



# Deep video models

Viorica Pătrăucean, Research scientist



DeepMind



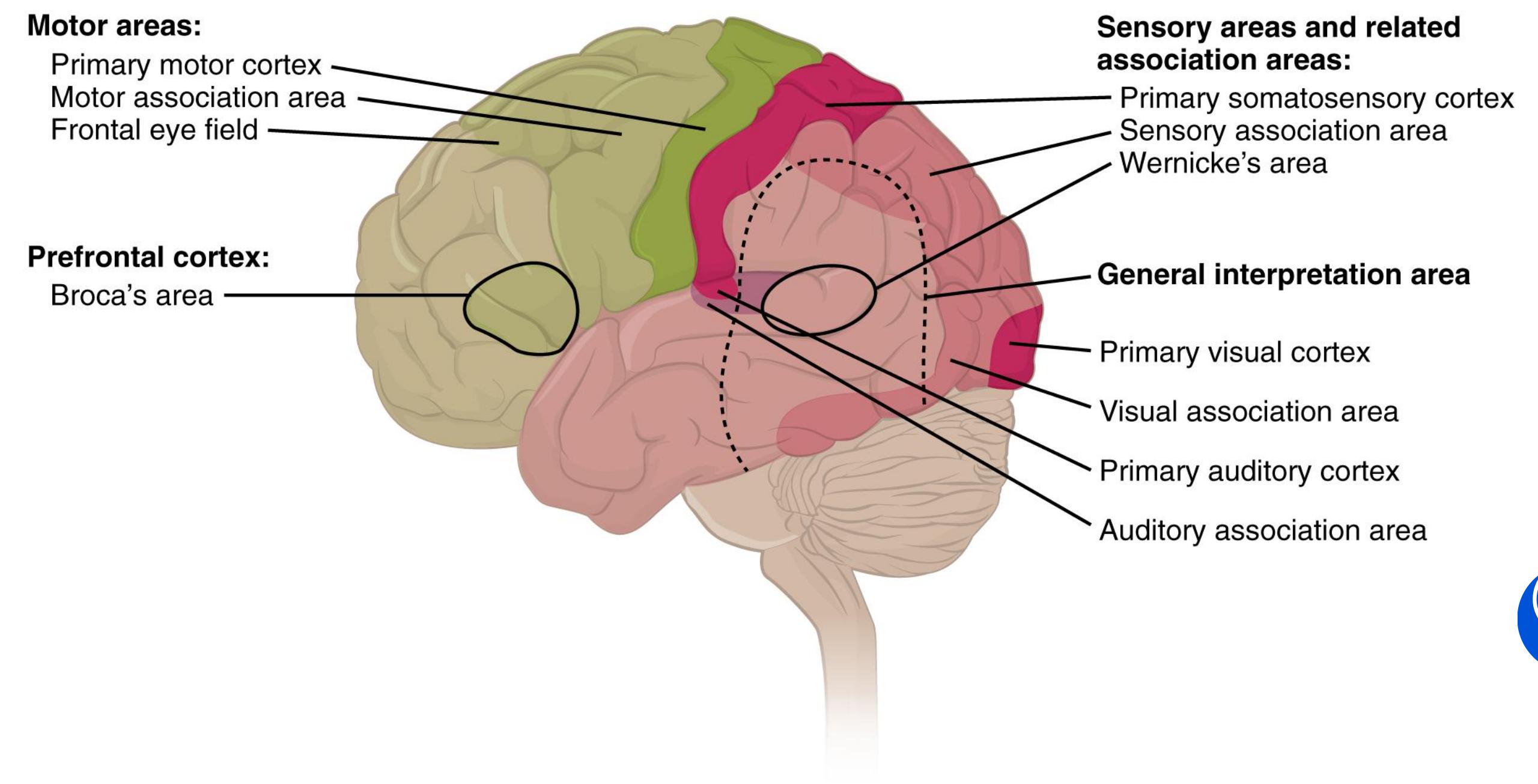
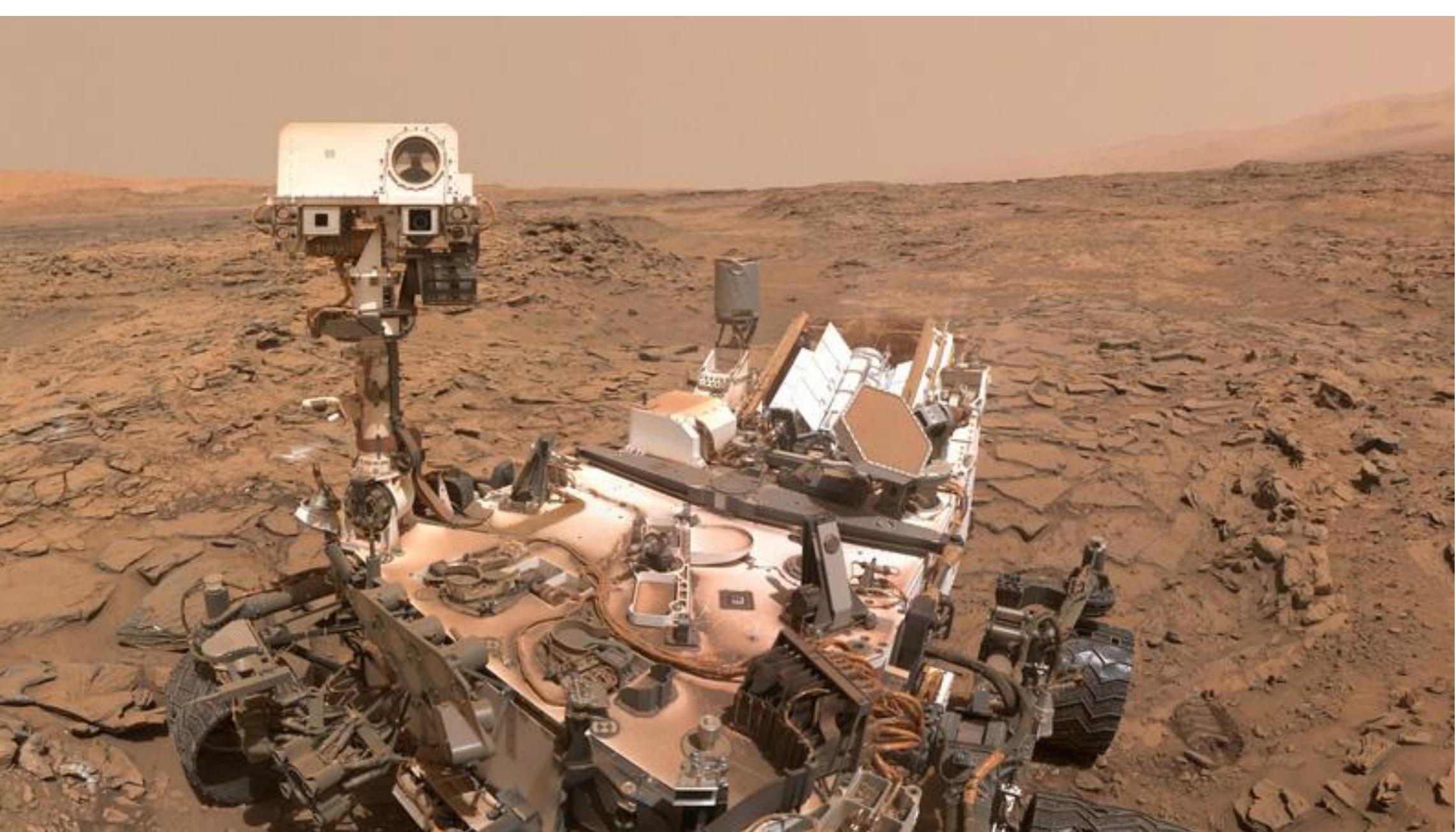
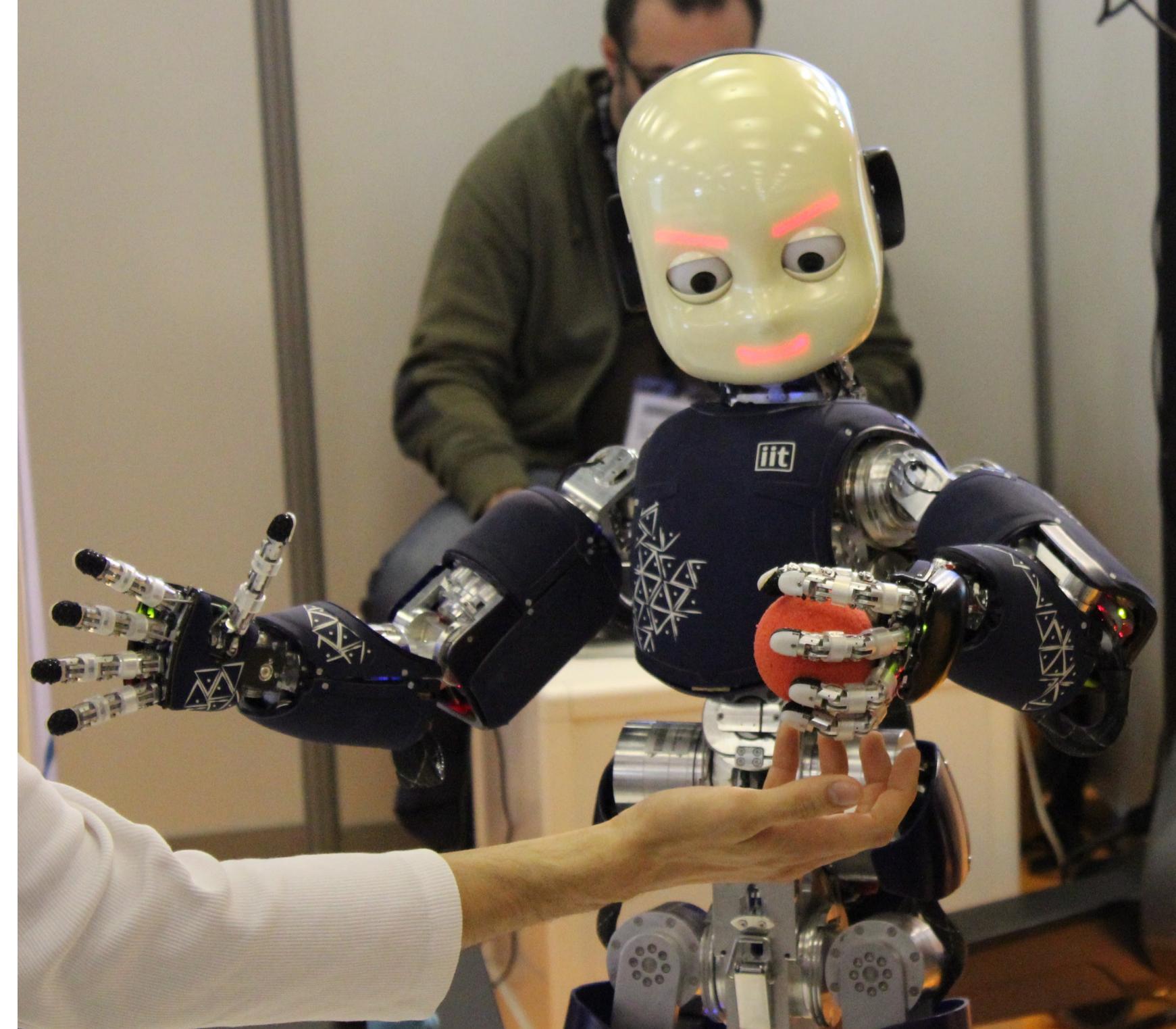
EEML

Vilnius Machine Learning Workshop 2021



# About myself

- Undergrad in Bucharest Romania – Computer Engineering
- Master's and PhD in Toulouse France – Applied Maths (2D image processing)
- PostDoc in Paris and Cambridge – Research on 3D shapes and videos
- Research scientist in DeepMind since 2016
- EEML organiser since 2017



# Outline

**01**

Overview

Why videos?

Popular video models

Efficiency challenges

**02**

Learning paradigms

Multimodal learning

Self-supervised learning

**03**

Depth parallelism

Inference

Training

Biological plausibility

**04**

Conclusion

Summary

My vision for vision

*Disclaimer:* There are many concepts and works in this space. I will cover only a subset of them.

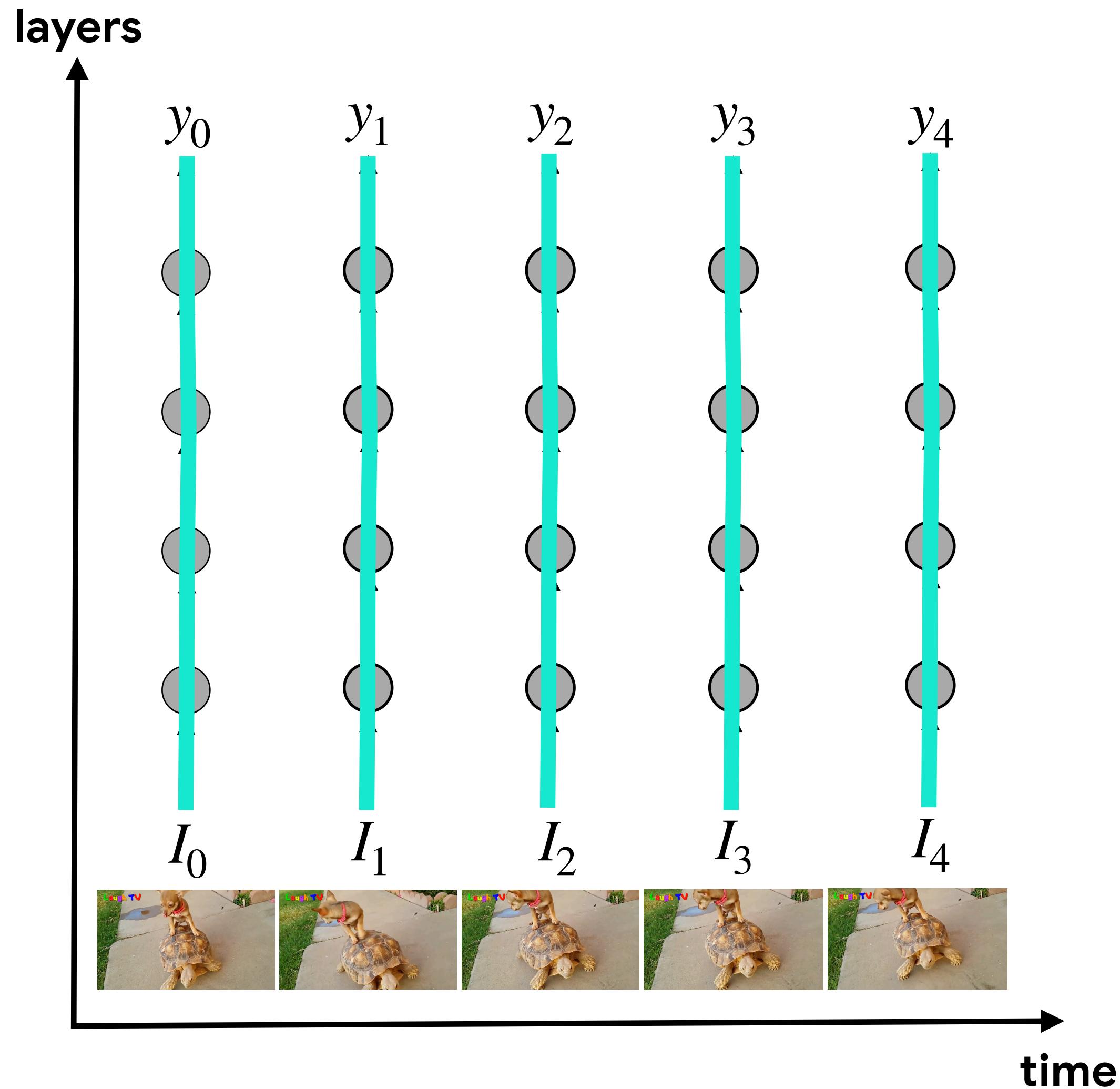


# 01

## Overview



# Why care about videos?



Frame-by-frame video model

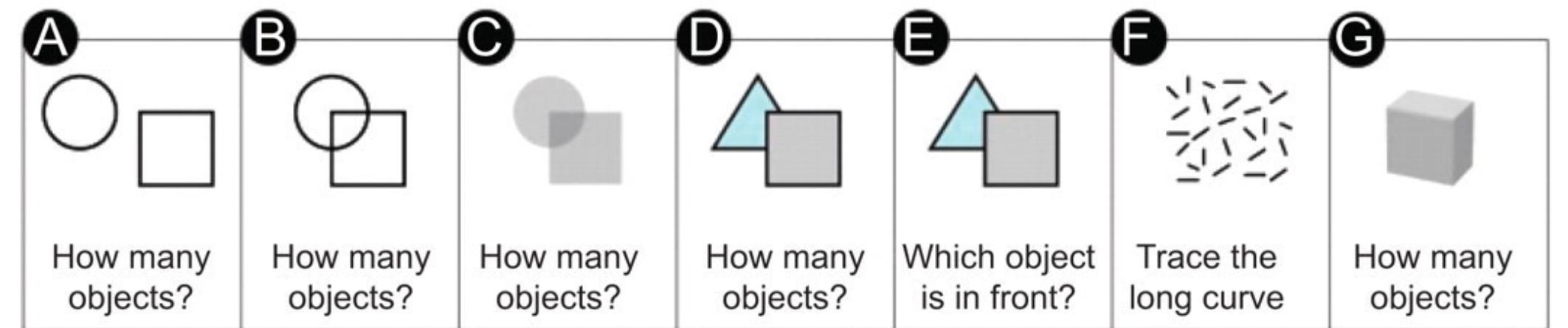
- Inference: run image model on consecutive frames
- Image model trained on single video frames
- Cannot integrate temporal information, no motion features
- Cannot exploit temporal redundancy to improve efficiency



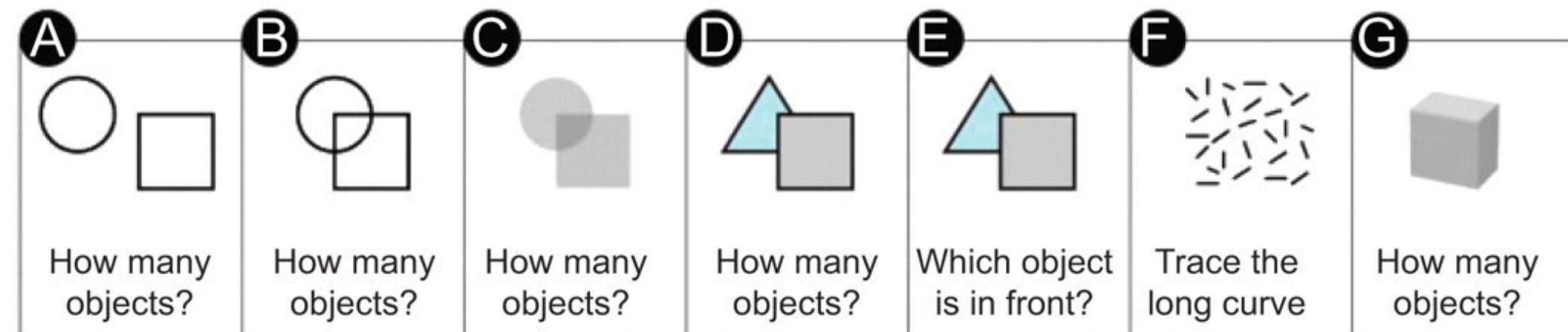
# 1. Frame-by-frame model



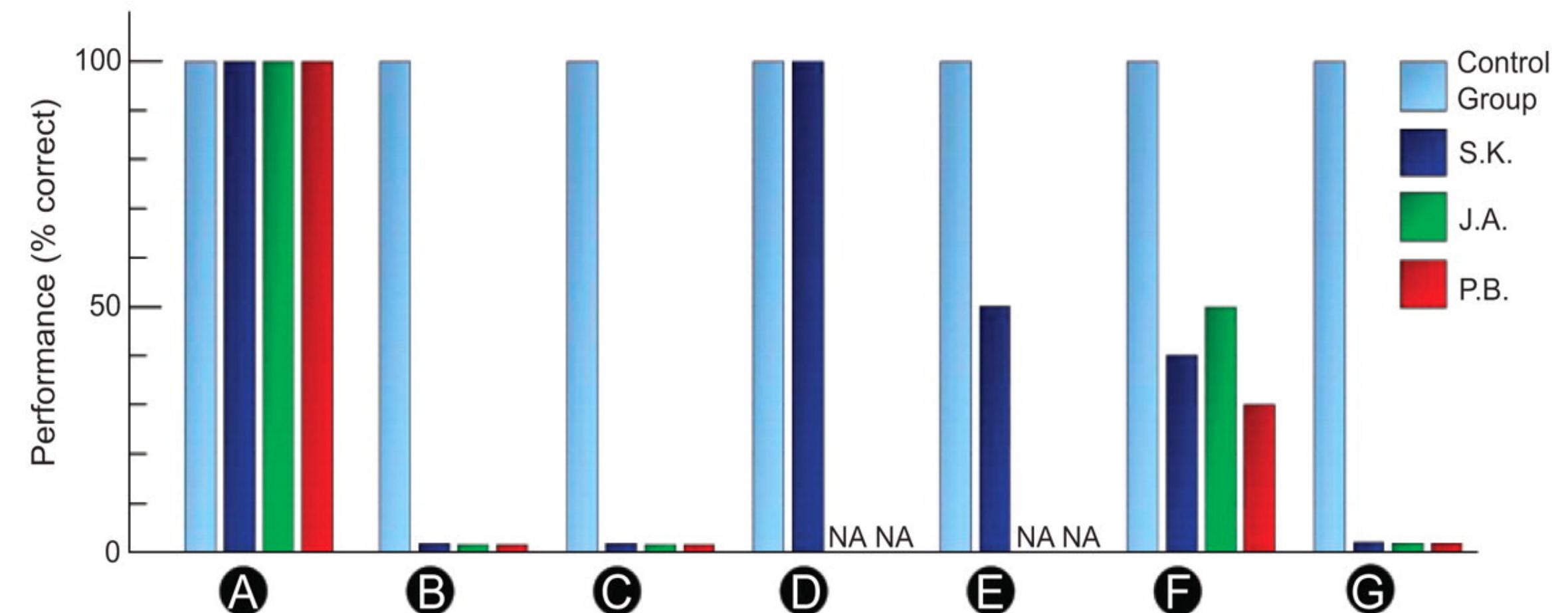
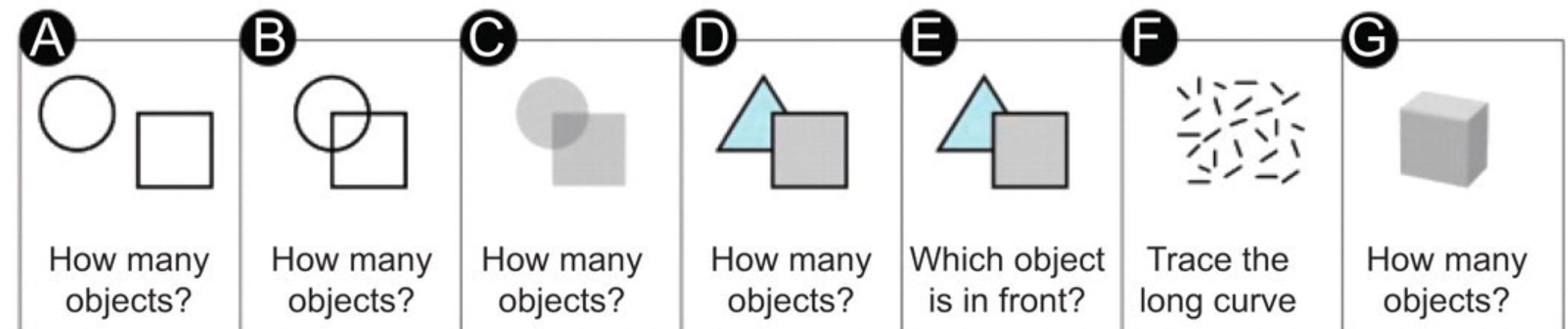
# Importance of motion



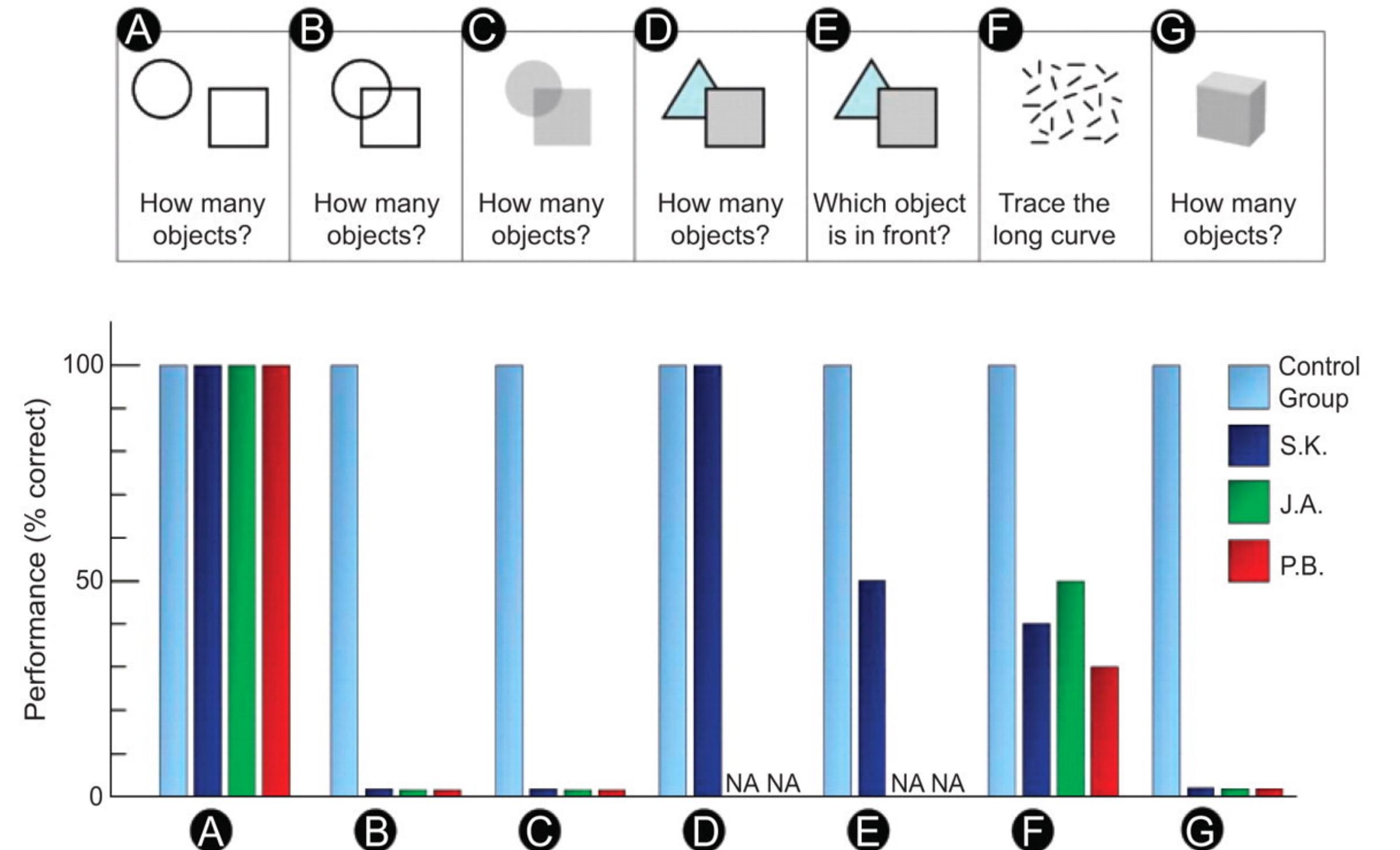
# Importance of motion



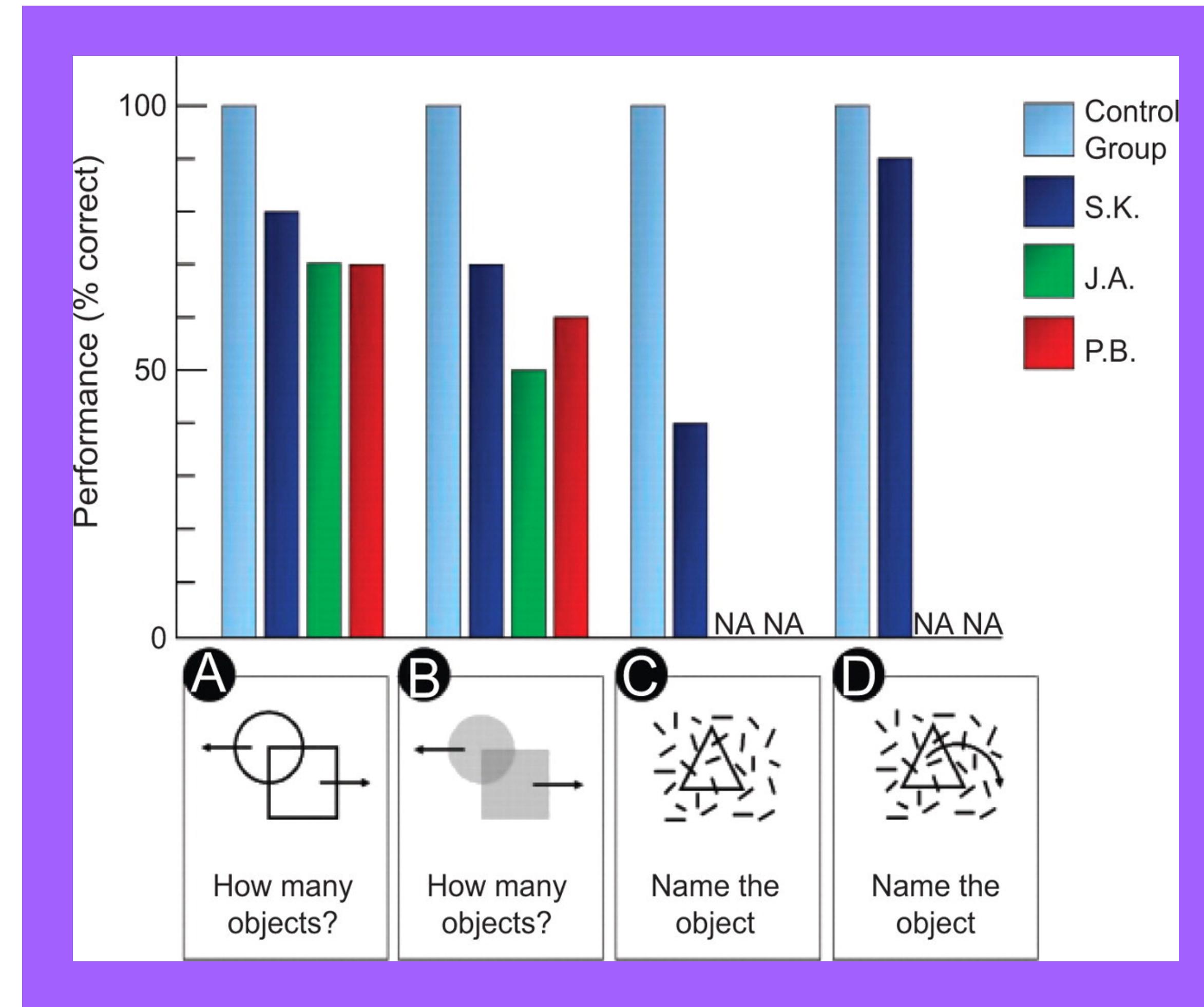
# Importance of motion



# Importance of motion



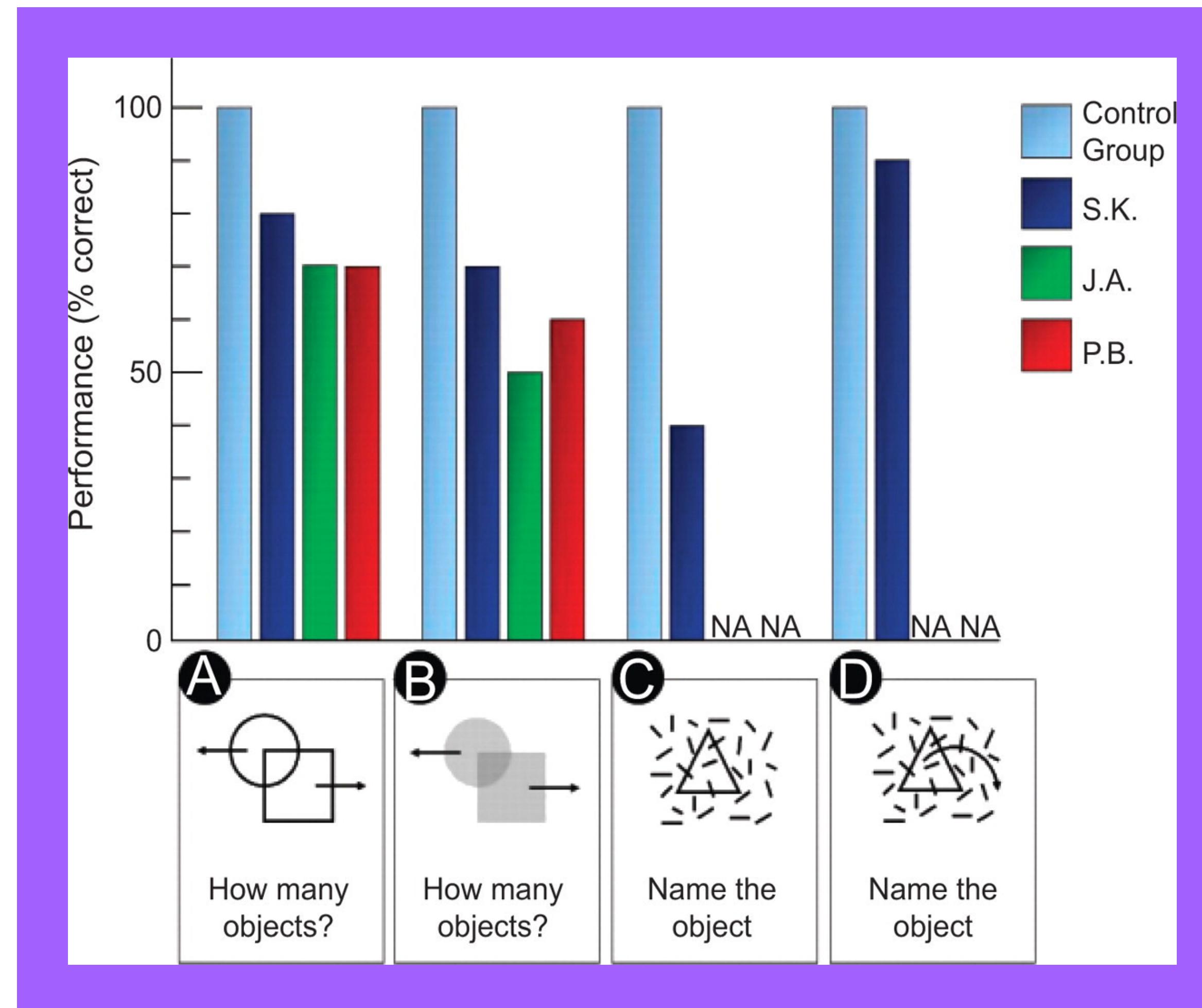
# Importance of motion



Visual Parsing After Recovery From Blindness, Ostrovsky et al, 2009



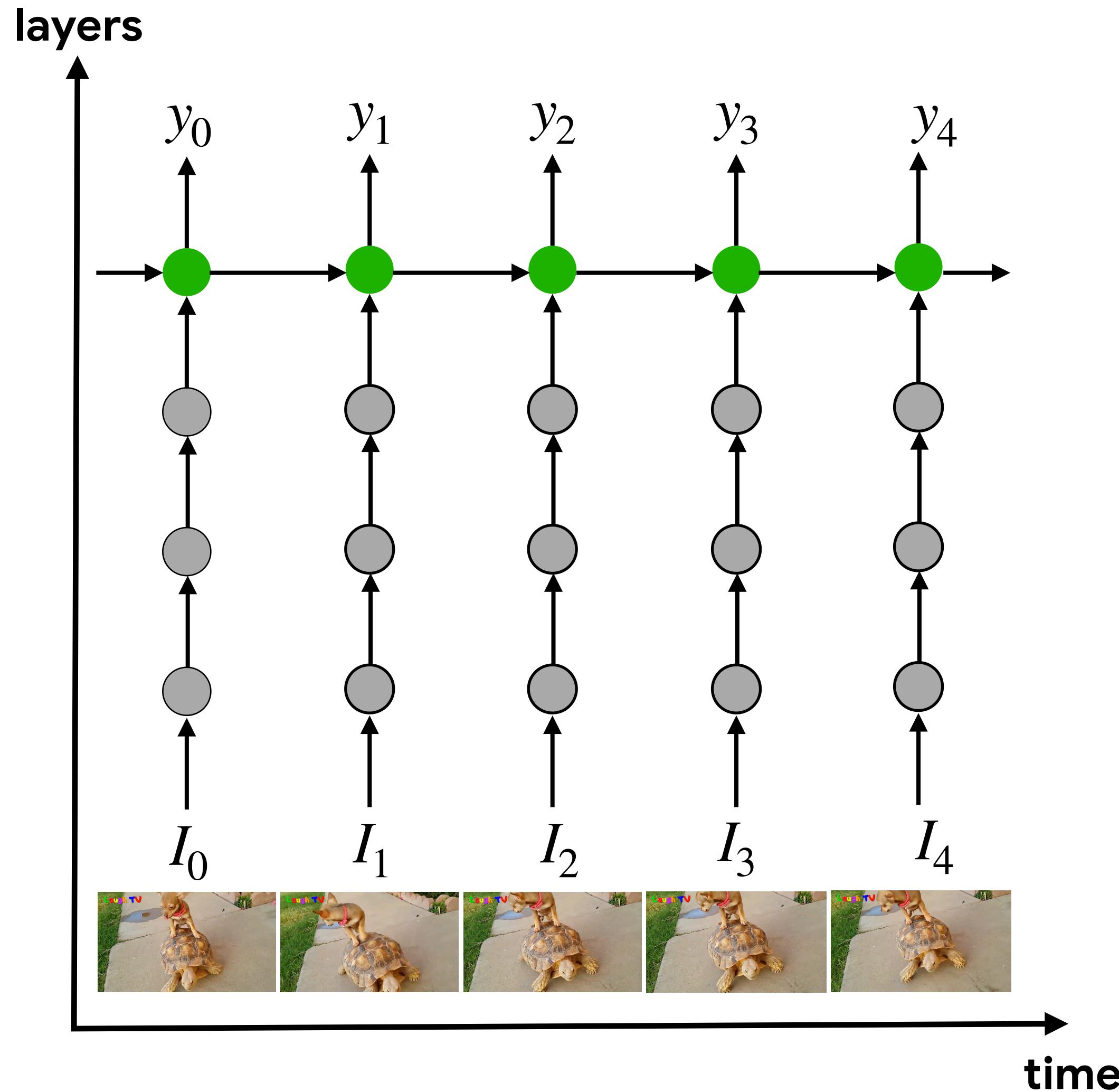
# Importance of motion



Motion helps object recognition when learning to see.



## 2. Image encoder + RNN



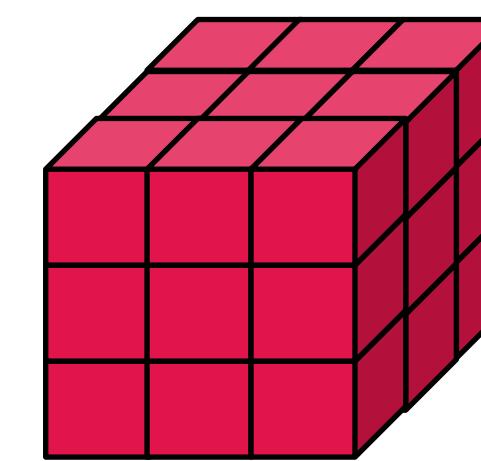
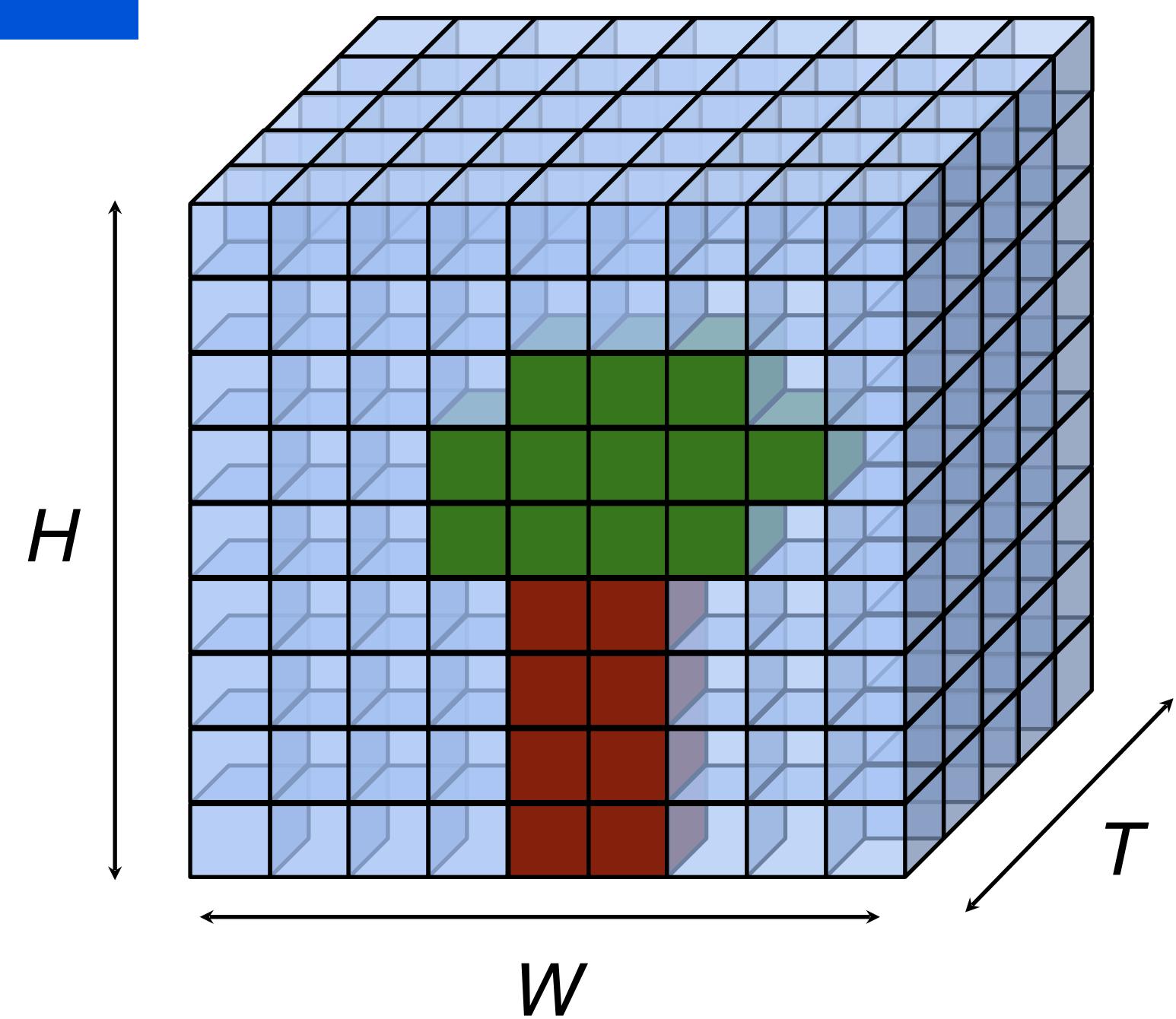
- Image encoder: standard image model (pre-trained on Imagenet)
- RNN: LSTM / GRU fully connected or convolutional
- Need large memory for training with BPTT



### 3. 3D convolutional models

Video as a volume

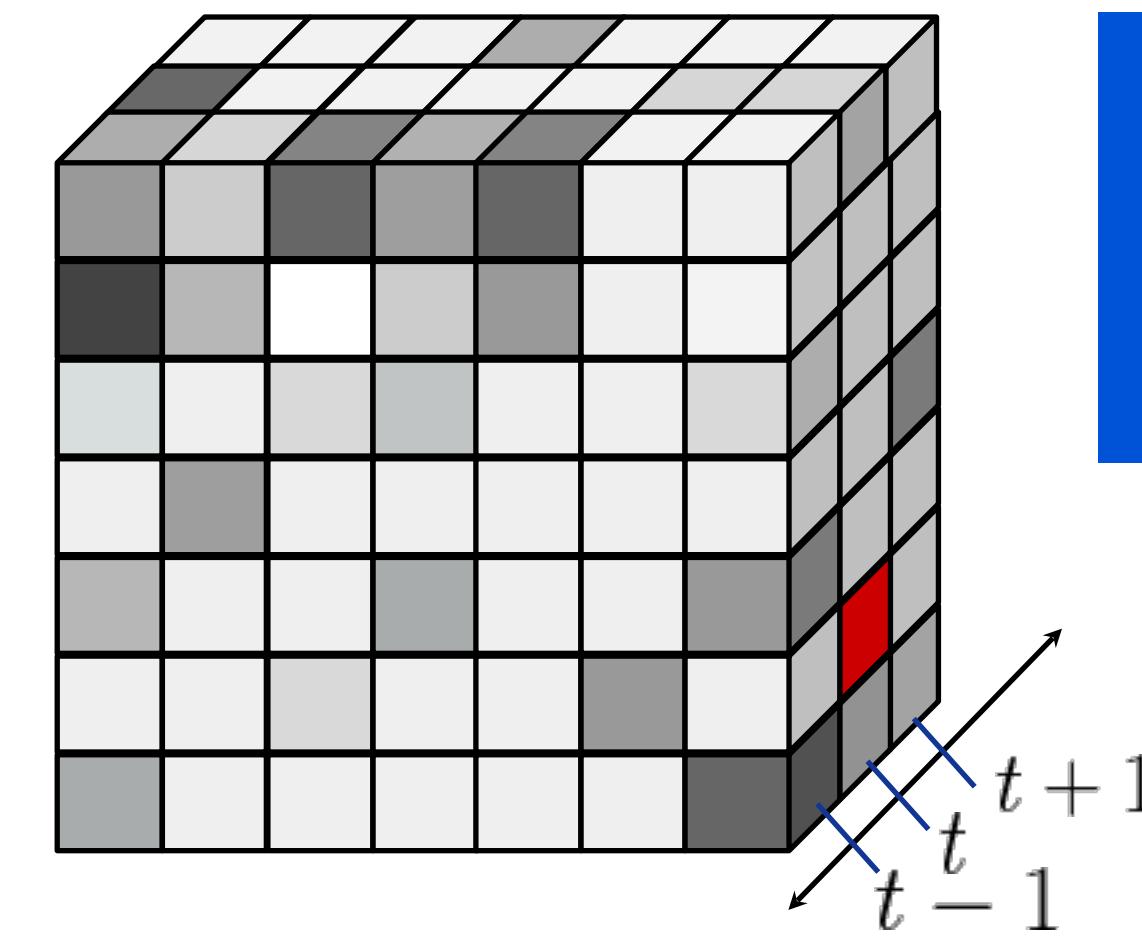
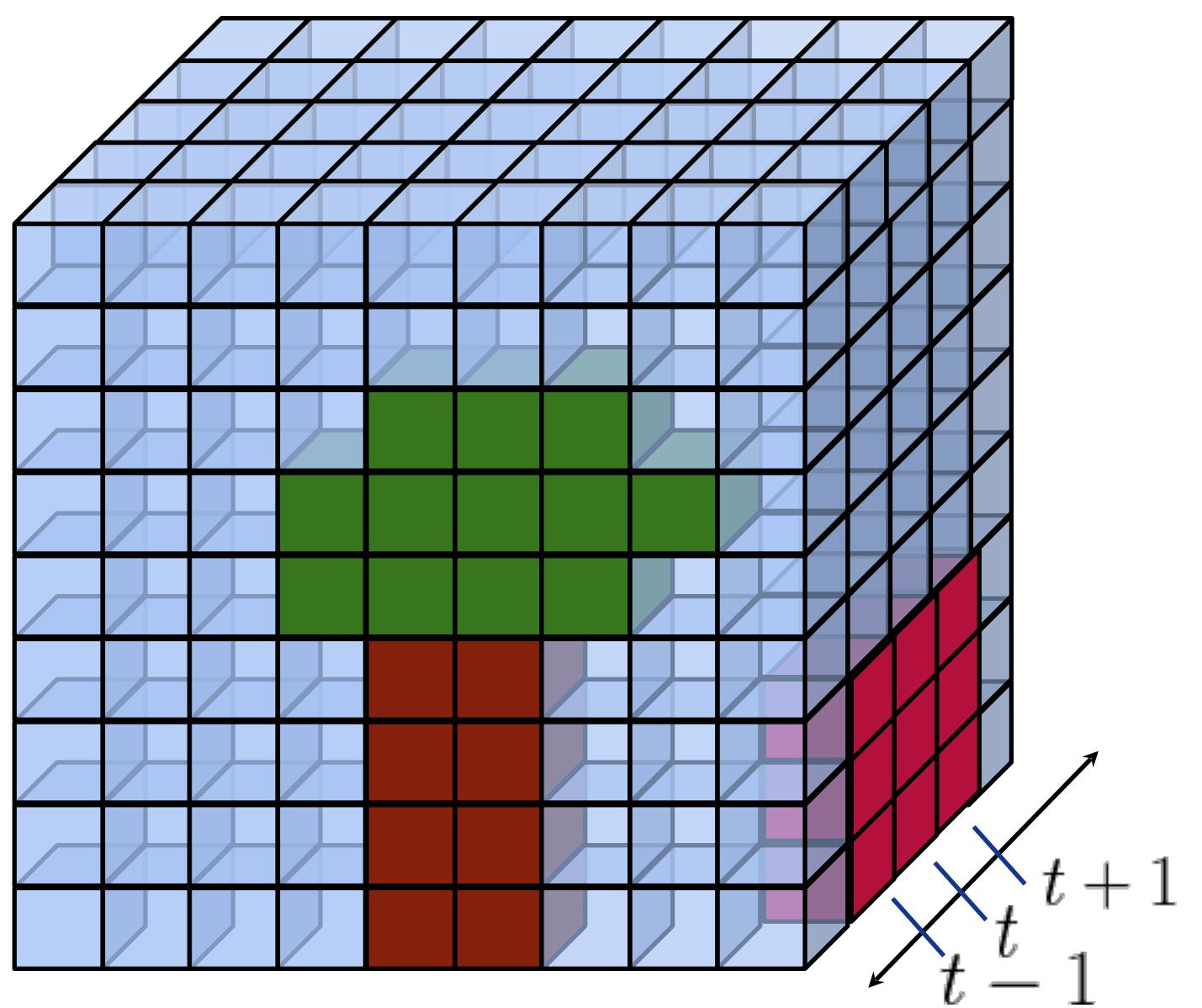
- ⇒ stack frames  $T \times H \times W \times 3$
- ⇒ apply 3D convolutions



$$y = \sum_{i \in 3 \times 3 \times 3} \mathbf{w}_i \mathbf{x}_i + b$$



### 3. 3D convolutional models

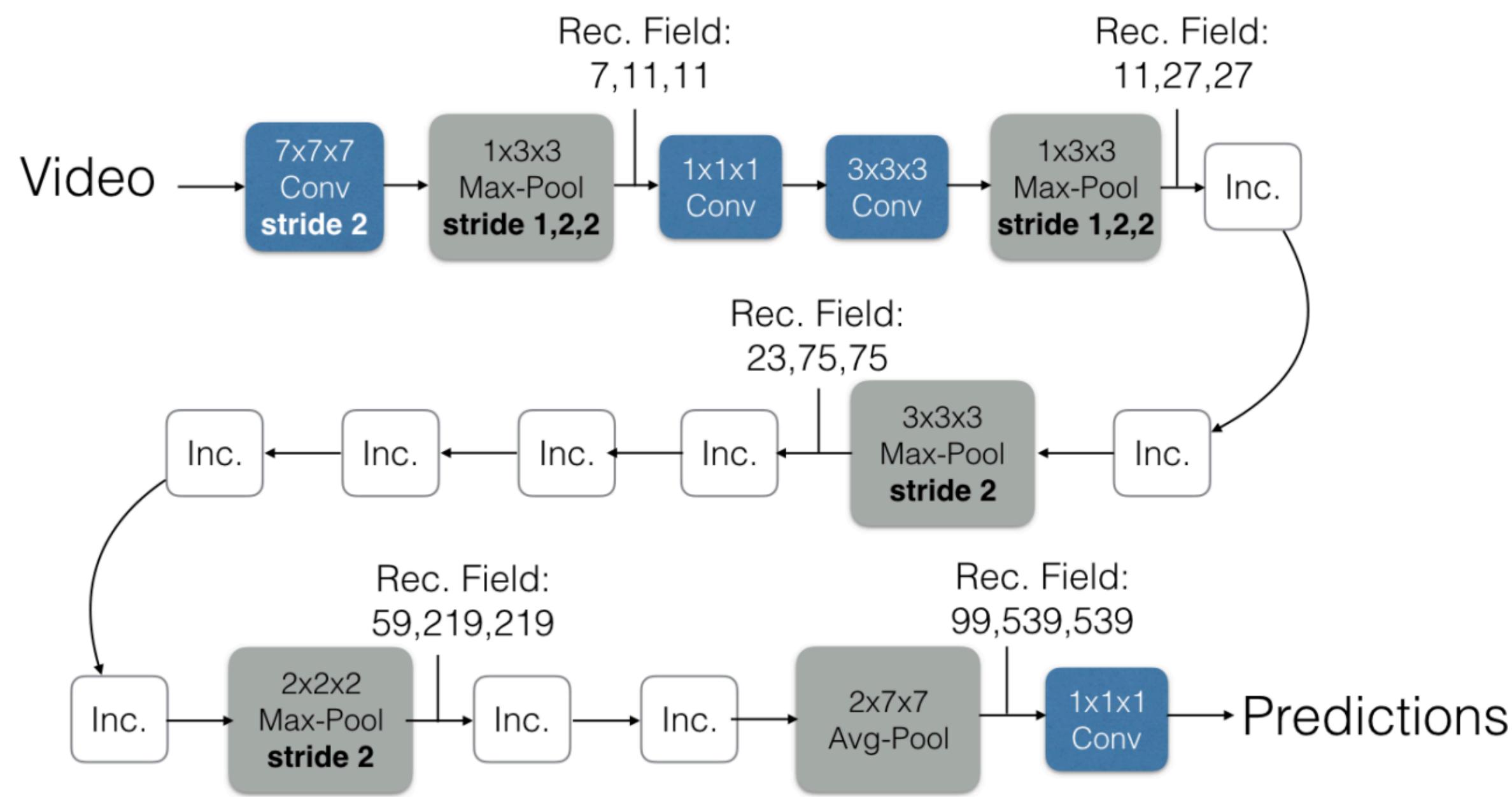


→ 3D convolutions are non-causal  
→ masked 3D convolutions are causal



# 3. 3D convolutional models

## Inflated Inception-V1



- 3D convnet
- Temporal striding
- Possibly inflated 2D model  
(transfer learning)



# Models for action recognition

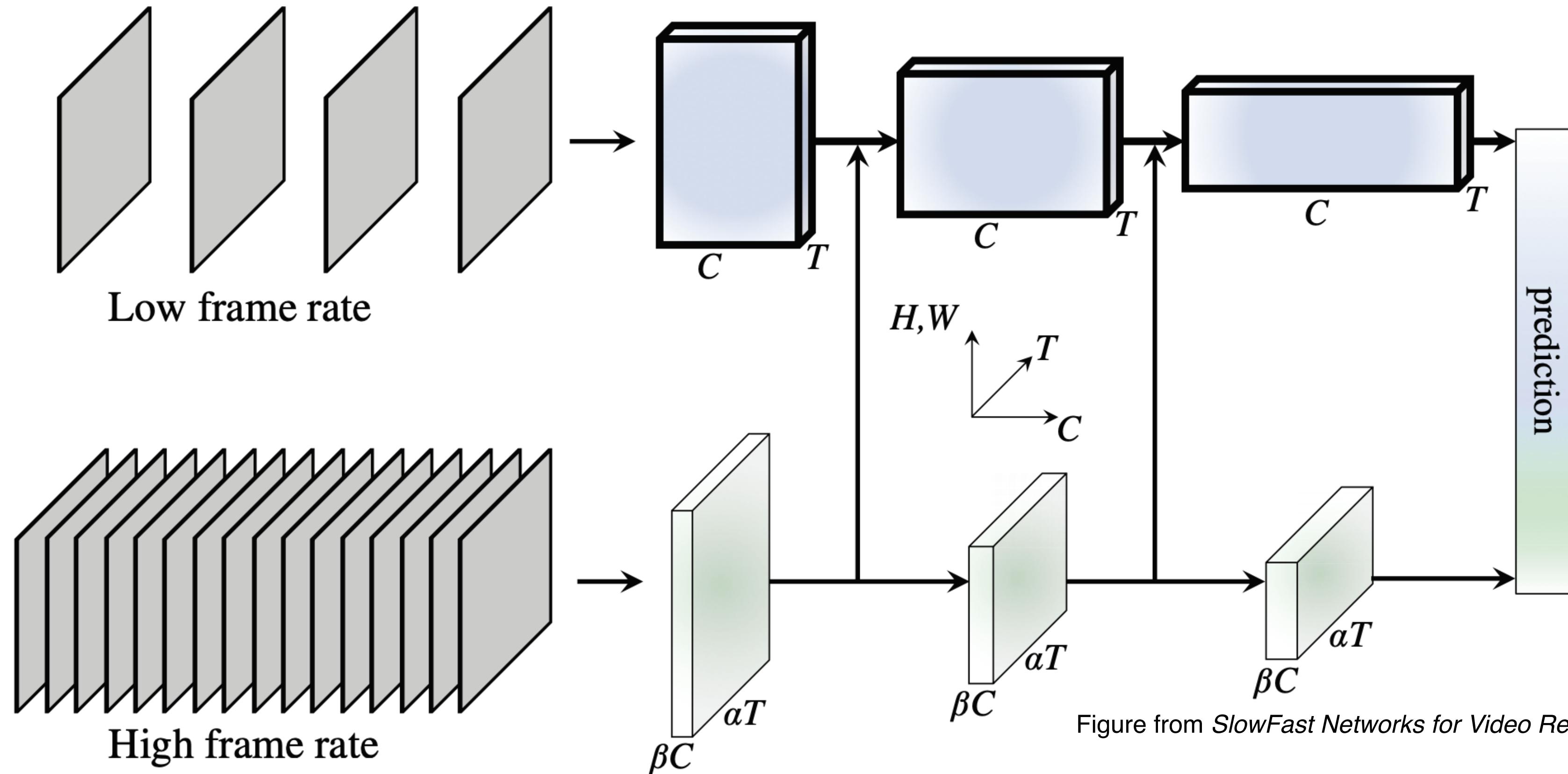
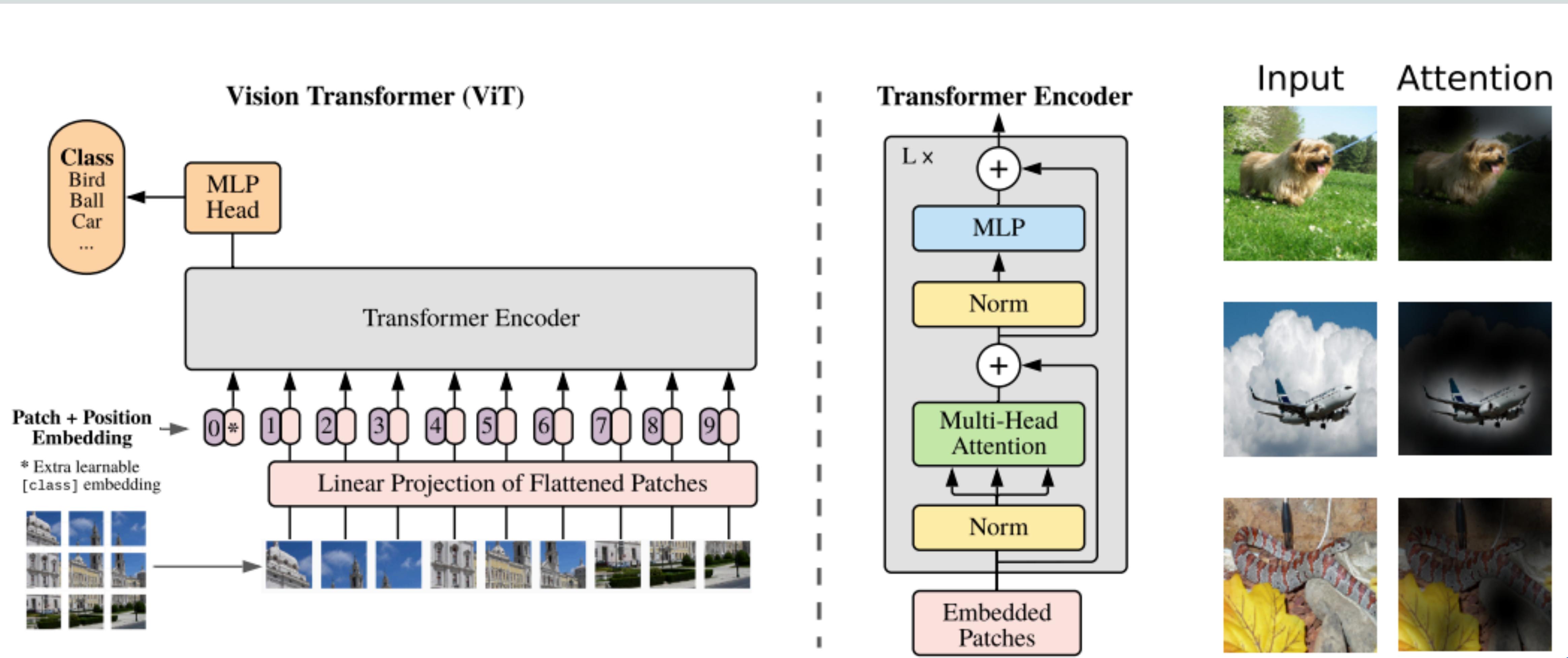


Figure from *SlowFast Networks for Video Recognition*, Feichtenhofer et al, 2019



# 4. Self-attention models – Vision Transformers



Figures from *An image is worth 16x16 words*, Dosovitskiy et al, 2021



# 4. Self-attention models – Video Transformers

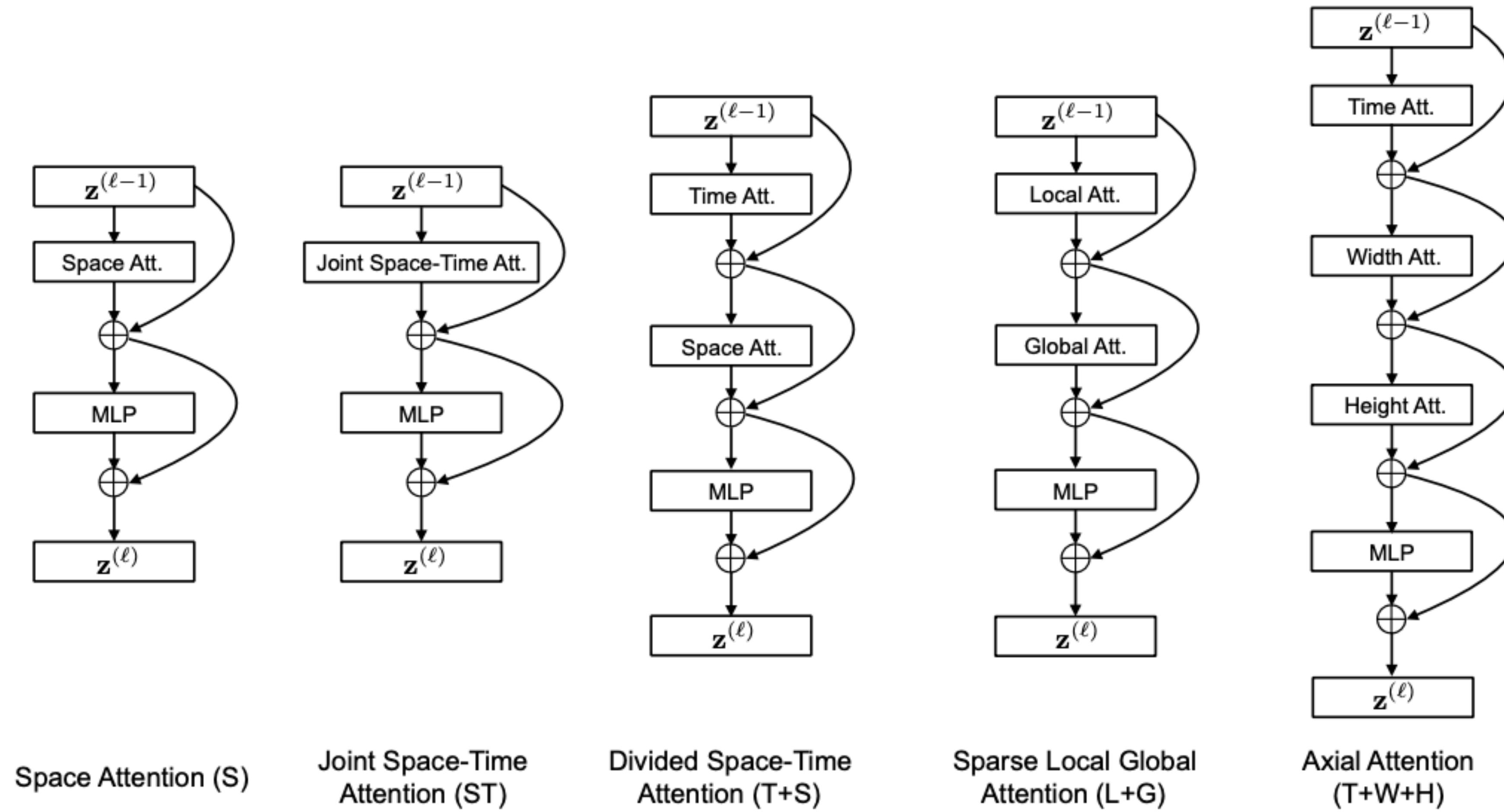
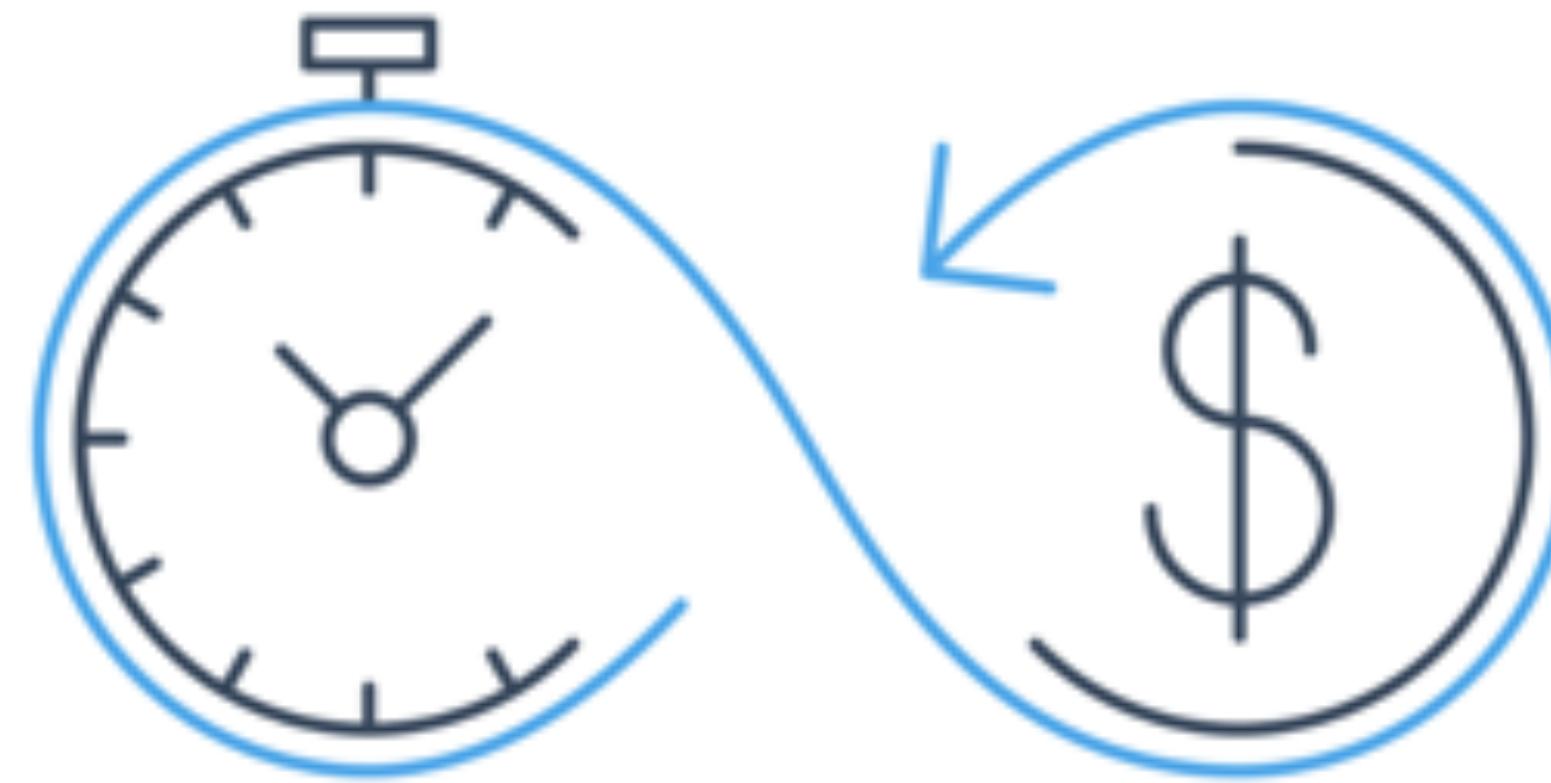


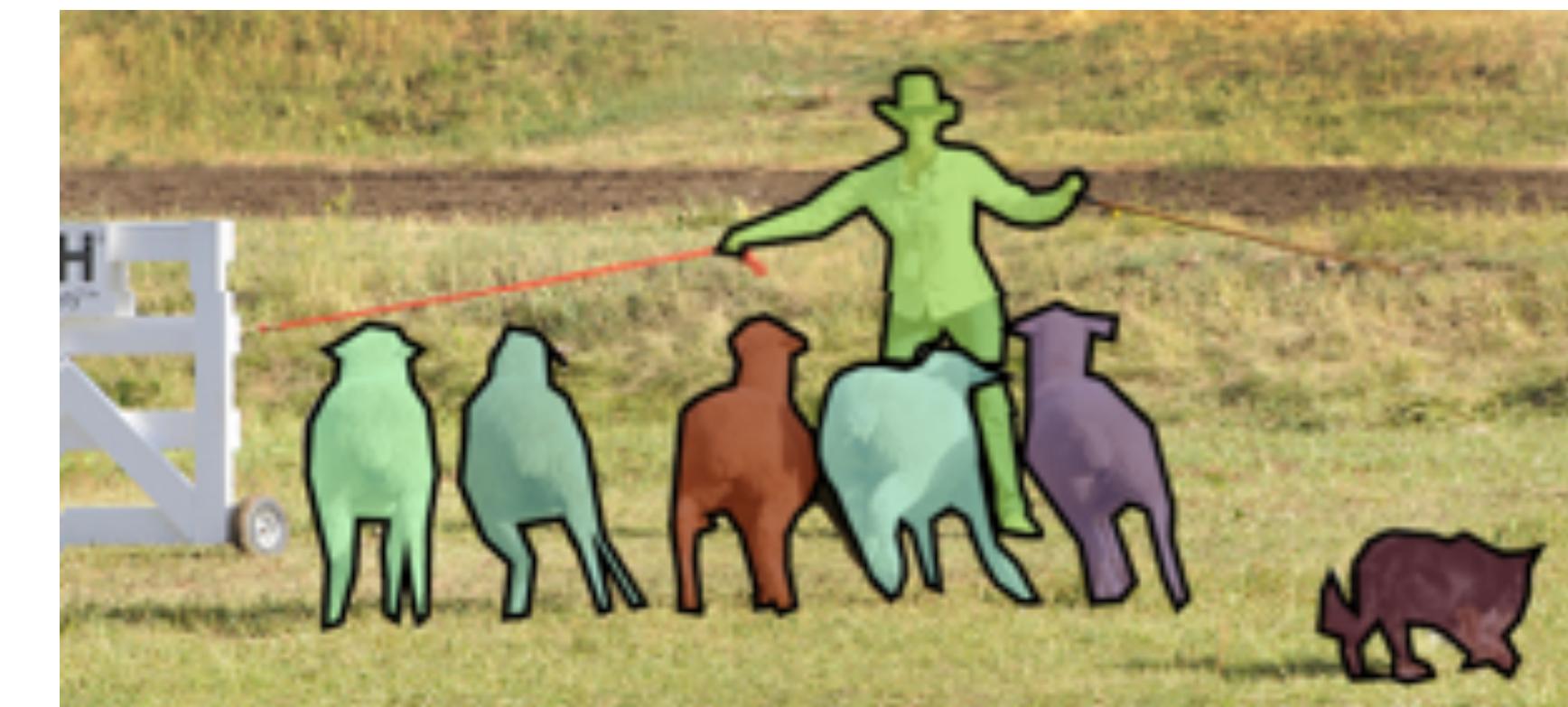
Figure from *Is Space-Time Attention All You Need for Video Understanding?* Bertasius et al, 2021



# Challenges in video processing



**Labels are expensive**



**Agreement: definition? granularity?**



# Challenges in video processing

- ✓ Ok
- ✗ Not great, not terrible
- ✗ Pretty bad

	Frame-by-frame (2D) models	2D-encoder + RNN	3D-conv models	Video transformer
Accuracy	✗	✗	✓	✓
Latency / throughput @ inference	✗	✗	✗	✗
Memory @ train	✓	✗	✗	✗
Memory @ test	✓	✓	✗	✗
Energy efficiency	✗	✗	✗	✗
Causal	✓	✓	✗	✗

Accuracy and latency/throughput dependent on the depth of the model  
 Video transformer - quadratic complexity in all dimensions

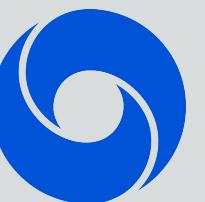


# Latency & throughput

Video from Davis2017 dataset @ 25fps



Same video @ 5fps



# Attempts to improve them

## Boost efficiency of image models

- Model compression / binarisation

Chen et al., Courbarieux et al.

- Distillation

Hinton et al.

- Lighter convolutional modules

MobileNet, Xception

- Budget methods

Karayev et al., Mathe et al.

## Improve efficiency of video models

- Temporal multi-scale models

I3D, SlowFast, TSN

- 2D *temporal* models

TSM, R(2d+1)

- Reuse of features (warping)

Zhu et al.

- Variable update rates

Shelhamer et al.

## Self-supervised methods



O2

# Learning paradigms

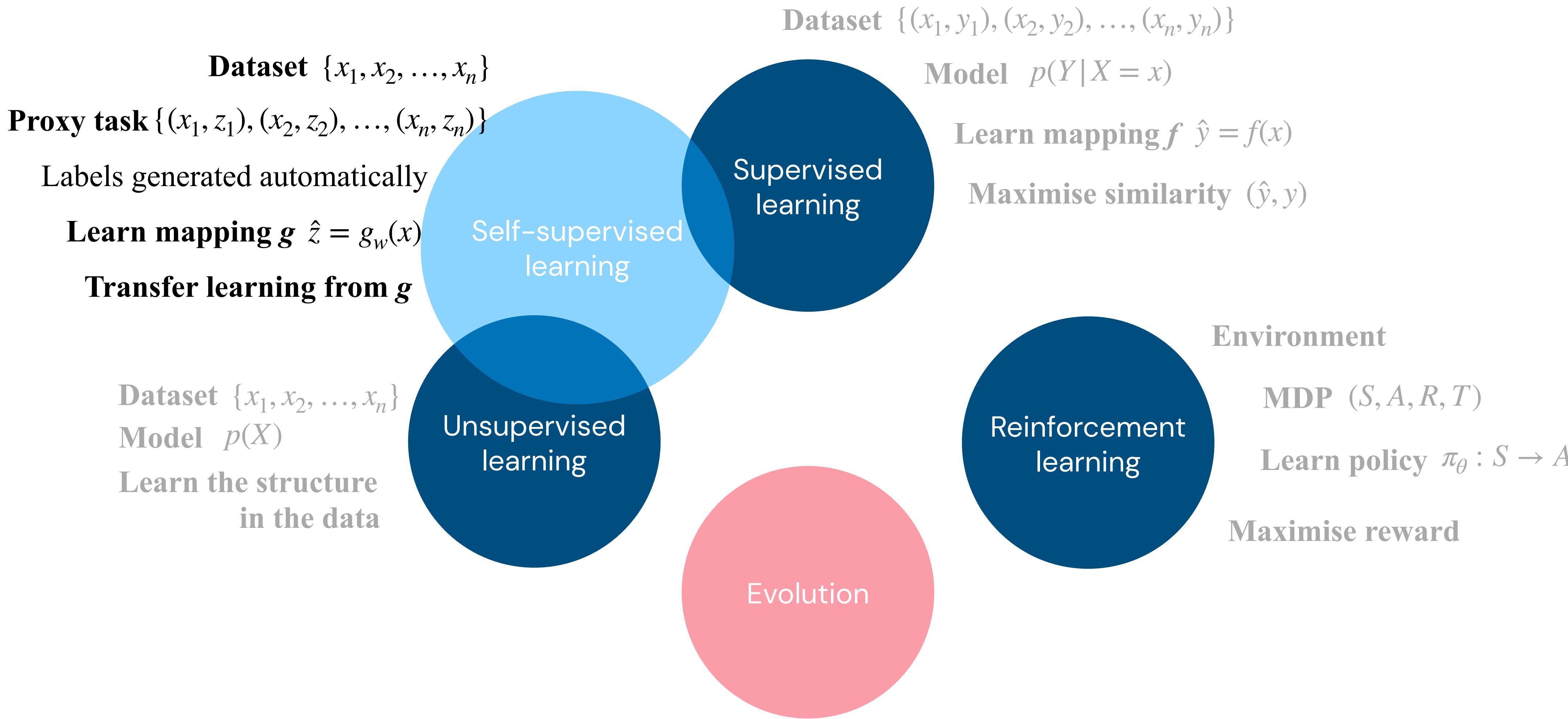


# Learning landscape

Want to learn more?



Zador. A critique of pure learning and what artificial neural networks can learn from animal brains (2019)



# Downstream tasks in videos

Tracking / segmentation



Video from Davis dataset

Action recognition



Video from Kinetics dataset



# Multimodal sensing of the environment

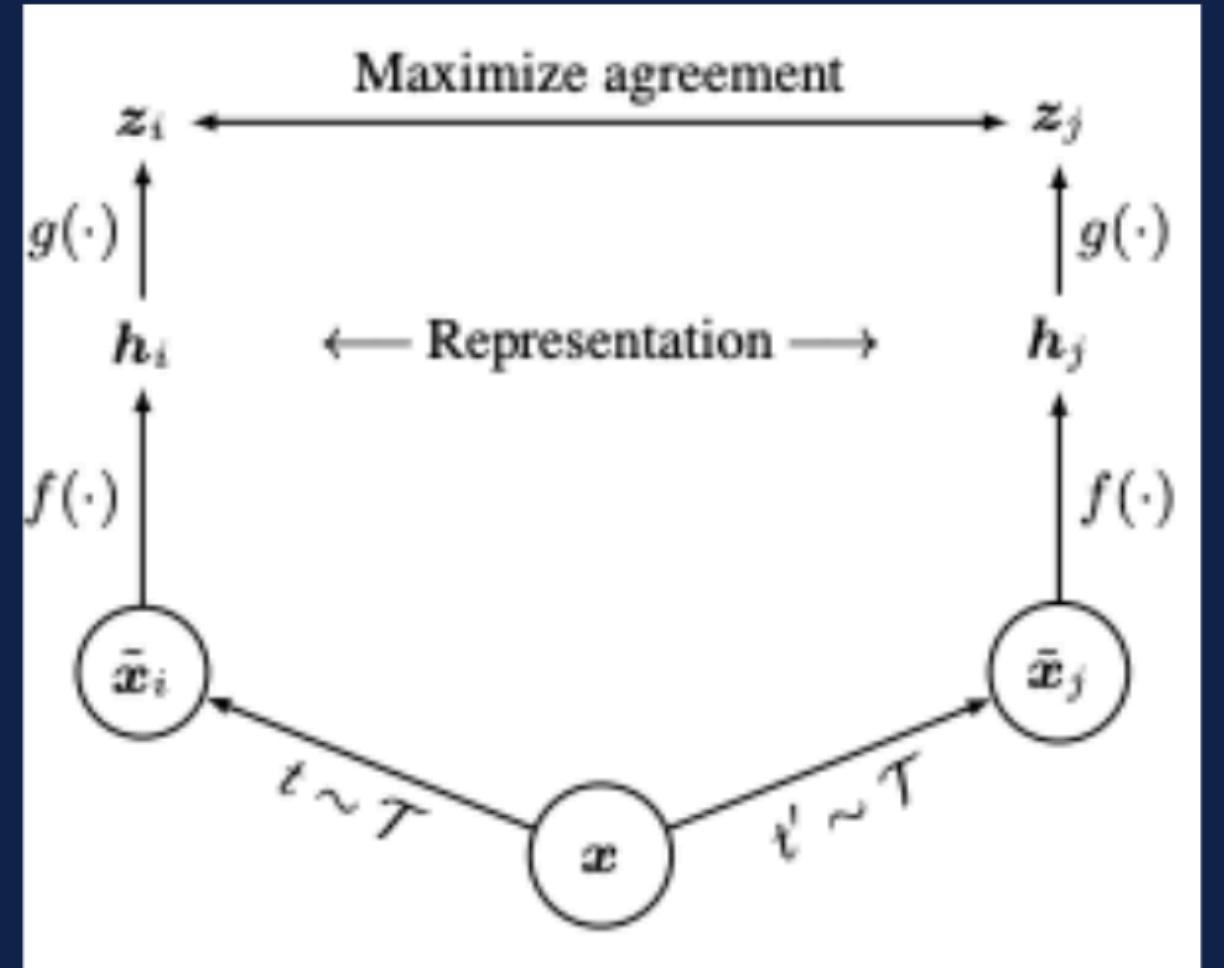


[...] towards the root and try to get as close to the root as possible, nice long strokes [...]

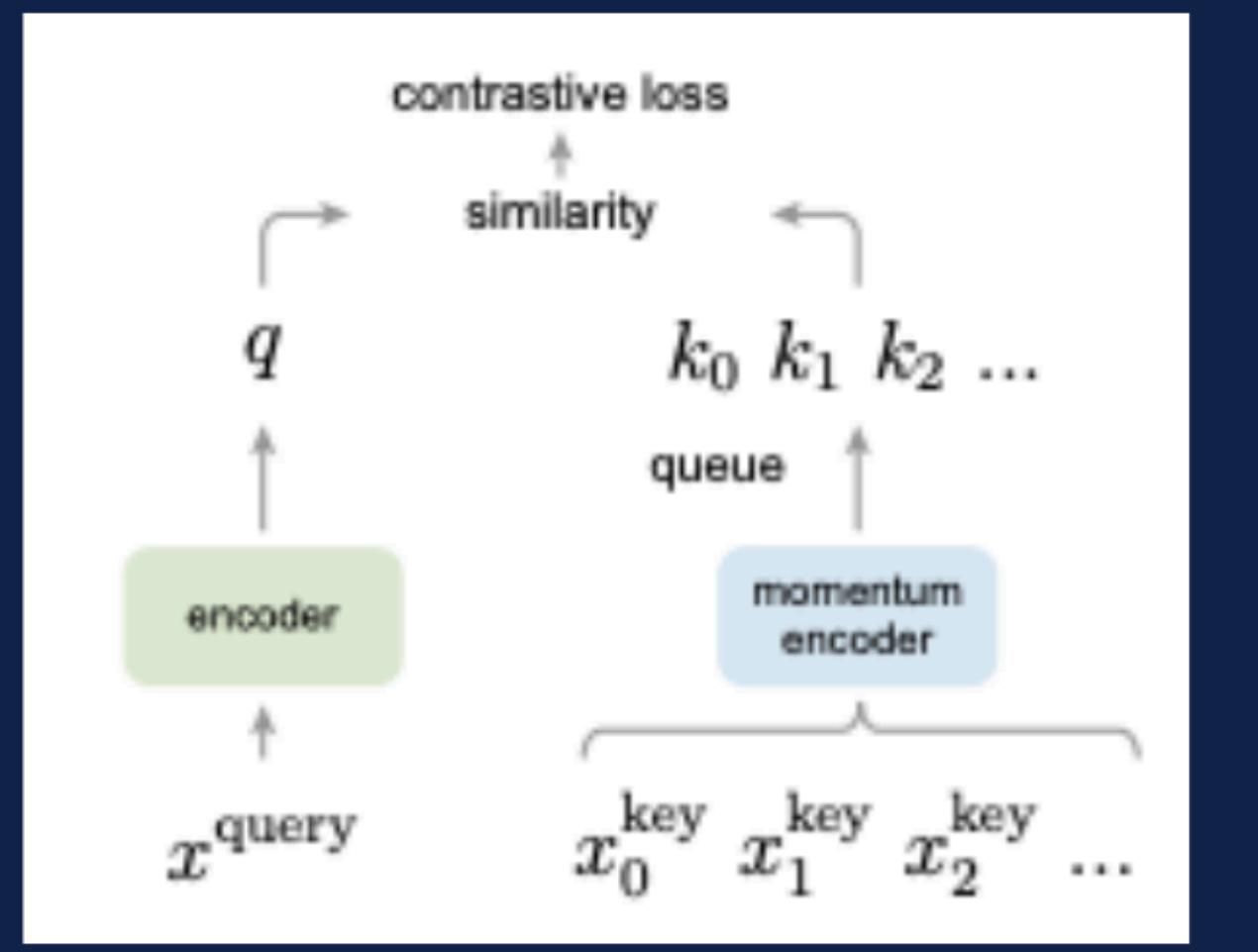


# Self-supervised learning for pre-training

## Vision



SimCLR: Chen et al, 2020

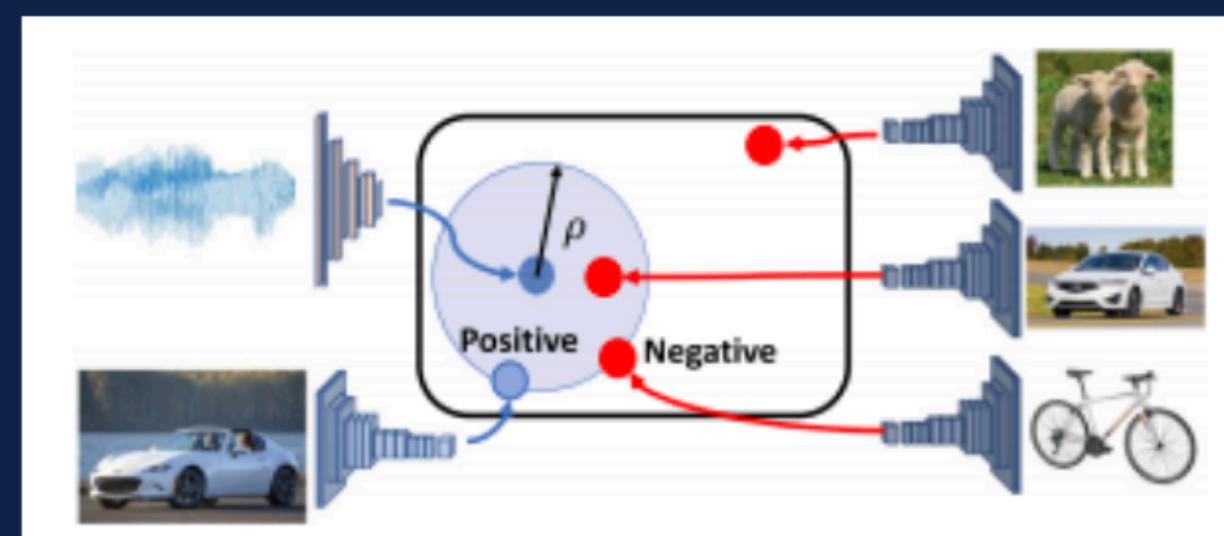


MOCO: He et al, 2020

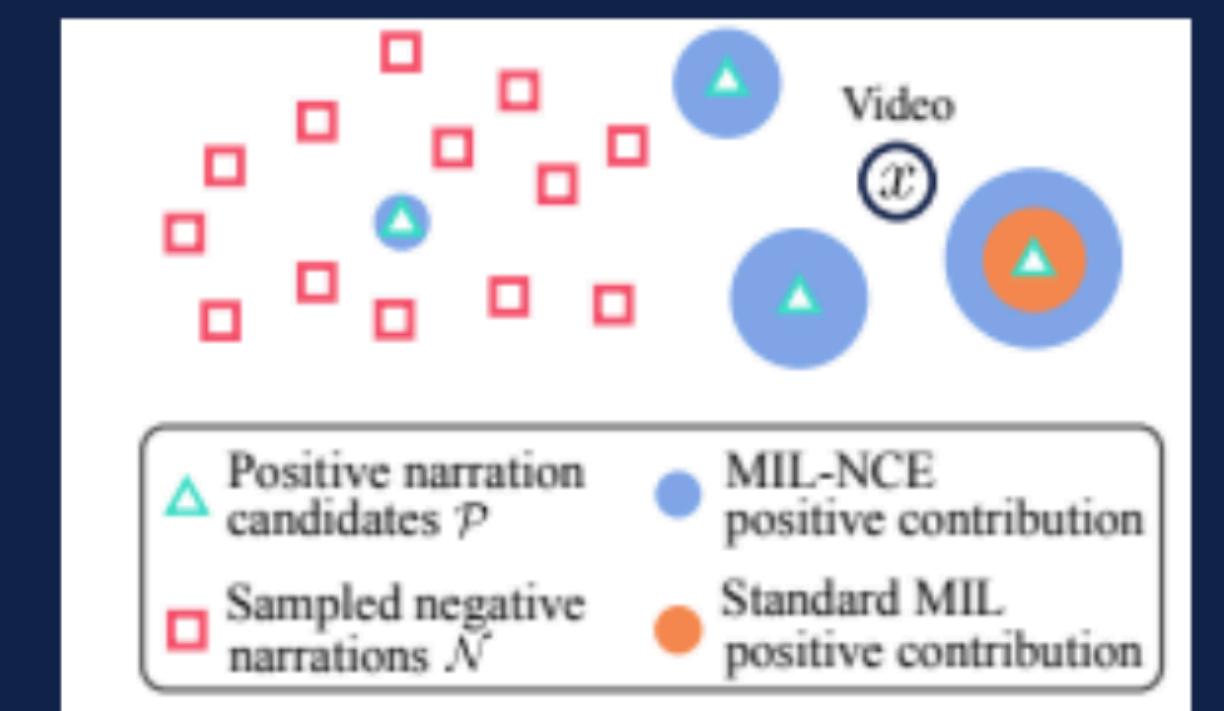
## Vision+Language



VideoBERT: Sun et al, 2019

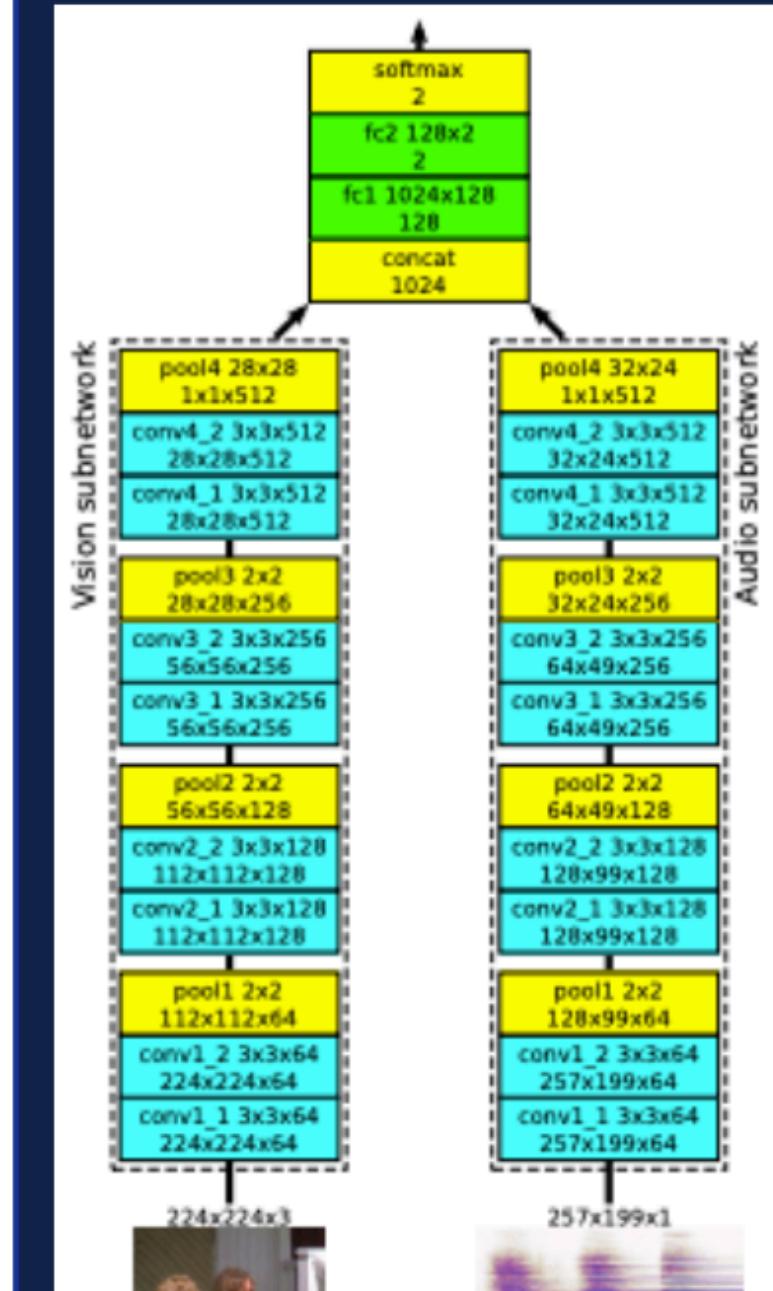


DaveNet: Harwath et al, 2018

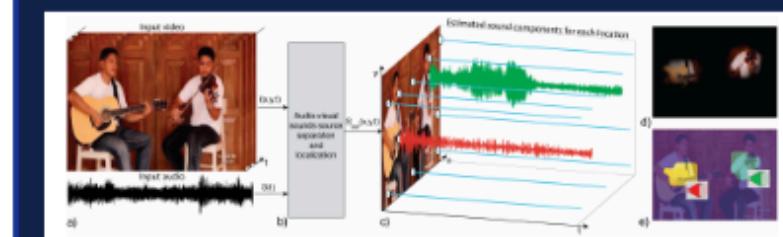


MIL-NCE: Miech, Alayrac et al, 2020

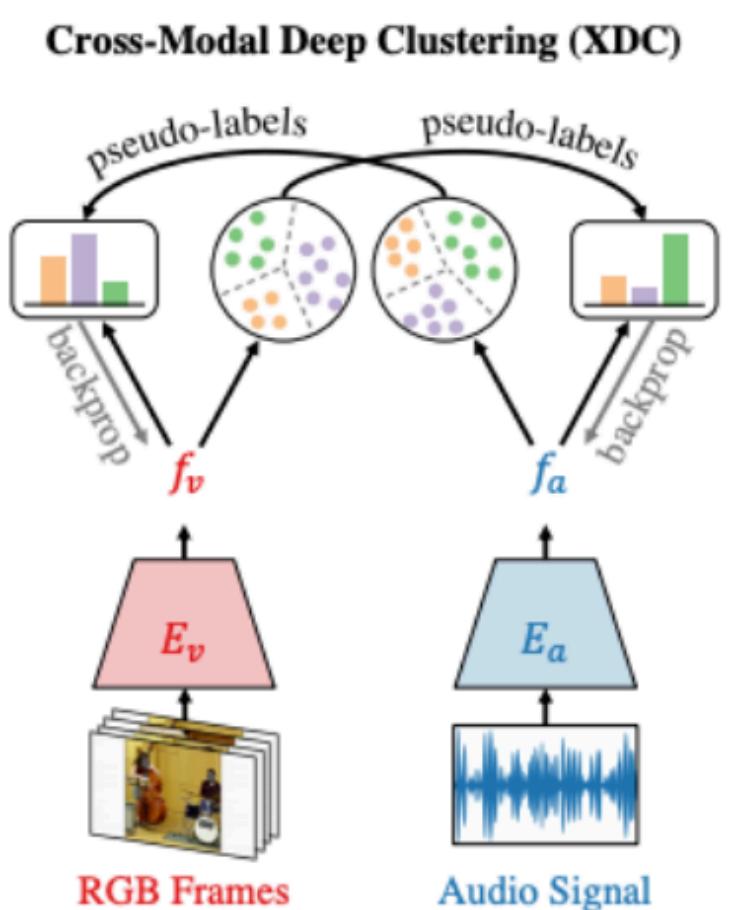
## Vision+Audio



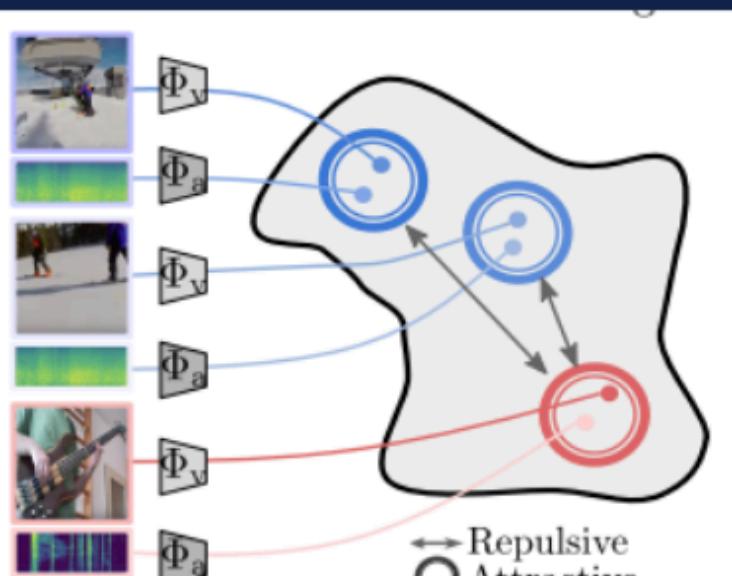
L3: Arandjelovic and Zisserman, 2017



Sound of Pixels: Zhao et al, 2018



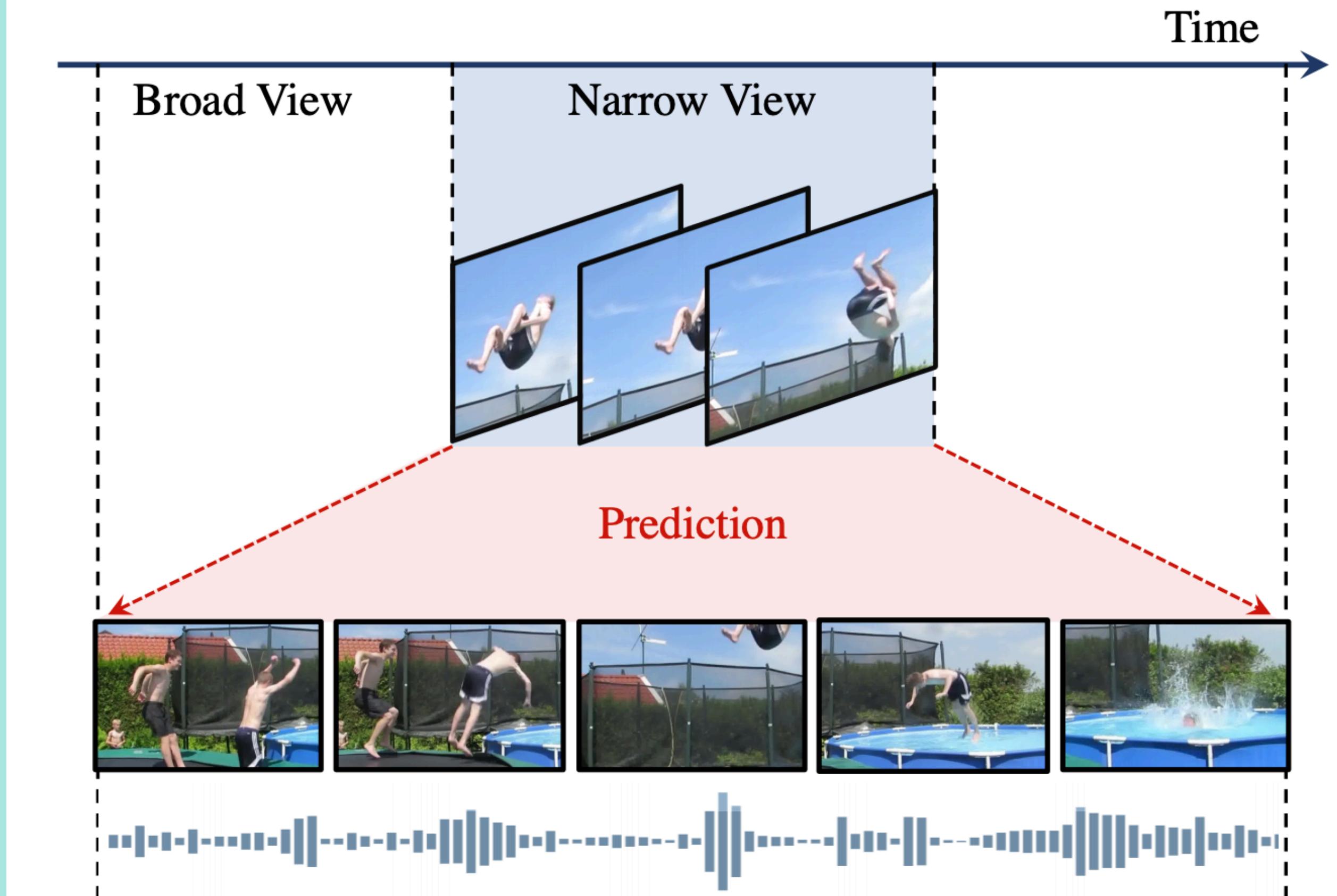
XDC: Alwassel at al, 2020



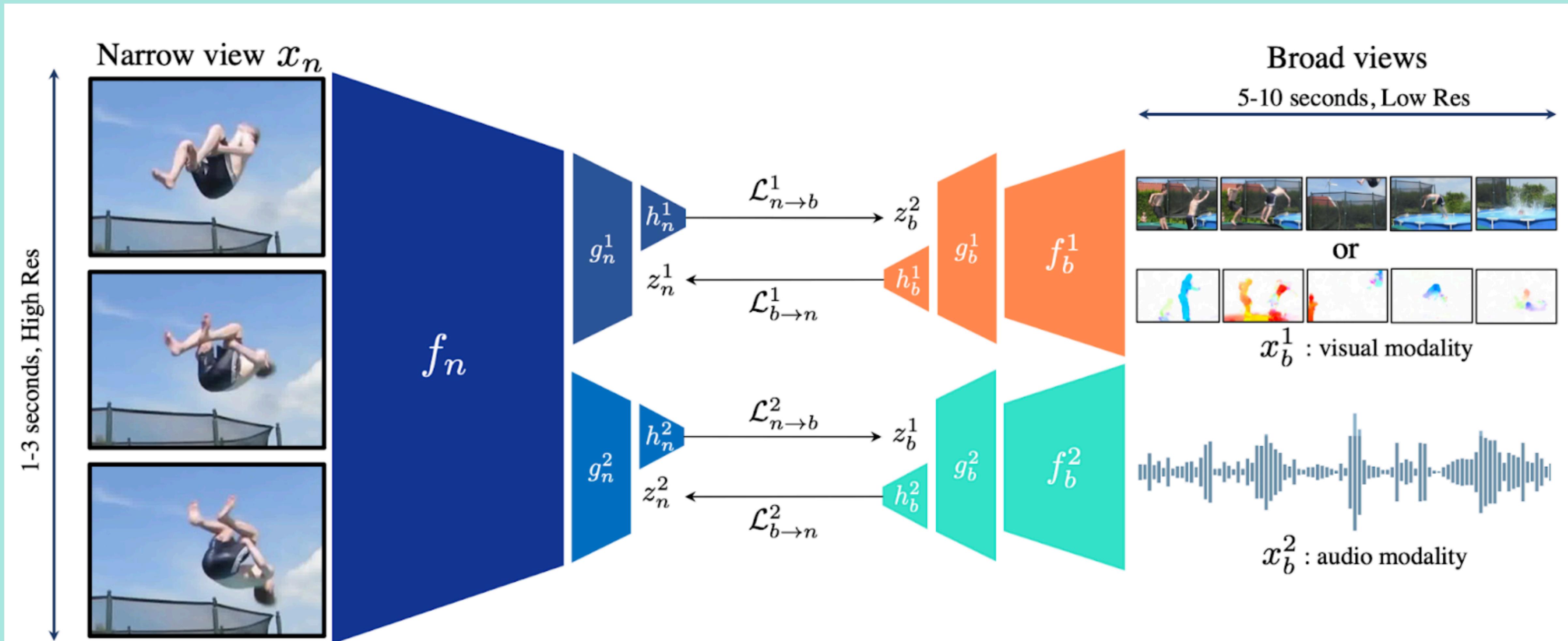
GDT: Patrick et al, 2020

# Broaden your views for self-supervised video learning

Recasens et al, ICCV2021



# BraVe



# BraVe training loss

**Global loss**

$$\mathcal{L}(x) = \underbrace{\mathcal{L}_{n \rightarrow b}(x)}_{\text{Narrow} \rightarrow \text{Broad}} + \underbrace{\mathcal{L}_{b \rightarrow n}(x)}_{\text{Broad} \rightarrow \text{Narrow}}$$

**Narrow to Broad Loss**

$$\mathcal{L}_{n \rightarrow b}(x) = \left\| \frac{h_n(z_n)}{\|h_n(z_n)\|_2} - \text{sg} \left[ \frac{z_b}{\|z_b\|_2} \right] \right\|_2^2$$

**Broad to Narrow Loss**

$$\mathcal{L}_{b \rightarrow n}(x) = \left\| \frac{h_b(z_b)}{\|h_b(z_b)\|_2} - \text{sg} \left[ \frac{z_n}{\|z_n\|_2} \right] \right\|_2^2$$



# BraVe results

- When using the same backbone and dataset, our model significantly beats the state-of-the-art.

Method	Backbone (#params)	Dataset	Years	$\mathcal{M}$	UCF101		HMDB51		K600	ESC-50	AS
					Linear	FT	Linear	FT	Linear	Linear	MLP
AVTS [45]	MC3 (11.7M)	AS	1	VA	89.0		61.6			80.6	
ELo [67]	R(2+1)D-50 (46.9M)	YT8M	13	VFA	93.8	64.5	67.4				
AVID [58]	R(2+1)D-50 (46.9M)	AS	1	VA	91.5		64.7			89.2	
GDT [64]	R(2+1)D-18 (33.3M)	AS	1	VA	92.5		66.1			88.5	
MMV [3]	R(2+1)D-18 (33.3M)	AS	1	VA	83.9	91.5	60.0	70.1	55.5	85.6	29.7
XDC [4]	R(2+1)D-18 (33.3M)	AS	1	VA	93.0		63.7			84.8	
XDC [4]	R(2+1)D-18 (33.3M)	IG65M	21	VA	95.5		68.9			85.4	
<b>BraVe:V↔A (ours)</b>	R(2+1)D-18 (33.3M)	AS	1	VA	89.9	94.1	64.8	71.1	63.6	90.4	34.7
<b>BraVe:V↔A (ours)</b>	TSM-50 (23.5M)	AS	1	VA	93.0	94.8	69.4	72.6	70.1	90.5	34.4
<b>BraVe:V↔FA (ours)</b>	TSM-50 (23.5M)	AS	1	VFA	93.1	95.4	70.0	74.6	69.3	90.1	34.5
<b>BraVe:V↔FA (ours)</b>	R(2+1)D-50 (46.9M)	AS	1	VFA	92.5	95.1	68.3	73.6	69.4	91.6	34.5
<b>BraVe:V↔FA (ours)</b>	TSM-NF-F0 (71.5M)	AS	1	VFA	94.1	95.8	71.4	73.1	72.6	90.2	34.5
<b>BraVe:V↔FA (ours)</b>	TSM-50x2 (93.9M)	AS	1	VFA	93.1	95.7	70.5	77.8	71.4	91.1	34.8
Supervised [11, 44, 67, 87]					96.8	71.5	75.9	82.4		94.7	43.9



03

Depth  
parallelism



# Attempts to improve efficiency of video models

## Boost efficiency of image models

- Model compression / binarisation

Chen et al., Courbarieux et al.

- Distillation

Hinton et al.

- Lighter convolutional modules

MobileNet, Xception

- Budget methods

Karayev et al., Mathe et al.

## Improve efficiency of video models

- Temporal multi-scale models

I3D, SlowFast, TSN

- 2D *temporal* models

TSM, R(2d+1)

- Reuse of features (warping)

Zhu et al.

- Variable update rates

Shelhamer et al.

All models are executed sequentially over depth



# Our direction: Is sequential mode optimal?

NUMBERS | NEUROSCIENCE

## Why Is the Human Brain So Efficient?

*How massive parallelism lifts the brain's performance above that of AI.*

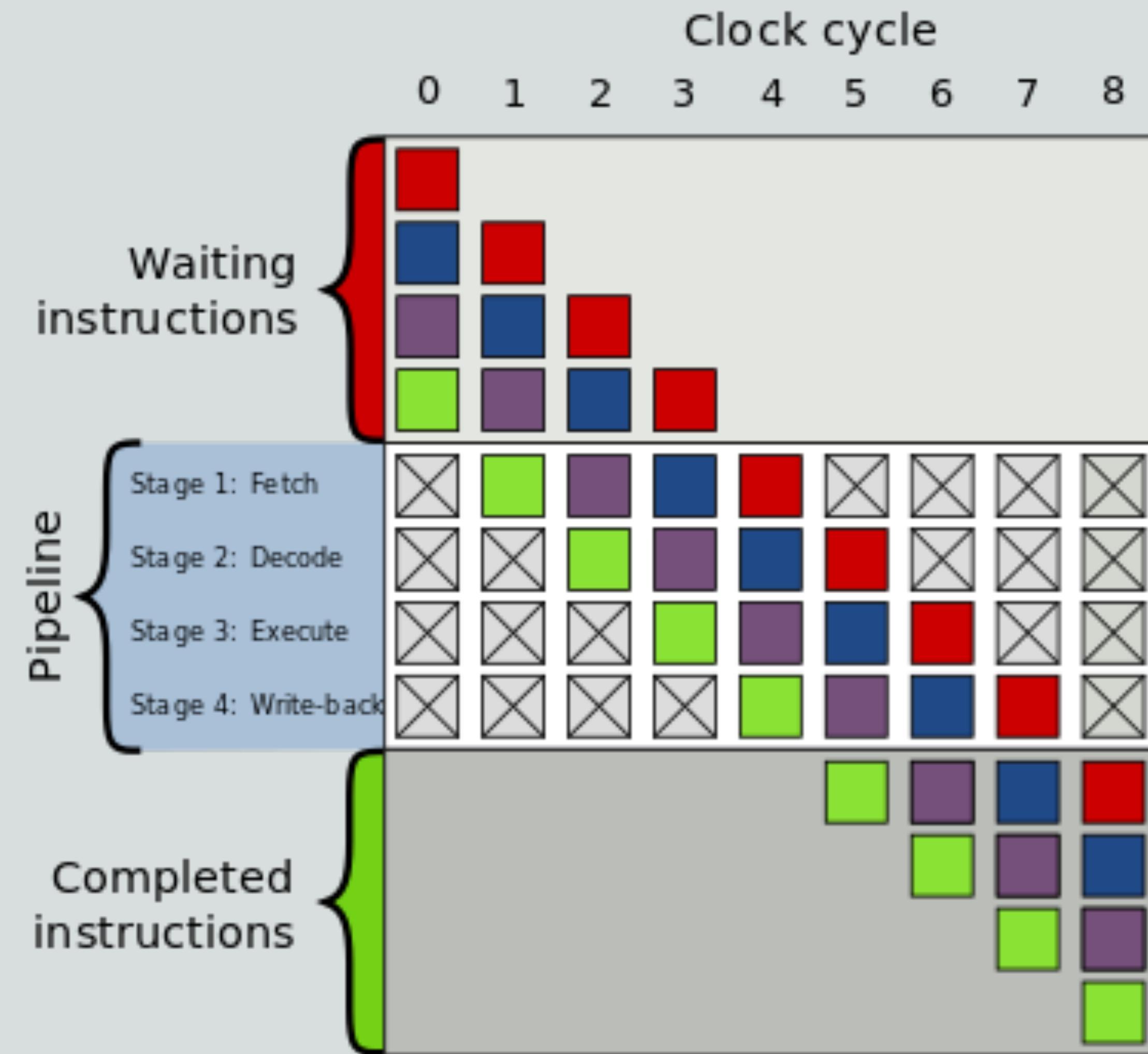
BY LIQUN LUO

APRIL 12, 2018

# Our direction: Is sequential mode optimal?

General-purpose processors use **pipelined** instructions enabling **parallel** processing.

**Maximise throughput**  
Efficient use of hw resources



# Outline

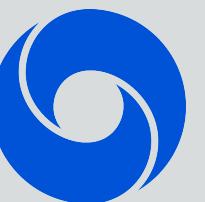
## Depth-parallel video models

A. Depth-parallel inference

B. Depth-parallel training

## Goals

- Increase throughput
- Reduce latency
- Reduce clock cycles (energy)
- Biological plausibility

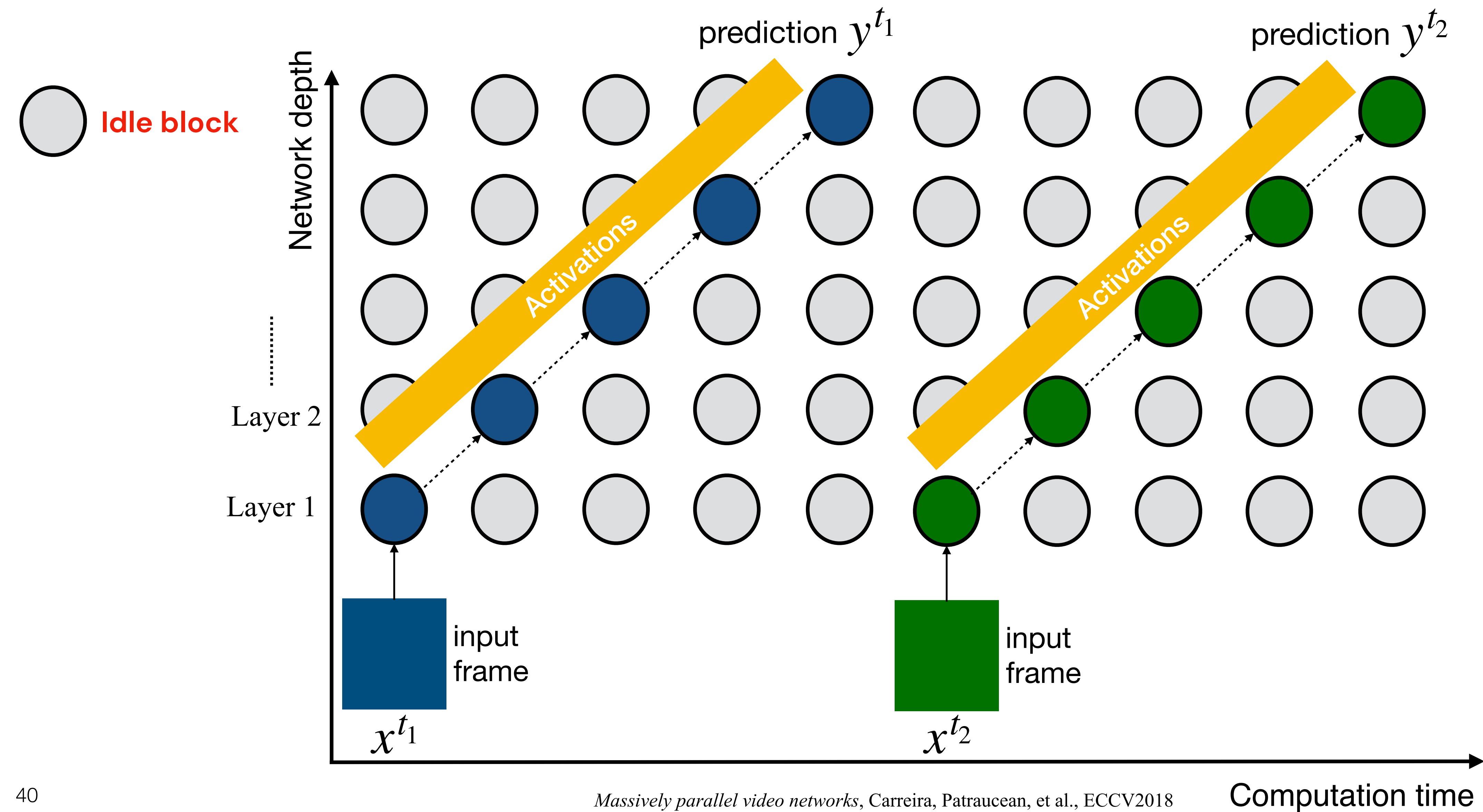




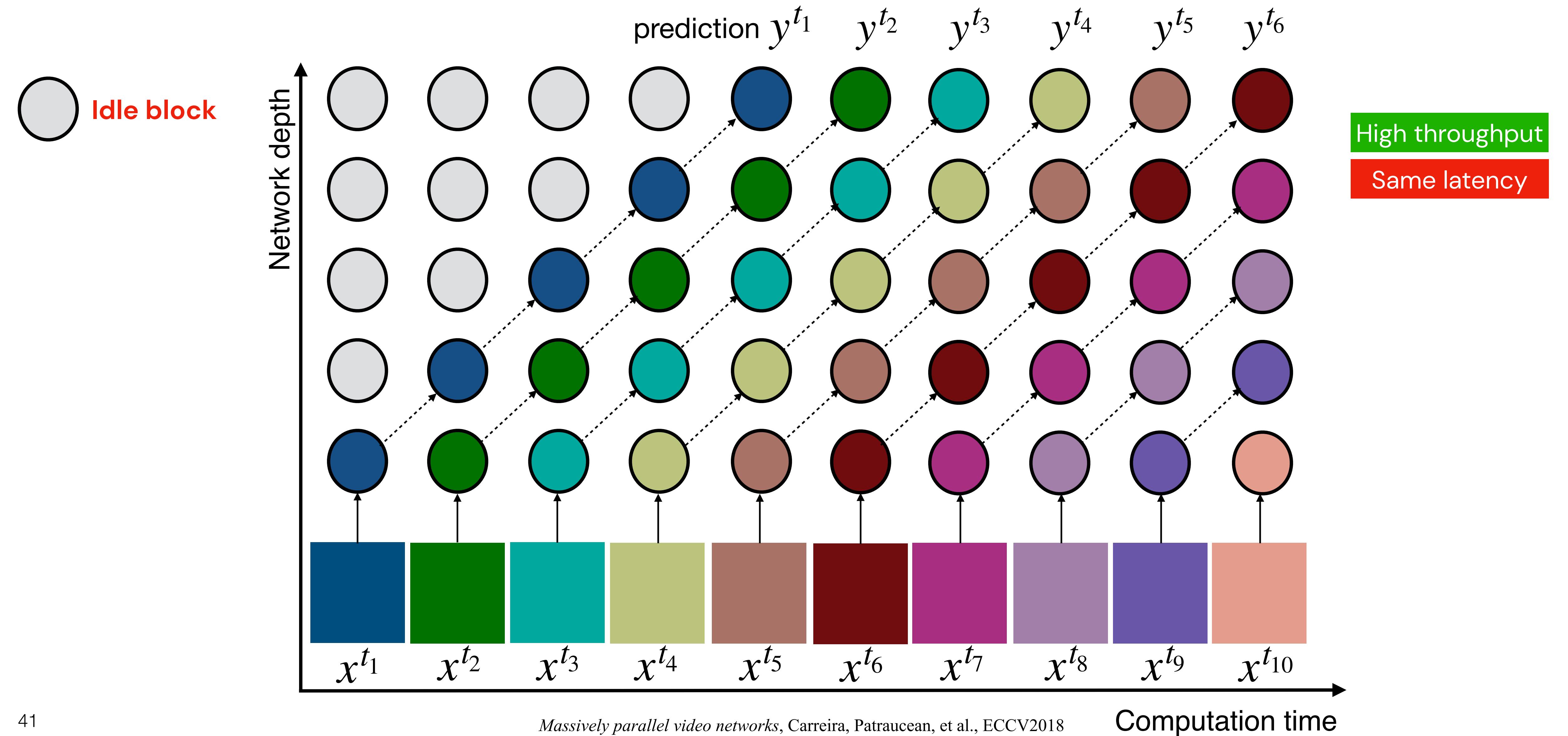
## A. Depth-parallel inference



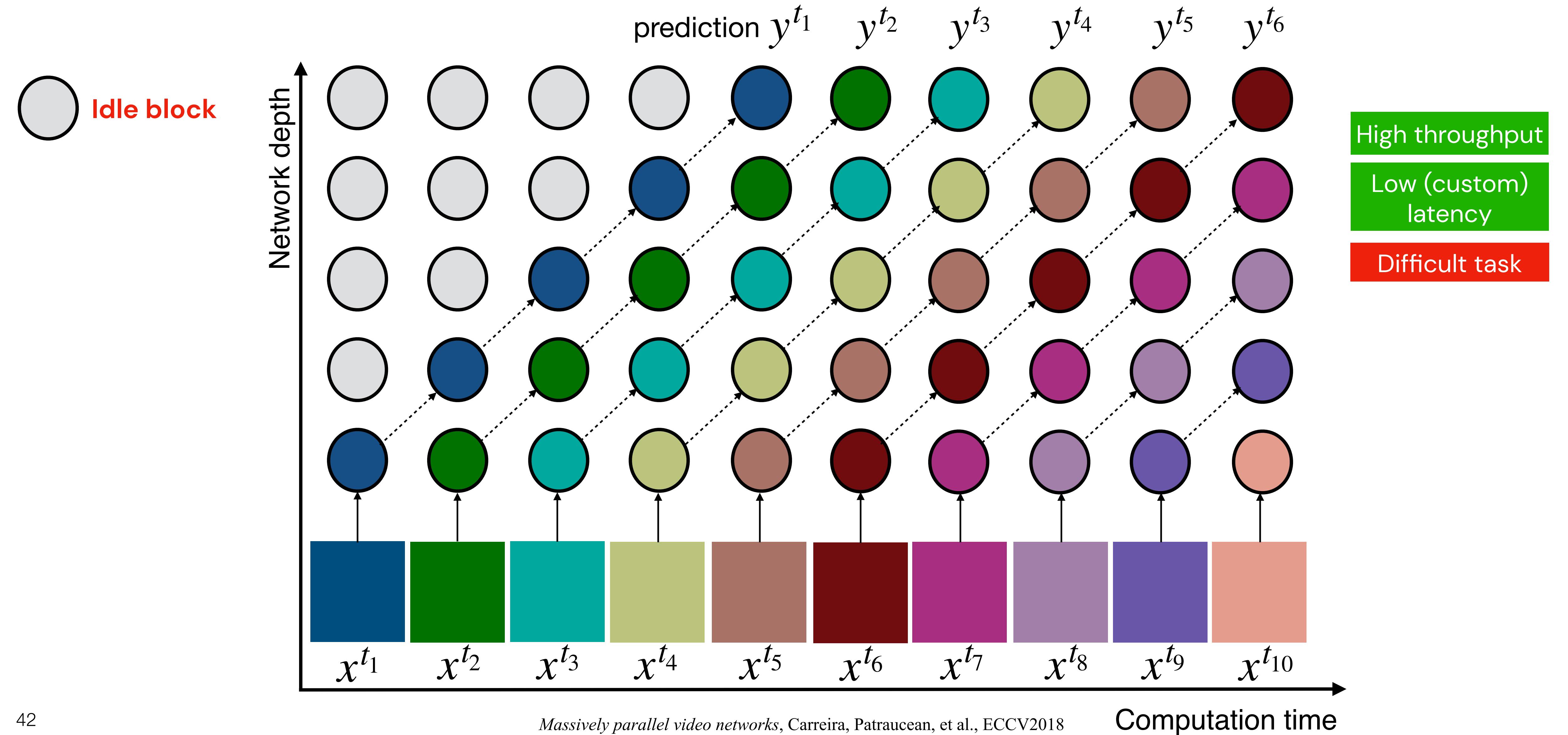
# Sequential (frame-by-frame) model



# Pipelined model

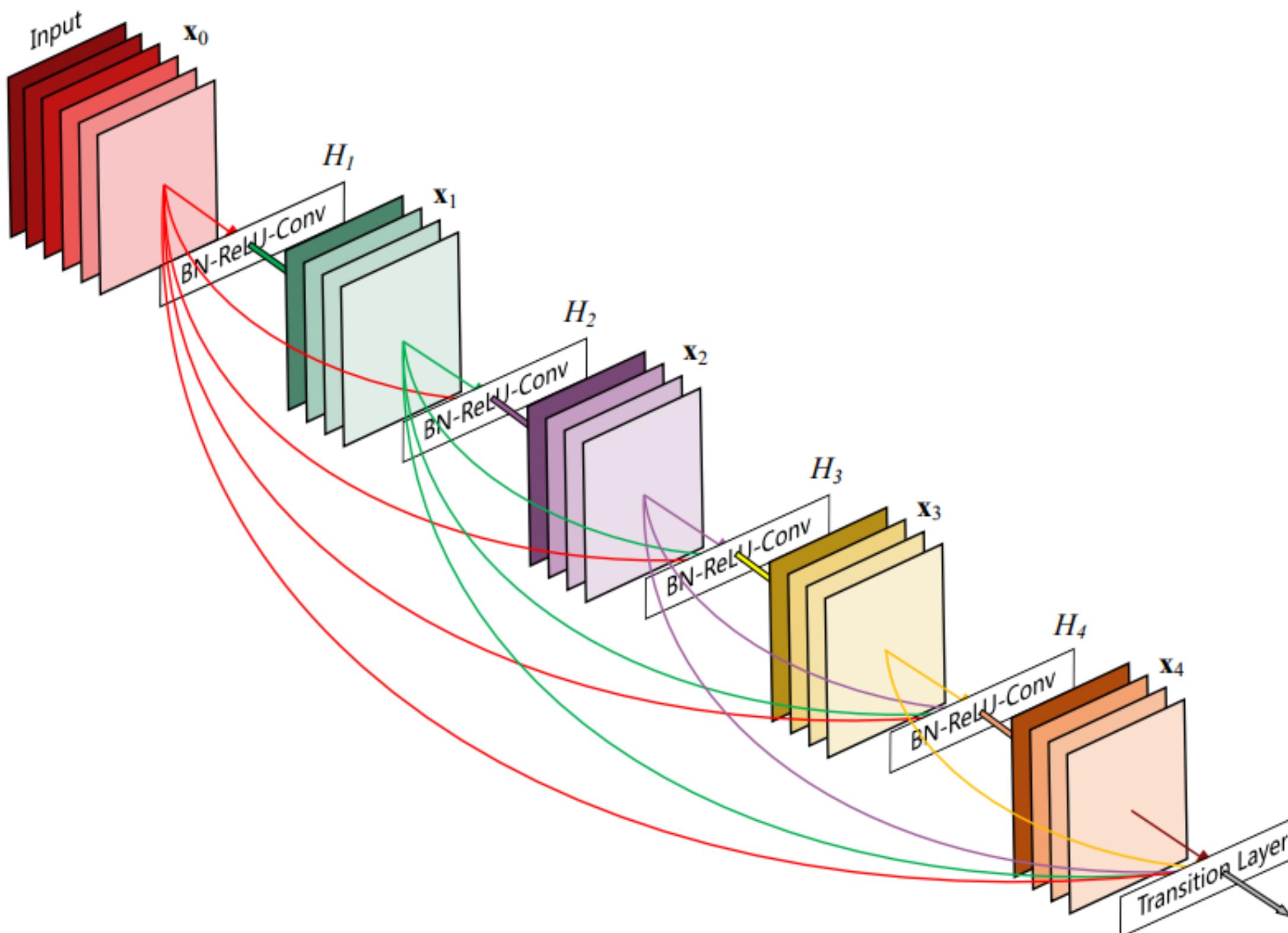


# Pipelined model



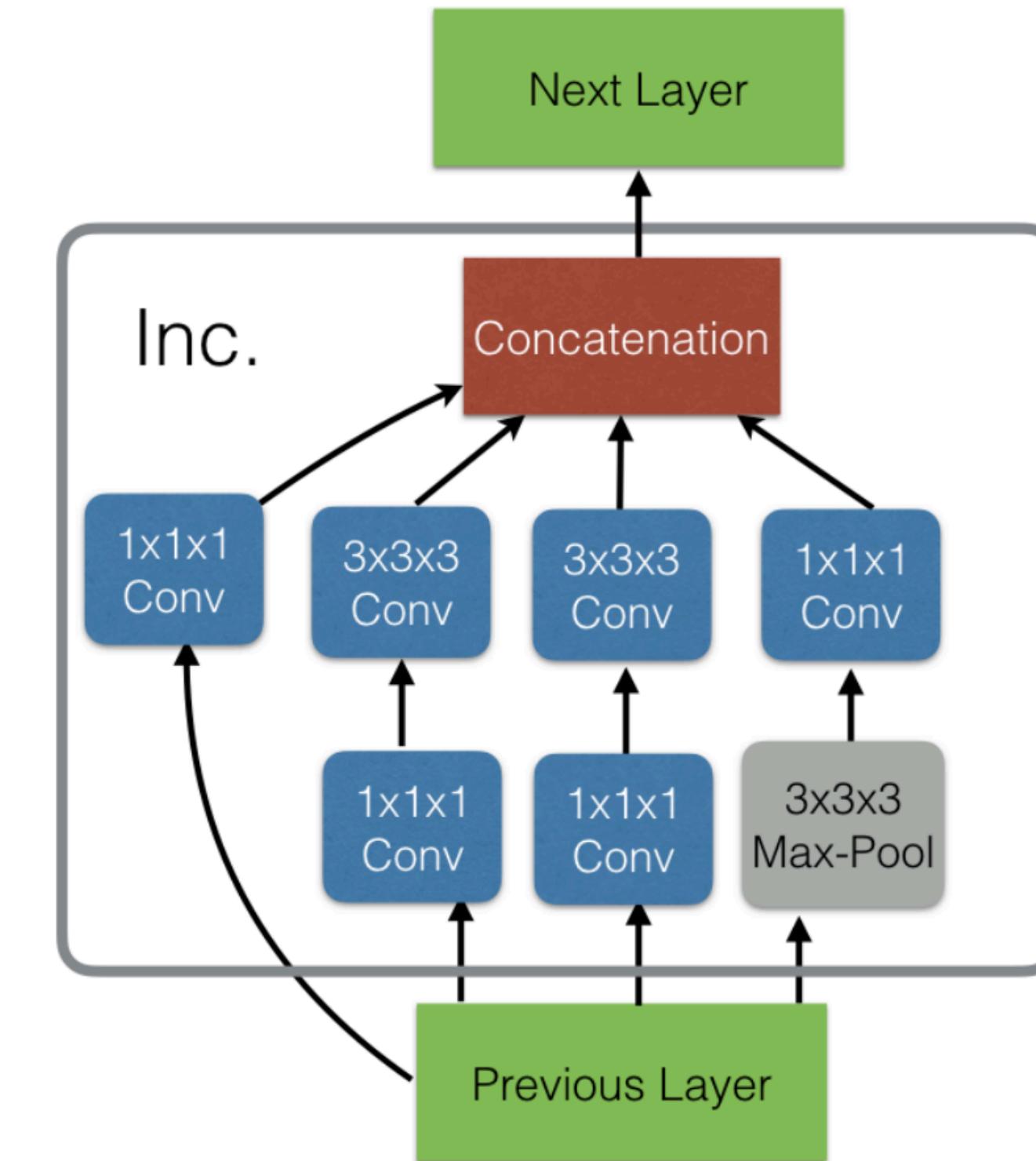
# Any existing model can be pipelined

DenseNet



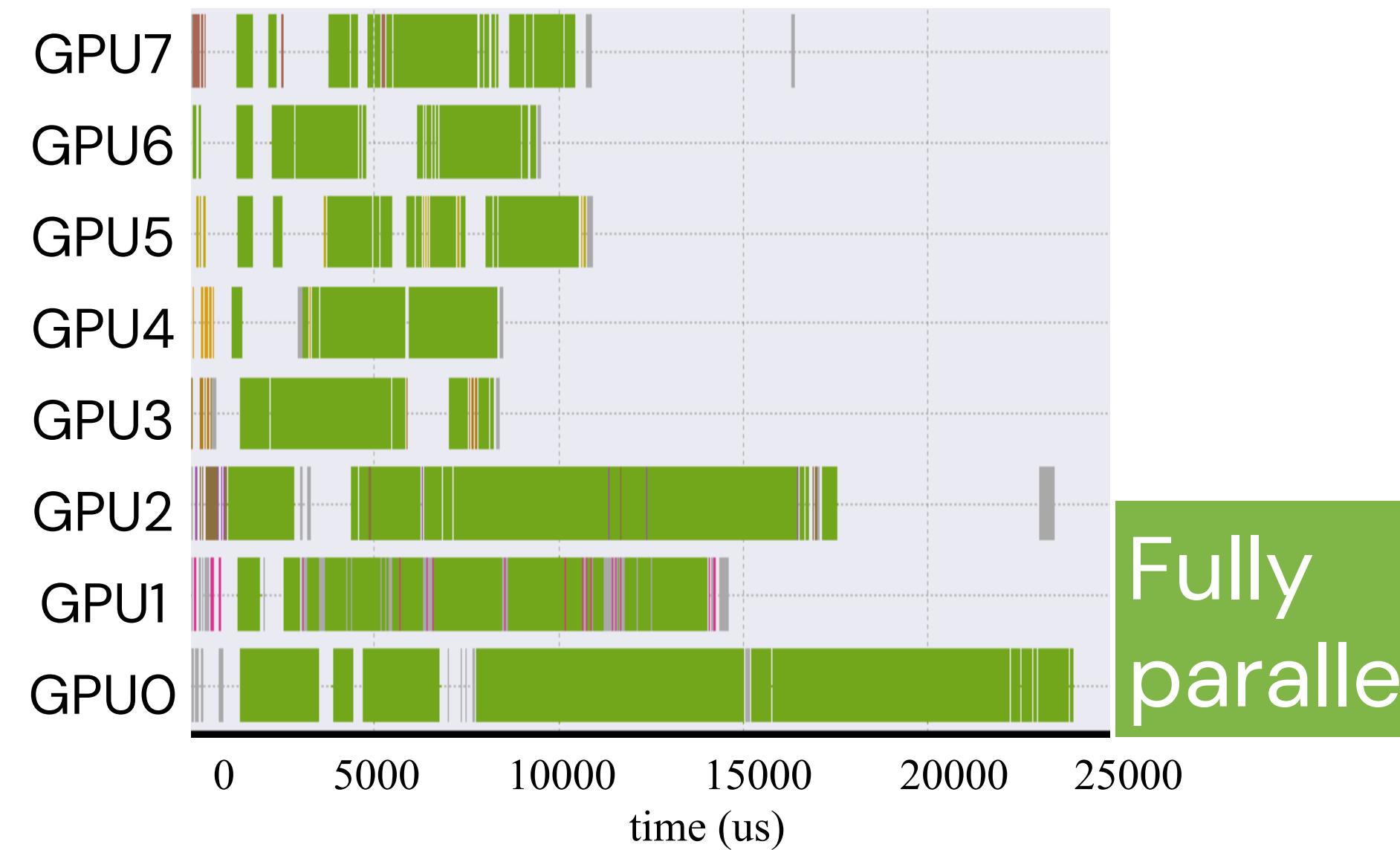
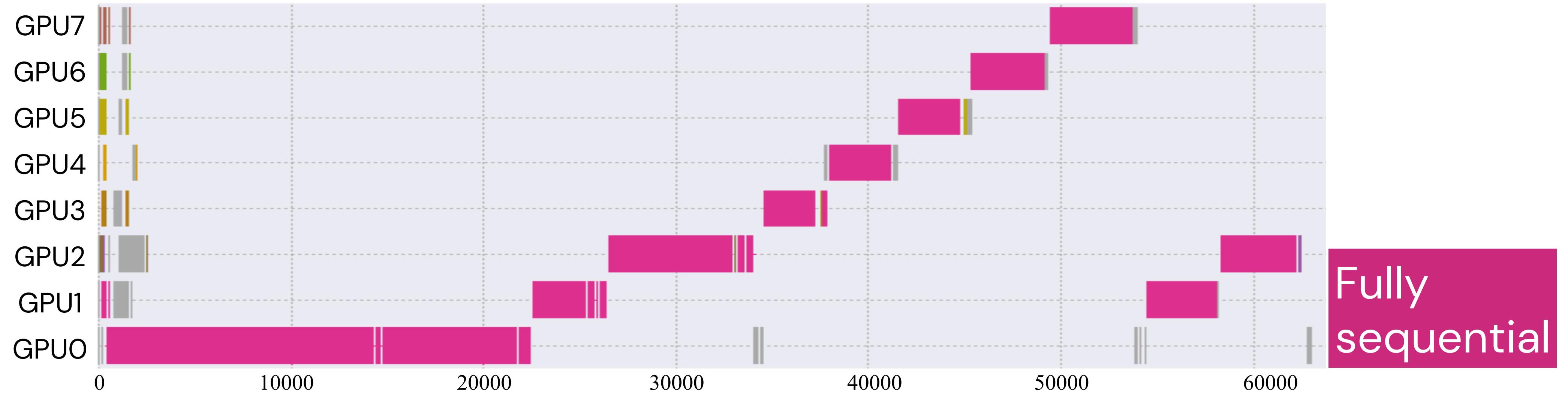
*Densely connected convolutional neural networks*, Huang et al., CVPR2017

Inception

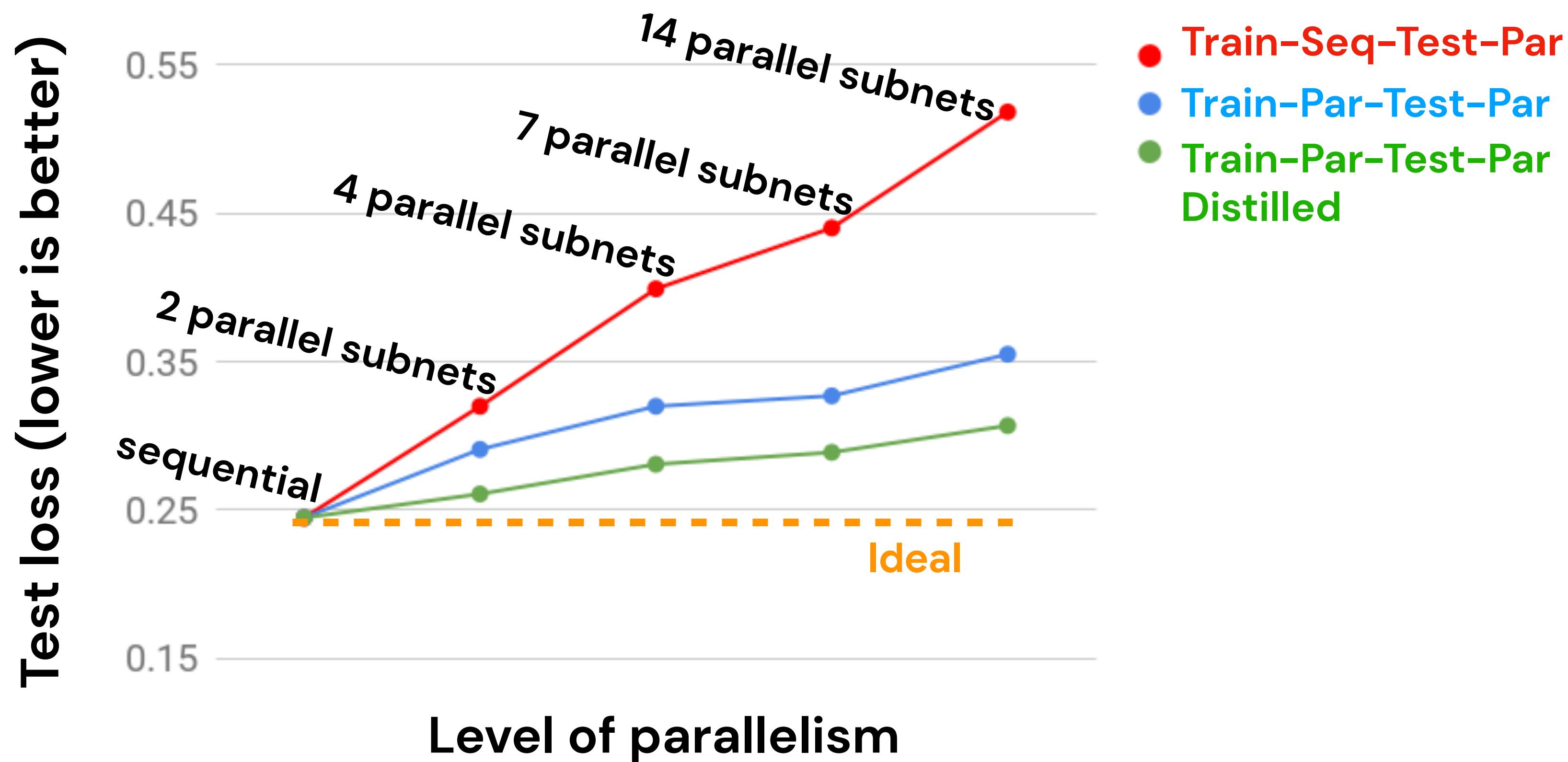


*Quo Vadis, action recognition?* Carreira and Zisserman, CVPR2017

# GPU utilisation

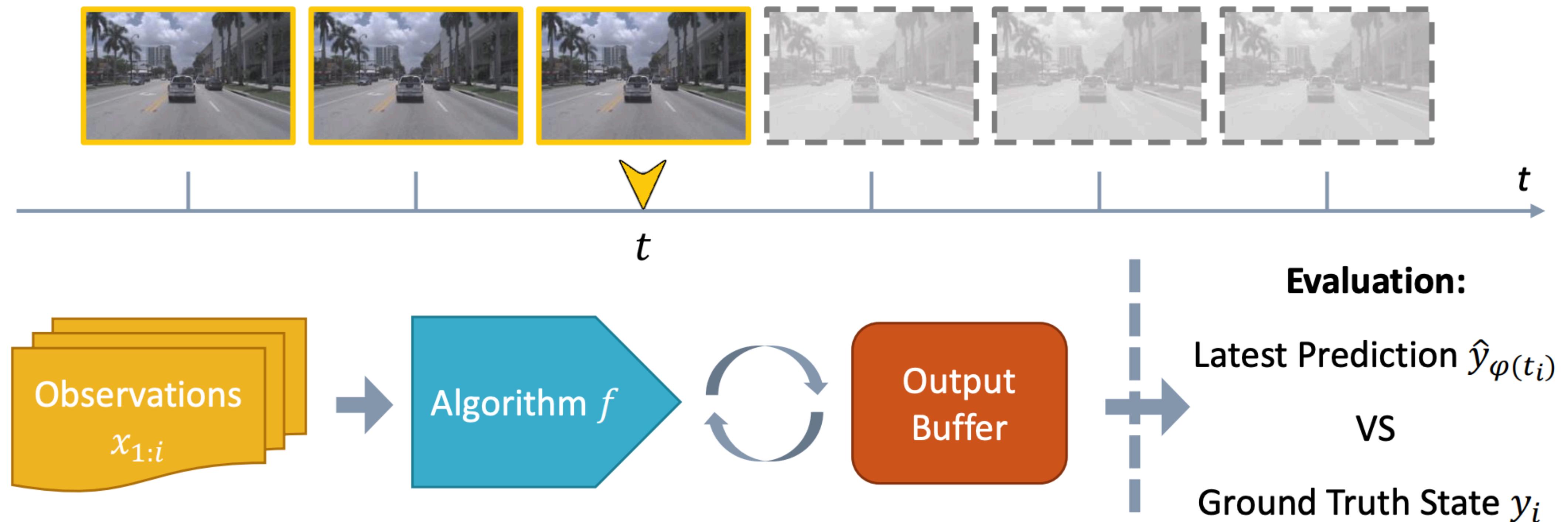


# Human keypoint localisation



# Fair(er) asynchronous evaluation

*Towards streaming image understanding, Li et al, ECCV2020*



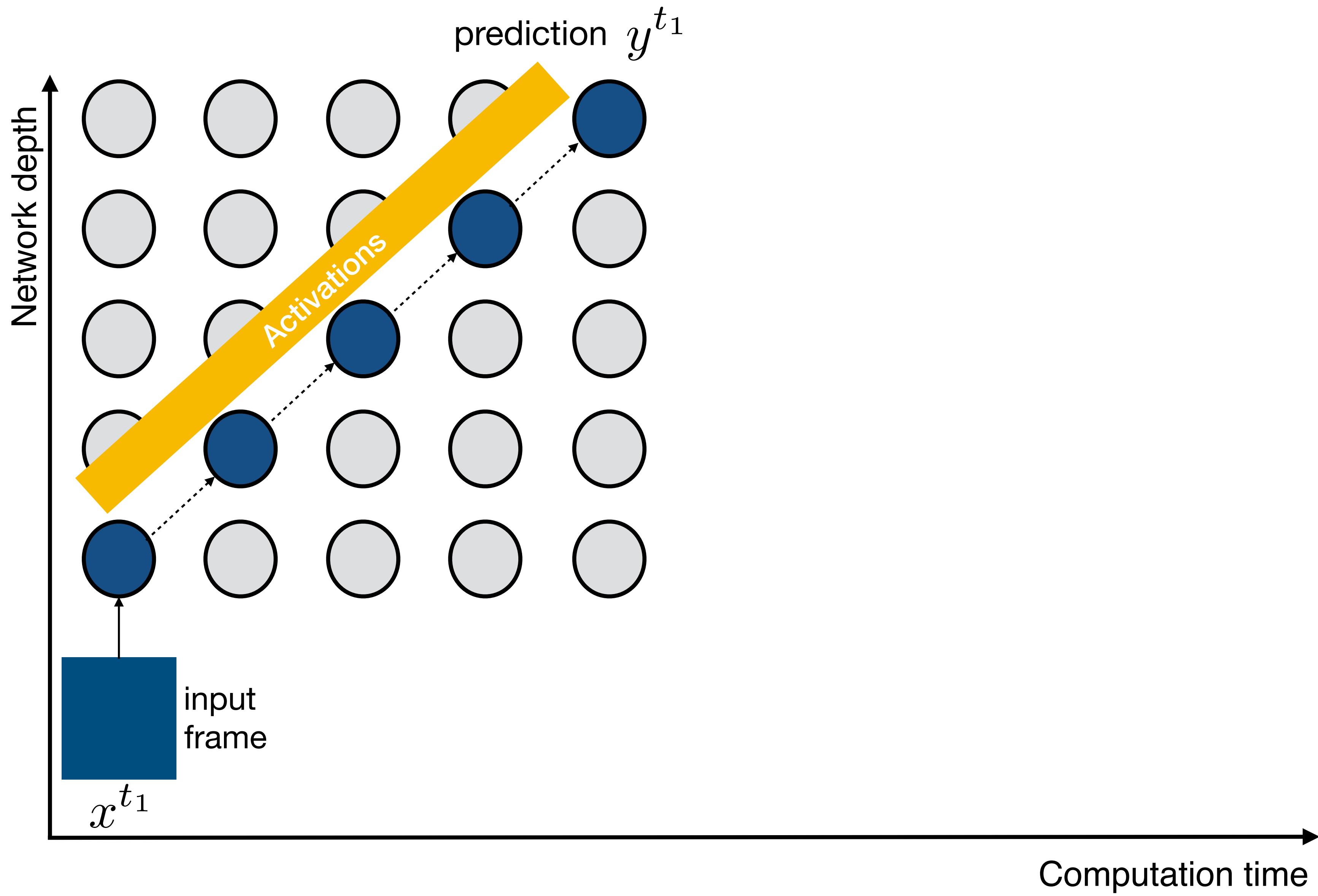
Object detector: offline precision 38.0, streaming average precision 20.3



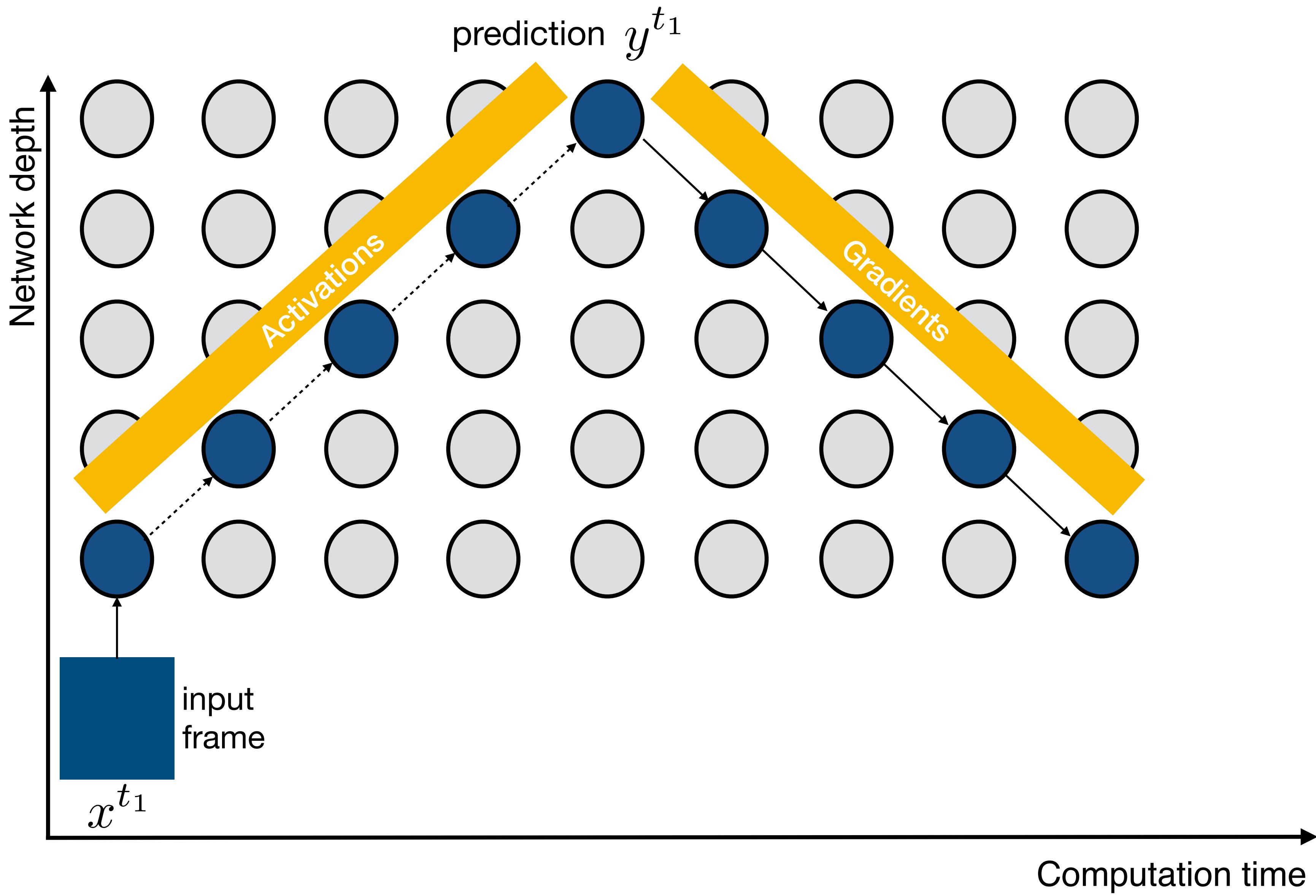
## B. Depth-parallel training



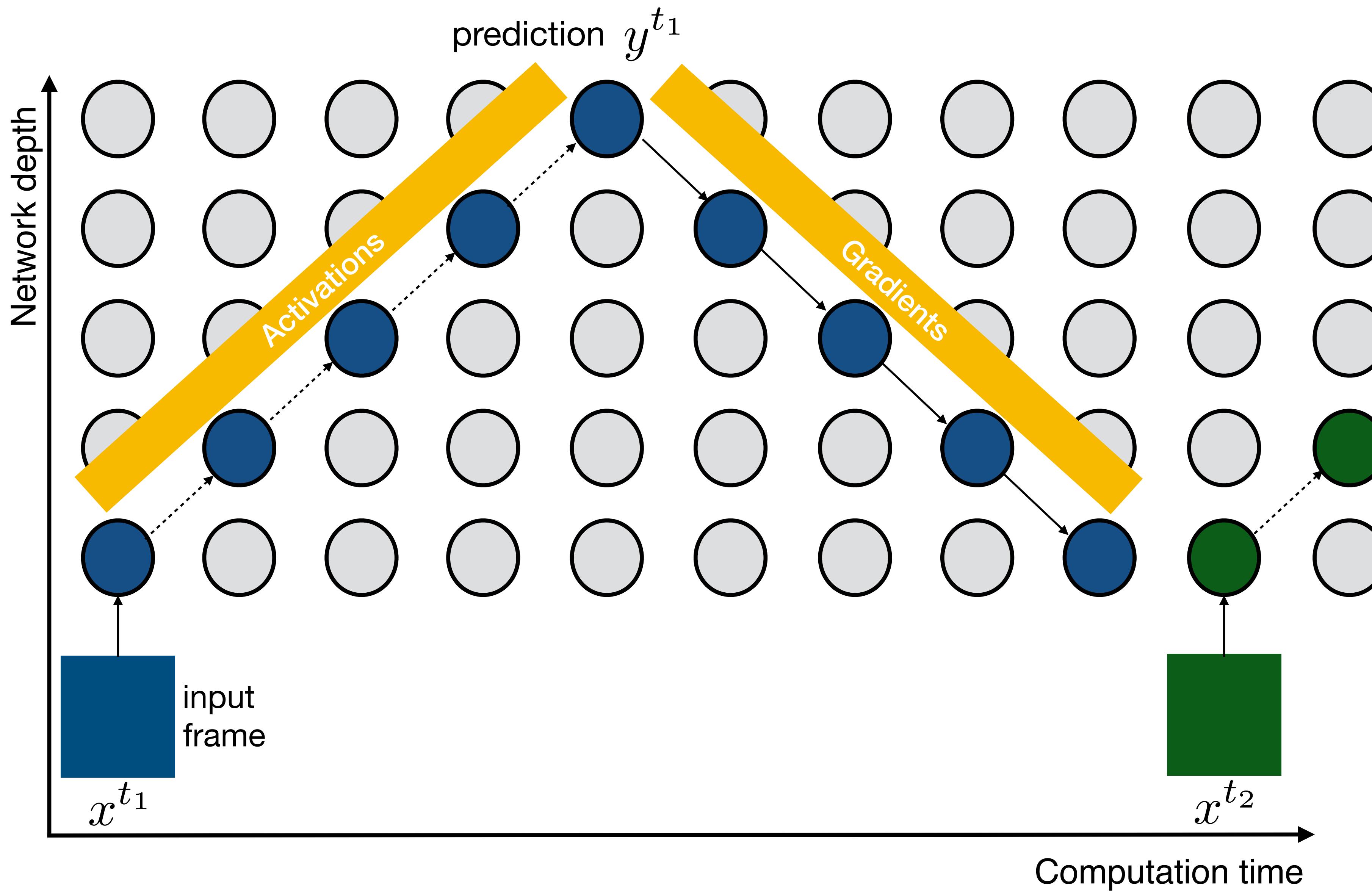
# Sequential (frame-by-frame) model



# Sequential (frame-by-frame) model

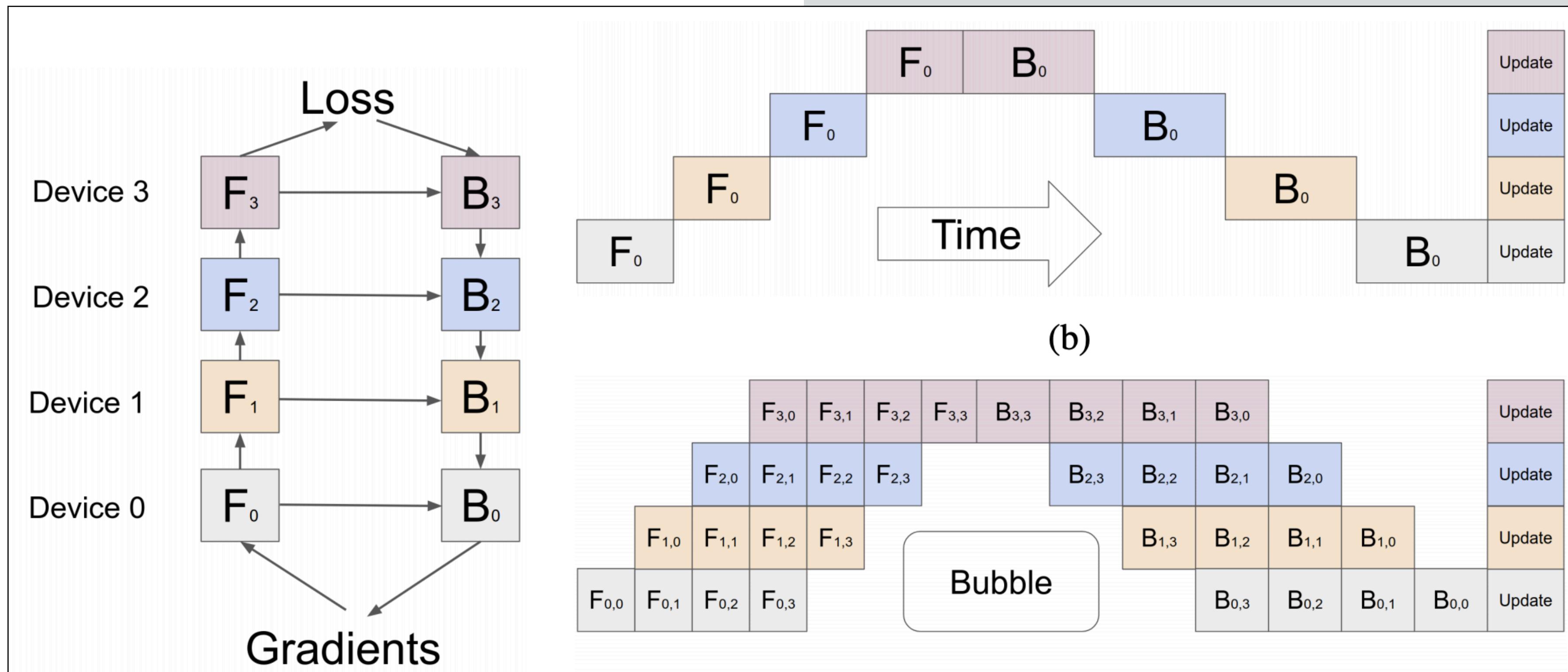


# Sequential (frame-by-frame) model



# Pipelined backpropagation: no blocking

*GPipe: Easy Scaling with Micro-Batch Pipeline Parallelism*, Yanping et al, 2019



**Pipelined training of image classifiers by buffering activations**

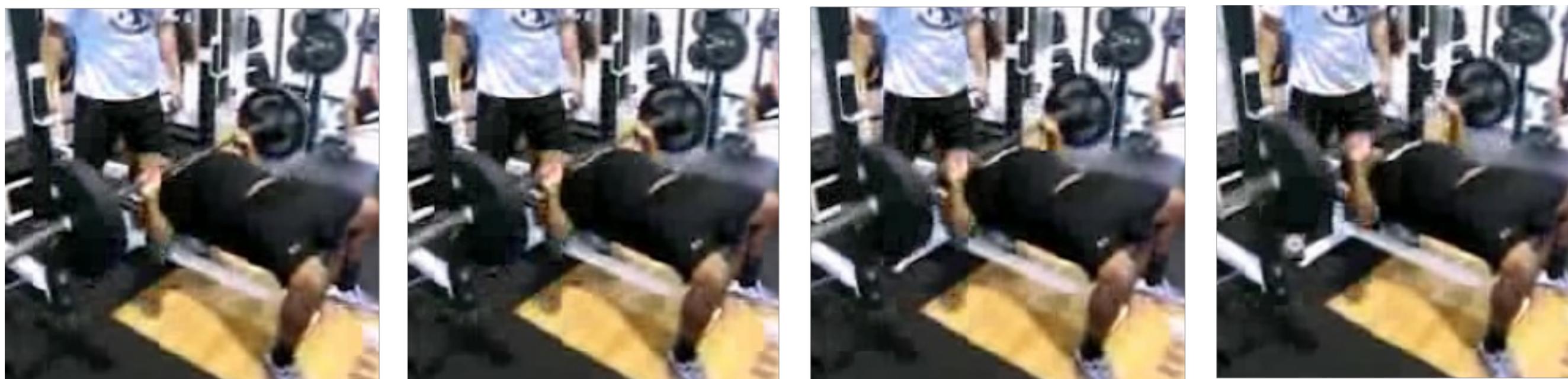


Figure from *GPipe: Easy Scaling with Micro-Batch Pipeline Parallelism*, Yanping et al, 2019

# Videos are temporally smooth: do we need buffering?

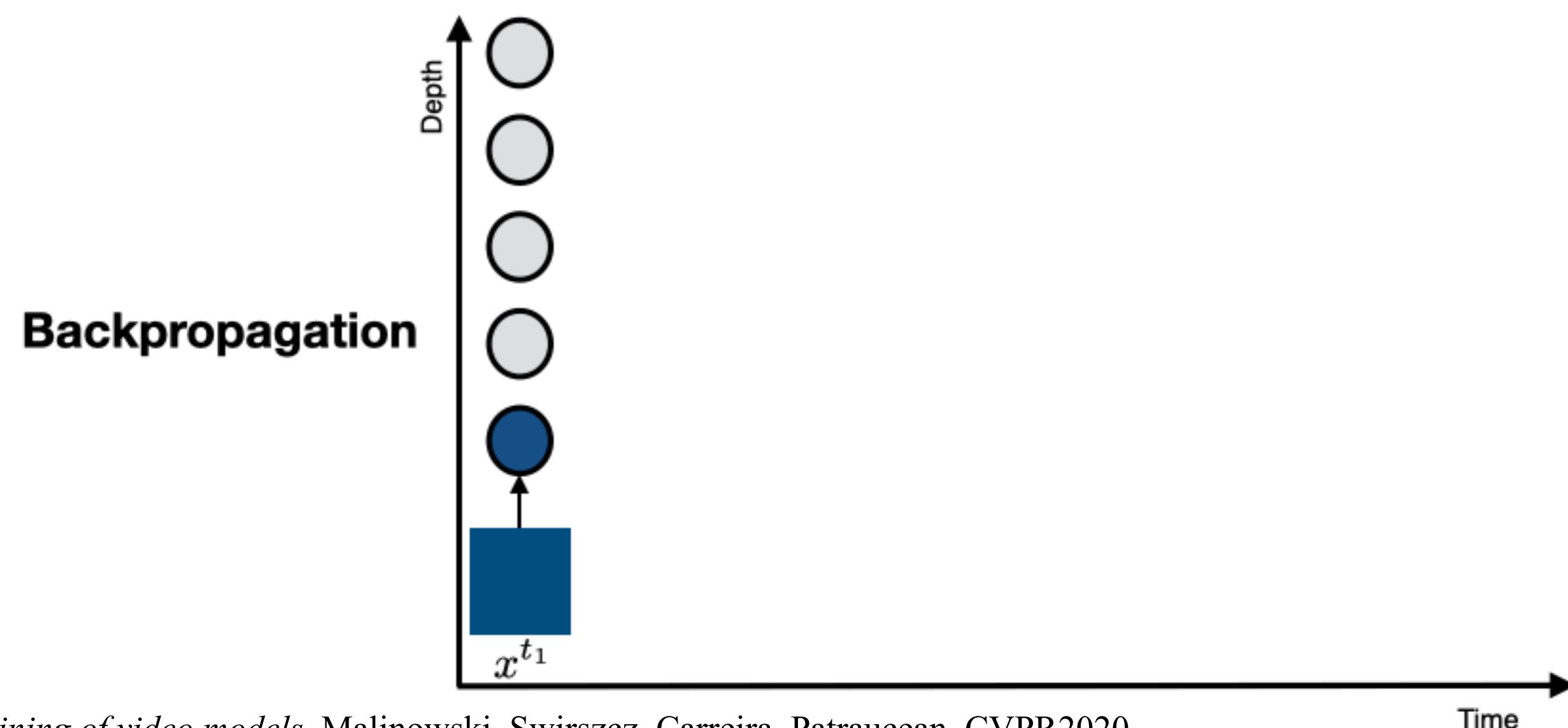
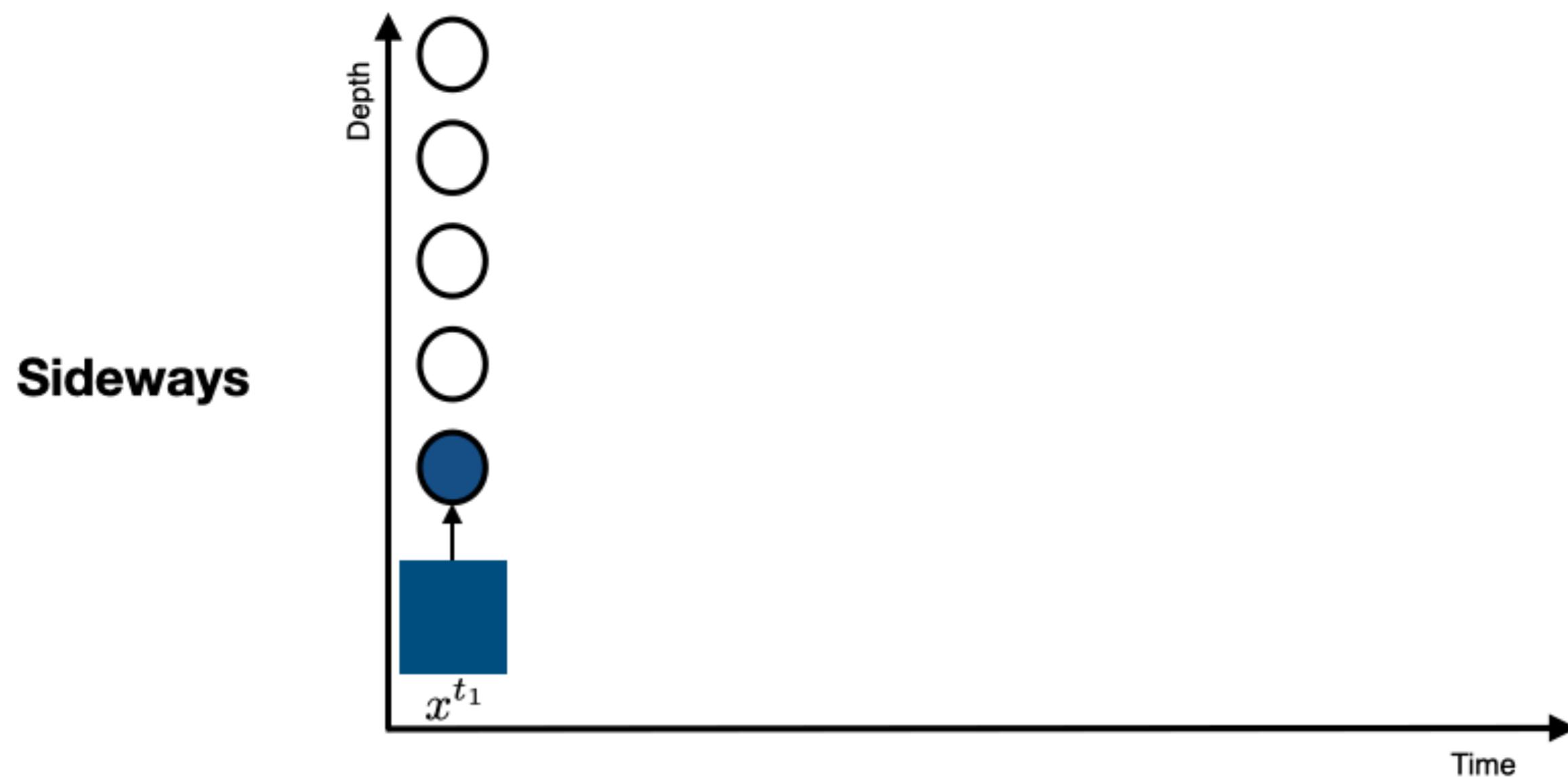


Video from Kinetics600 dataset

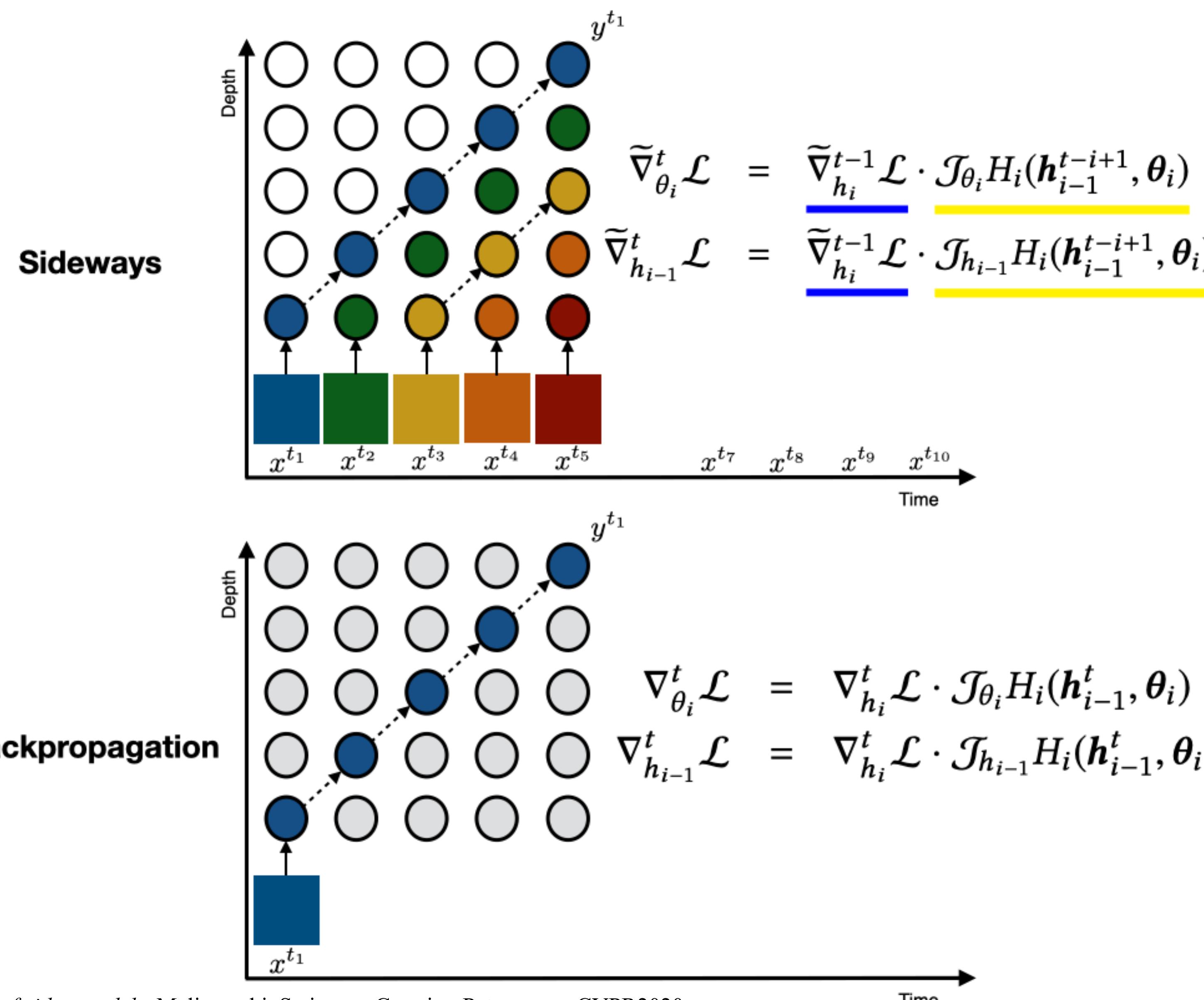


Sample consecutive frames extracted from the video

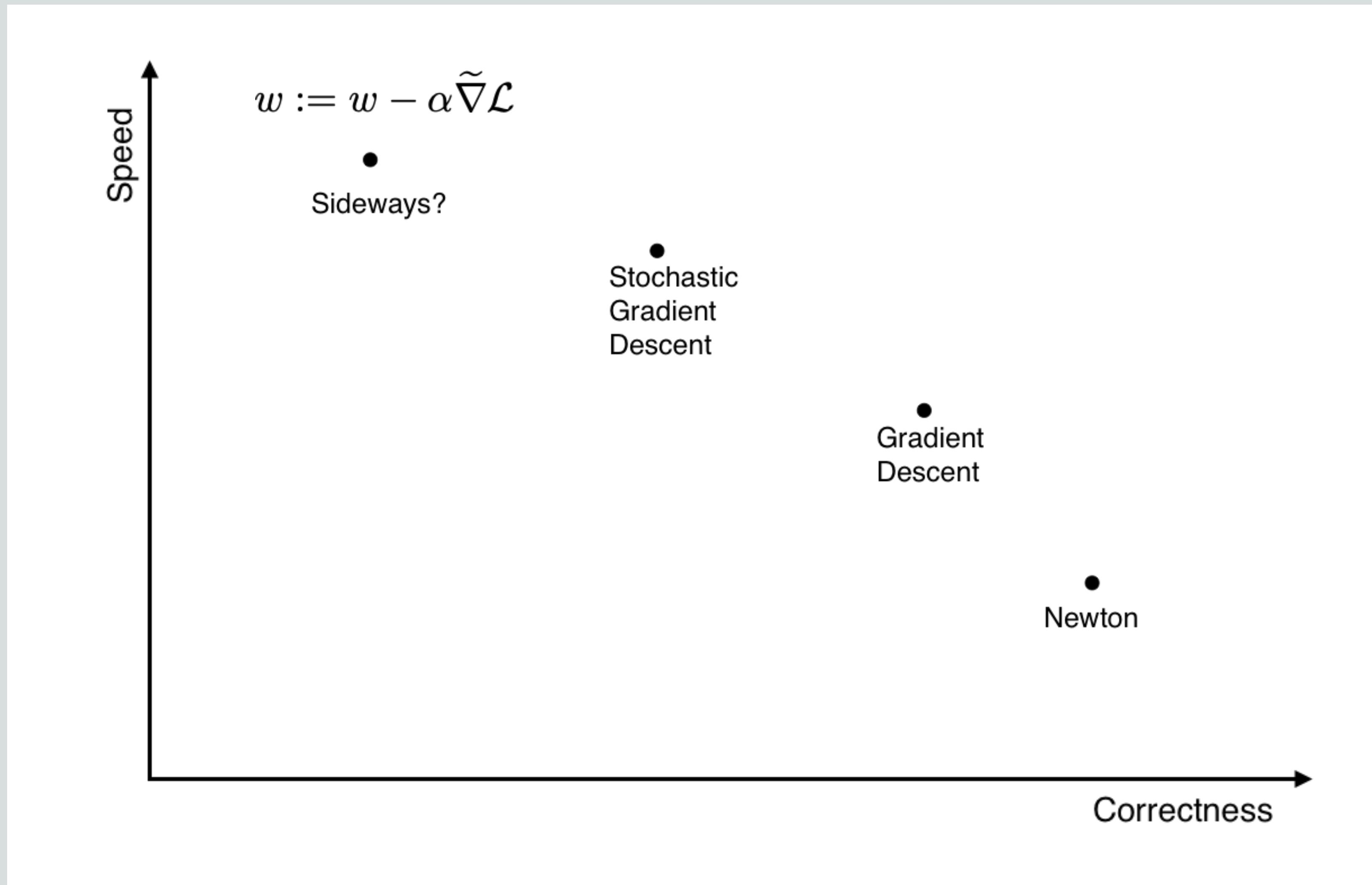
# Sideways vs. Backprop



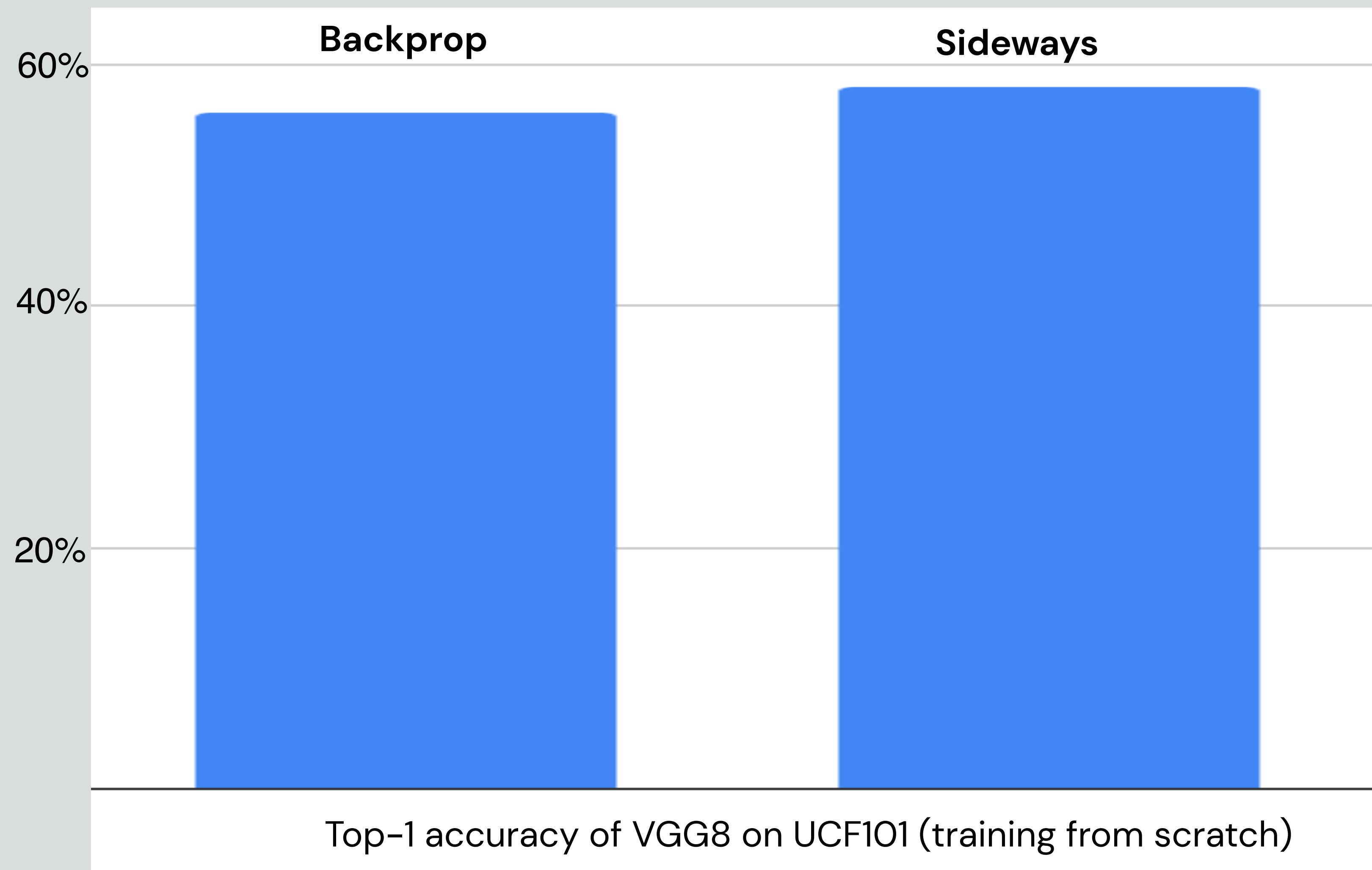
# Sideways vs. Backprop



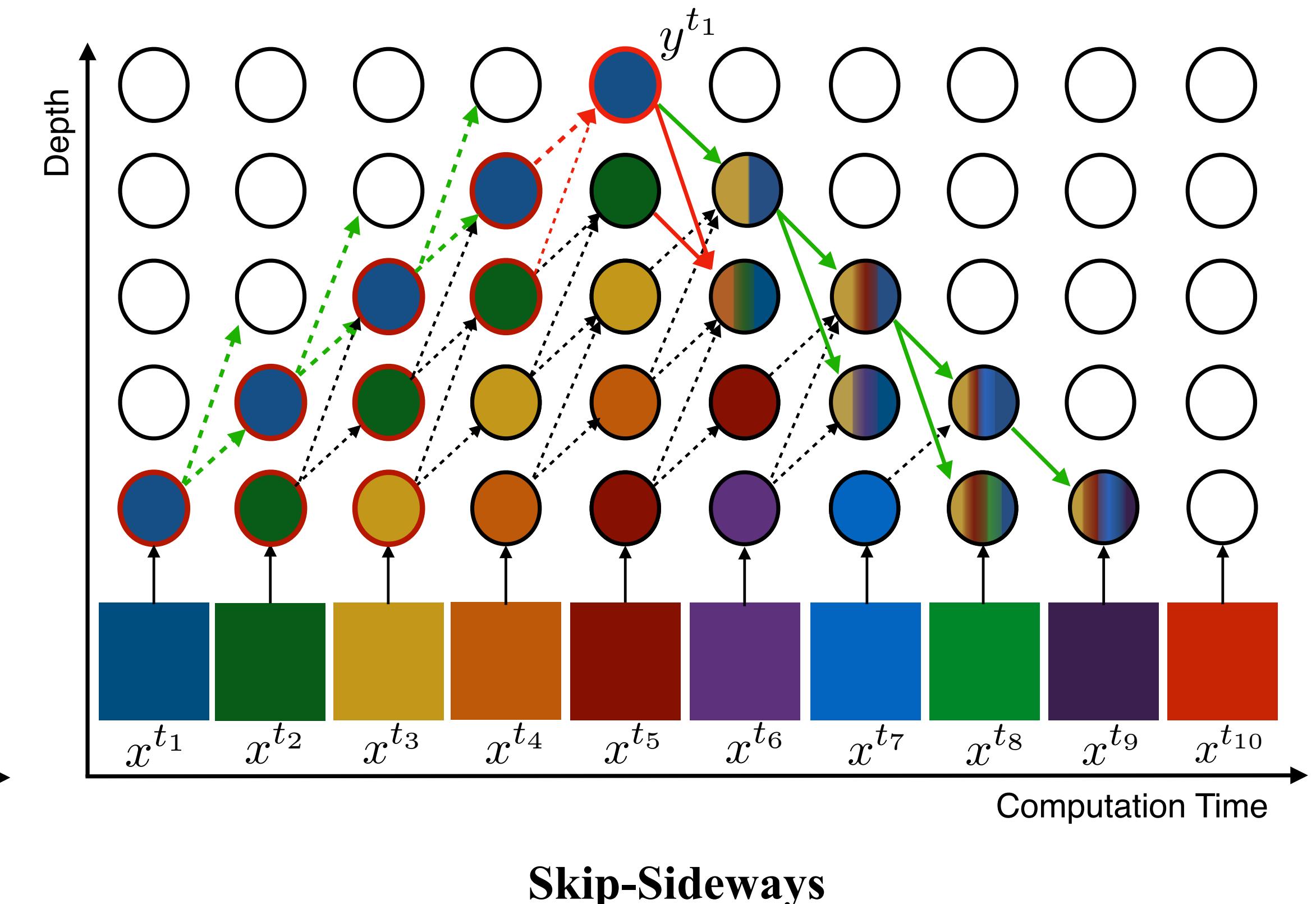
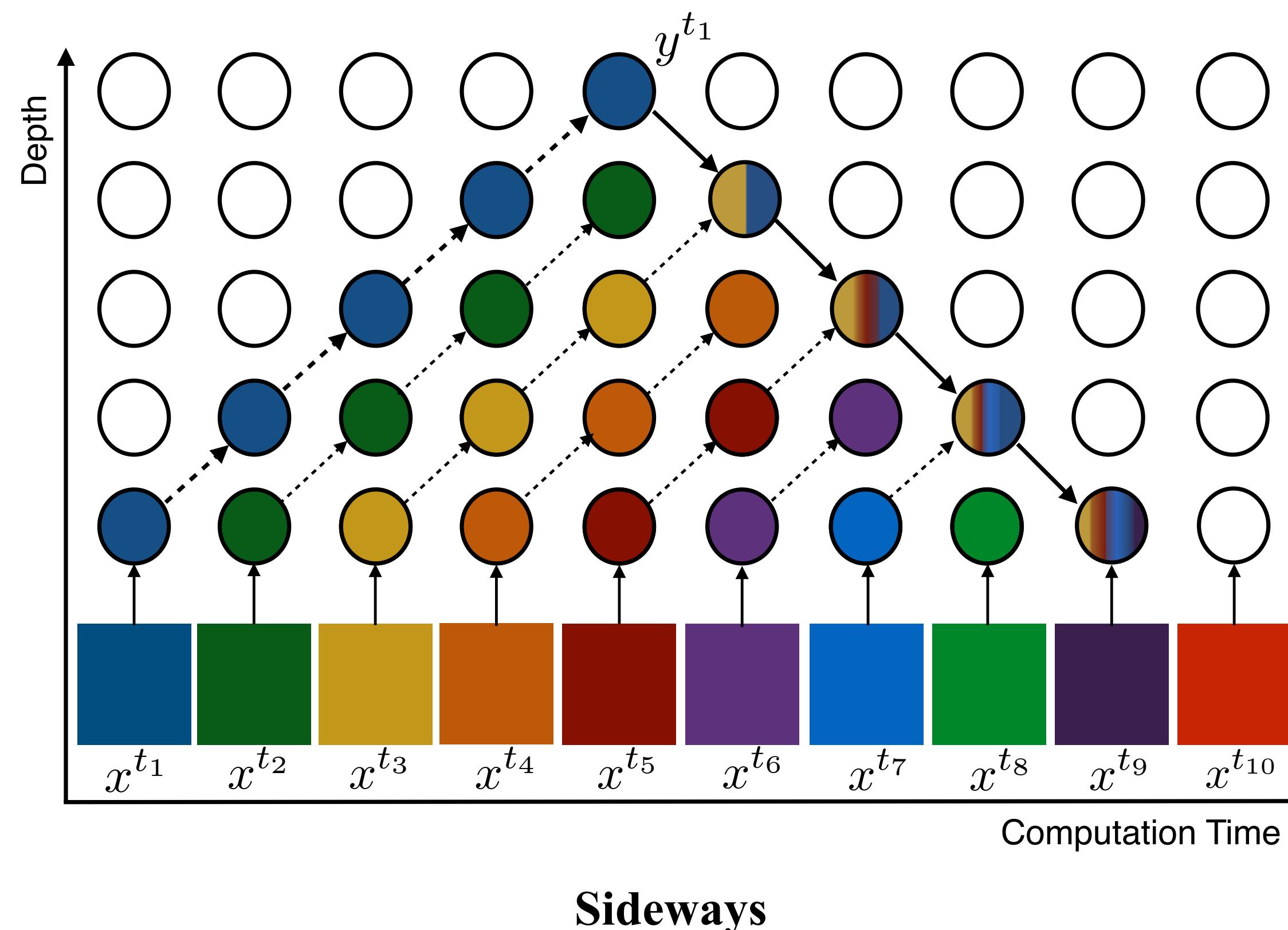
# Mathematical correctness of gradients vs. speed



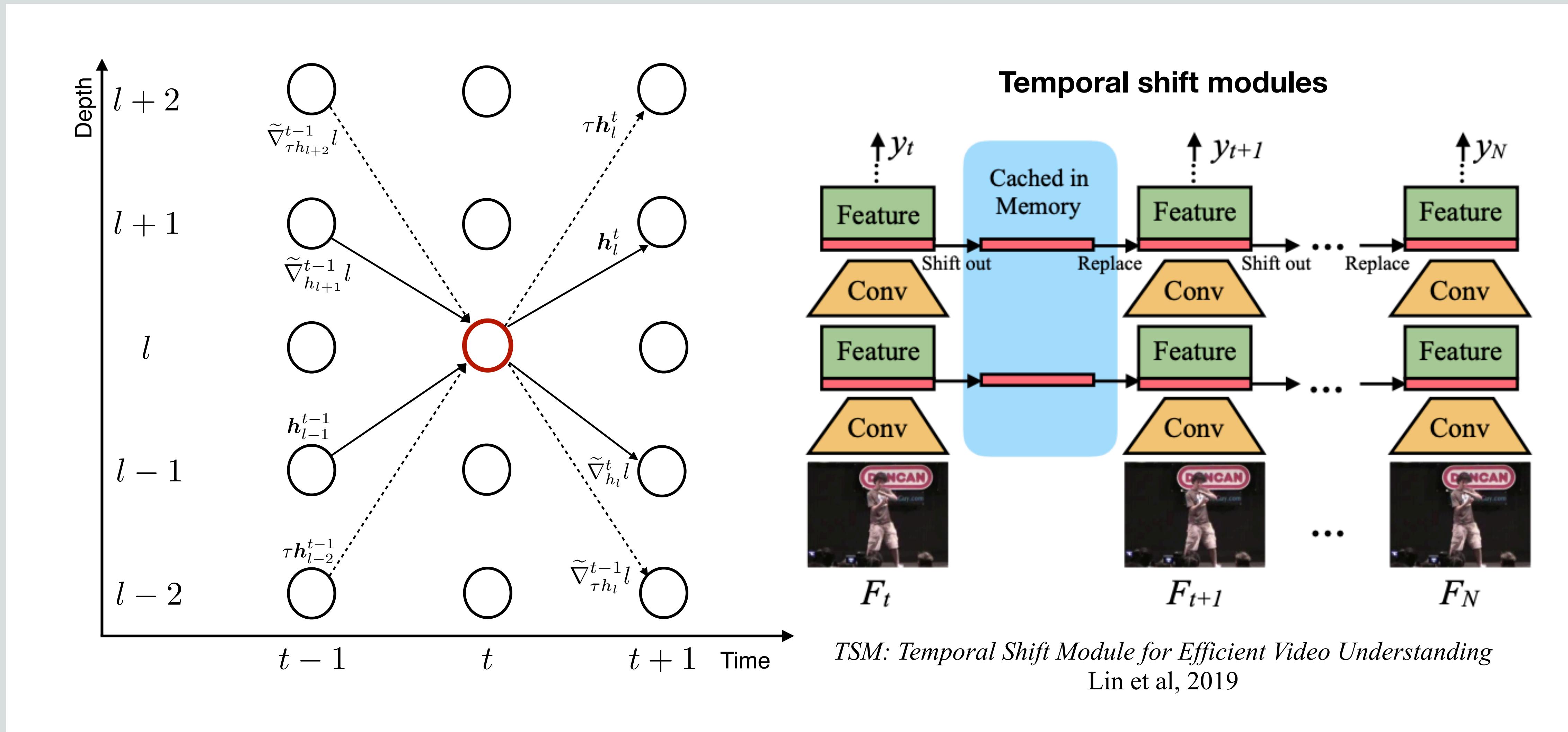
# Sideways vs. Backprop



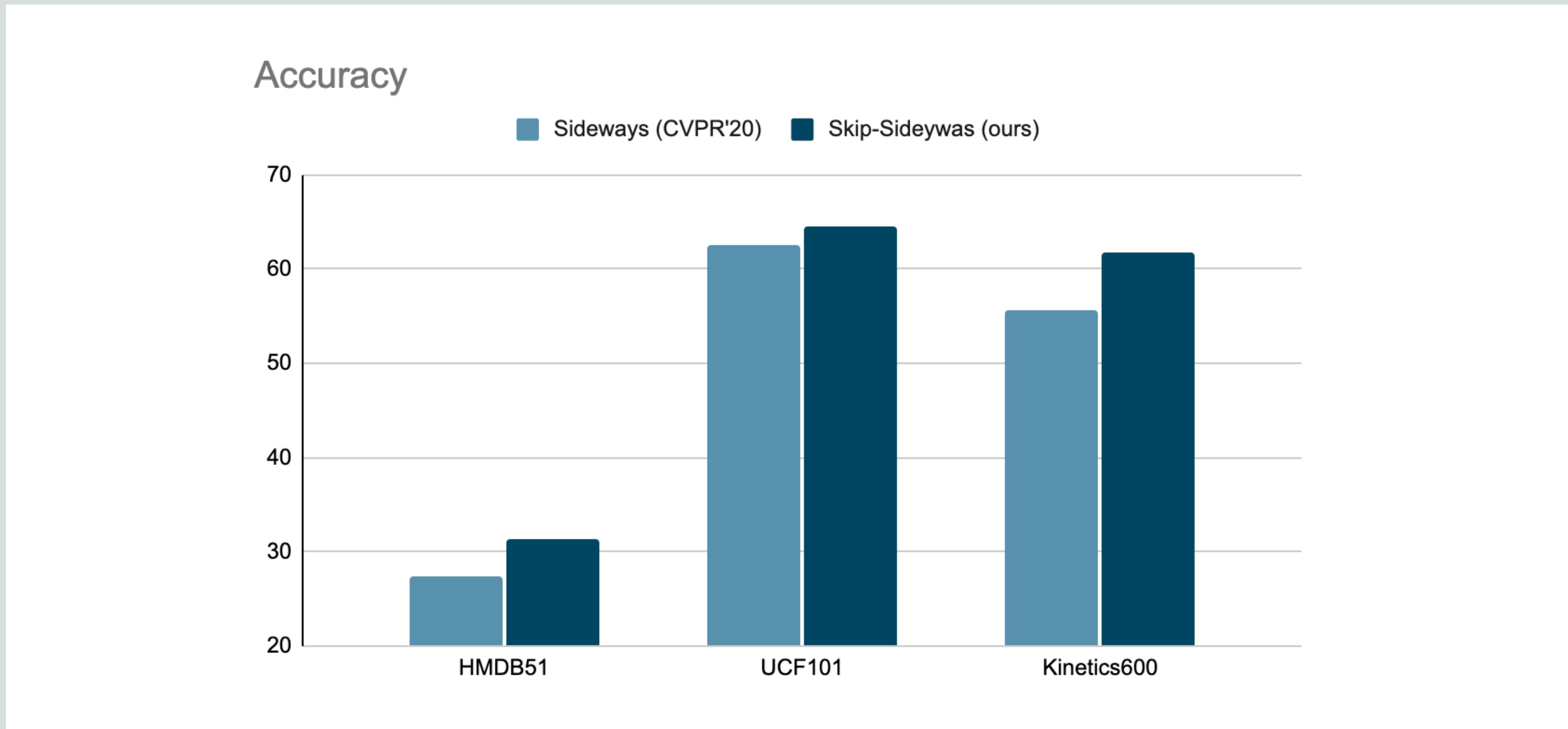
# Skip-Sideways: 2D *temporal* model



# Skip-Sideways vs. Temporal shift modules (TSM)



# Skip-Sideways: 2D *temporal* model



# (Skip-)Sideways vs. GPipe vs. Backprop

	Inputs	Uses	Blocking	Requires buffering	Correct gradients
BP	any	chain rule	yes	yes	yes
GPipe	any	pipelining of operations	partial	yes	yes
(Skip-)Sideways	temporally smooth sequences	pipelining and smoothness of inputs	no	no	approximate

# Biological (im)plausibility of backprop

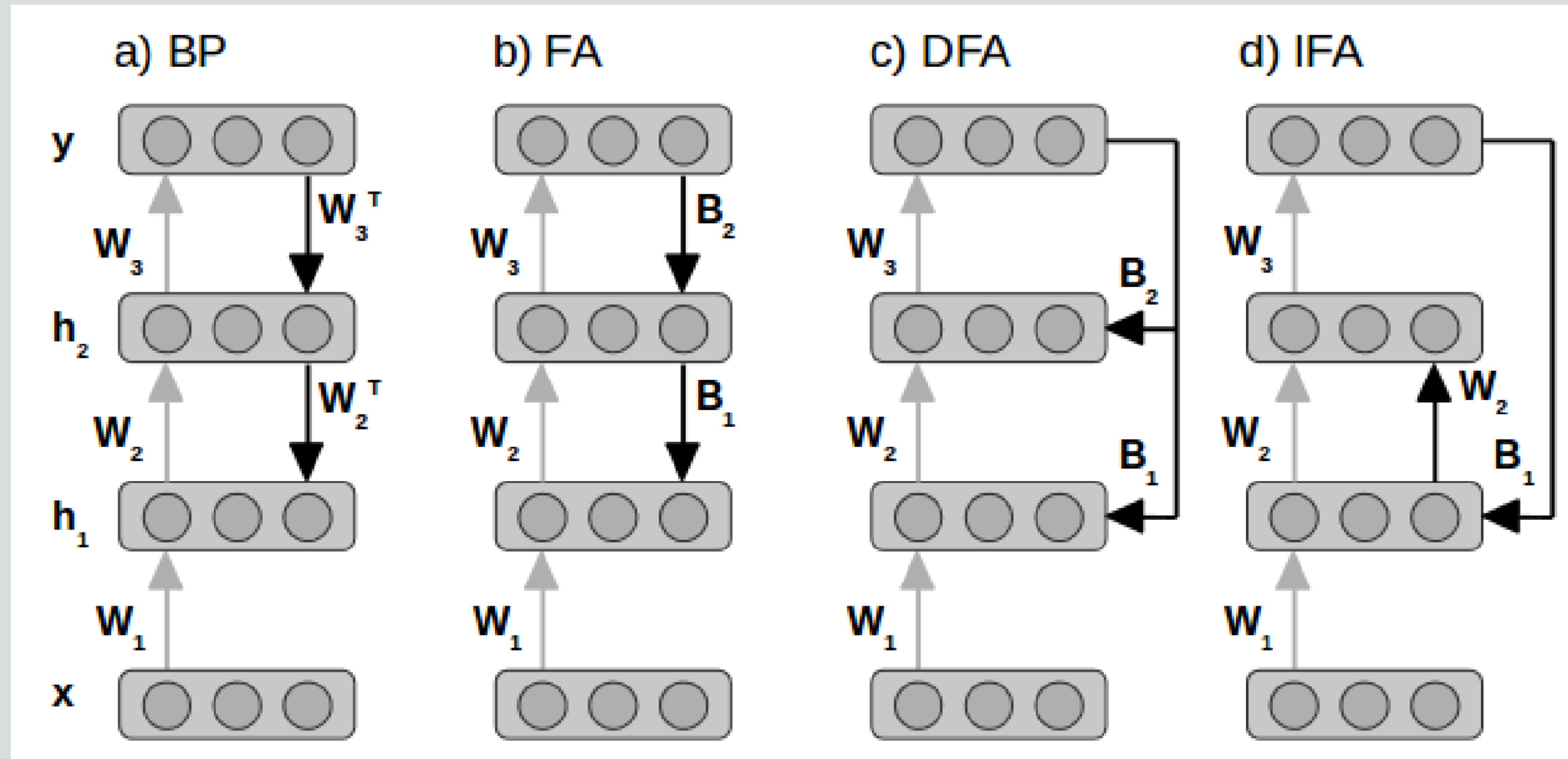


Figure from *Direct Feedback Alignment Provides Learning in Deep Neural Networks*, Nokland, 2016

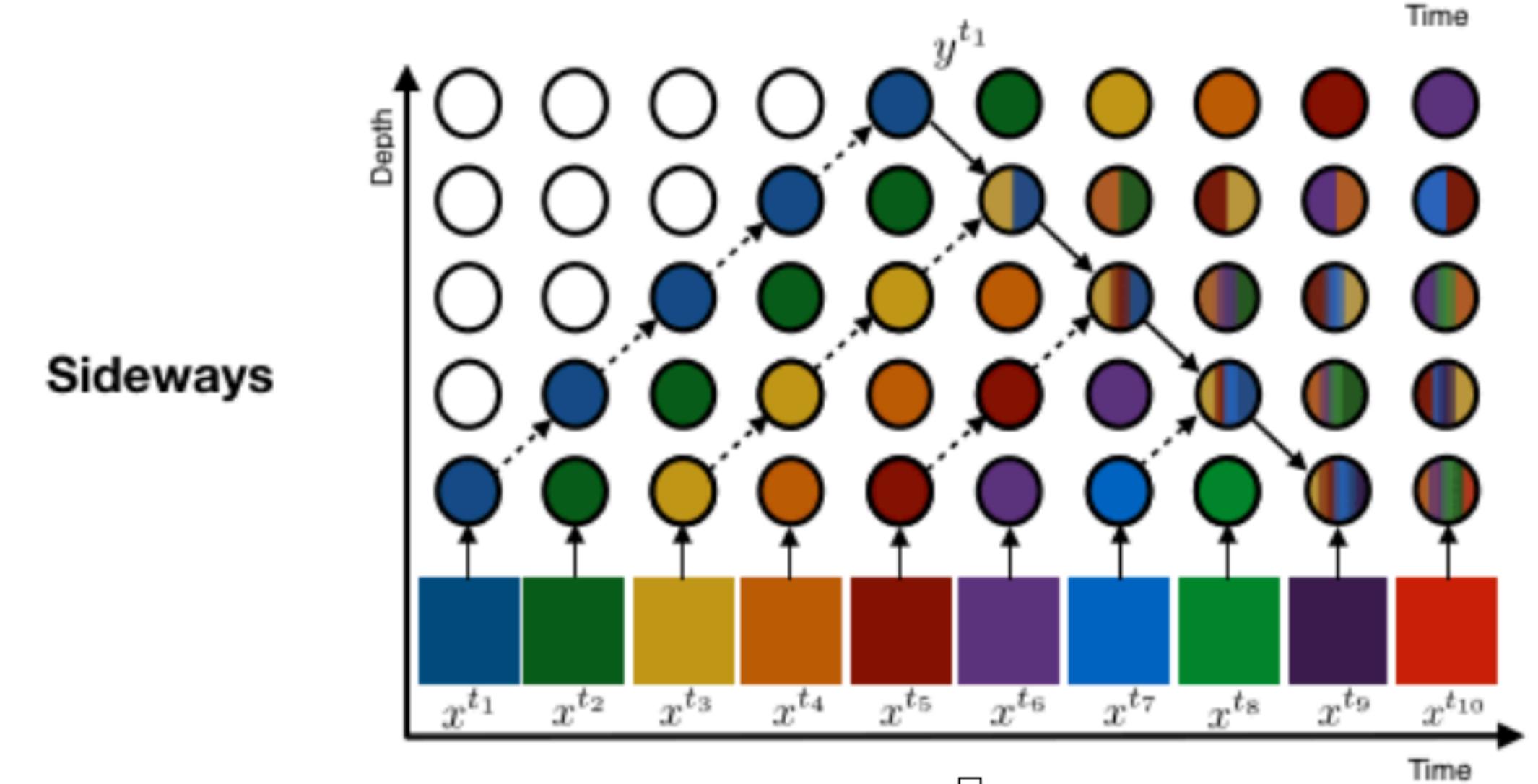
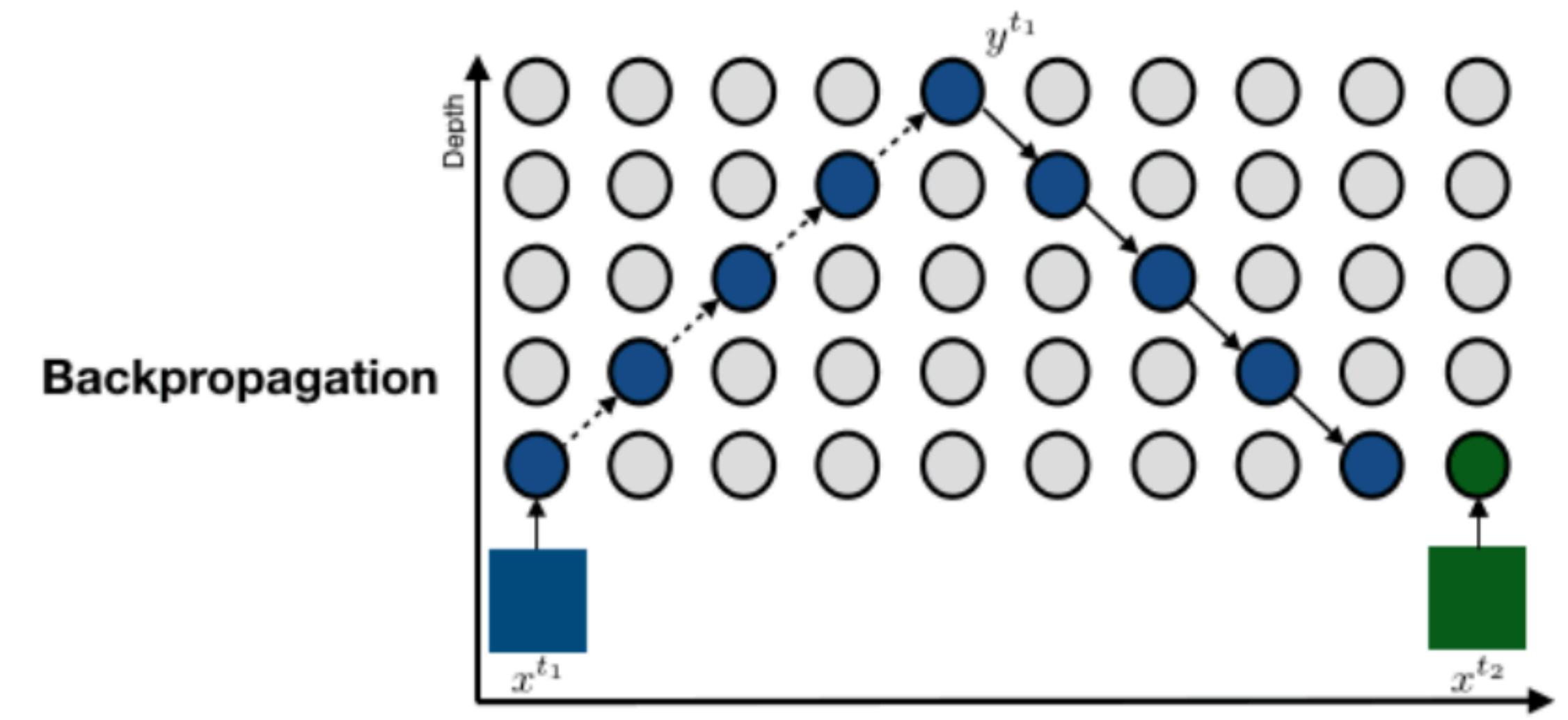
# Sideways: more biologically plausible

BP: instantaneous computation  
(blocking)

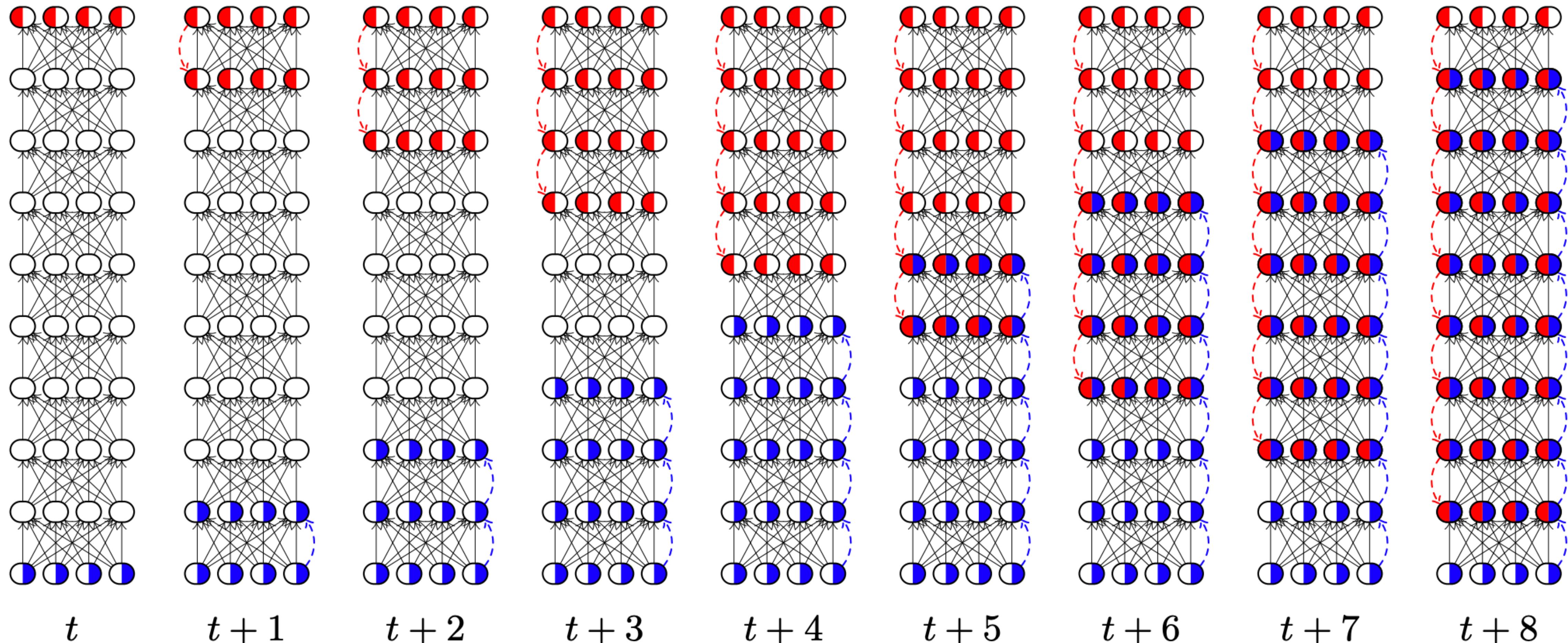
BPTT: sends gradients back in time

RTL: sends gradients forward, but  
computationally intractable

Sideways: sends *approximate*  
gradients forward



# Backprop as diffusion



# 04

# Conclusion



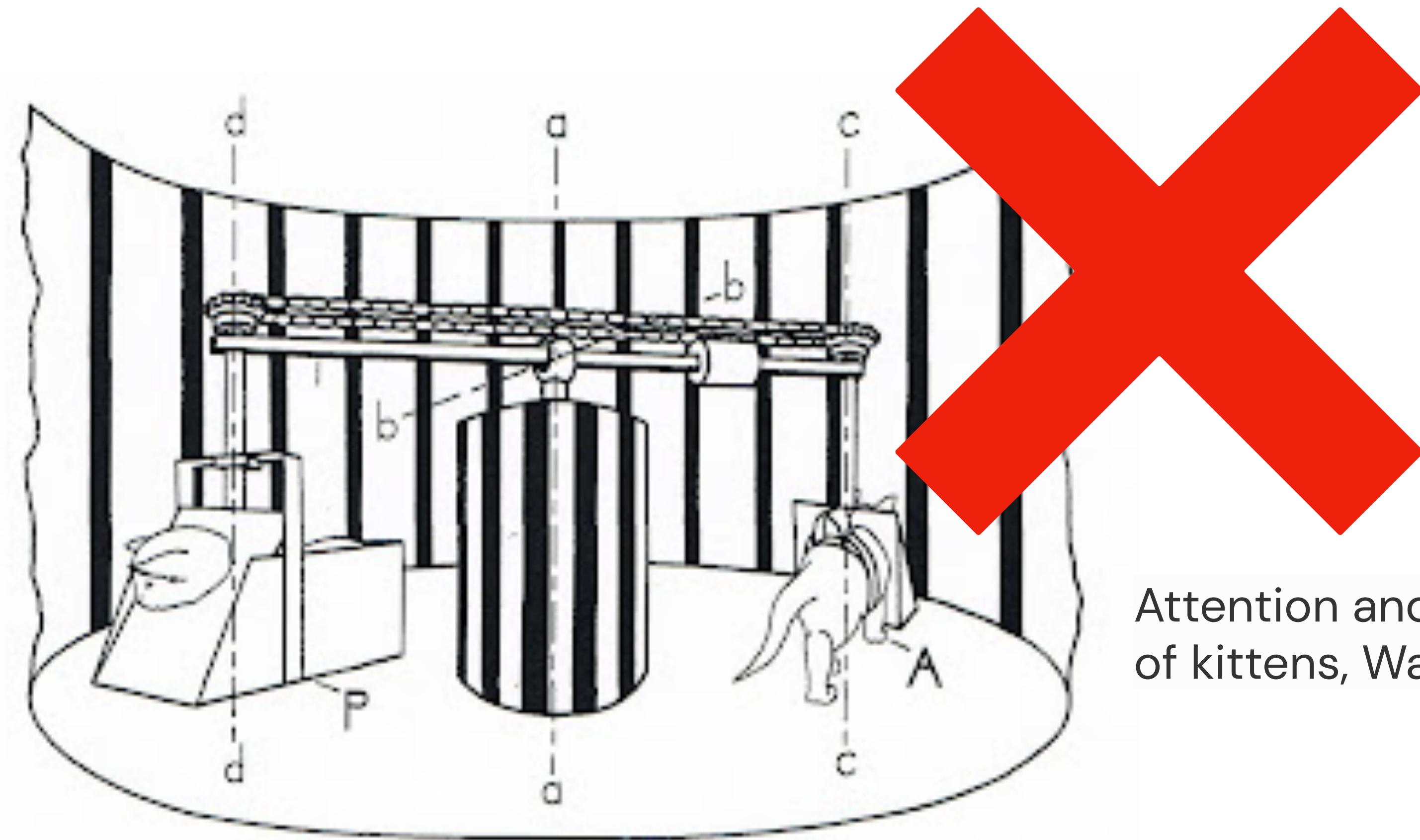
# Summary and Future work

- Design principles for efficient training and inference through
  - self-supervision
  - pipelining
- Exploit the temporal smoothness of videos to speed up training
- Latency and throughput are not dependent on depth: we can increase depth (and accuracy) without affecting the latency and throughput
- **Next steps:** theoretical framework for convergence with approximate gradients
- **Next steps:** Skip-Sideways + self-supervision -> local learning rule

# Roadmap/curriculum to solve vision



# Learning vision requires active locomotion



Attention and the depth perception  
of kittens, Walk et al, 1988

[Held, R. and Hein A. (1963). Movement-produced stimulation in the development of visually guided behavior. Journal of Comparative and Physiological Psychology 56(5): 872-876.]

# Final advice

**Question everything! Don't take things for granted!**