

Video

Computer Vision – Lecture 12

Further Reading

- Slides from [L Fei-Fei](#)
- Slides from [J Johnson](#)
- Slides from M Niessner & L Leal-Taixé ([2nd part](#))

Videos

- Sequence of frames: $T \times 3 \times H \times W$
- Frame rate: ~30 FPS (frames per second)
- Temporal coherence
- Sometimes: shot changes/skips



Video Tasks

- Activity Classification
- Temporal/Spatial Action Localisation
- Event Dense Captioning
- Active Speaker Recognition
- Sign Language Transcription
- ...



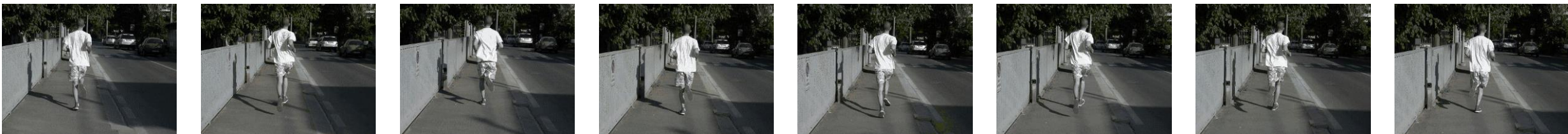
Videos

Problem: videos are big!

- HD (1080x1920):
 - $60 \times 30 \times 3 \times 1080 \times 1920$ bytes = 11.2GB / minute
- This is just the data, we still need to compute things.
- Use short, low-res clips: $T=16$, $FPS=5$, $H=W=112$: (60MB/min)

Windowed Video Processing

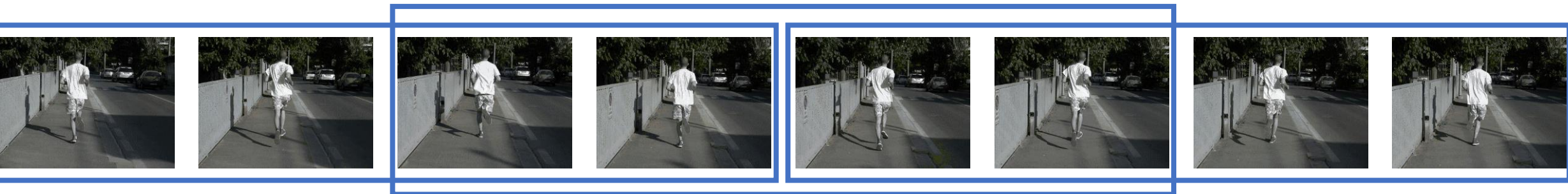
Raw Video: long, high resolution, high FPS



Training: short clips, low resolution, low FPS



Testing: run model on overlapping clips, average predictions



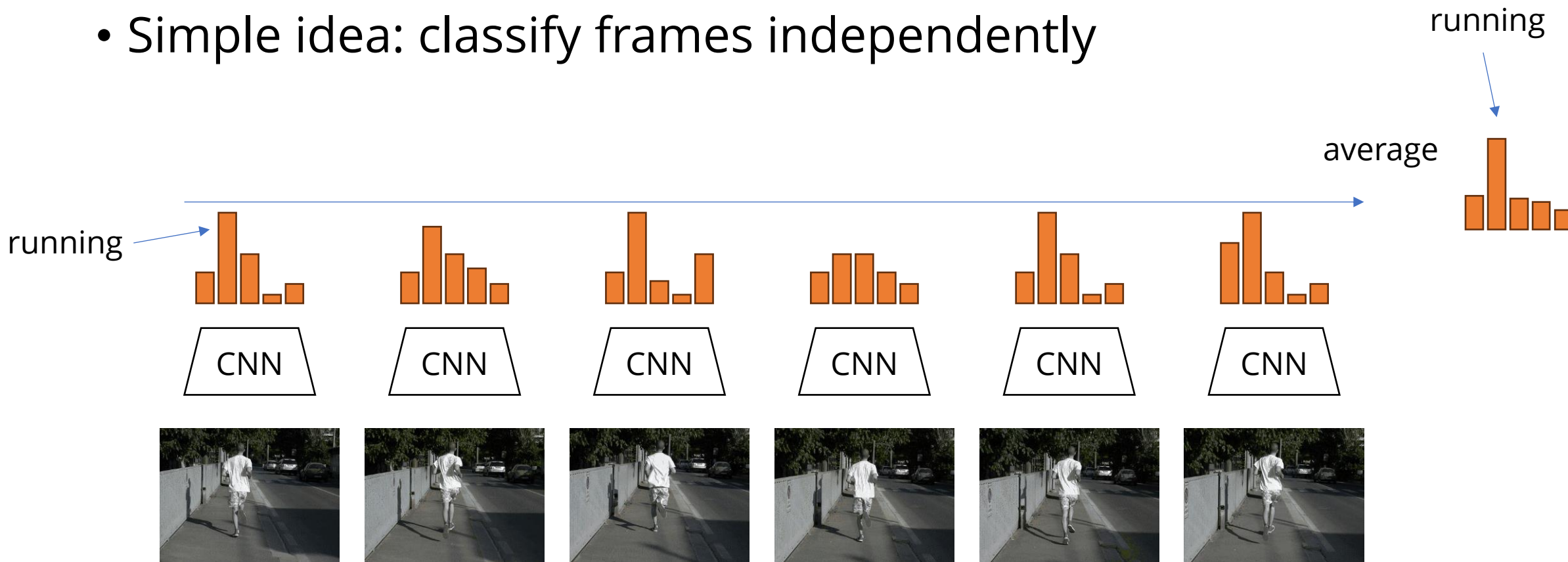
Windowed Video Processing

Familiar ideas:

- Subsampling (now in time and space)
- Sliding window (now mostly in time)
- Share as much computation as possible to increase the efficiency

Task: Video Classification

- Action recognition: running, knitting, basketball, etc.
- Simple idea: classify frames independently

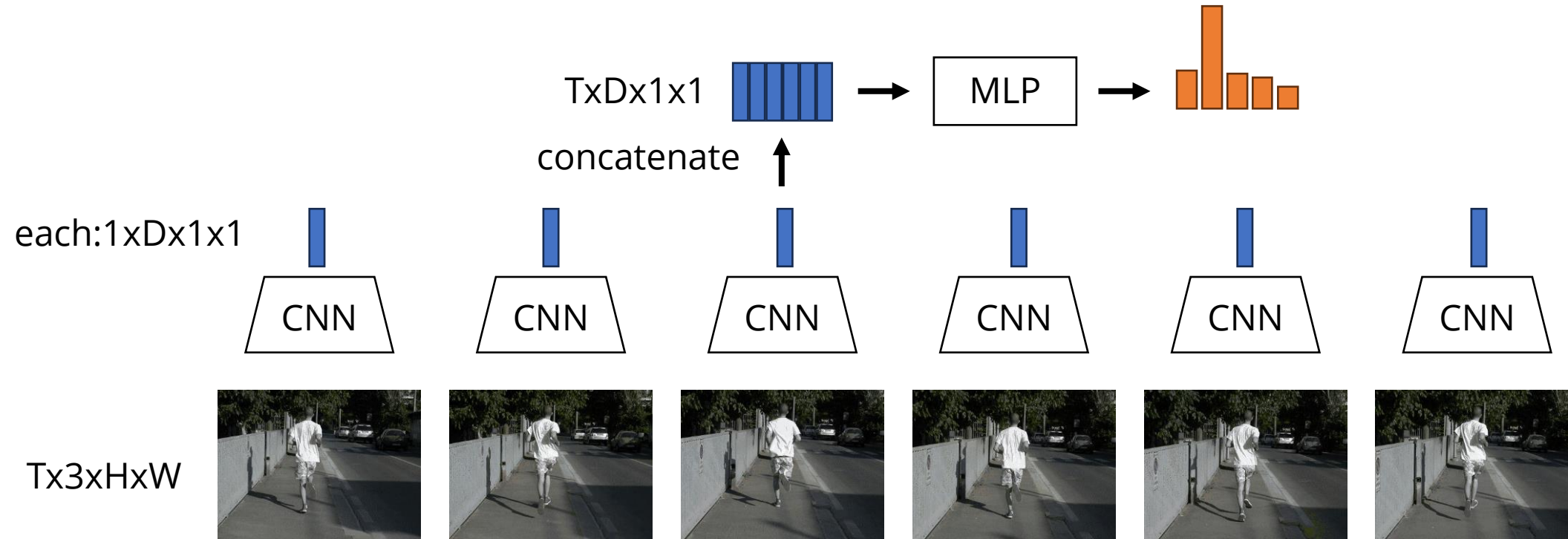


Per-Frame Models

- Predict for each frame independently with an image model.
- Average probabilities (mathematically questionable but works much better than multiplication).
- Often a very strong baseline!
- Intuition:
 - Only one frame is needed to differentiate between “running” and “swimming”.
 - This is often called: object-bias of action recognition.
- Depends on task: “sitting down” vs. “standing up” needs motion.

Sharing Computation

- Per-frame model cannot reason about time: we only average predictions.
- Add layers that have access across time.



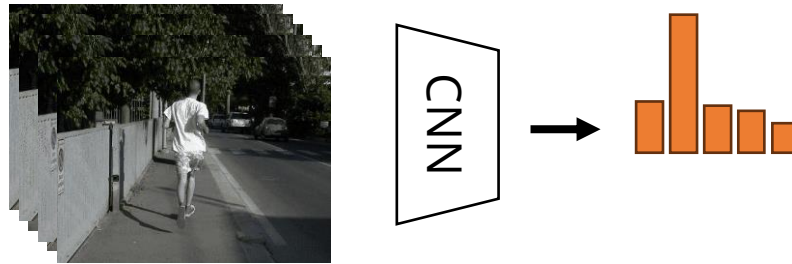
Late Fusion

- Here we decided to include temporal information “late” in the computation.
- Intuition: get per-frame high-level understanding, then combine them across time.
- Fusion mechanisms: concatenation, pooling, etc.
- Problem: low-level motion often lost after “compressing” a frame into a feature vector.

Early Fusion

- Fuse frames at the input level: $T \times 3 \times H \times W \rightarrow T \times 3 \times H \times W$
- Treat input as an image with many channels.

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



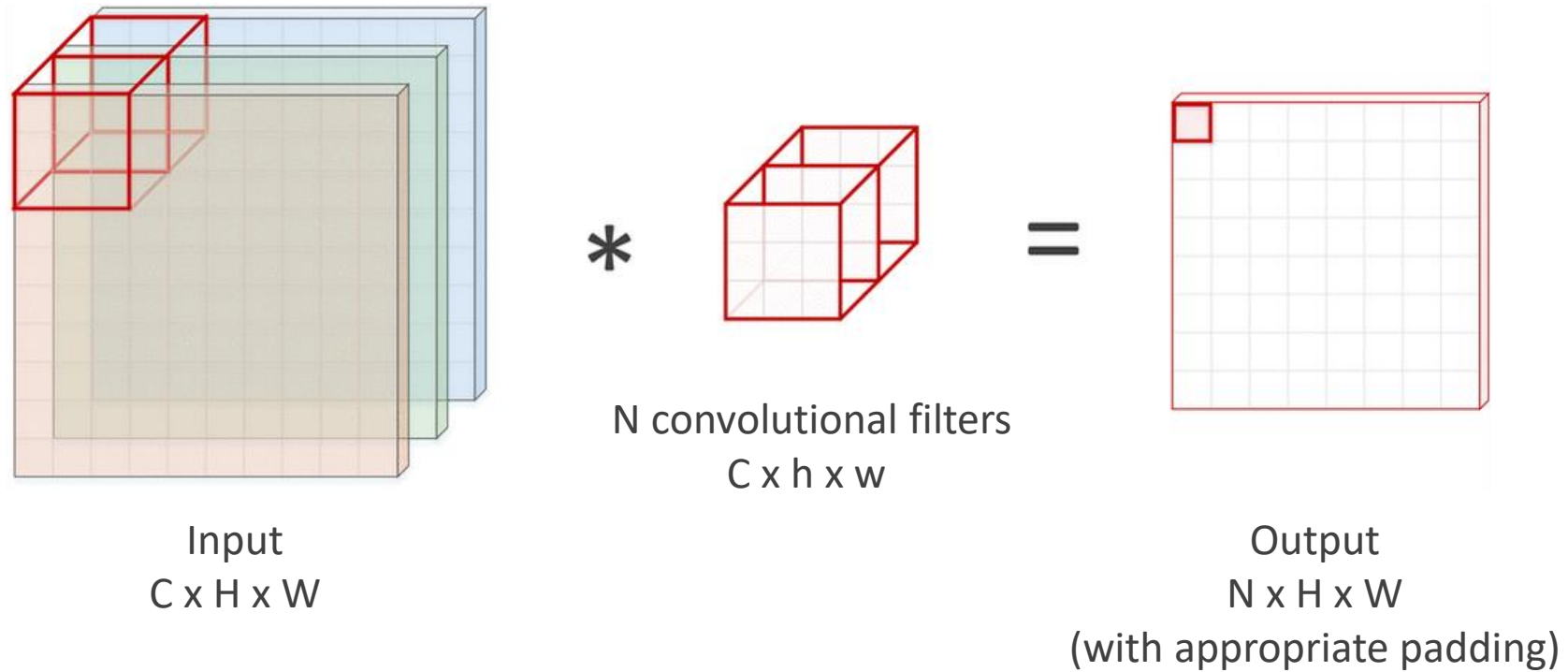
$T \times 3 \times H \times W$



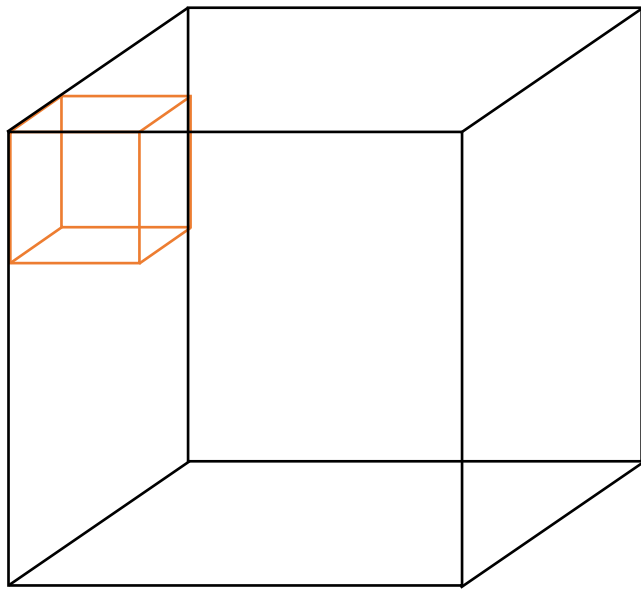
Early Fusion

- First layer takes as input the whole video stacked in the channel dimension.
- Problem: effectively only the first layer has access to temporal information
- First 2D convolution collapses all temporal information:
 $3T \times H \times W \rightarrow D \times H' \times W'$
- Remaining network is a standard 2D network.

2D Convolution

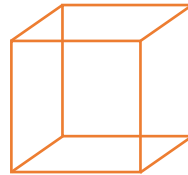


3D Convolution



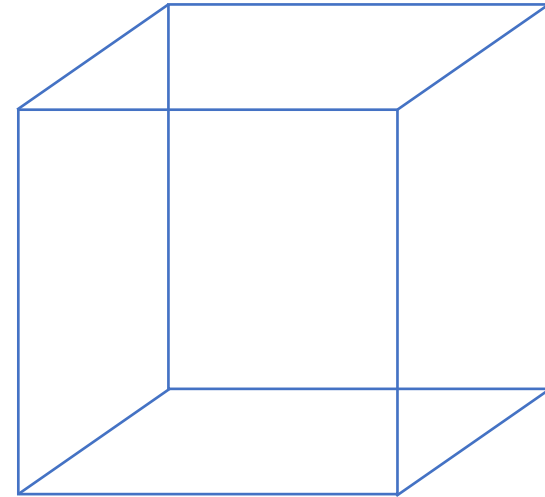
Input
 $C \times D \times H \times W$

*



N Filters
 $C \times d \times h \times w$

=



Output
 $N \times D \times H \times W$

3D Convolution

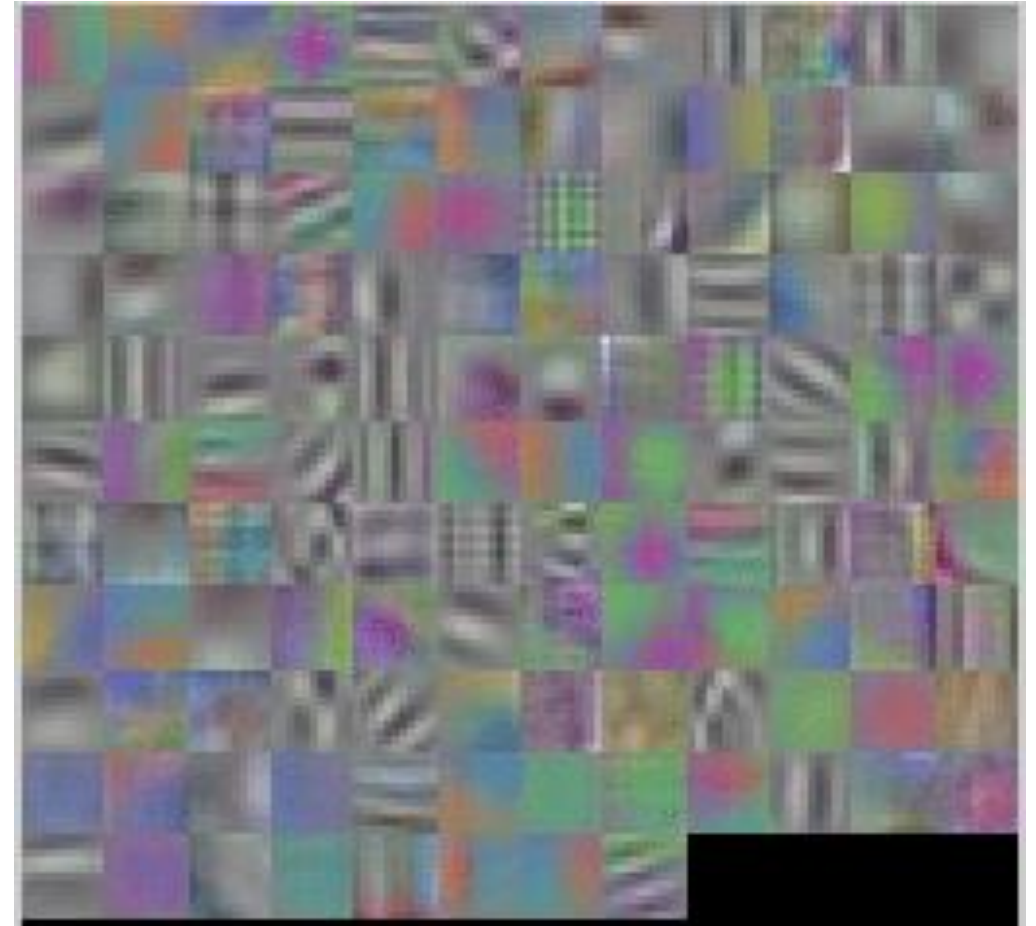
- The convolution now slides along 3 directions: width, height, and depth.
- There is still a channel dimension: each feature in the 3D feature map has C dimensions.
- The output is also a 3D feature map.
- Naming: we ignore the channel dimension. Input and output are actually 4D tensors. (weights: 5D)

Other operations

- 3D Pooling: works the same way. Filter size e.g. 2x2x2
- Activation: works element-wise. (no change)
- Fully connected layer: same as in 2D. Reshape into a vector before applying it. (or use a convolution that has the same size as the feature map)
- 3D CNN: swap 2D operations with 3D operations.

First Layer Filters

- Filters span space and time.
- We can visualise them by animating them through time.
- Moving edge filters.
- Not all filters change with time.

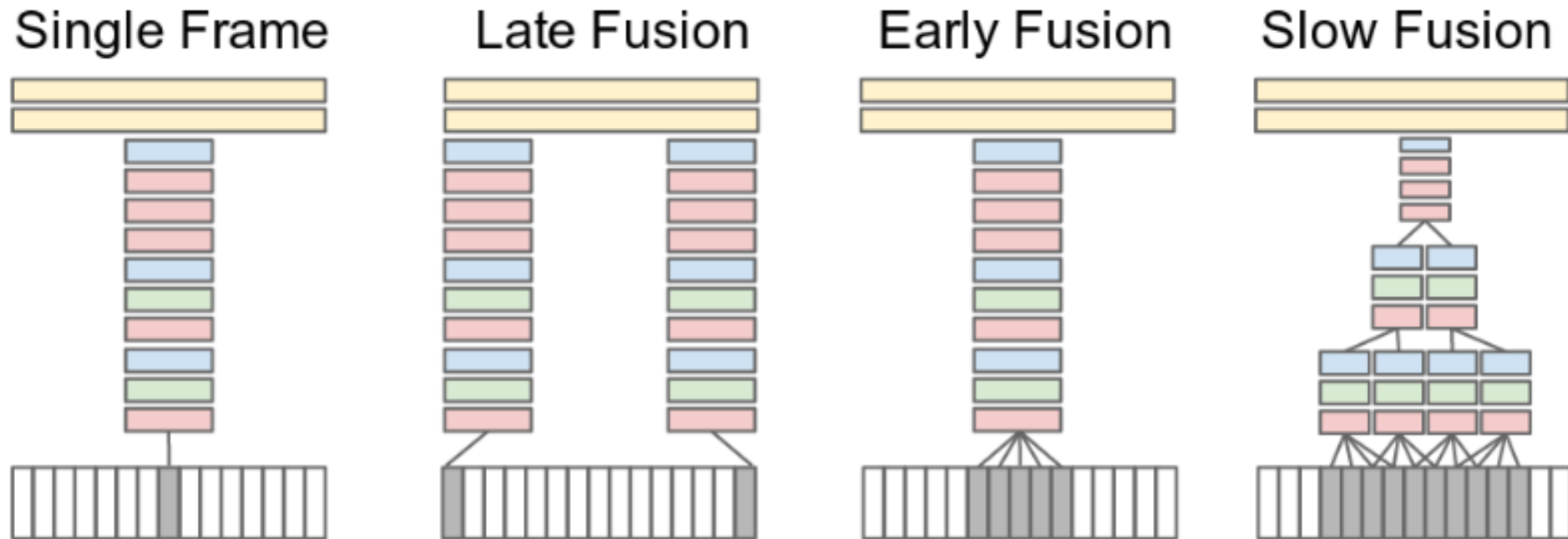


Large-scale Video Classification with Convolutional Neural Networks,
Karpathy et al., 2014

3D CNNs for Video Understanding

- Slow Fusion: slowly fuse temporal information over the course of the network
- Adds shift invariance in time (same motion at a different time).
- Subsampling in space and time gives larger and larger context to each successive layer.













Fusion Approaches



Video Classification Example

GT

pred

 <p>track cycling cycling track cycling road bicycle racing marathon ultramarathon</p>	 <p>ultramarathon ultramarathon half marathon running marathon inline speed skating</p>	 <p>heptathlon heptathlon decathlon hurdles pentathlon sprint (running)</p>	 <p>bikejoring mushing bikejoring harness racing skijoring carting</p>	 <p>longboarding longboarding aggressive inline skating freestyle scootering freeboard (skateboard) sandboarding</p>	 <p>ultimate (sport) ultimate (sport) hurling flag football association football rugby sevens</p>
 <p>demolition derby demolition derby monster truck mud bogging motocross grand prix motorcycle racing</p>	 <p>telemark skiing snowboarding telemark skiing nordic skiing ski touring skijoring</p>	 <p>whitewater kayaking whitewater kayaking rafting kayaking canoeing adventure racing</p>	 <p>arena football indoor american football arena football canadian football american football women's lacrosse</p>	 <p>reining barrel racing rodeo reining cowboy action shooting bull riding</p>	 <p>eight-ball nine-ball blackball (pool) trick shot eight-ball straight pool</p>

Video Classification Example

Single frame is very good
and even better with a
multi-resolution
approach

Model	Clip Hit@1
Feature Histograms + Neural Net	-
Single-Frame	41.1
Single-Frame + Multires	42.4
Single-Frame Fovea Only	30.0
Single-Frame Context Only	38.1
Early Fusion	38.9
Late Fusion	40.7
Slow Fusion	41.9

Slow fusion works better
than early and late
fusion

Motion



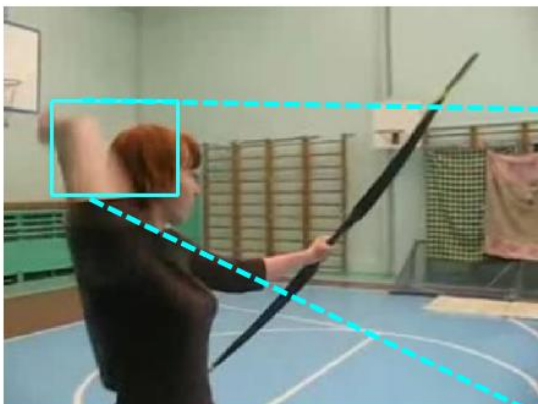
2-DIMENSIONAL MOTION PERCEPTION

2-DIMENSIONAL MOTION PERCEPTION

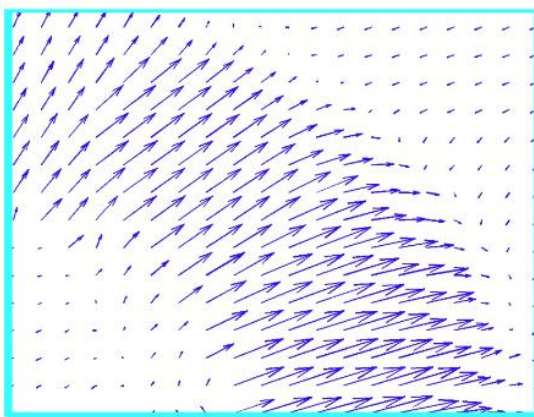


Recap: Optical Flow

Image at frame t



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



Tells where each pixel will move in the next frame:

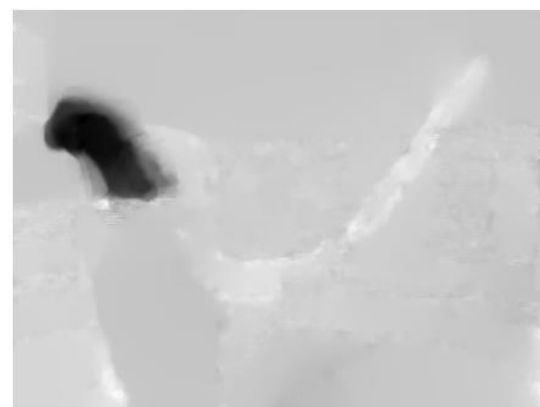
$$F(x, y) = (dx, dy)$$

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

Image at frame $t+1$

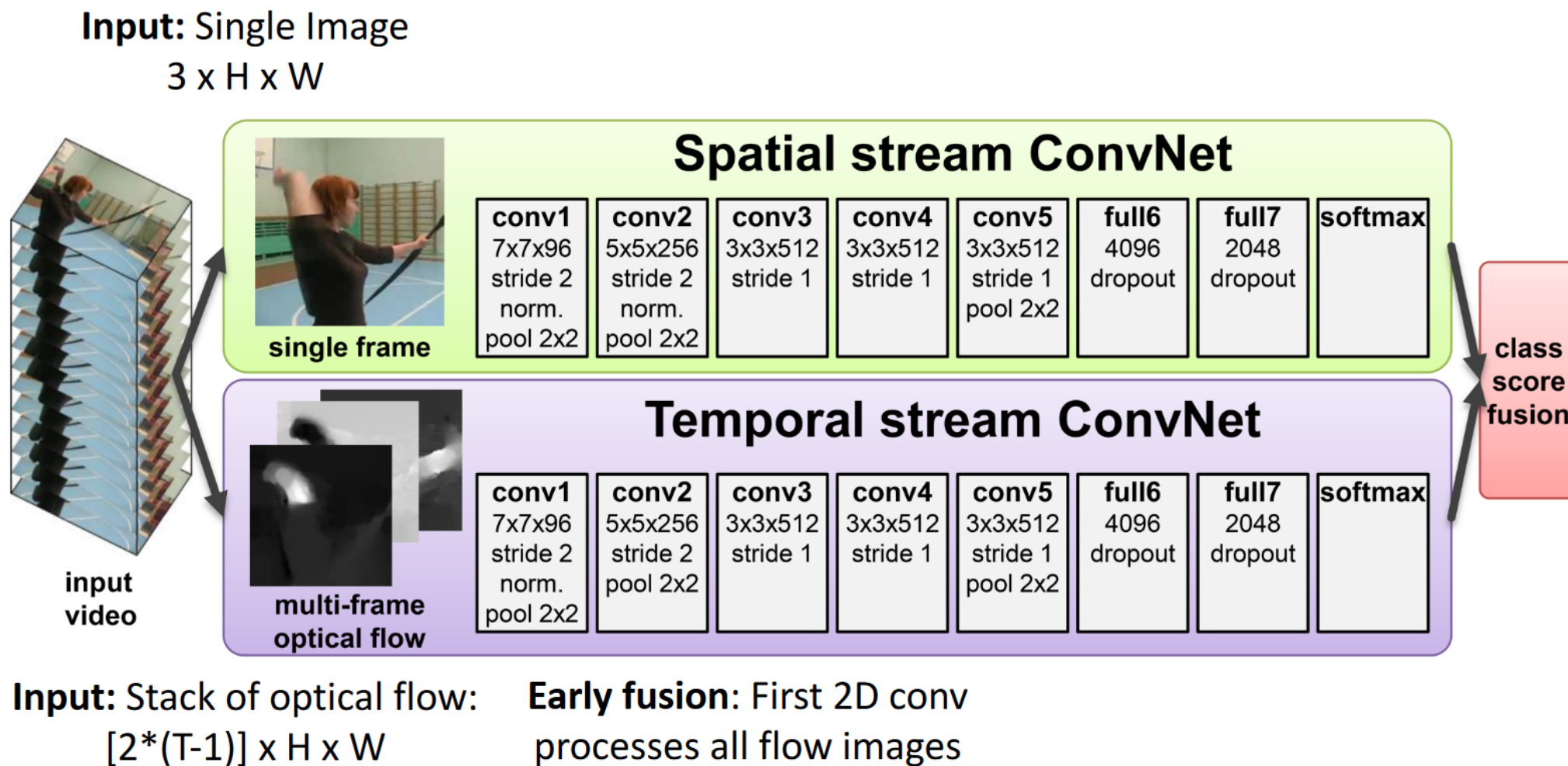


Horizontal flow dx

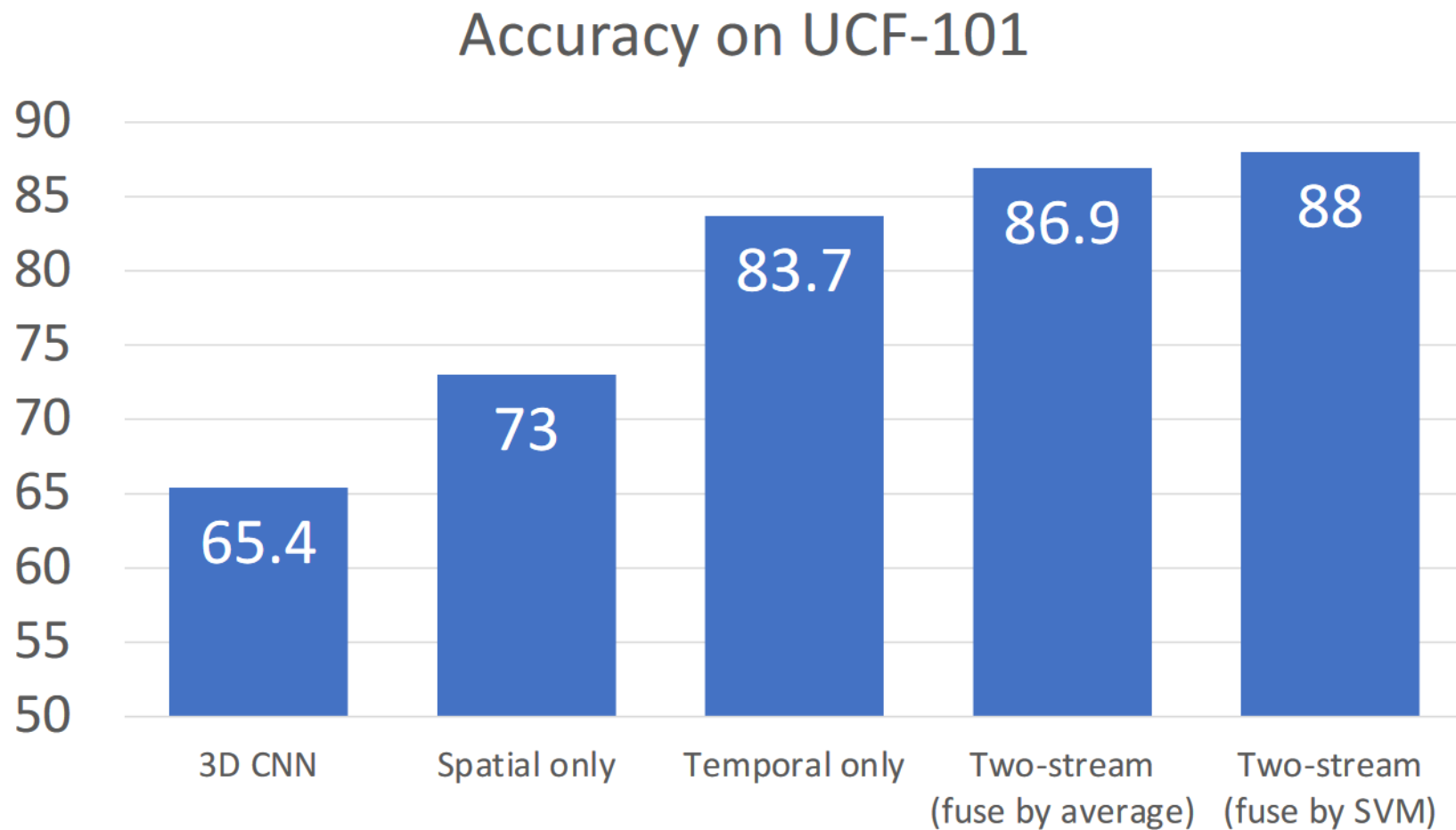


Vertical Flow dy

Action Recognition with Optical Flow



Two Stream Networks

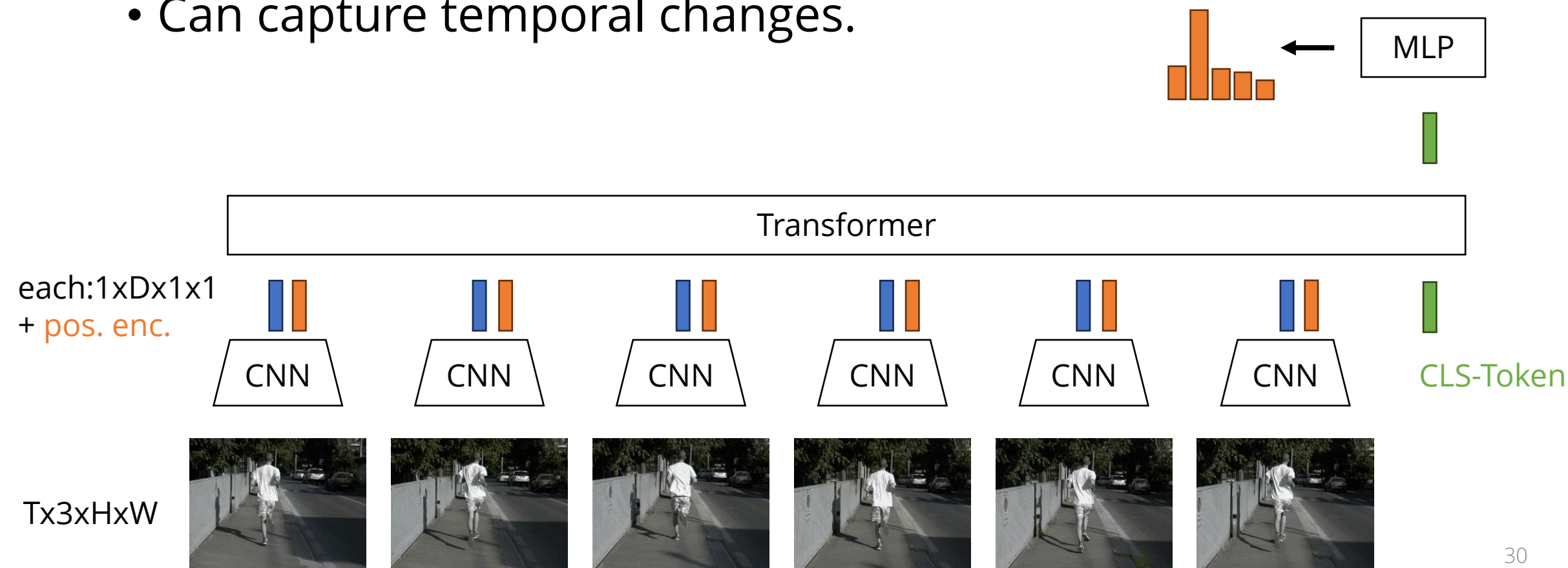


Two Stream Networks

- Can be used to fuse different modalities:
 - RGB and optical flow
 - Image and Audio
 - Image and Text
 - ...
- Typically late(-ish) fusion: process each modality separately before fusing.

Longer Temporal Context

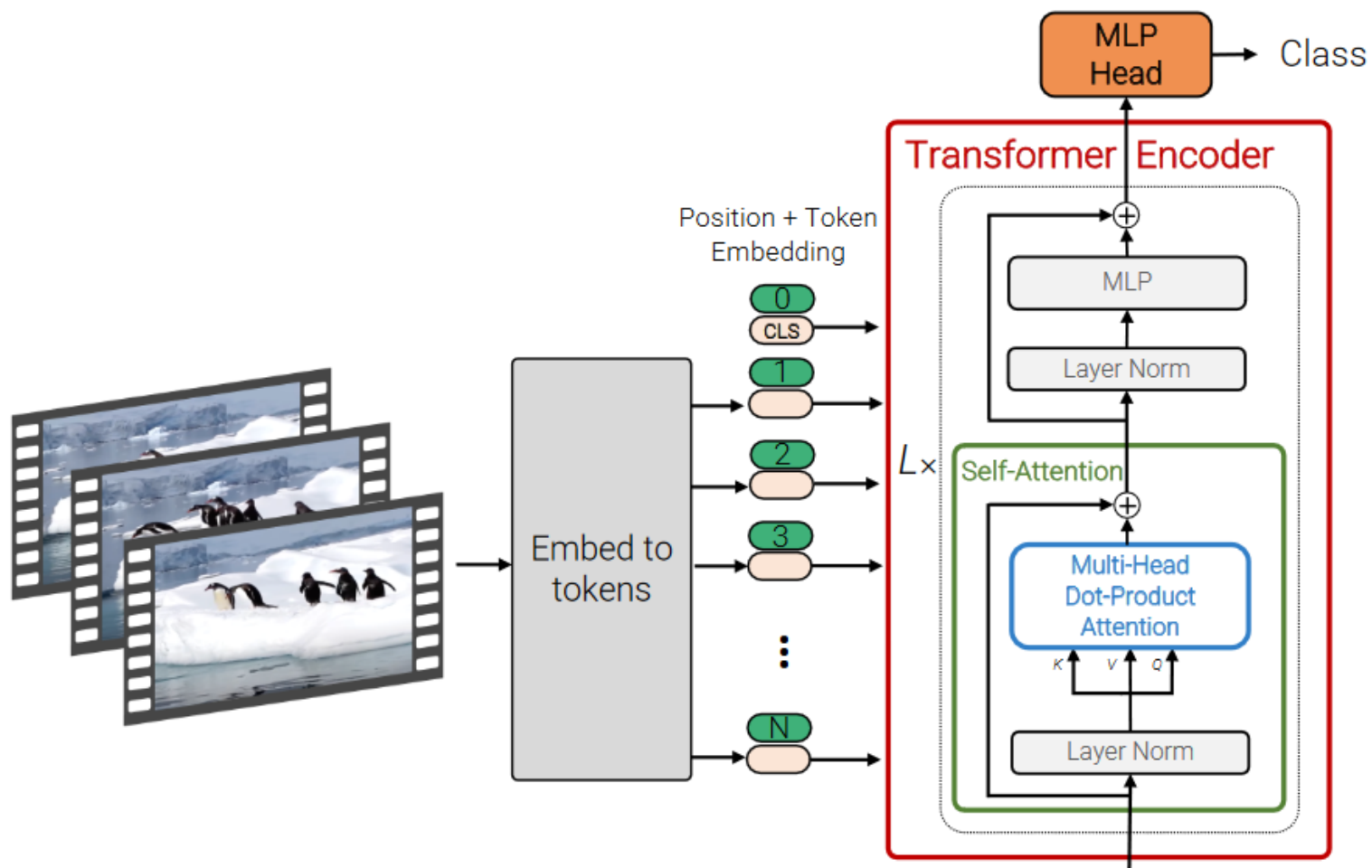
- Use a sequence model (e.g. Transformer) across time.
- Can capture temporal changes.



Longer Temporal Context

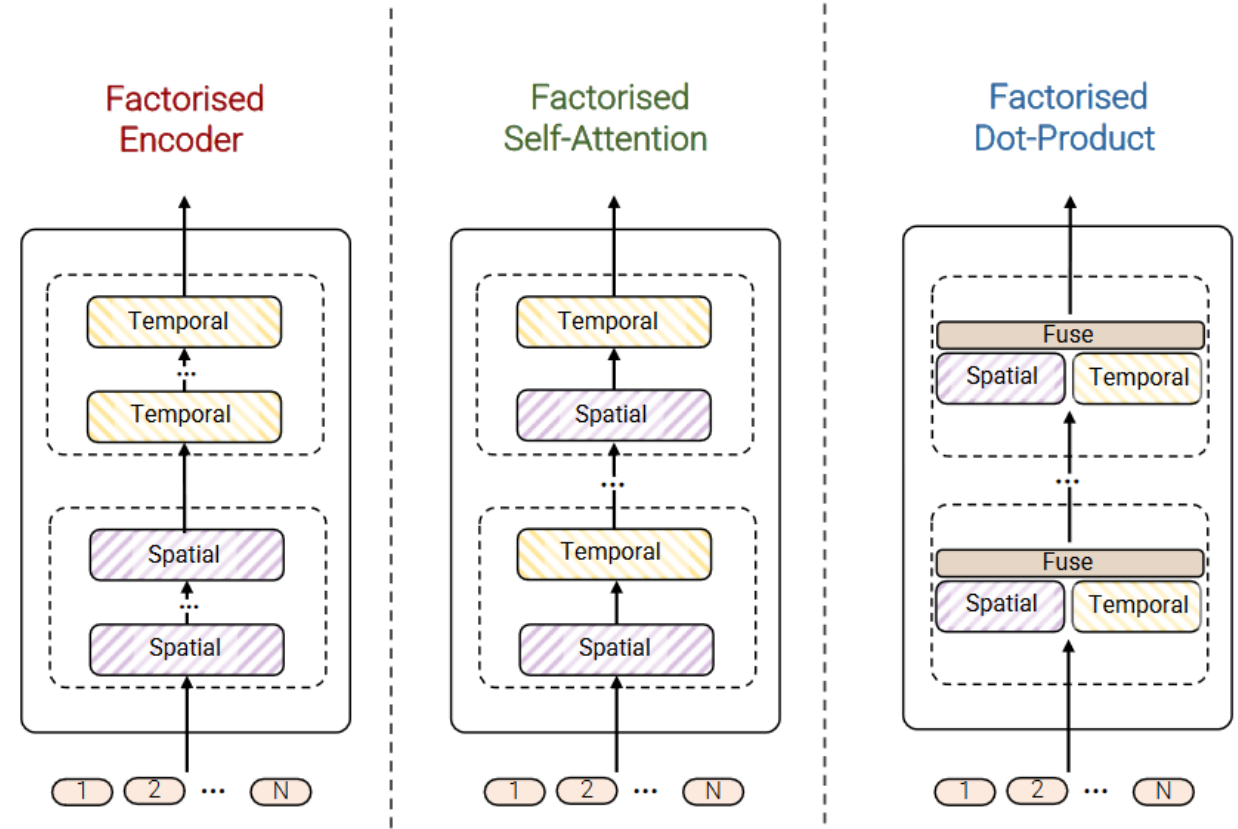
- Encoder can be per frame, or other architectures e.g. 3D CNN: produce temporal features.
- Sequence model reasons across time.
- Hybrid architecture: CNN for processing clips, transformer for processing video (composed of clips).
- Pure transformer architectures?

Transformer for Video Understanding



Video Vision Transformers

- Spatio-temporal tokens: each token comes from a “patch” in space and time.
- Attention is expensive $\mathcal{O}(n^2)$.
- Factorised attention: alternate spatial and temporal attention.



Factorised Attention

- Attention: input/output $B \times N \times D$ (batch, tokens, channels)
- Samples in batch dimension are processed separately.
- Our input: $B \times T \times D \times H \times W$
- Idea: use batch dimension to process dimensions we do not want to include in the attention.

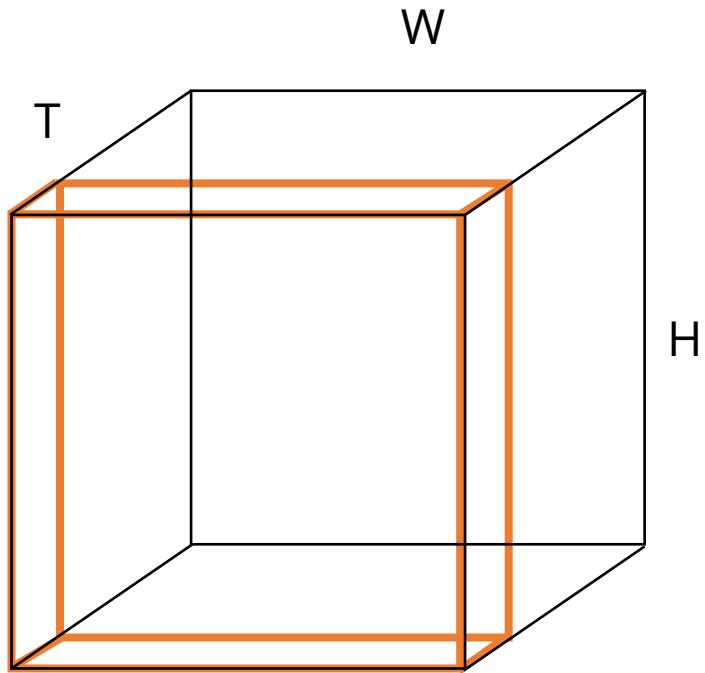
Temporal Attention

- Process each sample in the batch separately
- Process each spatial location separately
- Input $B \times T \times D \times H \times W$
 - Permute: $B \times H \times W \times T \times D$
 - Flatten: $BHW \times T \times D$
 - Attention: $BHW \times T \times D \rightarrow BHW \times T \times D$
 - Unflatten: $B \times H \times W \times T \times D$
 - Permute: $B \times T \times D \times H \times W$

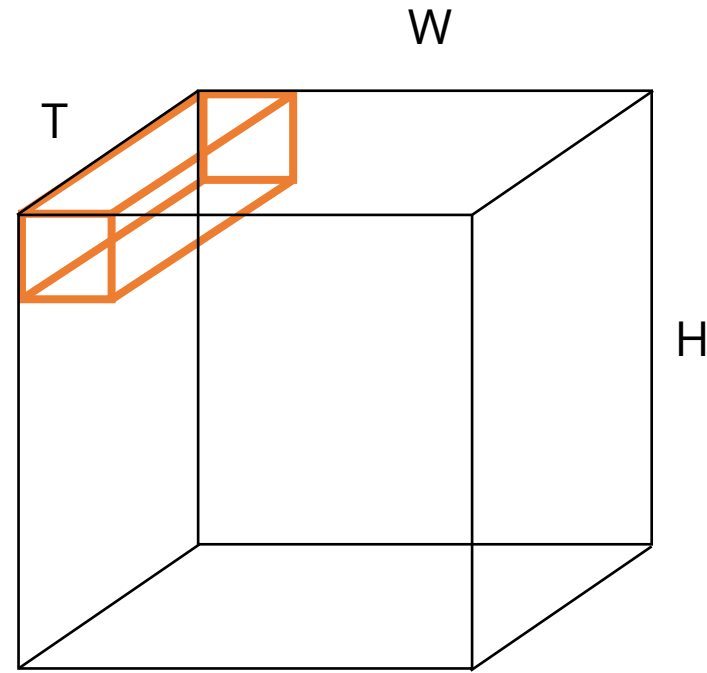
Spatial Attention

- Process each sample in the batch separately
- Process each time step separately
- Input $B \times T \times D \times H \times W$
 - Permute: $B \times T \times H \times W \times D$
 - Flatten: $BT \times HW \times D$
 - Attention: $BT \times HW \times D \rightarrow BT \times HW \times D$
 - Unflatten: $B \times T \times H \times W \times D$
 - Permute: $B \times T \times D \times H \times W$

Factorised Attention



Spatial Attention



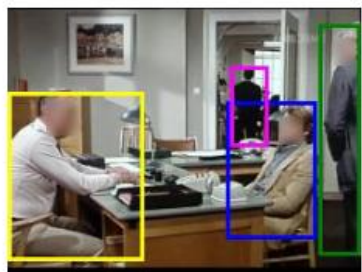
Temporal Attention

Factorised Attention

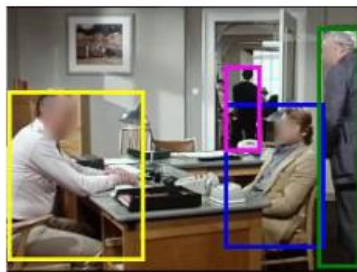
- Alternating spatial and temporal attention layers.
- Allows propagating information across time and space.
- Complexity:
 - Full attention: $\mathcal{O}((THW)^2)$
 - Factorised attention: $\mathcal{O}(T^2HW + T(HW)^2)$

Task: Spatio-Temporal Detection

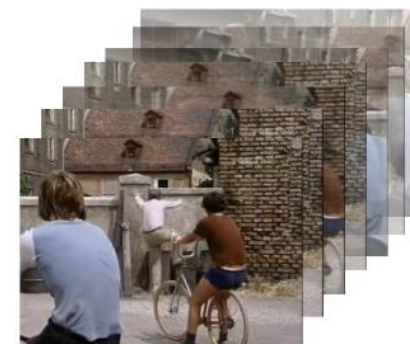
- Detect objects in a video and predict their actions.
- Tubelets: bounding box with time.



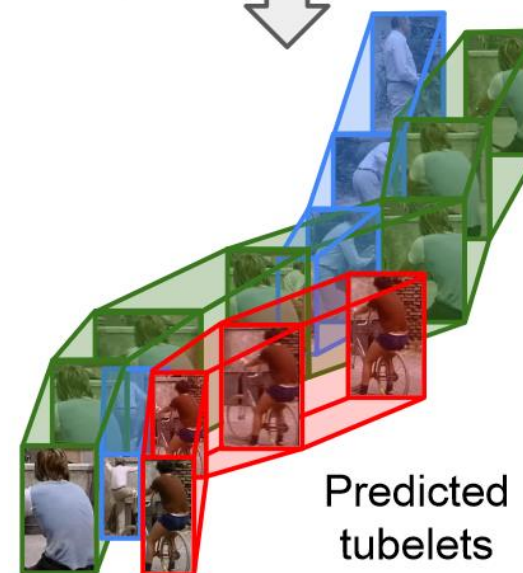
sit, talk to, watch, touch
watch, listen to, sit



stand, watch, listen to
walk, watch, listen to

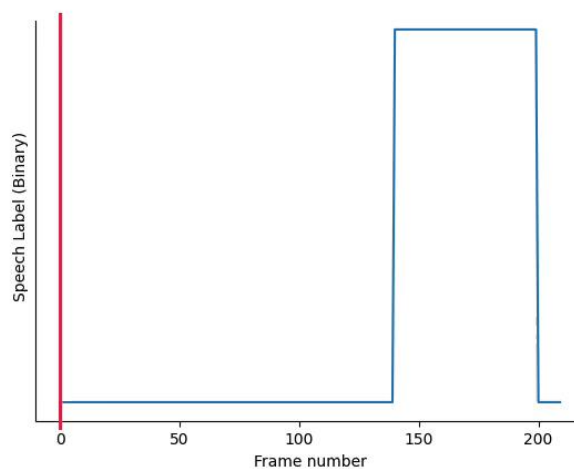


End-to-end
model



Task: Lip Reading

Step 1: Visual Speech Detection



Lip Reading with Attention



Figure 2. Visualization of the visual attention masks α from the VTP module superimposed on the input frames that produce them. The video clips used here are random samples from the LRS3 dataset. It is evident that the model follows the more discriminative mouth region.

Task: Audio Description

- Multi-modal task: audio-visual input -> text
- Difficult: long-range context understanding



Multi-Modal Learning

- Often, we have multiple modalities: image, audio, text, sensor signal, 3D, temperature, etc.
- Fusion-based architectures are good to combine different input types.
- Process each modality separately into a common shape.
- Fuse and process jointly.

Processing Videos

- Expensive task even after all the subsampling.
- High compute and memory cost.
- Many tasks: image models work surprisingly well.
- Architectures: combine image and time understanding.