

Massively parallel video networks

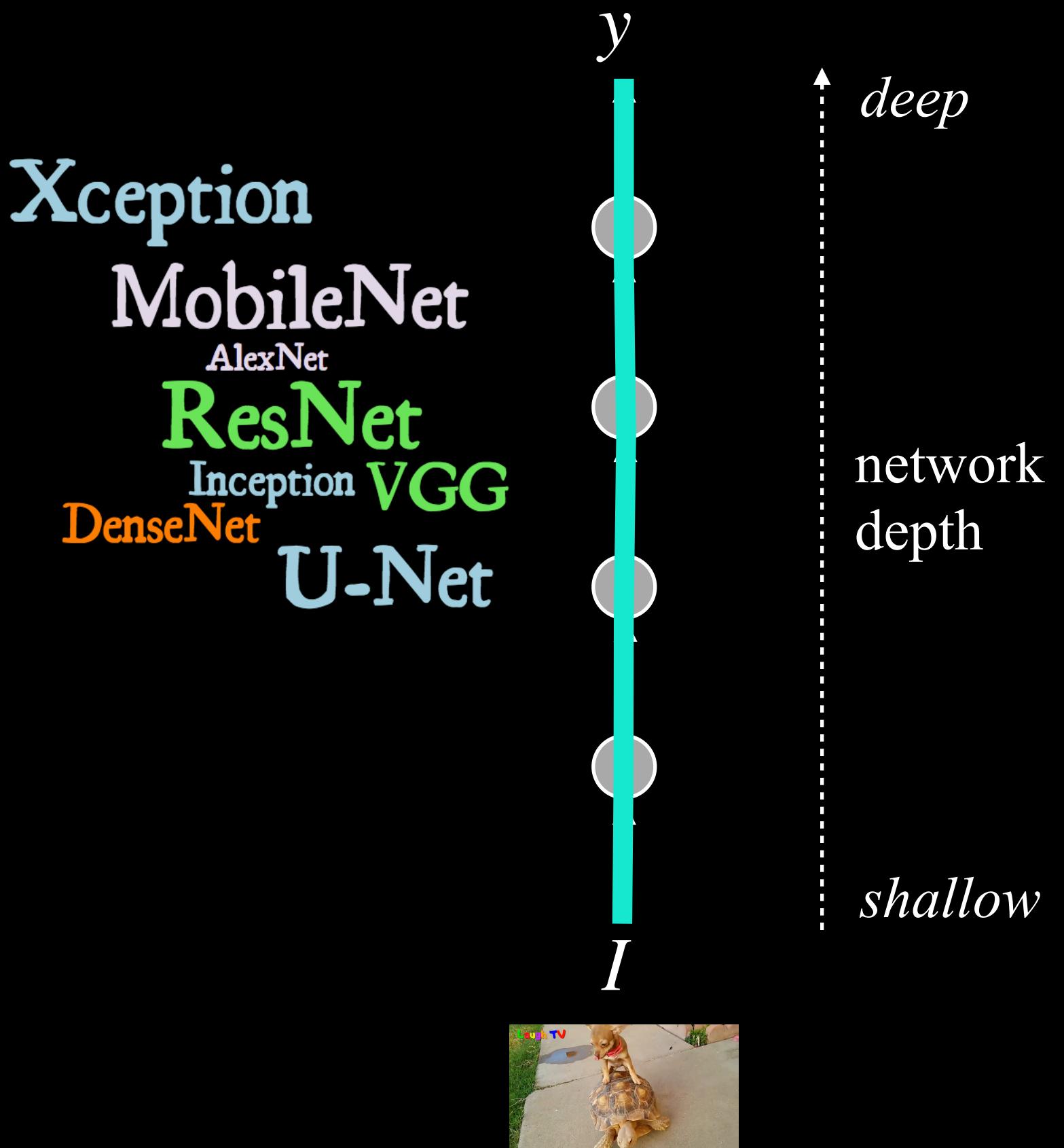
Viorica Pătrăucean
vierica@google.com

joint work with João Carreira, Laurent Mazare
Andrew Zisserman, Simon Osindero



Processing of visual input

- ▶ Vision: dominant sense for humans
- ▶ Most works focus on image processing
 - videos: image models applied frame-by-frame



Benefits from working with videos

→ better accuracy (temporal smoothness; reduced ambiguity)

- ▶ action recognition, optical flow estimation etc.

→ more efficient

- ▶ e.g. MotionJPEG vs. MPEG



Figure 1. A still from 'Quo Vadis' (1951). Where is this going?



Autonomous cars



The Duel: Timo Boll vs. KUKA Robot

Robotics applications

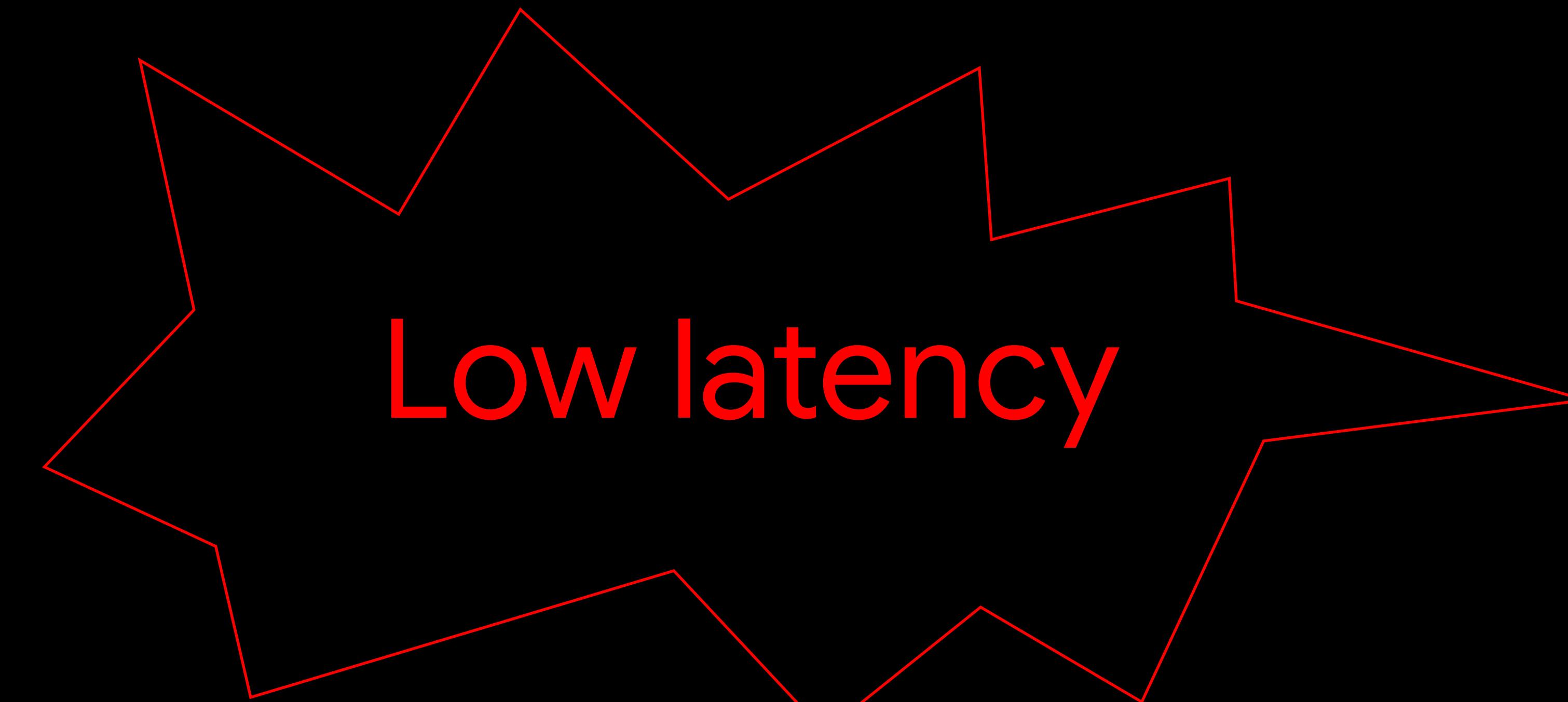


Minority report, Steven Spielberg (2002)

Augmented reality

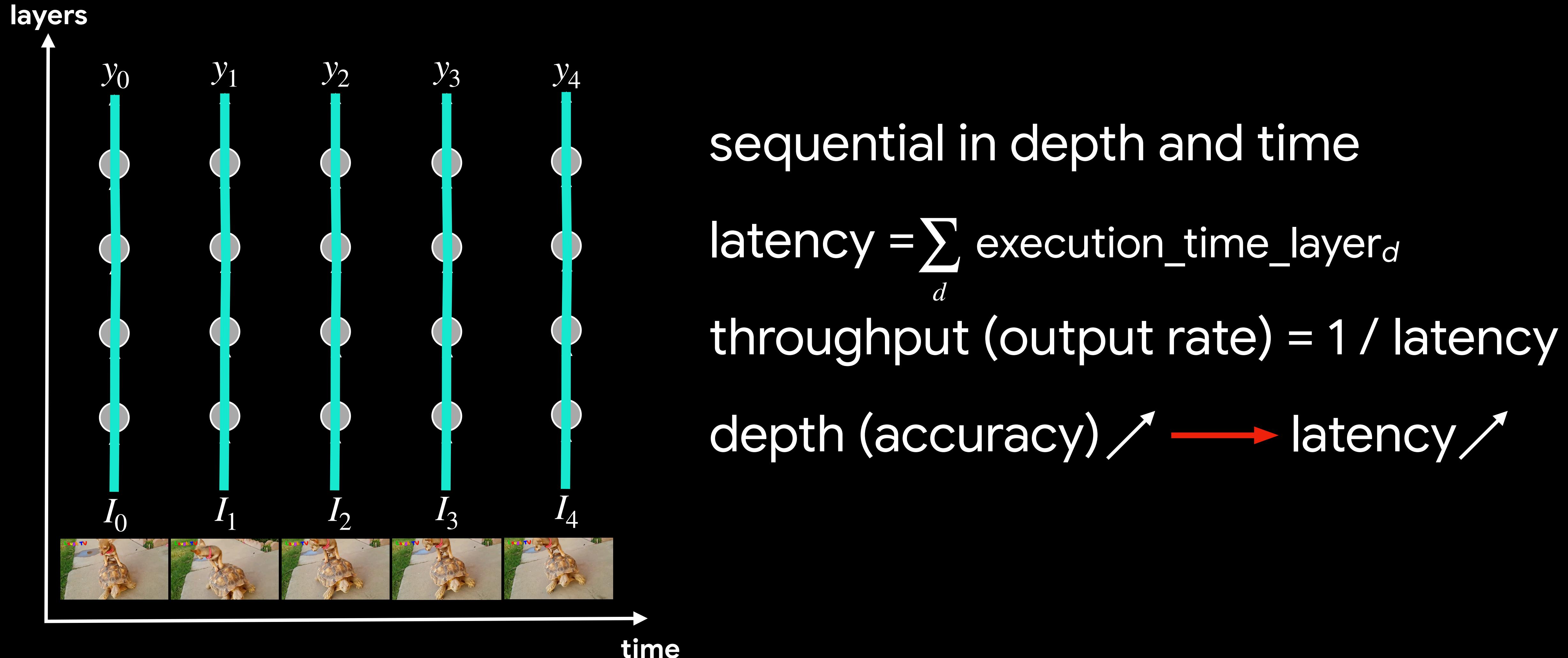
Setting

- online processing, one video stream, can't batch
- causal, no peeking into the future



Low latency

Deep video models = image models applied frame-by-frame



Current object detectors: ~5fps



25 fps



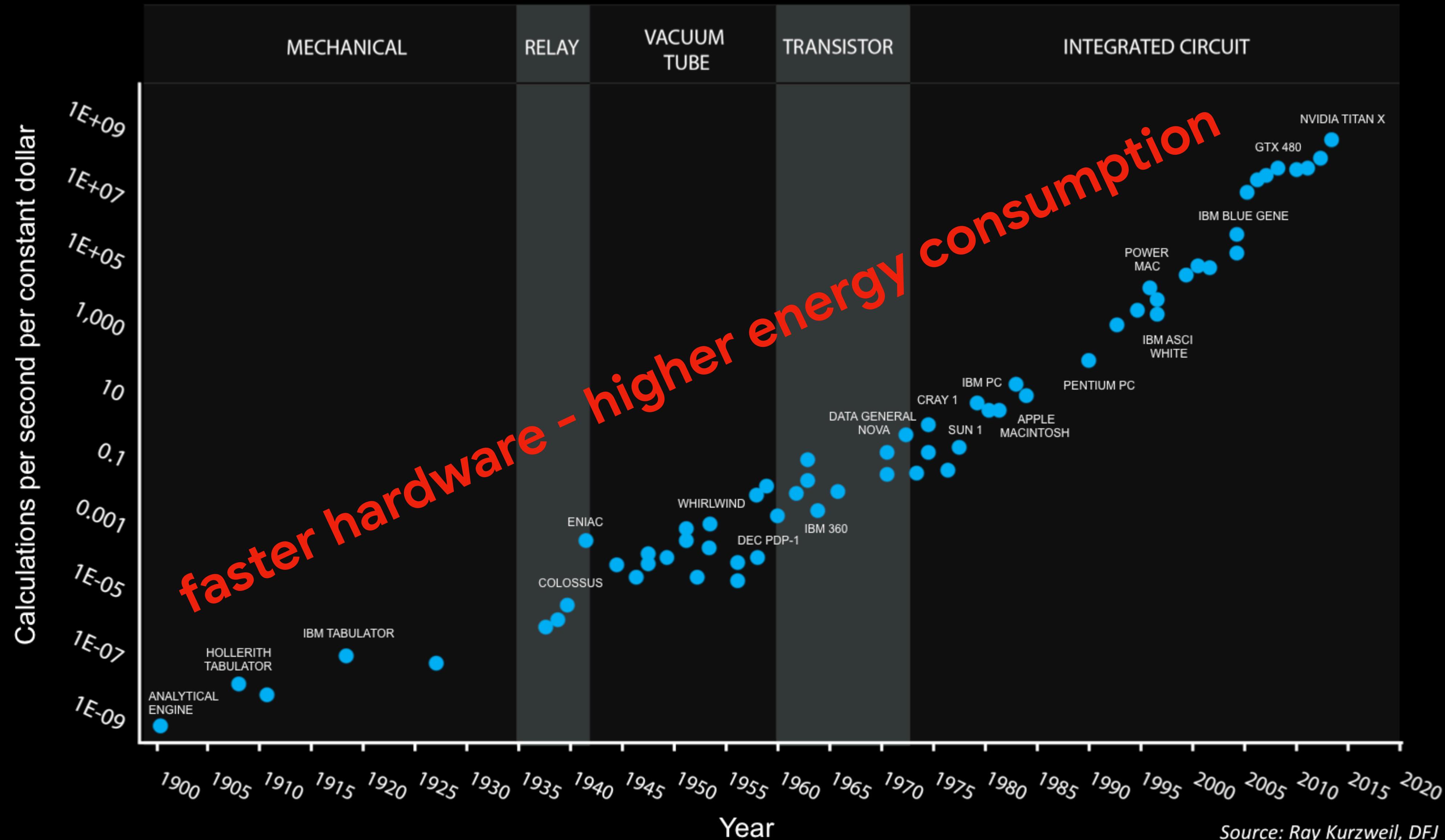
5 fps

Davis2017 dataset

Safe autonomous driving at 5fps



120 Years of Moore's Law



More sustainable ideas

Improve efficiency of image models

- model compression (Chen et al.), distillation (Hinton et al.), low representation format (Courbarieux et al), lighter convolutional architectures (Mobilenet, Xception)
- time budget methods: exit when time runs out (Karayev et al., Mathe et al.)

Improve efficiency of video models

- strided 3D conv (Tran et al., Carreira et al.), warping (Zhu et al.)
- different update rates for different features (Shelhamer et al.)

Our approach: orthogonal to these works

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.

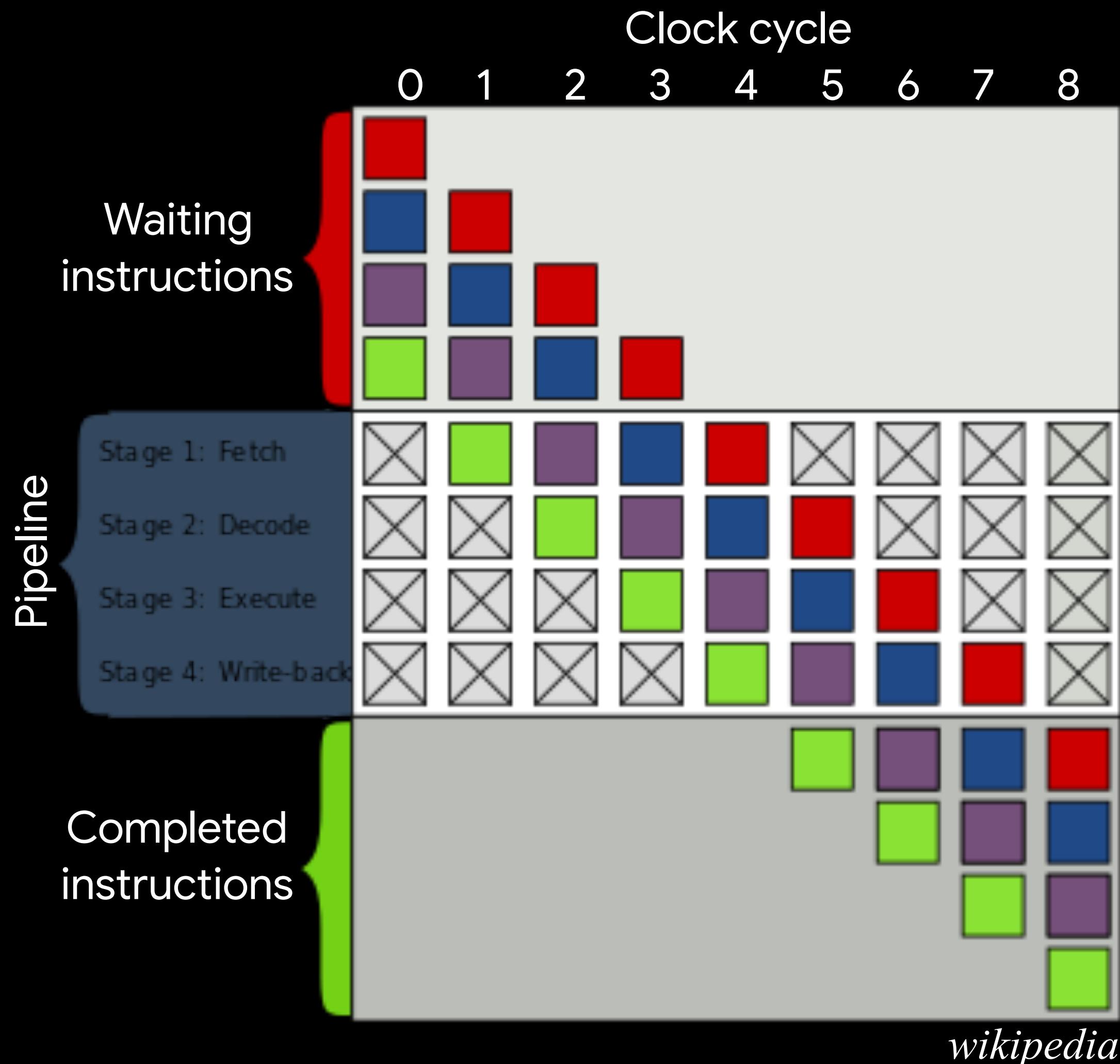
BY LIQUN LUO

APRIL 12, 2018

Is sequential mode optimal?

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

Maximise throughput
Efficient use of hw resources



Predictive deep learning

Increase throughput

Reduce latency

Reduce clock cycles

Increase throughput - option 1

Reduce latency

Reduce clock cycles

Setting

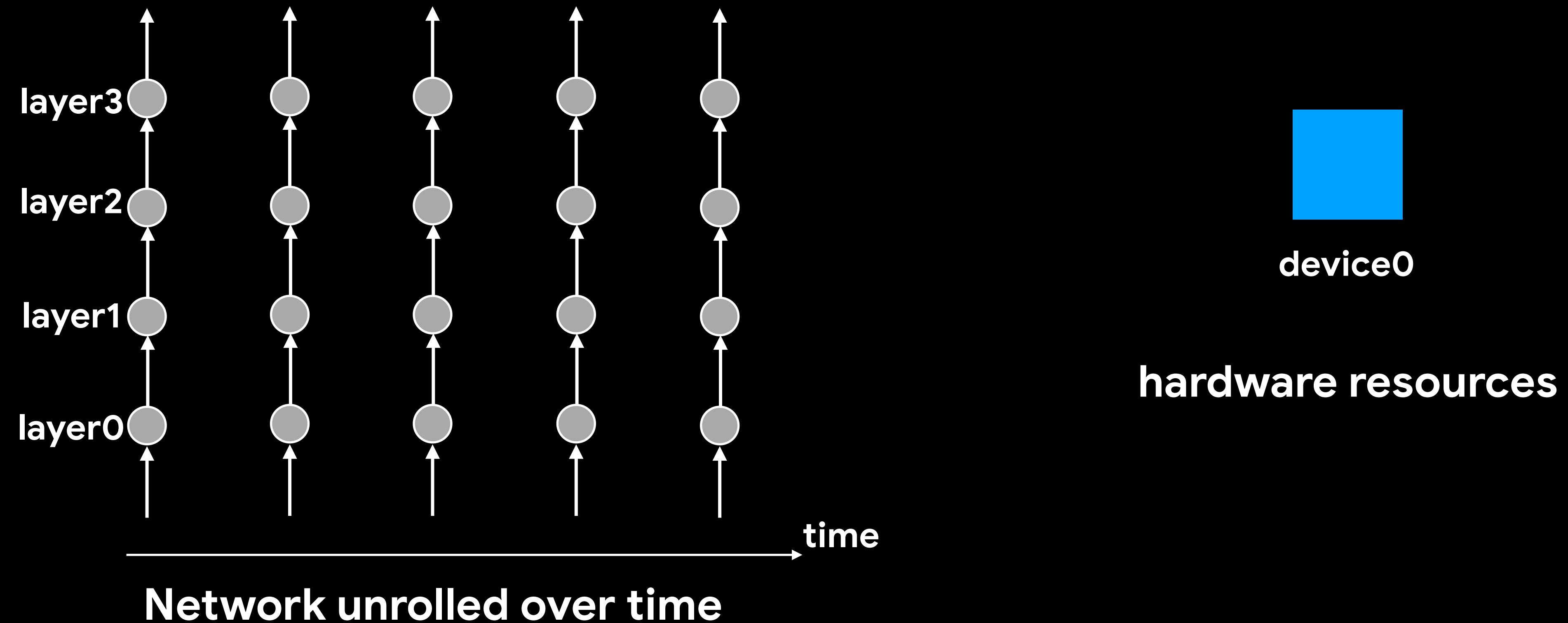
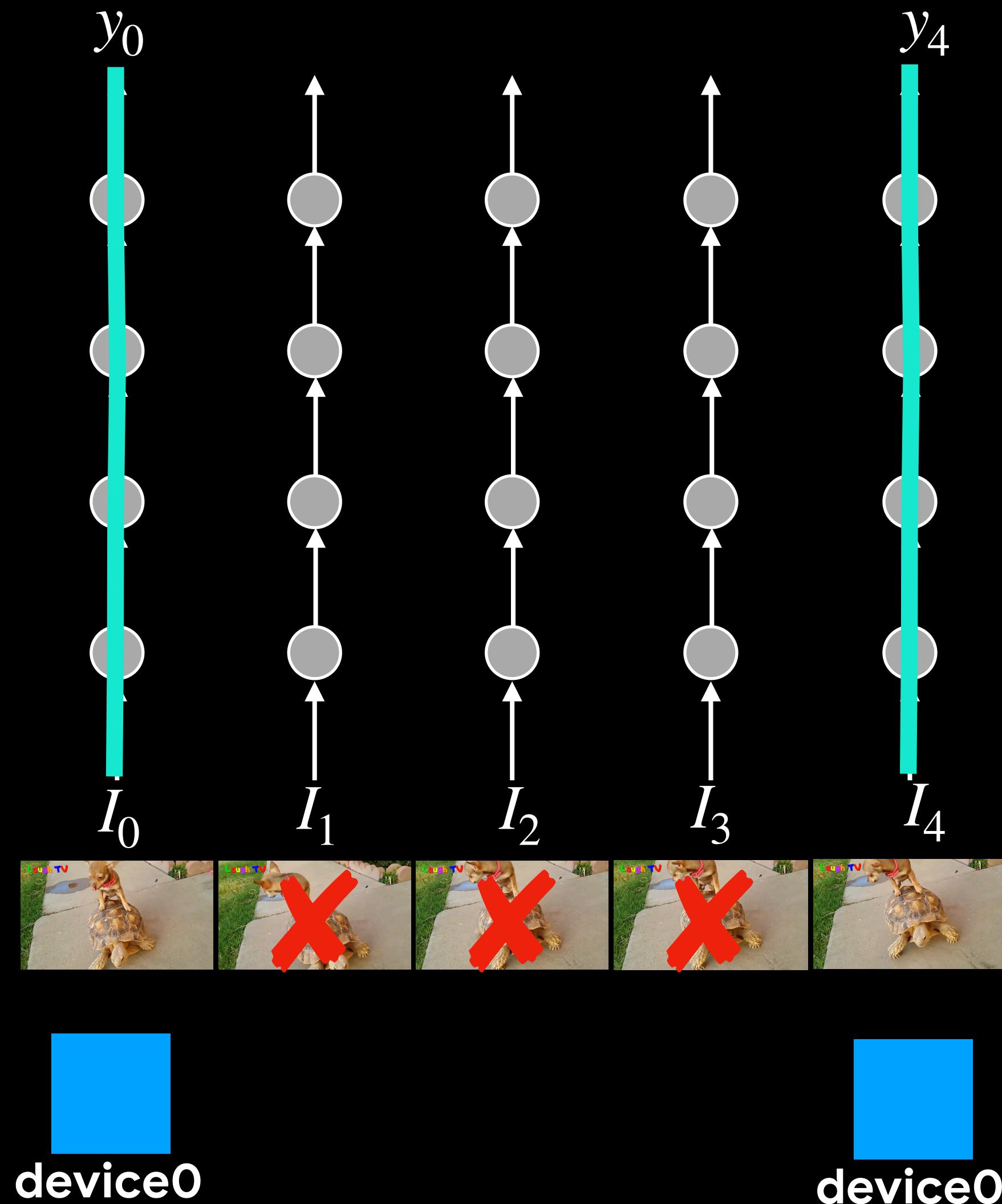


Image model frame-by-frame

non real-time model

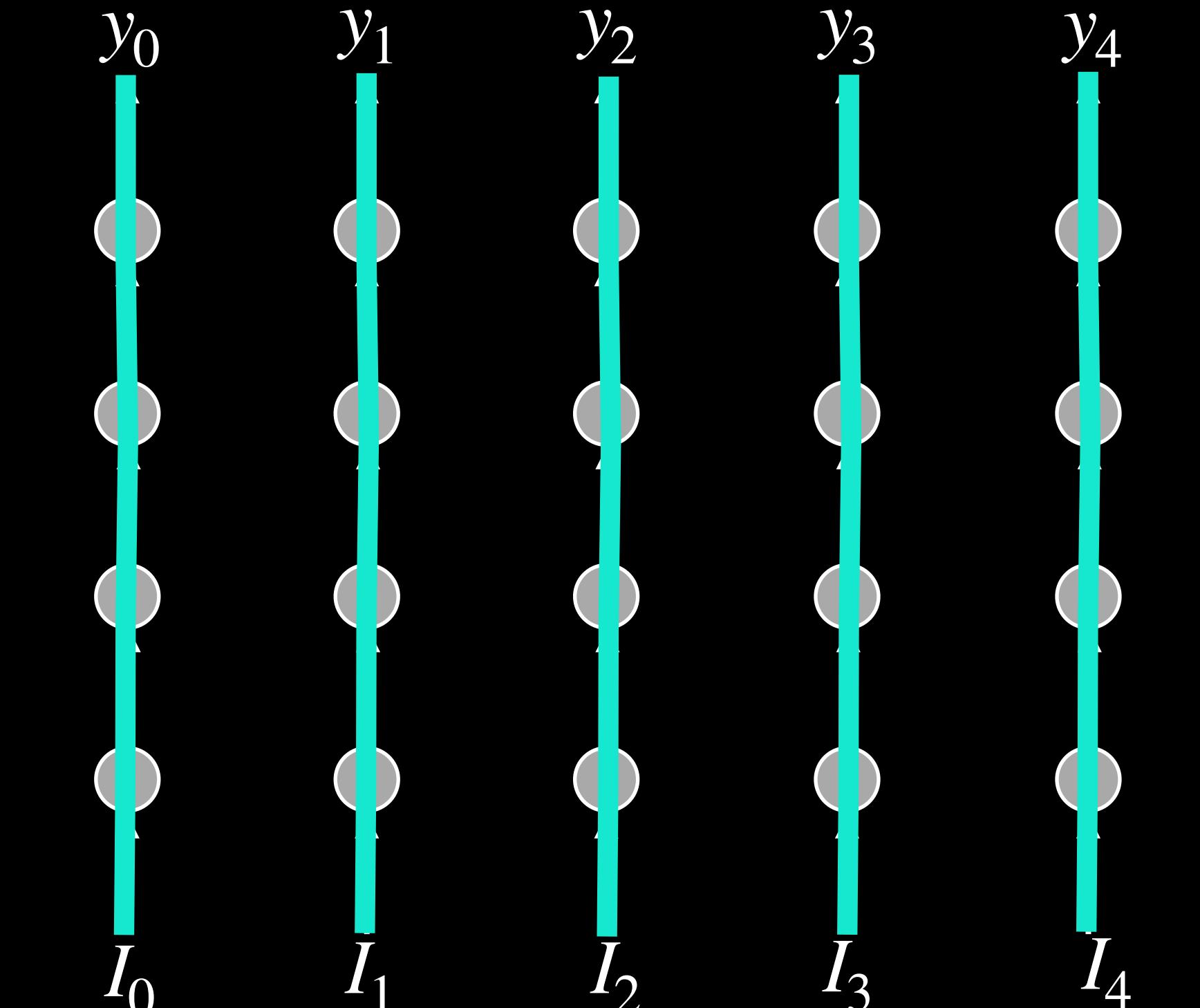


Con:
low throughput
high latency

Image model frame-by-frame

non real-time model

parallel resources



Pro:
high throughput

Con:
same (high) latency

round-robin over frames

Increase throughput

Reduce latency

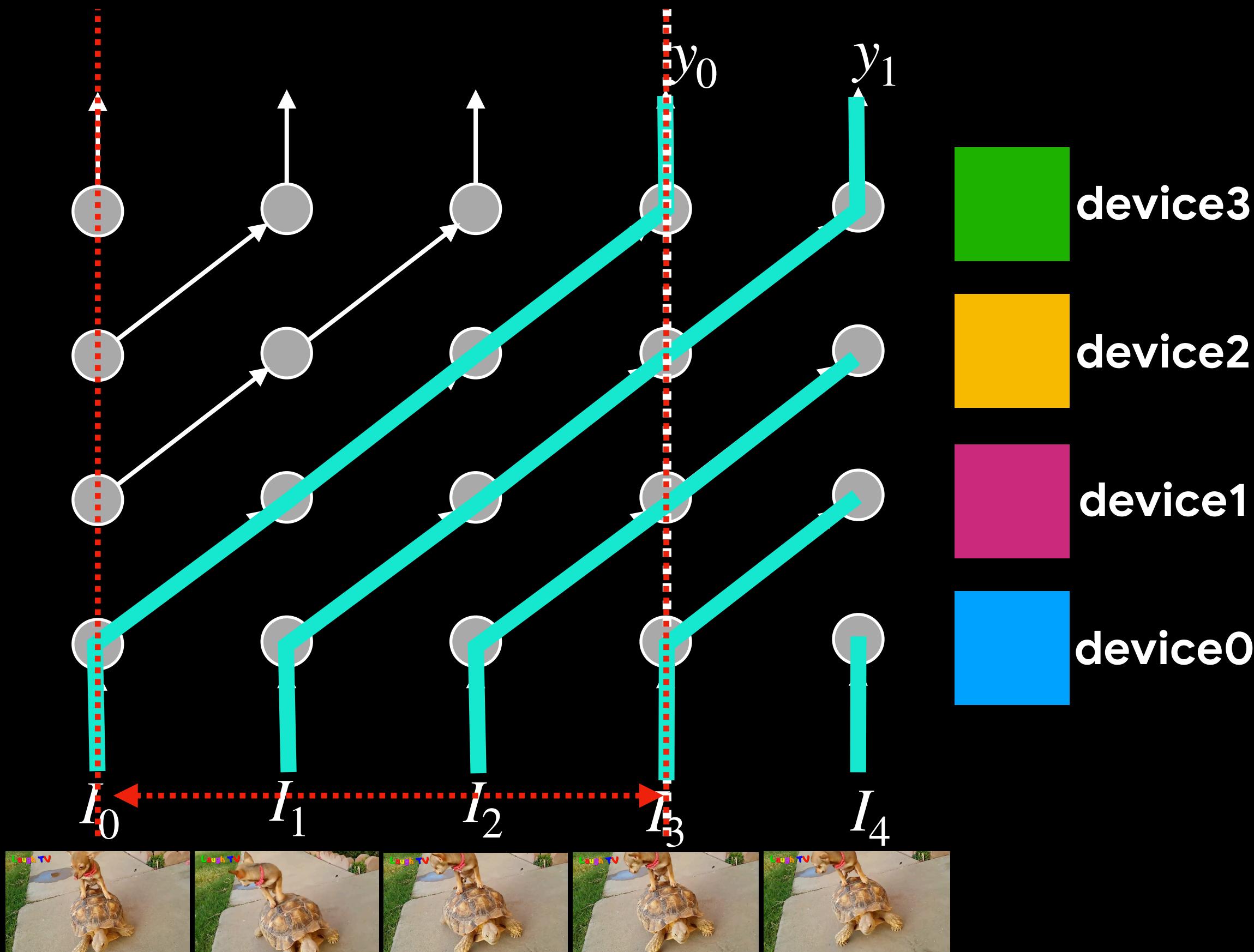
Reduce clock cycles

} our approach

Predictive pipelined model

Pro:
high throughput

Con:
same latency

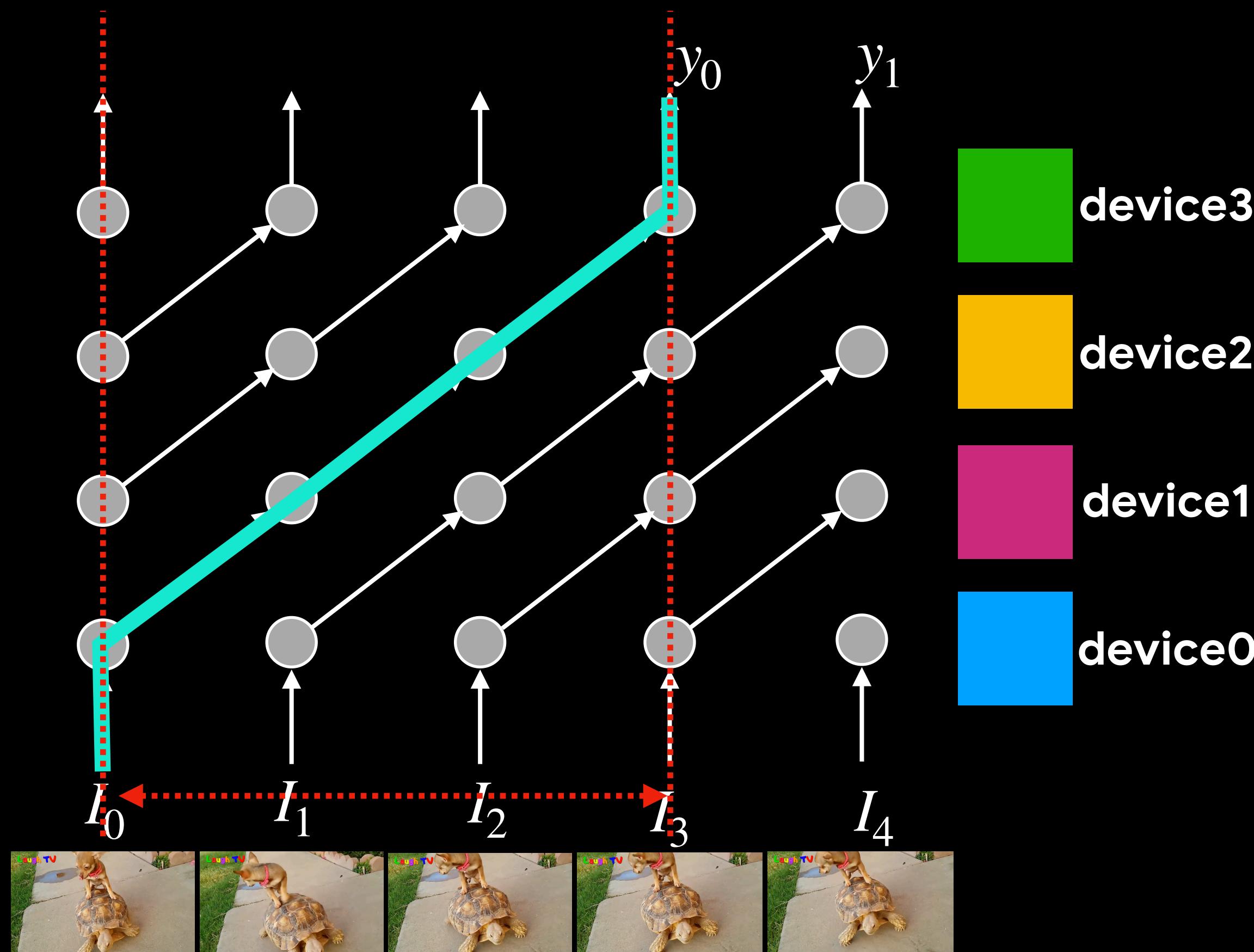


round-robin over layers

Predictive pipelined model

Pro:
high throughput
low latency

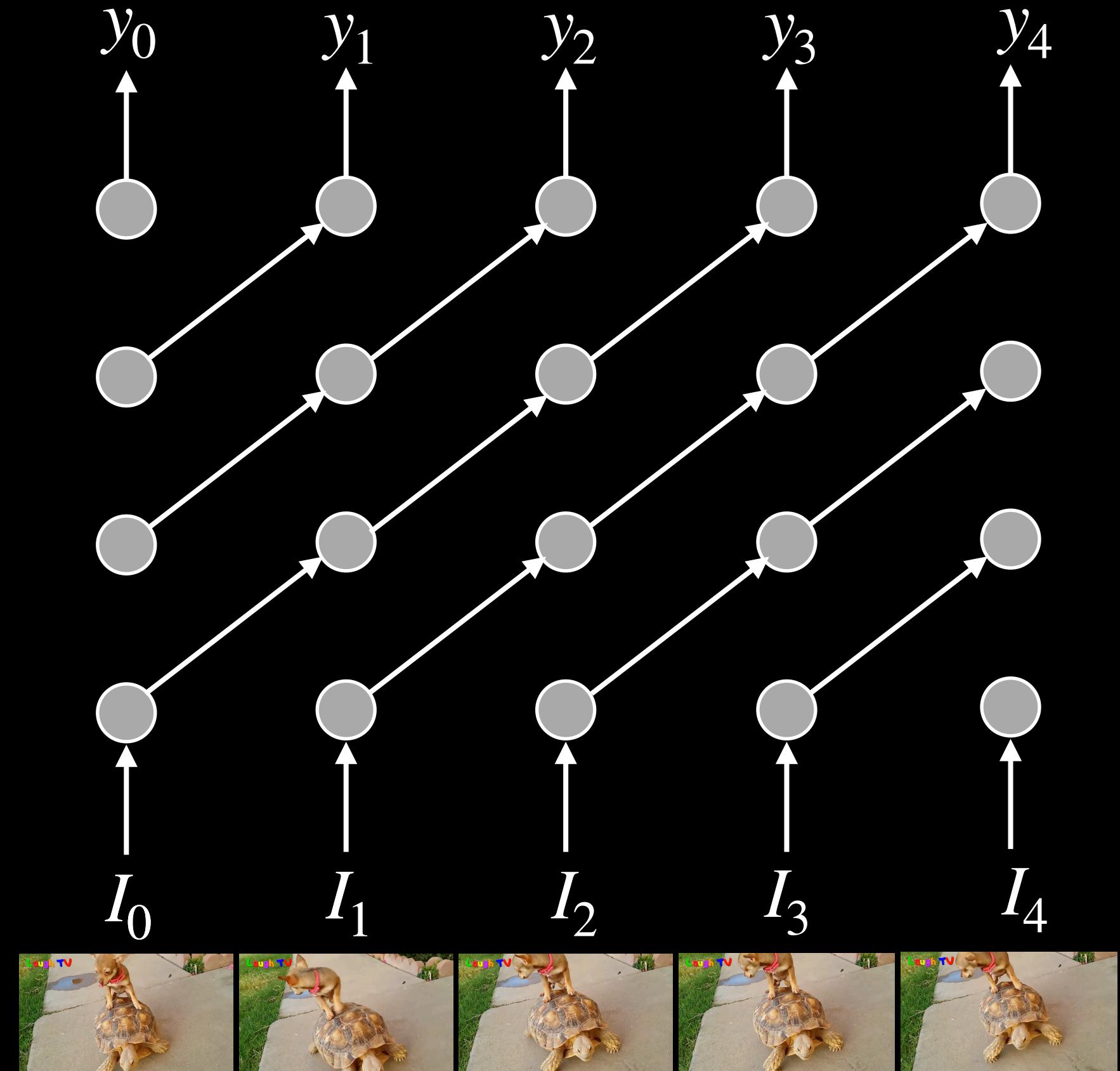
Con:
more difficult task
reduced accuracy



round-robin over layers

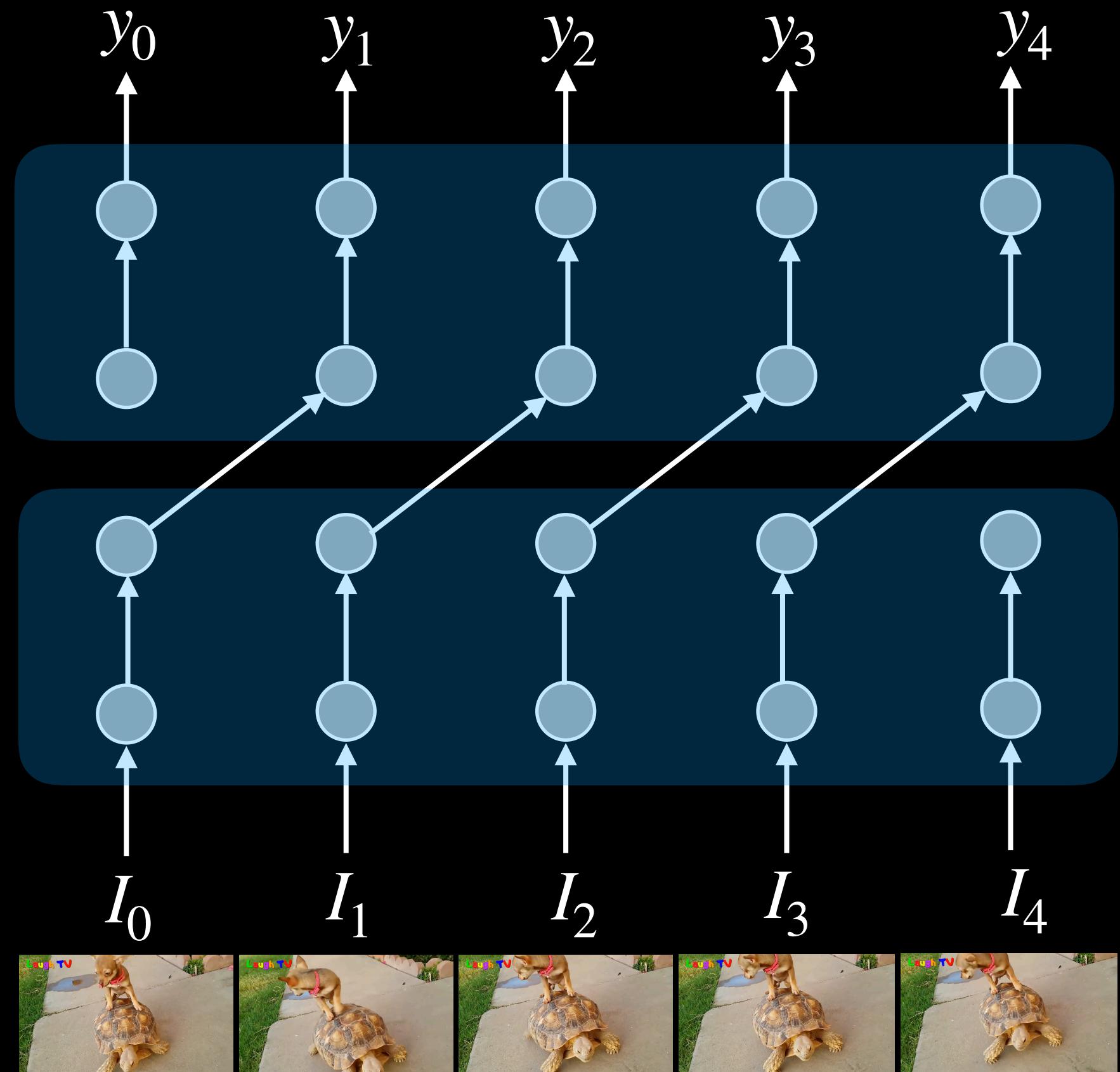
1. Partial pipelining

Fully-parallel model



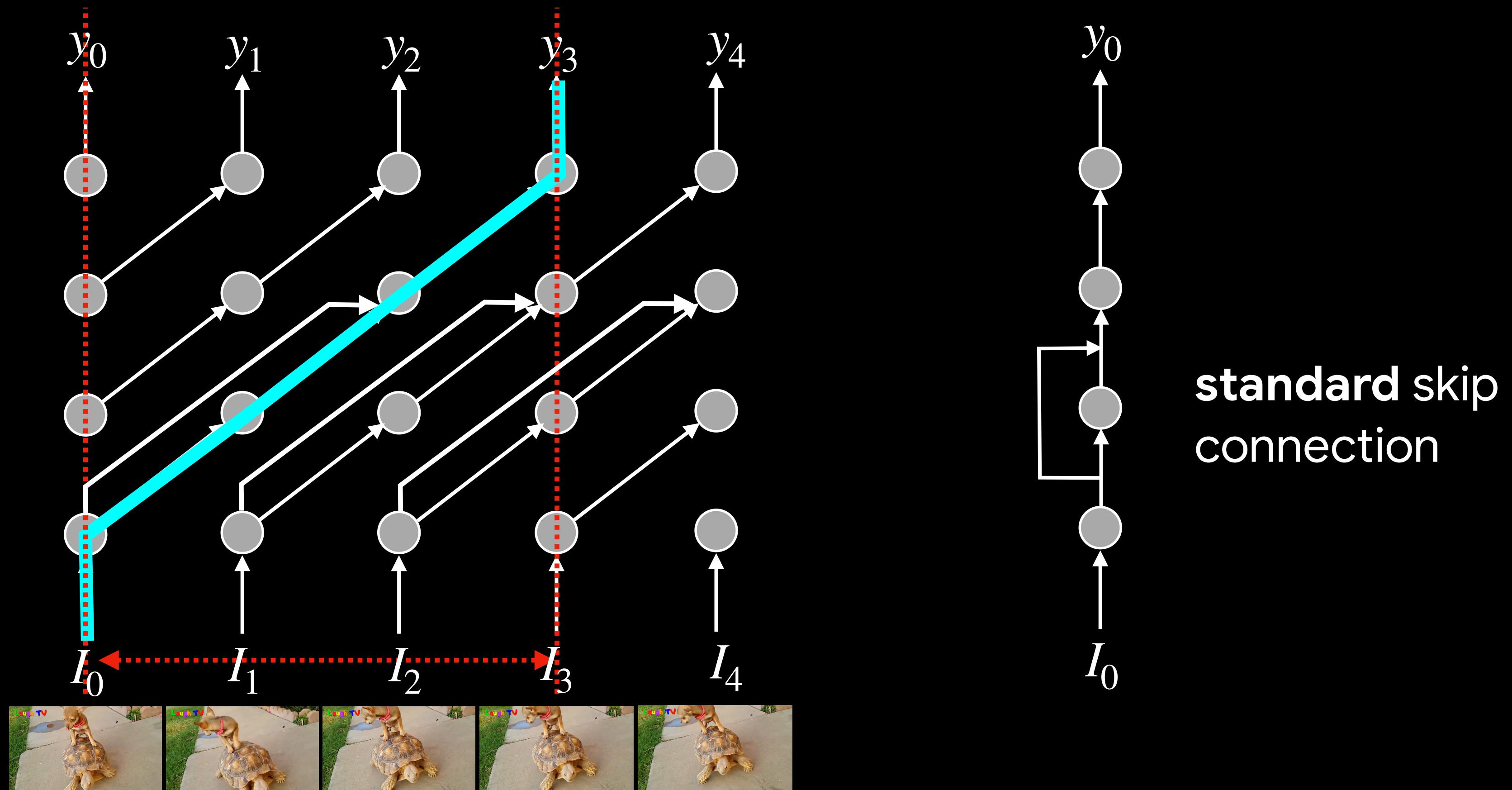
1. Partial pipelining

2 parallel subnetworks



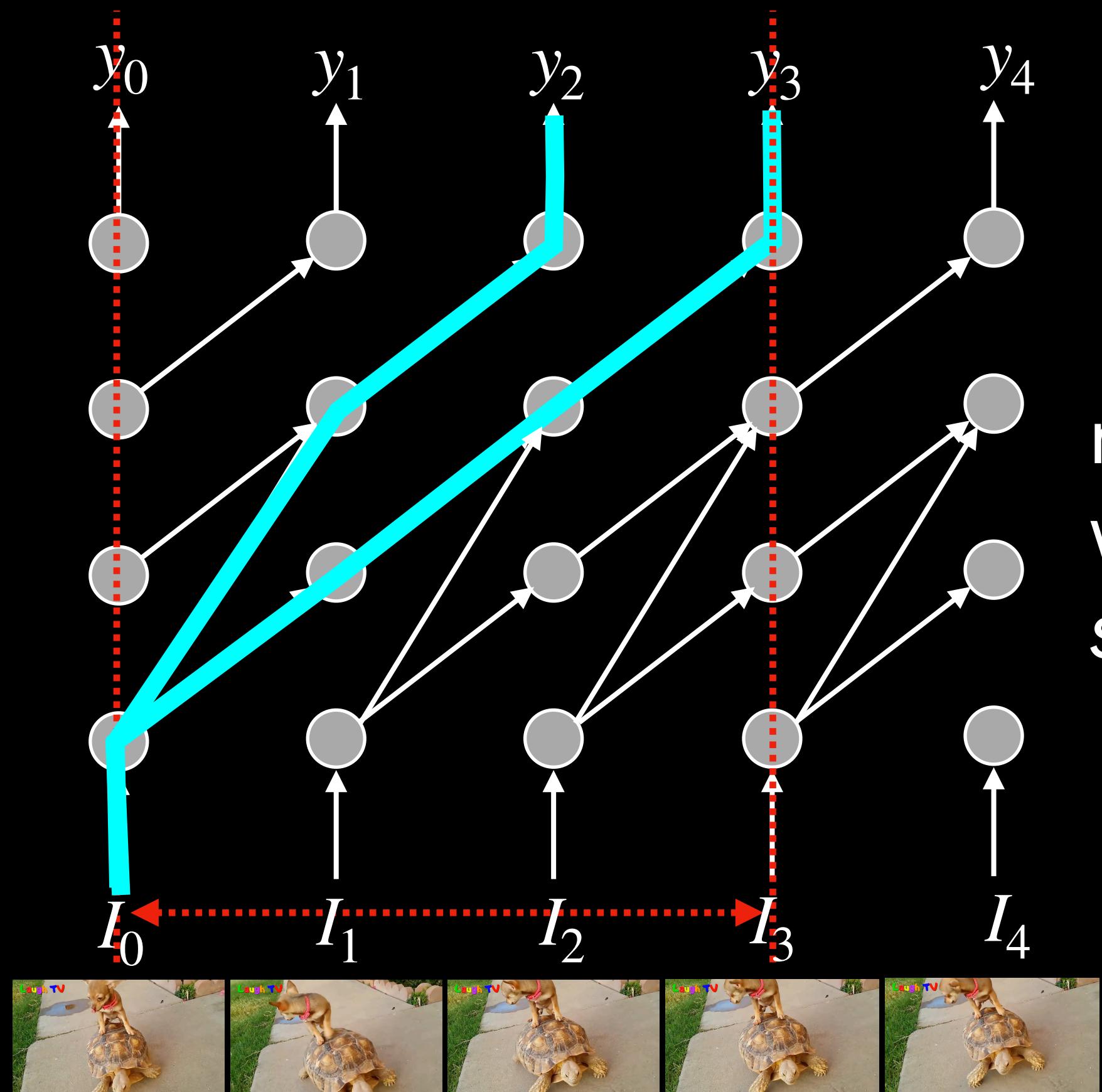
**Pro: trade-off
latency vs accuracy**

2. Skip connections



2. Skip connections

Temporal
skip connections



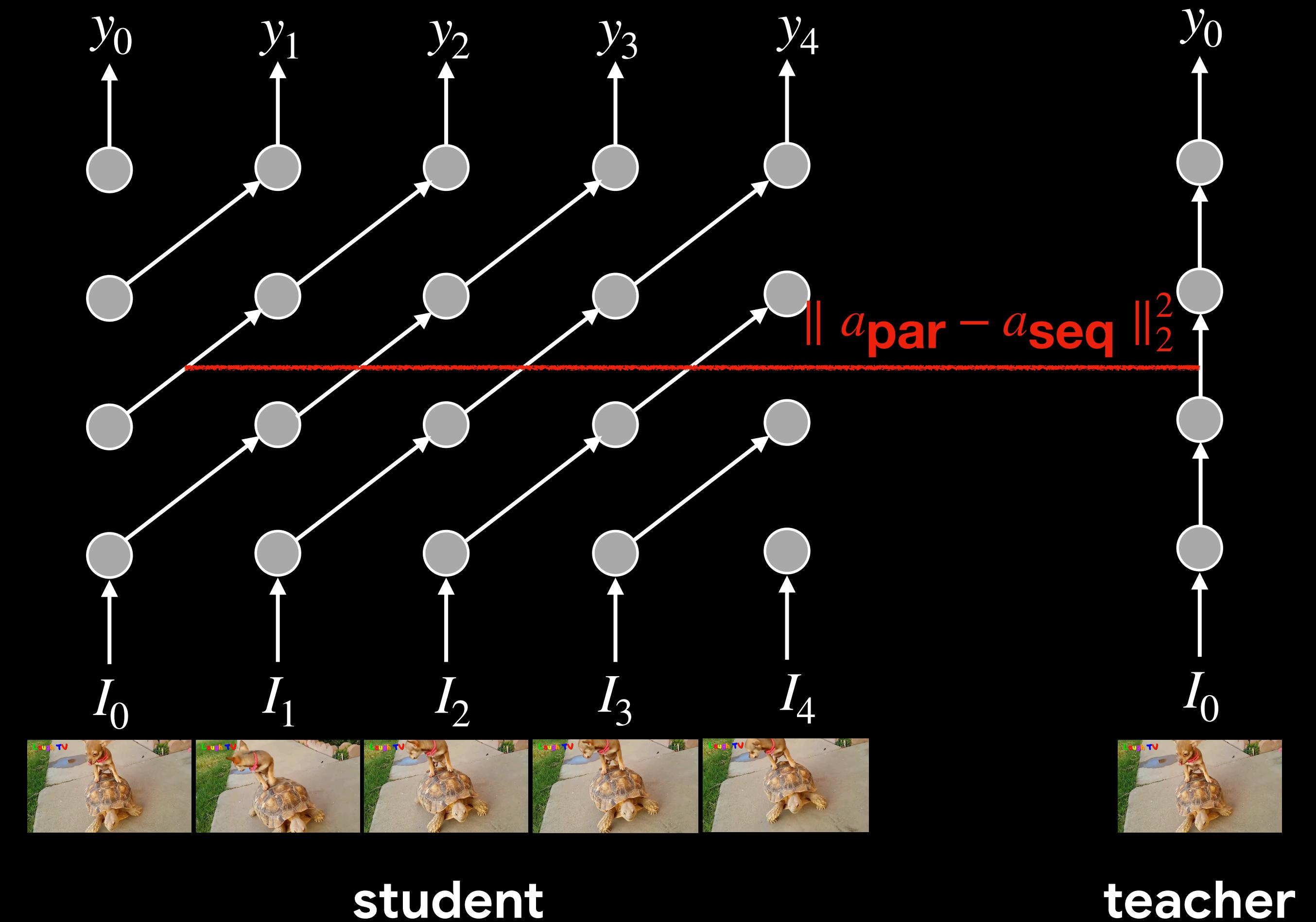
merge fresh shallow features
with old deep features
slowness principle [1]

[1] Wiskott and Sejnowski, Slow feature analysis: Unsupervised learning of invariances. Neural Computation 14(4) (2002)

3. Distillation

teacher - sequential model,
accurate but slow
student - pipelined model,
fast but inaccurate

Additional loss term:
intermediate activations
should match

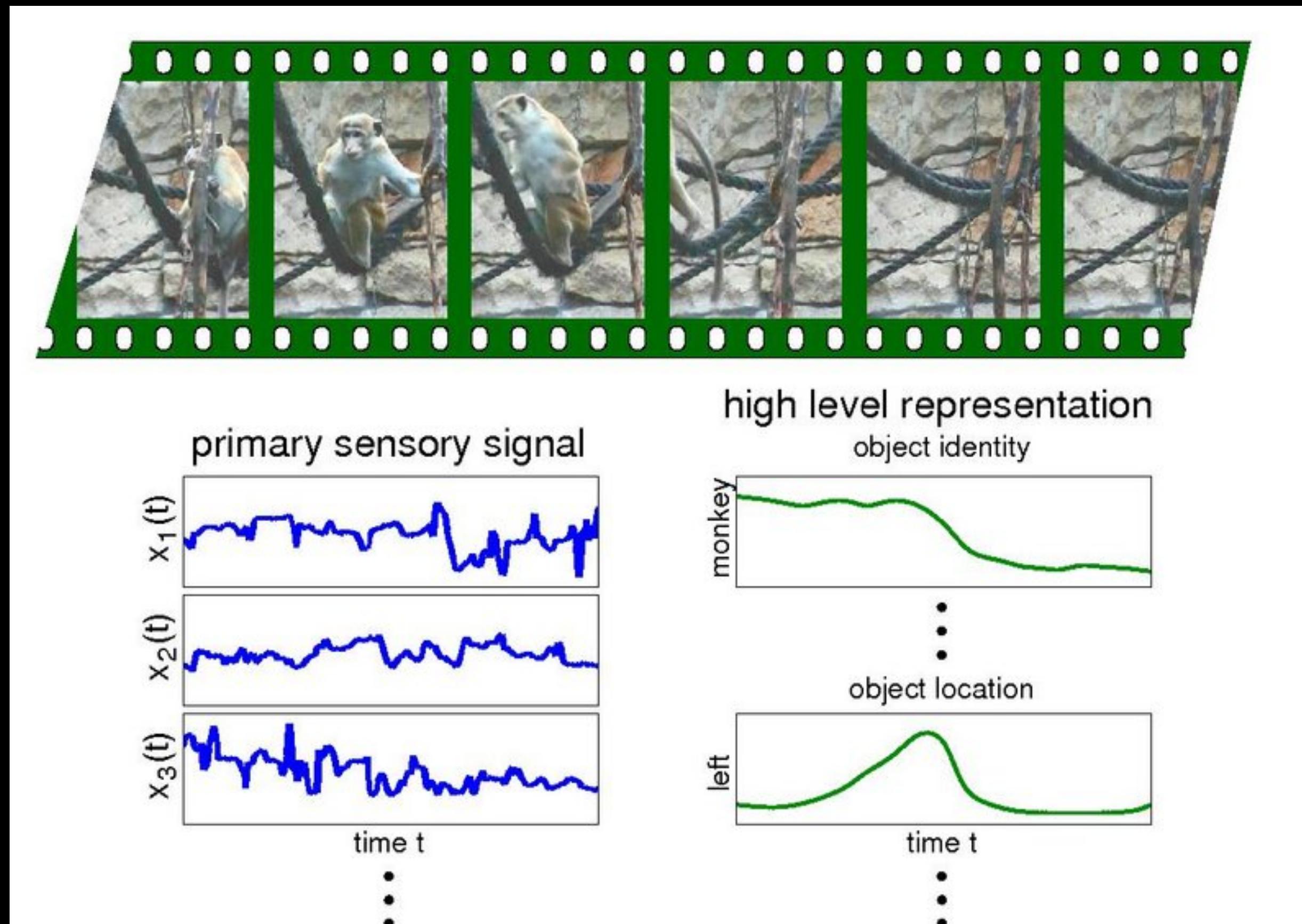


Increase throughput

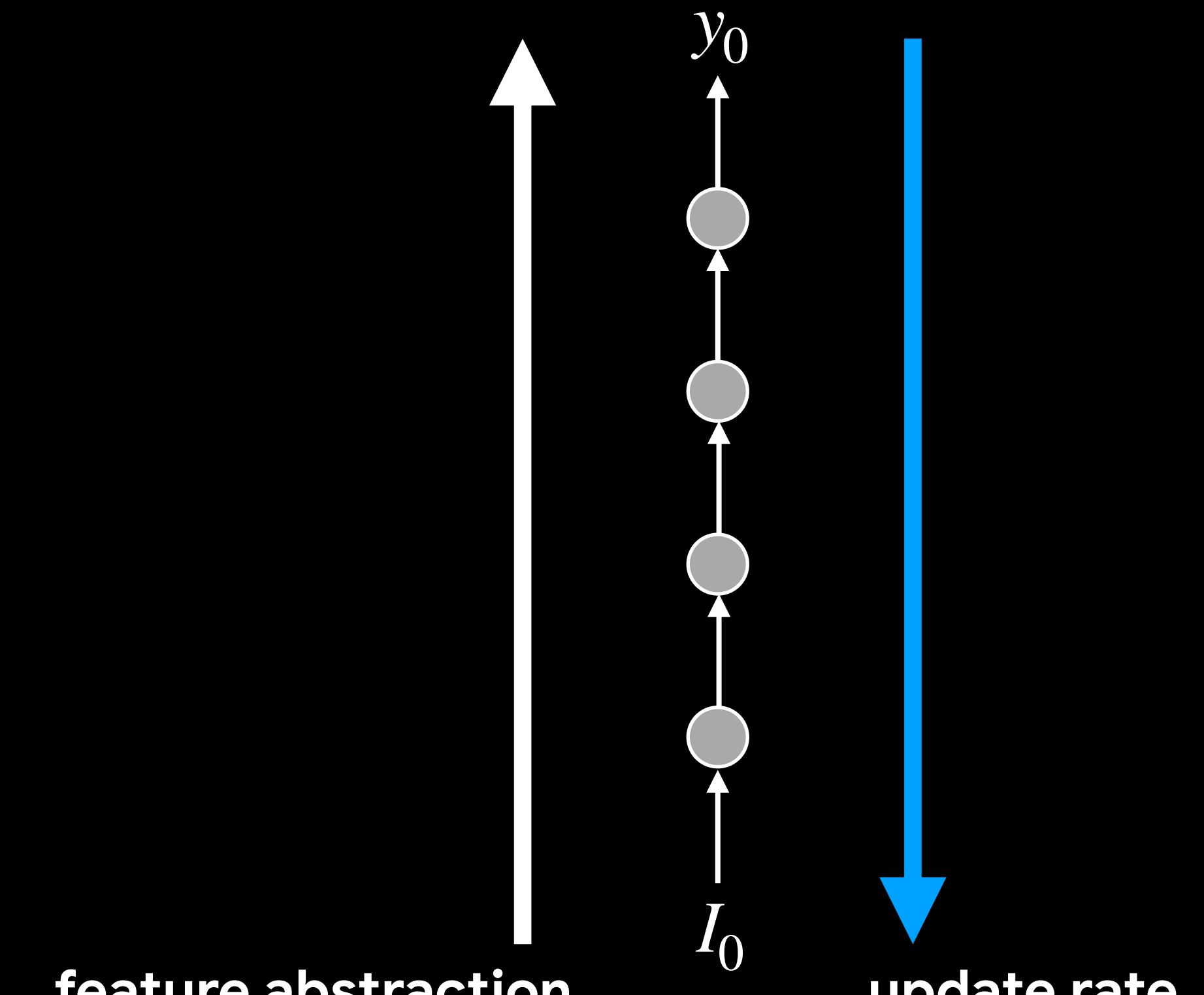
Reduce latency

Reduce clock cycles

Slowness principle



http://www.scholarpedia.org/article/Slow_feature_analysis



Exponentially reduced update rates (Wavenet style) \Leftrightarrow causal conv3D temporal stride > 1

Summary: predictive depth-parallelism

- 1. pipelined ops
 - 2. temporal skip connections
 - 3. distillation **change in train loss**
 - 4. multi-rate clocks **change in graph connectivity**
- } **changes in temporal connectivity**

Latency not dependent on depth

Experiments

Investigate the effect of pipelining on:

- accuracy
- speedup (throughput)

Tasks

Slow task

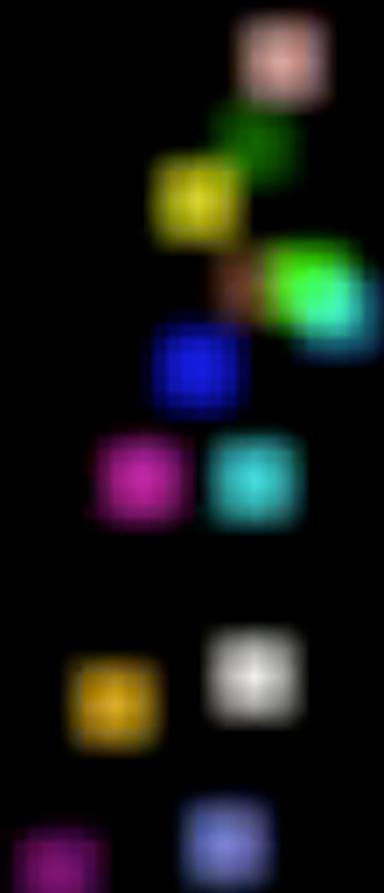
action recognition



video label “play basketball”

Fast task

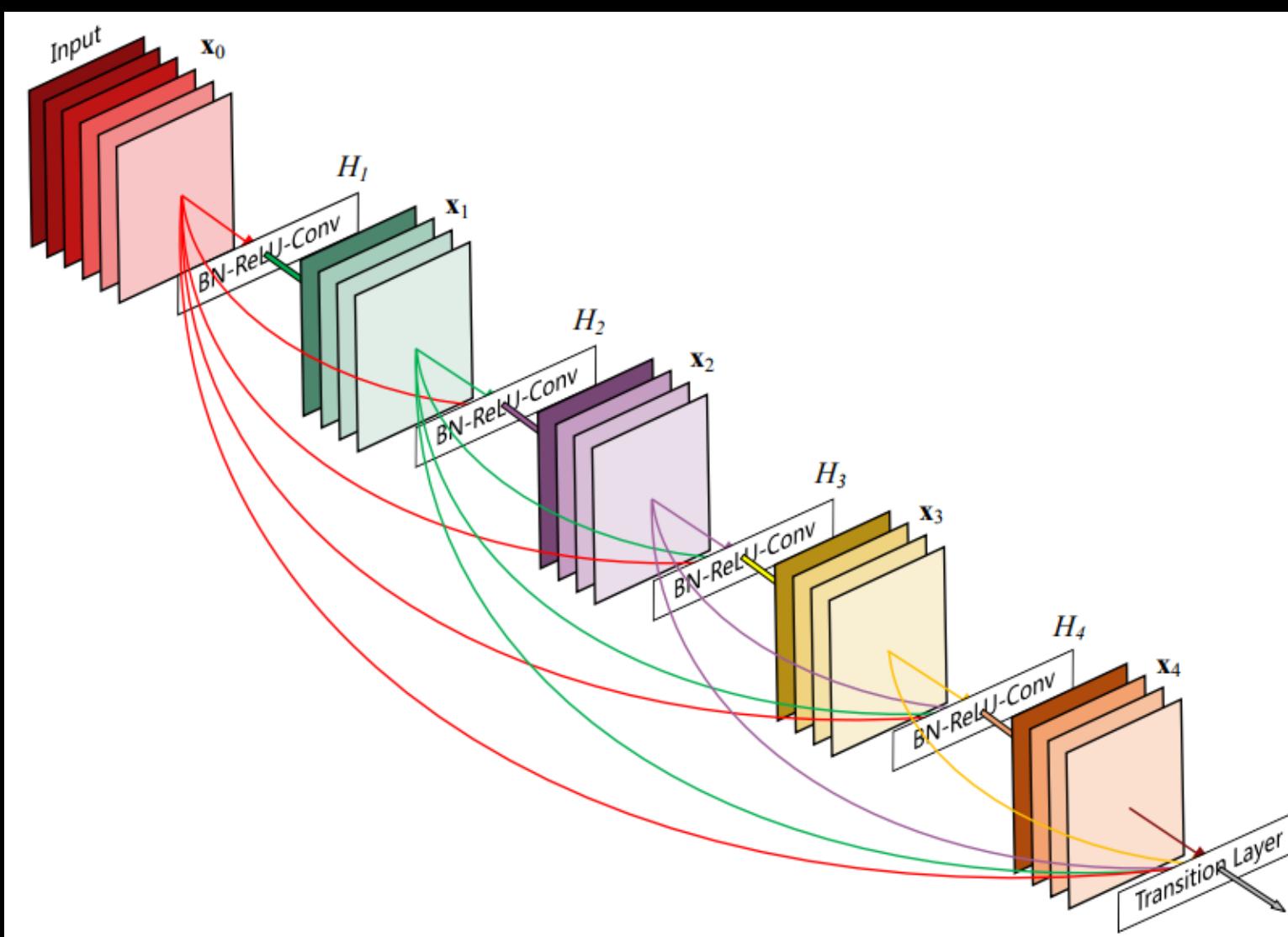
human keypoint localisation



per-frame labels: 13D joints heatmaps

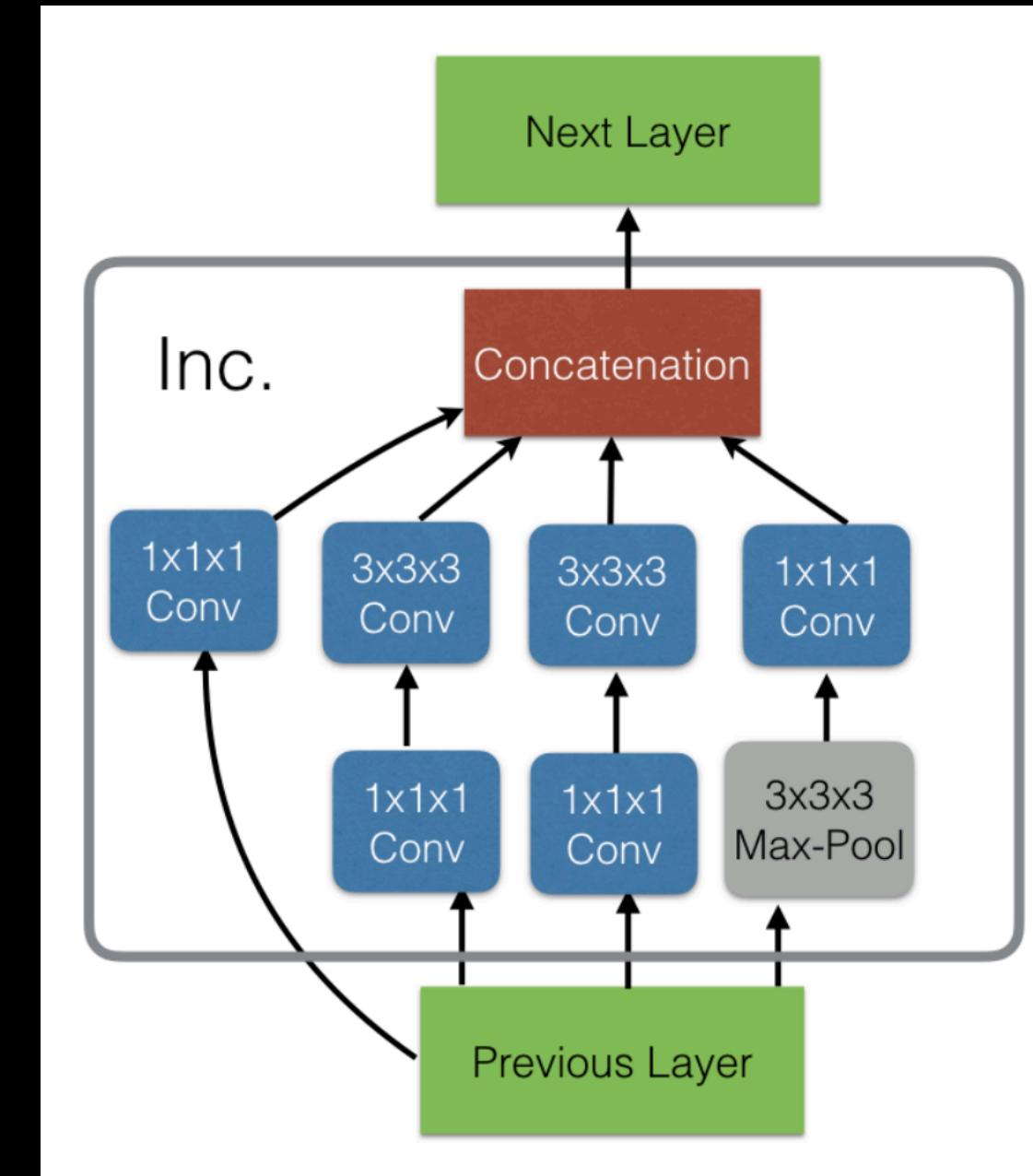
Models

Par-DenseNet



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

Par-Inception



Quo Vadis, action recognition? Carreira and Zisserman, CVPR2017

Dataset

Mini-Kinetics

200 action classes

80k training videos

5k validation

automatically generated
“ground truth” keypoints [1]

archery



country line dancing



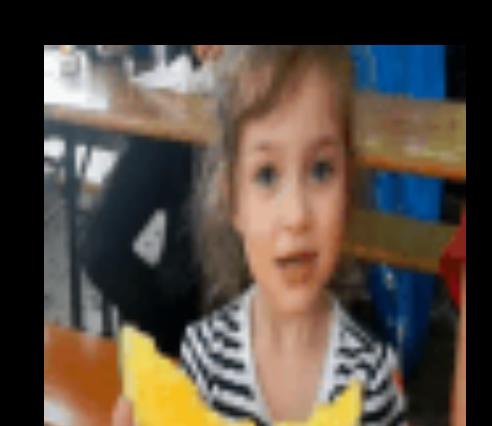
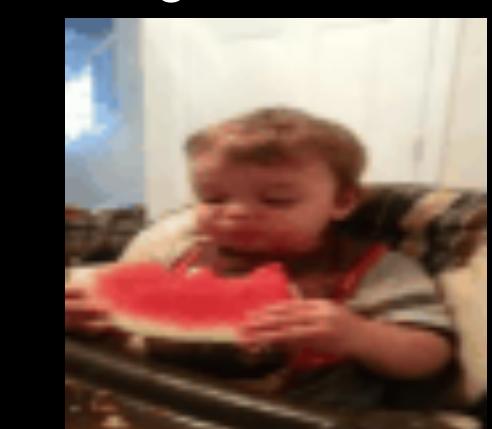
riding or walking with horse



playing violin



eating watermelon

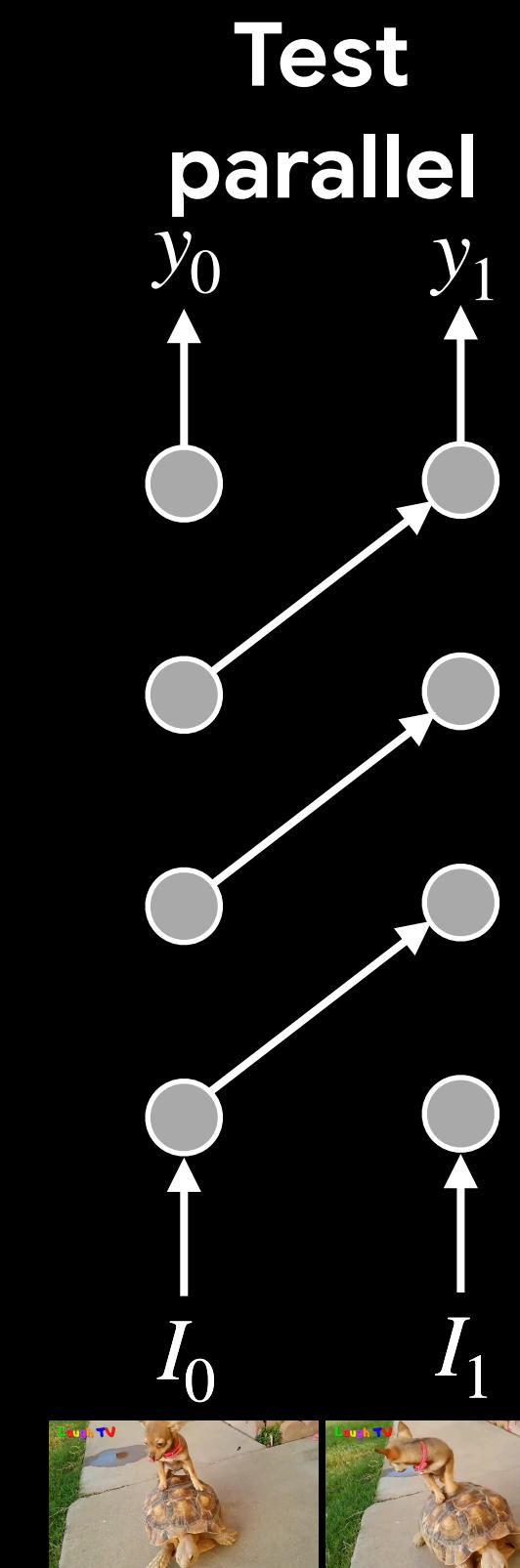
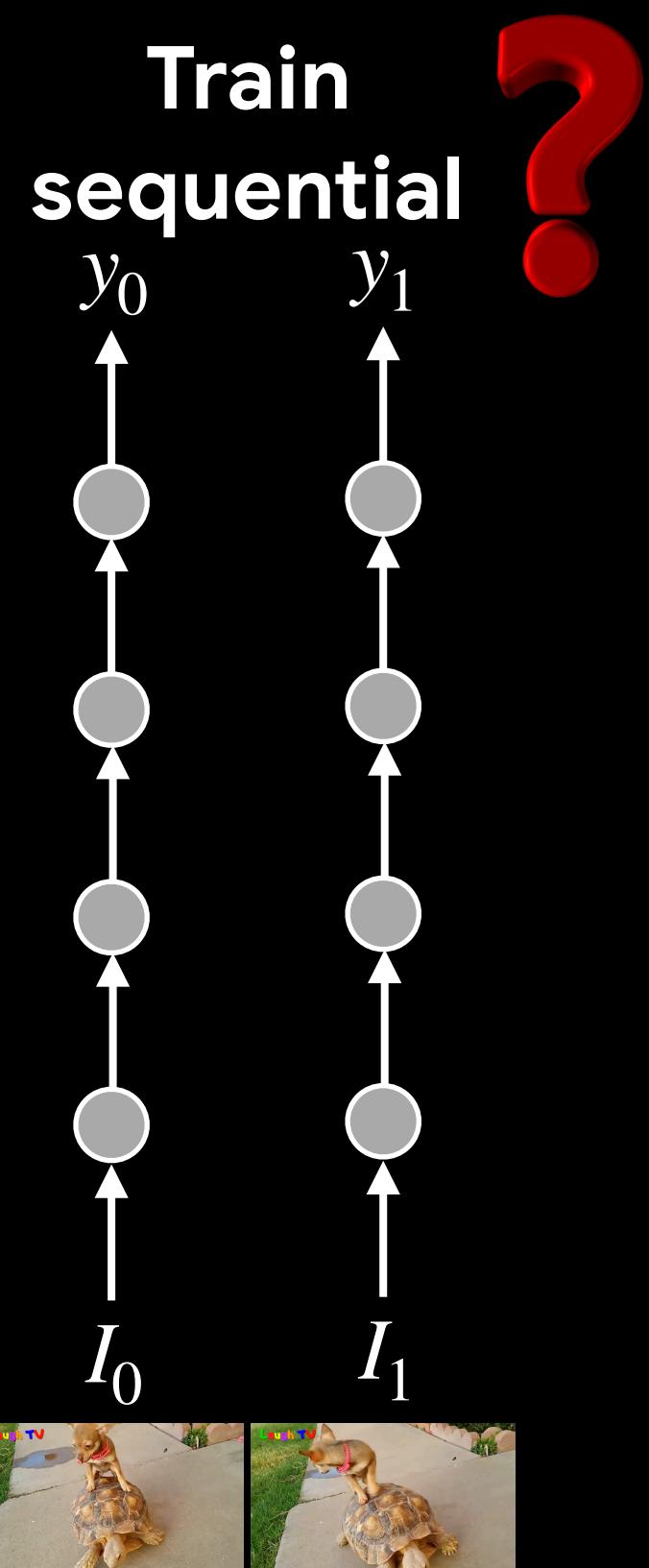
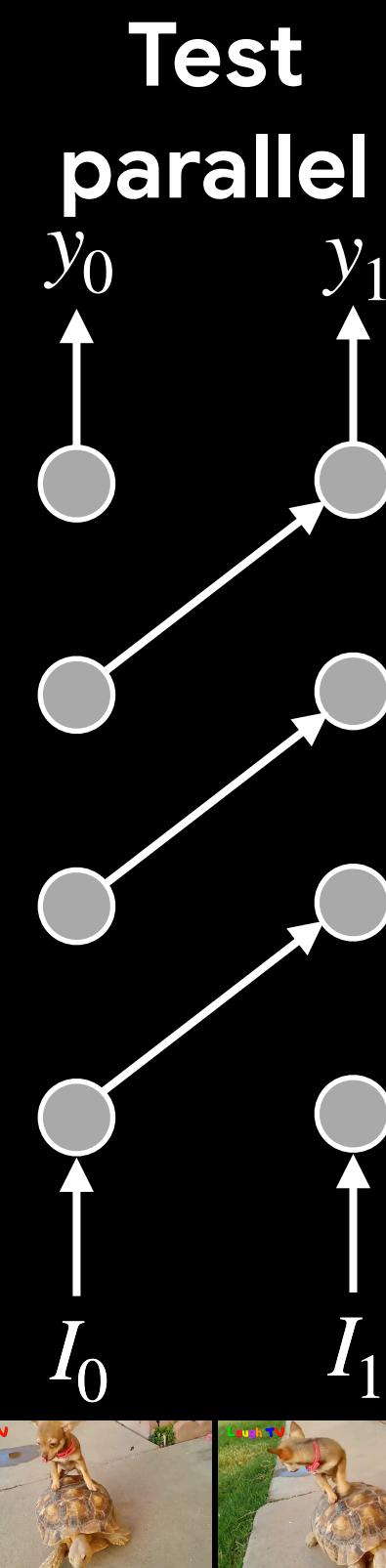
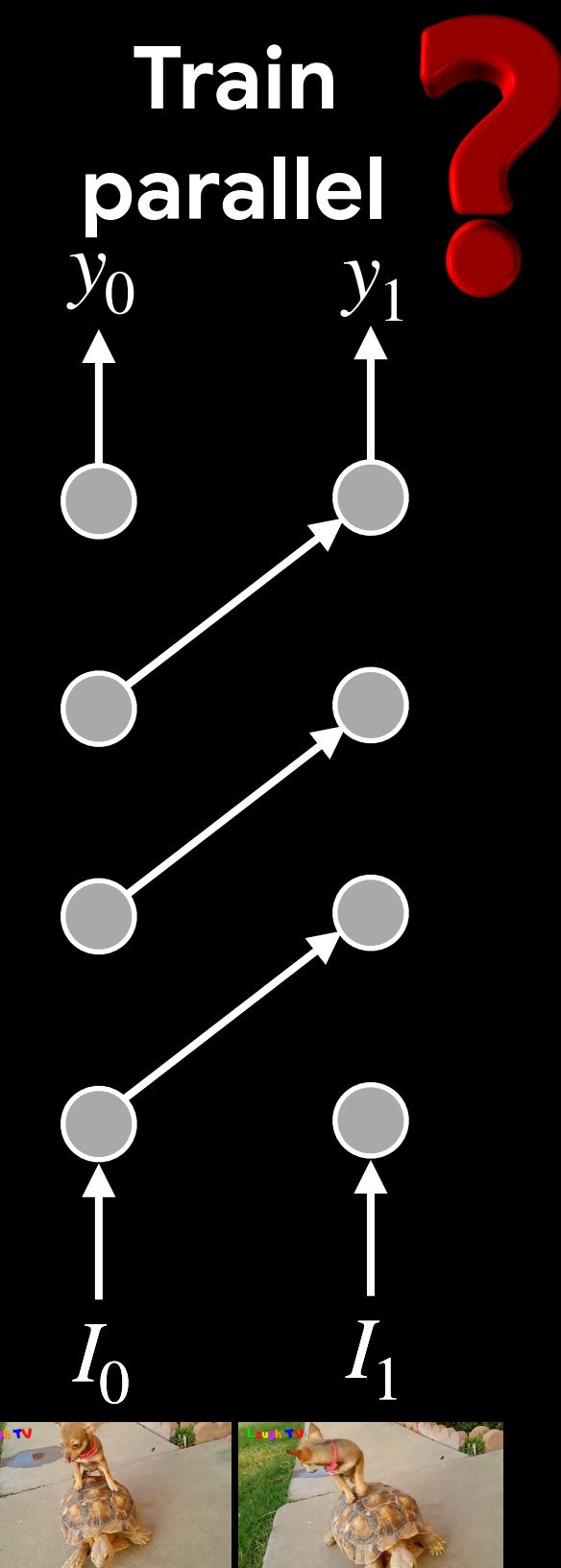


[1] Papandreou et al., Towards Accurate Multi-person Pose Estimation in the Wild, CVPR2017

Results: Human keypoint localisation

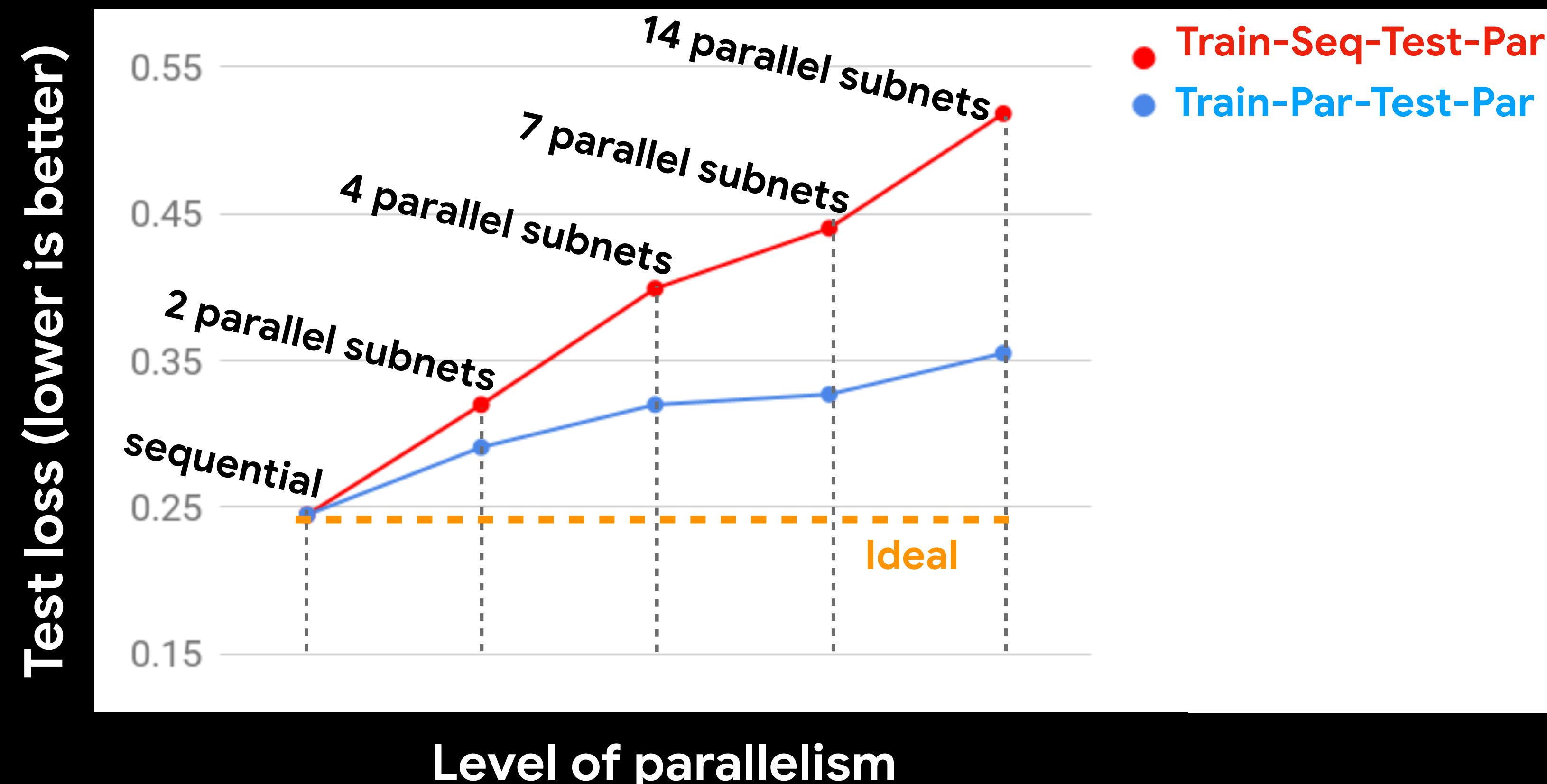
Results: Predictive parallelism

Q: Can we just parallelise at test time?



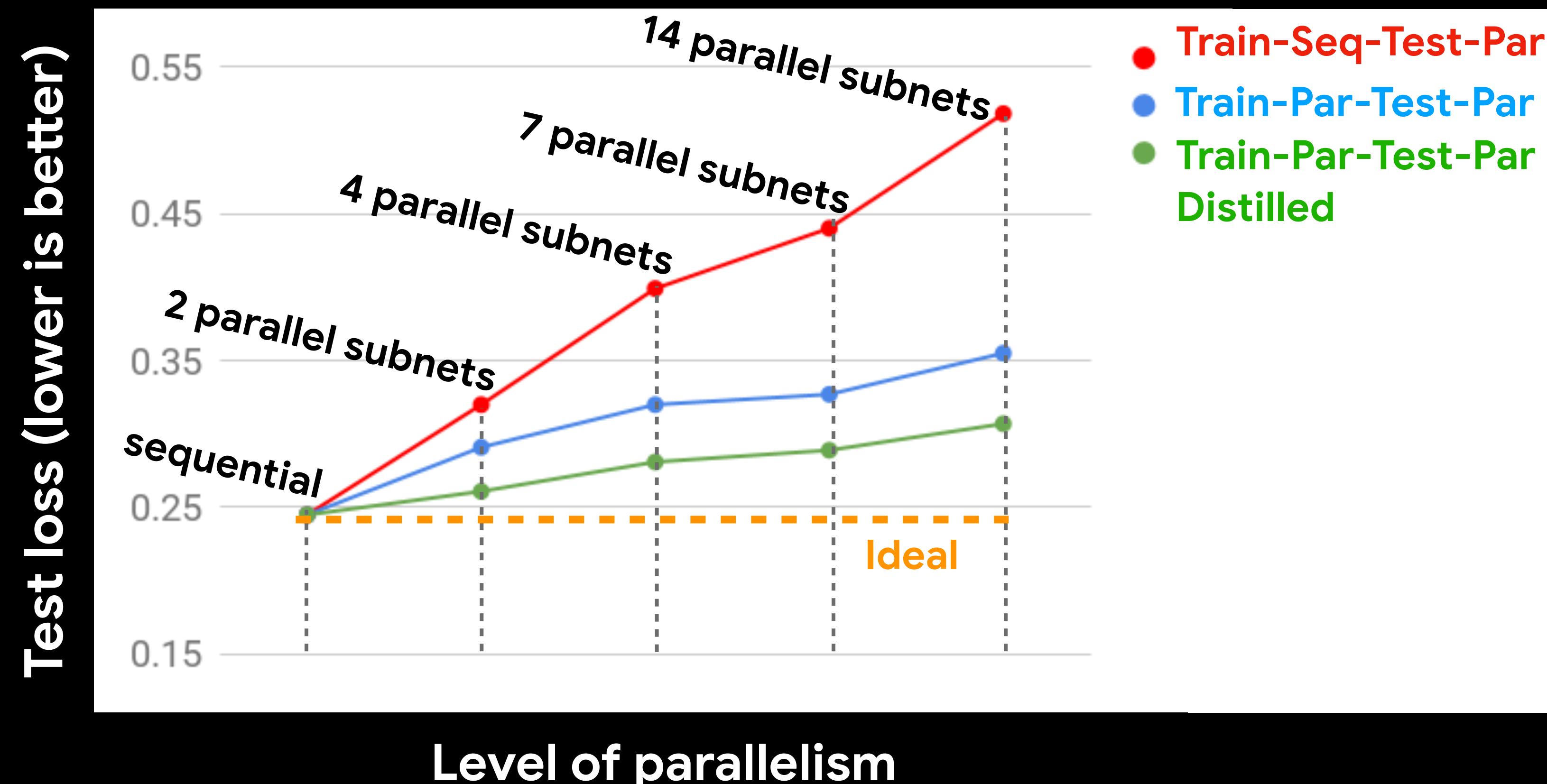
Results: Predictive parallelism

Q: Can we just parallelise at test time? A: **NO**

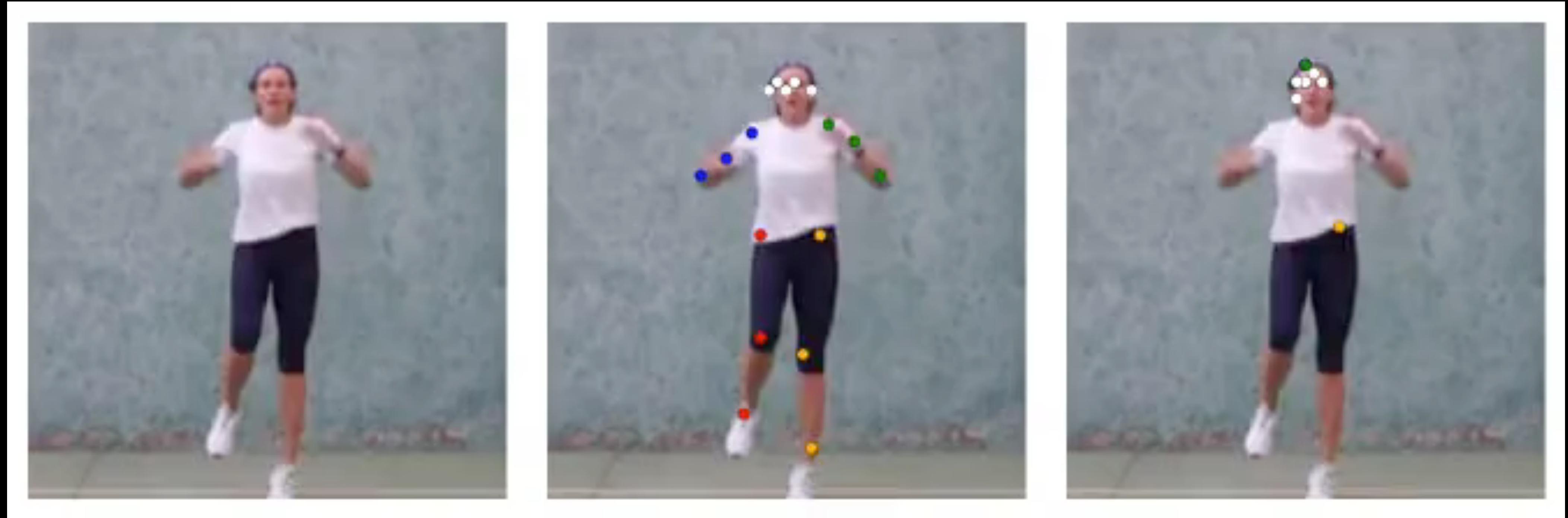


Results: Predictive parallelism

Q: Can we just parallelise at test time? A: **NO**



Results: Human keypoint localisation

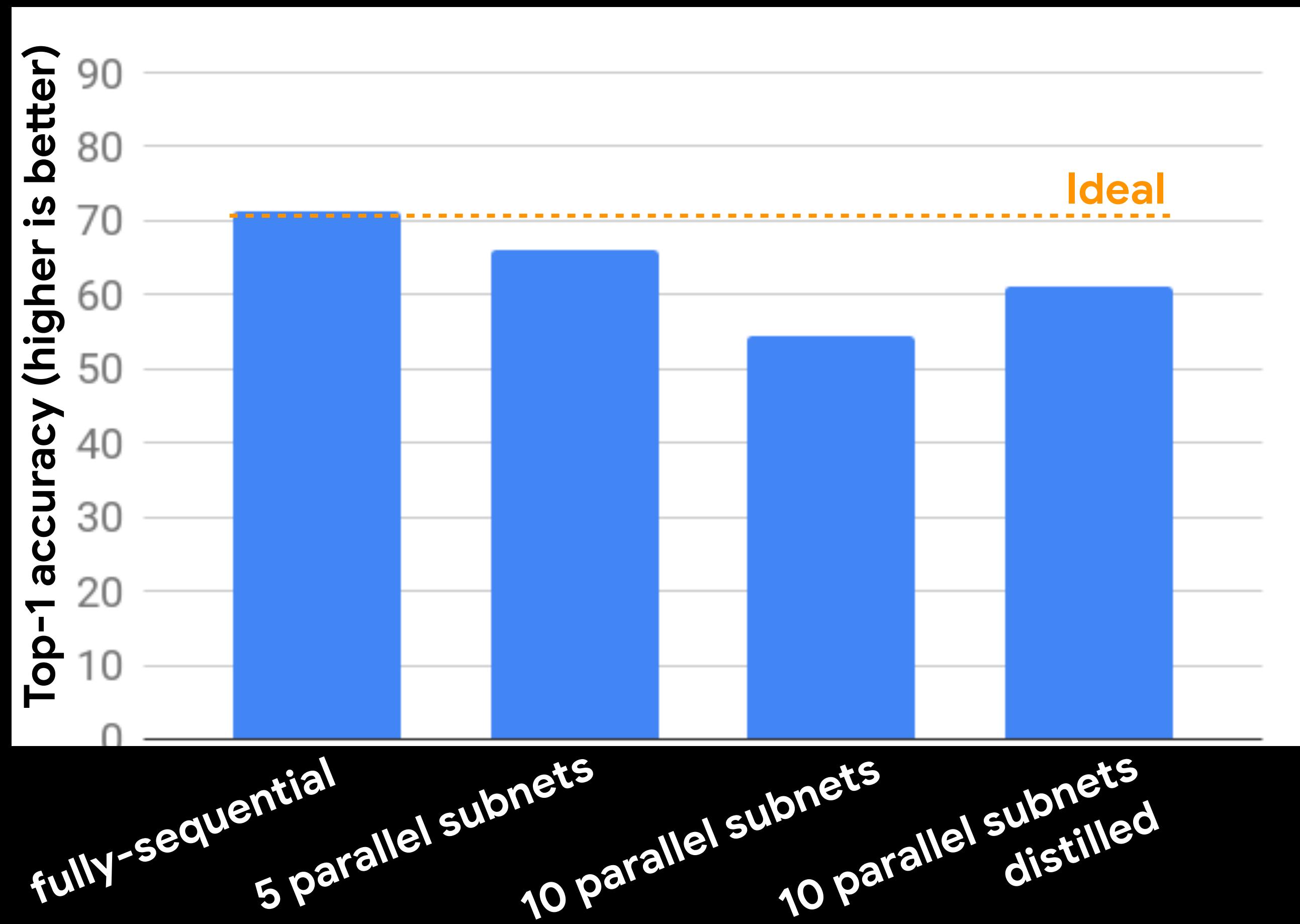


Input

Fully-sequential model
running offline

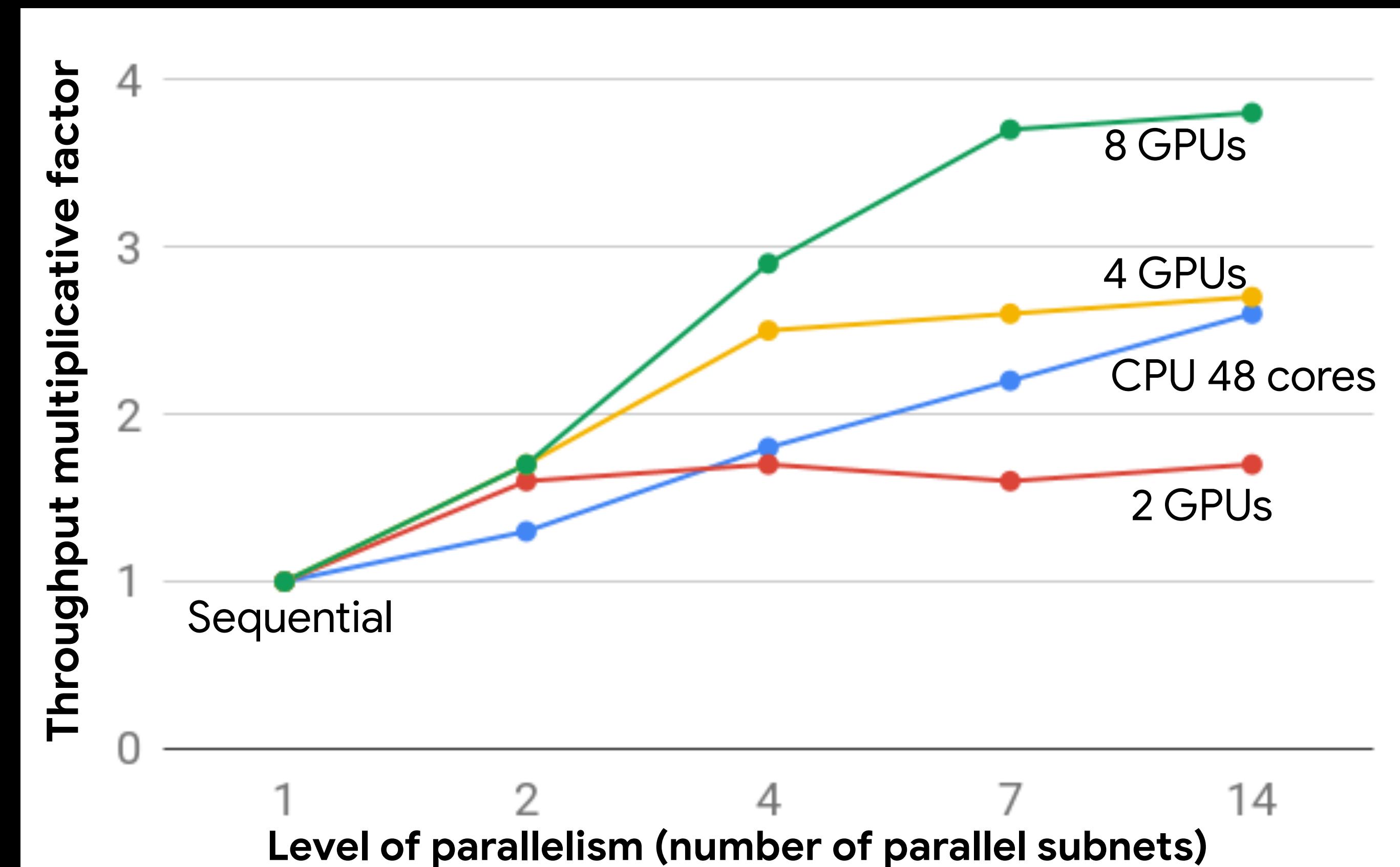
Fully-parallel model
running online

Results: Action recognition



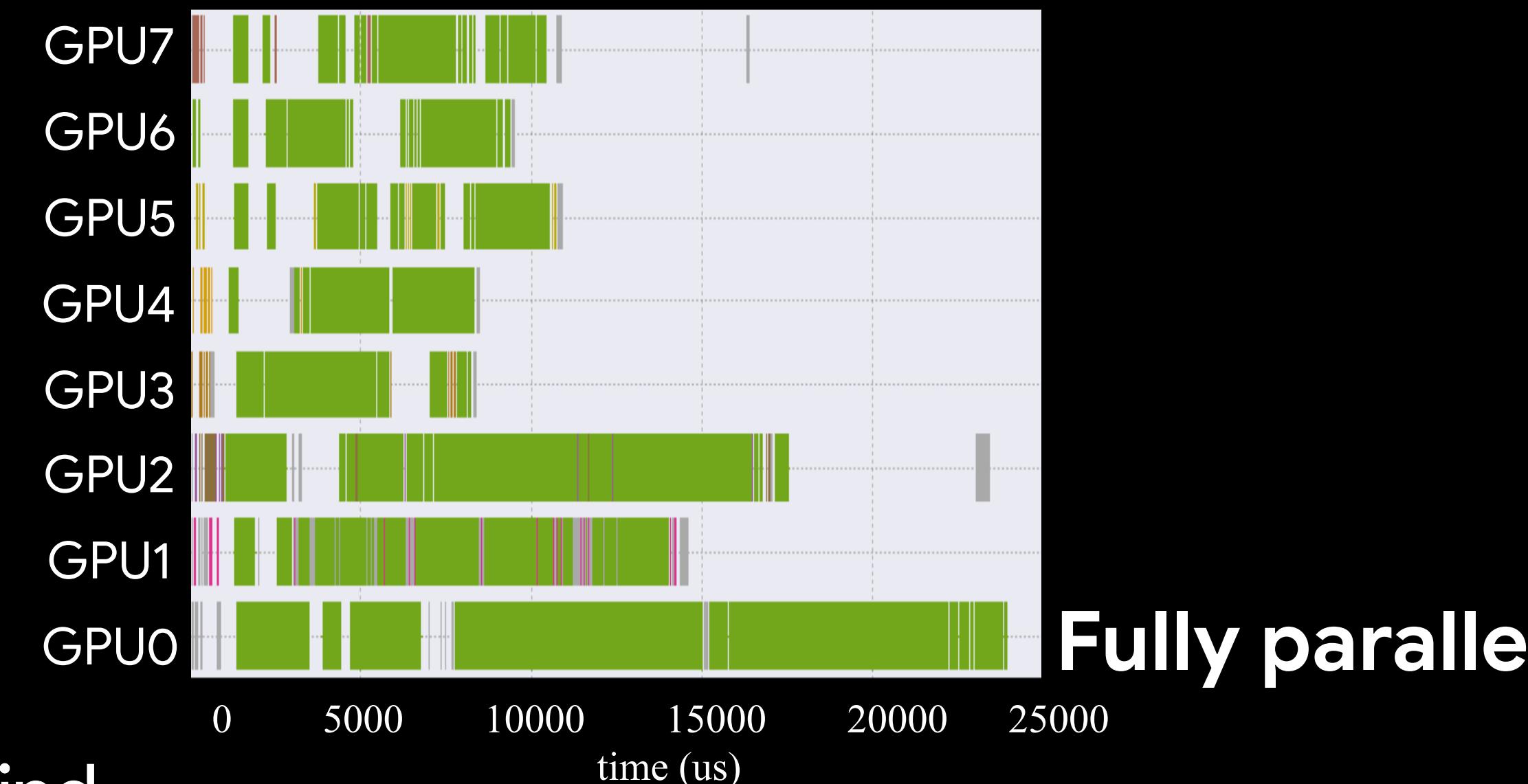
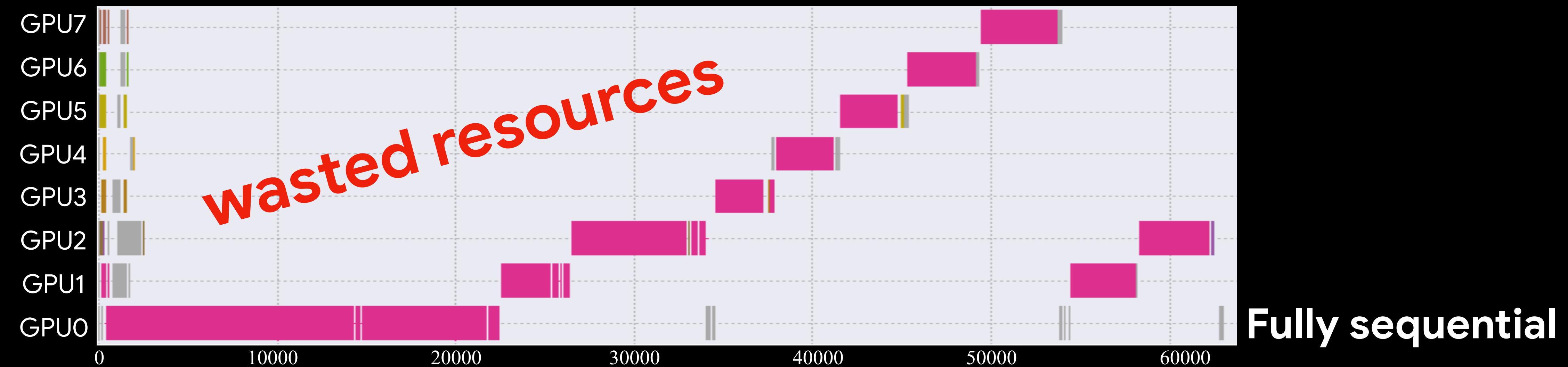
Speedup

Speedup - CPU or GPU box



Physically colocated hardware

GPU utilisation



Conclusion

Promising new avenue for scalable video networks design

Important speed-up gain: 4x - 8x

Interesting aspect: latency reduction shapes representations

More investigation of architecture space needed (evolution strategies, Hebbian learning)