

Lecture 13: Video and temporal models

Announcements

- PS6 due next Weds.
- Project proposal information out
 - Rolling deadline, due Nov. 13
- Discussion section: project office hours + GANs

Project proposal

- Due November 13th
- Rolling deadline
- What we want: 1 page summary of what you'd like to do.
- Not worth much of total grade. Graded pass/fail. We just want to know whether it's an acceptable project.

Project proposal

2 Project ideas

To help you think of projects, we've provided a few ideas below. Please note that these projects only cover a very small portion of the possible things you can do — most involve reimplementing and extending a paper. We encourage you to propose your own, creative project ideas, and to use these as a starting point! We may also add new project ideas to this list in the coming weeks.

Applications of vision. Apply computer vision to a task that's important to you! We highly encourage this option if it appeals to you, since it's often the most fun option. For example, students in previous EECS 442 classes have applied computer vision to Settlers of Catan, measured the volume of liquid that could be held by a teacup, and analyzed the coffee coming out of an espresso machine. Often these projects will involve applying a few different computer vision models to a task, and analyzing the results.

Image synthesis. Implement a (small) version of a generative image model, such as VQ-GAN [1] or a diffusion model [2].

Extend an existing image synthesis model. Extend an existing image synthesis model, such as Stable Diffusion [3] in an interesting way (see here [4] for architecture details).

Video magnification. Implement a motion magnification algorithm, such as the method of Wadhwa et al. [5]. Try running it on your own videos, too.

Stereo. Implement a system that can estimate depth from a collection of photos using stereo. An easy-to-implement reference point is Goesele et al. [6].





kindergarten classroom



television

person



chair

“What will the girl do next?”





Simple video task: action recognition



Making latte art



Jaywalking



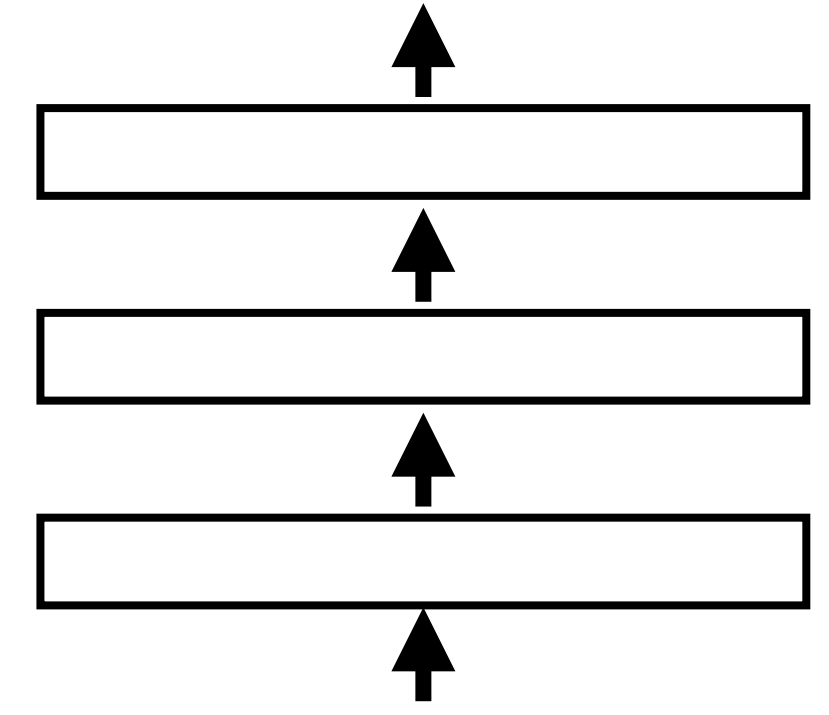
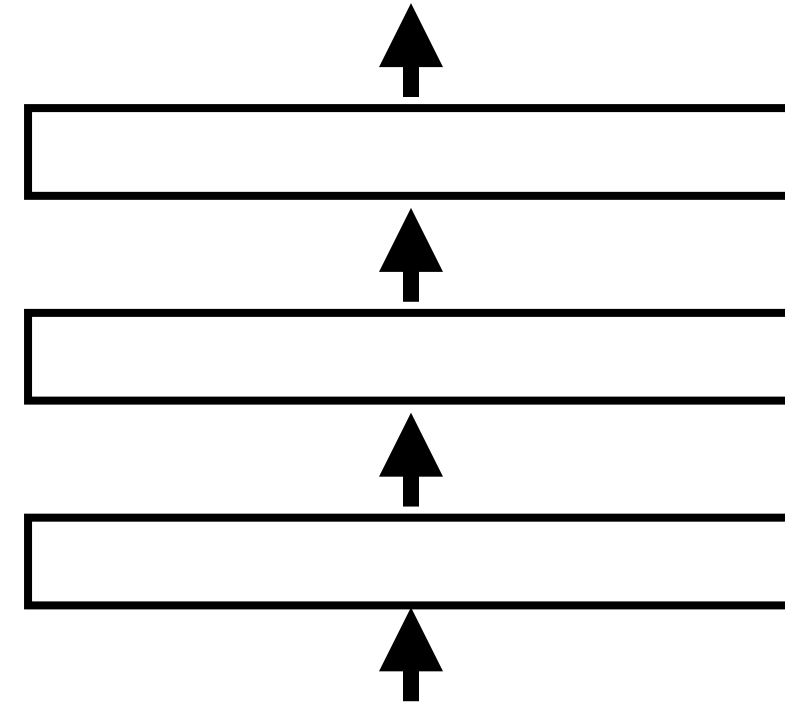
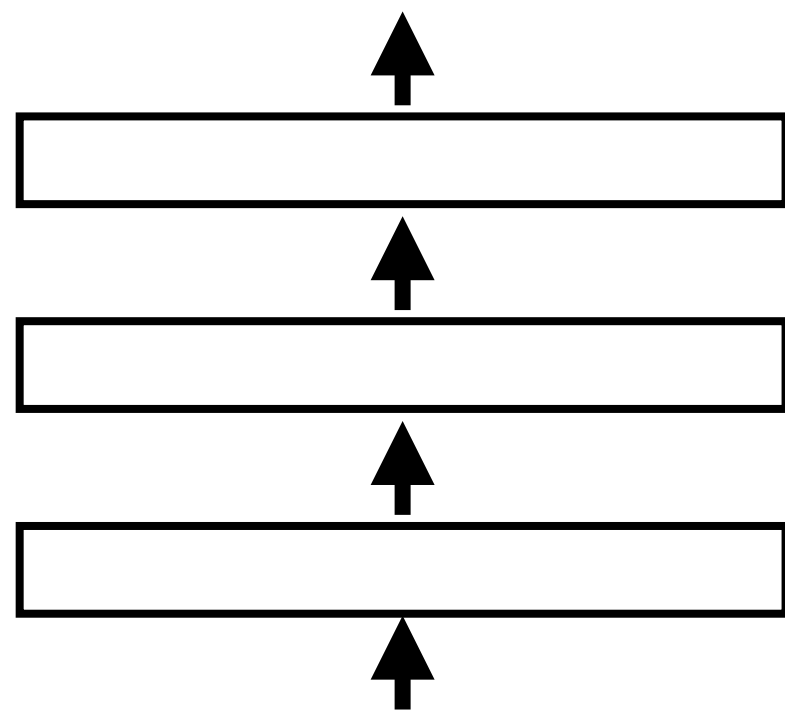
Grooming dog

Analogous to image recognition. Useful testbed for designing video models.

Examples from the Kinetics dataset [Carreira et al. 2017 - 2019]
700 human activity classes, 650K 10-sec videos

Simple model: averaging

$$\frac{1}{3} p(\text{making latte art} \mid I_1) + \frac{1}{3} p(\text{making latte art} \mid I_2) + \frac{1}{3} p(\text{making latte art} \mid I_3)$$



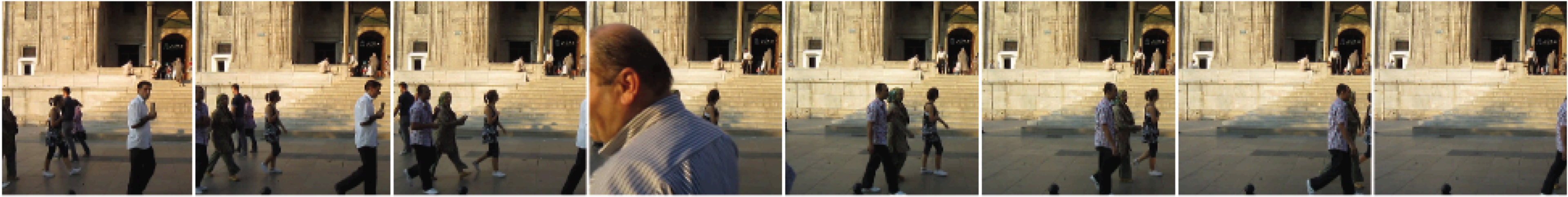
Can't analyze motion.

Temporal filtering

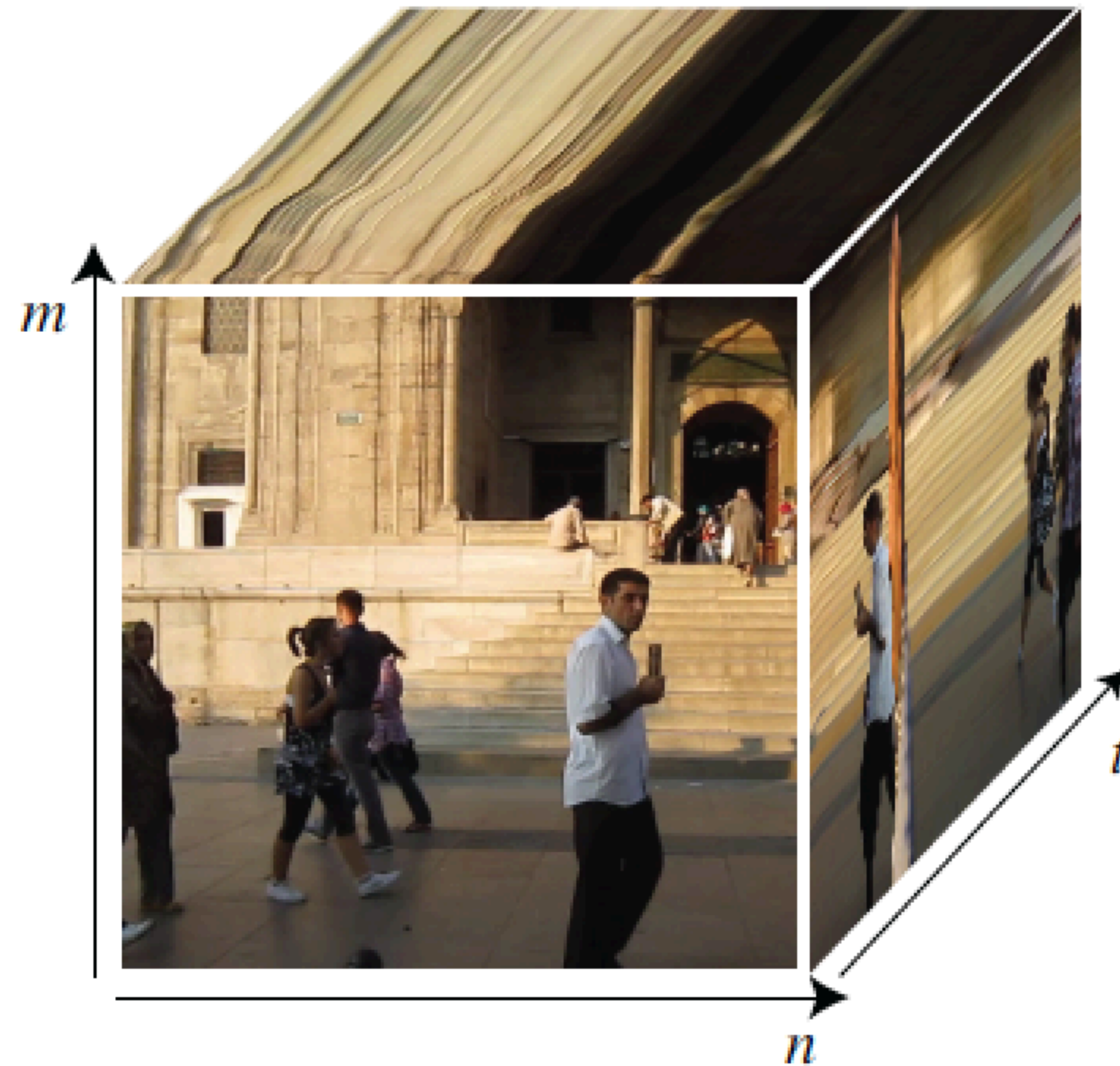
Temporal filtering



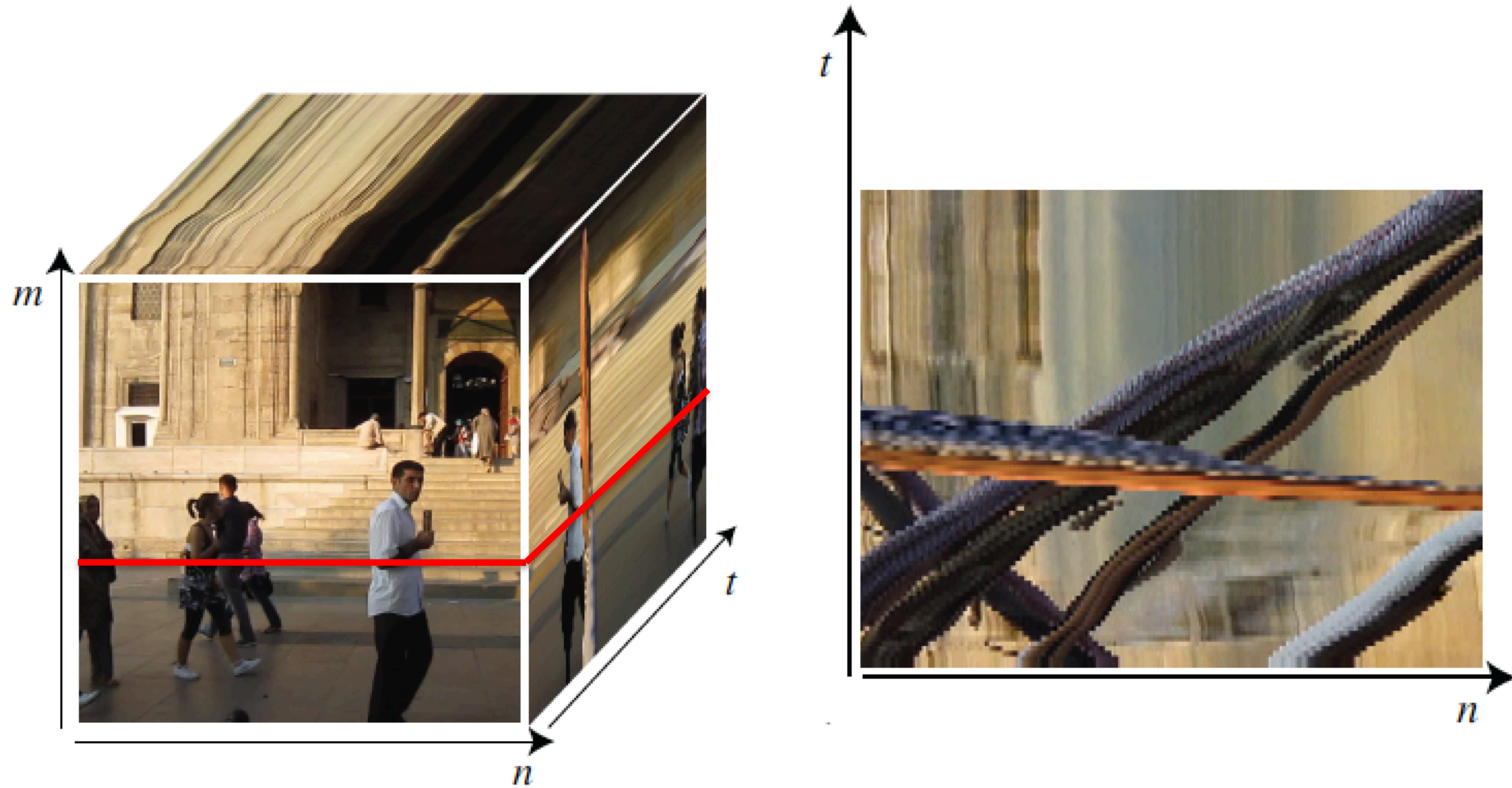
Videos



time →

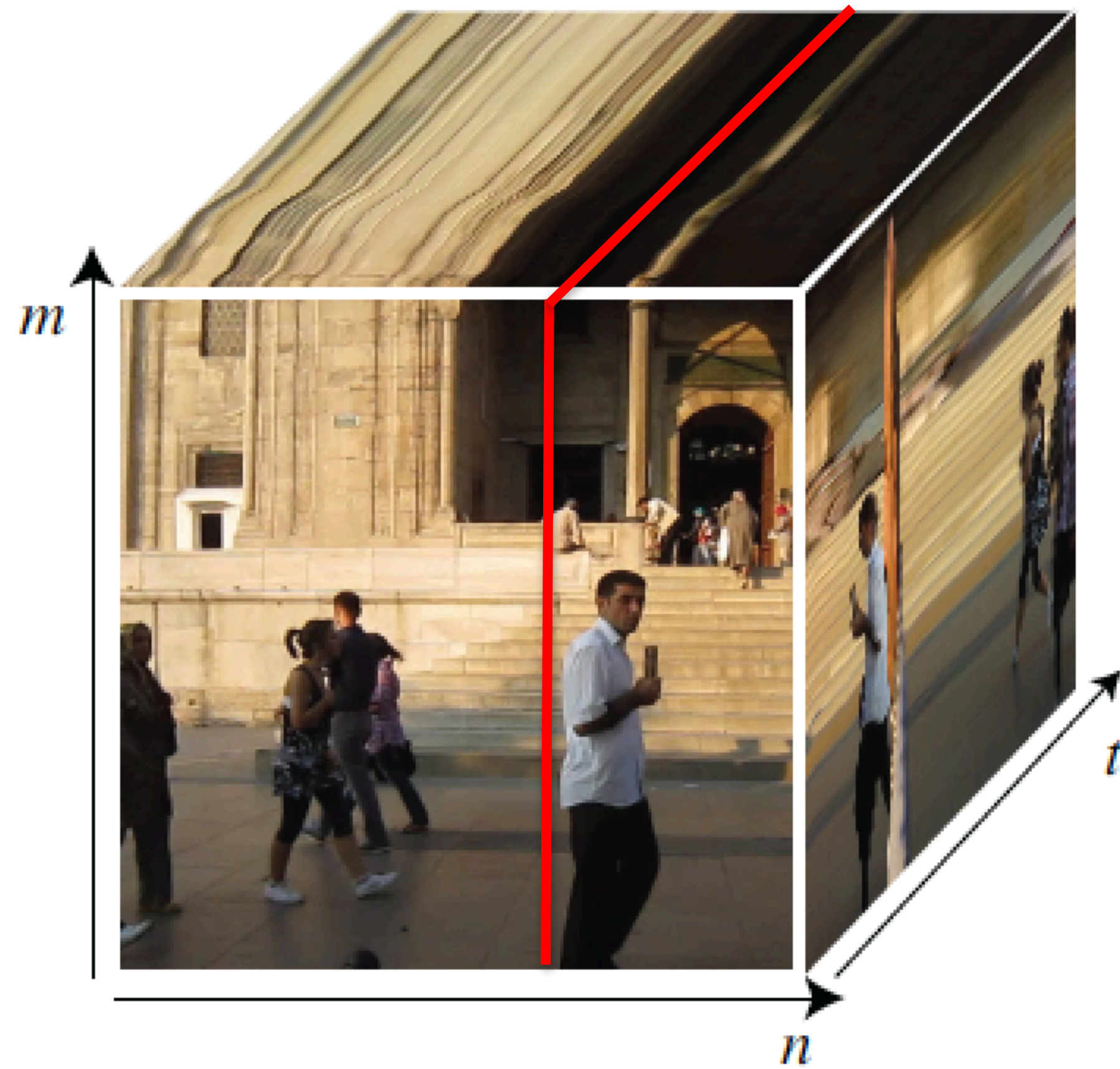


Videos



Cube size = 128x128x90

Videos



Cube size = 128x128x90



Examples of temporal filtering

Temporal median filtering



Background subtraction



-

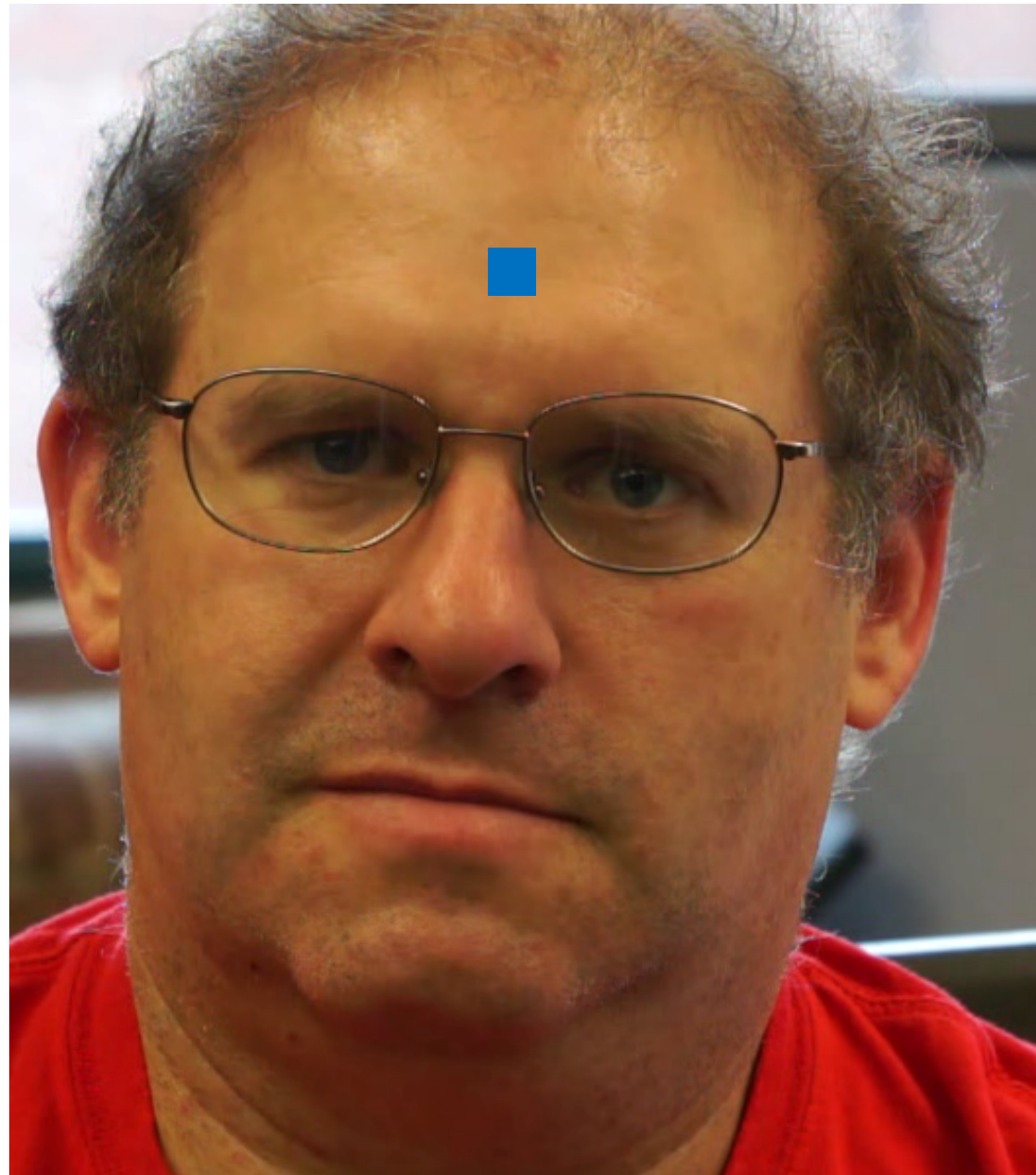


=

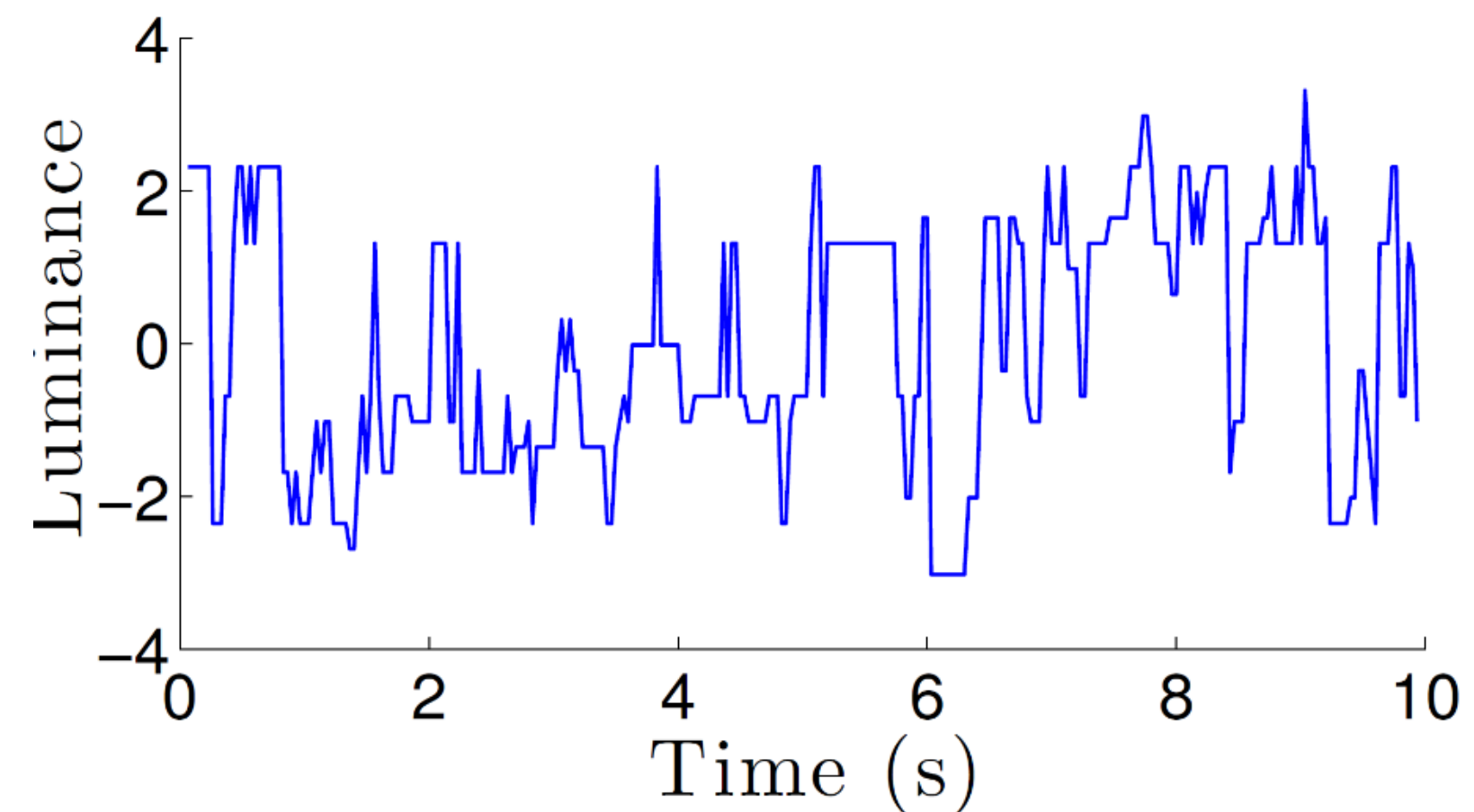


Finding subtle color variations

- The face gets slightly redder when blood flows

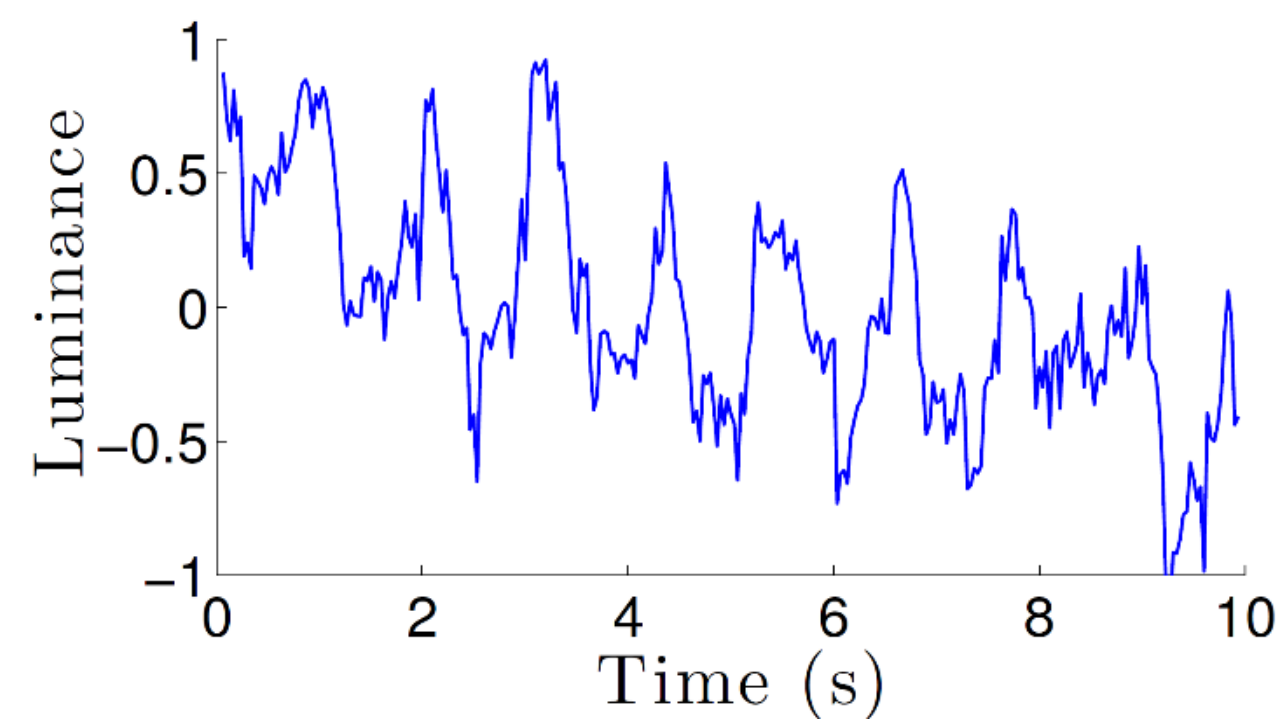
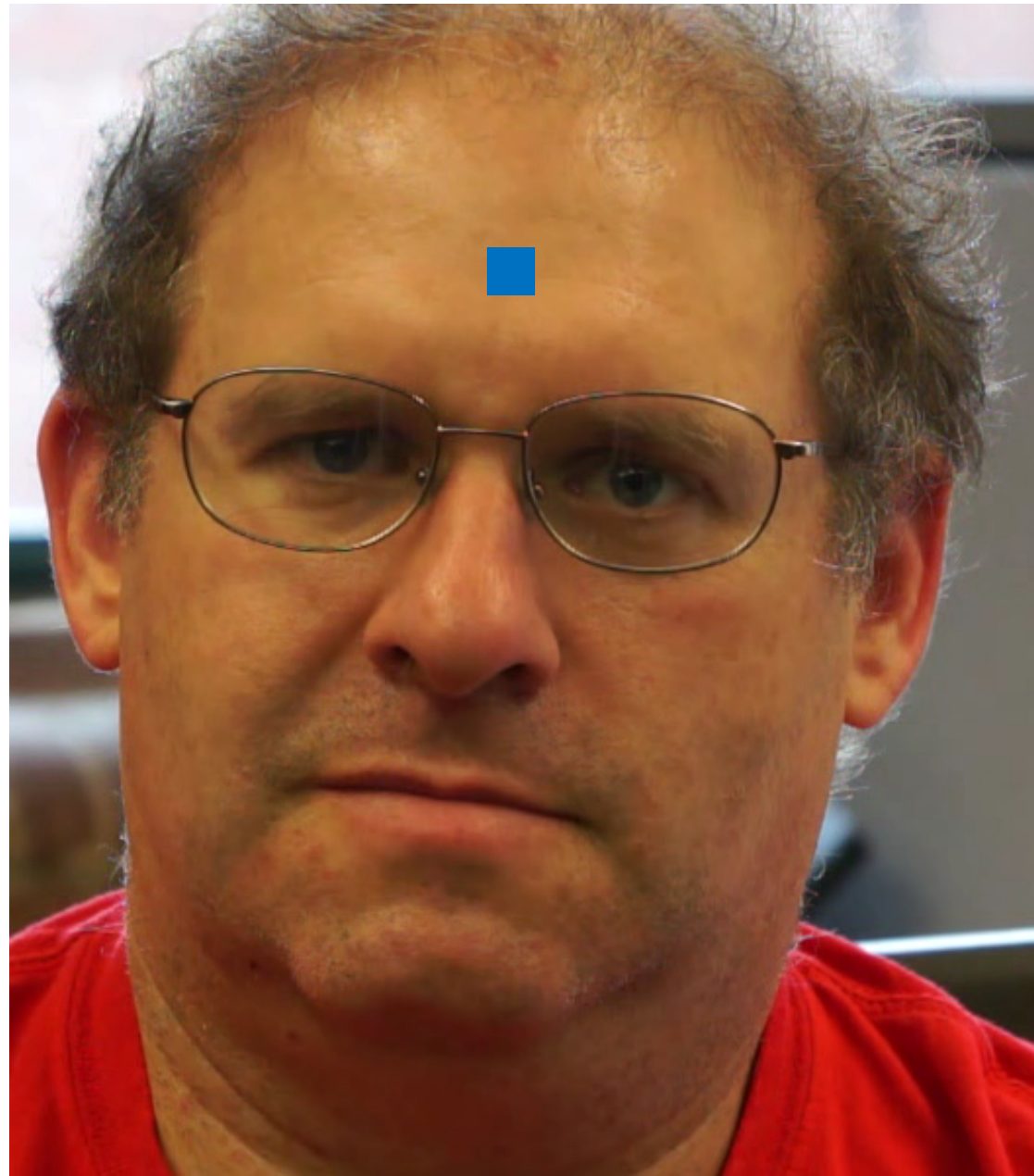


Input frame



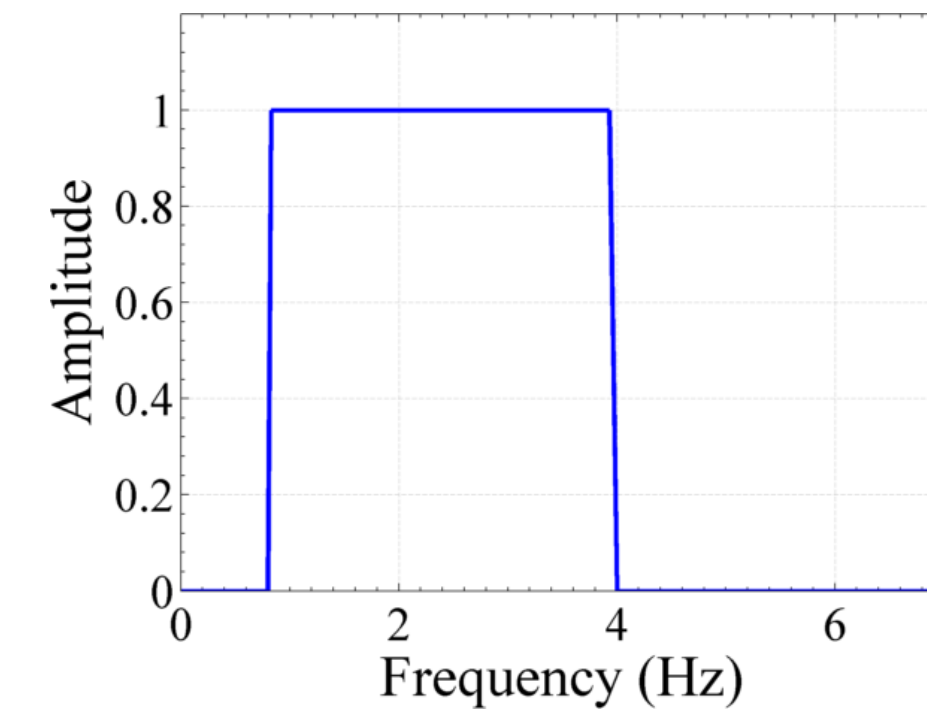
Luminance trace (zero mean)

Amplifying Subtle Color Variations



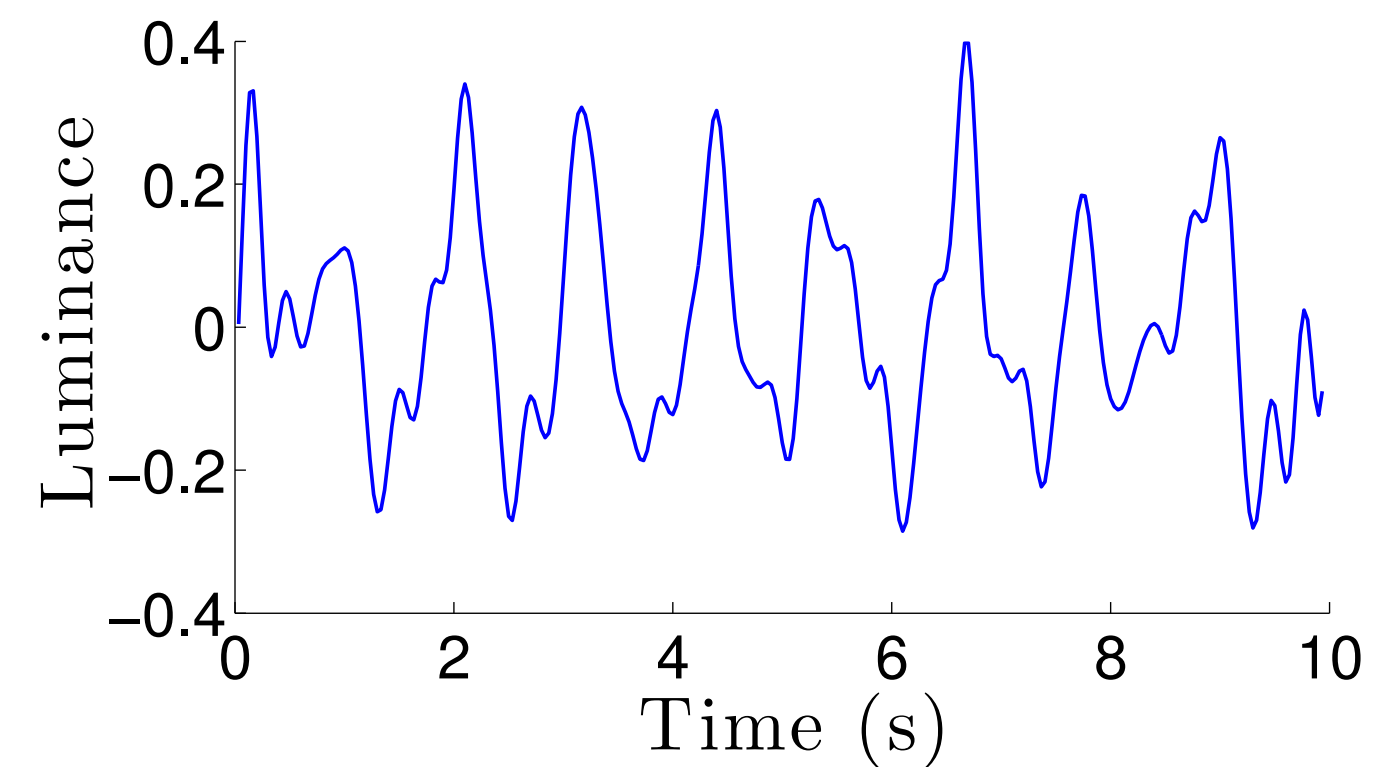
Spatially averaged luminance

\otimes



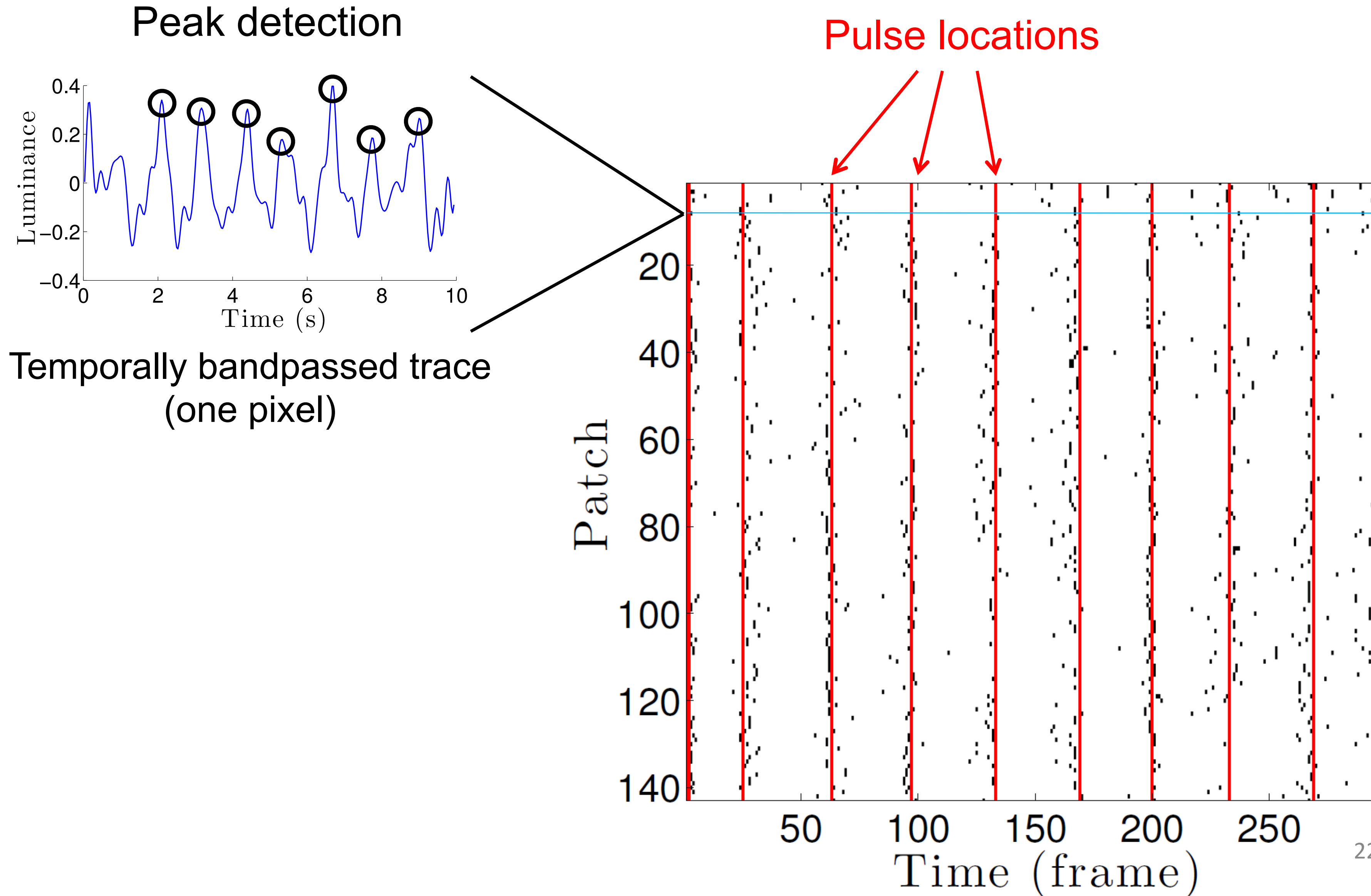
Temporal filter

$=$

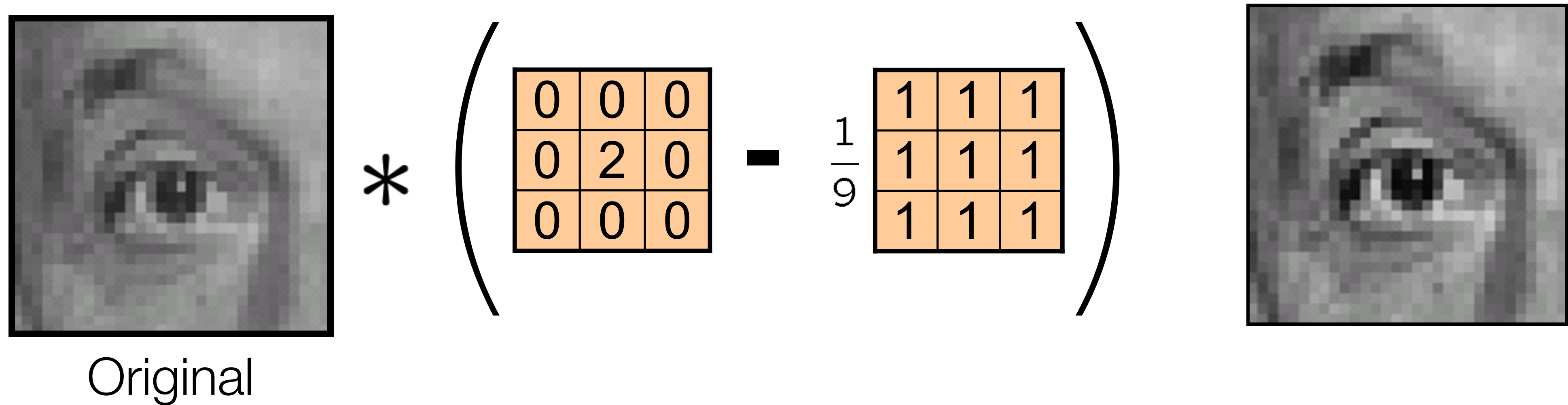


Result

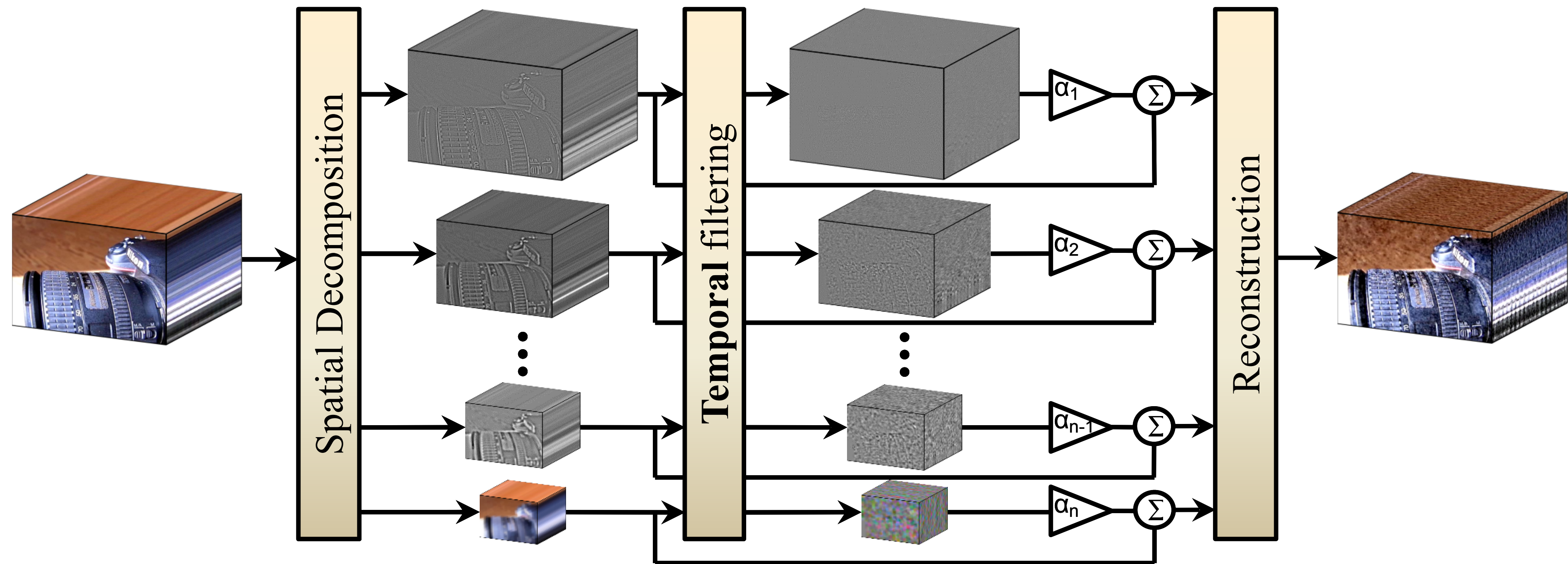
Heart Rate Extraction



Recall: sharpening filter



Magnifying tiny motions

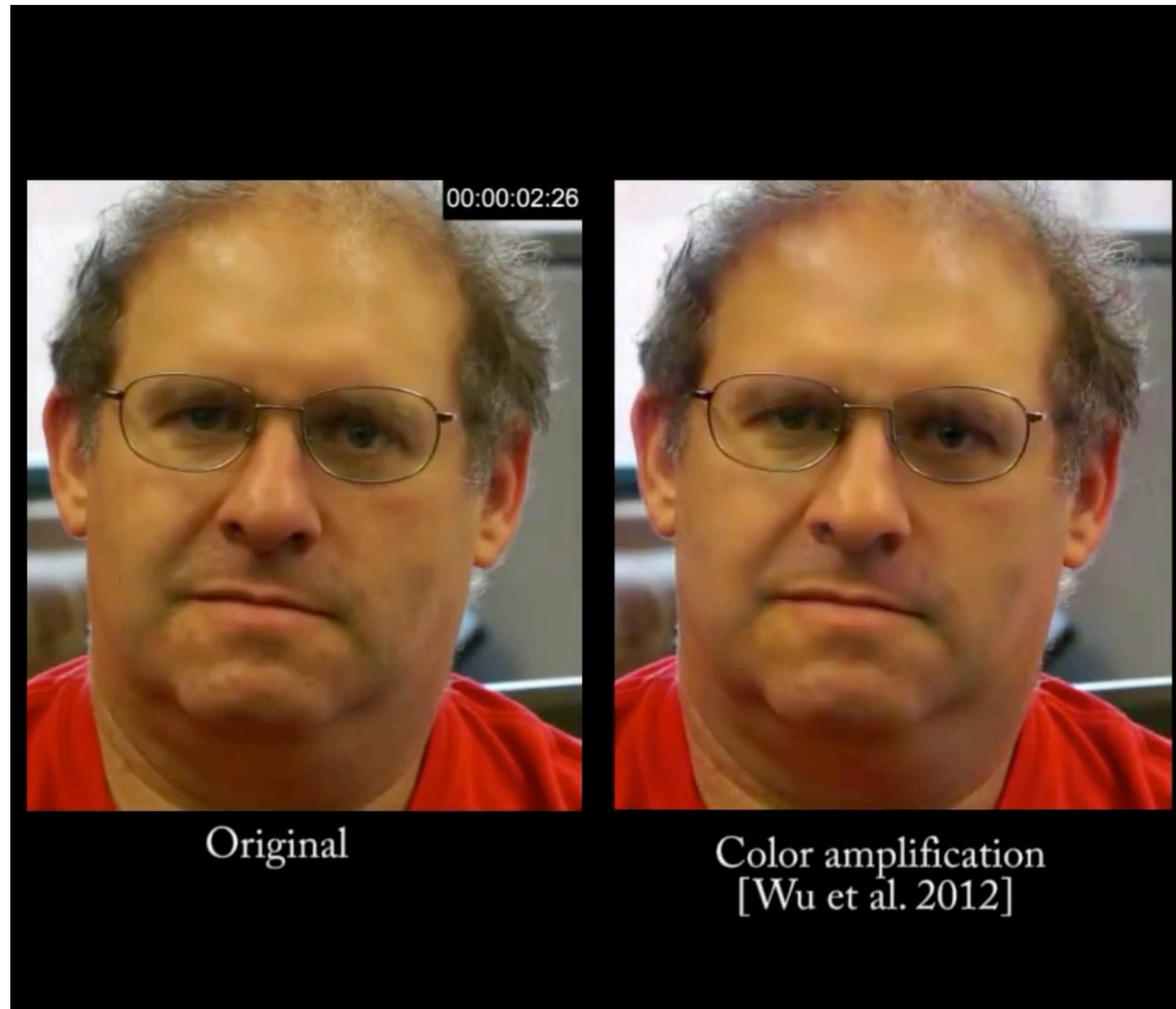


Laplacian
Pyramid

Bandpass filter
intensity at each
pixel over time

Amplify
bandpassed
signal and add
back to original

Color amplification



Motion magnification

Original



Original

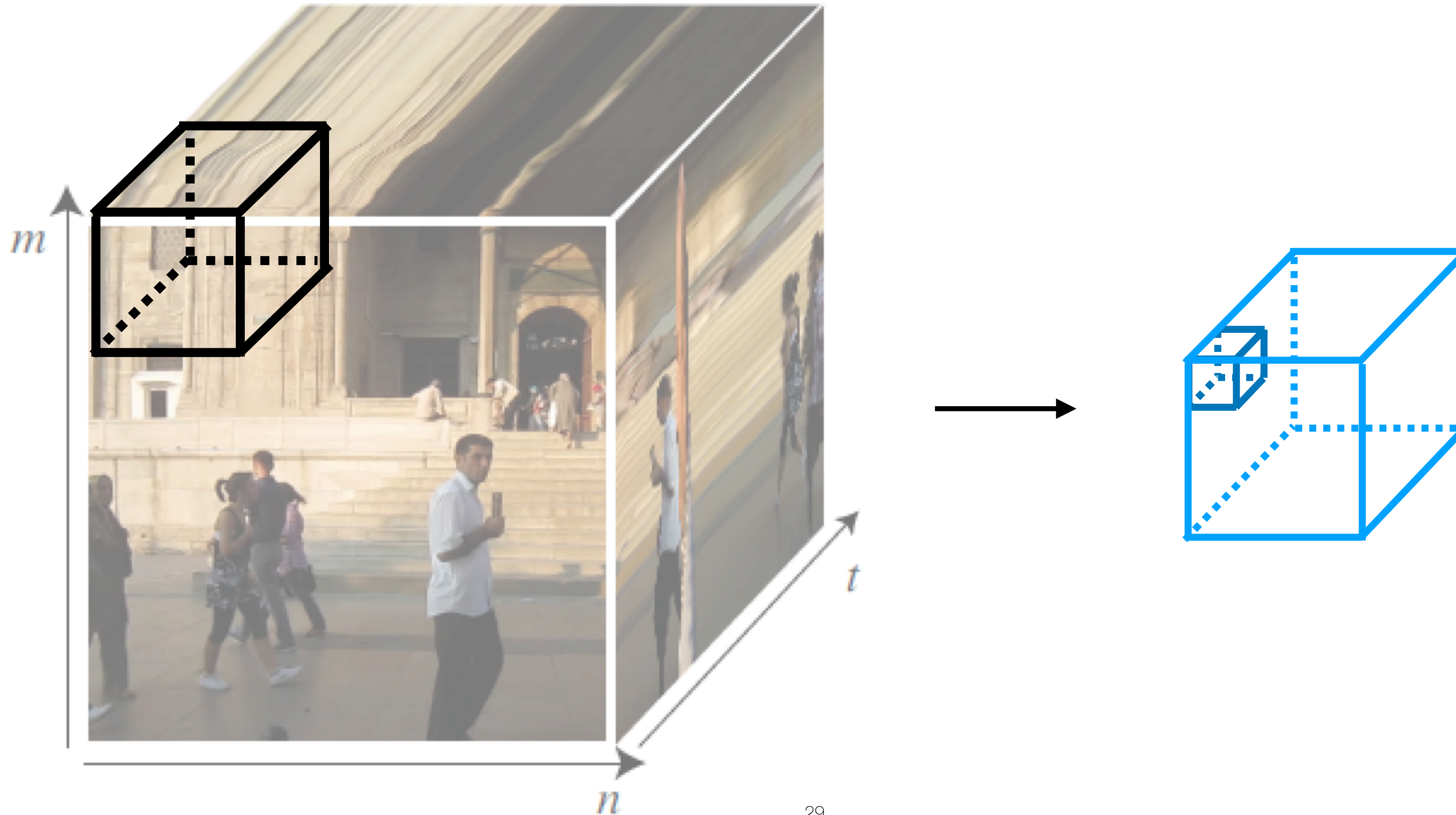


Source

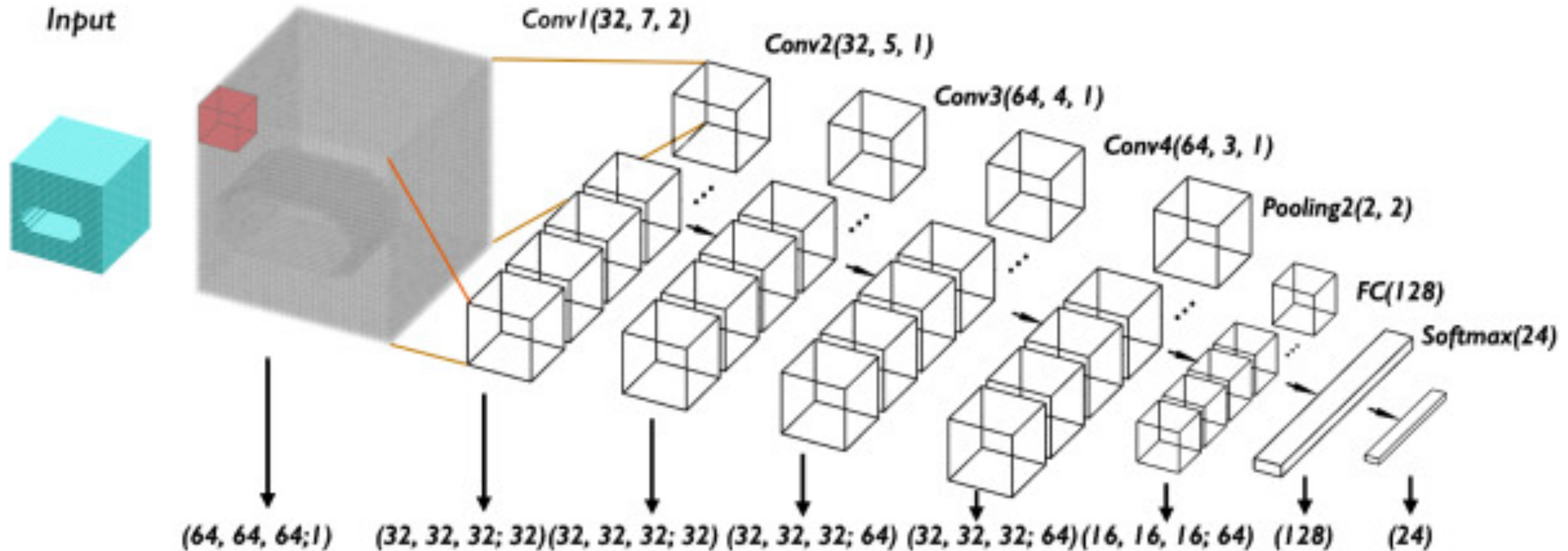
[Wadhwa et al., SIGGRAPH 2013]

Space-time convolutions

3D space-time convolution



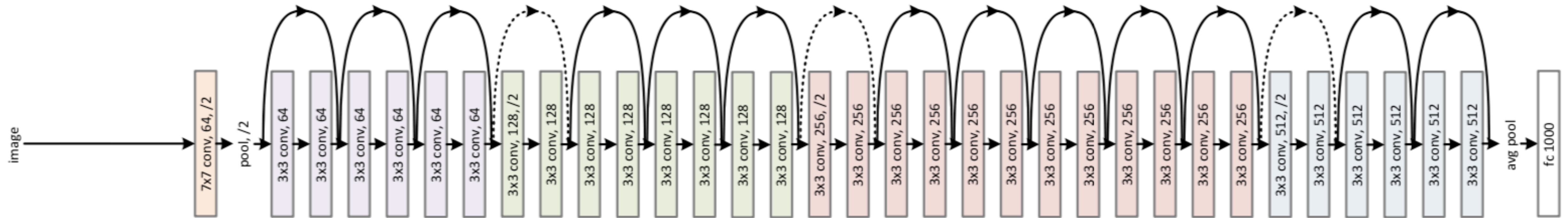
3D space-time CNN



[Source: FeatureNet: Machining feature recognition based on 3D Convolution Neural Network]

Designing a 3D CNN architecture

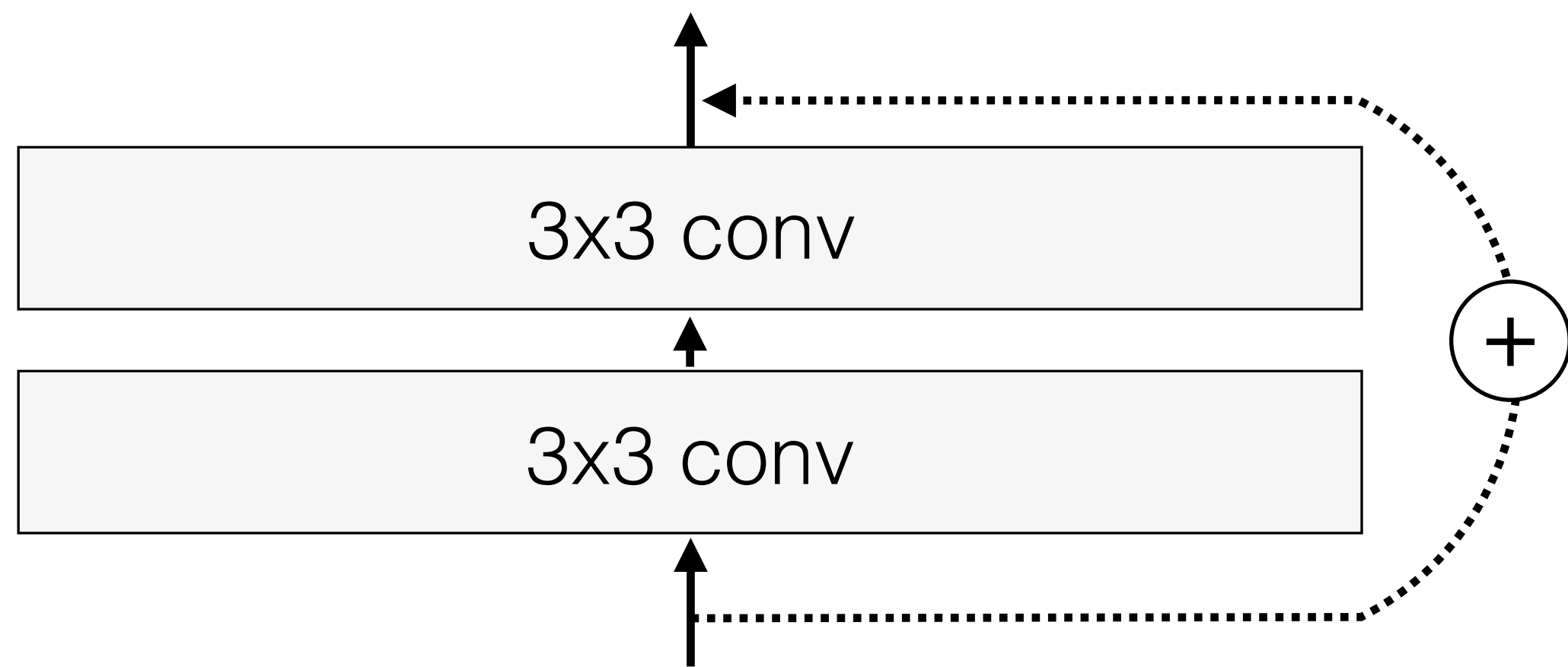
Starting point: 2D image CNNs



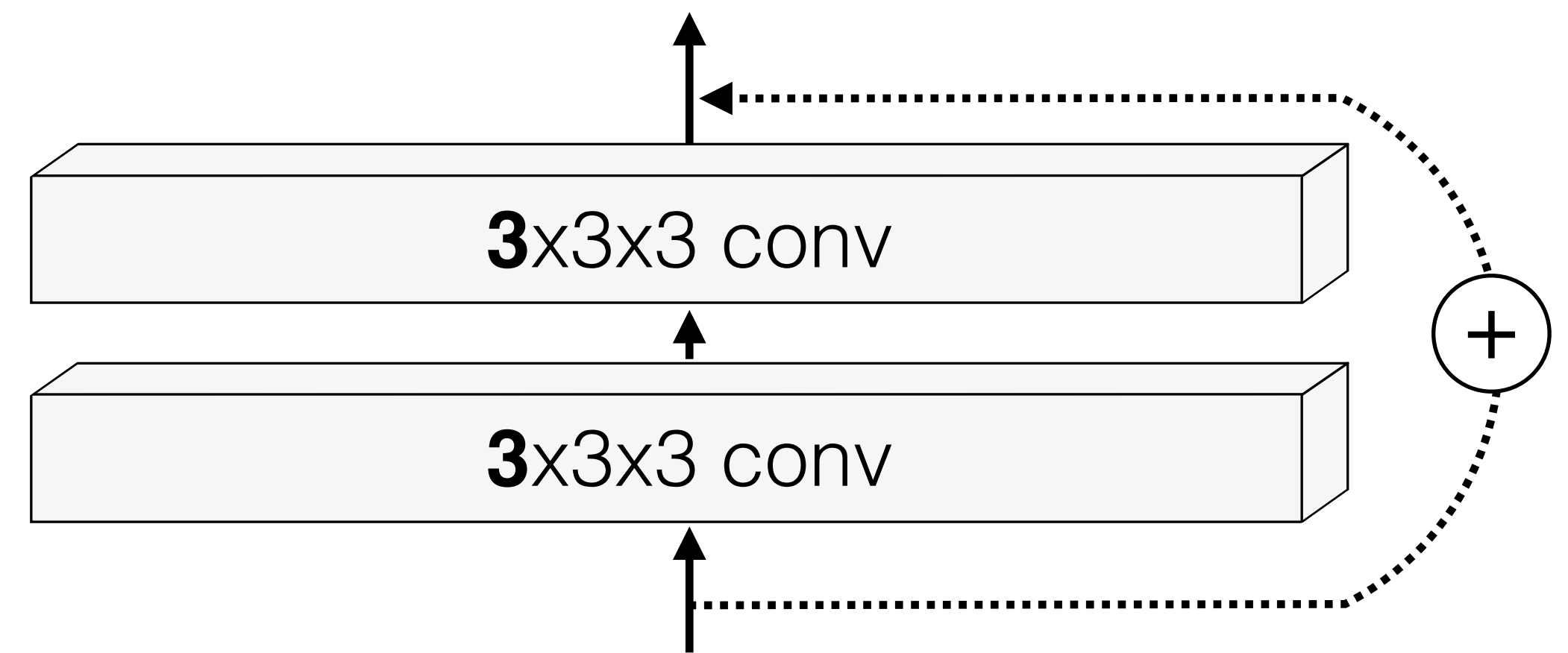
ResNet [Kaiming He et al. 2016]

Inflated convolutions

2D ResNet block



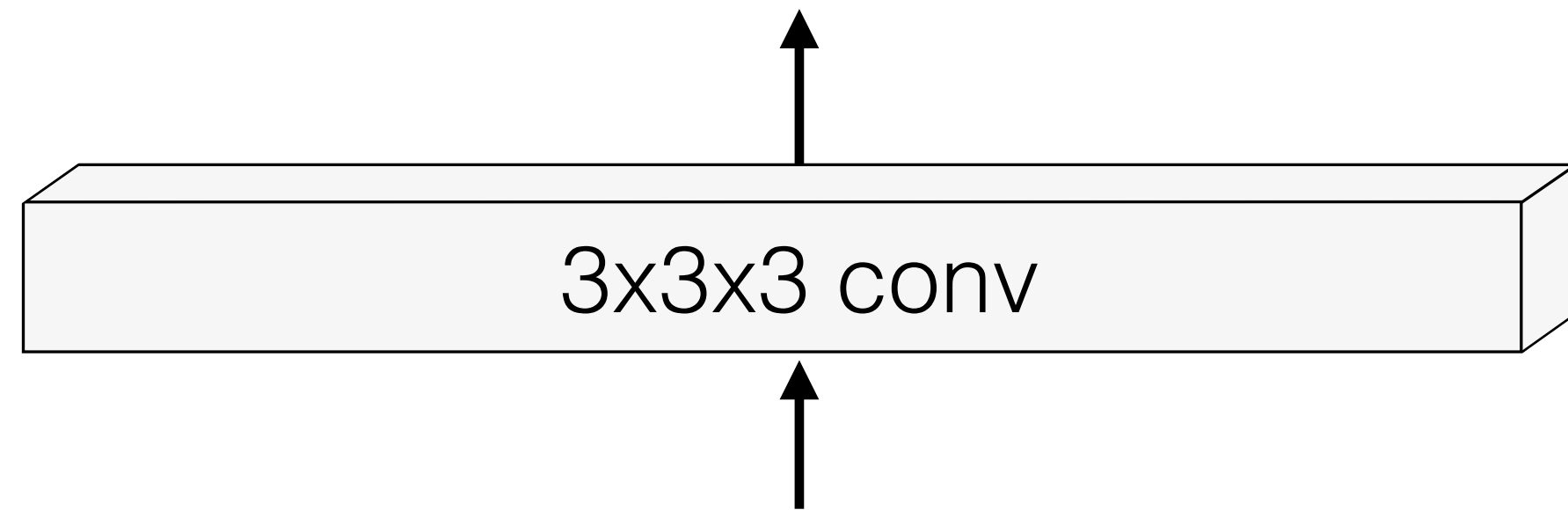
3D ResNet block



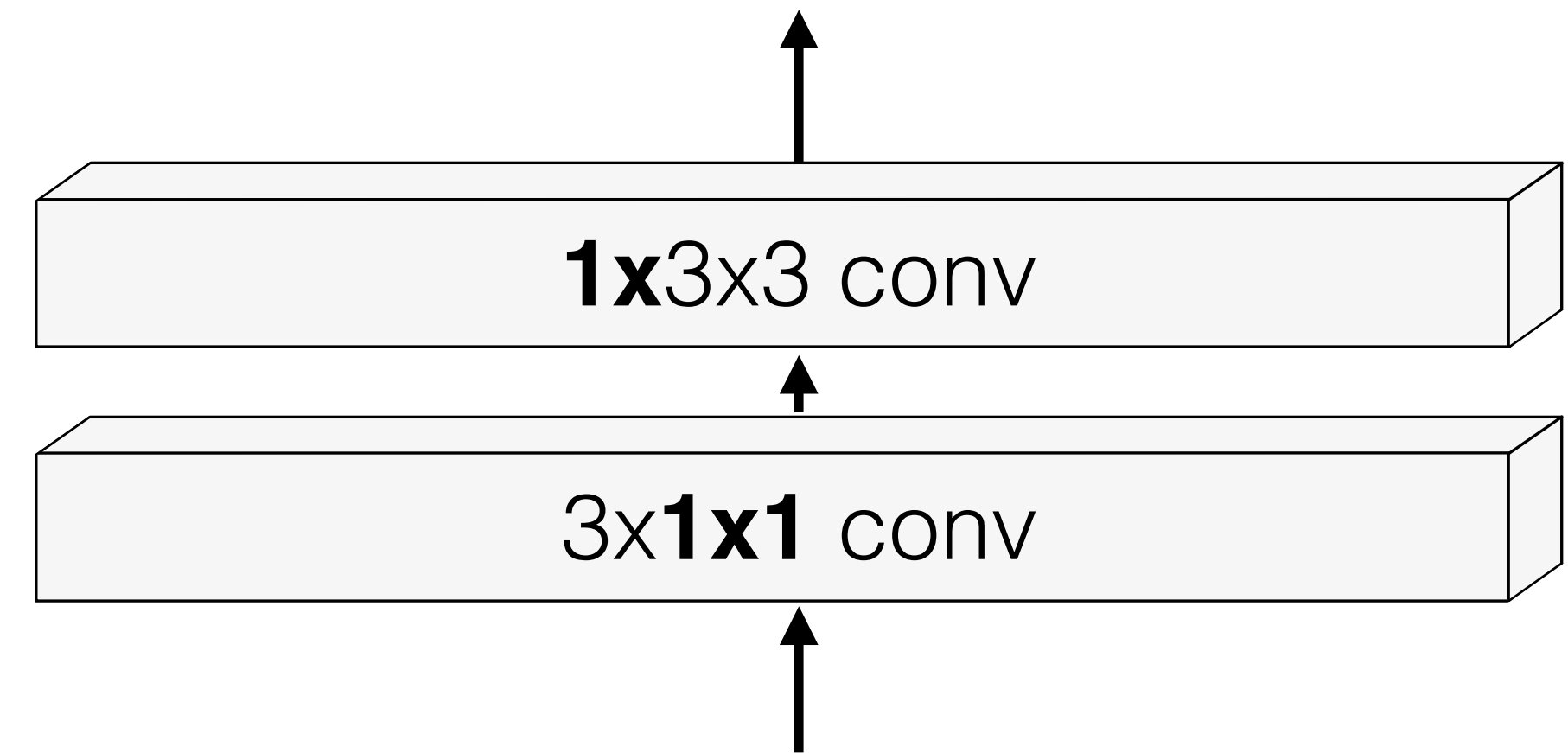
- Can reuse 2D architectures. [Carreira et al. 2017]
- Pretrain with 2D nets (“inflating” 2D filter to 3D)

Separable convolutions

3D convolution



Separate space/time



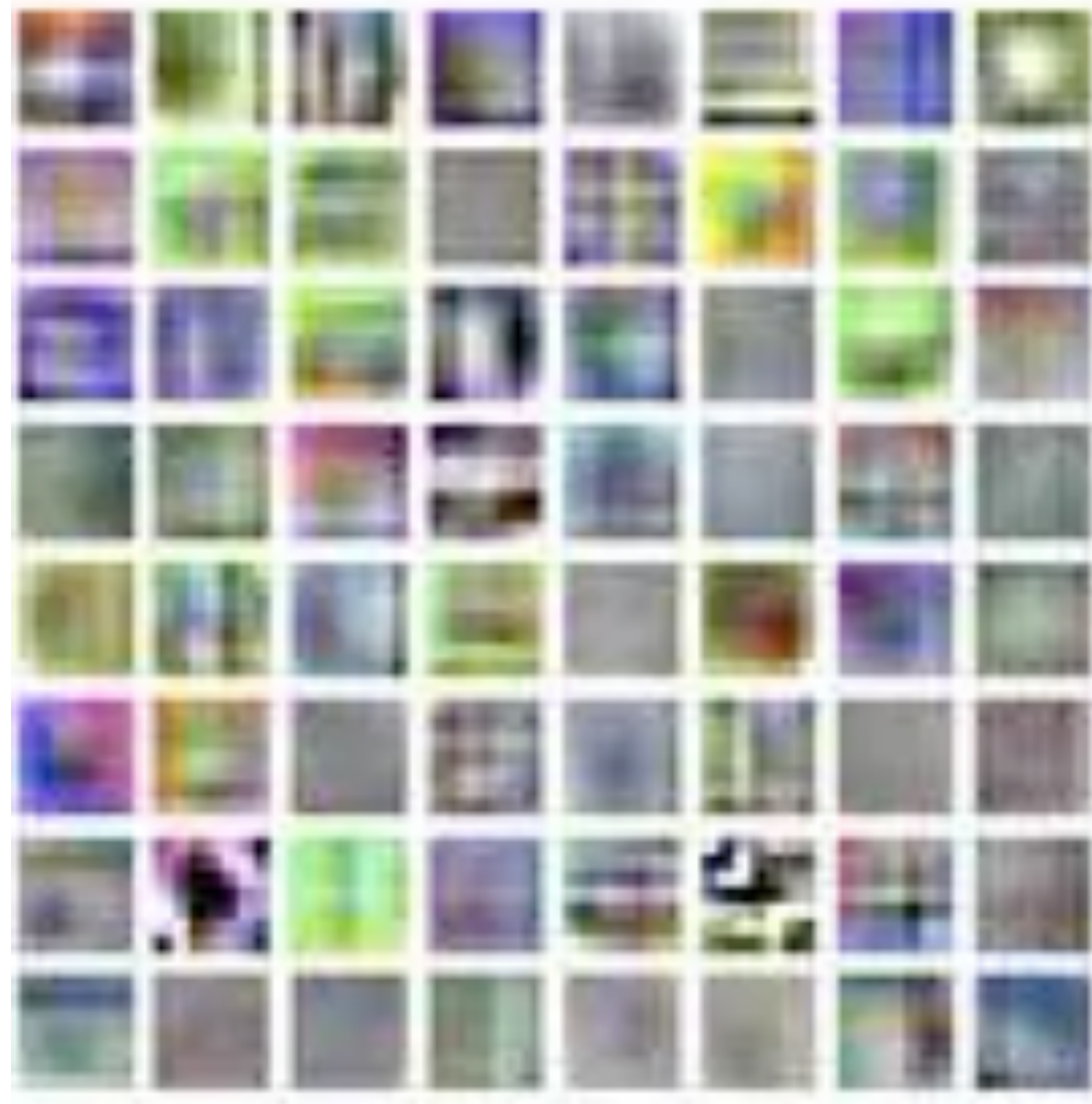
Often works well. Faster and fewer parameters.

$$3 \times 3 \times 3 = 27$$

vs.

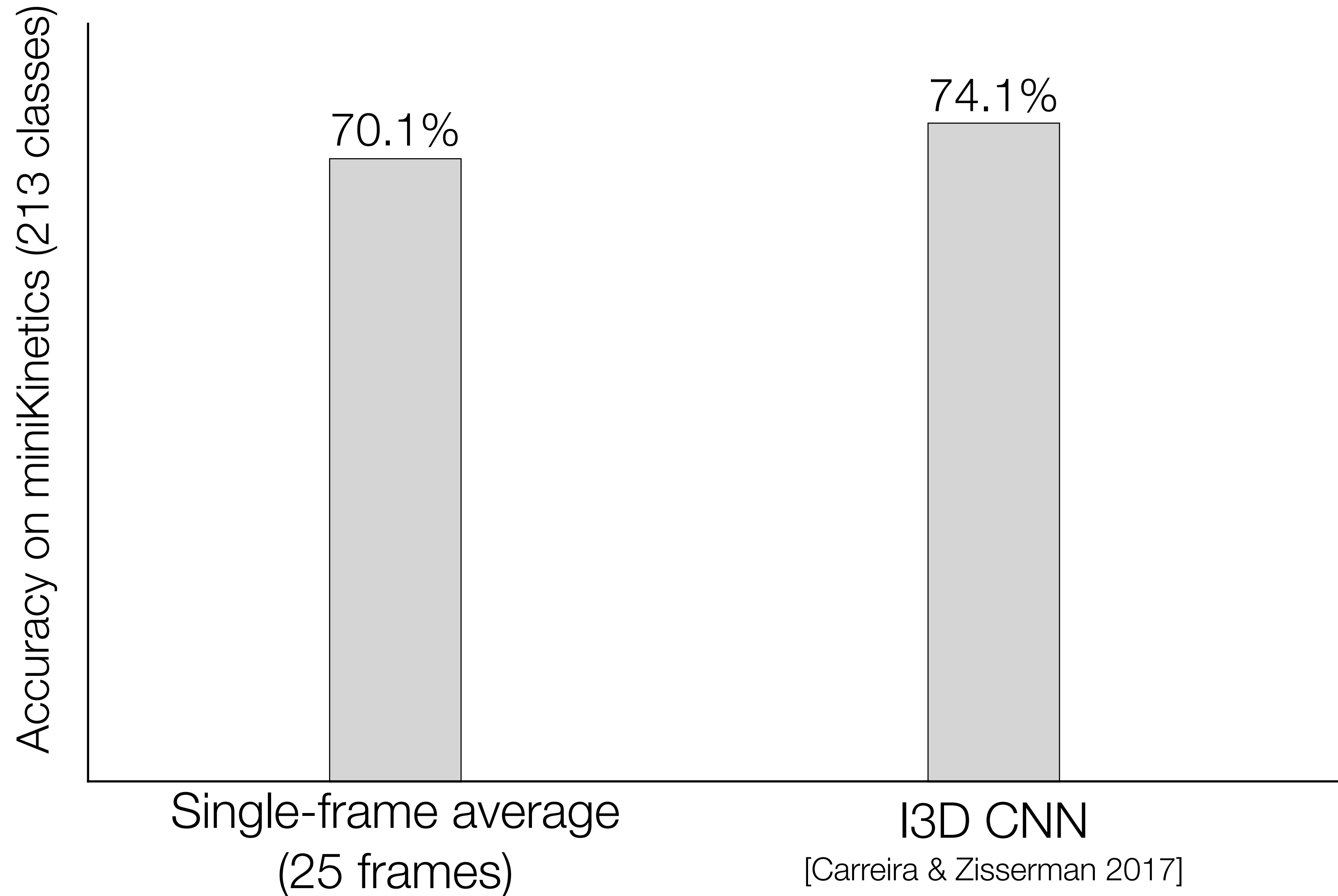
$$3 \times 3 + 3 = 12$$

Learned space-time filters



$(7 \times 7 \times 7)$ I3D conv1 filters, [Carreira & Zisserman 2017]

When do we actually need motion?



When do we actually need motion?

Let's look at these again:



Making latte art



Jaywalking



Grooming dog

When do we actually need motion?

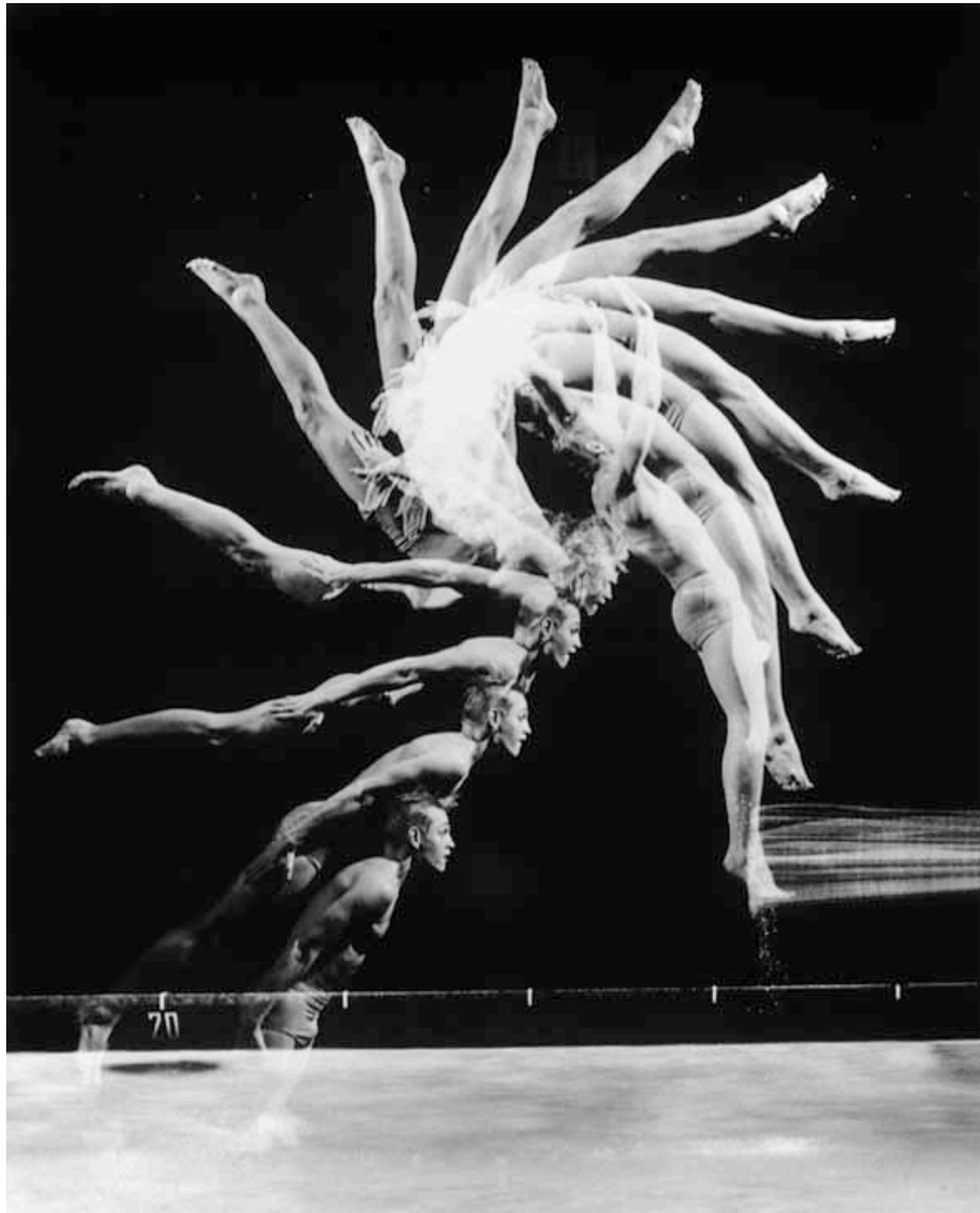


Photo by H. Edgerton

- Single-frame models usually don't work well.
- But single-image model + temporal pooling often is surprisingly competitive.
- In many tasks that use video, time often provides **extra samples**, rather than motion.
- Later in the course, we'll see tasks where motion is essential, such as 3D reconstruction.

Recurrent networks

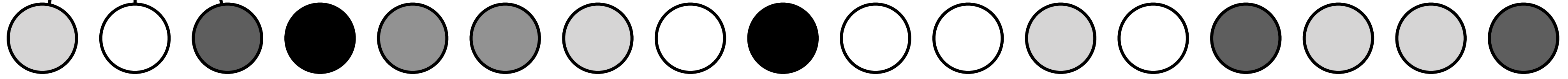
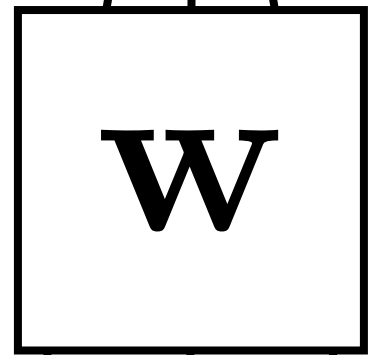
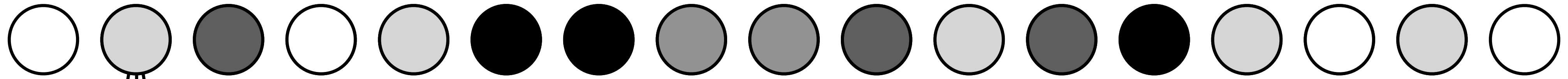
mihaifrancu



39

[<https://www.youtube.com/watch?v=wxfgT-kKxiM>]

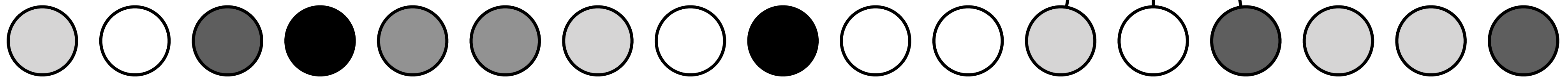
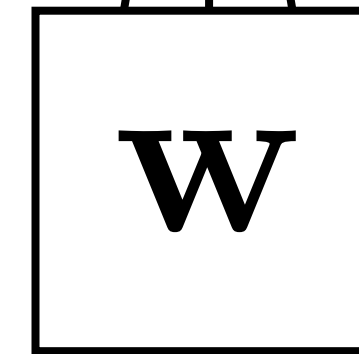
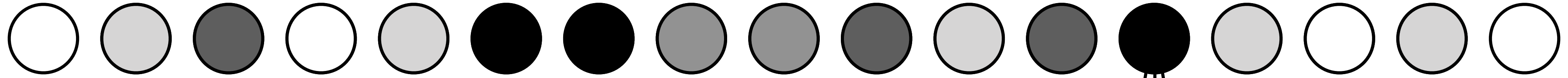
Rufus



time



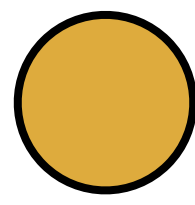
Douglas



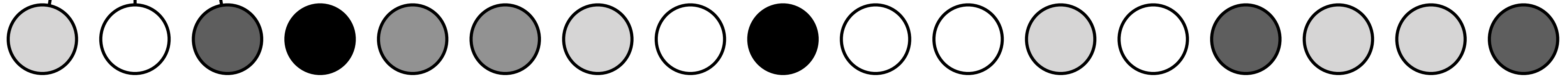
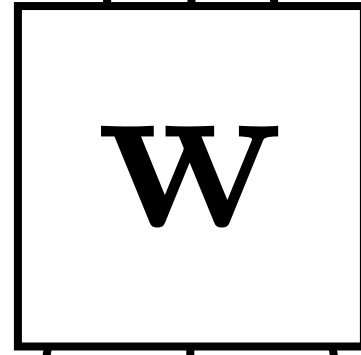
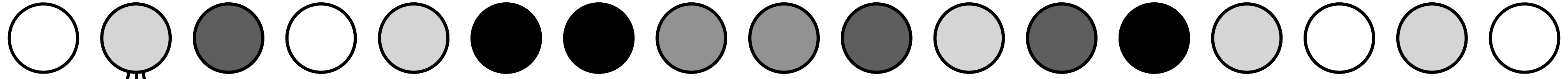
time



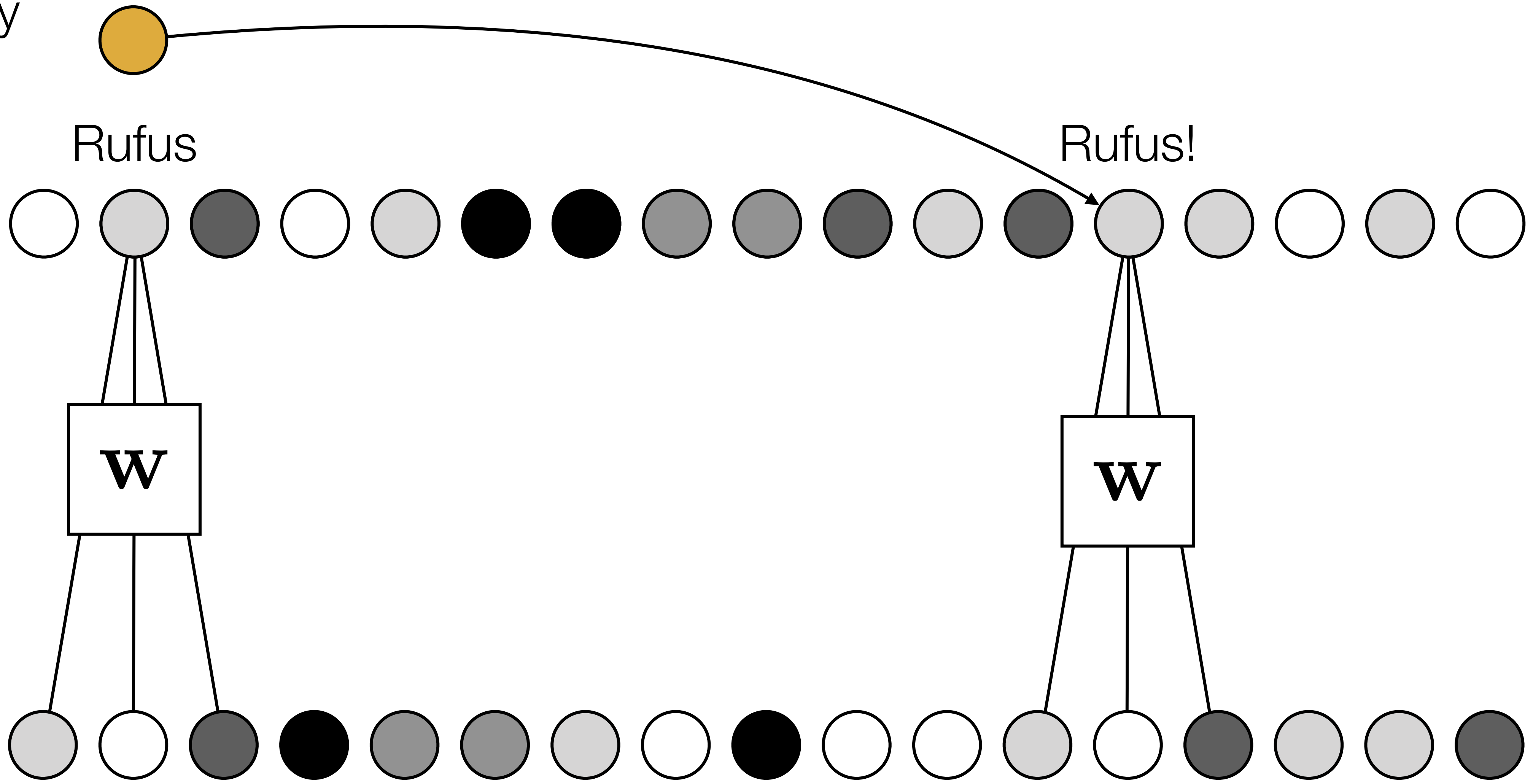
Memory unit



Rufus

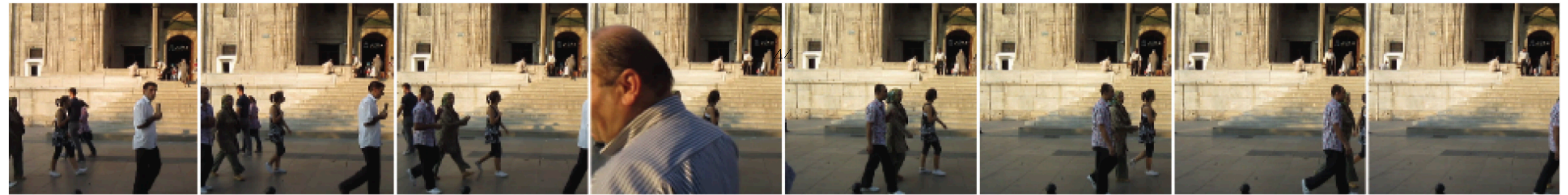
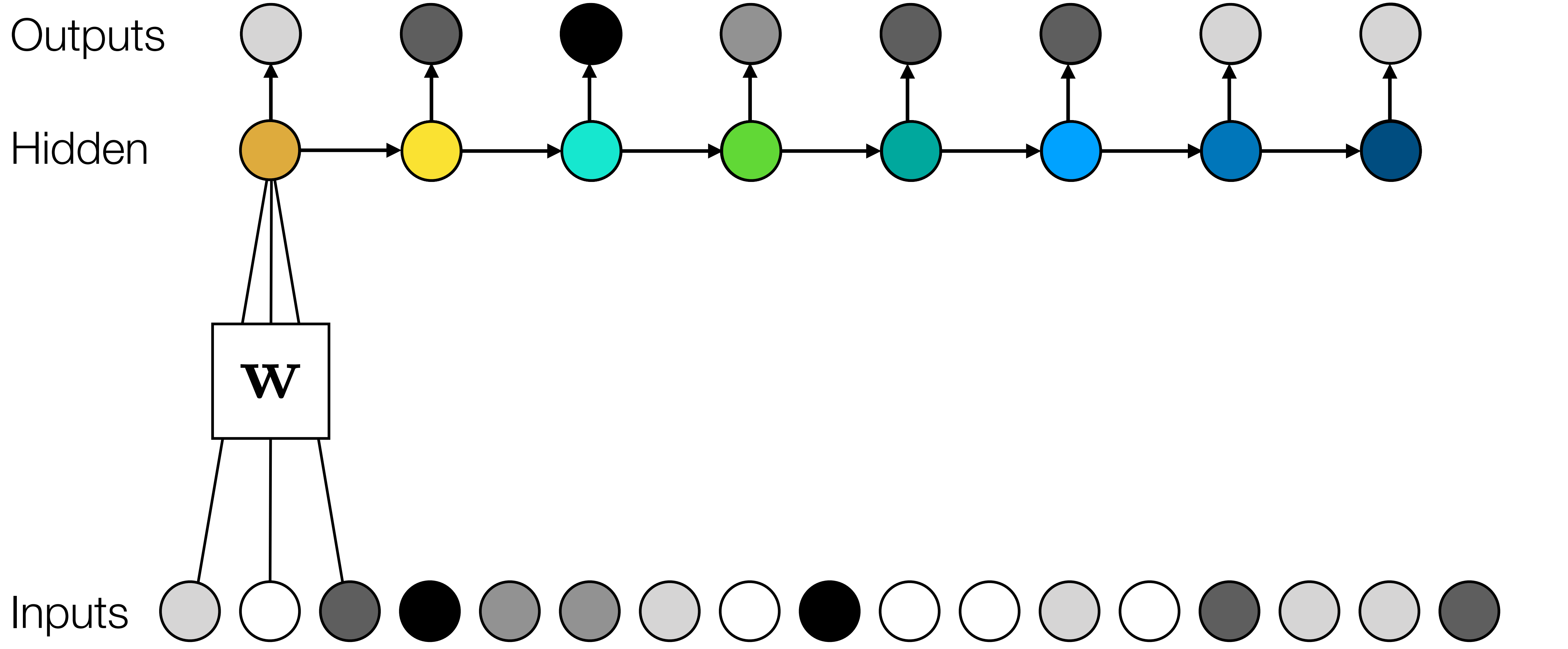


Memory unit

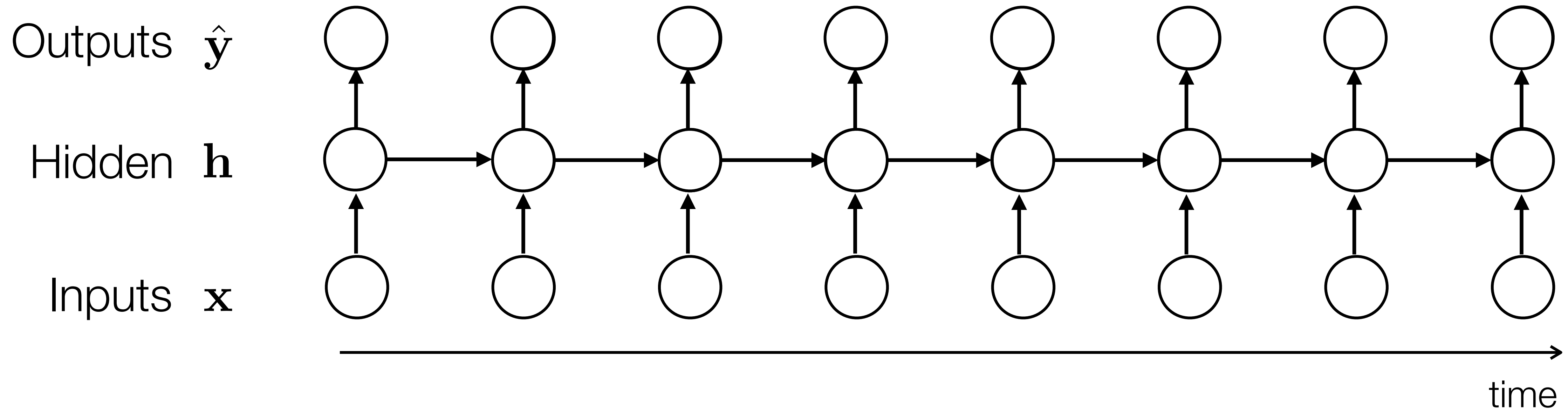


time

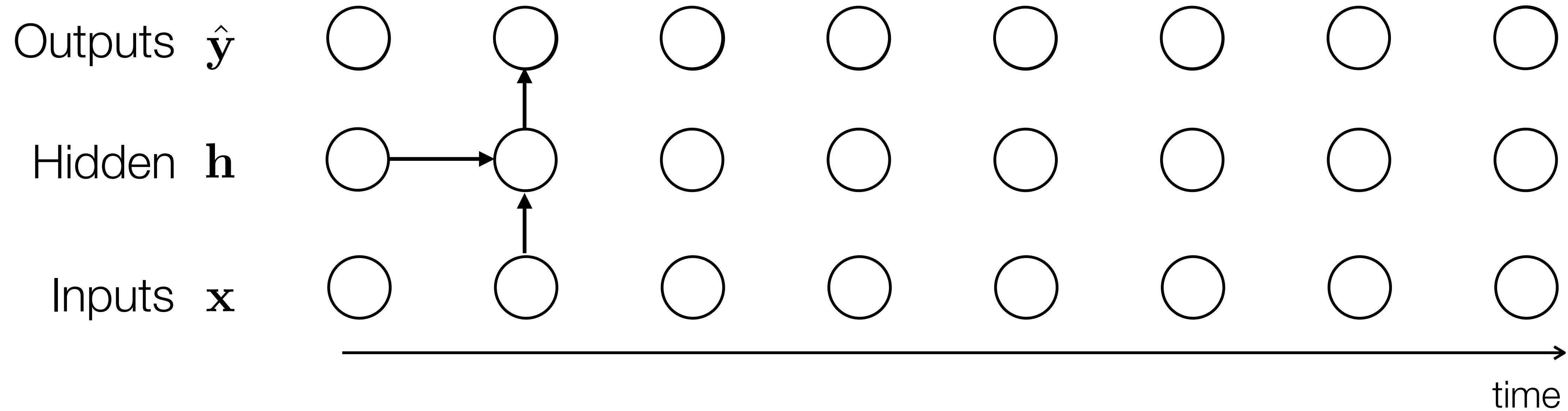
Recurrent Neural Networks (RNNs)



Recurrent Neural Networks (RNNs)



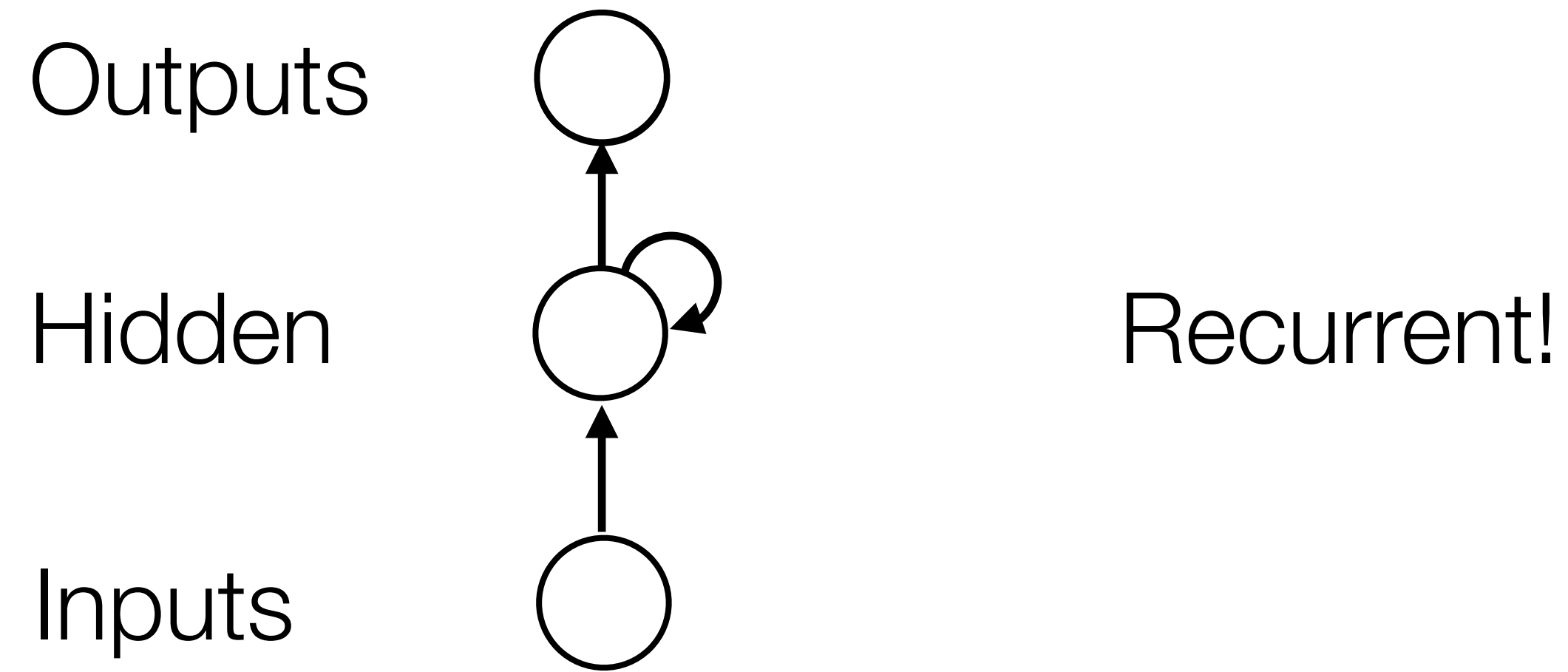
Recurrent Neural Networks (RNNs)



$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)})$$

$$\mathbf{y}^{(t)} = g(\mathbf{h}^{(t)})$$

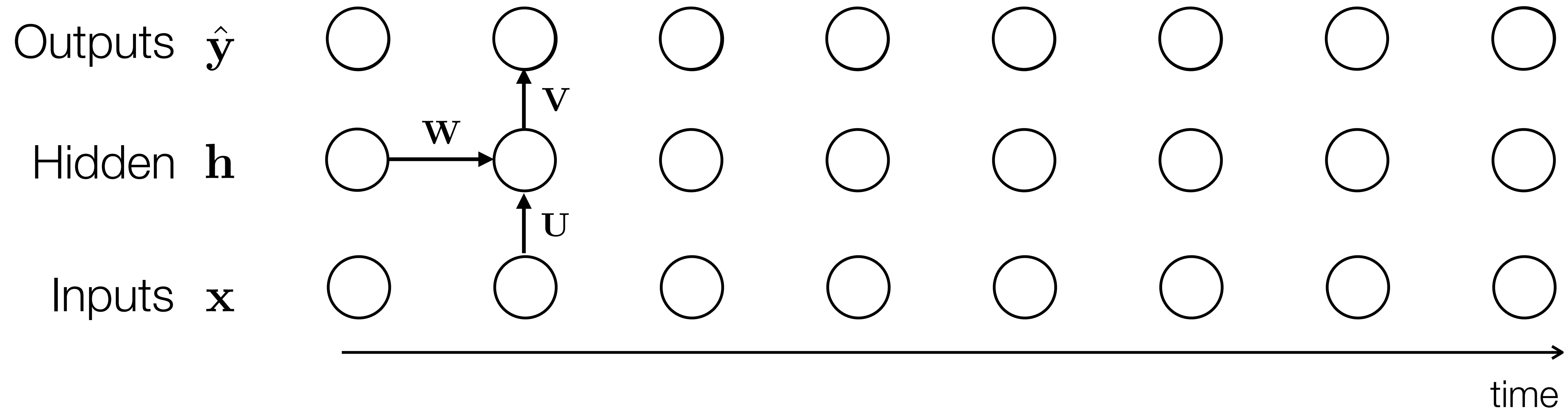
Recurrent Neural Networks (RNNs)



$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)})$$

$$\mathbf{y}^{(t)} = g(\mathbf{h}^{(t)})$$

Recurrent Neural Networks (RNNs)



$$\mathbf{a}^{(t)} = \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b}$$

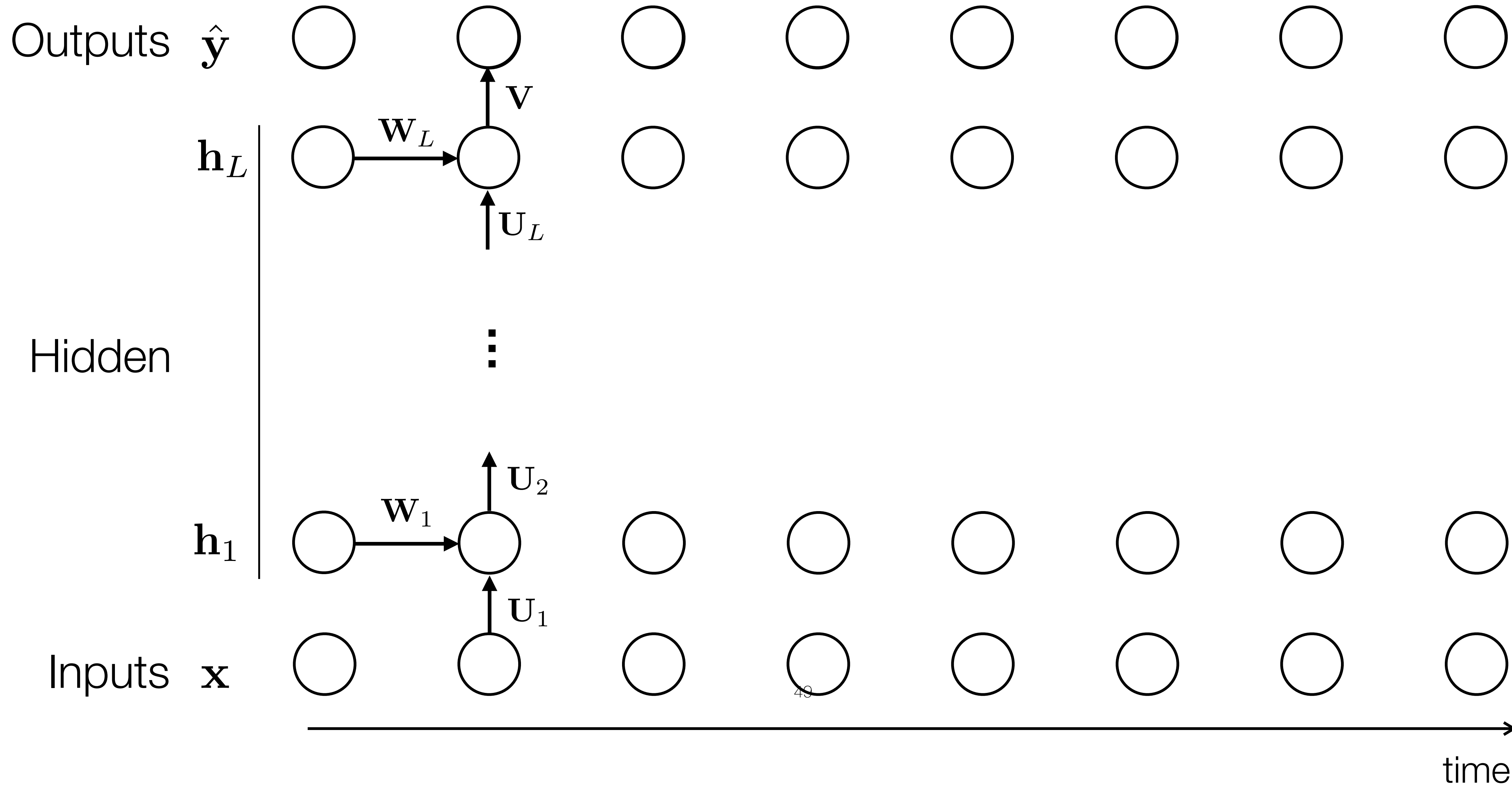
$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)})$$

$$\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)} + \mathbf{c}$$

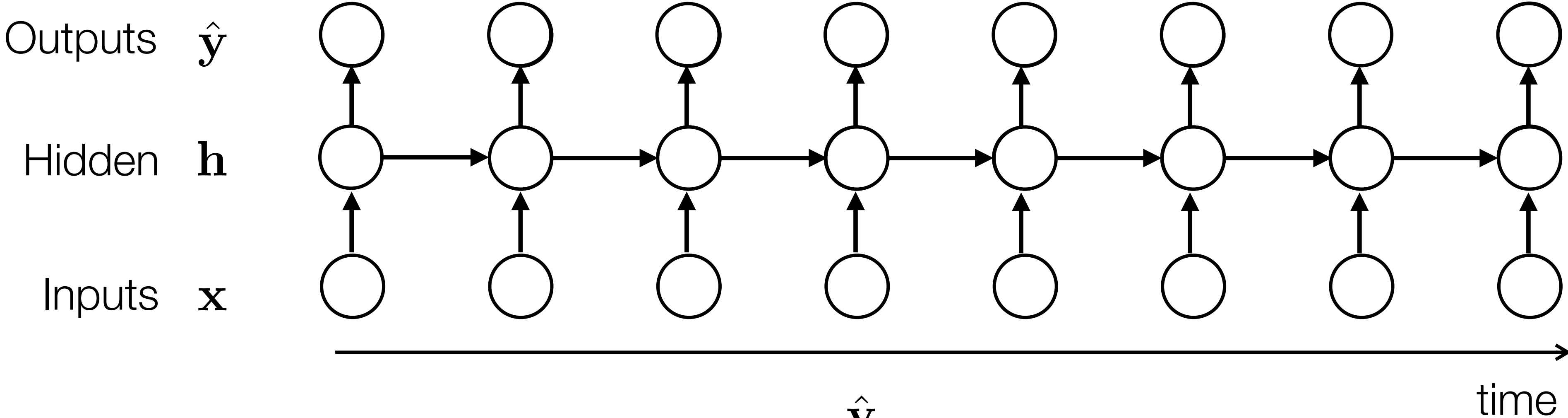
48

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)})$$

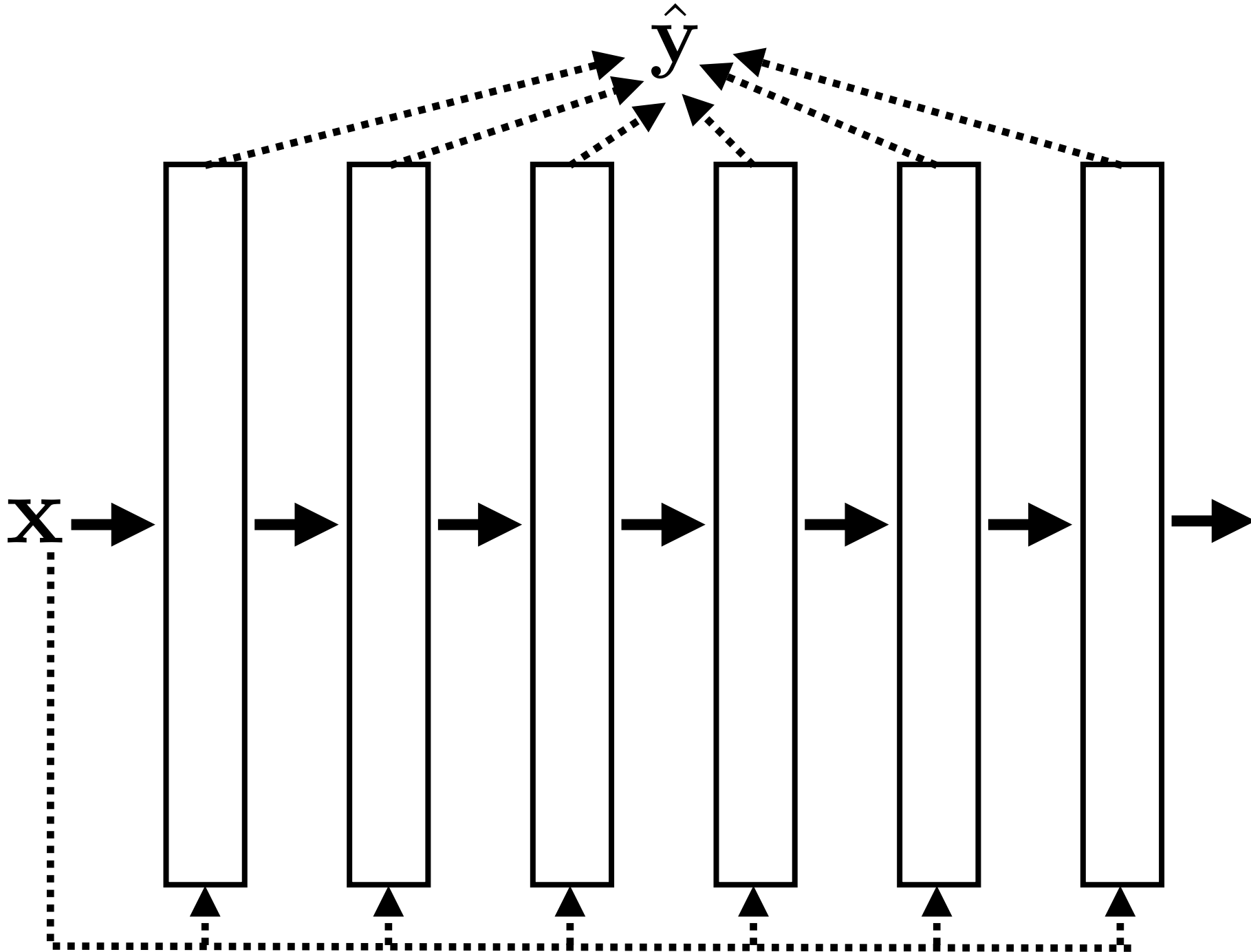
Deep Recurrent Neural Networks (RNNs)



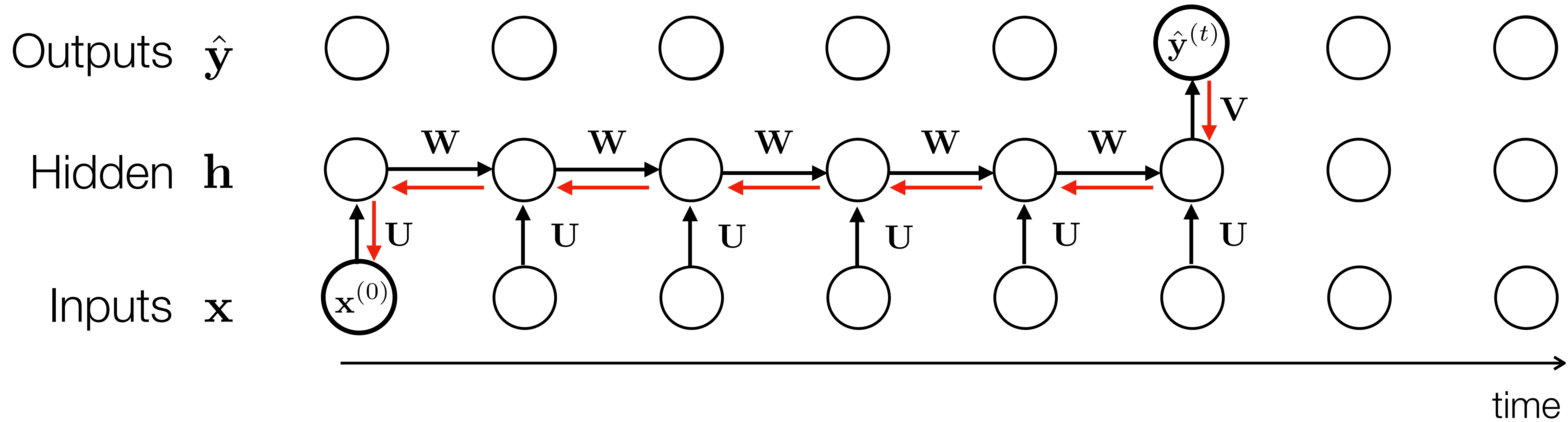
Unrolling an RNN



Equivalent to
"unrolled" network
with shared weights

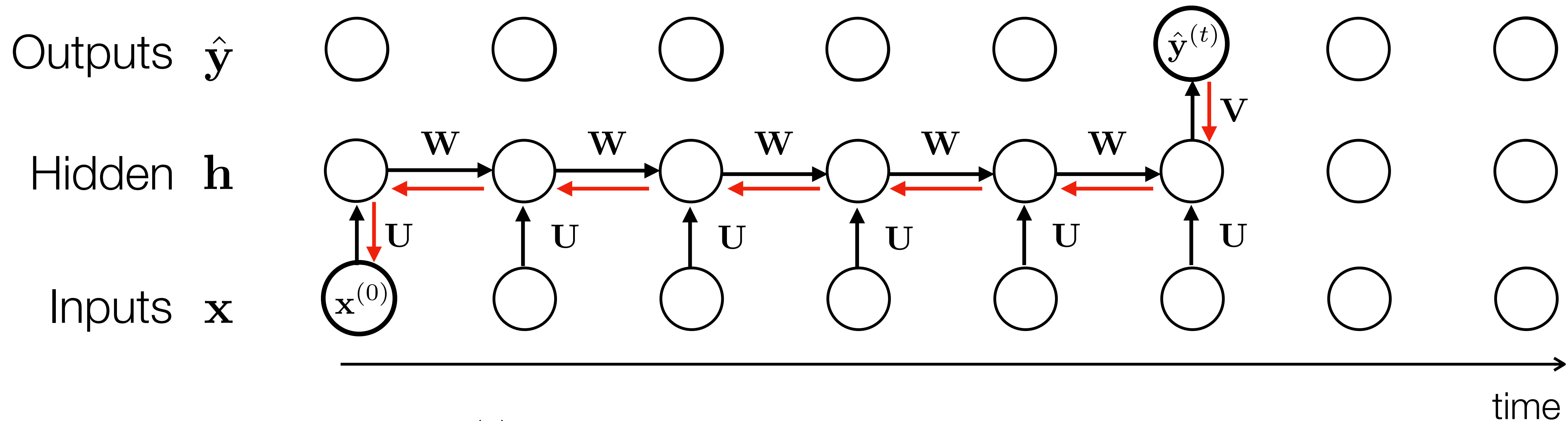


Backprop in RNNs



$$\frac{\partial \hat{\mathbf{y}}^{(t)}}{\partial \mathbf{x}^{(0)}} = \frac{\partial \hat{\mathbf{y}}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \cdots \frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{h}^{(0)}} \frac{\partial \mathbf{h}^{(0)}}{\partial \mathbf{x}^{(0)}}$$

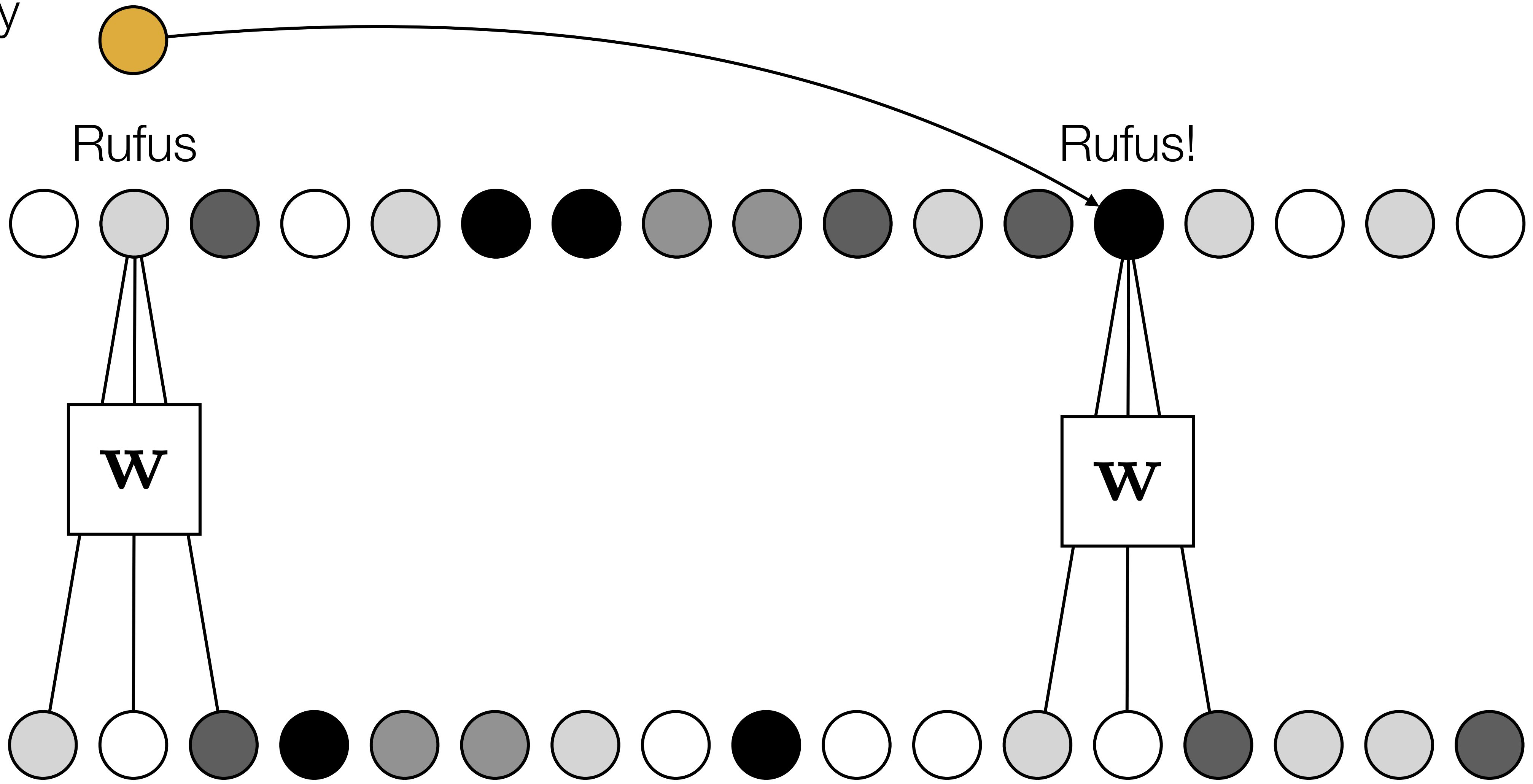
The problem of long-range dependences



$$\frac{\partial \hat{\mathbf{y}}^{(t)}}{\partial \mathbf{x}^{(0)}} = \frac{\partial \hat{\mathbf{y}}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \cdots \frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{h}^{(0)}} \frac{\partial \mathbf{h}^{(0)}}{\partial \mathbf{x}^{(0)}}$$

- Capturing long-range dependences requires propagating information through a long chain.
- Old observations are forgotten
- Stochastic gradients become high variance (noisy), and gradients may **vanish** or **explode**

Memory unit

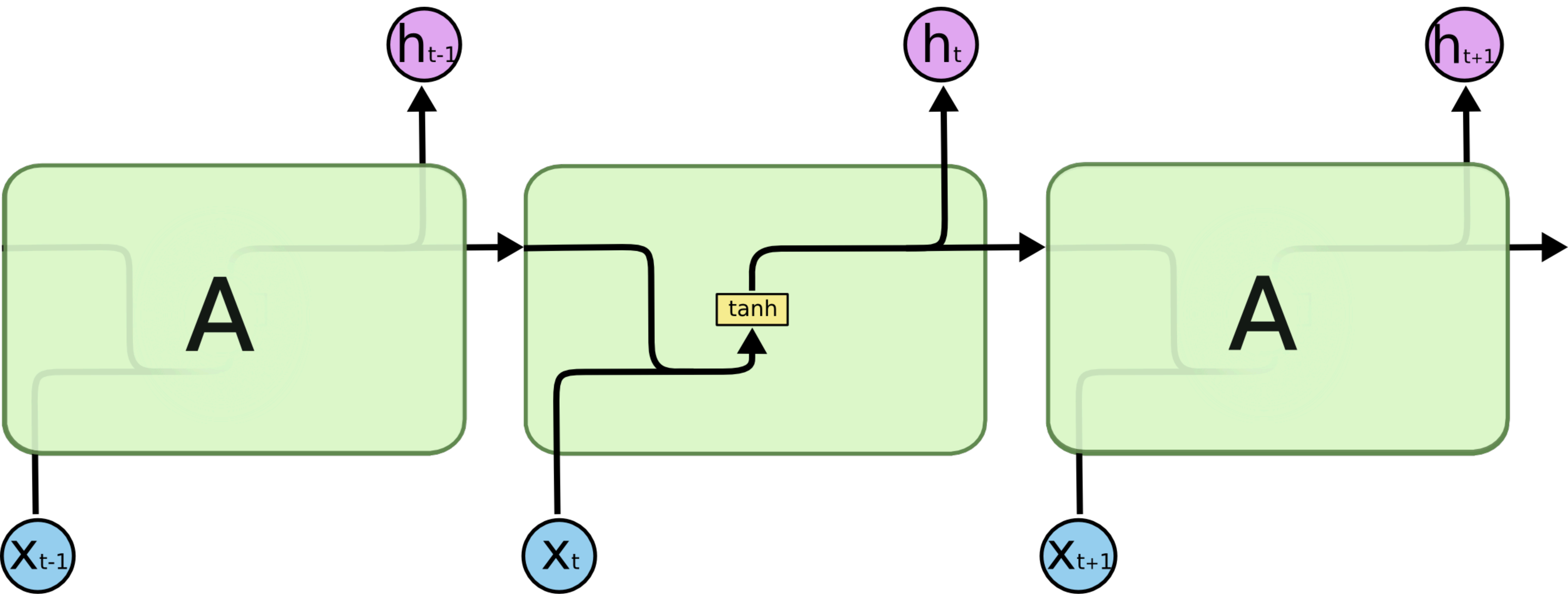


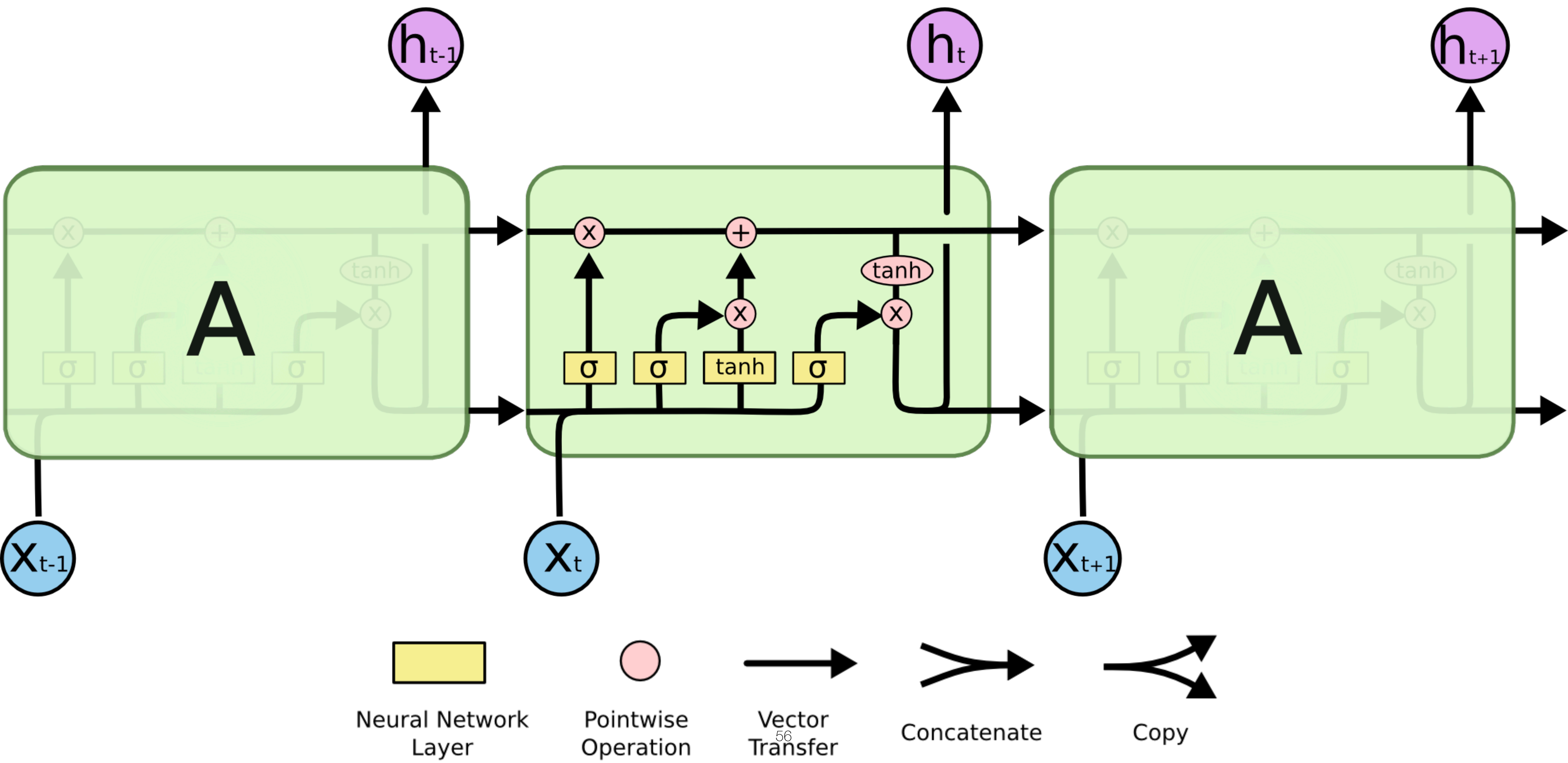
time

LSTMs

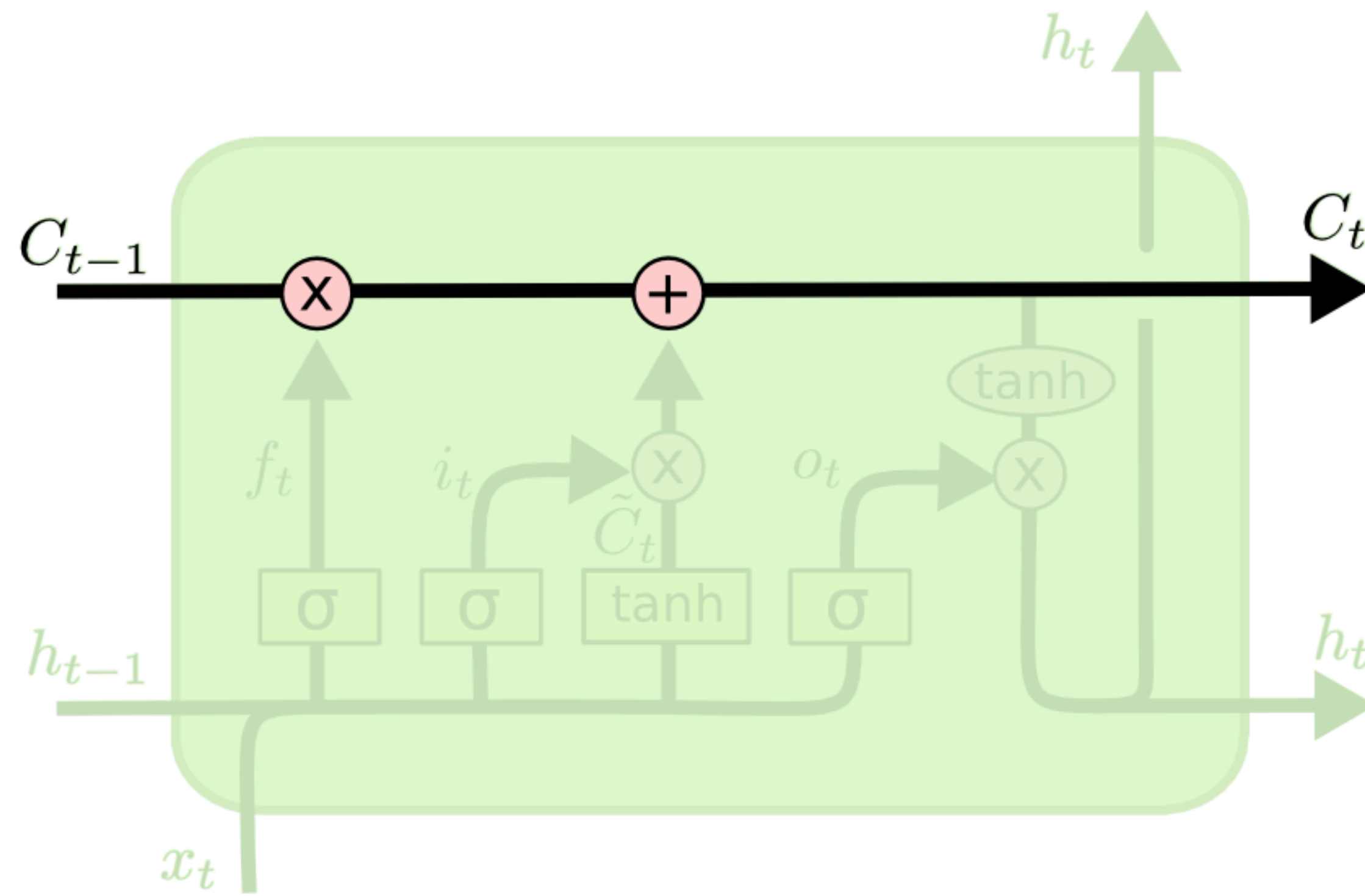
Long Short Term Memory

- A special kind of RNN designed to avoid forgetting [Hochreiter & Schmidhuber 1995].
- Related to ResNets' bias is that state transition is an identity function.
- This way the default behavior is not to forget an old state. Instead of forgetting by default, the network has to *learn to forget*.
- Bit of a complex design. Works well but simpler methods like Gated Recurrent Unit (GRU) are competitive [Jozefowicz et al. 2015].

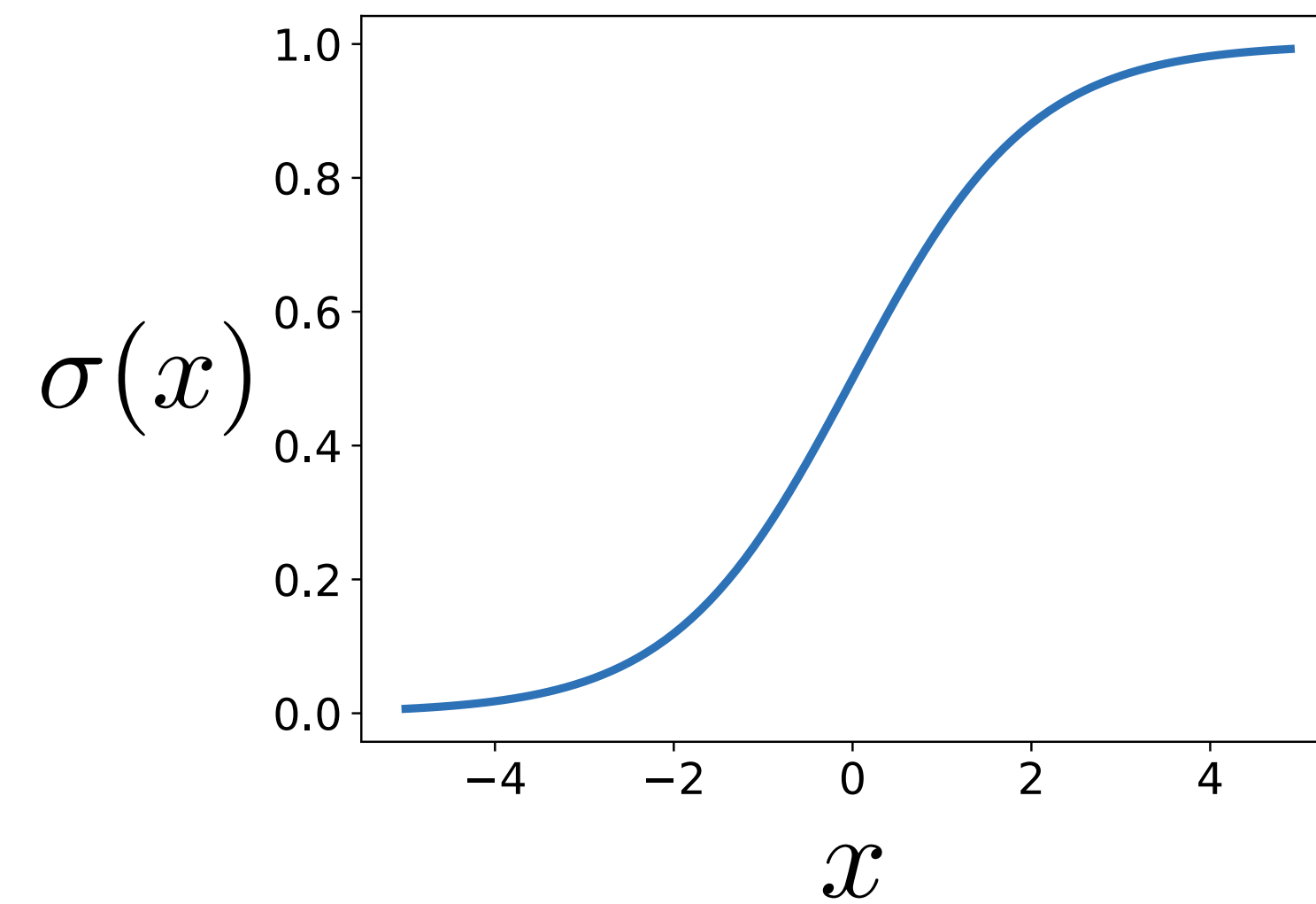
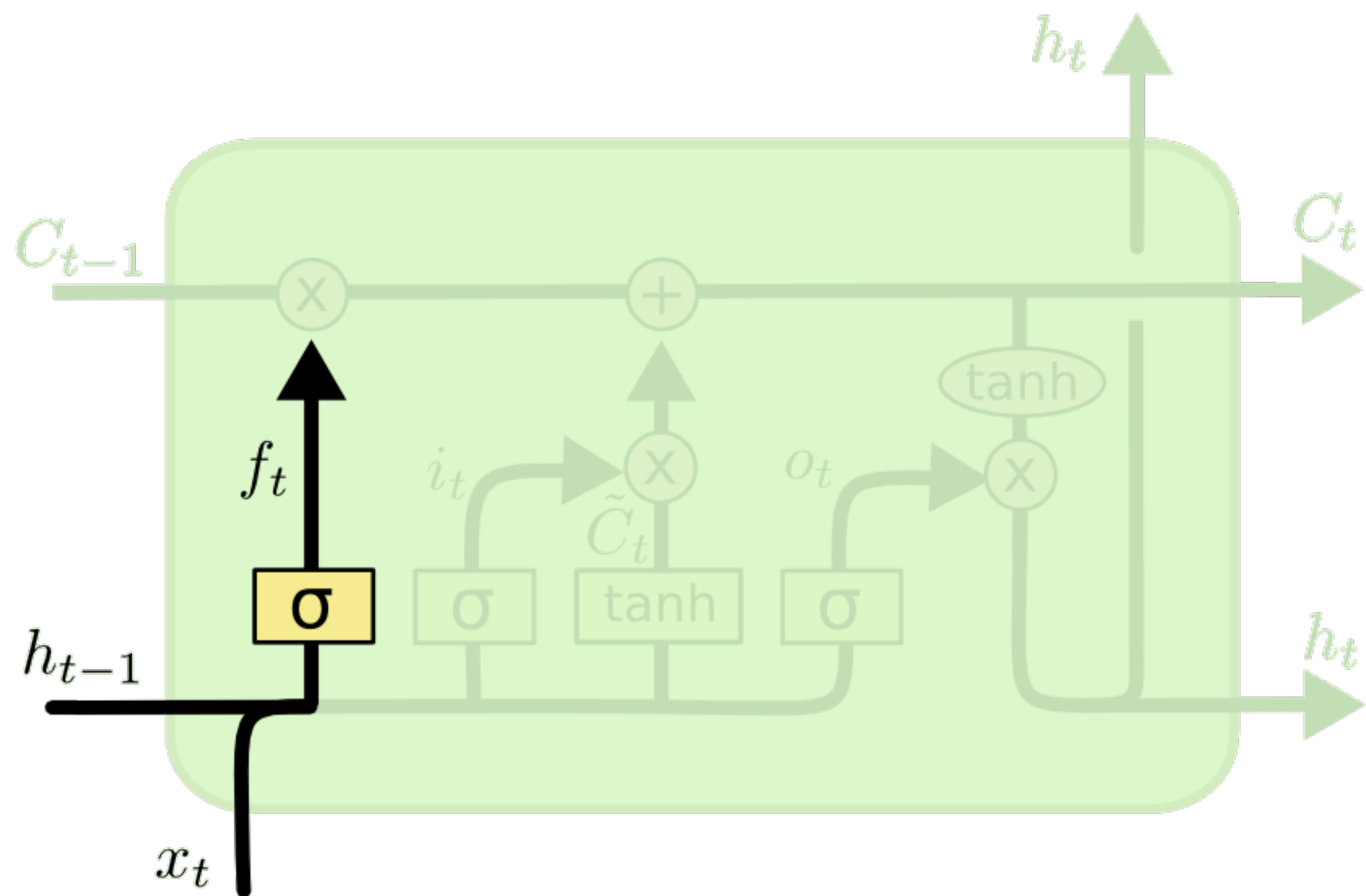




[Slide derived from Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>]



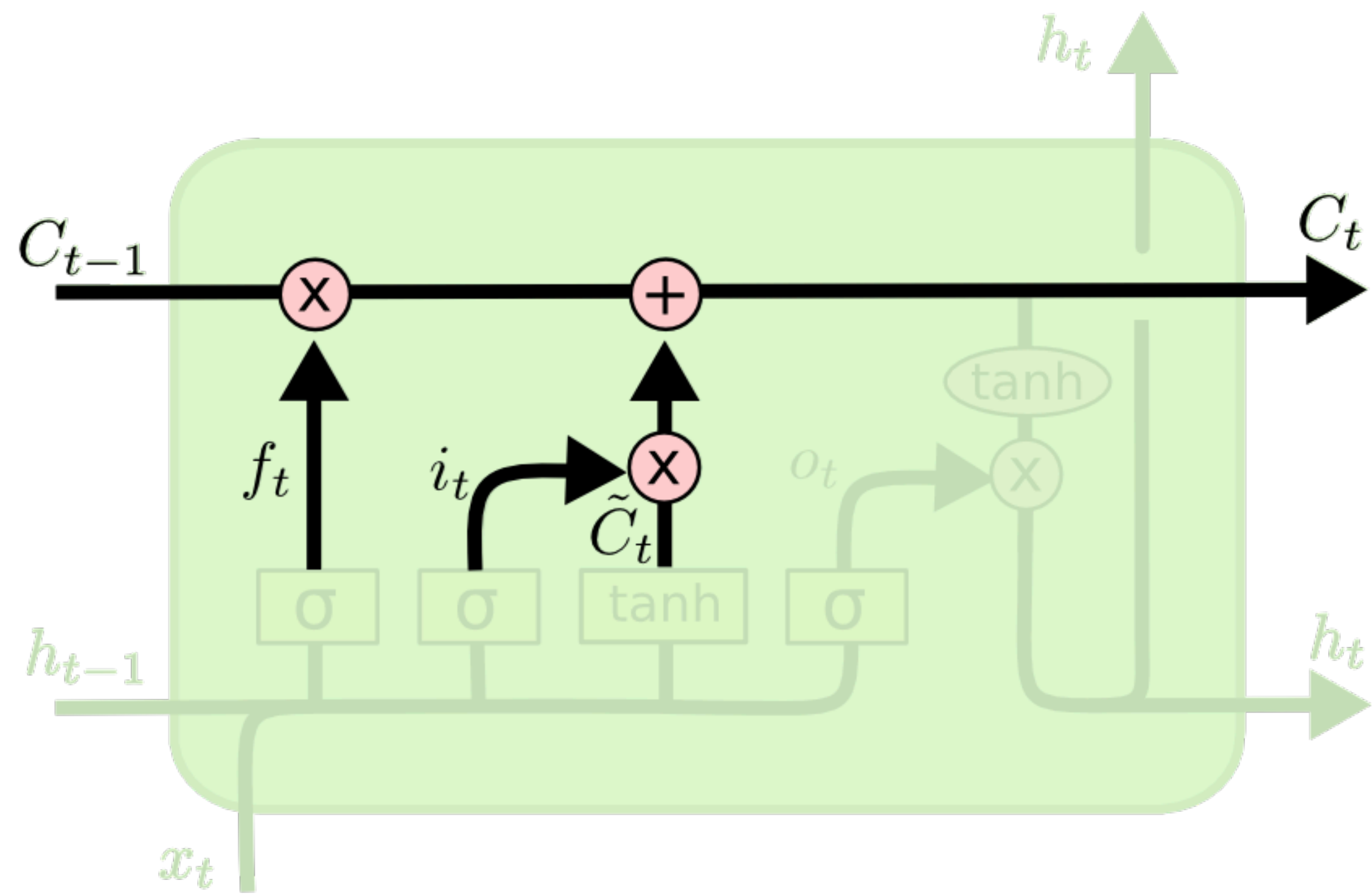
C_t = Cell state



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Decide what information to throw away from the cell state.

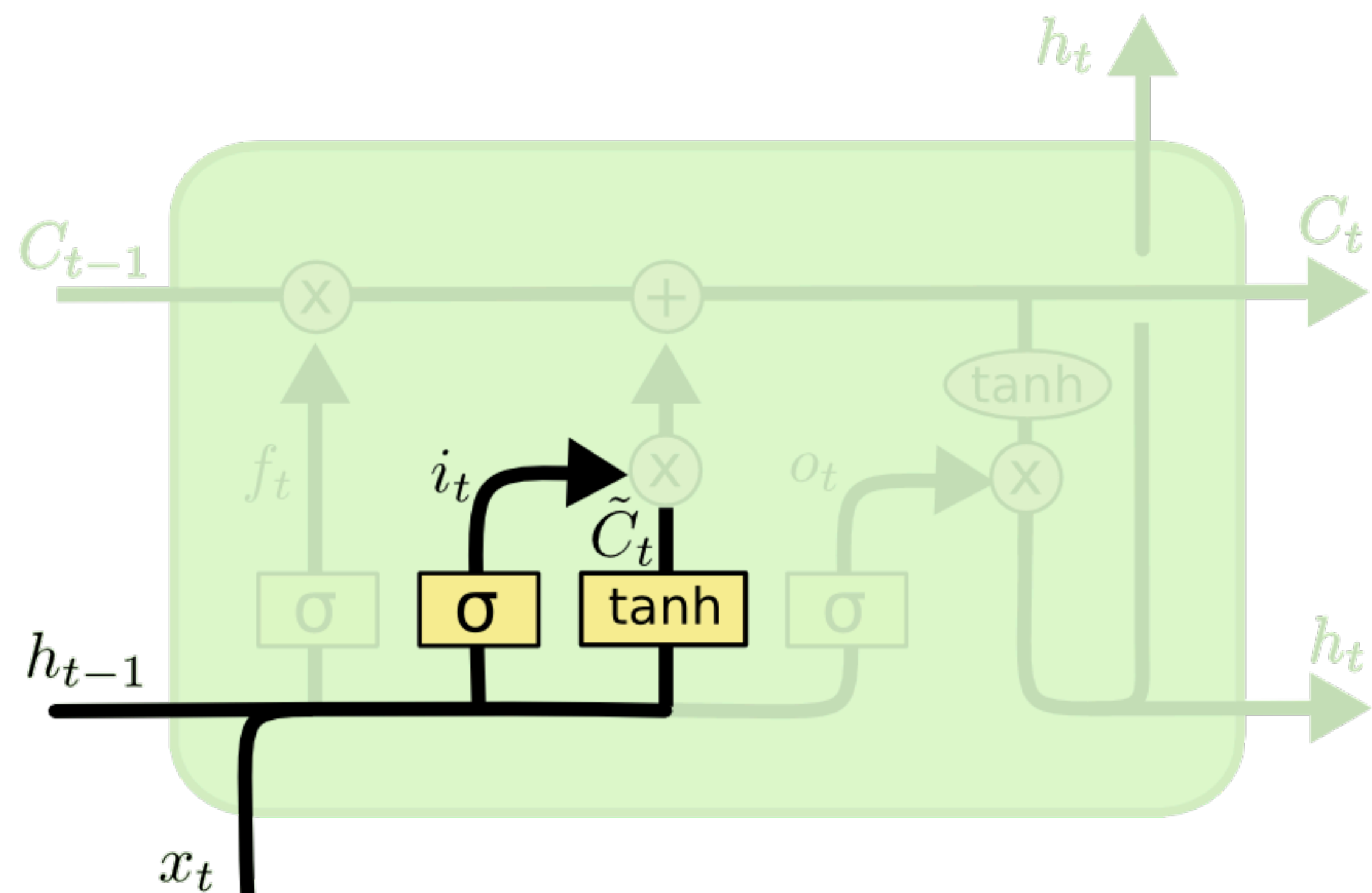
Forget gate: each element of cell state is multiplied by:
 ~ 1 (remember) or ~ 0 (forget).



Forget New values to write

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Forget selected old information, write selected new information.



which components to write to

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

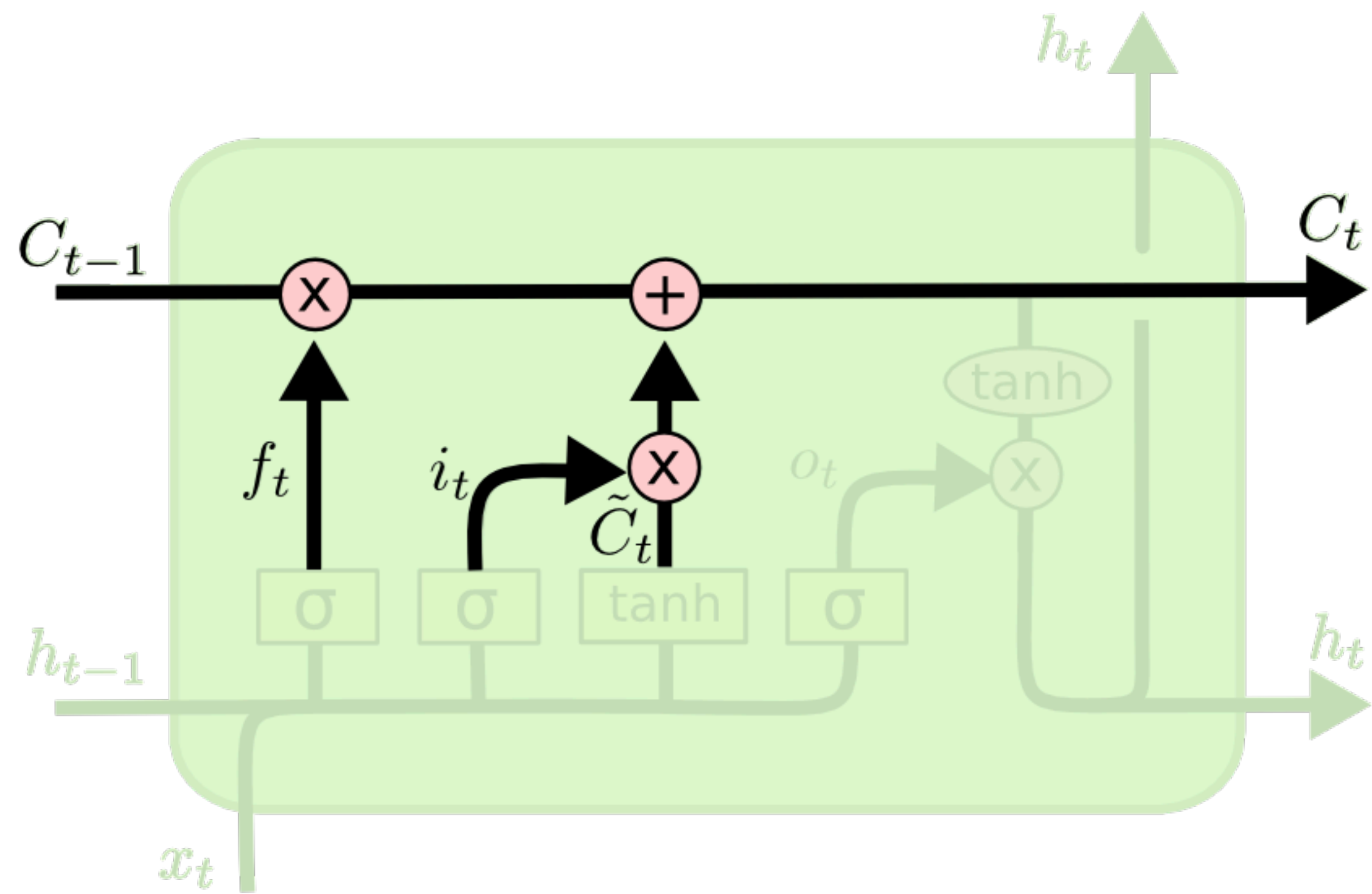
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

what to write into those components

Decide what new information to add to the cell state.

60

[Slide derived from Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>]

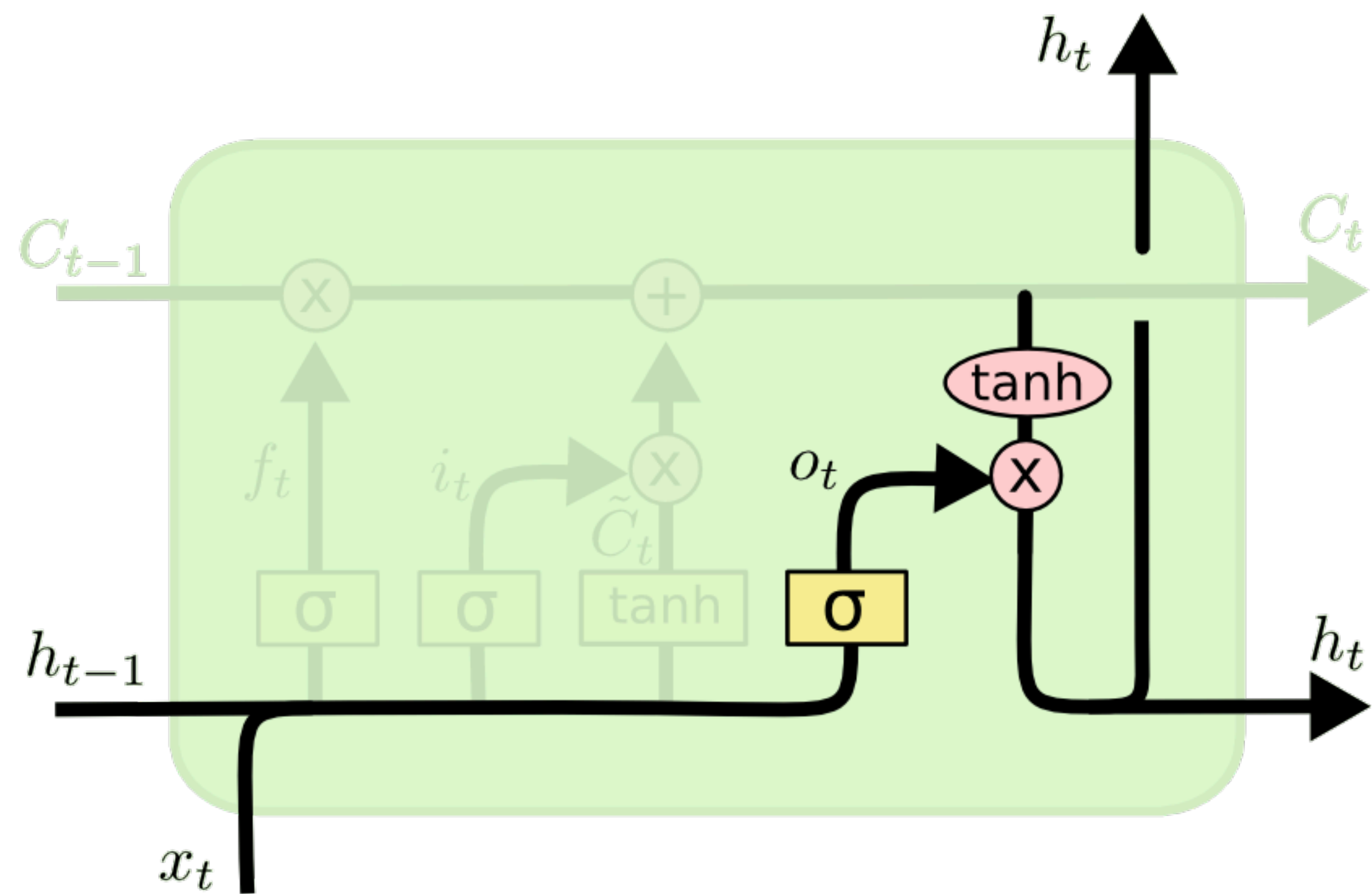


Forget Write new values

↓ ↓

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Forget selected old information, write selected new information.



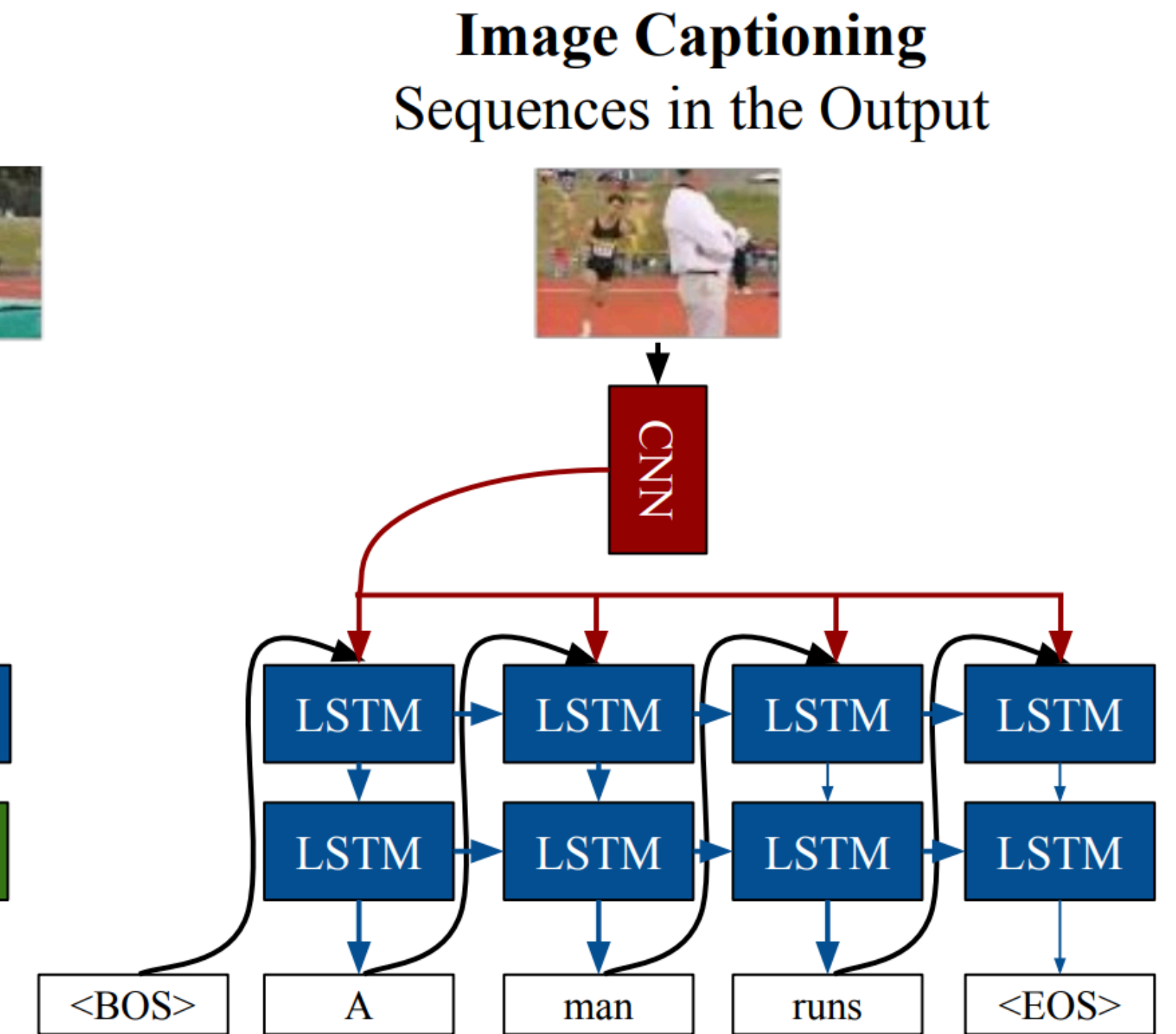
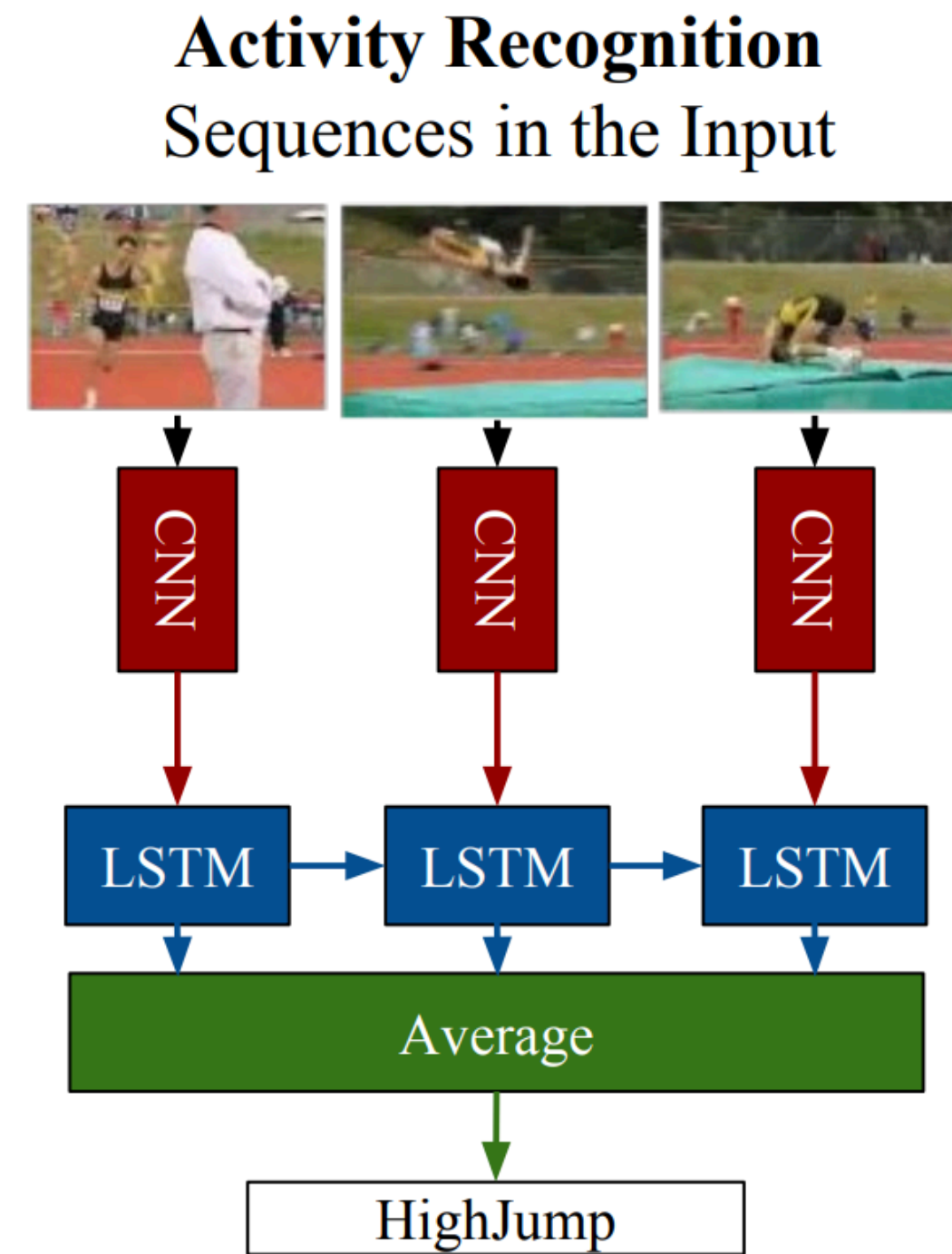
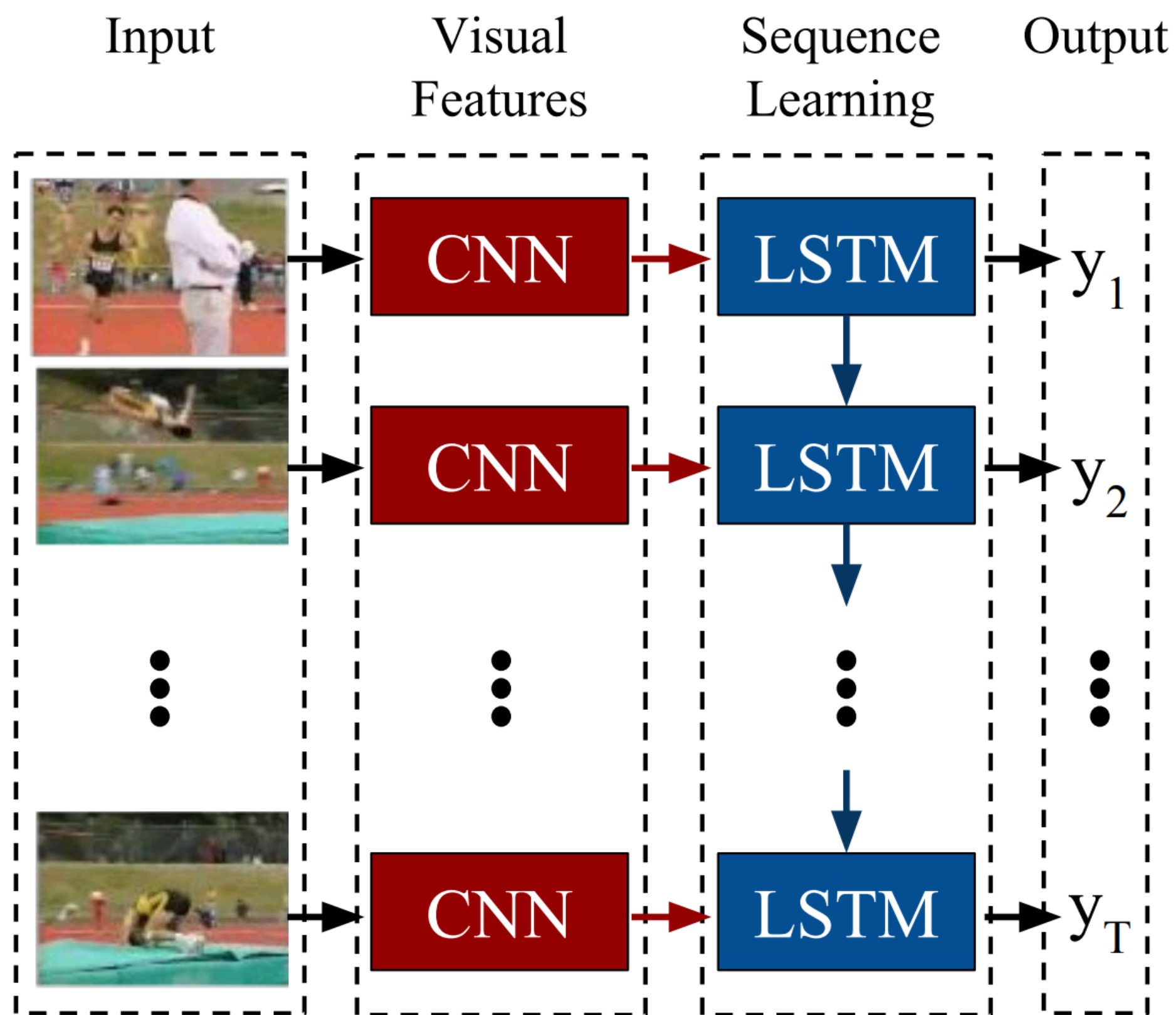
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

After having updated the cell state's information, decide what to output.

[Slide derived from Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>]

Some uses for RNNs



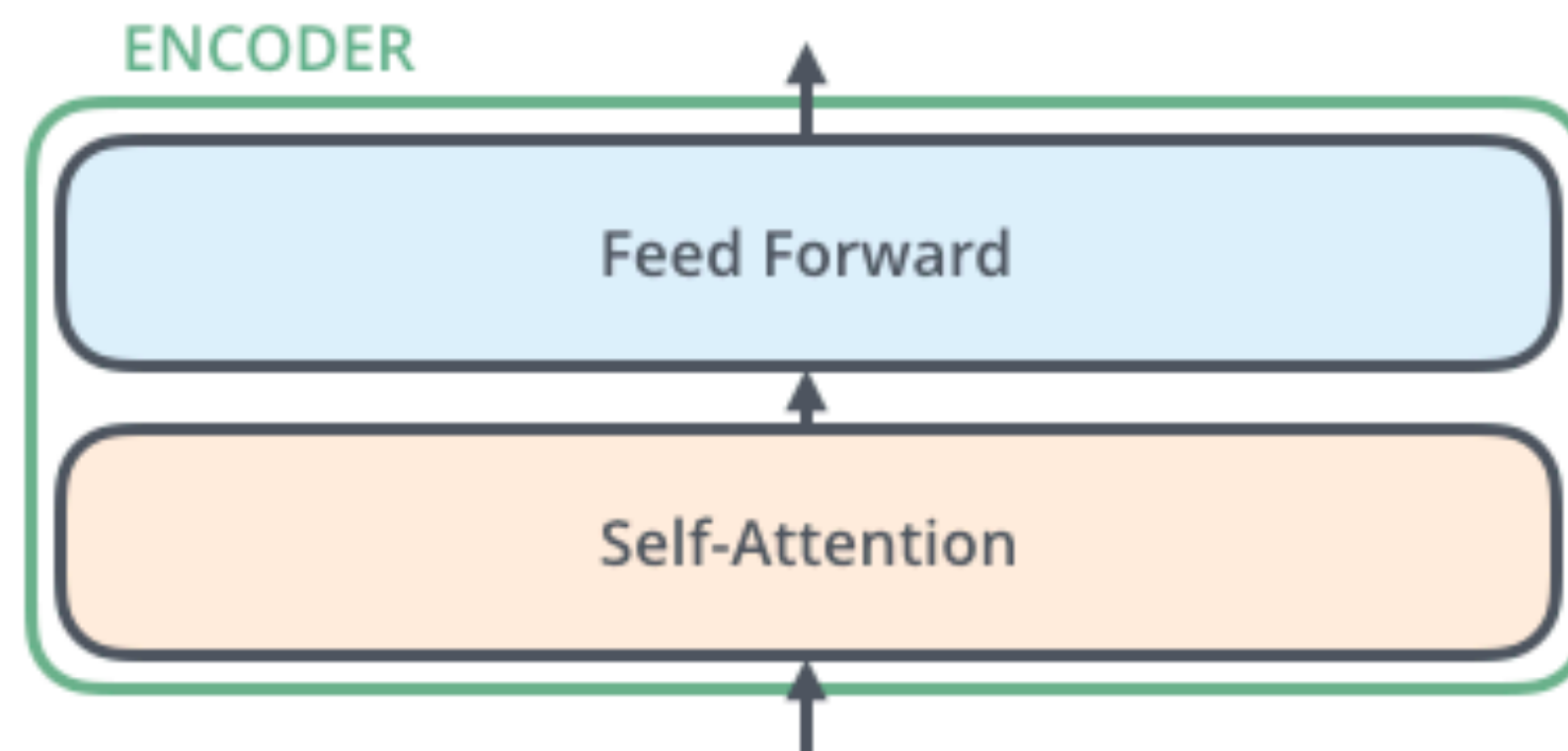
Problems with RNNs

- Hard to obtain motion information.
- Results in very deep networks (often slow, hard to train)
- Doesn't parallelize well
- Depth of network = length of the sequence

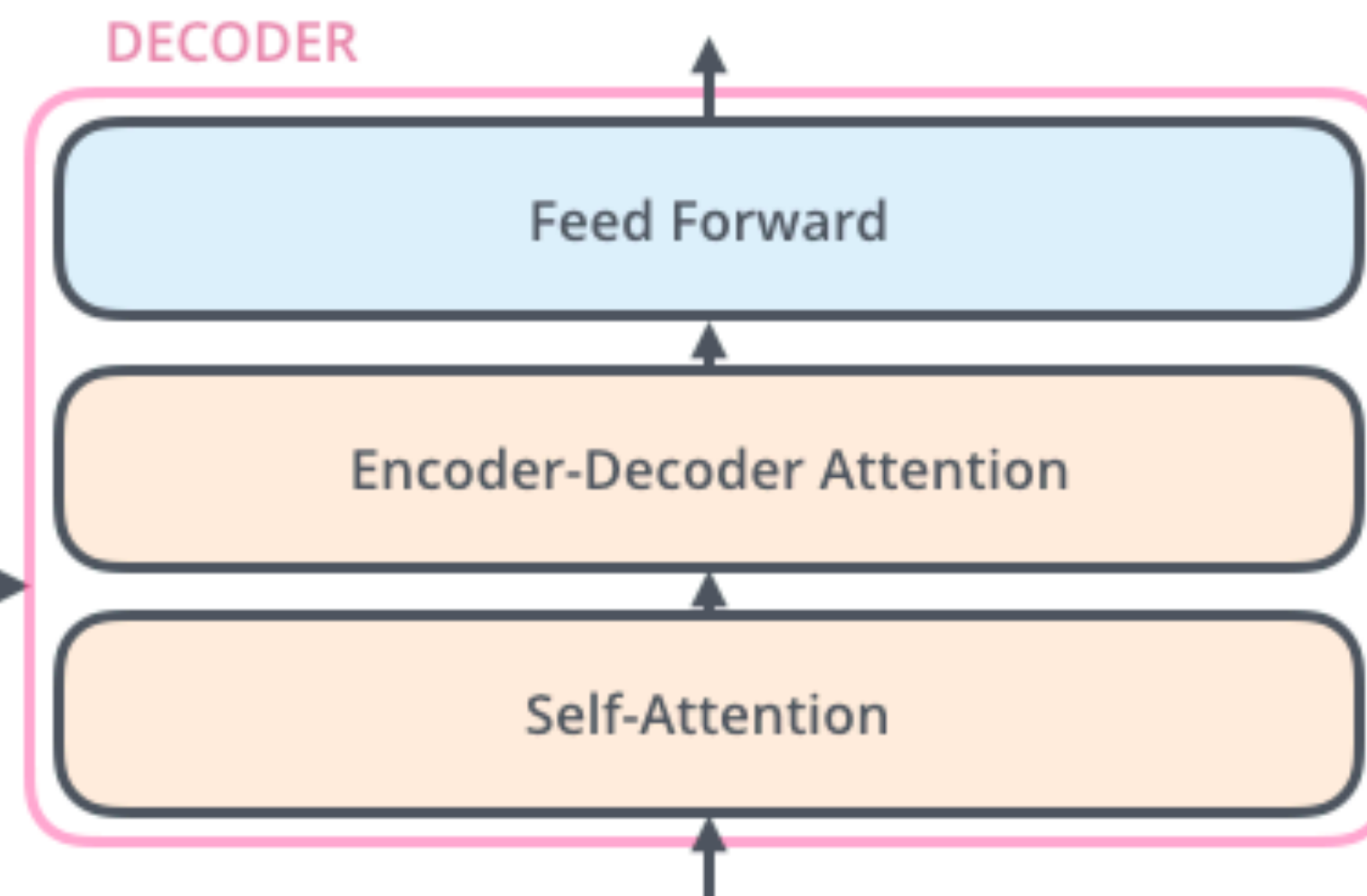
Basic transformer model

- Do we really need these sequence models?
- Sequence-to-sequence architecture using only point-wise processing and attention (no recurrent units or convolutions)

Encoder: receives entire input sequence and outputs encoded sequence of the same length



Decoder: predicts next token conditioned on encoder output and previously predicted tokens



A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, I. Polosukhin, [Attention is all you need](#), NeurIPS 2017

Self-attention

- Used to capture context within the sequence

The animal didn't cross the street because **it** was too tired .

As we are encoding “it”, we should focus on “the animal”

The animal didn't cross the street because **it** was too wide .

As we are encoding “it”, we should focus on “the street”

Self-attention layer

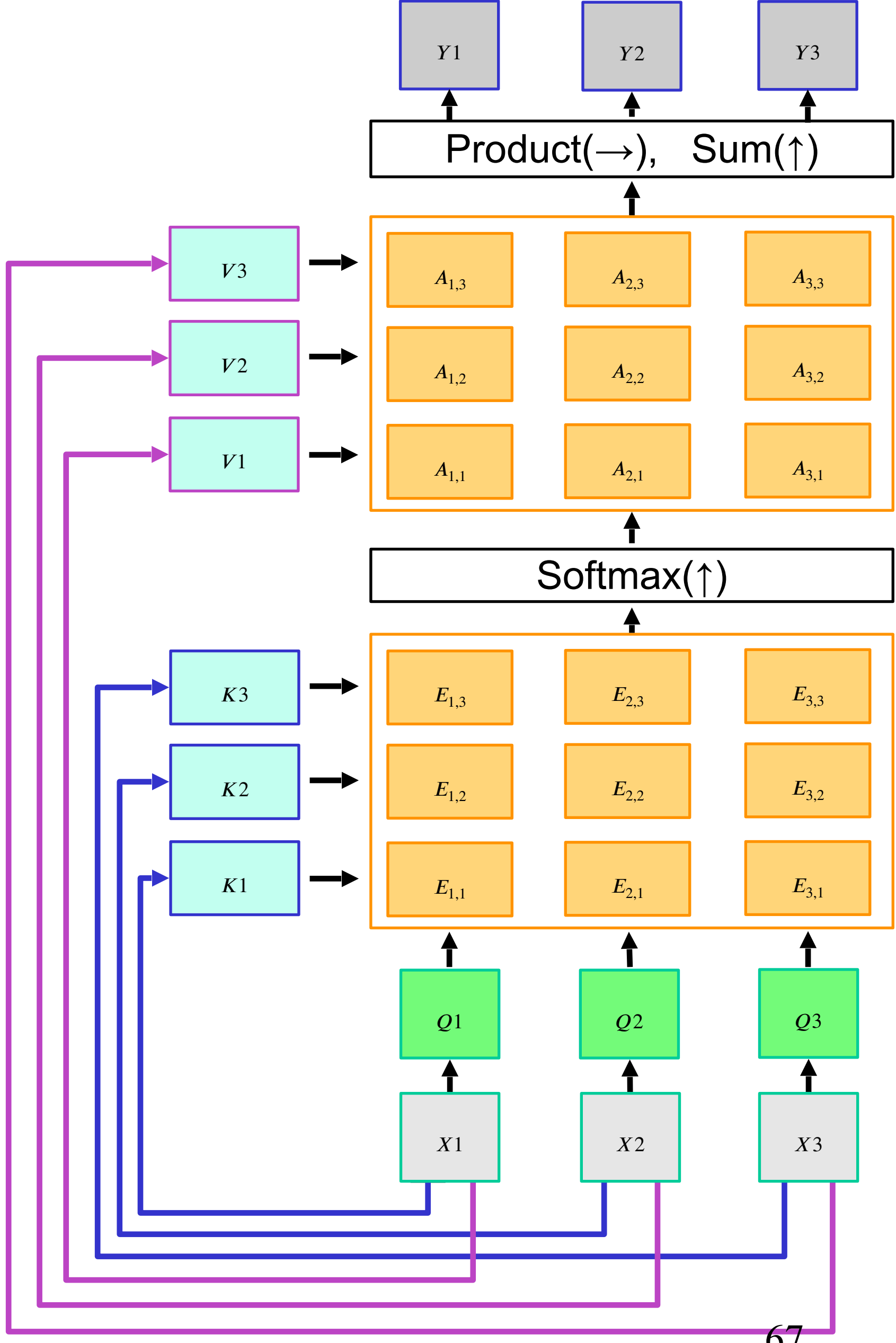
- Query vectors: $Q = XW_Q$
- Key vectors: $K = XW_K$
- Value vectors: $V = XW_V$
- Similarities: *scaled dot-product attention*

$$E_{i,j} = \frac{(Q_i \cdot K_j)}{\sqrt{D}} \quad \text{or} \quad E = QK^T / \sqrt{D}$$

(D is the dimensionality of the keys)

- Attn. weights: $A = \text{softmax}(E, \text{dim} = 1)$
- Output vectors:

$$Y_i = \sum_j A_{i,j} V_j \quad \text{or} \quad Y = AV$$

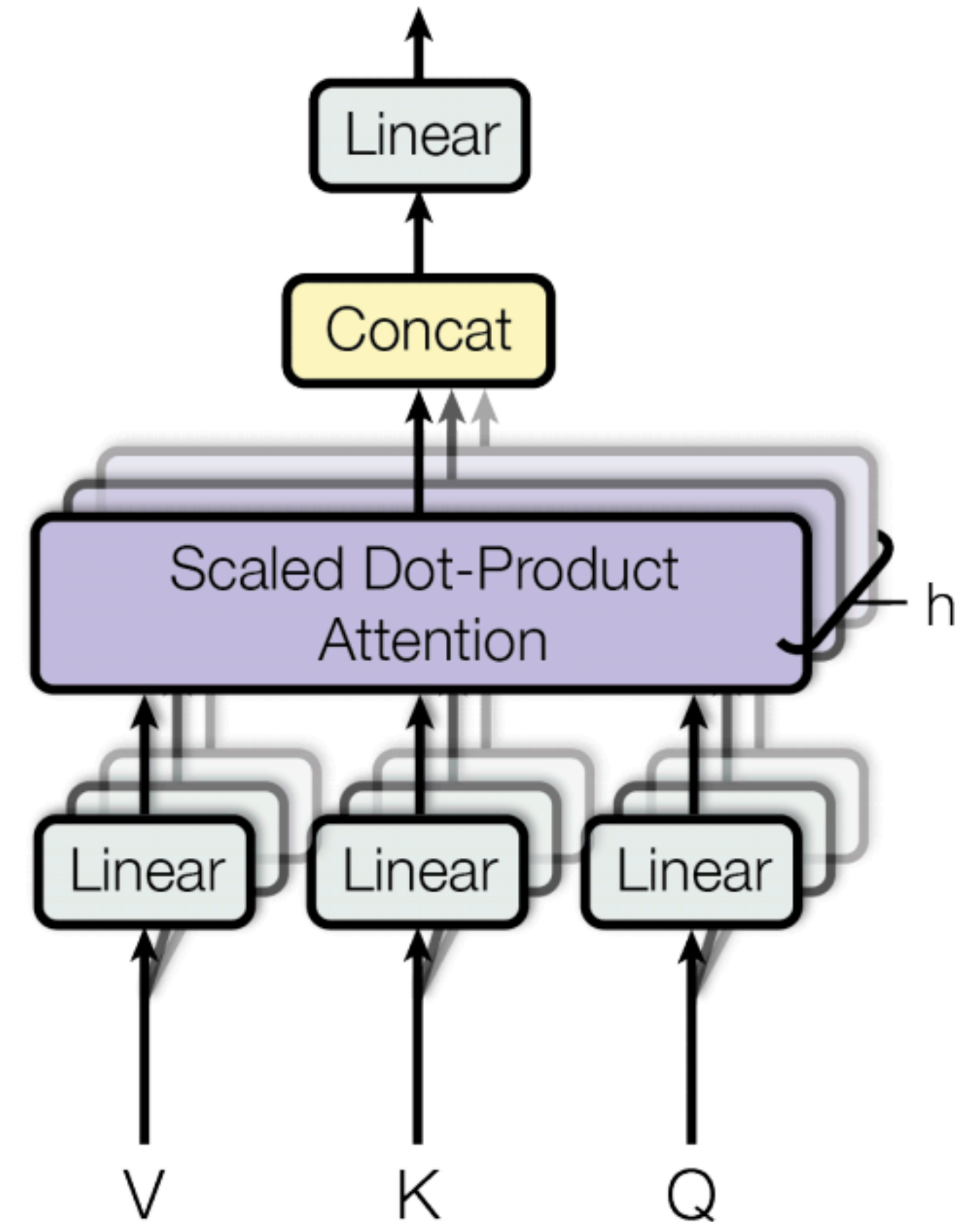


One query per input vector

Adapted from J. Johnson and S. Lazebnik.

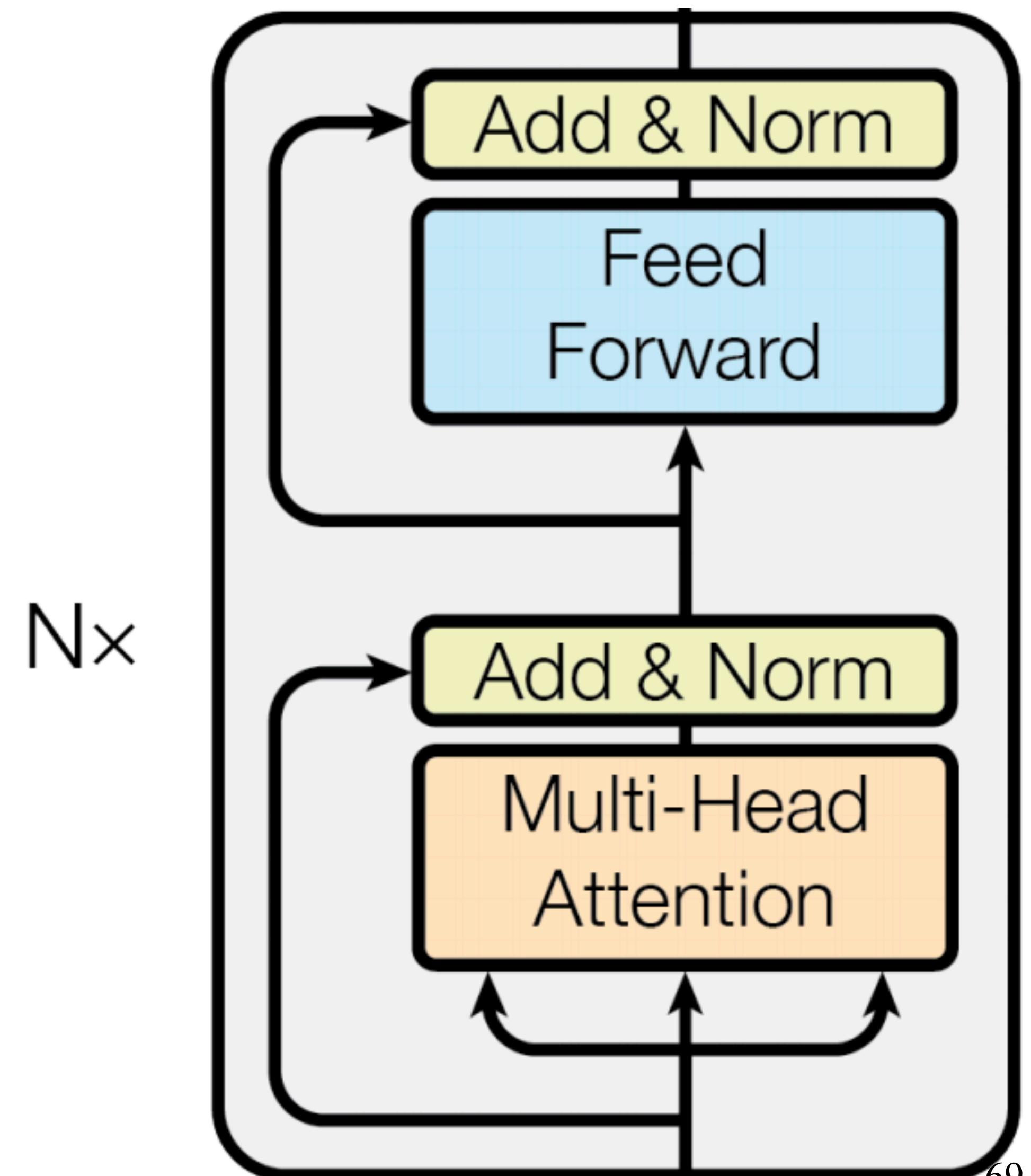
Multi-head attention

- Run h attention models in parallel on top of different linearly projected versions of Q , K , V ; concatenate and linearly project the results
- Intuition: enables model to attend to different kinds of information at different positions



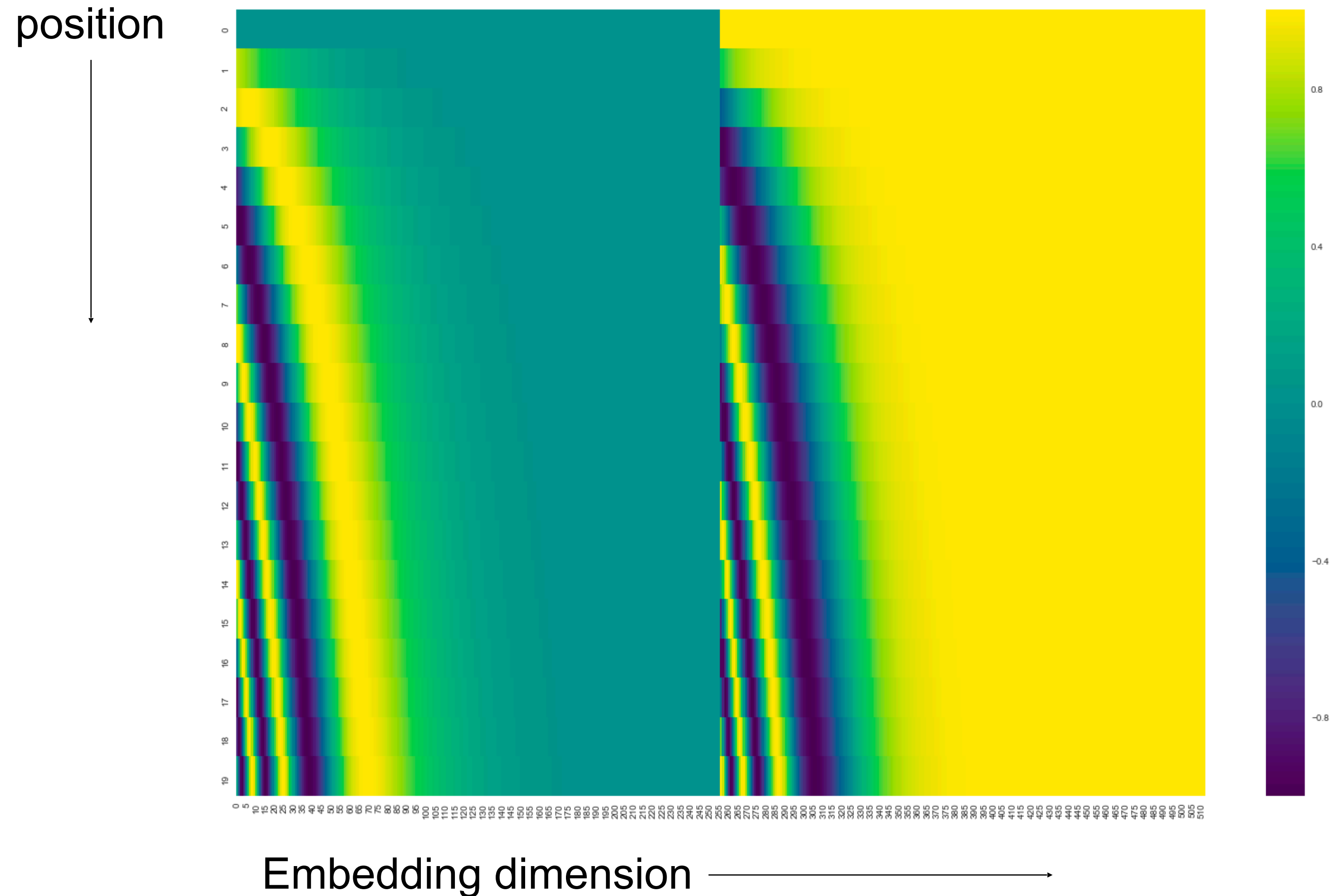
Transformer blocks

- A **Transformer** is a sequence of transformer blocks
- Vaswani et al.: $N=12$ blocks, embedding dimension = 512, 6 attention heads
- **Add & Norm**: residual connection followed by [layer normalization](#)
- **Feedforward**: two linear layers with ReLUs in between, applied independently to each vector
- Attention is the only interaction between inputs!

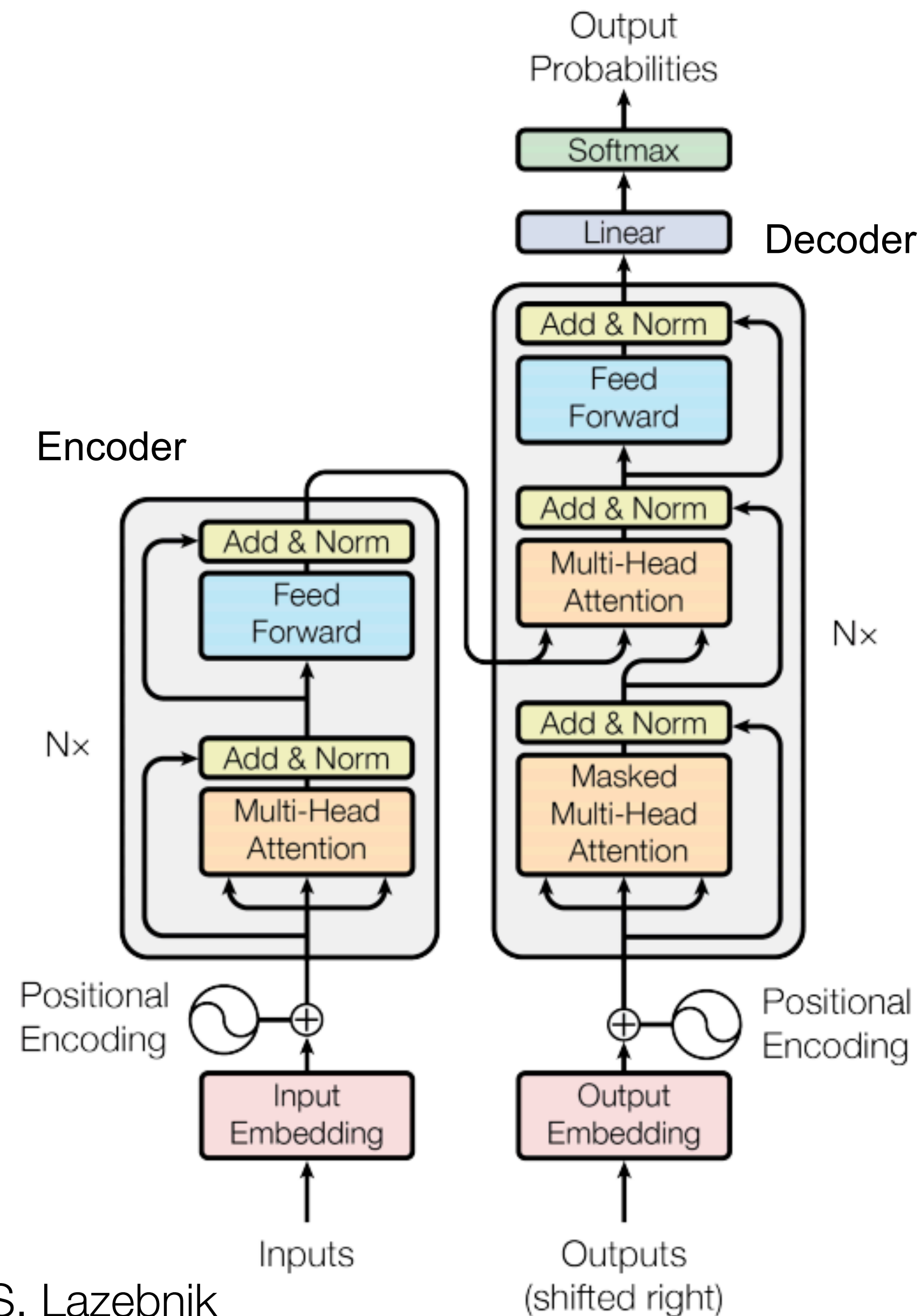


Positional encoding

- To give transformer information about ordering of tokens, add function of position (based on sines and cosines) to every input

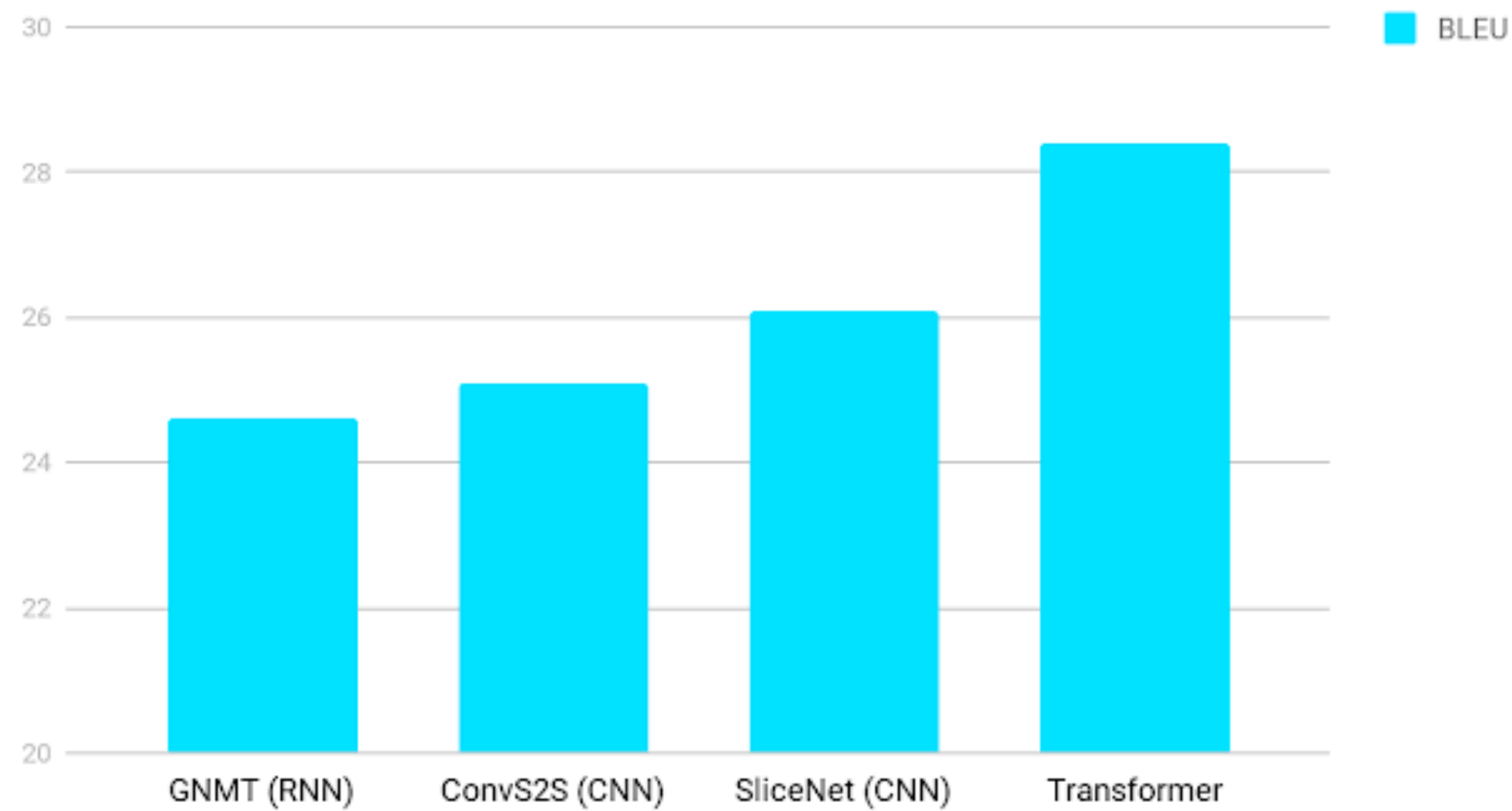


Transformer architecture: Zooming back out

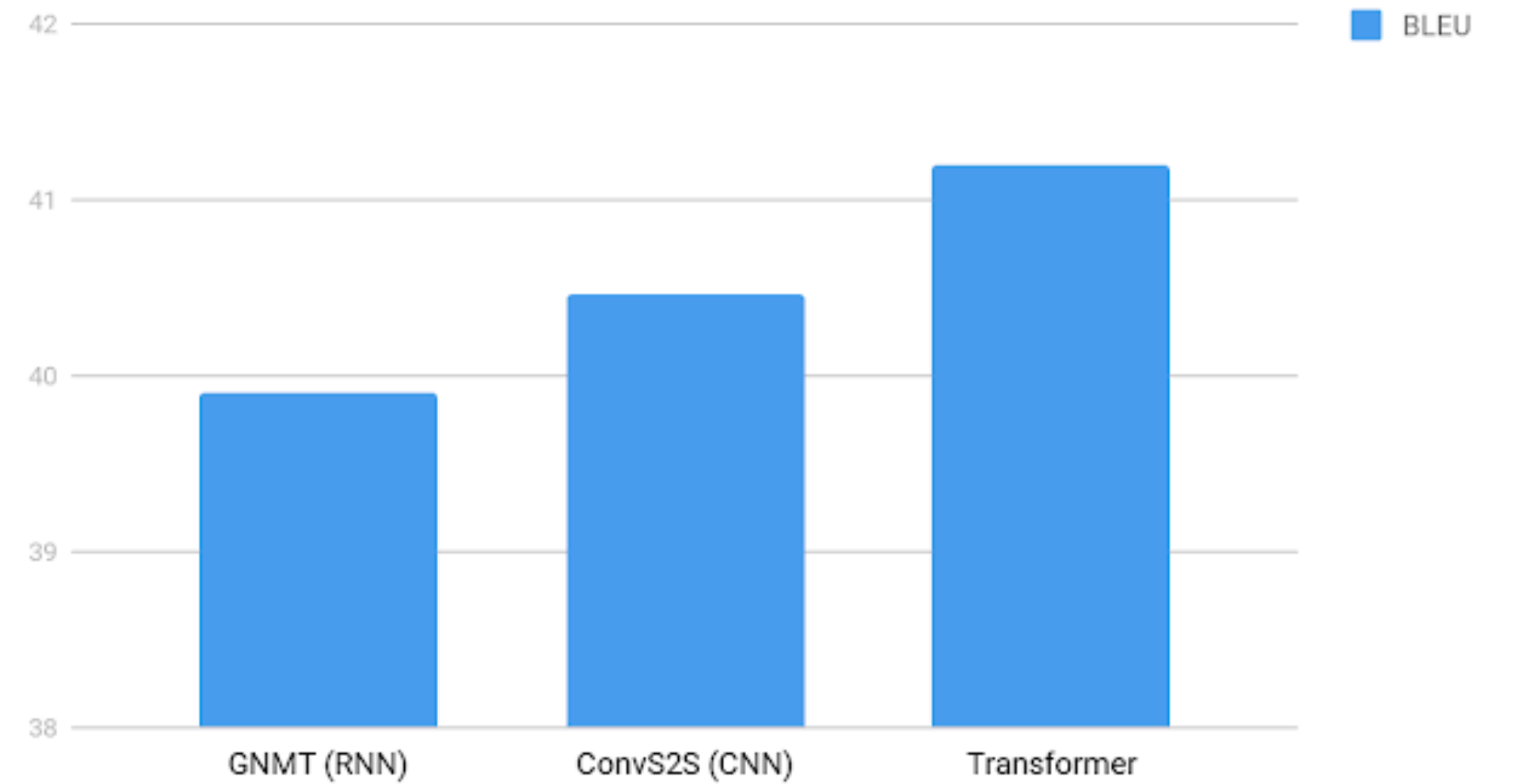


Language translation results

English German Translation quality



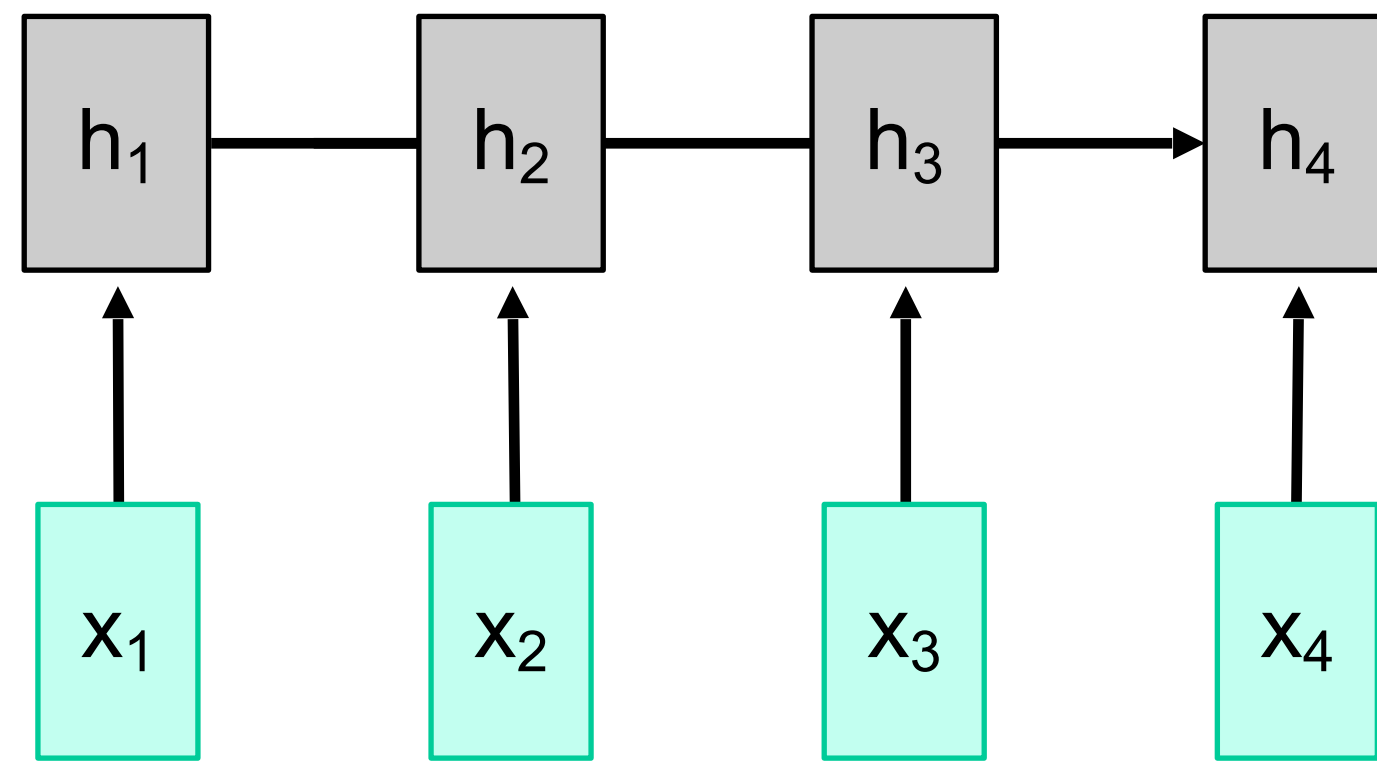
English French Translation Quality



<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

Different ways of processing sequences

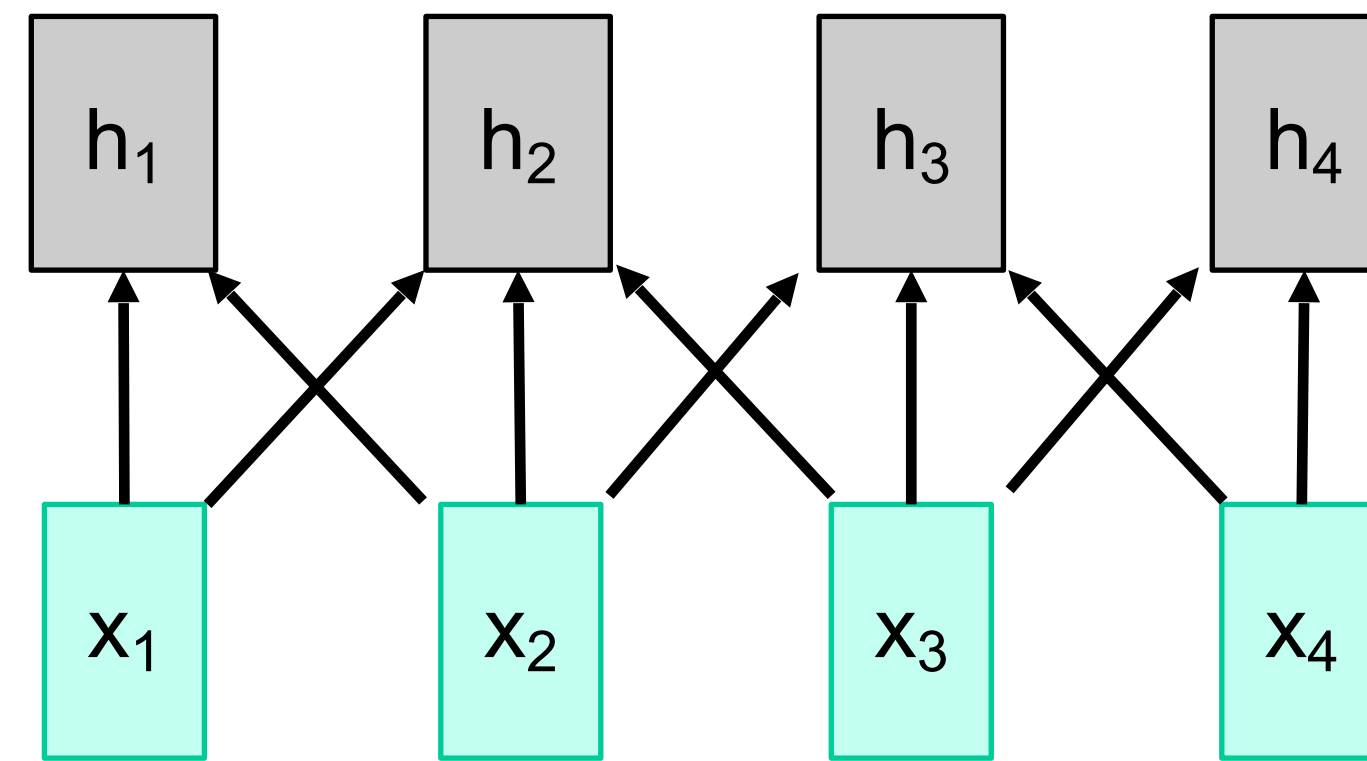
RNN



Works on **ordered sequences**

- Pros: Good for long sequences: After one RNN layer, h_T "sees" the whole sequence
- Con: Not parallelizable: need to compute hidden states sequentially. Very deep.
- Con: Hidden states have limited expressive capacity

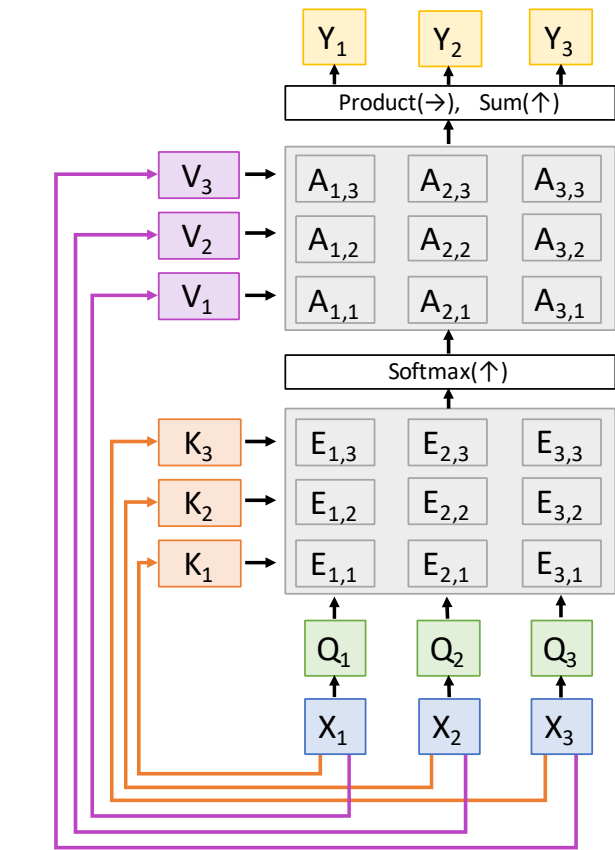
Convolutional network



Works on **multidimensional grids**

- Pro: Each output can be computed in parallel (at training time)
- Con: Bad at long sequences: Need to stack many conv layers for outputs to "see" the whole sequence

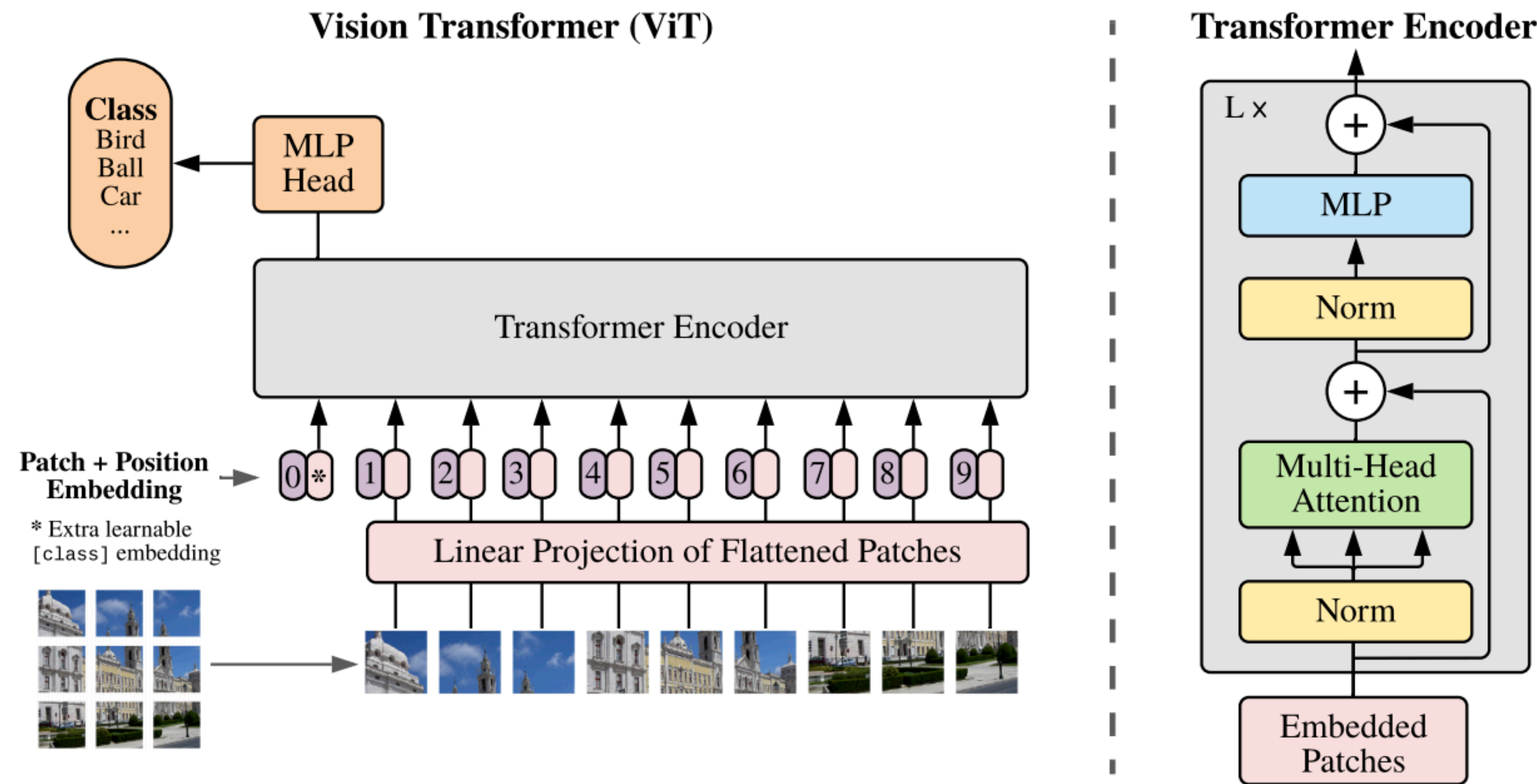
Transformer



- Works on **sets of vectors**
- Pro: Good at long sequences: after one self-attention layer, each output "sees" all inputs!
- Pro: Each output can be computed in parallel (at training time)
- Con: Memory-intensive. $O(N^2)$ without modifications.

Preview for later in the course: Vision transformer (ViT)

- Split an image into patches, feed linearly projected patches into standard transformer encoder
 - With patches of 14x14 pixels, you need $16 \times 16 = 256$ patches to represent 224×224 images



A. Dosovitskiy et al. [An image is worth 16x16 words: Transformers for image recognition at scale](#). ICLR 2021

Next class: representation learning