velocity at each pixel by solving Lucas-Kanade equations
2. Warp H towards I using the estimated flow field
 - *use image warping techniques*
3. Repeat until convergence

Optical Flow: Iterative Estimation

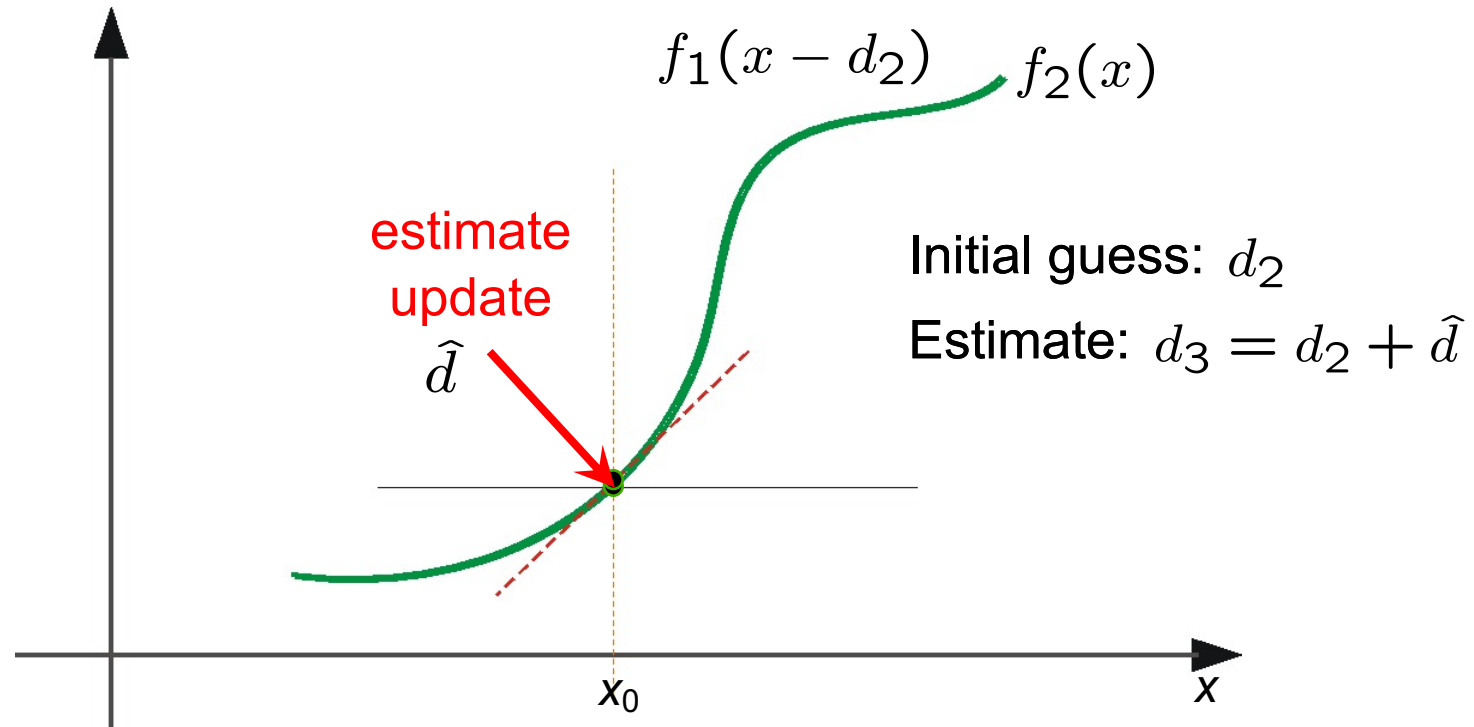


(using d for *displacement* here instead of u)

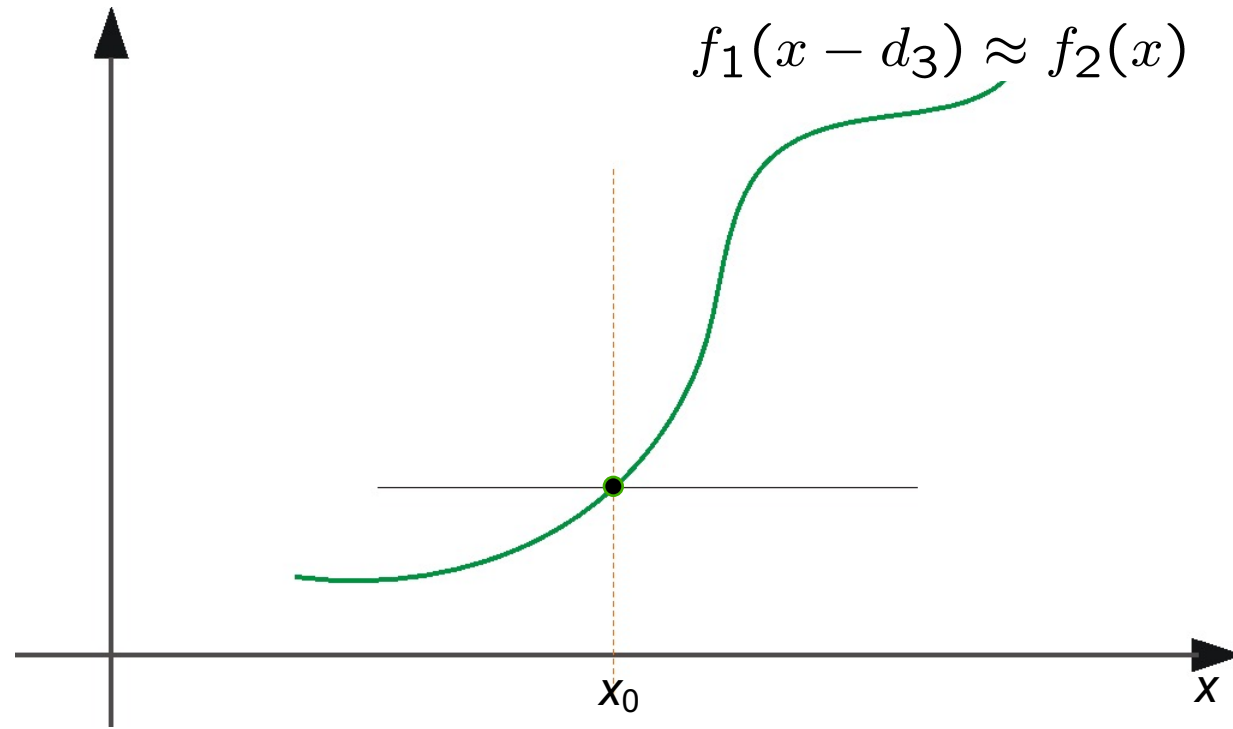
Optical Flow: Iterative Estimation



Optical Flow: Iterative Estimation



Optical Flow: Iterative Estimation



Optical Flow: Iterative Estimation

Some Implementation Issues:

- Warping is not easy (ensure that errors in warping are smaller than the estimate refinement)
- Warp one image, take derivatives of the other so you don't need to re-compute the gradient after each iteration.
- Often useful to low-pass filter the images before motion estimation (for better derivative estimation, and linear approximations to image intensity)

Revisiting the small motion assumption



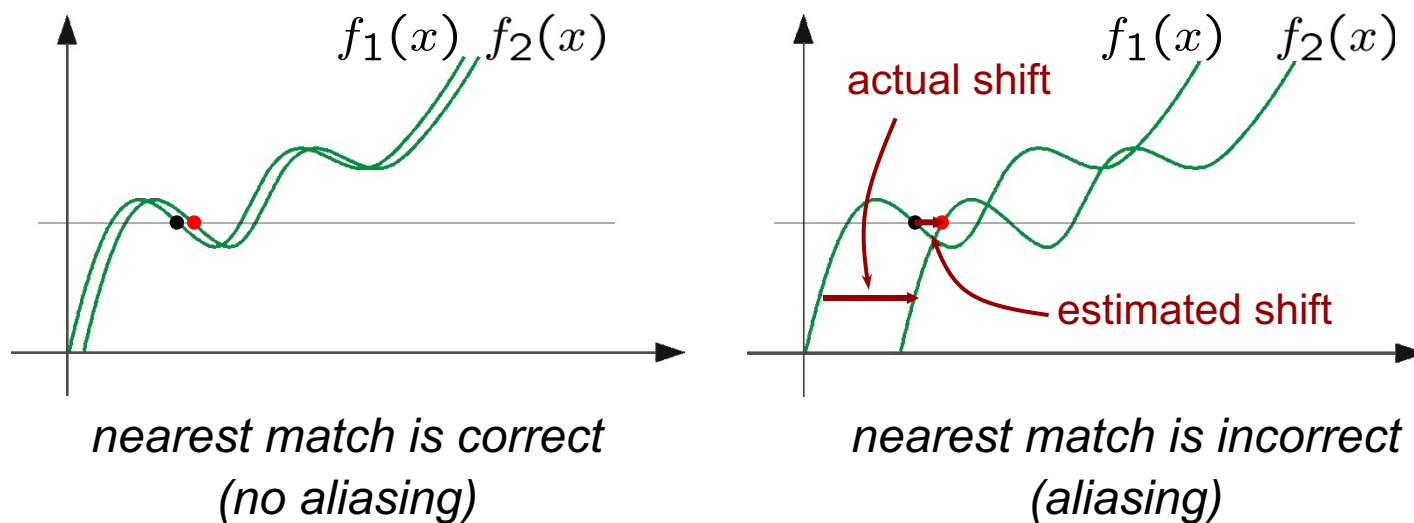
Is this motion small enough?

- Probably not—it's much larger than one pixel (2nd order terms dominate)
- How might we solve this problem?

Optical Flow: Aliasing

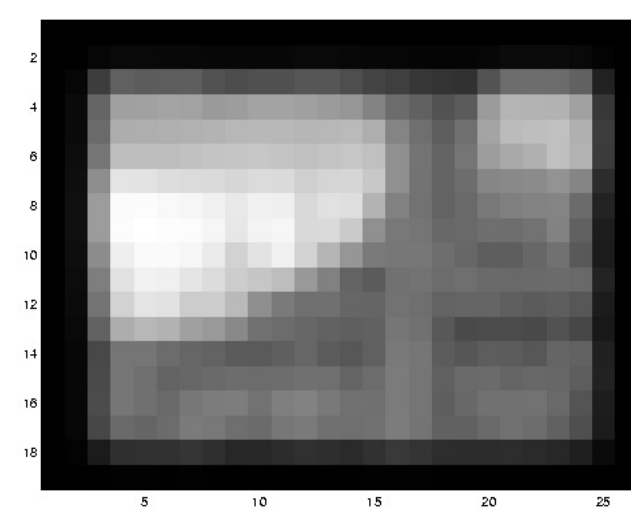
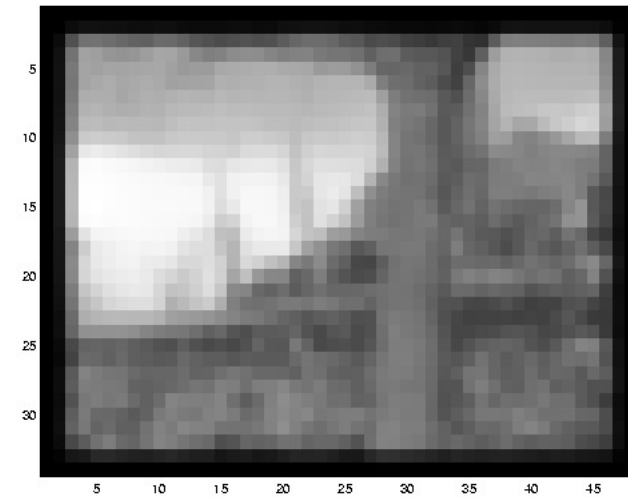
Temporal aliasing causes ambiguities in optical flow because images can have many pixels with the same intensity.

I.e., how do we know which 'correspondence' is correct?

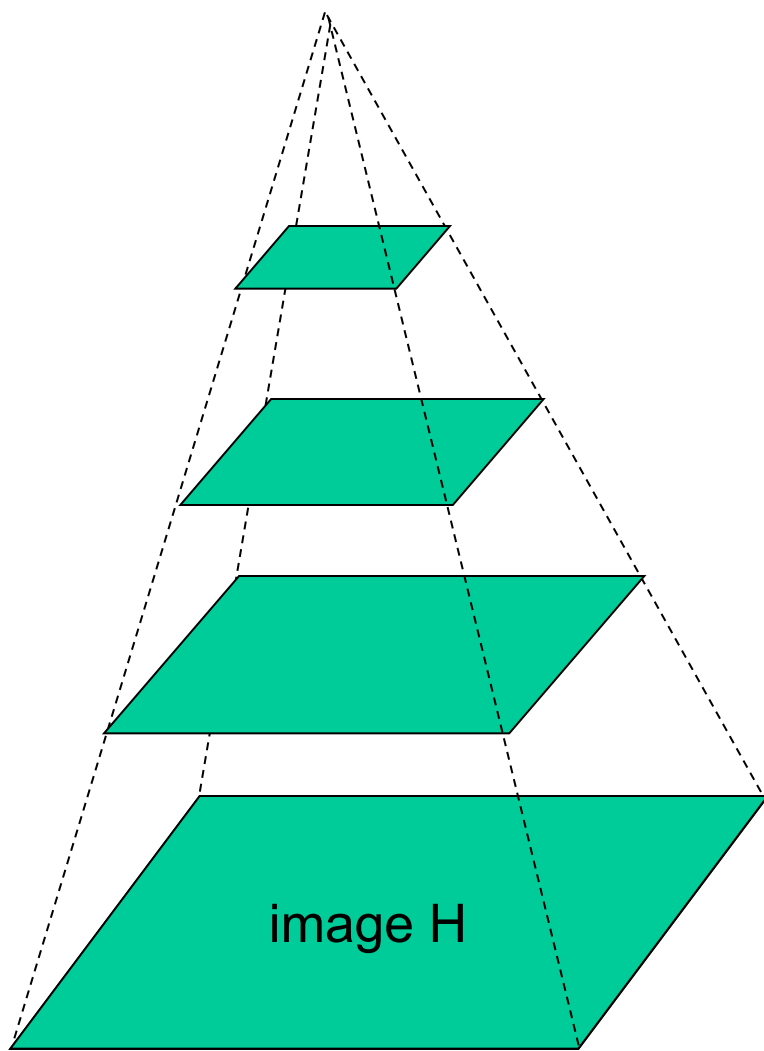


To overcome aliasing: coarse-to-fine estimation.

Reduce the resolution!



Coarse-to-fine optical flow estimation



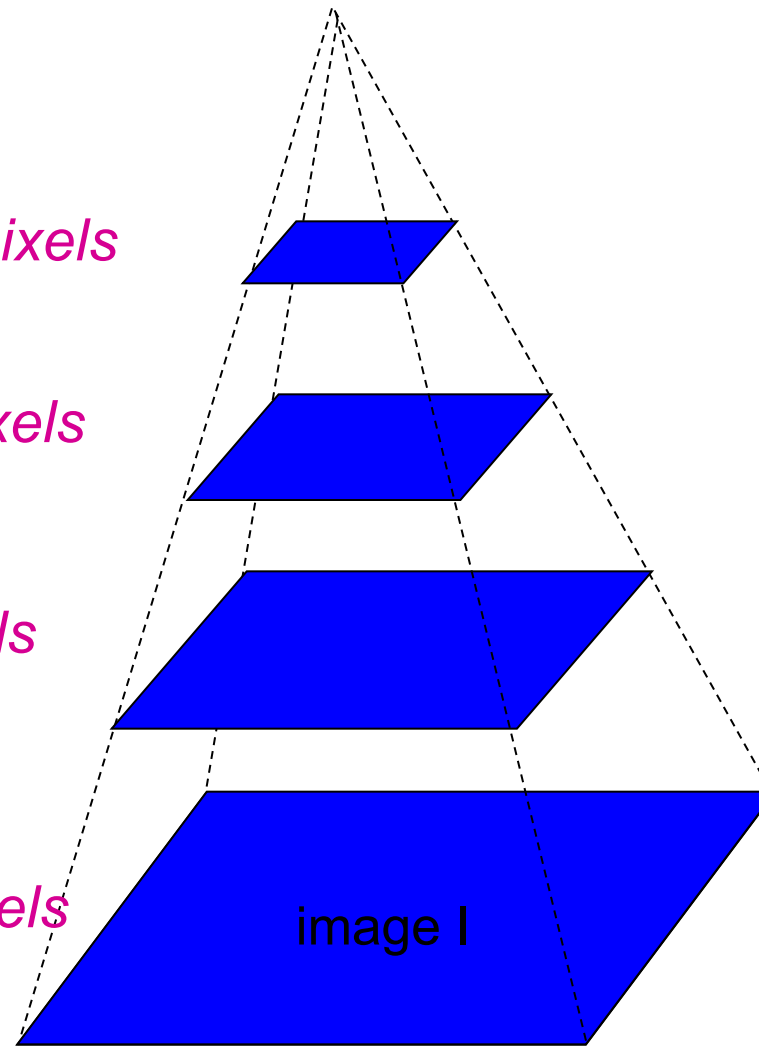
Gaussian pyramid of image H

$u=1.25$ pixels

$u=2.5$ pixels

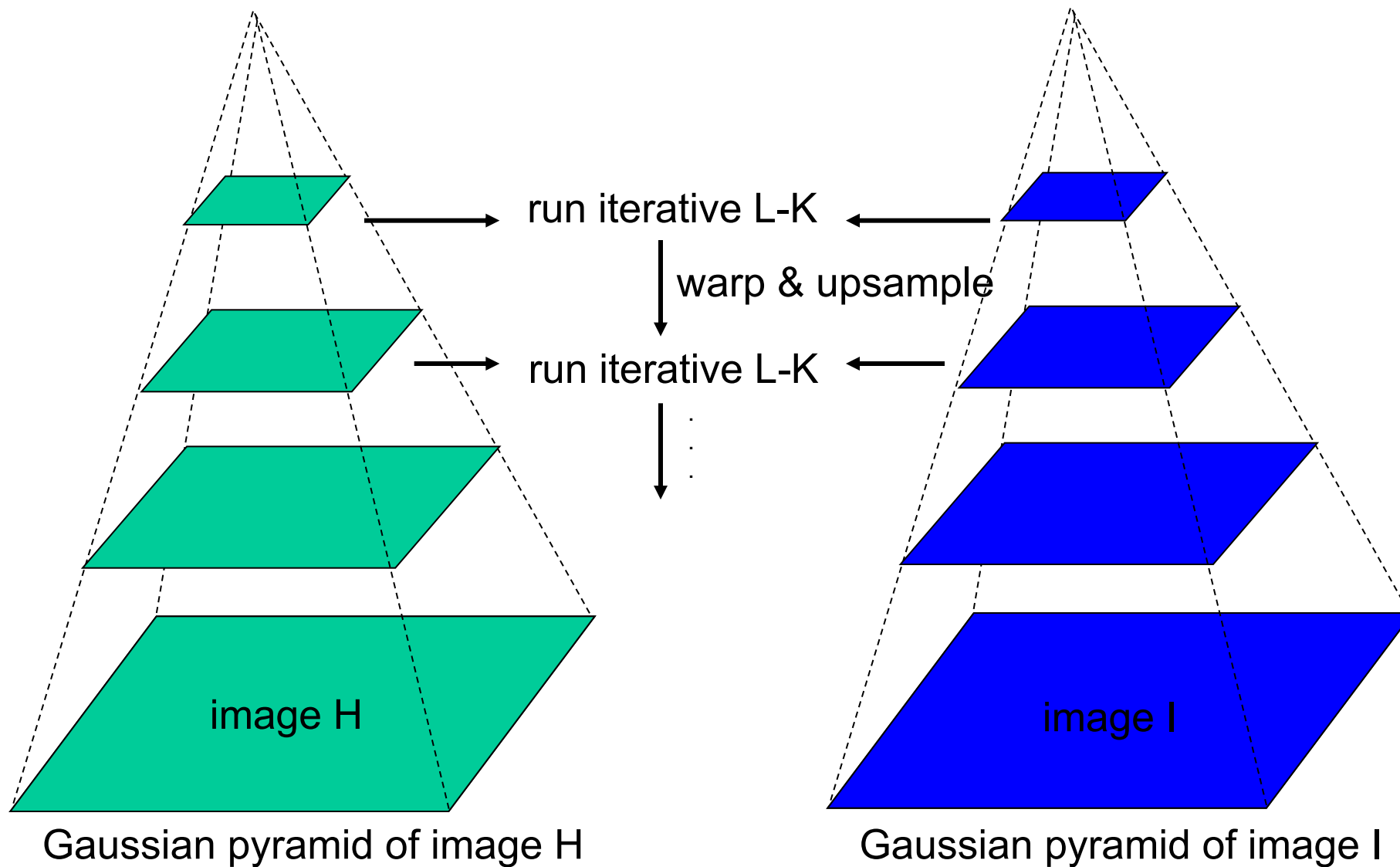
$u=5$ pixels

$u=10$ pixels



Gaussian pyramid of image I

Coarse-to-fine optical flow estimation



Recap: Classes of Techniques

Direct-methods (e.g. optical flow)

- Directly recover image motion from spatio-temporal image brightness variations
- Global motion parameters directly recovered without an intermediate feature motion calculation
- Dense motion fields, but more sensitive to appearance variations
- Suitable for video and when image motion is small (< 10 pixels)

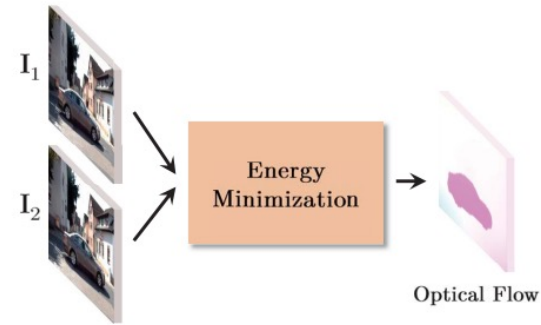
Feature-based methods (e.g. SIFT+Ransac+regression)

[To be covered next lecture]

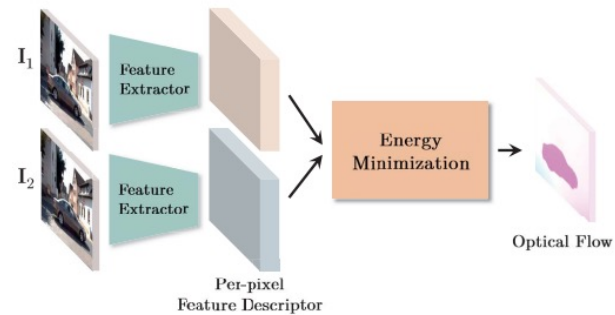
- Extract visual features (corners, textured areas) and track them over multiple frames
- Sparse motion fields, but possibly robust tracking
- Suitable especially when image motion is large (10-s of pixels)

Optical Flow Estimation in the Deep Learning Age

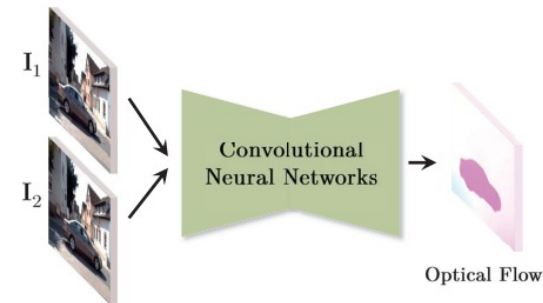
[Hur & Roth, 2020] Survey paper
<https://arxiv.org/pdf/2004.02853.pdf>



(a) Classical energy-based approach



(b) Using CNNs as a feature extractor



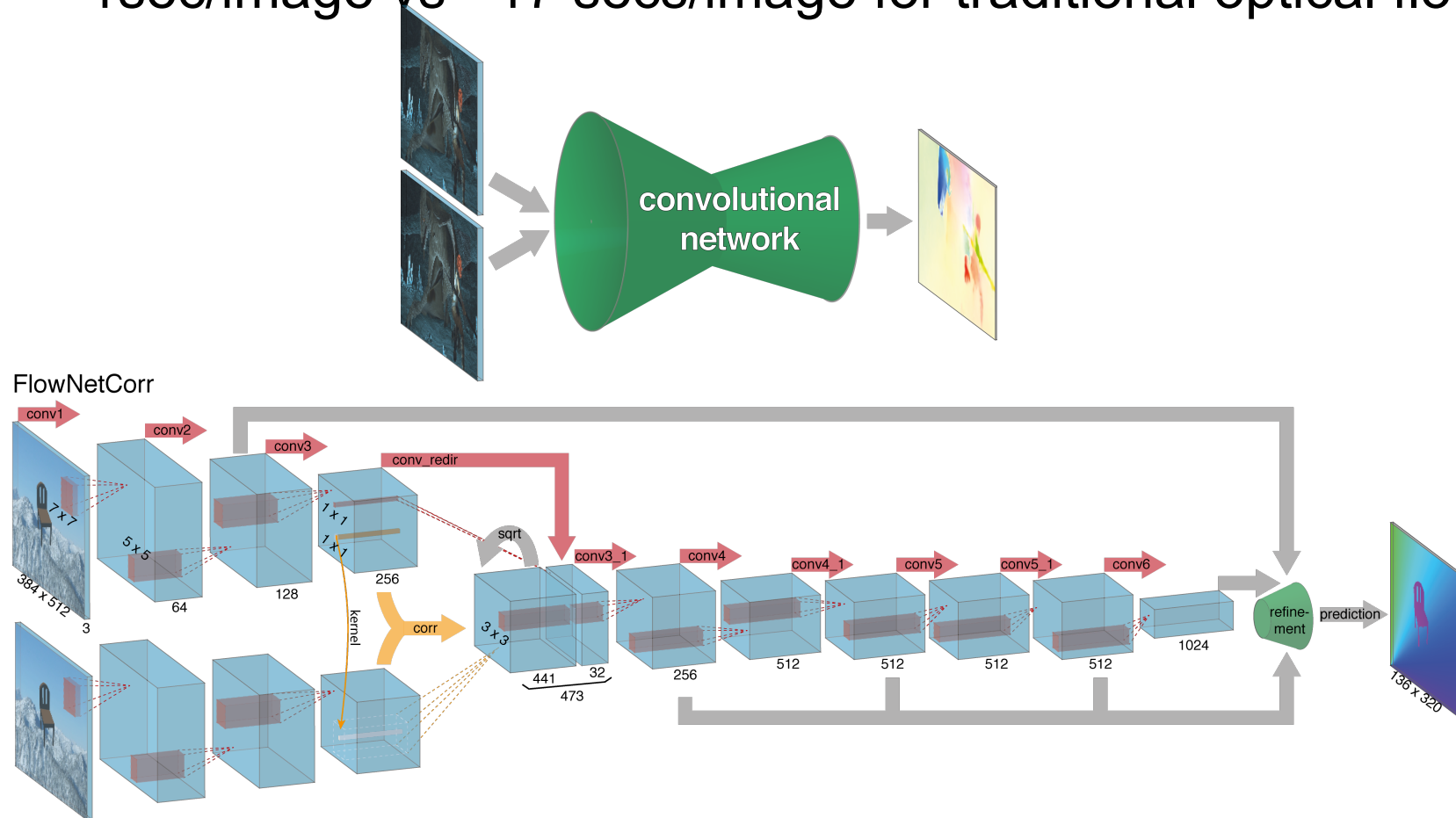
(c) CNN regression architecture

Fig. 1: Transition from (a) classical energy-based approaches to (b) CNN-based approaches that use CNNs as a feature extractor or to (c) end-to-end trainable CNN regression architectures.

FlowNet

FlowNet: Learning Optical Flow with Convolutional Networks [Fischer et al. 2015]

~ 1sec/image vs ~17 secs/image for traditional optical flow



FlowNet

- FlowNet: Learning Optical Flow with Convolutional Networks [Fischer et al. 2015]
End-to-end regression approach

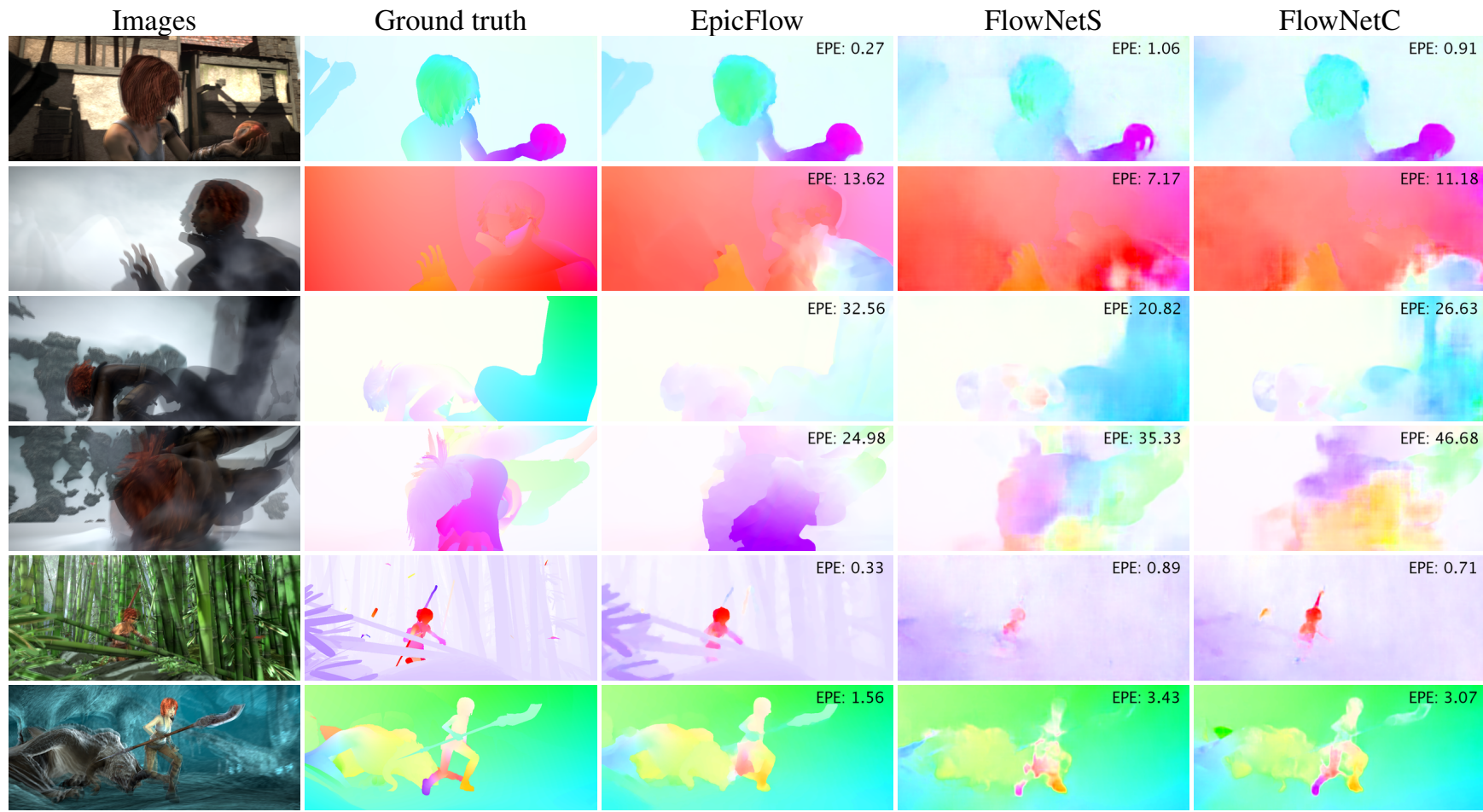


Table 2: Quantitative comparison on public benchmarks: MPI Sintel [8] and KITTI [14, 37].

Methods	MPI Sintel ^a		KITTI ^b	
	Clean	Final	2012	2015
FlowNetS [10]	6.158	7.218	37.05 %	–
FlowNetC [10]	6.081	7.883	–	–
SPyNet [41]	6.640	8.360	12.31 %	35.07 %
FlowNet2 [25]	3.959	6.016	4.82 %	10.41 %
PWC-Net [48]	4.386	5.042	4.22 %	9.60 %
LiteFlowNet [23]	3.449	5.381	3.27 %	9.38 %
IRR-PWC [24]	3.844	4.579	3.21 %	7.65 %
HD ³ [58]	4.788	4.666	2.26 %	6.55 %
VCN [56]	2.808	4.404	–	6.30 %

^a Evaluation metric: end point error (EPE).

^b Evaluation metric: outlier rate (*i.e.* less than 3 pixel or 5% error is considered an inlier)

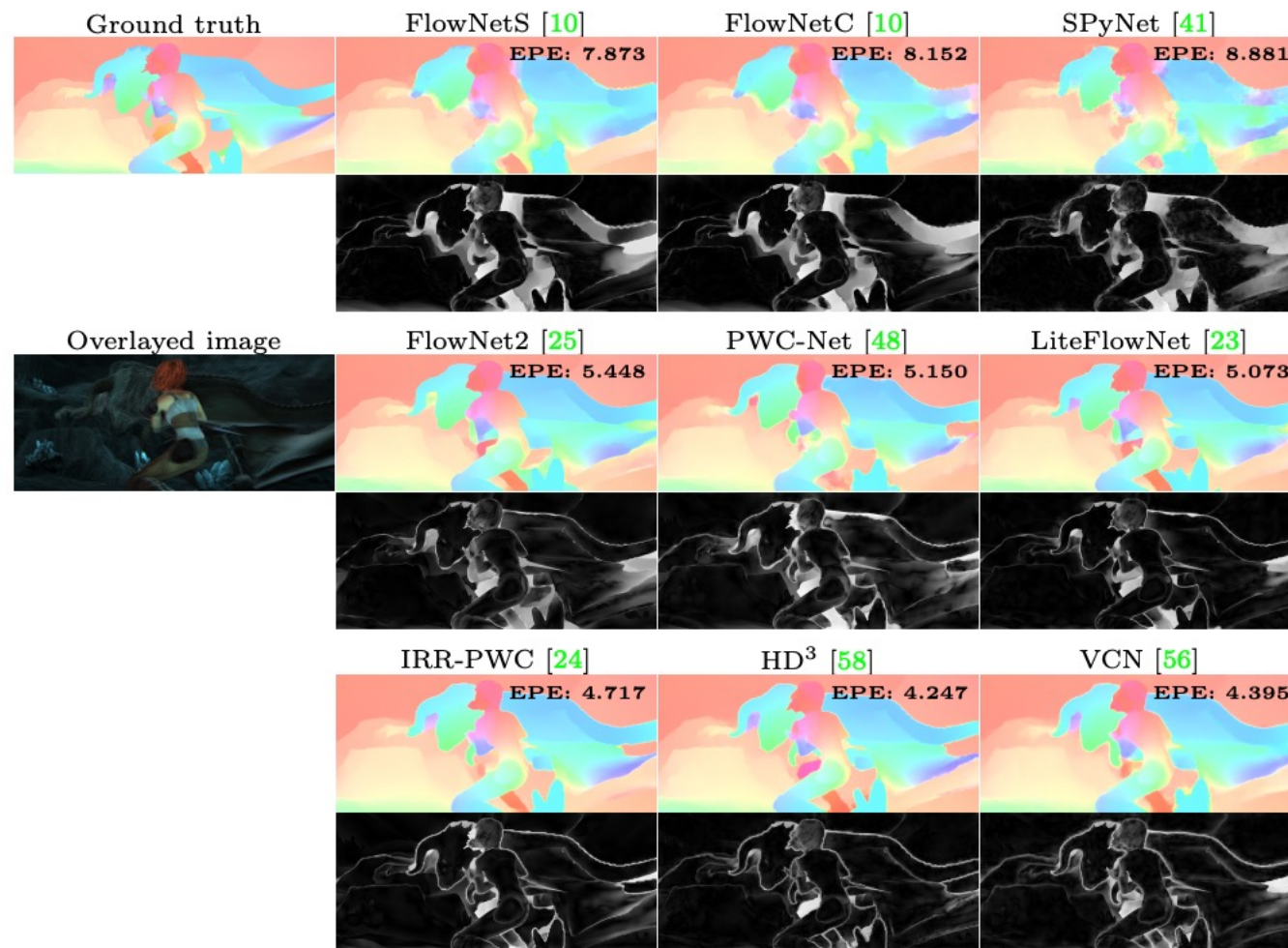


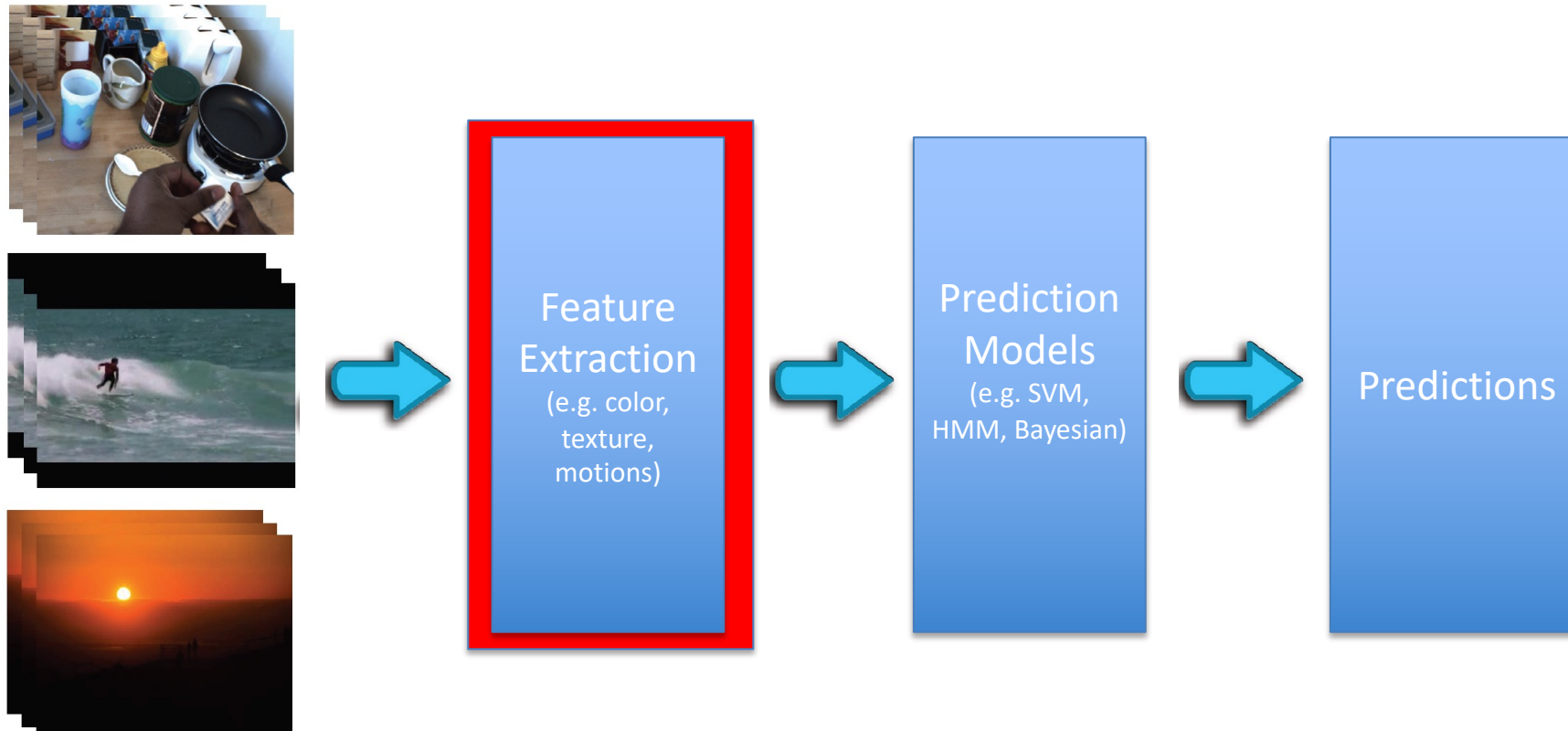
Fig. 3: **Qualitative comparison of end-to-end architectures:** Example from Sintel Final Test [8]. The first column shows the ground-truth flow and the overlaid input images. In the further columns, we show the color-coded flow visualization of each method, overlaid with the end point error (EPE) and their error maps (the brighter a pixel, the higher its error).

[Hur & Roth, 2020] Survey paper
<https://arxiv.org/pdf/2004.02853.pdf>

Overview

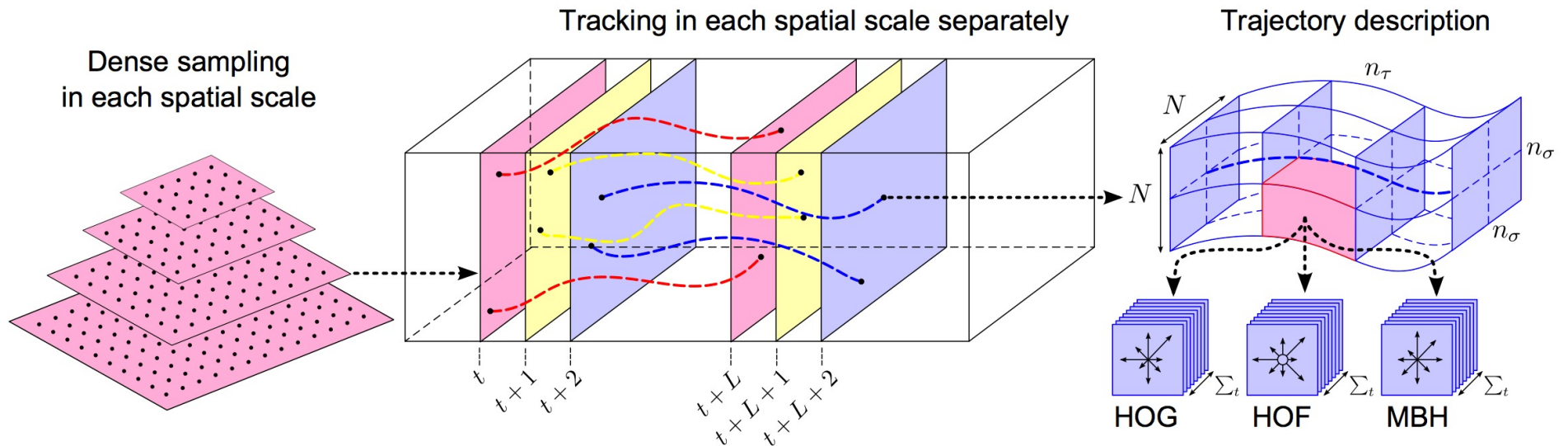
- Optical Flow
- ConvNets for Video

Traditional Computer Vision Pipeline



Best (non-DL) Video Features

- improved Dense Trajectories (iDT)



Wang et al. IJCV'13

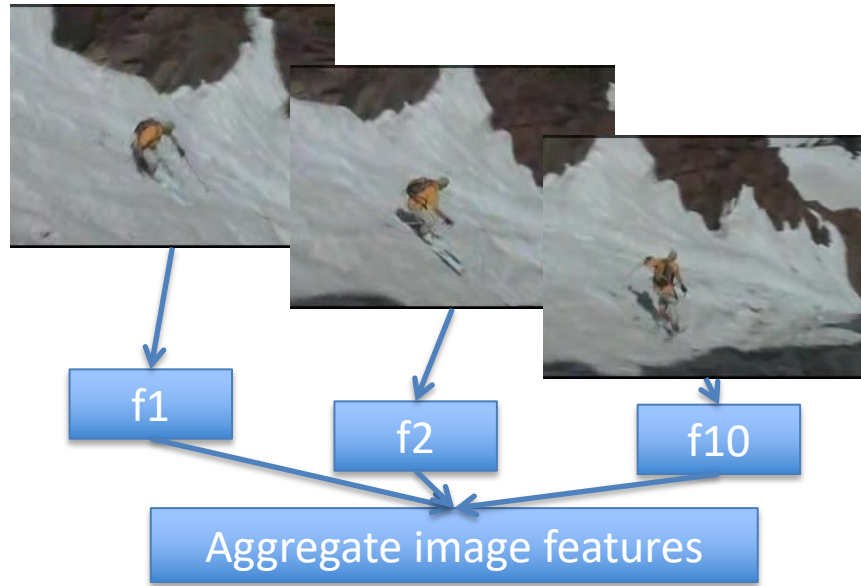
Pros:

- Don't need to learn
- Don't need large-scale training data

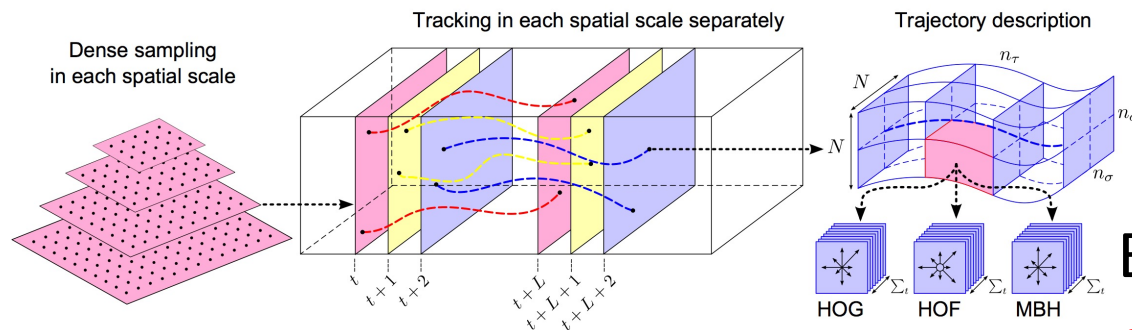
Cons:

- Highly hand-crafted
- Computational intensive
- Hard to parallelize

Spatiotemporal Feature Learning



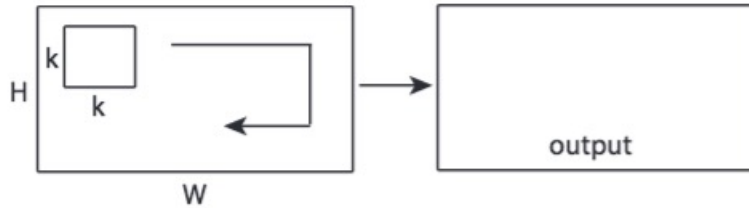
No explicit **motion modeling**



Biased to **human design** & **computationally expensive**

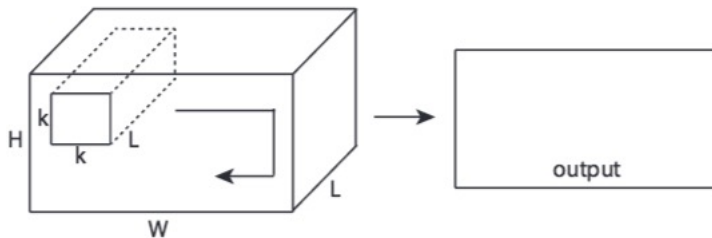
- Wang et al. IJCV'13

3D ConvNets



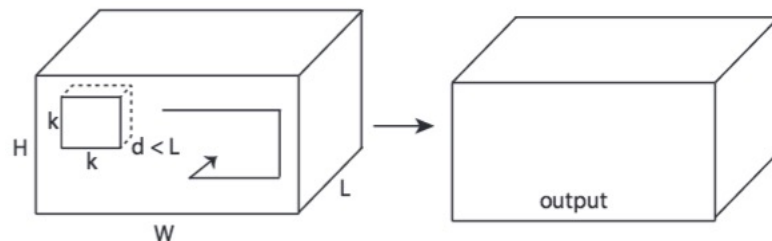
2D convolve on an image

-> no motion modeling



2D convolve on multiple images as channels

-> collapse temporal signal after one convolution layer



Spatial-temporally convolve on multiple frames

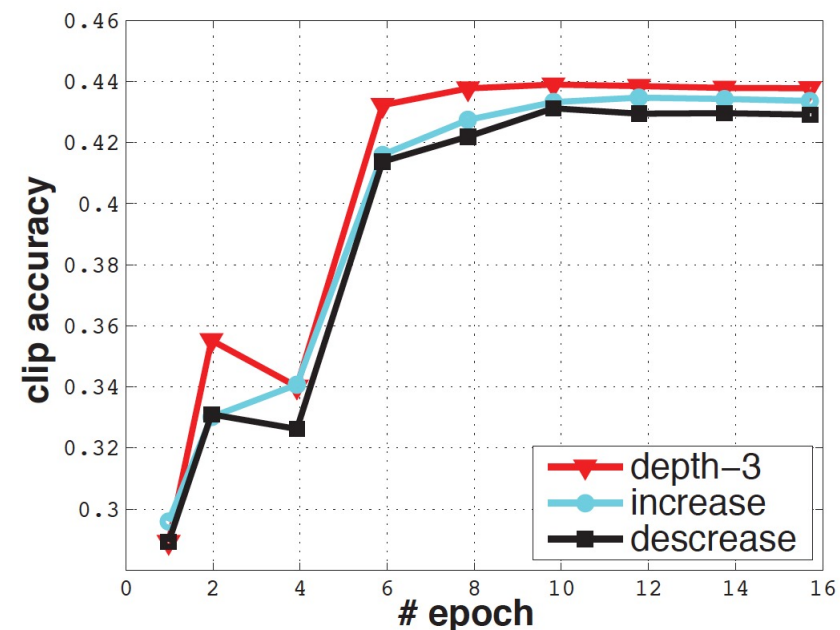
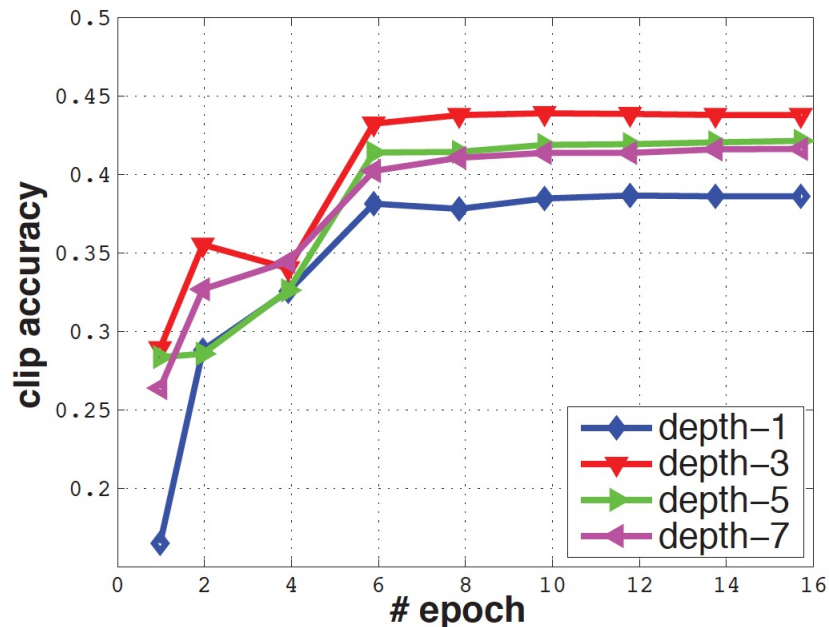
-> hierarchically group temporal signal



What is a Good Architecture for 3D ConvNets?

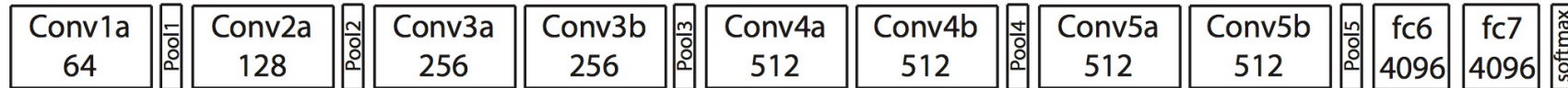
D. Tran, L. Bourdev, R. Fergus, L. Torresani, M. Paluri, *Learning Spatiotemporal Features with 3D Convolutional Networks*, **ICCV15**.

- Dataset: UCF101
- Use VGG-similar architecture, varying kernel temporal length



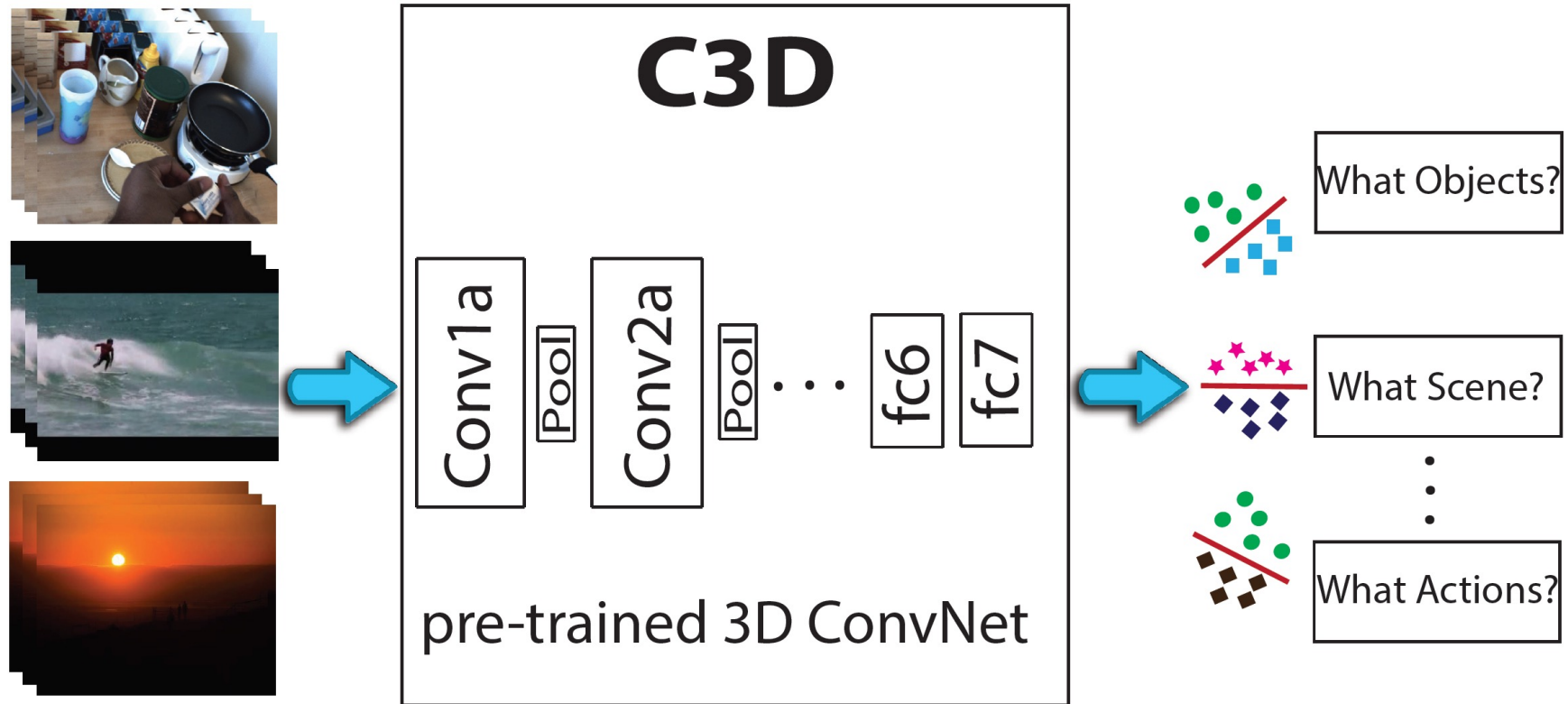
What is a Good Architecture for 3D ConvNets?

D. Tran, L. Bourdev, R. Fergus, L. Torresani, M. Paluri, *Learning Spatiotemporal Features with 3D Convolutional Networks*, **ICCV15**.



- C3D architecture
 - 8 convolution, 5 pool, 2 fully-connected layers
 - 3x3x3 convolution kernels
 - 2x2x2 pooling kernels
- Dataset: Sports-1M [Karpathy et al. CVPR14]
 - 1.1M videos of 487 different sport categories
 - Train/test splits are provided

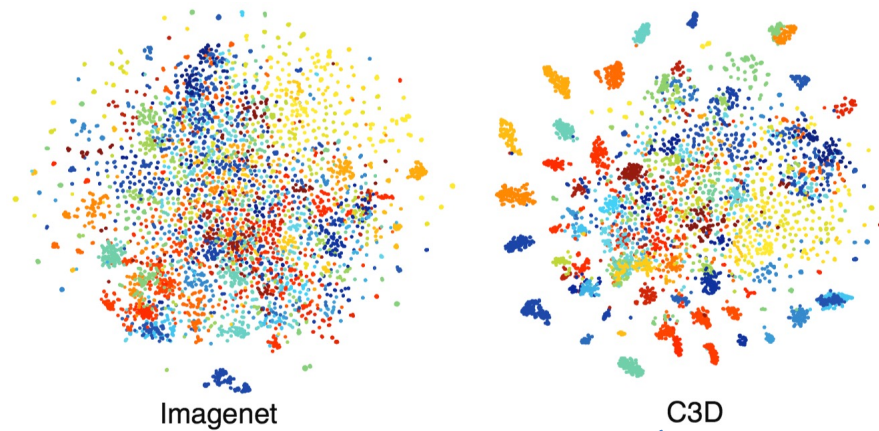
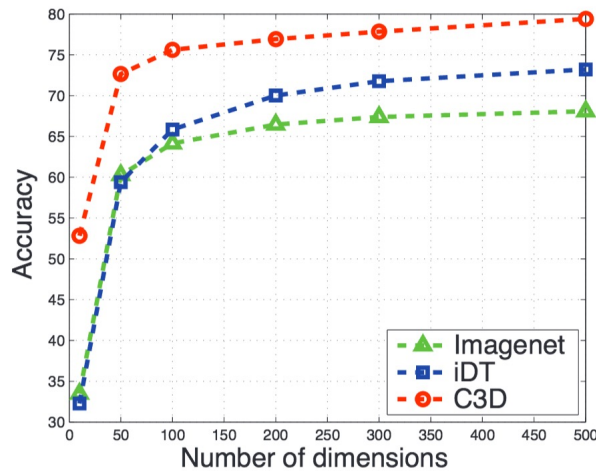
C3D as Generic Features



Simple recipe: C3D + linear SVM = good performance

Video Classification with C3D

Dataset Task	Sport1M action recognition	UCF101 action recognition	ASLAN action similarity labeling	YUPENN scene classification	UMD scene classification	Object object recognition
Method	[19]	[39]([26])	[31]	[10]	[10]	[32]
Result	80.2	75.8 (89.1)	68.7	96.2	77.7	12.0
C3D	85.2	85.2 (90.4)	78.3	98.1	87.7	22.3
Δ	5.0	9.4 (1.3)	9.6	1.9	10.0	10.3



C3D is discriminative and compact!

C3D code/model is publicly available

Action Recognition Task

UCF101



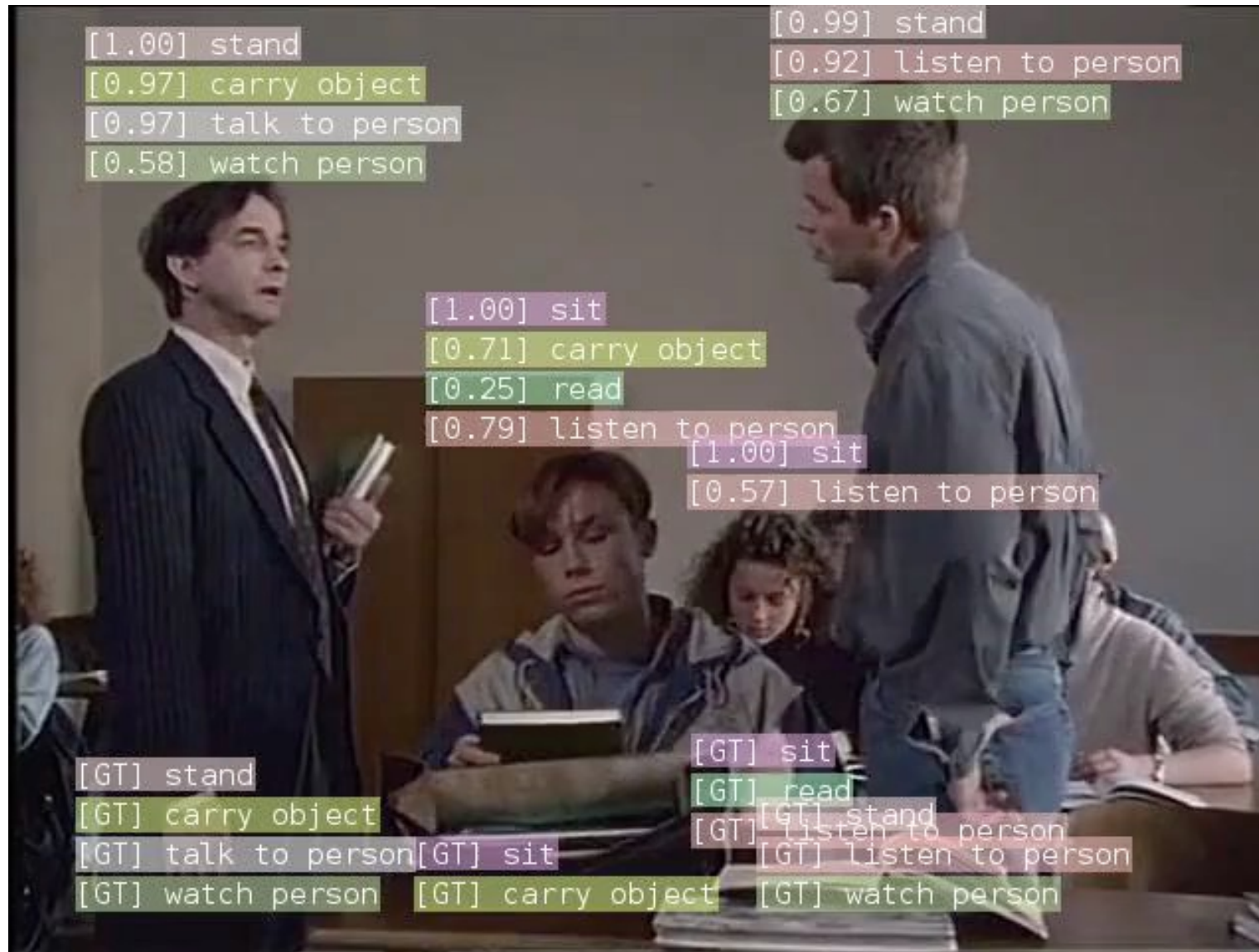
Recognition in Video

ICCV 2019 Tutorial

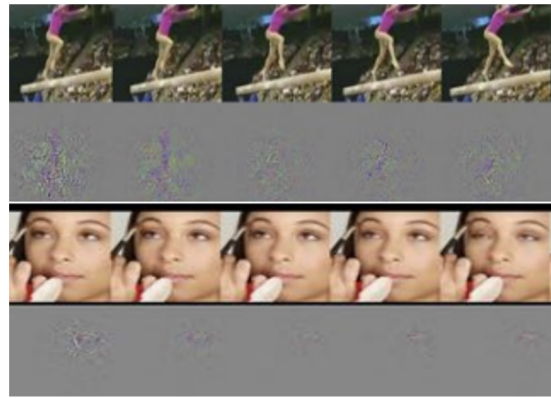
Christoph Feichtenhofer

Facebook AI Research (FAIR)

Task: Human action classification & detection



Outline: Components for state-of-the-art video understanding

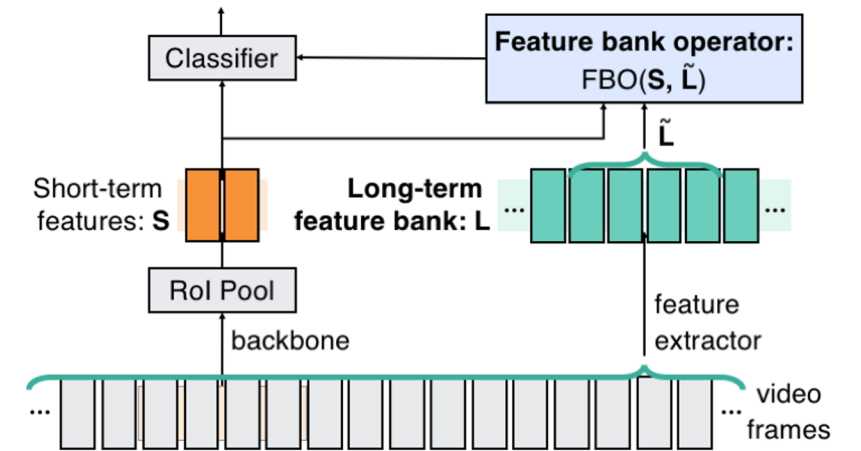


3D ConvNets

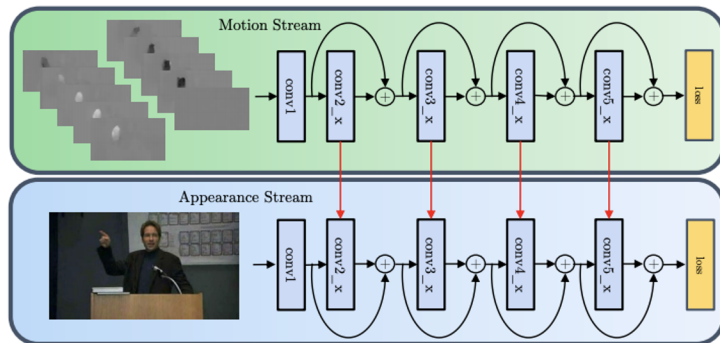
[Taylor et al. 2010, Karpathy et al. 2014, Tran et al. 2015,...]



Attention-based models, Non-local network blocks, [Wang et al., 2018 2019, Girdhar et al. 2019 ,...]

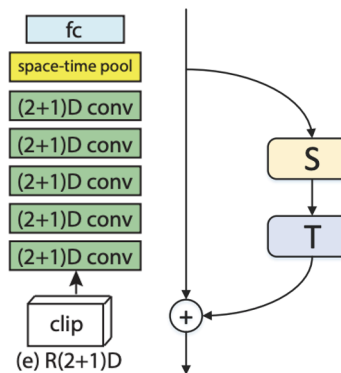


Long-term Models [Varol et al. 2017, Wu et al. 2019, ...]

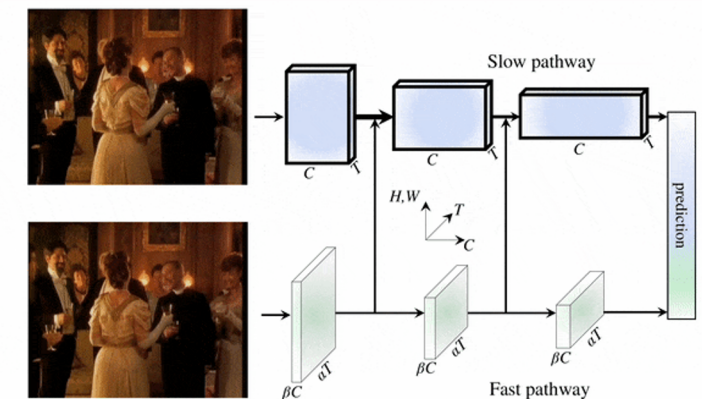


Two-stream ConvNets (RGB+optical flow)

[Simonyan et al. 2014, Feichtenhofer et al. 2016, Wang et al. 2016, ...]

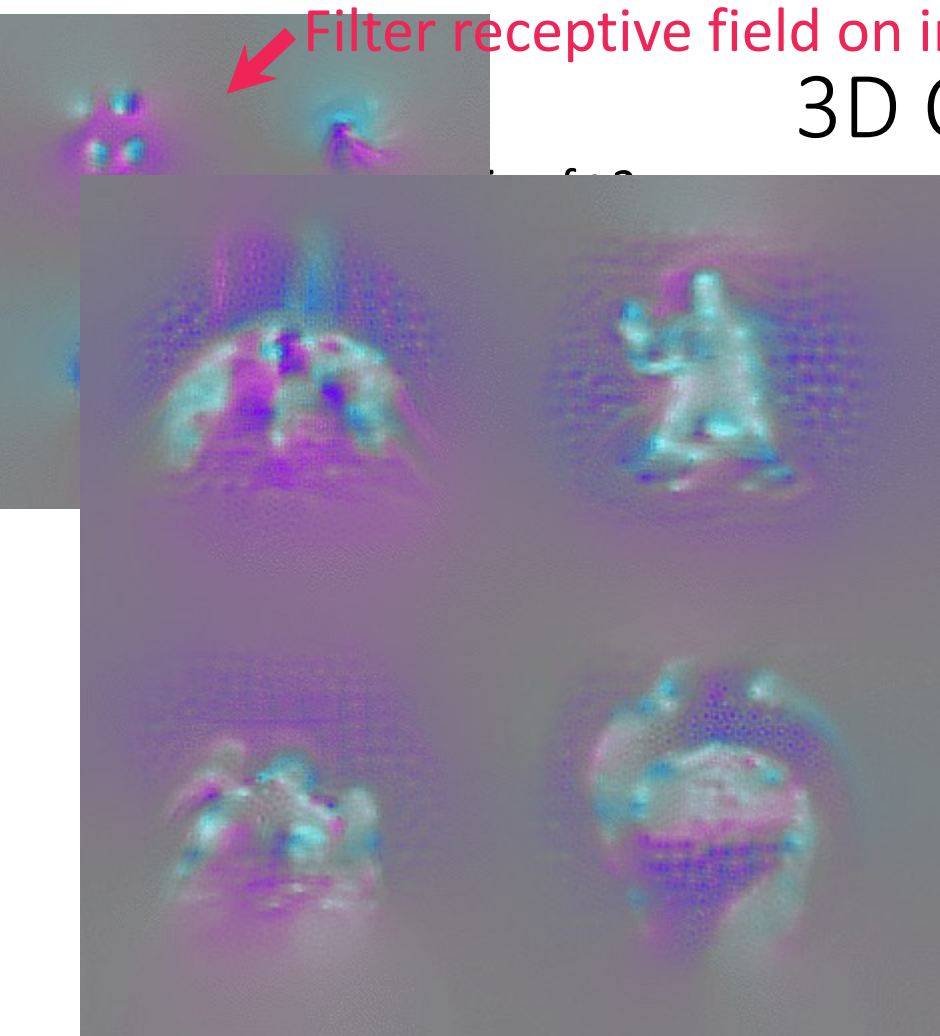


Local decomposition spatial and temporal information [Feichtenhofer et al. 2016, Qiu et al. 2017, Tran et al. 2018, Xie et al. 2018, ...]

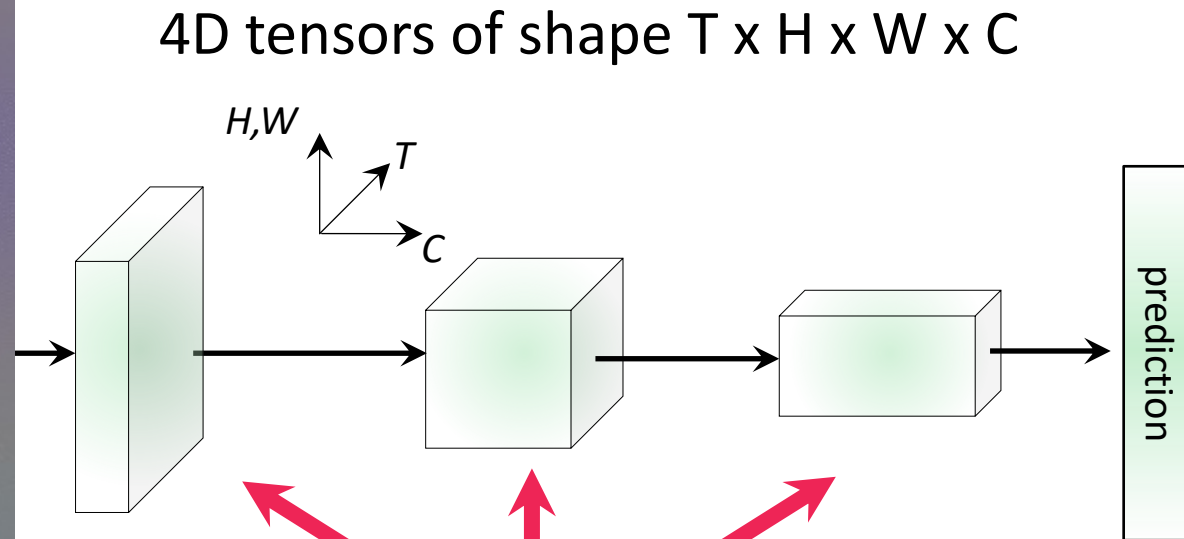


Global decomposition spatial and temporal information
SlowFast networks contrast features of different framerate and channel capacity [Feichtenhofer et al. 2019]

3D Convolutional Networks



(Kinetics classification annotation)



Intermediate filters only capture *local* information (in x, y, t) with a growing receptive field size

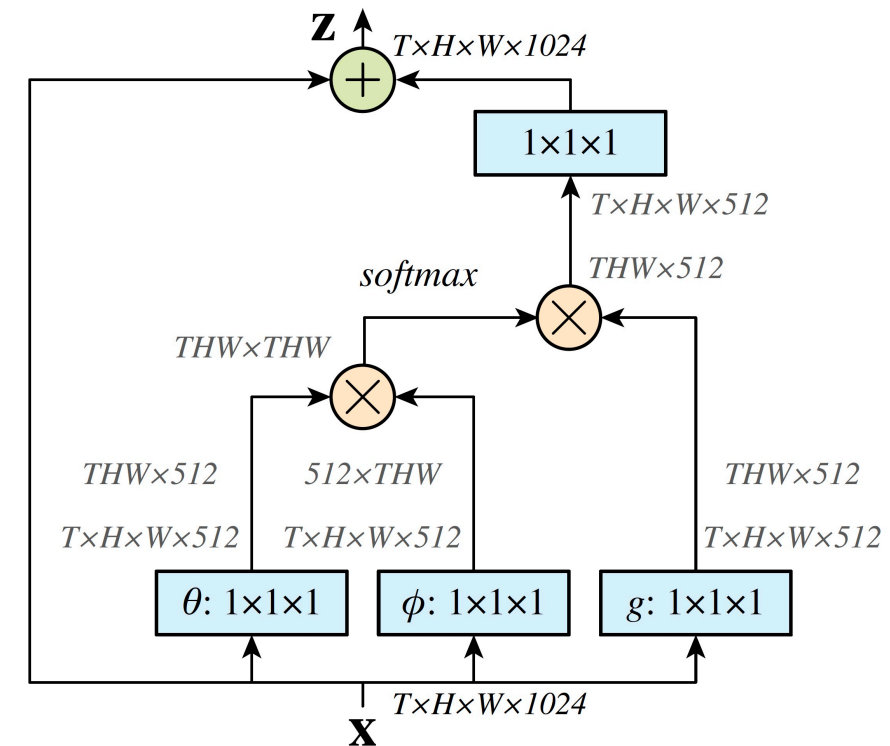
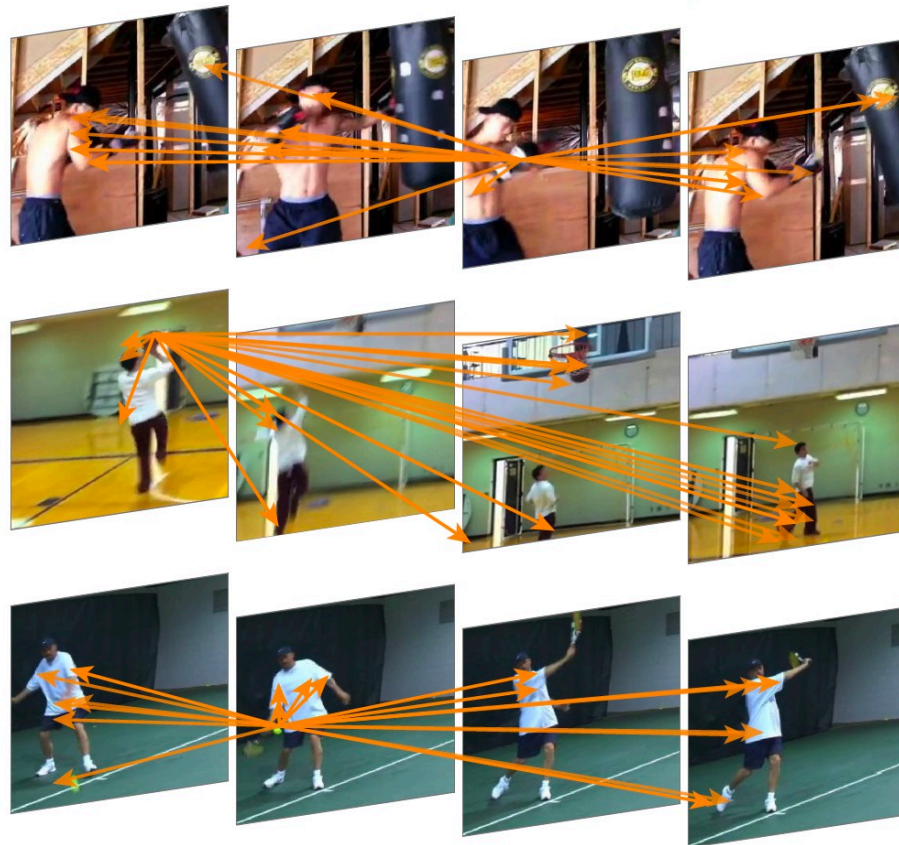
G. W. Taylor, R. Fergus, Y. LeCun, and C. Bregler. Convolutional learning of spatio-temporal features. In Proc. ECCV, 2010.

D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3D convolutional networks. In Proc. ICCV, 2015.

J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In Proc. CVPR, 2017.

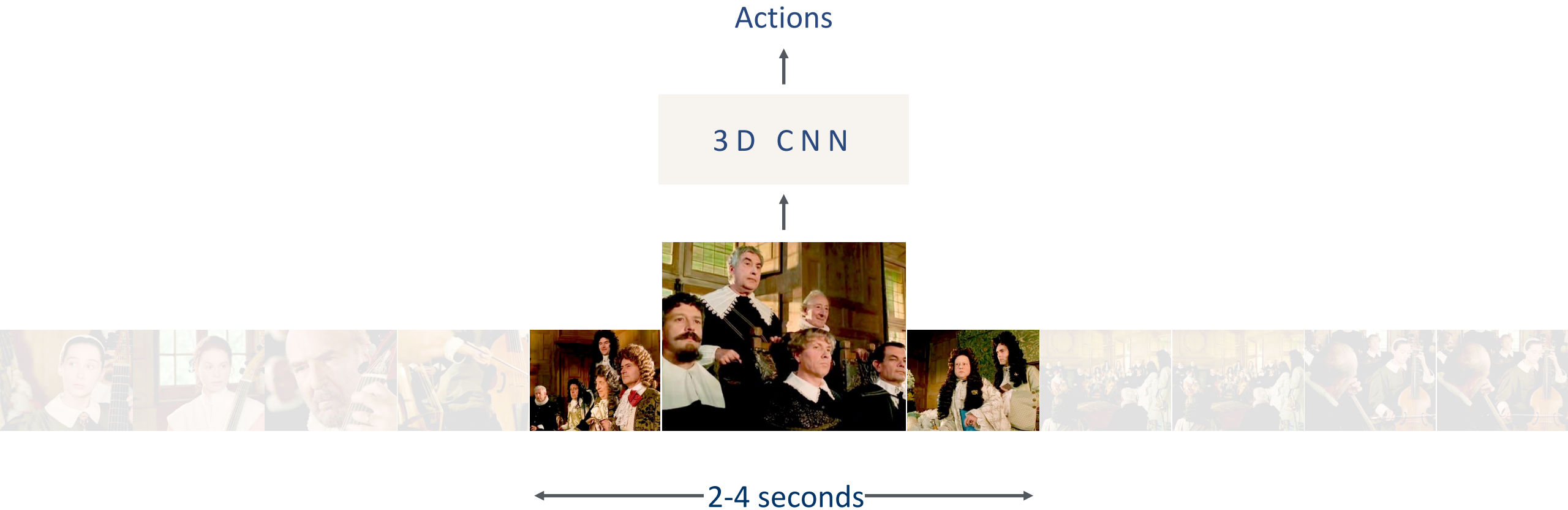
Pytorch code now available:

Non-Local Blocks

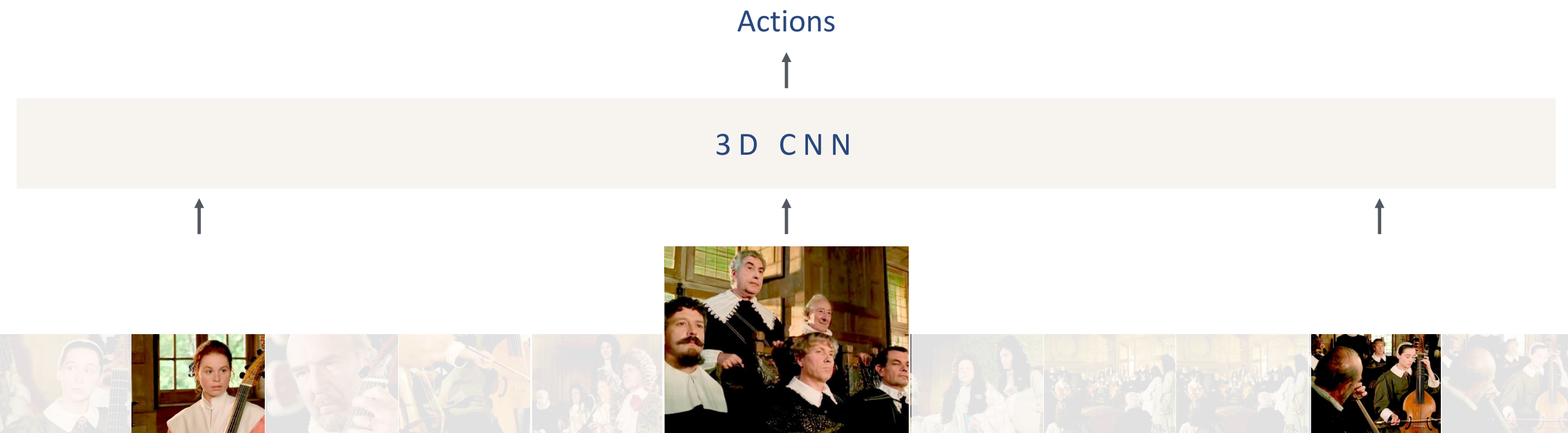
<https://github.com/facebookresearch/SlowFast>

○ Self-attention in the spatiotemporal domain allows long-range feature aggregation

Limited temporal input length of 3D ConvNets

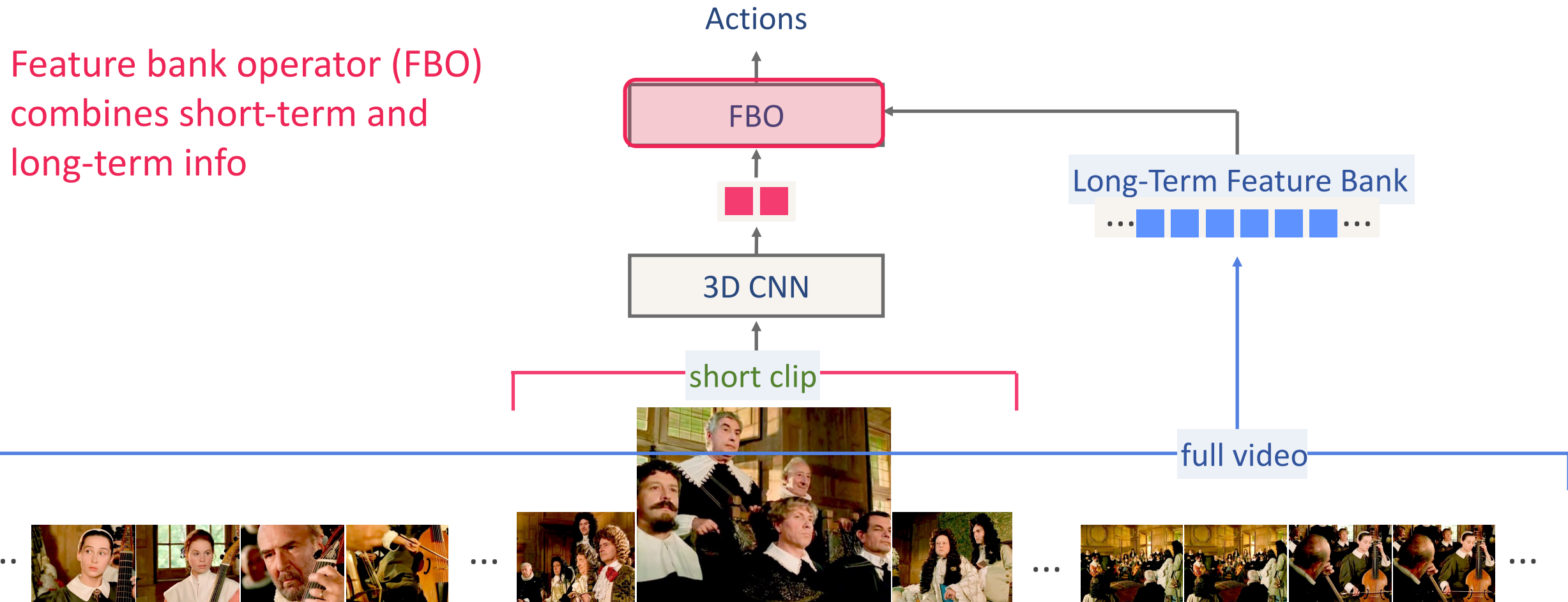


Temporal striding (subsampling)



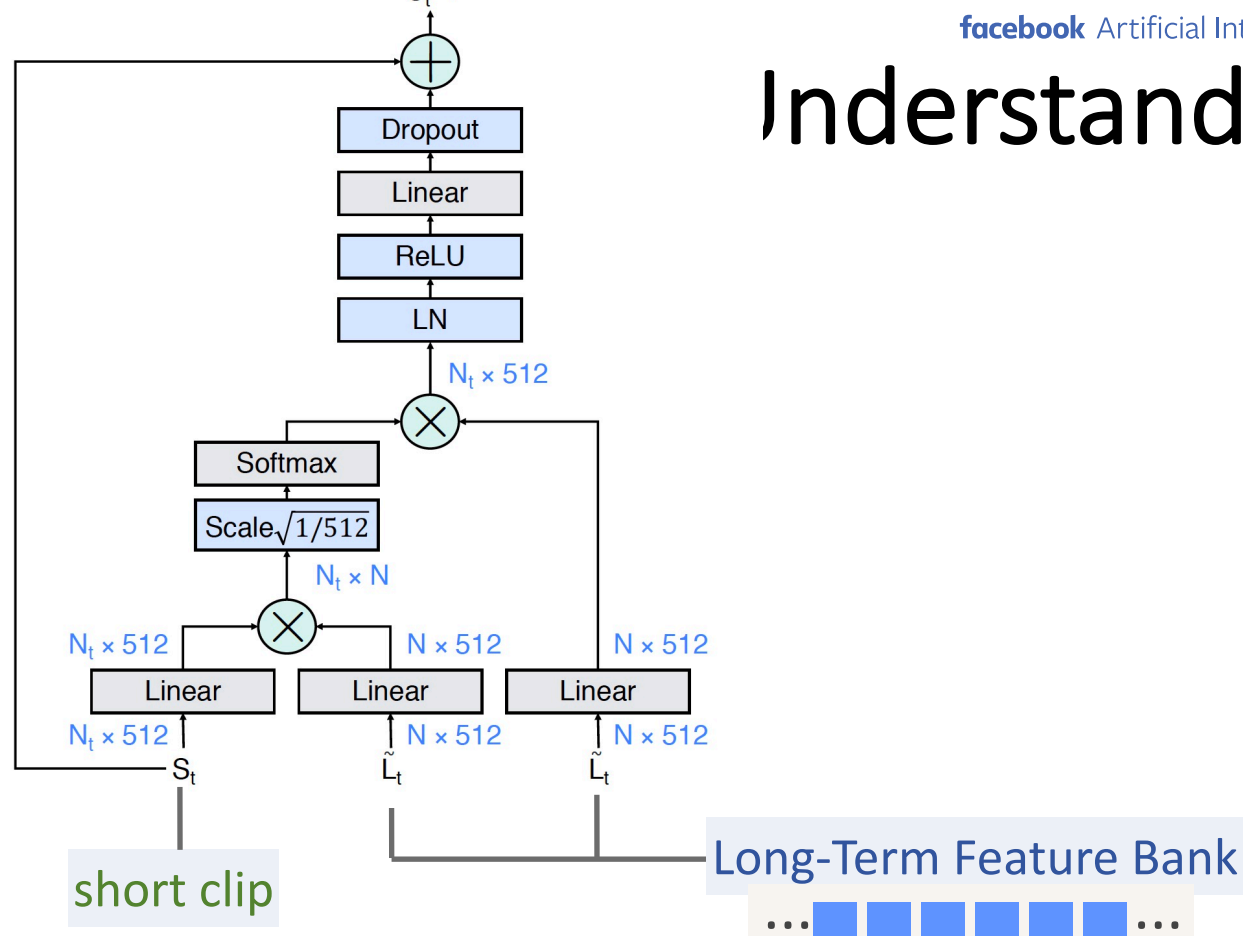
Long-Term Feature Banks for Video Understanding

Feature bank operator (FBO) combines short-term and long-term info



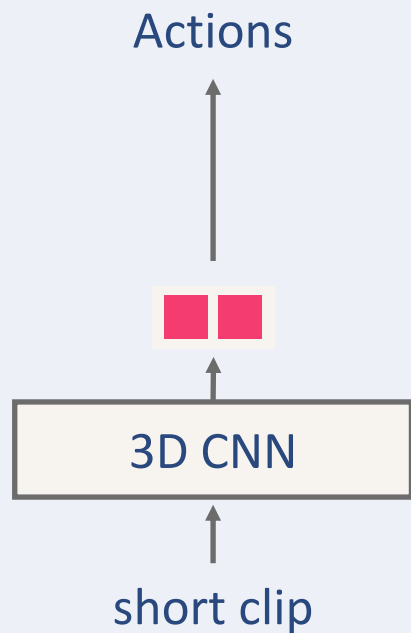
Long-Term Feature Bank

Understanding

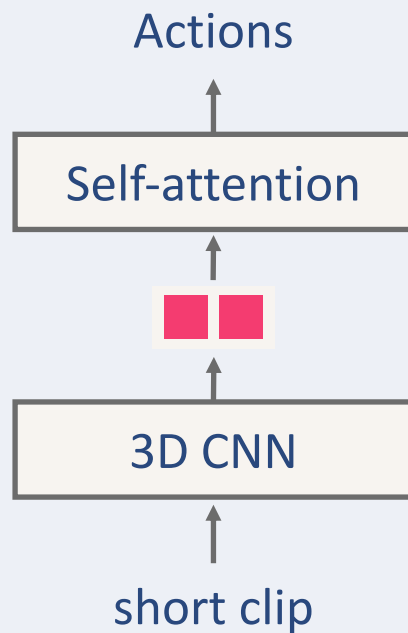


full video

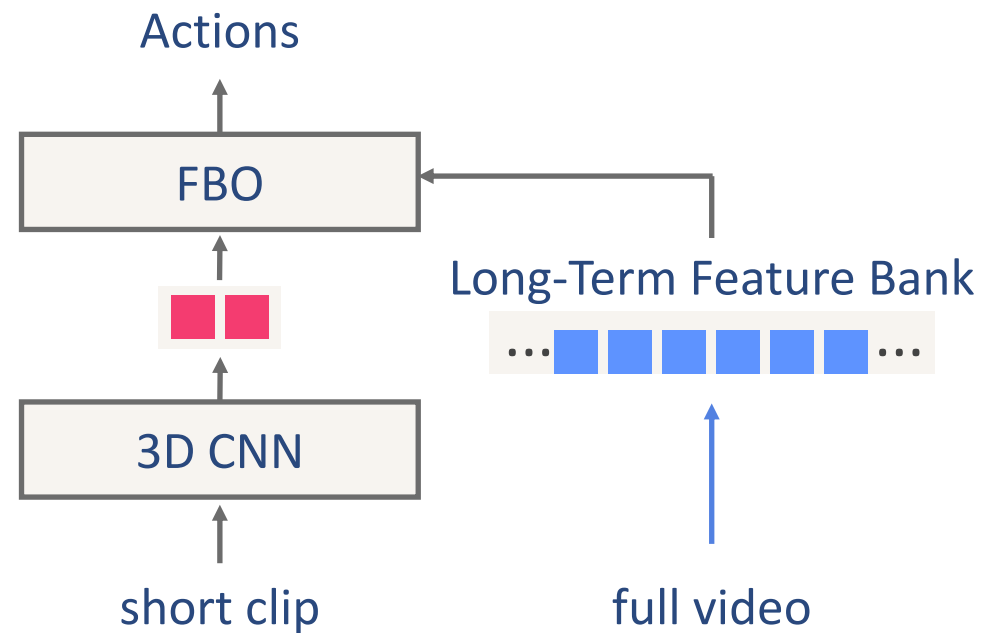
Baseline 1: 3D CNN x2



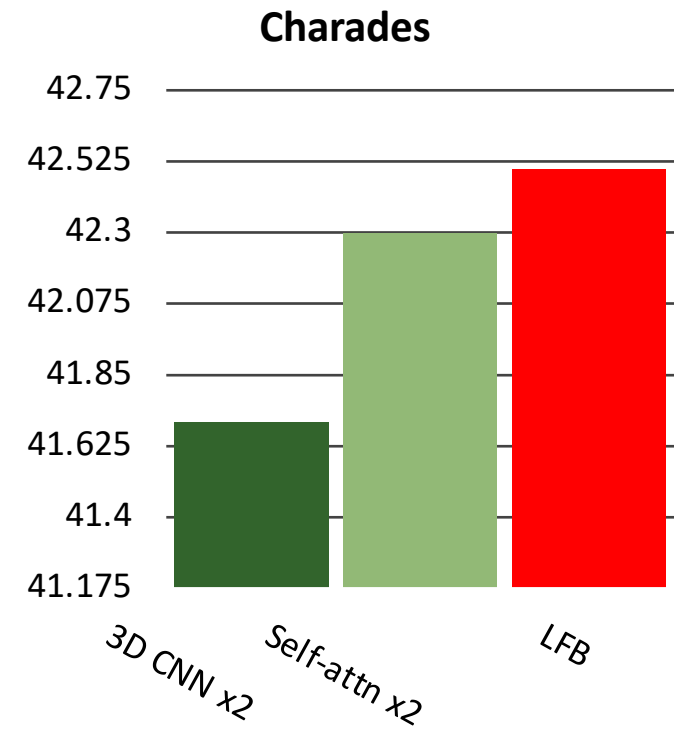
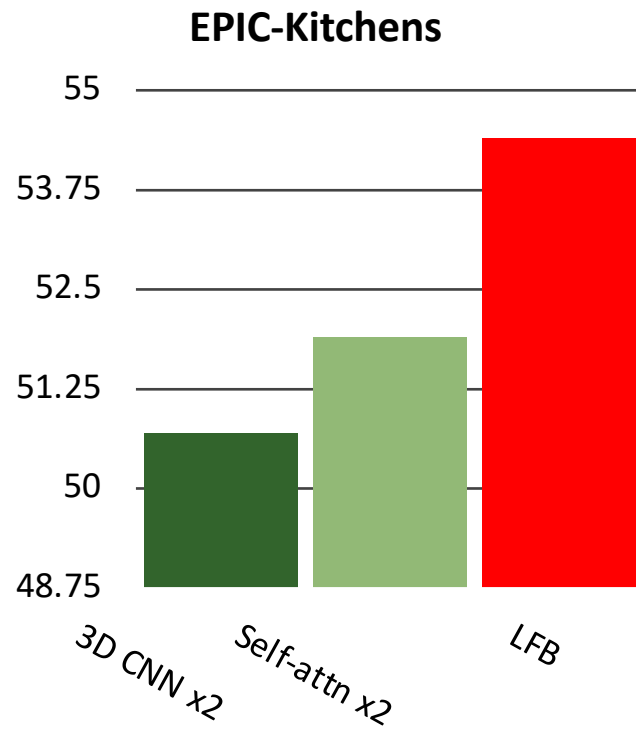
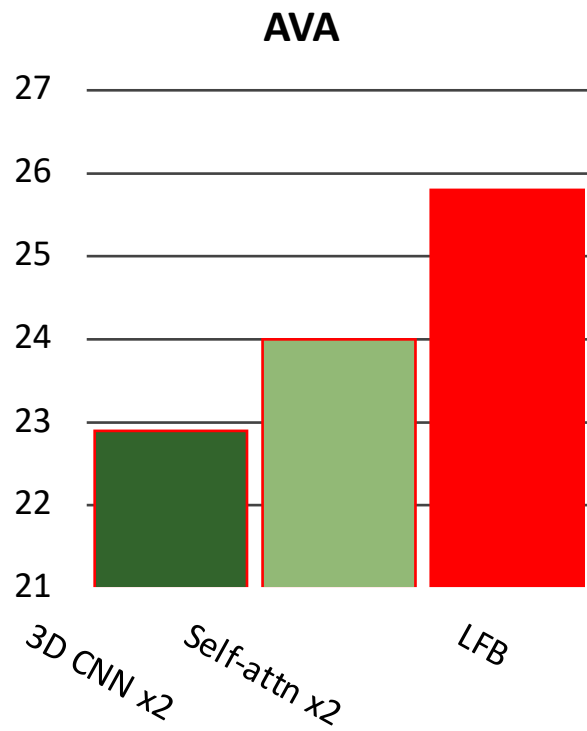
Baseline 2: Self-attention x2



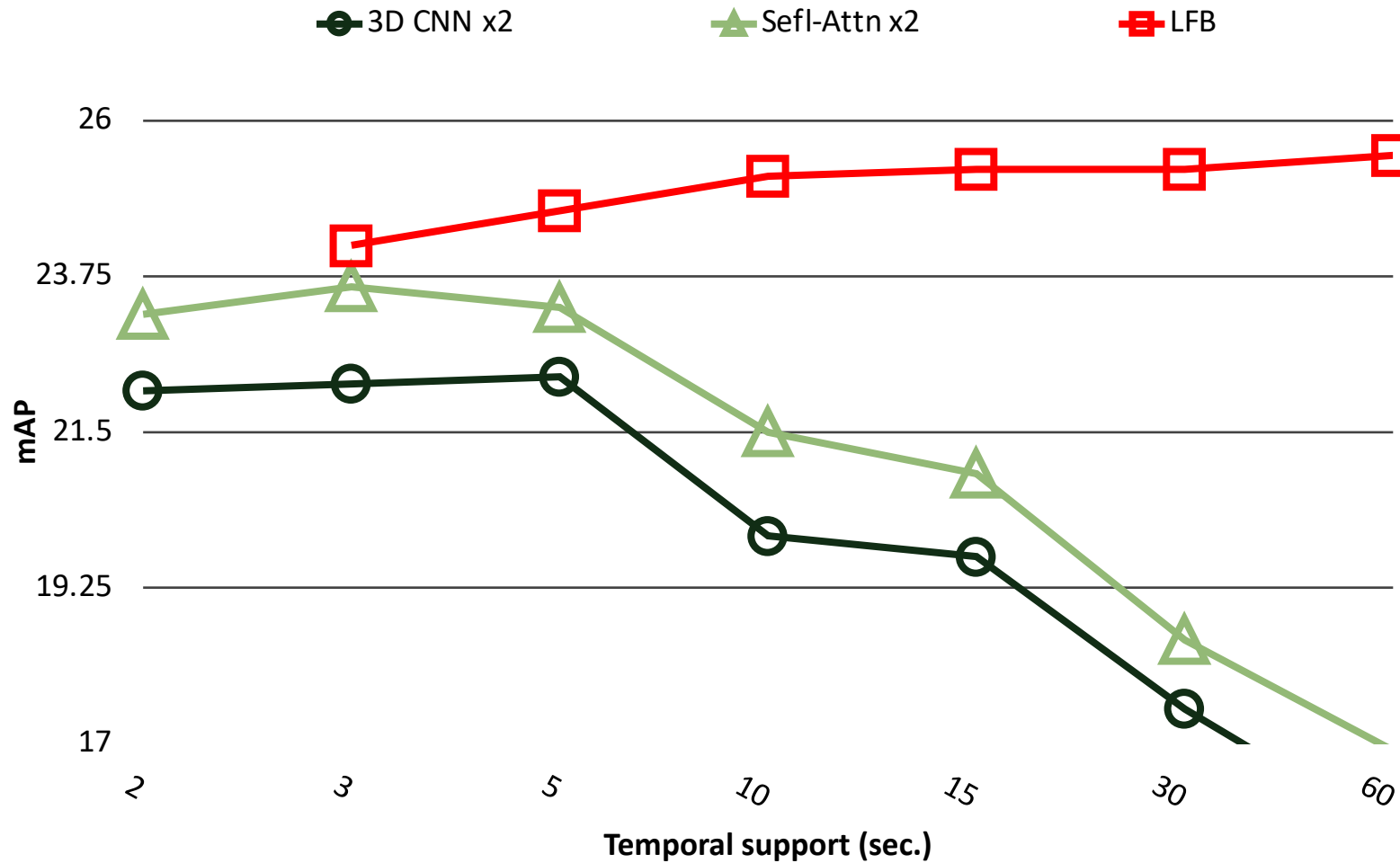
LFB



Ablation study on short-term vs. long-term



Ablation on input duration: subsampling vs LFB



Code/models:

<https://github.com/facebookresearch/video-long-term-feature-banks>

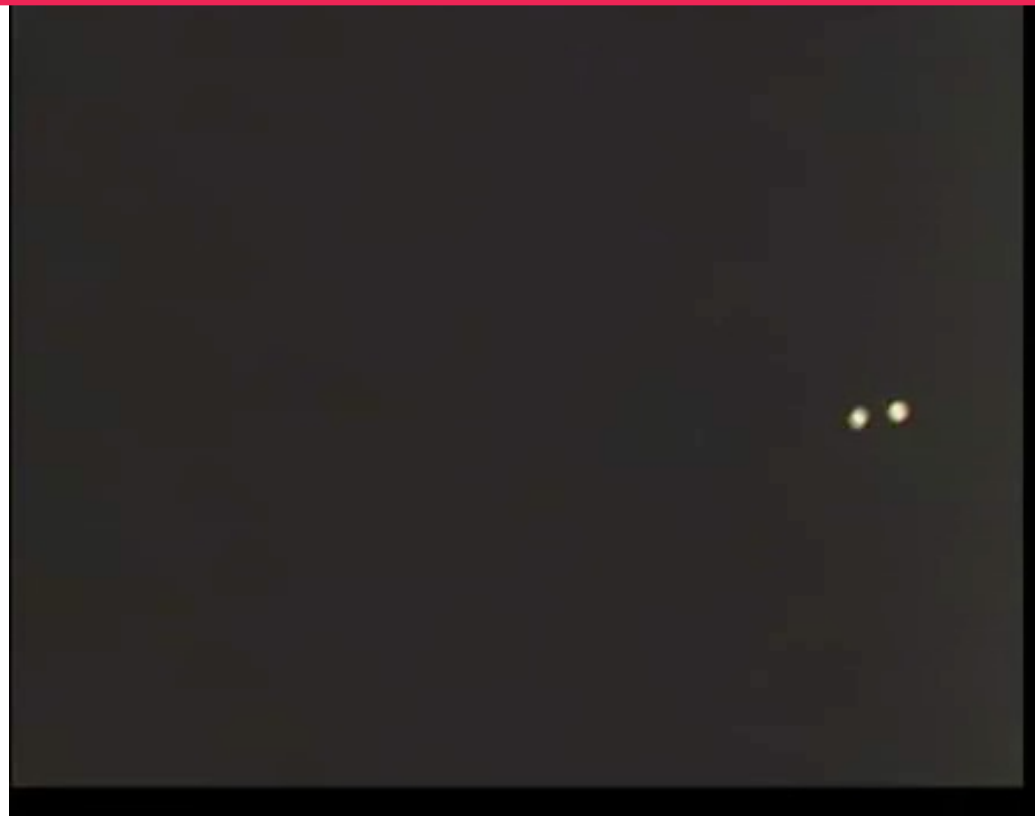


P(holding an object)?

Johansson

→ Amazing what a human brain can do
without appearance information

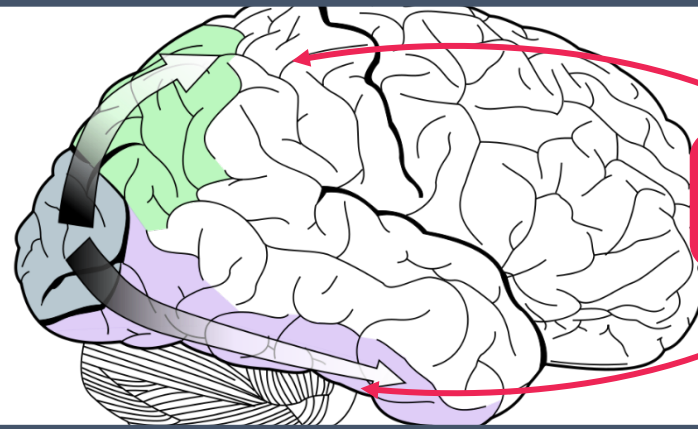
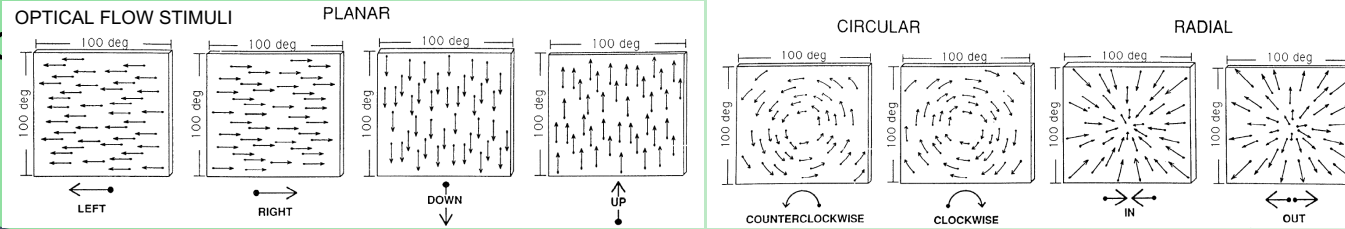
otion



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

Motivation: Separate visual pathways in nature

→ Dorsal stream ('where') recognizes motion and locates



→ "Interconnection"
e.g. in STS area

→ Ventral ('what') stream performs object recognition



Sources: "Sensitivity of MST neurons to optic flow stimuli. I. A continuum of response selectivity to large-field stimuli." *Journal of neurophysiology* 65.6 (1991).

"A cortical representation of the local visual environment", *Nature*. 392 (6676): 598–601, 2009

https://en.wikipedia.org/wiki/Two-streams_hypothesis

Two-Stream Convolutional Networks

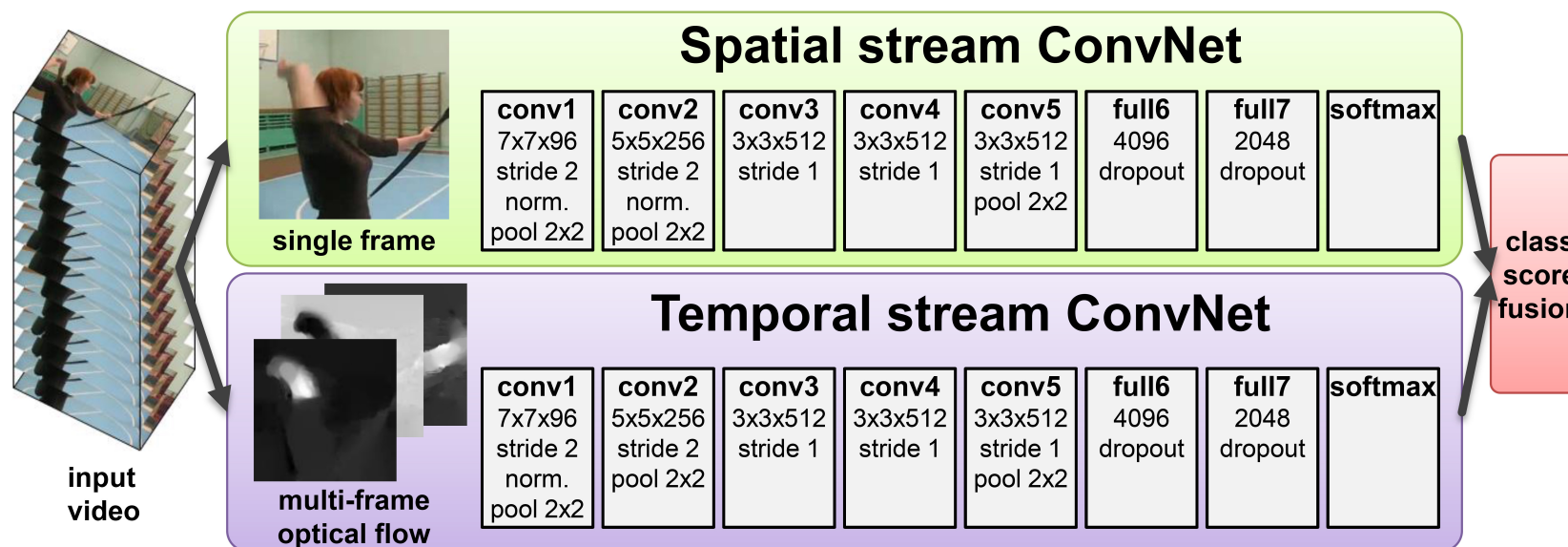
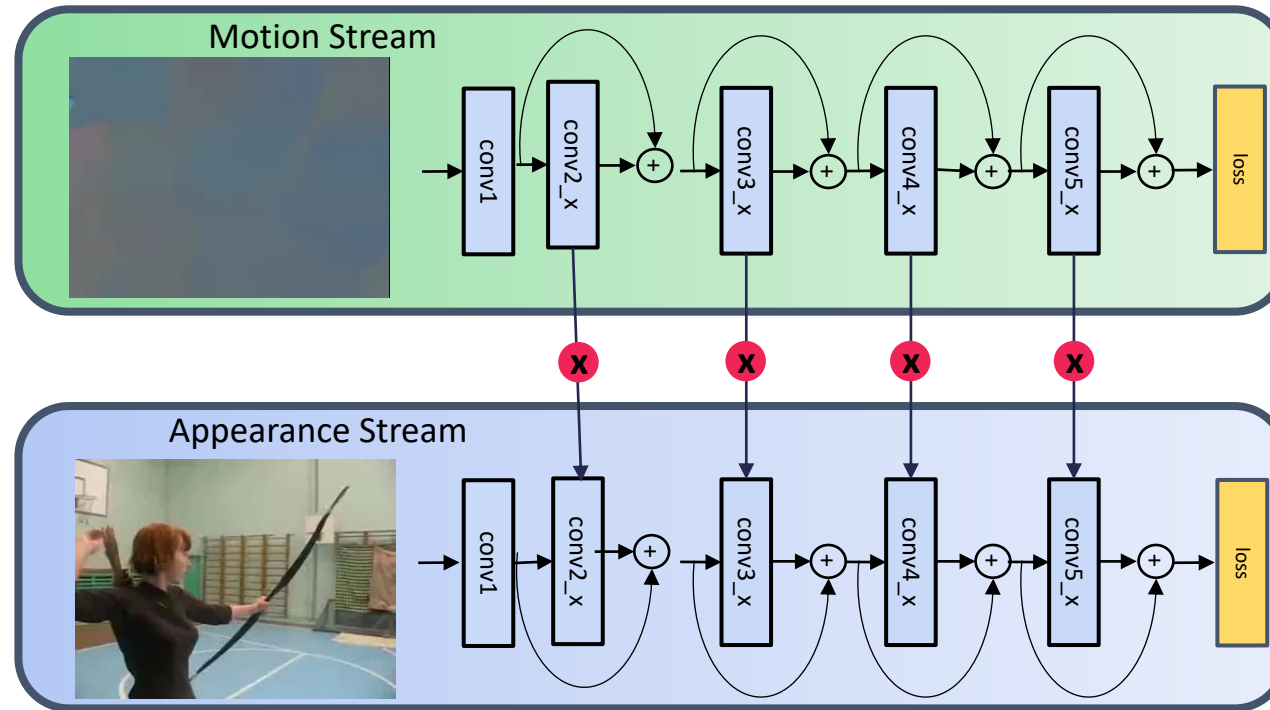


Figure 1: **Two-stream architecture for video classification.**

Individual processing of spatial and temporal information

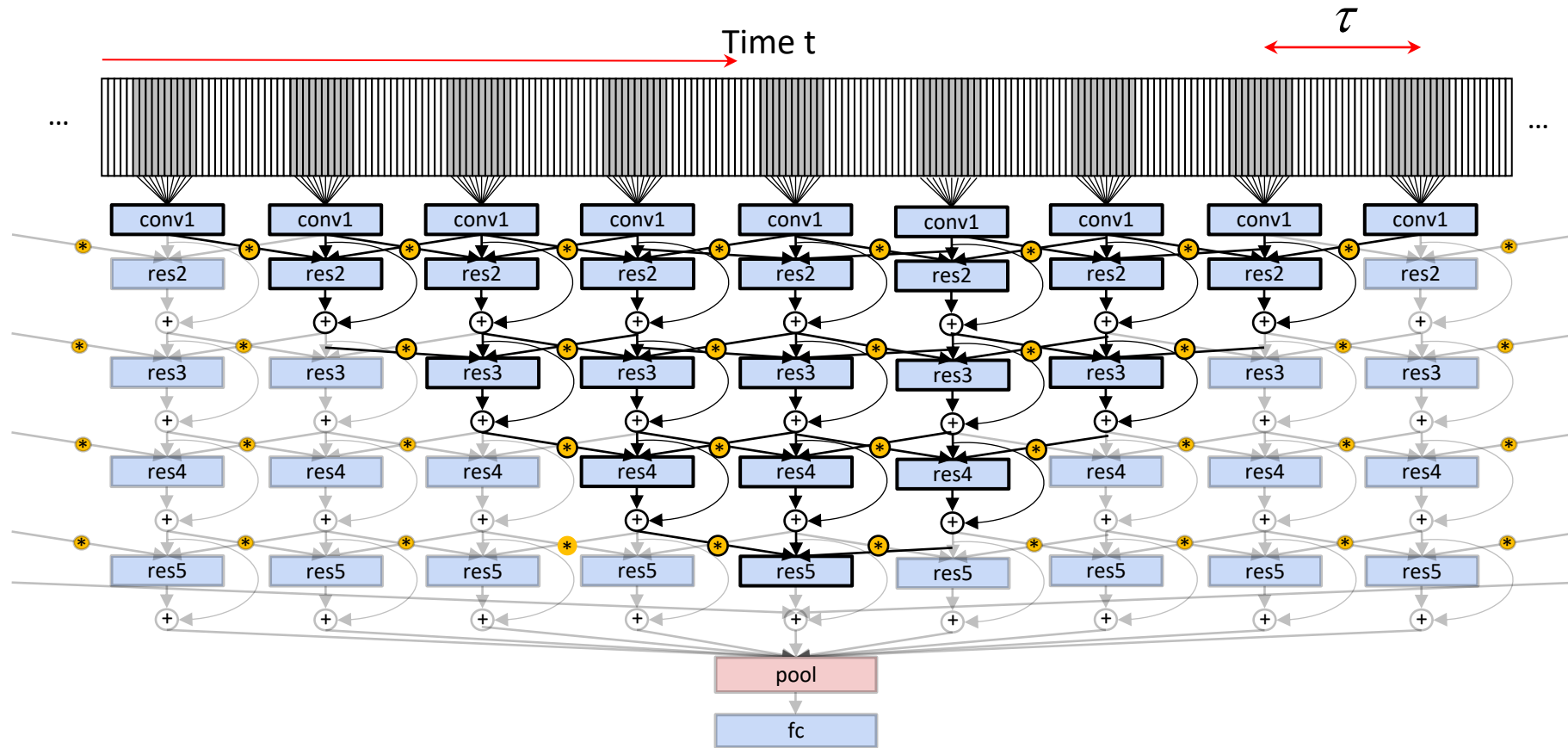
- Using a separate **2D (x,y)** ConvNet recognition stream for each
- Late fusion via softmax score averaging

Two-Stream Network Fusion and Long-term Two-Stream networks



- ST-ResNet allows the hierarchical learning of spacetime features by **connecting** the appearance and motion channels of a two-stream architecture.

Long-term Two-Stream networks and transforming filters by Inflation



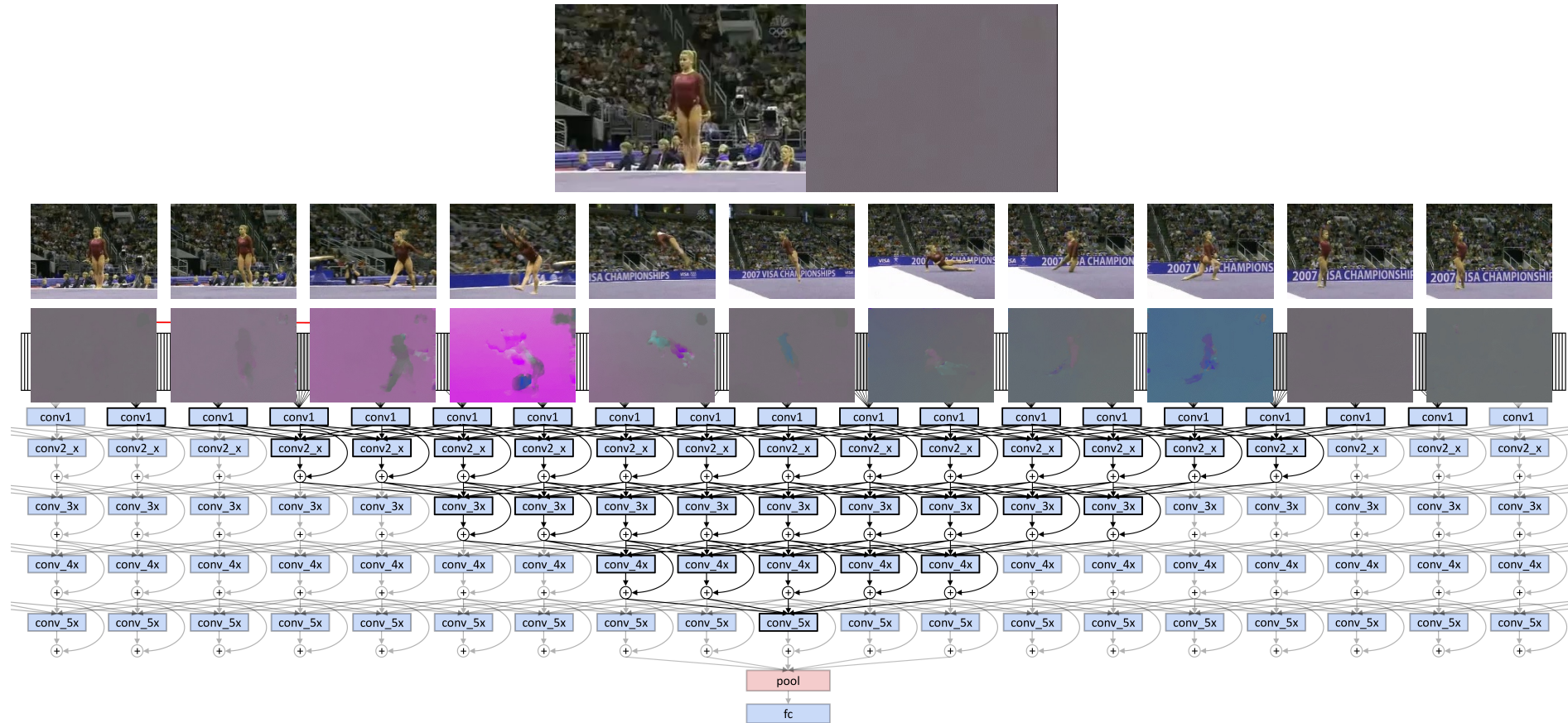
- Inflation allows to transform spatial filters to spatiotemporal ones (3D or 2D spatial +1D temporal)

C. Feichtenhofer, A. Pinz, and R. Wildes. Spatiotemporal residual networks for video action recognition. In NIPS, 2016.

Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X. and Van Gool, L., Temporal segment networks: Towards good practices for deep action recognition. ECCV 2016

J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In Proc. CVPR, 2017.

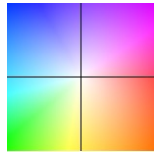
Two-Stream Network Fusion and Long-term Two-Stream networks



C. Feichtenhofer, A. Pinz, and R. Wildes. Spatiotemporal residual networks for video action recognition. In NIPS, 2016.

Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X. and Van Gool, L., Temporal segment networks: Towards good practices for deep action recognition. ECCV 2016

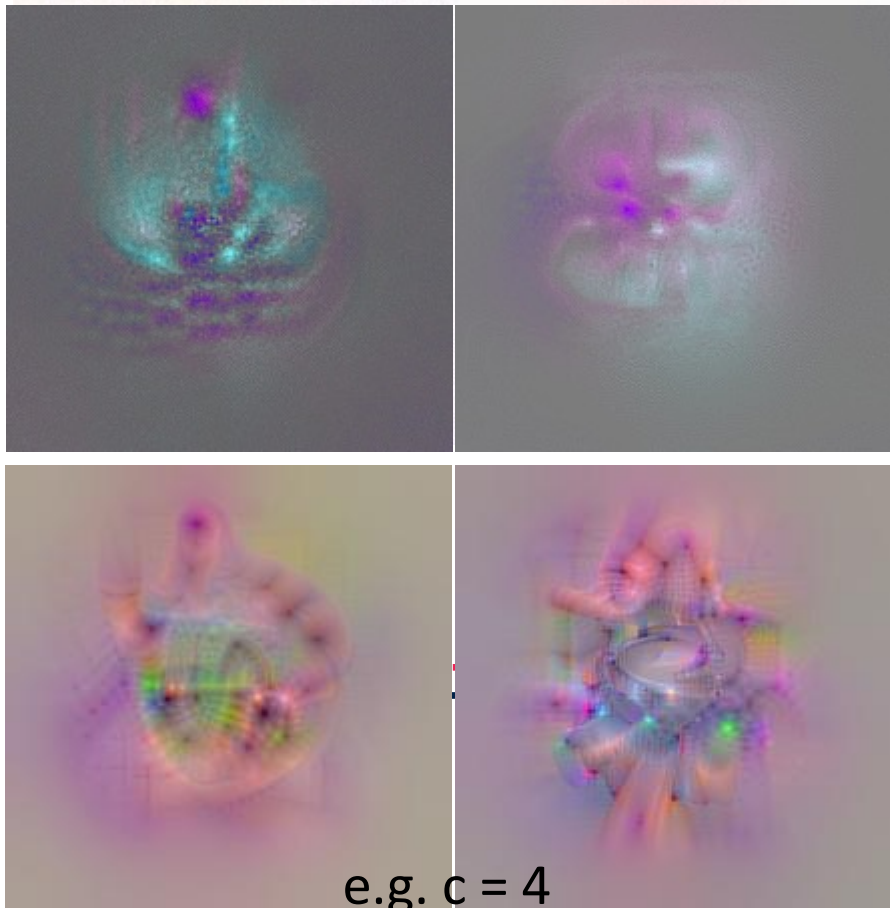
J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In Proc. CVPR, 2017.



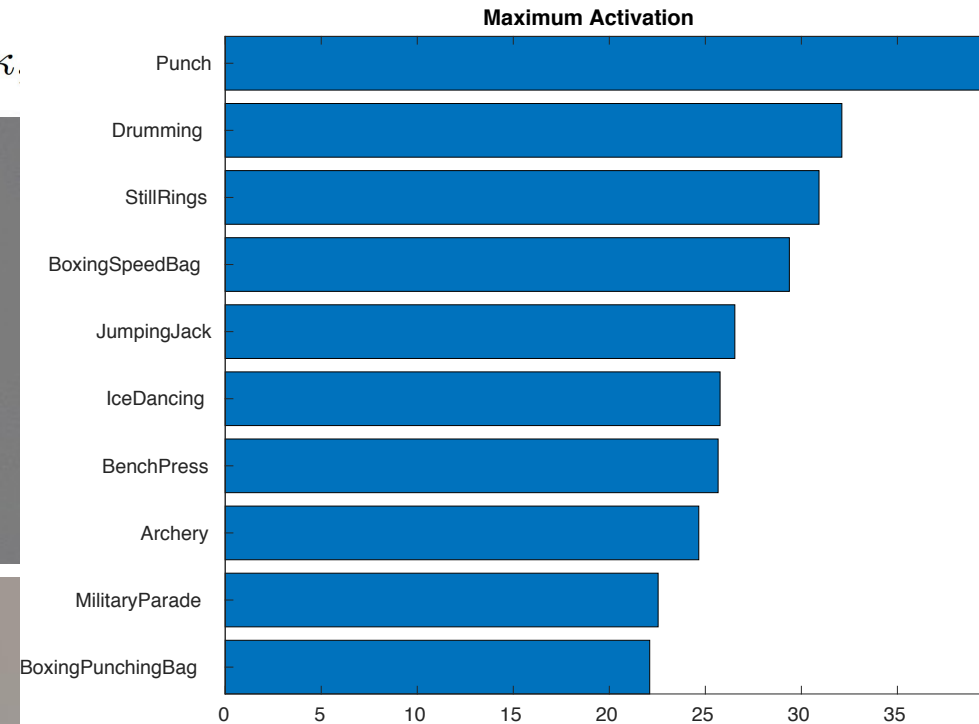
Visualizing the learned representation

Slow motion
(high temporal reg.)

Fast motion
(low temporal reg.)



e.g. $c = 4$



C. Feichtenhofer, A. Pinz, and R. Wildes, A. Zisserman. What have we learned from video recognition?. In CVPR, 2018.

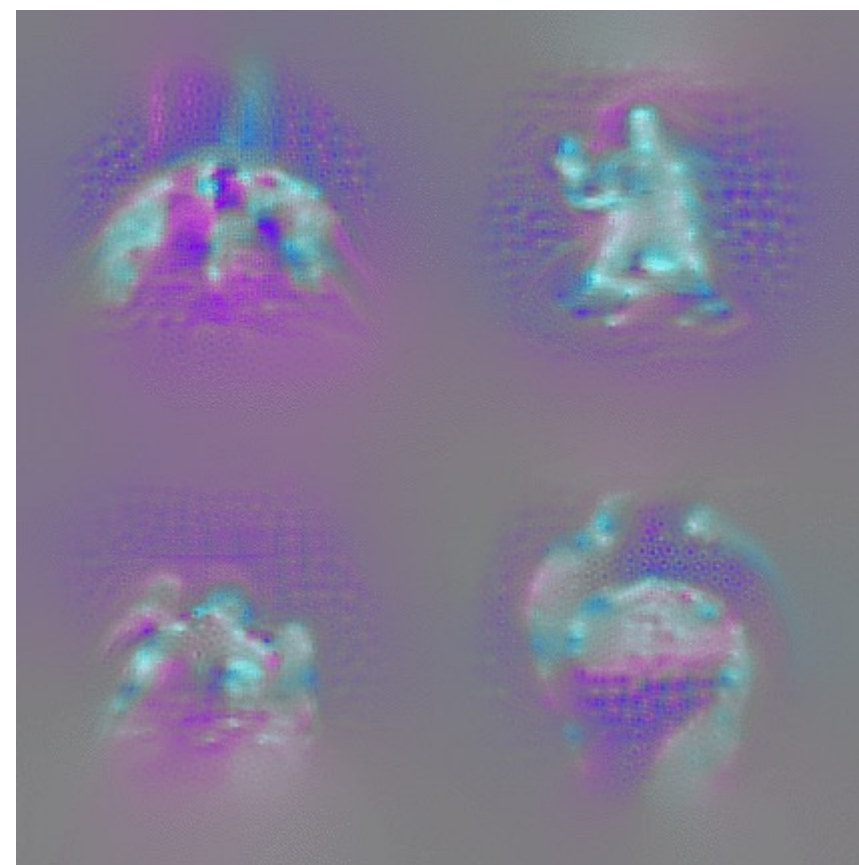
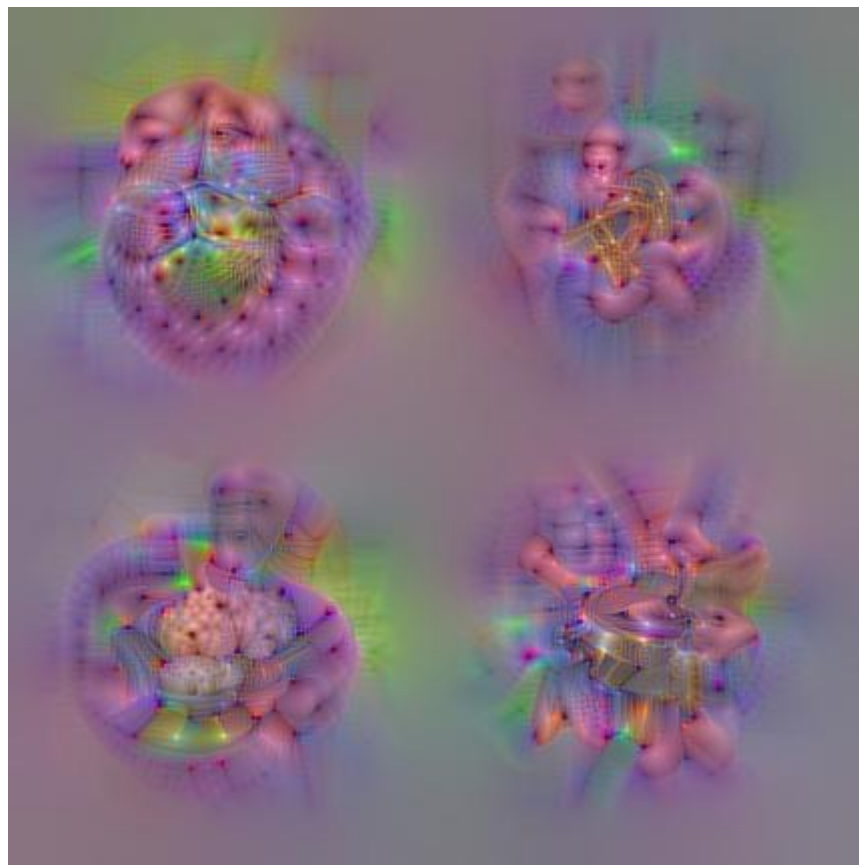
C. Feichtenhofer, A. Pinz, and R. Wildes, A. Zisserman. Deep insights into convolutional networks for video recognition?. In IJCV, 2019.

Going through the conv layers of VGG-16 (first four filters of each layer are shown)

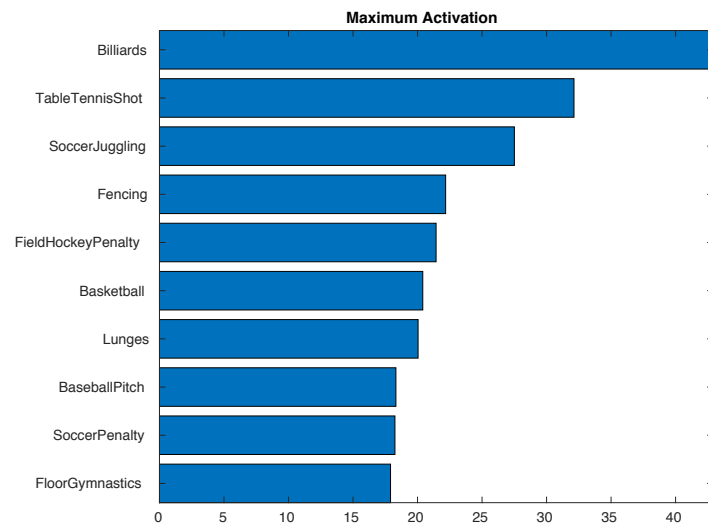
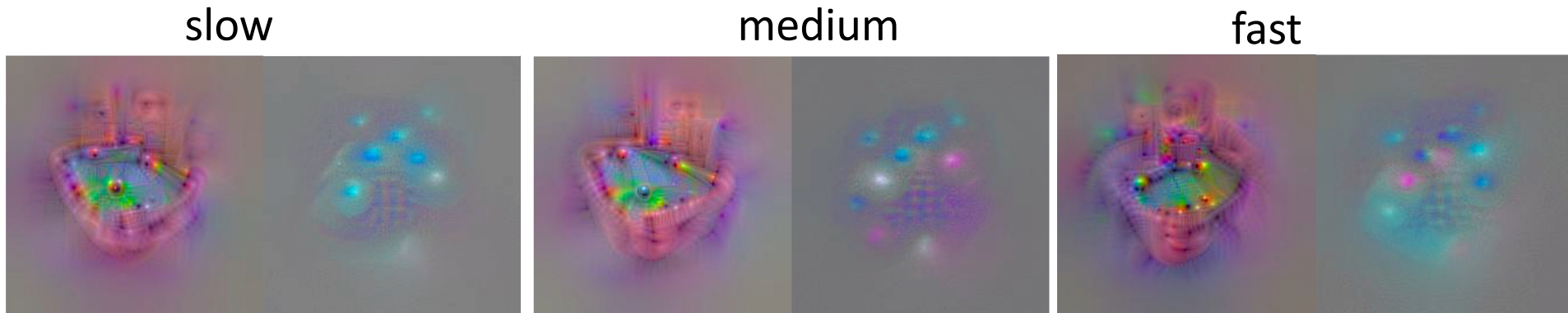
Appearance

conv4_3 f1-4

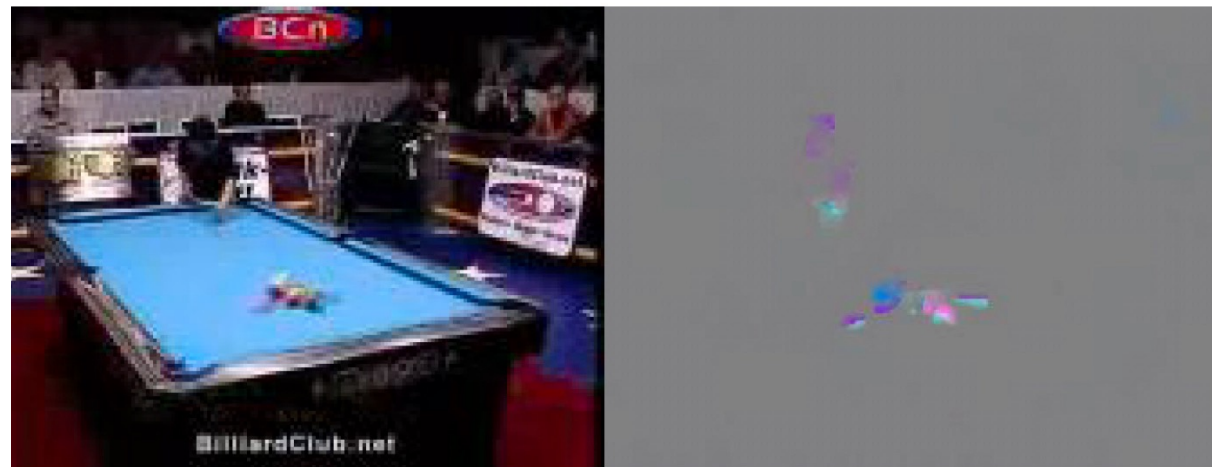
Motion



Filter #251 at conv5 fusion – the strongest local Billiards unit



(c) test set activity



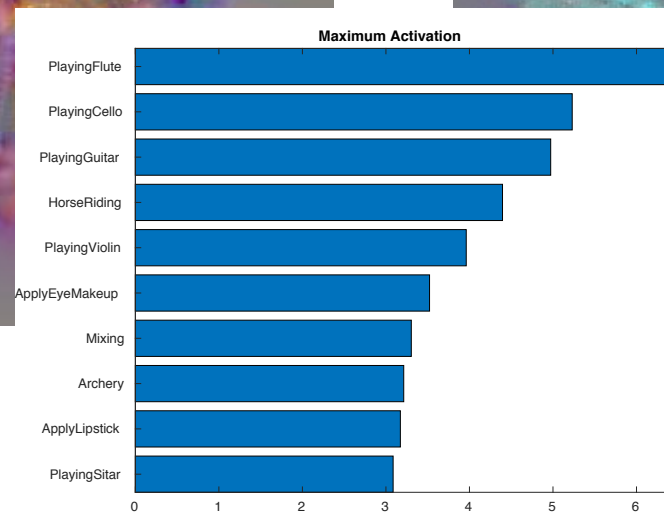
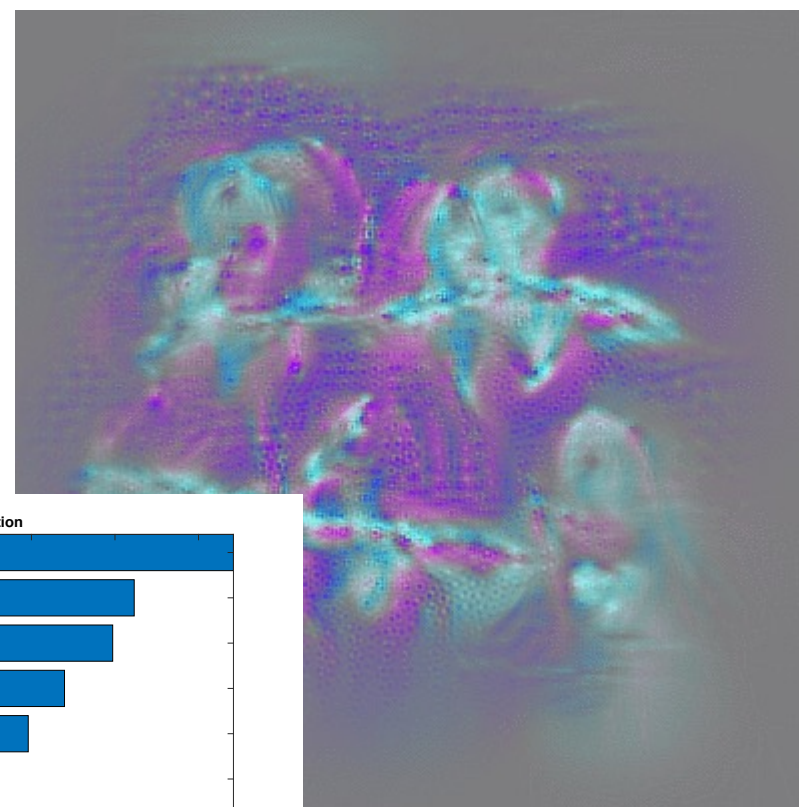
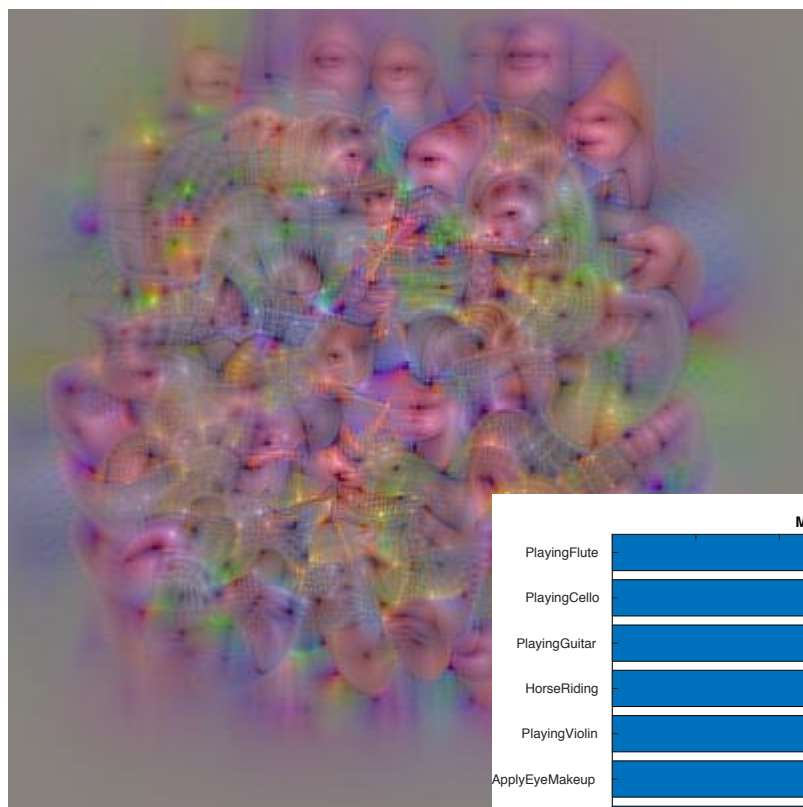
(d)

(e)

FC 6 (4096 features; RF 404x404)

Appearance

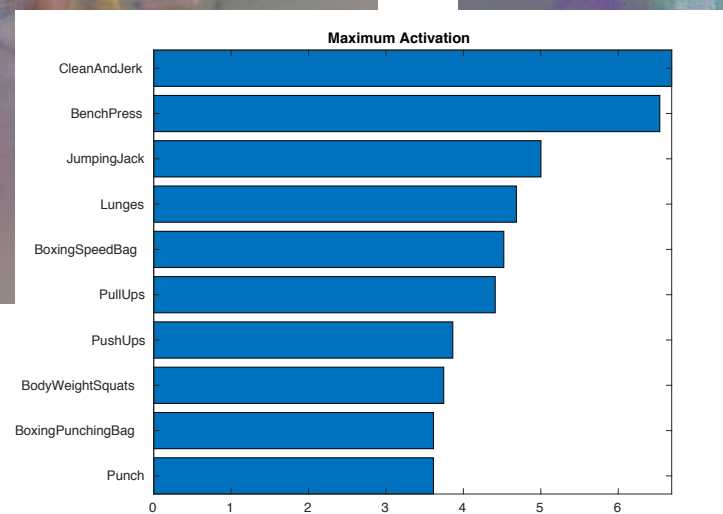
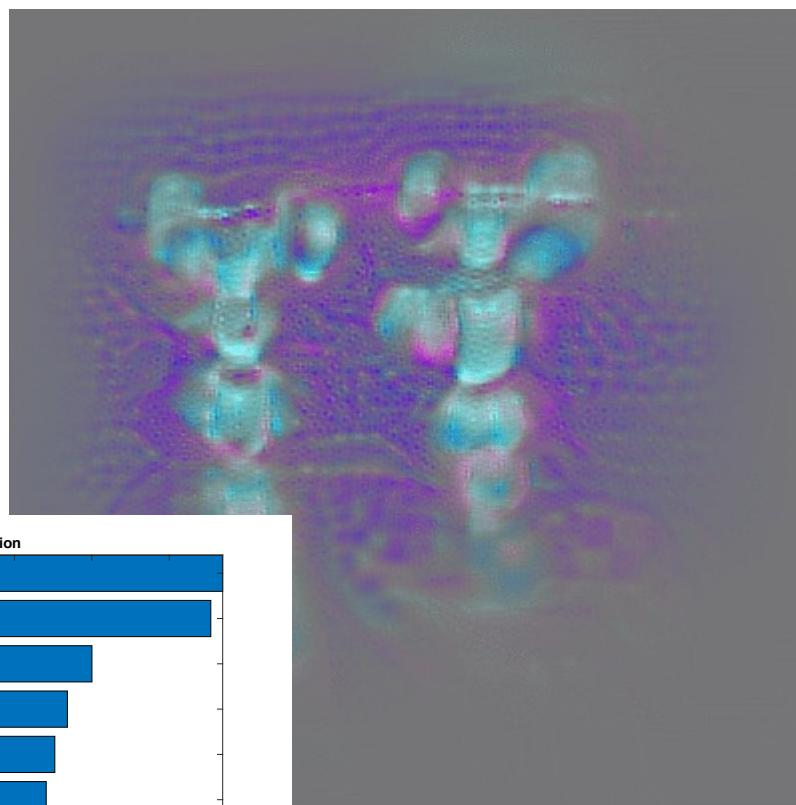
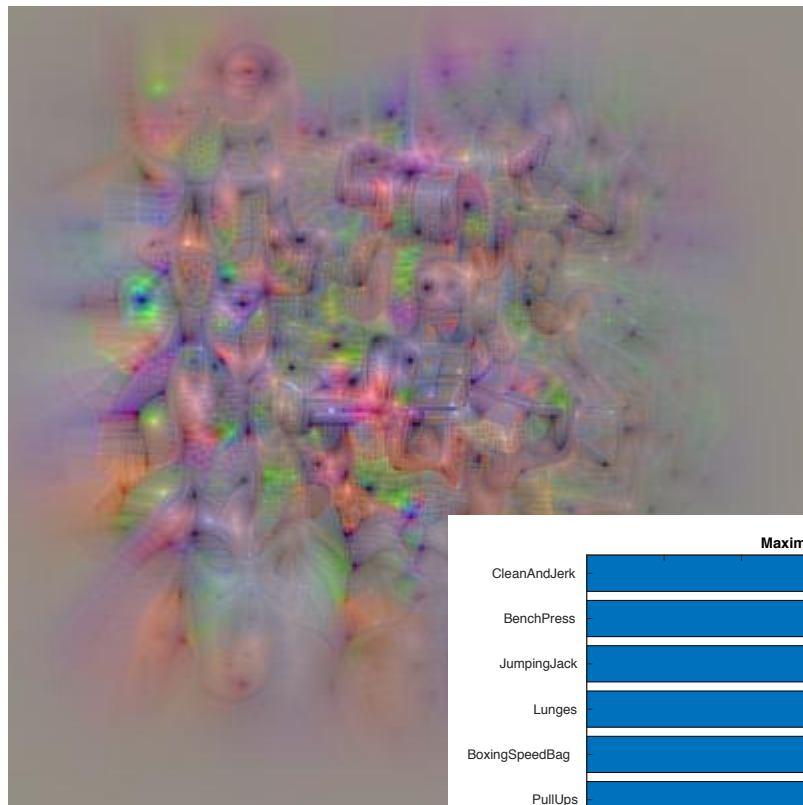
Slow motion



FC 7 (4096 features; RF 404x404)

Appearance

Slow motion



Last layer



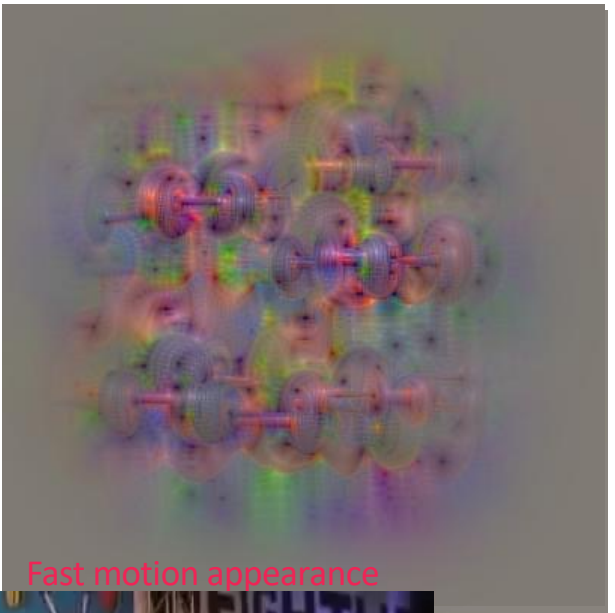
Appearance

→
"CleanAndJerk"

Slow motion

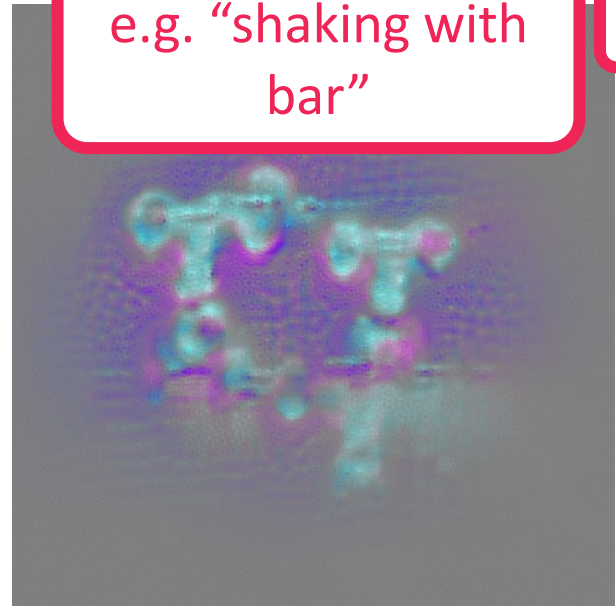


Fast motion

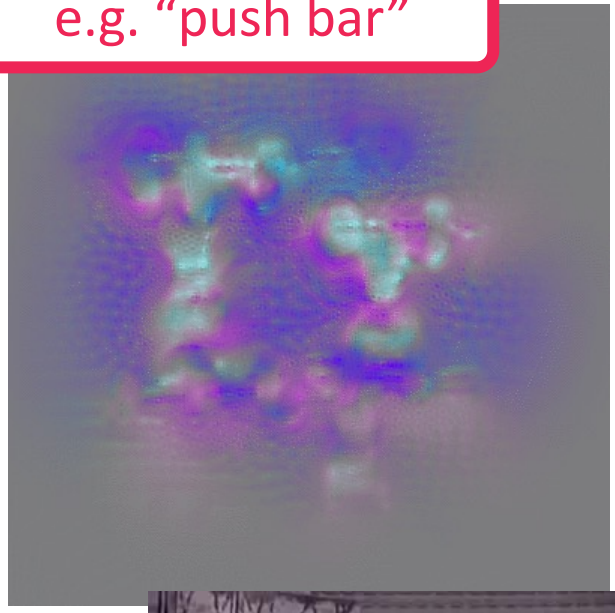


Fast motion appearance

e.g. "shaking with bar"



e.g. "push bar"



ing idiosyncracies in dat

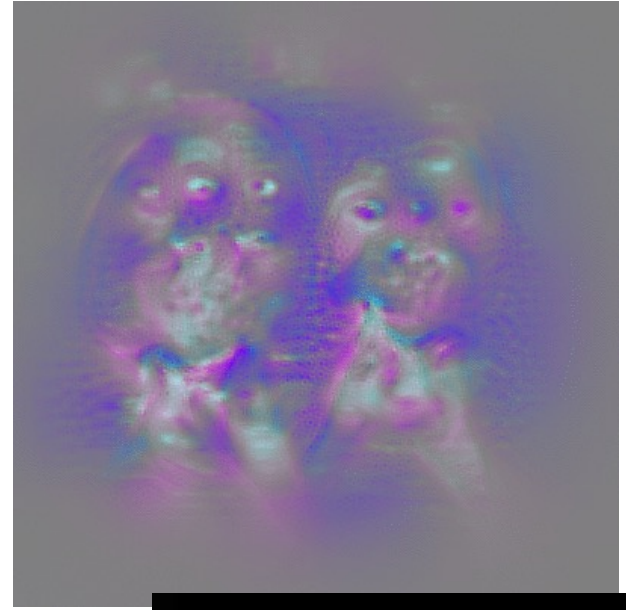
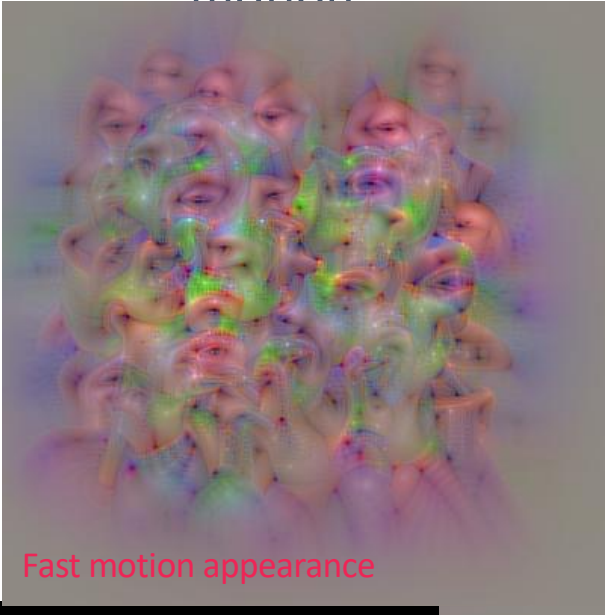


→ "ApplyLipstick"

Appearance
motion

Slow motion

Fast



Revealing idiosyncracies in data



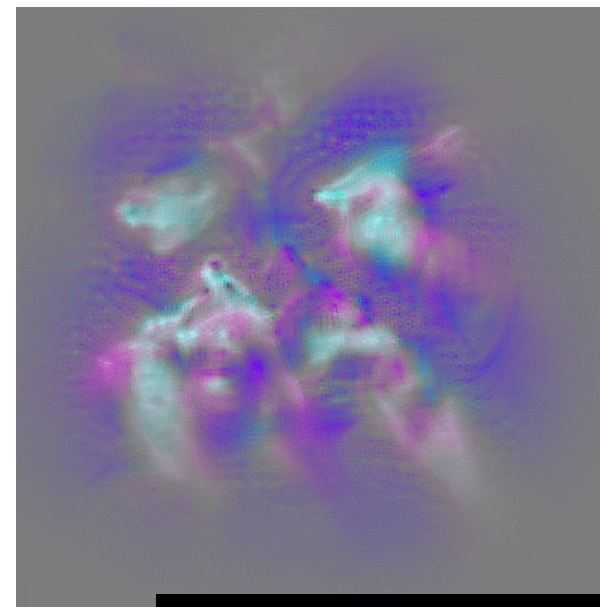
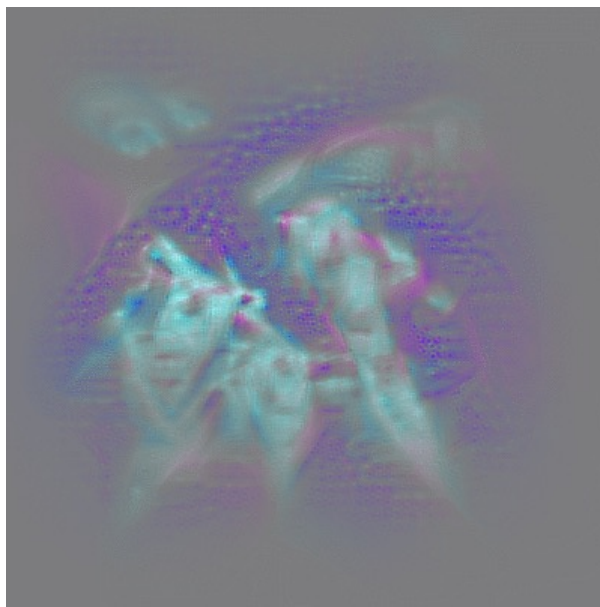
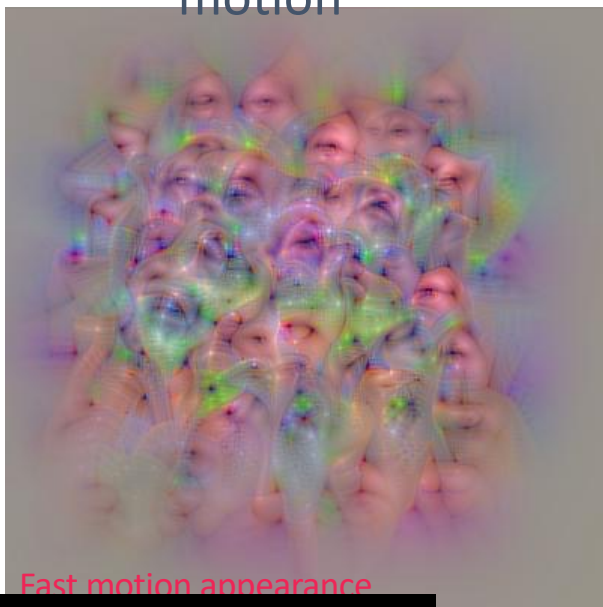
Appearance
motion



Slow motion



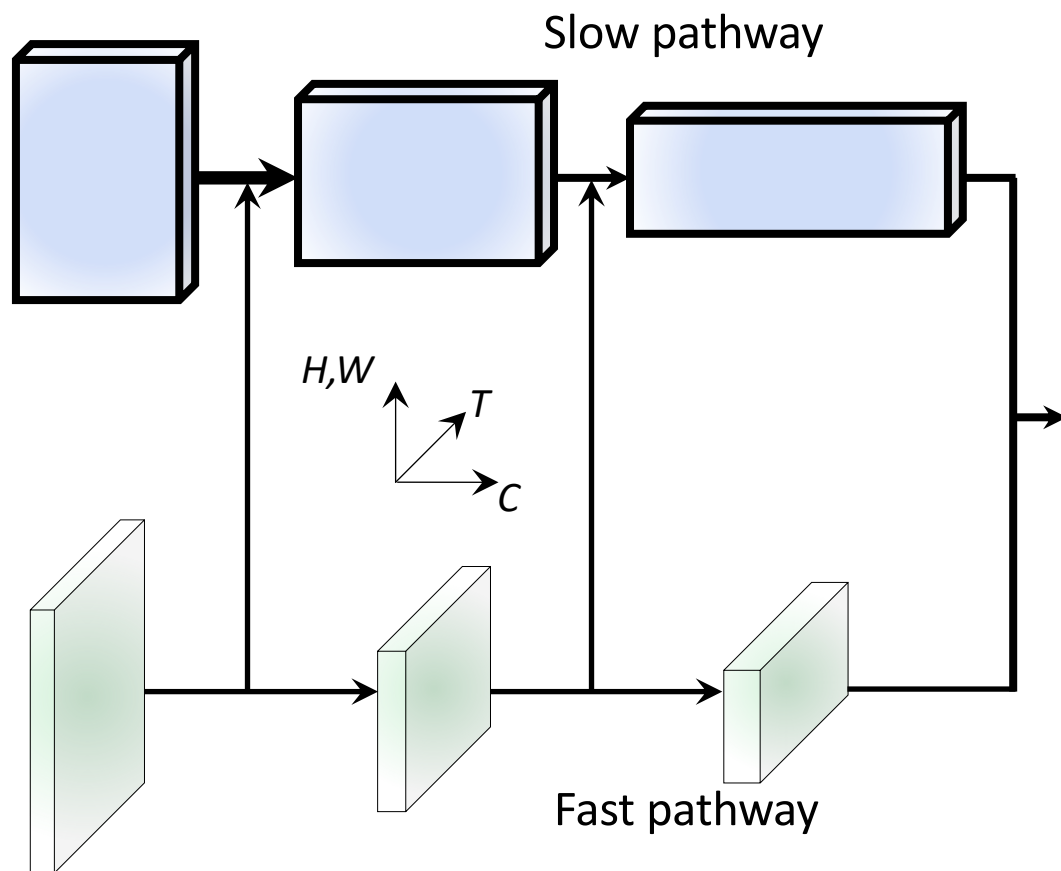
Fast



SlowFast Networks for Video Recognition

Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik and Kaiming He

- New backbone network for human action classification & detection



Slow frame rate

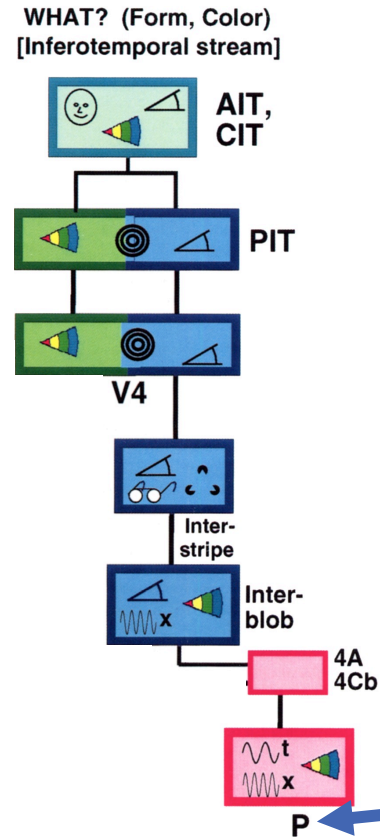


← “Hand-clap”
(action detection annotation)

Human brain: Separate visual pathways

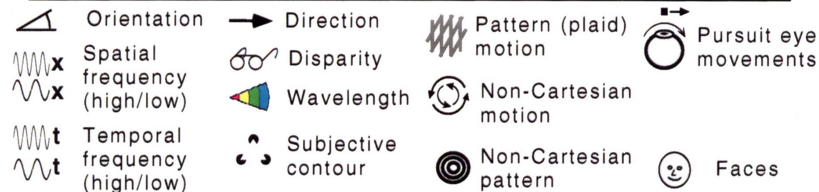
- ➔ Minority: $\approx 20\%$
- ➔ Fast conduction rate (more myelin)
- ➔ Grayscale
- ➔ Processes information about depth & motion
- ➔ Large receptive field

Magno cells



- ➔ Majority: $\approx 80\%$
- ➔ Slow conduction rate (less myelin)
- ➔ Color
- ➔ Processes information about color & detail
- ➔ Small receptive field

Parvo cells

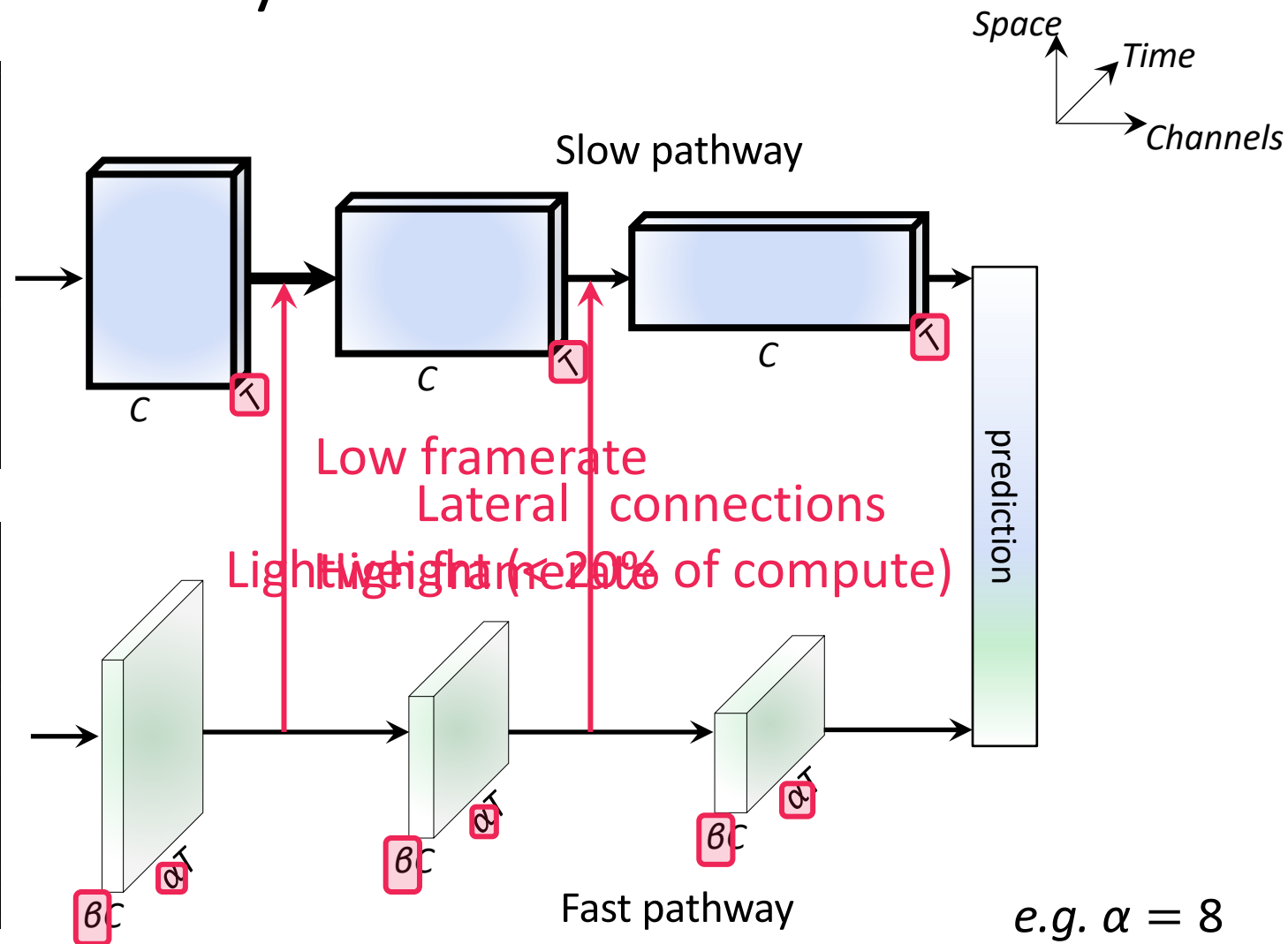


Basic idea: Two pathways

- **Slow** pathway
 - Low frame rate
 - Capturing spatial semantics
- **Fast** pathway
 - High frame rate
 - Capturing motion information

Basic idea: Two pathways

Slow



Fast



e.g. $\alpha = 8$
 $\beta = 1/8$

Example instantiation of a SlowFast network

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

stage	Slow pathway
raw clip	-
data layer	stride 16, 1^2
conv ₁	$1 \times 7^2, 64$ stride 1, 2^2
pool ₁	1×3^2 max stride 1, 2^2
res ₂	$\begin{bmatrix} 1 \times 1^2, 64 \\ 1 \times 3^2, 64 \\ 1 \times 1^2, 256 \end{bmatrix} \times 3$
res ₃	$\begin{bmatrix} 1 \times 1^2, 128 \\ 1 \times 3^2, 128 \\ 1 \times 1^2, 512 \end{bmatrix} \times 4$
res ₄	$\begin{bmatrix} 3 \times 1^2, 256 \\ 1 \times 3^2, 256 \\ 1 \times 1^2, 1024 \end{bmatrix} \times 6$
res ₅	$\begin{bmatrix} 3 \times 1^2, 512 \\ 1 \times 3^2, 512 \\ 1 \times 1^2, 2048 \end{bmatrix} \times 3$
global average pool, concate, fc	
# classes	

SlowFast ablations: Individual paths

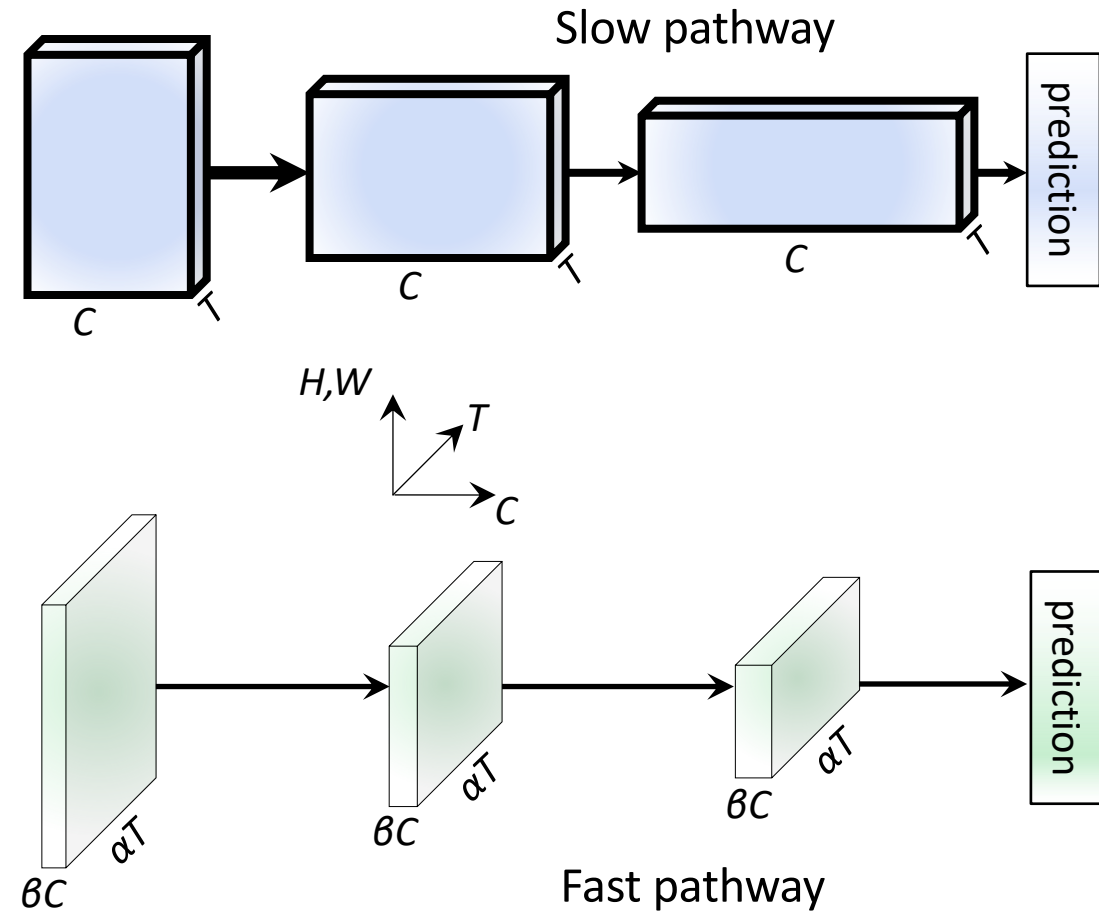
- Kinetics action classification dataset has 240k training videos and 20k validation videos in 400 classes

model	$T \times \tau$	t-reduce	top-1	top-5	GFLOPs
3D R-50	8×8	2^1	73.5	90.8	28.1
3D R-50	8×8	1	74.6	91.5	44.9
our Slow-only, R-50	4×16	1	72.6	90.3	20.9
our Fast-only, R-50	32×2	1	51.7	78.5	4.9

(b) **Individual pathways:** Training our Slow-only or Fast-only pathway alone, using the structure specified in Table 1. “t-reduce” is the total temporal downsampling factor within the network.

$$\alpha = 8$$

$$\beta = 1/8$$



SlowFast ablations: Learning curves

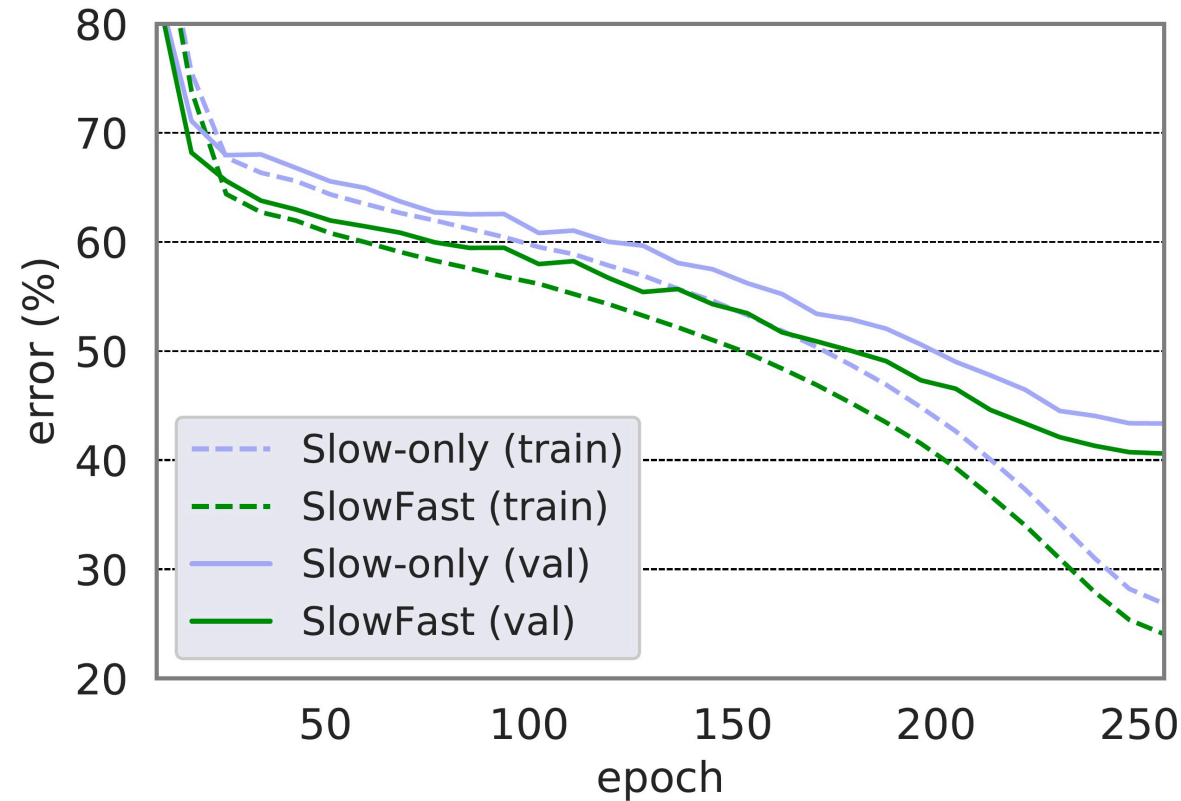
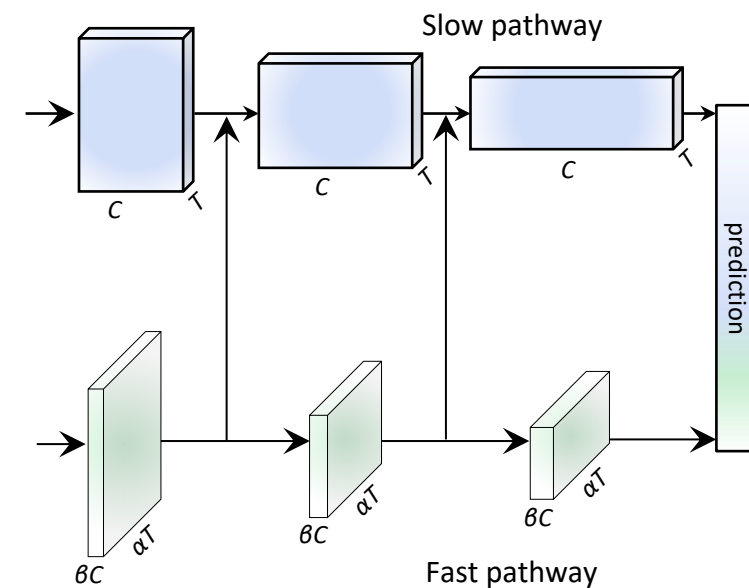
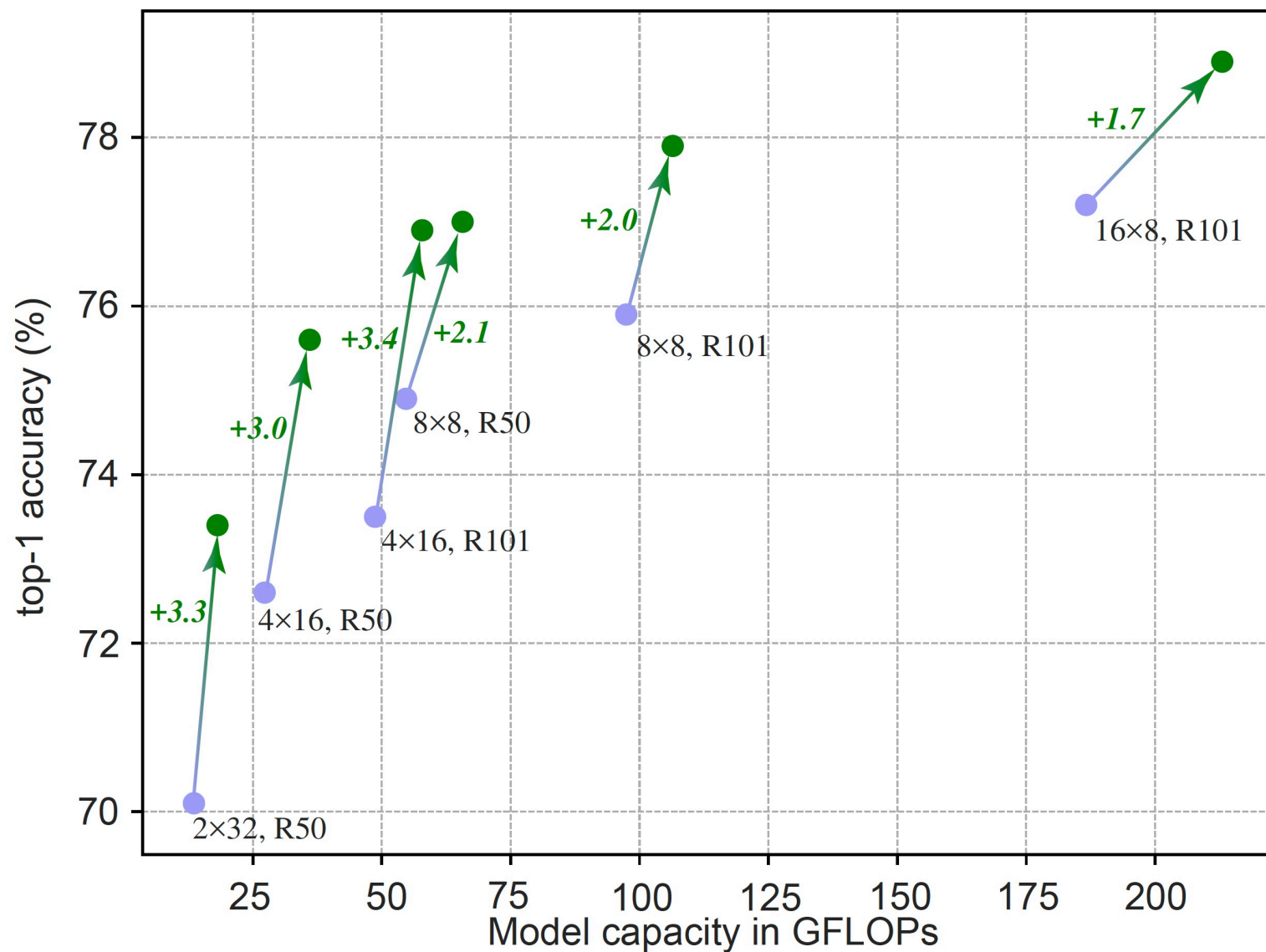


Figure 2. Training procedure on Kinetics for Slow-only (blue) vs. SlowFast (green) network. We show the top-1 training error (dash) and validation error (solid). The curves are single-crop *errors*; the video *accuracy* is 72.6% vs. 75.6% (see also Table 2c).

SlowFast ablations: Video action classification

Kinetics action classification performance

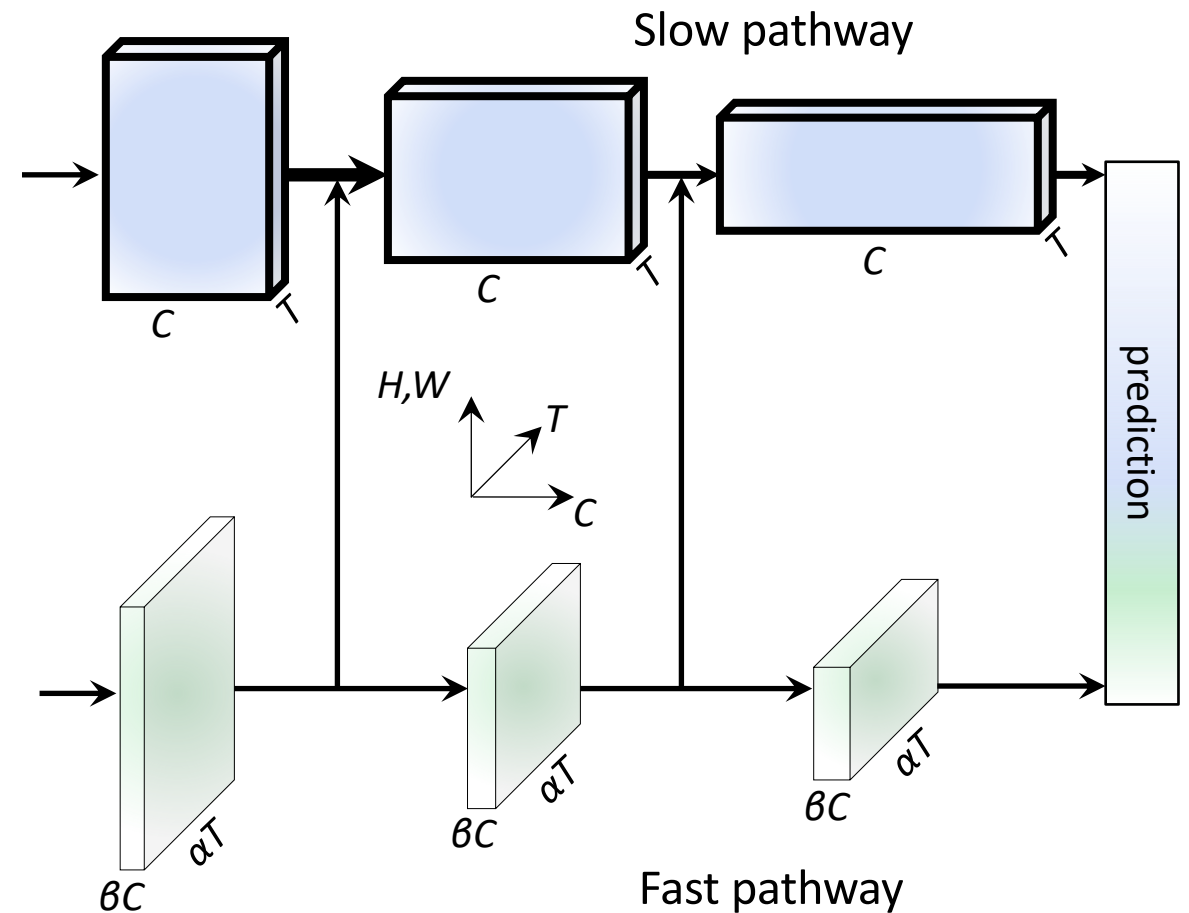


SlowFast ablations: Making the Fast path thin in channel dimension

- Kinetics dataset has 240k training videos and 20k validation videos in 400 classes

	top-1	top-5	GFLOPs
Slow-only	72.6	90.3	20.9
$\beta = 1/4$	75.6	91.7	41.7
1/6	75.8	92.0	32.0
1/8	75.6	92.1	27.6
1/12	75.2	91.8	25.1
1/16	75.1	91.7	23.4
1/32	74.2	91.3	21.9

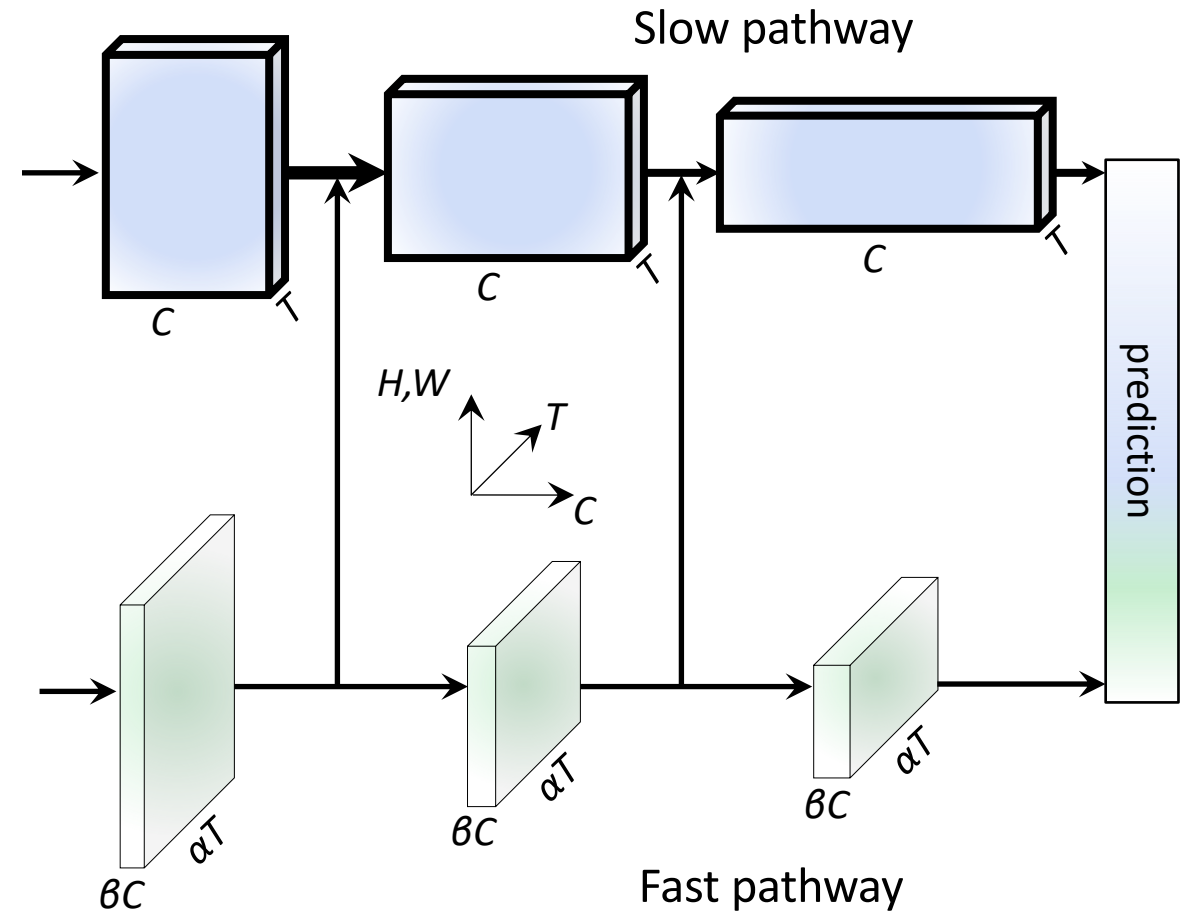
(d) **Channel capacity ratio:** Varying values of β , the channel capacity ratio of the Fast pathway. Backbone: R-50.



SlowFast ablations: Weaken the Fast appearance information

Fast pathway	spatial	top-1	top-5	GFLOPs
RGB	-	75.6	92.1	27.6
RGB, $\beta=1/4$	<i>half</i>	74.7	91.8	26.3
gray-scale	-	75.5	91.9	26.1
time diff	-	74.5	91.6	26.2
optical flow	-	73.8	91.3	26.9

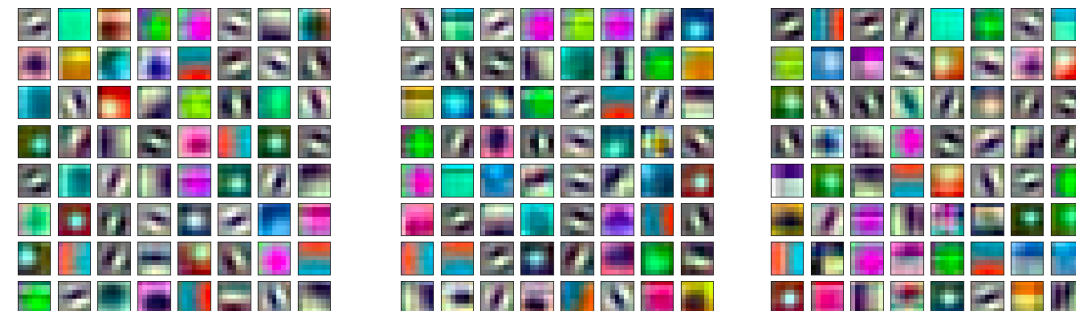
(e) **Weaker spatial input to Fast pathway:** Various ways of weakening spatial inputs to the Fast pathway in SlowFast models. $\beta=1/8$ unless specified otherwise. Backbone: R-50.



SlowFast ablations: Weaker input & reduced channels: conv1 filters

Fast pathway	spatial	top-1	top-5	GFLOPs
RGB	-	75.6	92.1	27.6

Slow



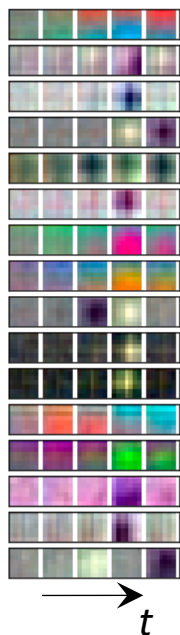
rgb

grayscale

time diff

$\beta = 1/4$

gray-scale	-	75.5	91.9	26.1
time diff	-	74.5	91.6	26.2

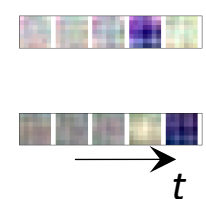
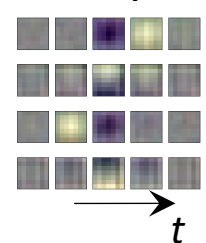
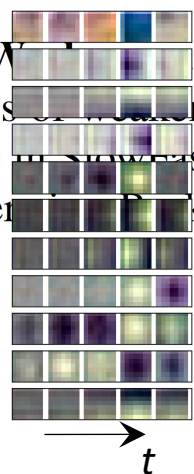


$\beta = 1/6$

(e) **Weak spatial input to Fast pathway:** Various ways of weakening spatial inputs to the Fast pathway in slowfast models. $\beta=1/8$ unless specified otherwise. Backbone: $\beta R=50$

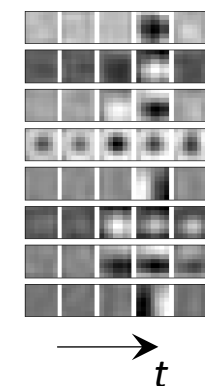
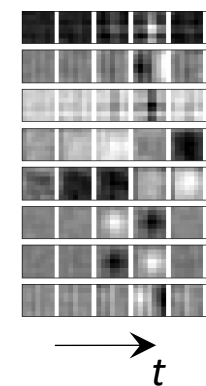
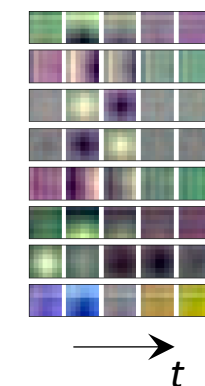
$\beta = 1/16$

$\beta = 1/32$



$\beta = 1/8$

Fast



rgb



grayscale



dt

SlowFast: State-of-the-art comparison on Kinetics

model	flow	pretrain	top-1	top-5	inference GFLOPs×crops
I3D [1]		ImageNet	72.1	90.3	108 × N/A
Two-Stream I3D [1]	✓	ImageNet	75.7	92.0	216 × N/A
S3D-G [6]	✓	ImageNet	74.7	93.4	142.8 × N/A
Nonlocal R-50 [5]		ImageNet	76.5	92.6	282 × 30
Nonlocal R-101 [5]		ImageNet	77.7	93.3	359 × 30
R(2+1)D Flow [3]	✓	-	67.5	87.2	152 × 115
STC [2]		-	68.7	88.5	N/A × N/A
ARTNet [4]		-	69.2	88.3	23.5 × 250
S3D [6]		-	69.4	89.1	66.4 × N/A
ECO [7]		-	70.0	89.4	N/A × N/A
I3D [1]	✓	-	71.6	90.0	216 × N/A
R(2+1)D [3]		-	72.0	90.0	152 × 115
R(2+1)D [3]	✓	-	73.9	90.9	304 × 115
SlowFast, R50 (4×16)		-	75.6	92.1	36.1 × 30
SlowFast, R50		-	77.0	92.6	65.7 × 30
SlowFast, R50 + NL		-	77.7	93.1	80.8 × 30
SlowFast, R101		-	77.9	93.2	106 × 30
SlowFast, R101 + NL		-	79.0	93.6	115 × 30

+ 5.1%
top-1

at 10%
of FLOPs

Table 1. **Comparison with the state-of-the-art on Kinetics-400.** In the column of computational cost, we report the cost of a single spacetime crop and the numbers of such crops used. “N/A” indicates the numbers are not available for us. The SlowFast models are the $T \times \tau = 8 \times 8$ versions, unless specified.

- [1] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proc. CVPR*, 2017.
- [2] A. Diba, M. Fayyaz, V. Sharma, M. M. Arzani, R. Yousefzadeh, J. Gall, and L. Van Gool. Spatio-temporal channel correlation networks for action classification. In *Proc. ECCV*, 2018.
- [3] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proc. CVPR*, 2018.
- [4] L. Wang, W. Li, W. Li, and L. Van Gool. Appearance-and-relation networks for video classification. In *Proc. CVPR*, 2018.
- [5] X. Wang, R. Girshick, A. Gupta, and K. He. Non-local neural networks. In *Proc. CVPR*, 2018.
- [6] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy. Rethinking spatiotemporal feature learning for video understanding. *arXiv preprint arXiv:1712.04851*, 2017.
- [7] M. Zolfaghari, K. Singh, and T. Brox. ECO: efficient convolutional network for online video understanding. In *Proc. ECCV*, 2018.

SlowFast: State-of-the-art comparison Charades¹

- Charades has 9.8k training videos and 1.8k validation videos in 157 classes
- Multi-label classification setting of longer activities spanning 30 seconds on average

model	pretrain	mAP	inference GFLOPs × views
CoViAR, R-50 [55]	ImageNet	21.9	N/A
Asyn-TF, VGG16 [39]	ImageNet	22.4	N/A
MultiScale TRN [58]	ImageNet	25.2	N/A
Nonlocal, R101 [52]	ImageNet+Kinetics400	37.5	544 × 30
STRG, R101+NL [53]	ImageNet+Kinetics400	39.7	630 × 30
our baseline (Slow-only)	Kinetics-400	39.0	187 × 30
SlowFast	Kinetics-400	41.8	213 × 30
SlowFast, +NL	Kinetics-400	42.5	234 × 30
SlowFast, +NL	Kinetics-600	45.2	234 × 30

Table 4. **Comparison with the state-of-the-art on Charades.** All our variants are based on $T \times \tau = 16 \times 8$, R101.

Annotated Actions: (gray if not active)

Turning on a light
 Walking through a doorway
 Taking a box from somewhere
 Holding a box
 Opening a box
 Taking a pillow from somewhere
 Taking something from a box
 Closing a box
 Holding a pillow
 Snuggling with a pillow
 Putting something on a shelf
 Putting a box somewhere

Video 21 of 50: (3x Speed)



Annotated Objects:

Box, Closet, Doorway, Light, Pillow, Shelf

Script:

A person turns on the light in a closet, opens a large container, then grasps a pillow from it.

¹G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In ECCV, 2016. , CVPR 2016

Experiments: AVA¹ Action Detection

- Fine-scale localization of 80 different physical actions
- Data from 437 different movies and spatiotemporal labels are provided in a 1Hz interval
- 211k training and 57k validation video segments
- We follow the standard protocol of evaluating on 60 most frequent classes
- Every person is annotated with a bounding box and (possibly multiple) actions

SlowFast detector output

AVA validation set videos

We show:

- Detected boxes in **green**, with predictions (if confidence > 0.5) on **top**
- Ground-Truth (GT) boxes in **red**, with annotated labels on the **bottom**

Detections and GT are shown every second, with reduced playback speed

¹Gu et al. AVA: A Video Dataset of Spatio-temporally Localized Atomic Visual Actions, CVPR 2018

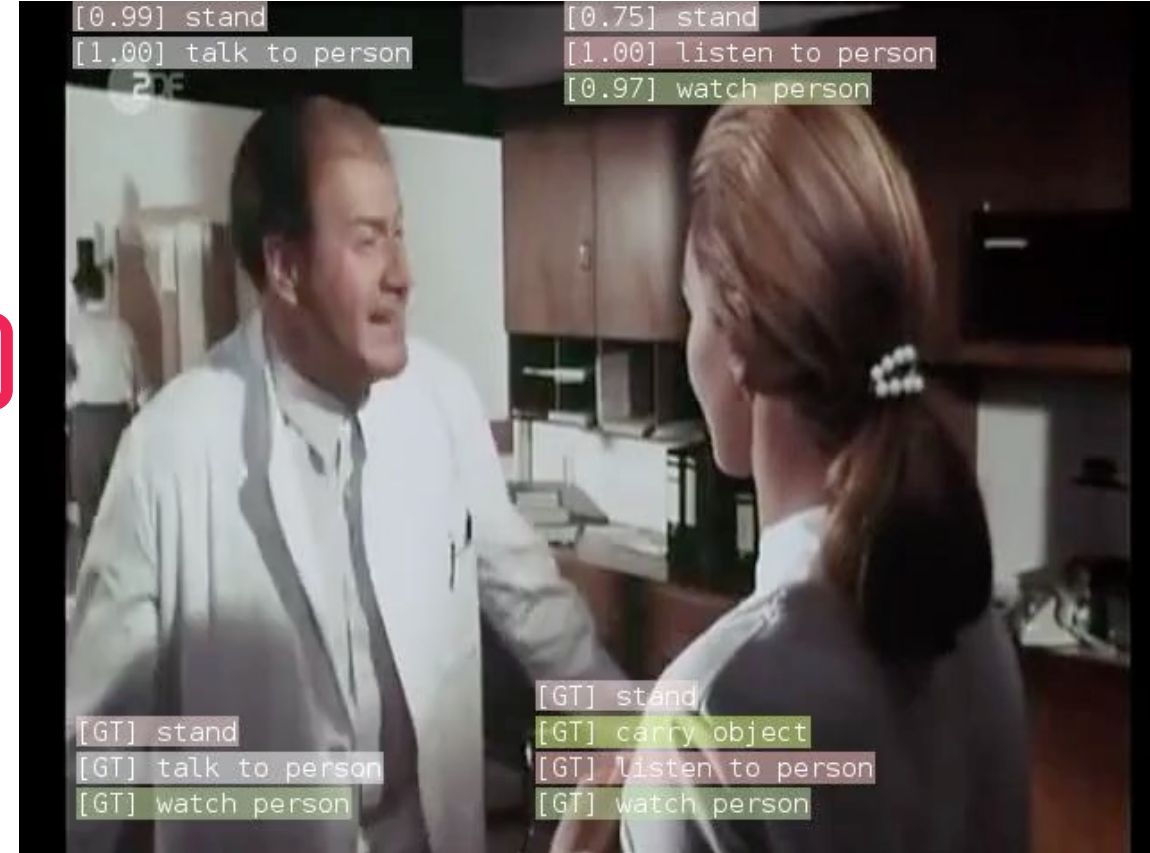
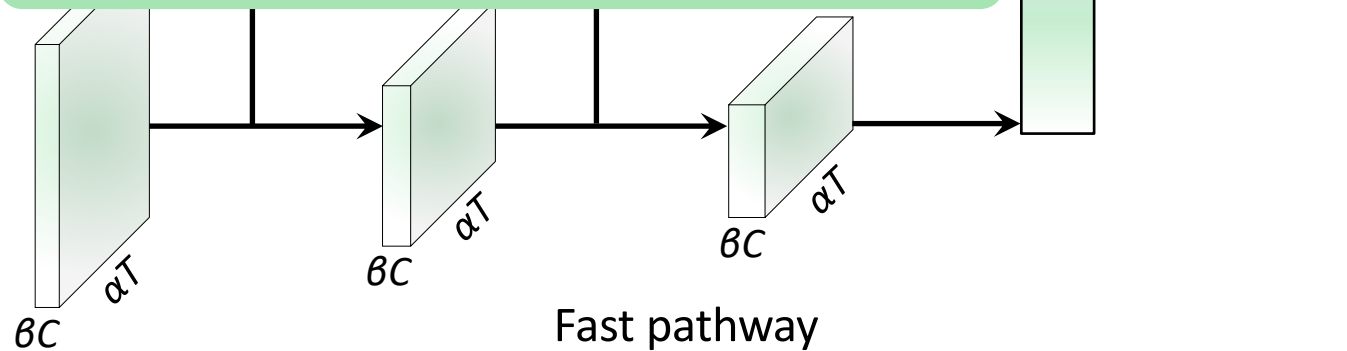
SlowFast: AVA action detection



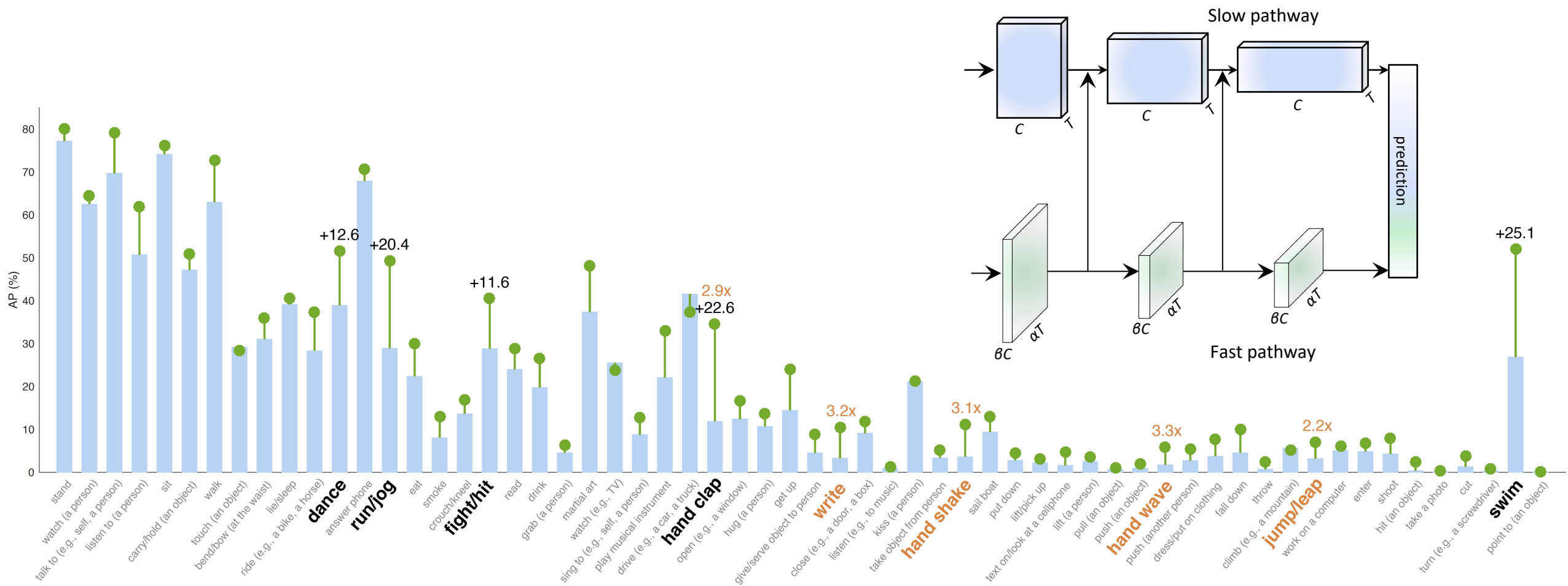
2018: Challenge winner: 21.1 mAP

2019: SlowFast winner: 34.3 mAP

Top-3 ranked teams used SlowFast



SlowFast ablations: AVA class level performance



Experiments: AVA Qualitative results



Experiments: AVA Qualitative results

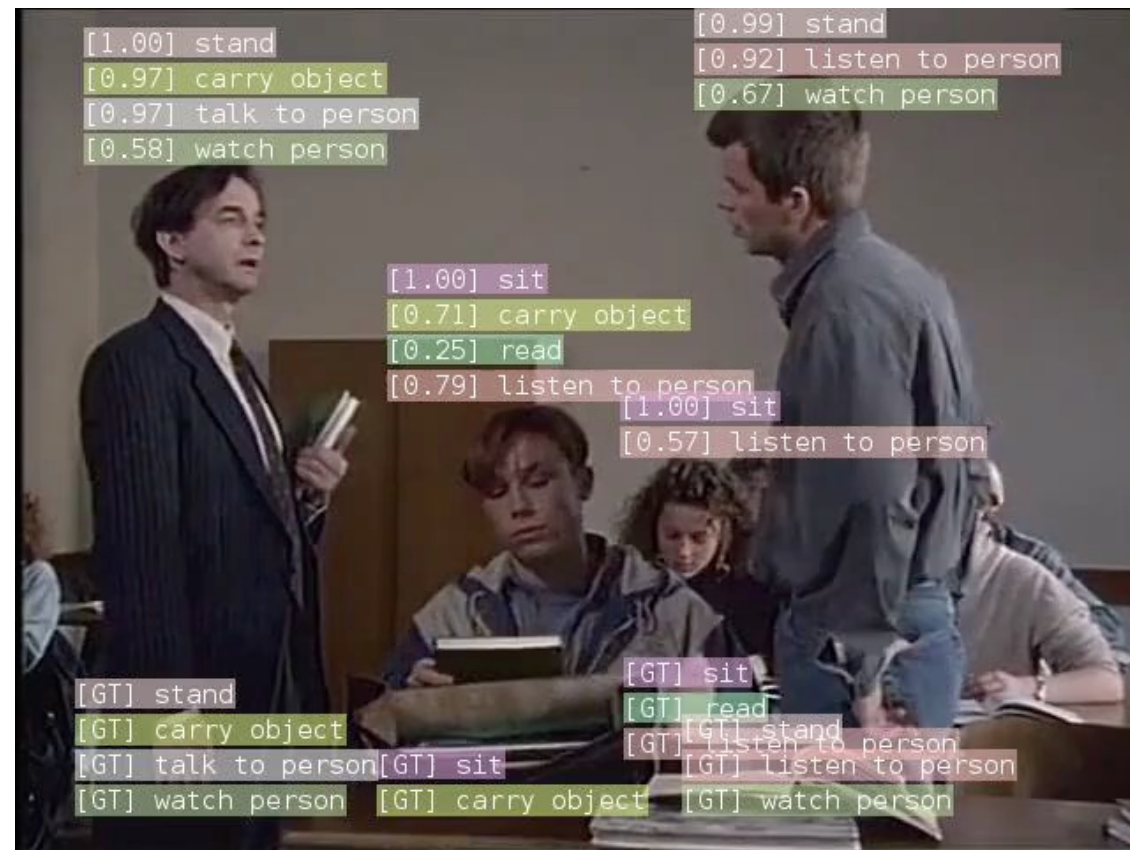


Pytorch code available:

Conclusion

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

- The time axis is a special dimension of video
- *3D ConvNets* treat space and time uniformly
- *Non-local networks* and *Long-term feature banks* aggregate long-term spatiotemporal information
- *SlowFast* & *Two-Stream* networks treat space and time differently and share motivation from neuroscience
- The *SlowFast* architecture design focuses on contrasting the speed along the temporal axis
- Given the mutual benefits of jointly modeling video with different temporal speeds, we hope that this concept can foster further research in video analysis



Overview

- Optical Flow
- ConvNets for Video
- Video Generation

Video Generation



Samples from a text-conditioned video diffusion model, conditioned on the string *fireworks*.

(video from: Ho et al., “Video Diffusion Models”, *arXiv*, 2022,
<https://video-diffusion.github.io/>)

[Ho et al., “Video Diffusion Models”, *arXiv*, 2022](#)

[Harvey et al., “Flexible Diffusion Modeling of Long Videos”, *arXiv*, 2022](#)

[Yang et al., “Diffusion Probabilistic Modeling for Video Generation”, *arXiv*, 2022](#)

[Höppe et al., “Diffusion Models for Video Prediction and Infilling”, *arXiv*, 2022](#)

[Voleti et al., “MCVD: Masked Conditional Video Diffusion for Prediction, Generation, and Interpolation”, *arXiv*, 2022](#)

[Kreis, Gao, Vahdat CVPR 2022]

Video Generation

Video Generation Tasks:

- Unconditional Generation (Generate all frames)
- Future Prediction (Generate future from past frames)
- Past Prediction (Generate past from future frames)
- Interpolation (Generate intermediate frames)

➔ Learn a model of the form:

$$p_{\theta}(\mathbf{x}^{t_1}, \dots, \mathbf{x}^{t_K} | \mathbf{x}^{\tau_1}, \dots, \mathbf{x}^{\tau_M})$$

Given frames: $\mathbf{x}^{\tau_1}, \dots, \mathbf{x}^{\tau_M}$

Frames to be predicted: $\mathbf{x}^{t_1}, \dots, \mathbf{x}^{t_K}$

[Ho et al., "Video Diffusion Models", arXiv, 2022](#)

[Harvey et al., "Flexible Diffusion Modeling of Long Videos", arXiv, 2022](#)

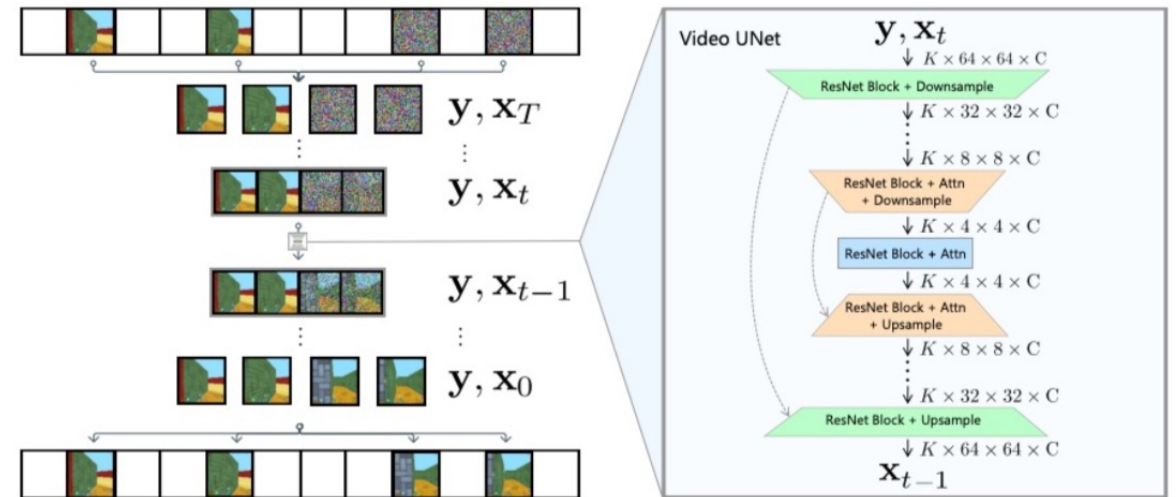
[Yang et al., "Diffusion Probabilistic Modeling for Video Generation", arXiv, 2022](#)

[Höppe et al., "Diffusion Models for Video Prediction and Infilling", arXiv, 2022](#)

[Voleti et al., "MCVD: Masked Conditional Video Diffusion for Prediction, Generation, and Interpolation", arXiv, 2022](#)

➔ Learn one model for everything:

- Architecture as **one diffusion model over all frames concatenated**.
- Mask frames to be predicted; provide conditioning frames; vary applied masking/conditioning for different tasks during training.
- Use **time position encodings** to encode times.



(image from: Harvey et al., "Flexible Diffusion Modeling of Long Videos", arXiv, 2022)

Video Generation

Architecture Details

Architecture Details:

Data is 4D (image height, image width, #frames, channels)

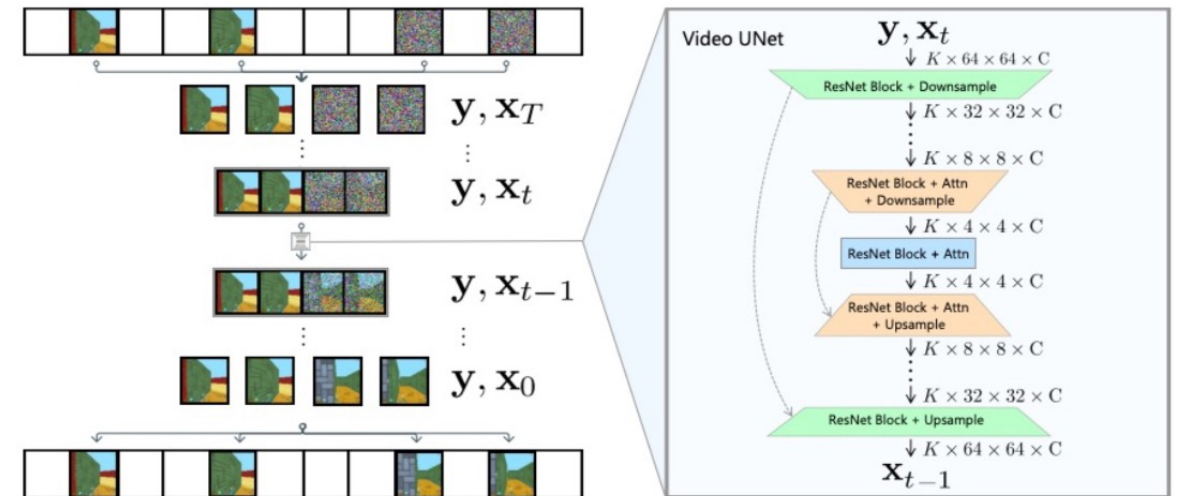
- Option (1): 3D Convolutions. Can be computationally expensive.
- Option (2): Spatial 2D Convolutions + Attention Layers along frame axis.

➔ Additional Advantage:

Ignoring the attention layers initially, the model can be pre-trained on pure image data!

➔ Learn one model for everything:

- Architecture as **one diffusion model over all frames concatenated**.
- Mask frames to be predicted; provide conditioning frames; vary applied masking/conditioning for different tasks during training.
- Use **time position encodings** to encode times.



(image from: Harvey et al., "Flexible Diffusion Modeling of Long Videos", arXiv, 2022)

[Ho et al., "Video Diffusion Models", arXiv, 2022](#)

[Harvey et al., "Flexible Diffusion Modeling of Long Videos", arXiv, 2022](#)

[Yang et al., "Diffusion Probabilistic Modeling for Video Generation", arXiv, 2022](#)

[Höppe et al., "Diffusion Models for Video Prediction and Infilling", arXiv, 2022](#)

[Voleti et al., "MCVD: Masked Conditional Video Diffusion for Prediction, Generation, and Interpolation", arXiv, 2022](#)

[Kreis, Gao, Vahdat CVPR 2022]

Video Generation

Results

Long term video generation in hierarchical manner:

- 1. Generate future frames in sparse manner, conditioning on frames far back
- 2. Interpolate in-between frames



1+ hour coherent video generation possible!

Test Data:



Generated:



(video from: Harvey et al., “Flexible Diffusion Modeling of Long Videos”, *arXiv*, 2022, <https://plai.cs.ubc.ca/2022/05/20/flexible-diffusion-modeling-of-long-videos/>)

[Ho et al., “Video Diffusion Models”, *arXiv*, 2022](#)

[Harvey et al., “Flexible Diffusion Modeling of Long Videos”, *arXiv*, 2022](#)

[Yang et al., “Diffusion Probabilistic Modeling for Video Generation”, *arXiv*, 2022](#)

[Höppe et al., “Diffusion Models for Video Prediction and Infilling”, *arXiv*, 2022](#)

[Voleti et al., “MCVD: Masked Conditional Video Diffusion for Prediction, Generation, and Interpolation”, *arXiv*, 2022](#)

[Kreis, Gao, Vahdat CVPR 2022]

Video Generation demos

- <https://video-diffusion.github.io/>
- <https://makeavideo.studio/>