

Lecture 12: Learning temporal models

Announcements

- PS1 regrade requests due today!
- PS2 grades out. Regrades due next Tues.
- PS6 out.
- Project proposal will be next “assignment”.
We’ll give you more info next week

Project

- Open-ended! Example projects:
 - Implement and extend a recent computer vision paper
 - Use computer vision in your research
 - We'll also provide a list of project ideas
- Work in small groups (up to 4 people)
- We'll allow updates to project proposal if you change ideas

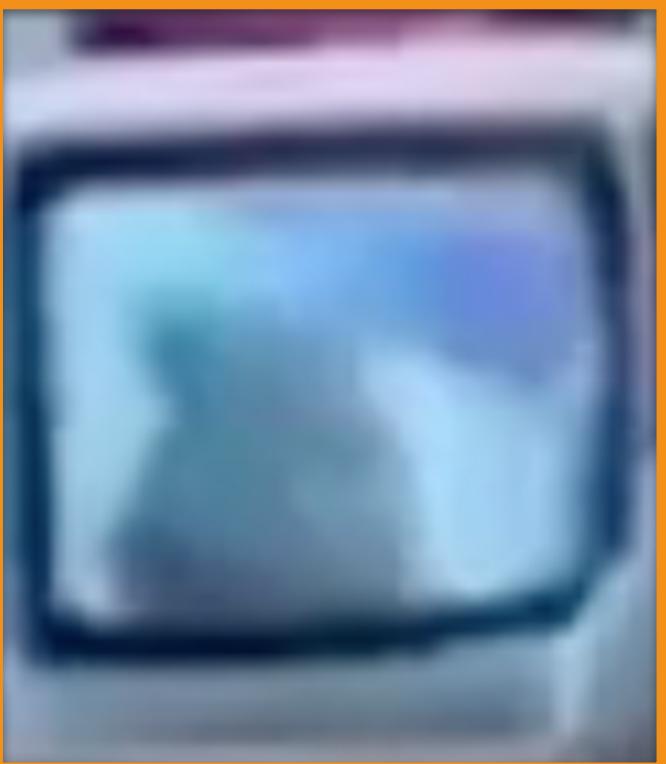


Source: Torralba,
Freeman, Isola



kindergarten classroom

Source: Torralba,
Freeman, Isola



television

person



chair

“What color is the chair?”



Source: Torralba,
Freeman, Isola

“What color is the chair?”
red



Source: Torralba,
Freeman, Isola



“What will the girl do next?”

Source: Torralba,
Freeman, Isola



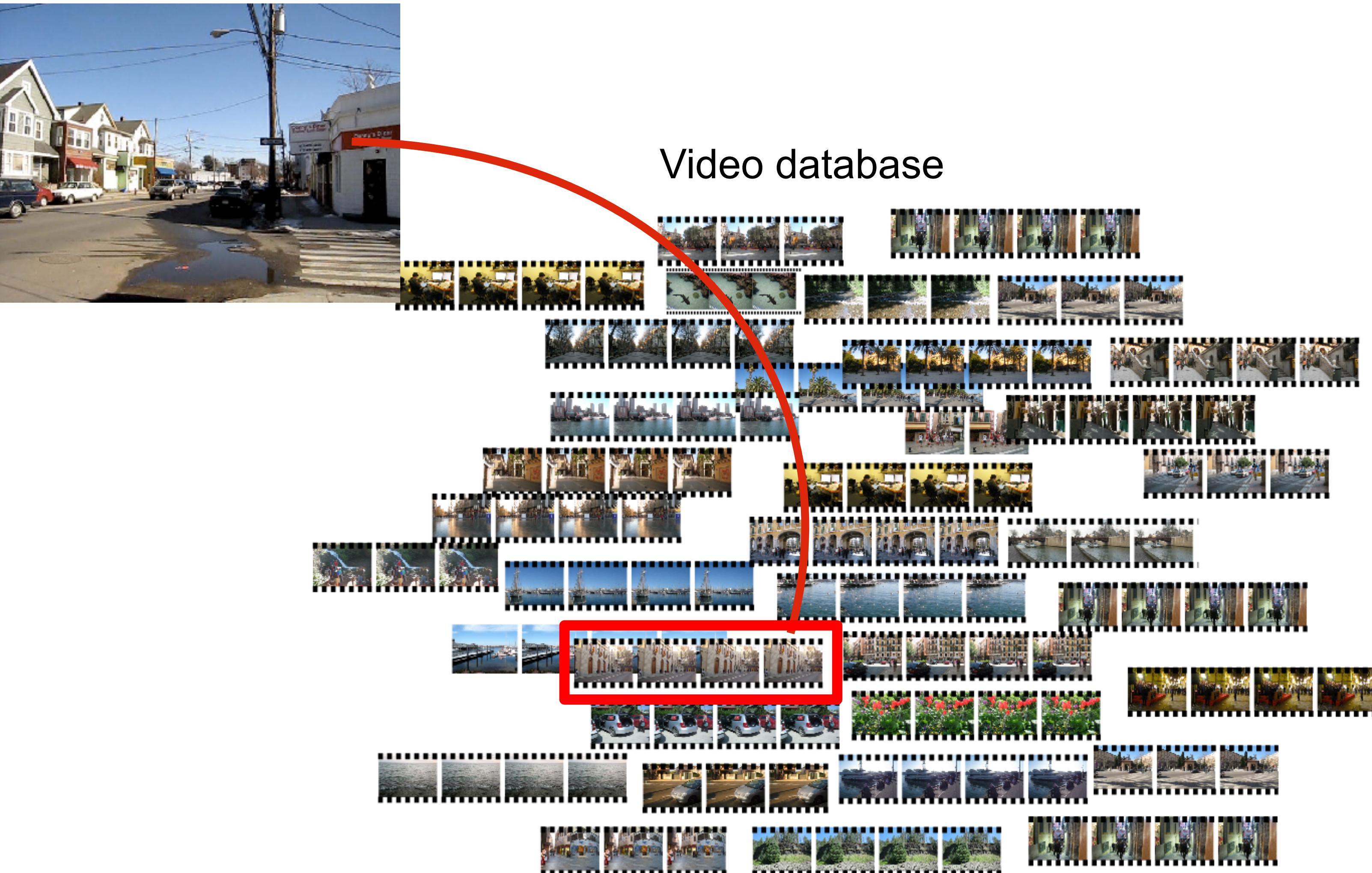
Event prediction

What can happen here?



Event prediction

What can happen here?

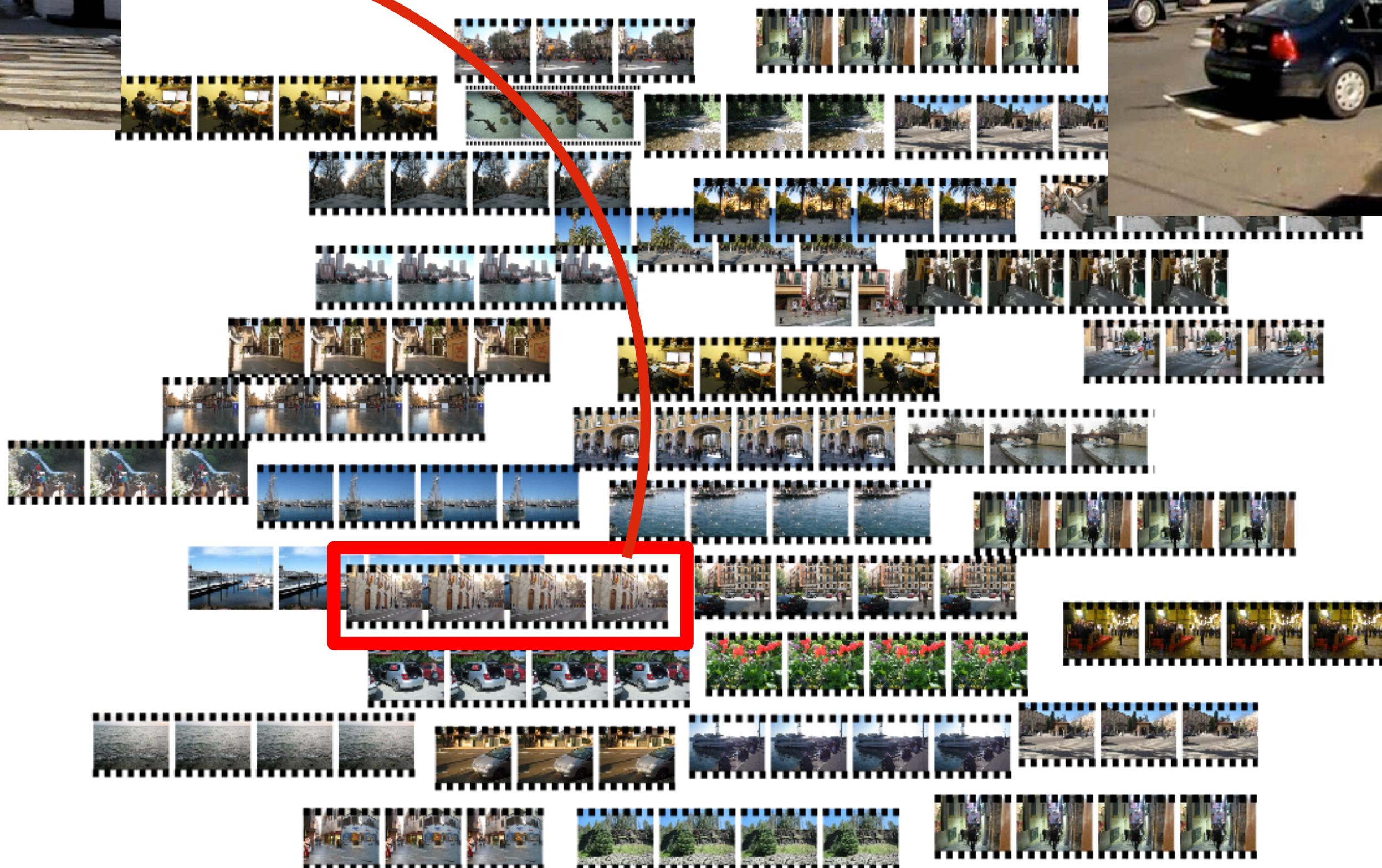


Event prediction

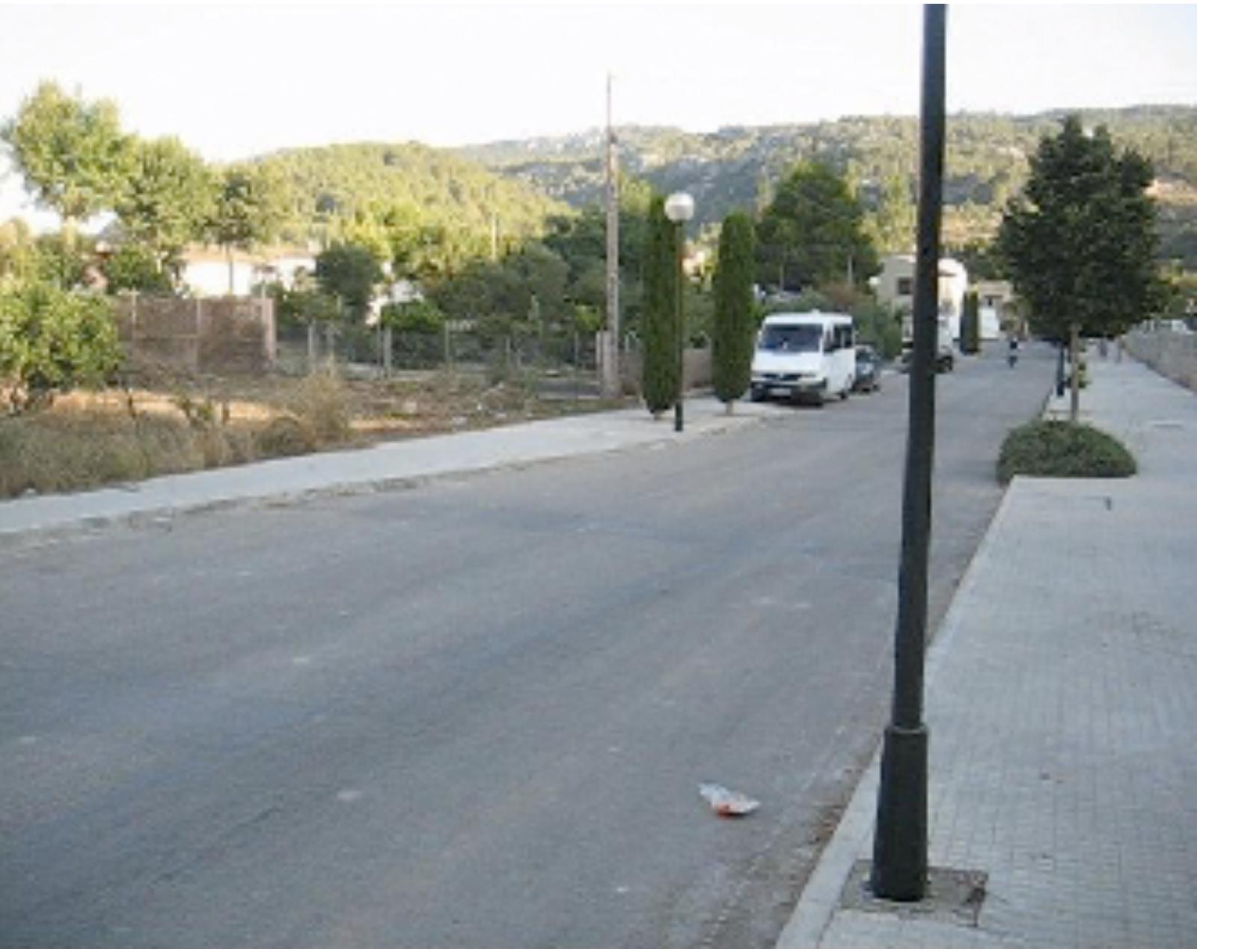
What can happen here?



Video database



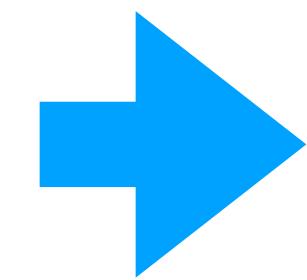
What can happen here?



What can happen here?

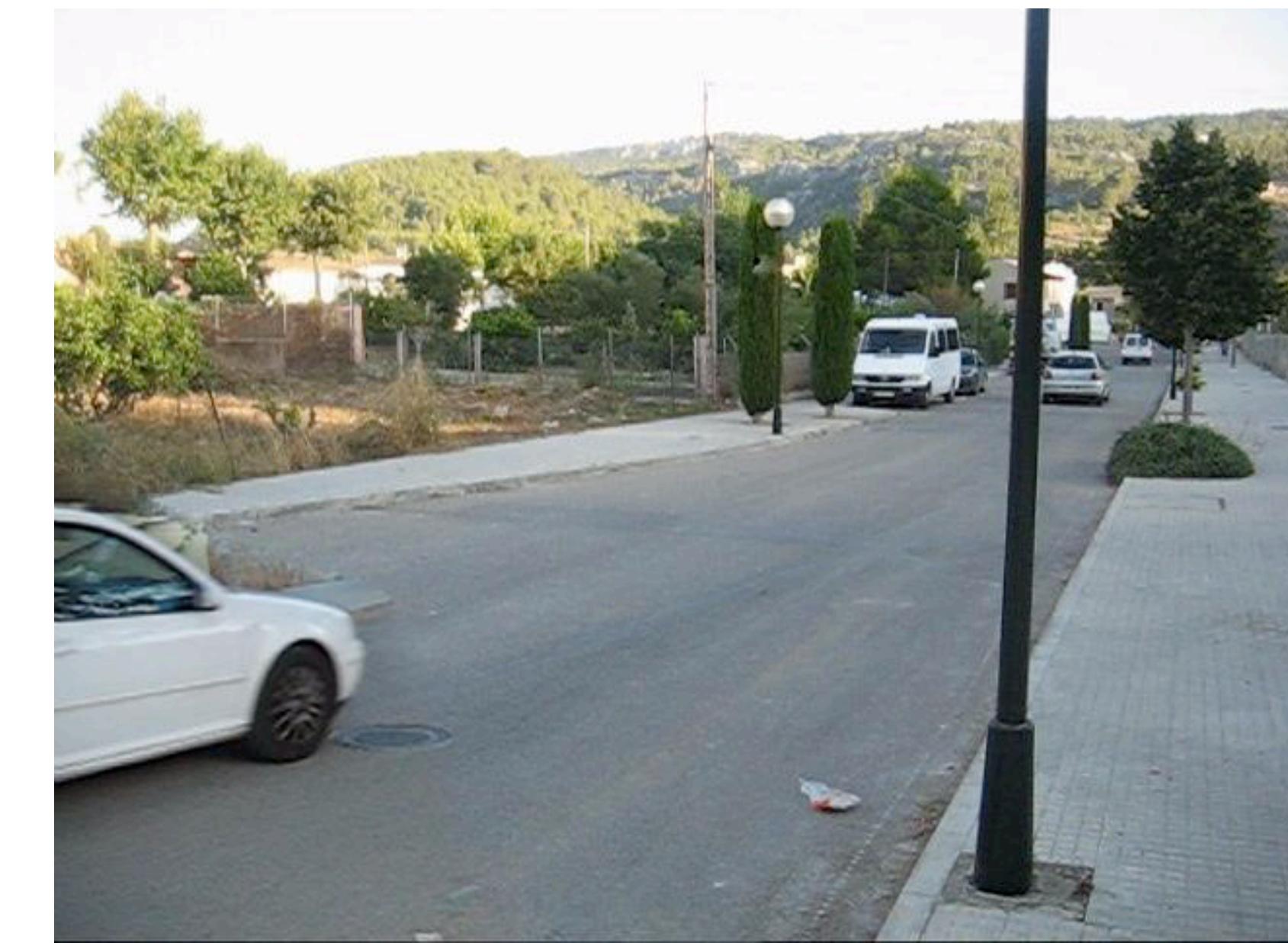
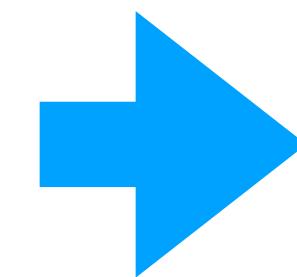


Prediction



What can happen here?

Prediction



Nearest neighbor



16

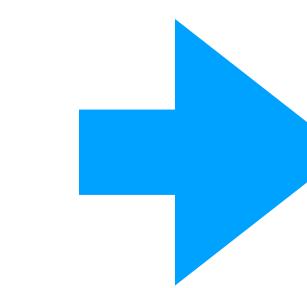
What can happen here?



What can happen here?



Prediction



Action recognition



Making latte art



Jaywalking

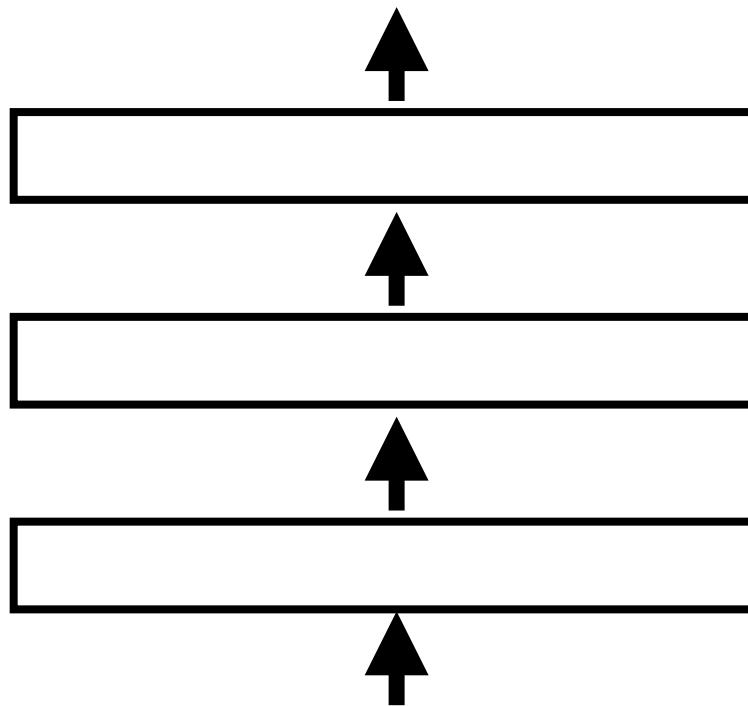


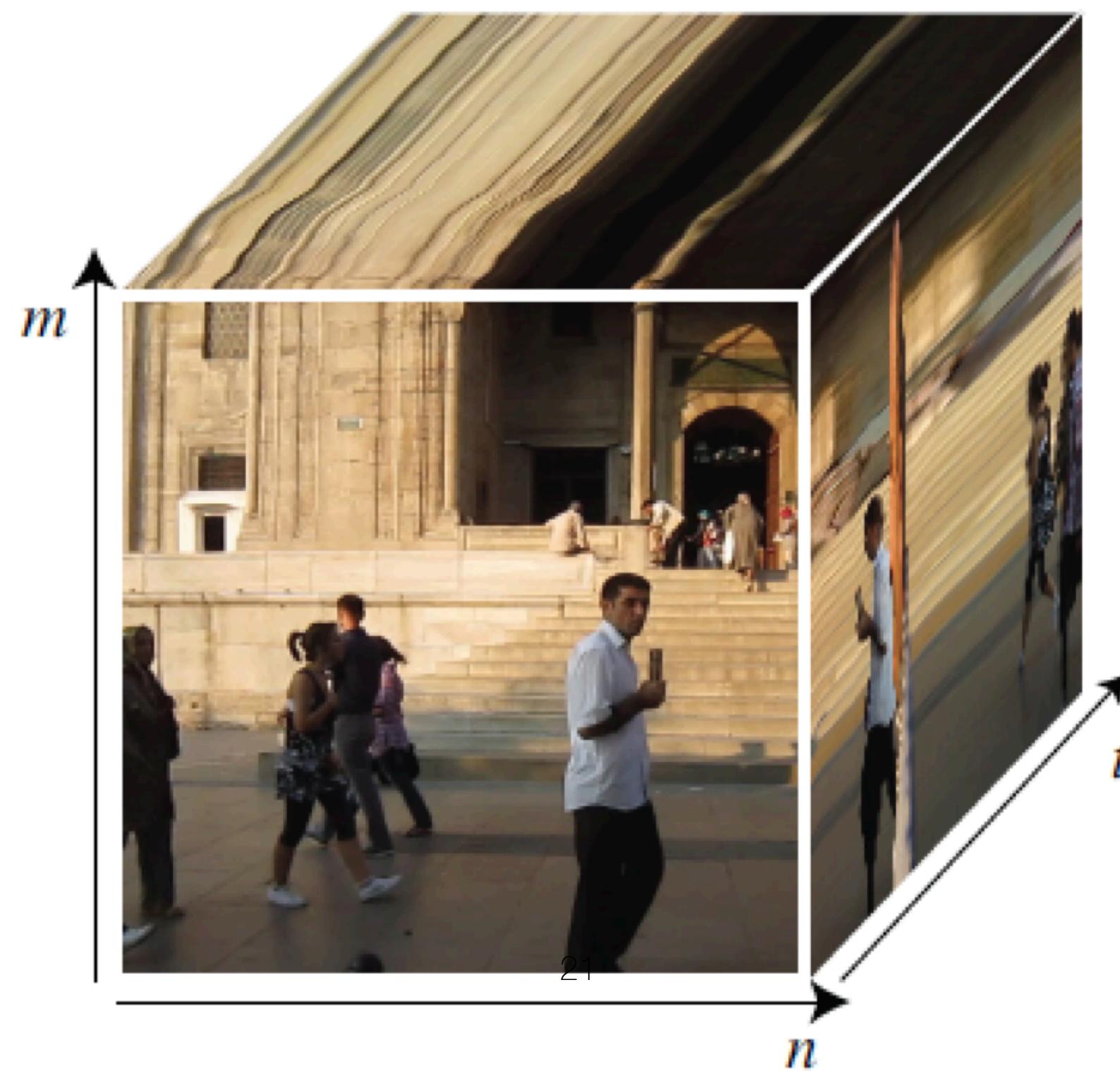
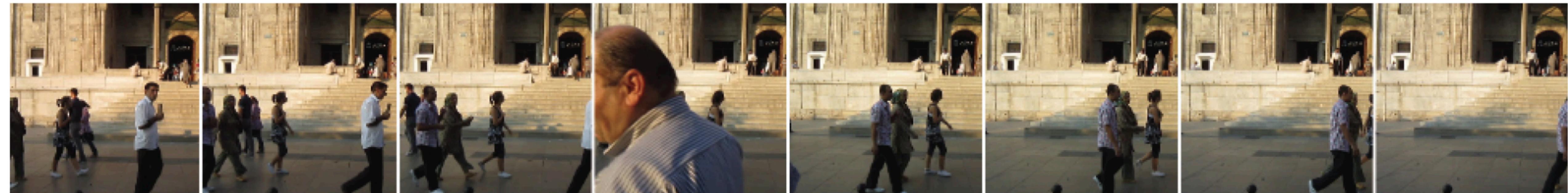
Grooming dog

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

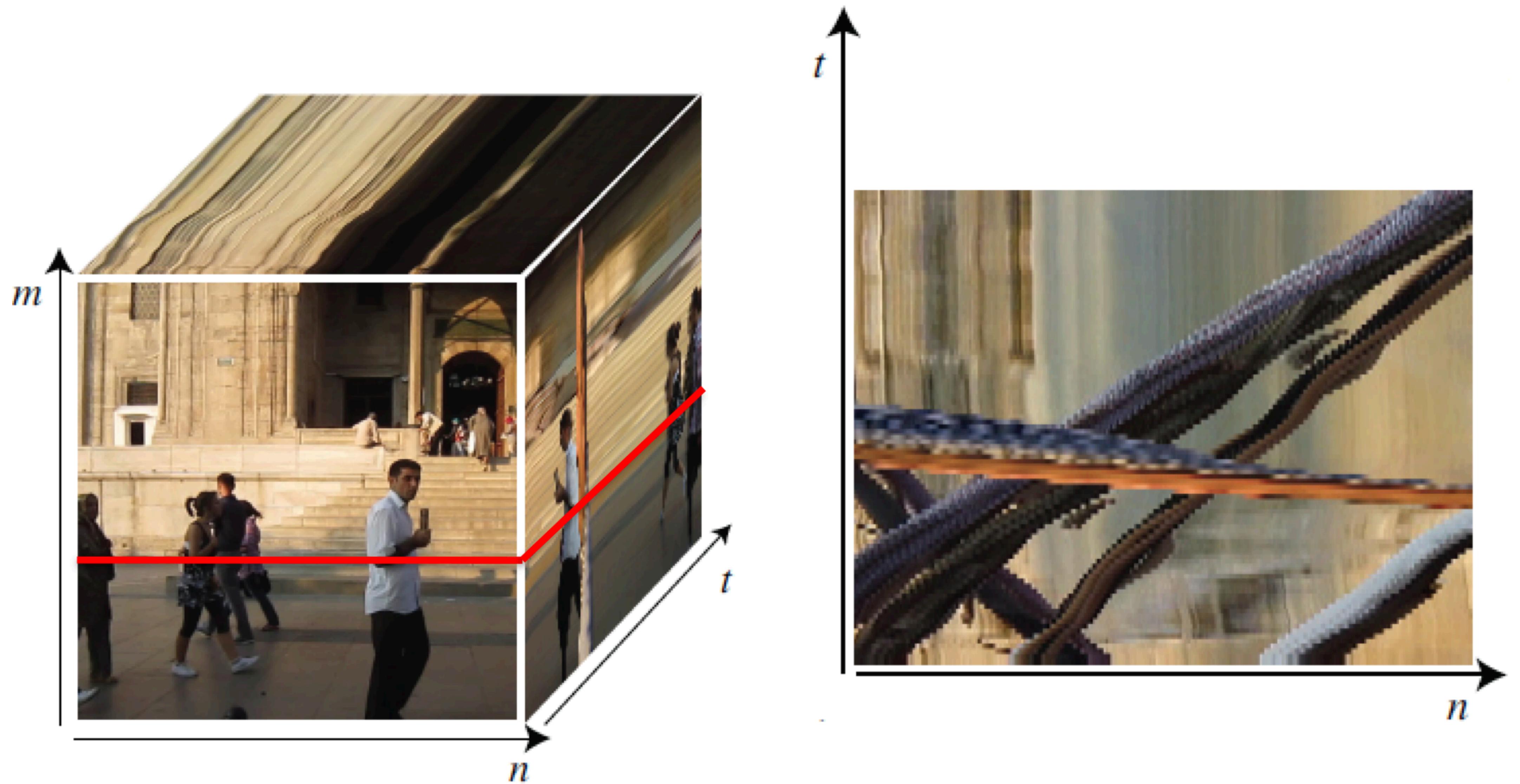
Simple model: averaging

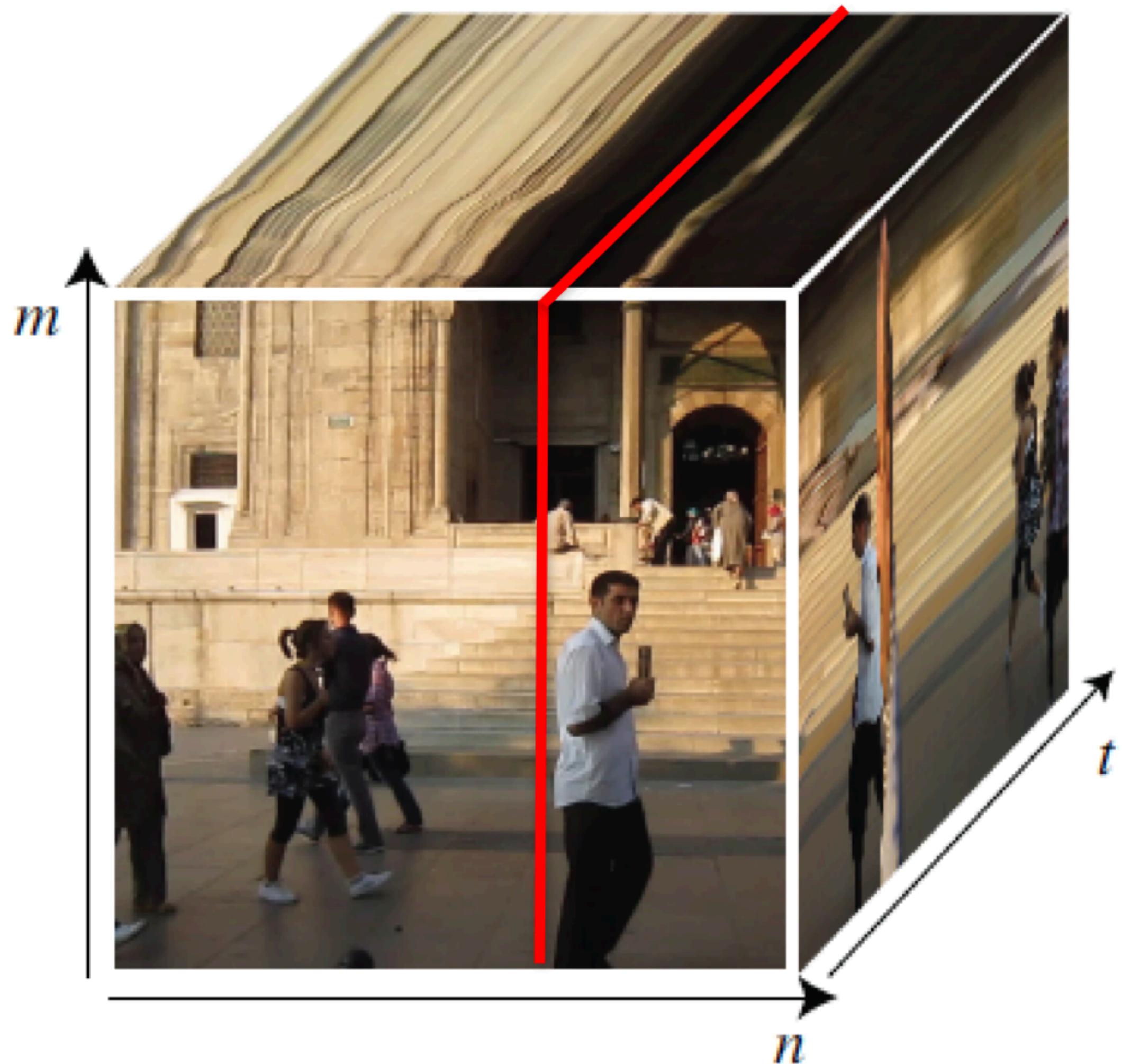
$$\frac{1}{3} p(\text{making latte art} \mid I_1) + \frac{1}{3} p(\text{making latte art} \mid I_2) + \frac{1}{3} p(\text{making latte art} \mid I_3)$$



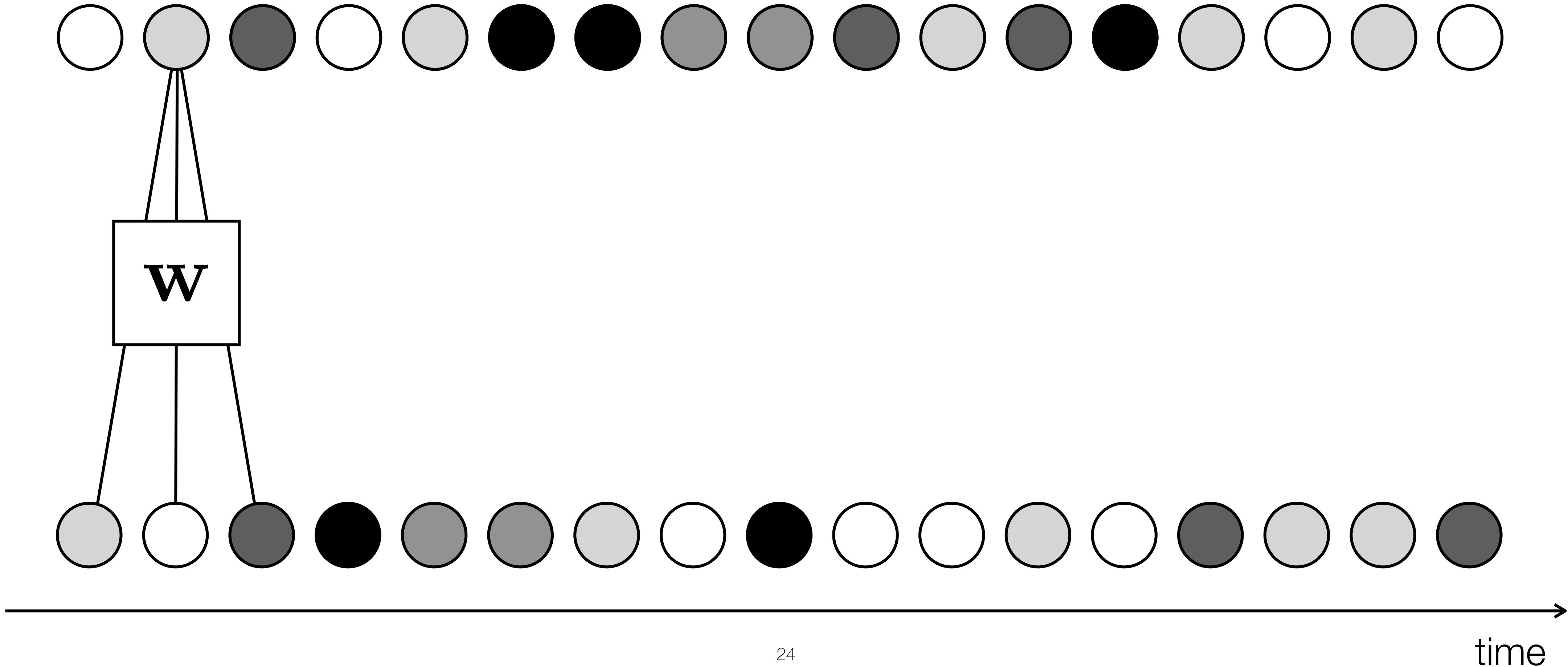


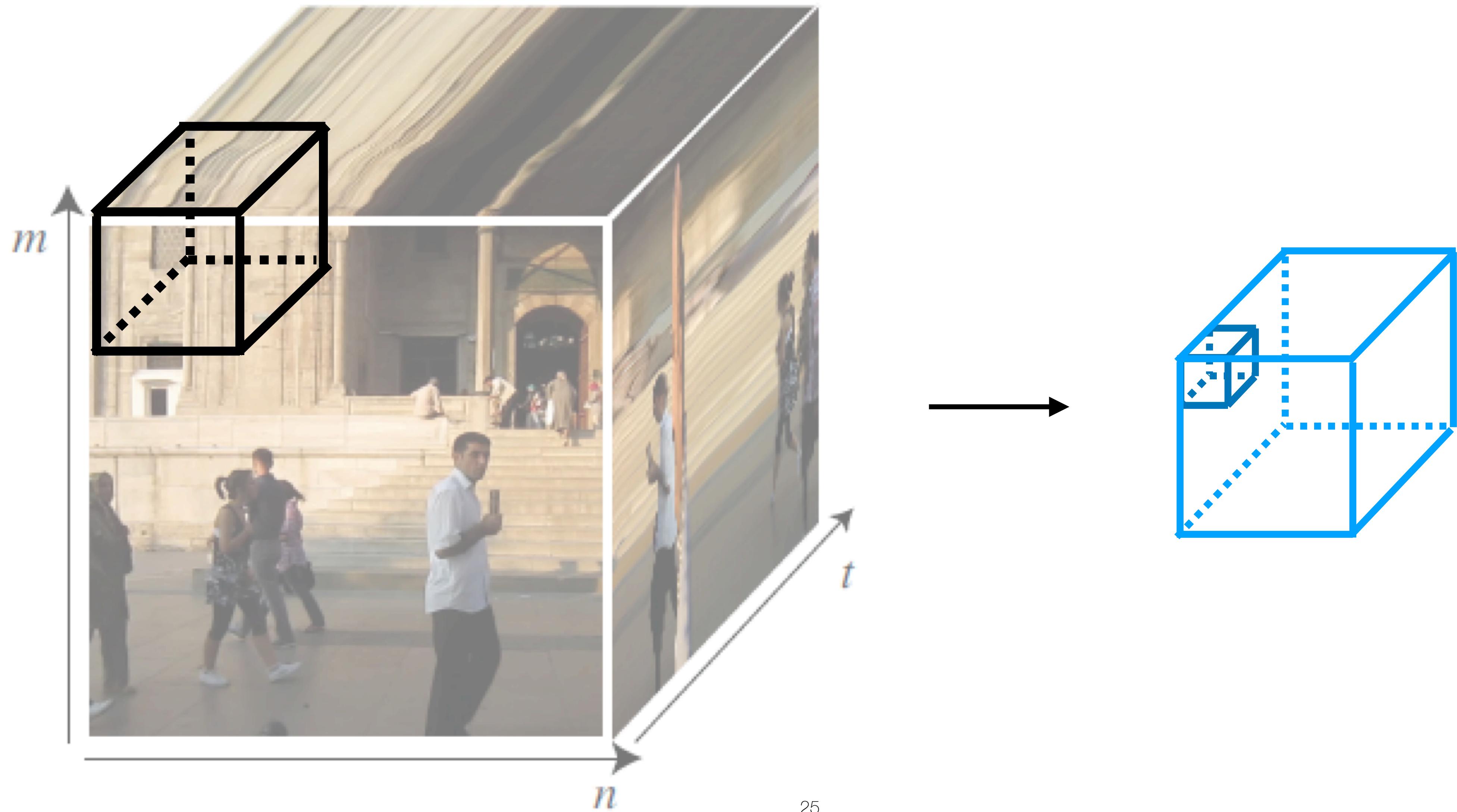
Source: Torralba, Freeman, Isola



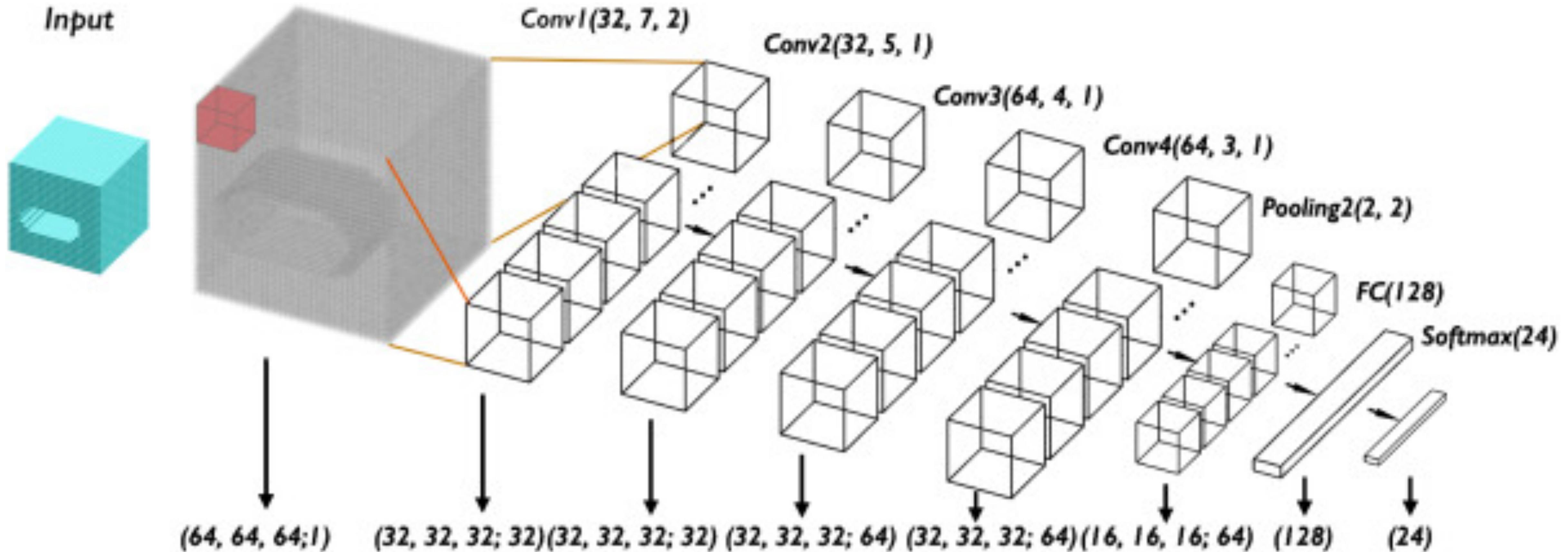


Convolutions in time



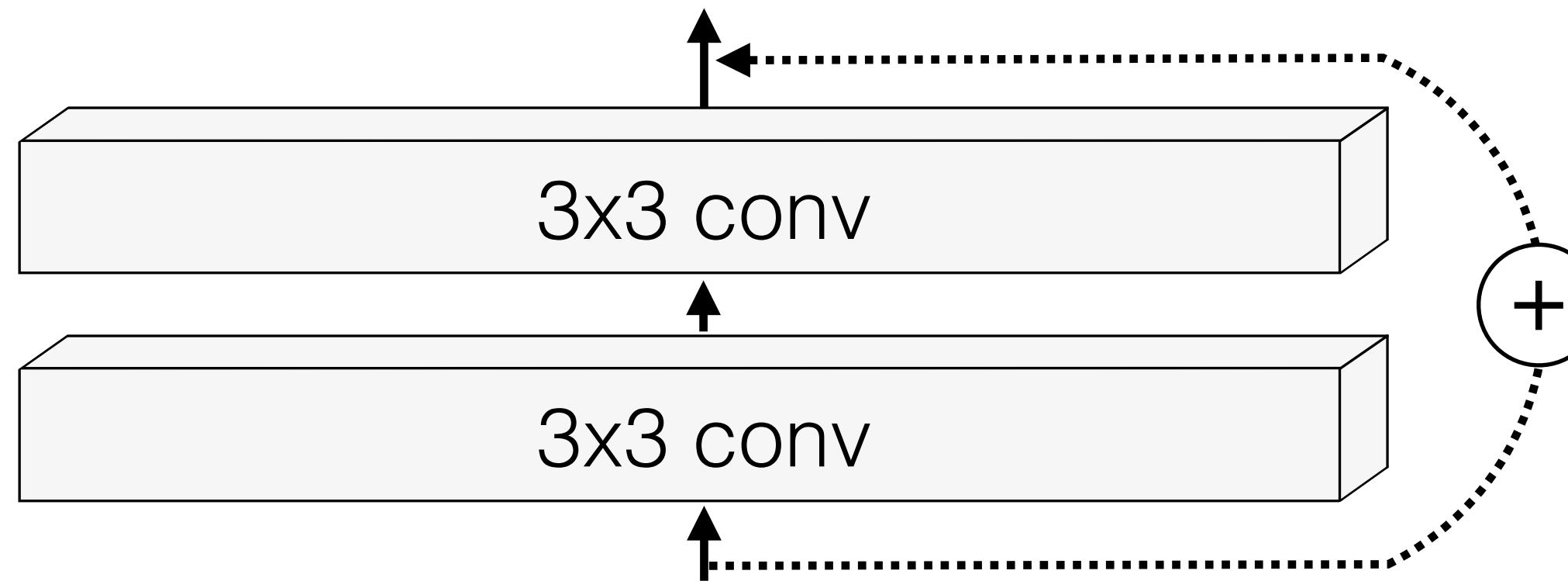


25



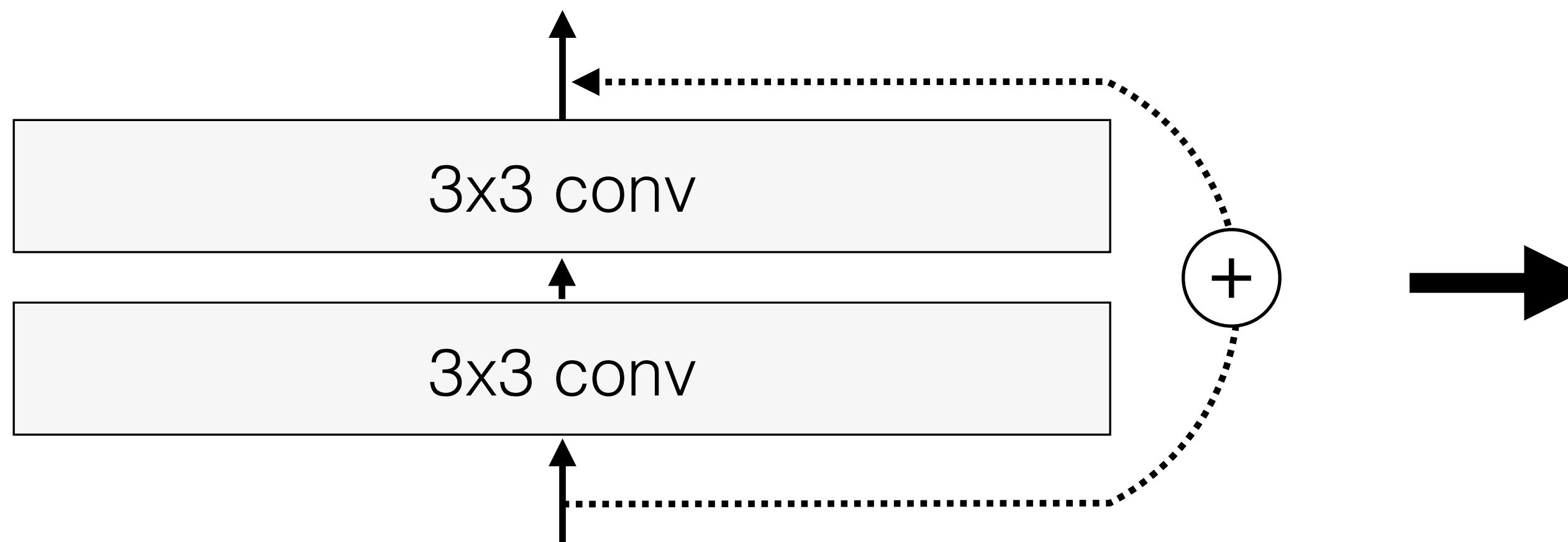
[fig from FeatureNet: Machining feature recognition based on 3D Convolution Neural Network]

Inflated convolutions

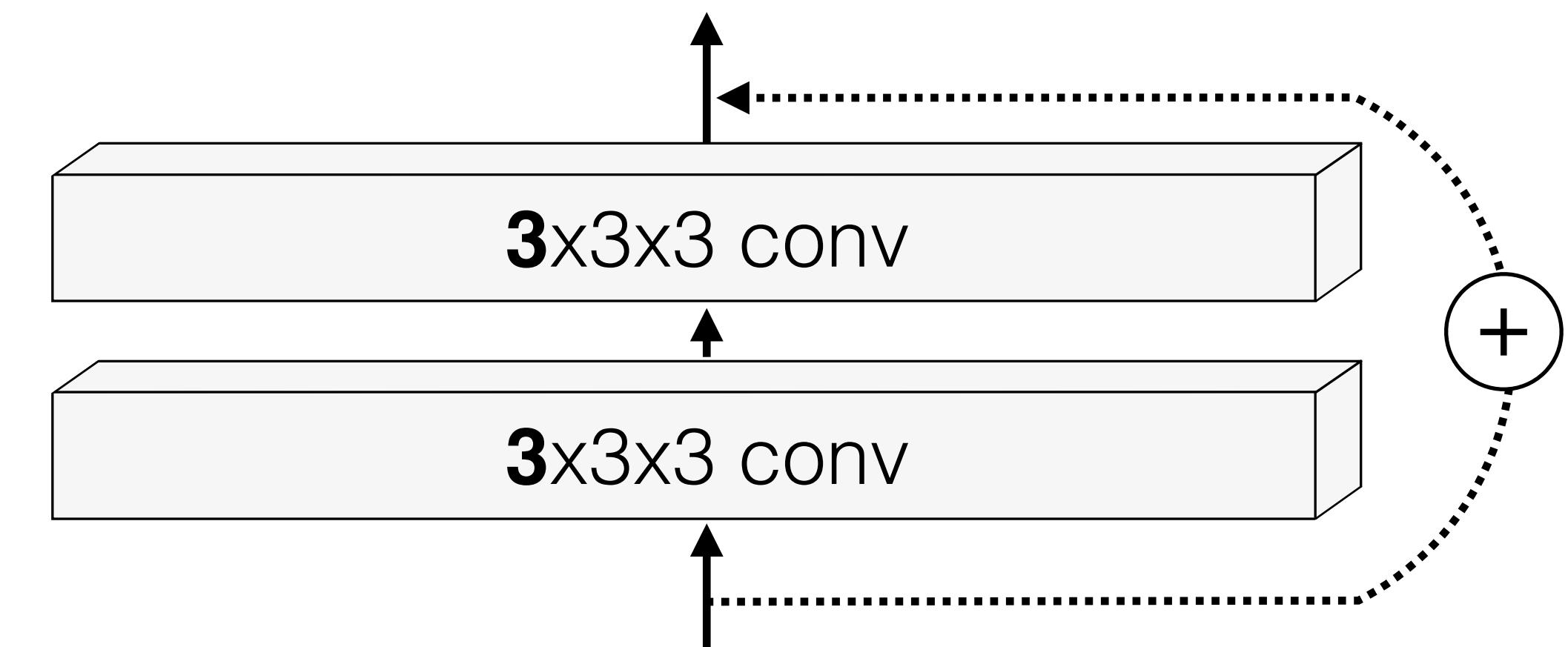


Inflated convolutions

2D ResNet block



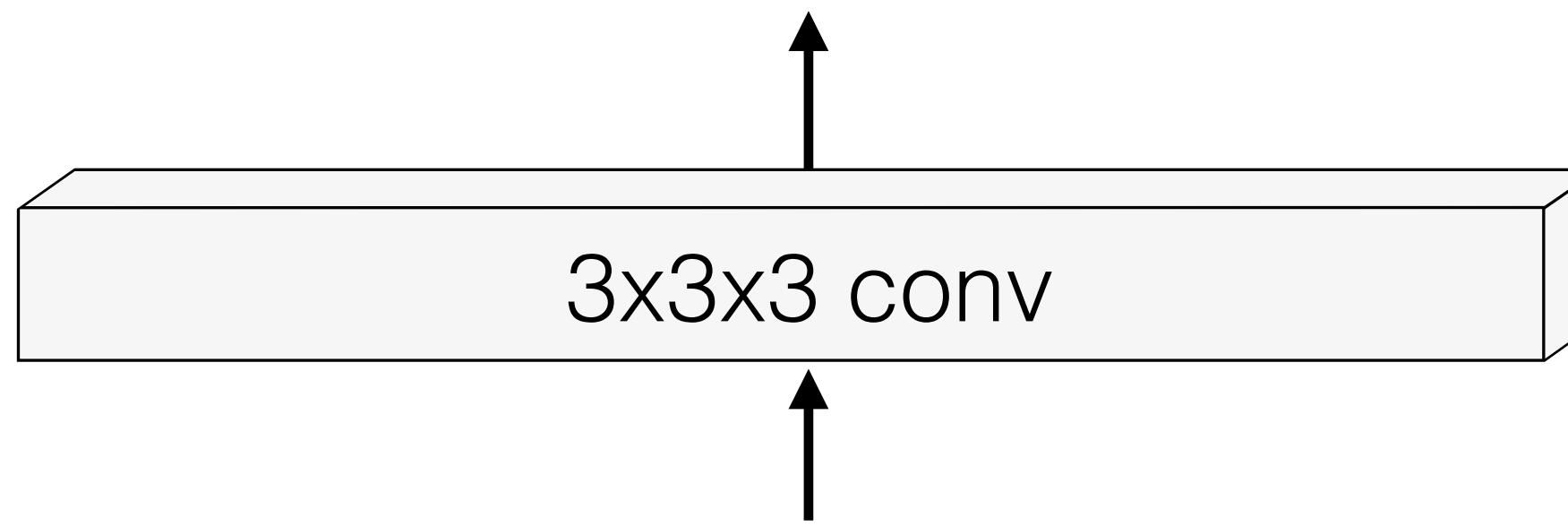
3D ResNet block



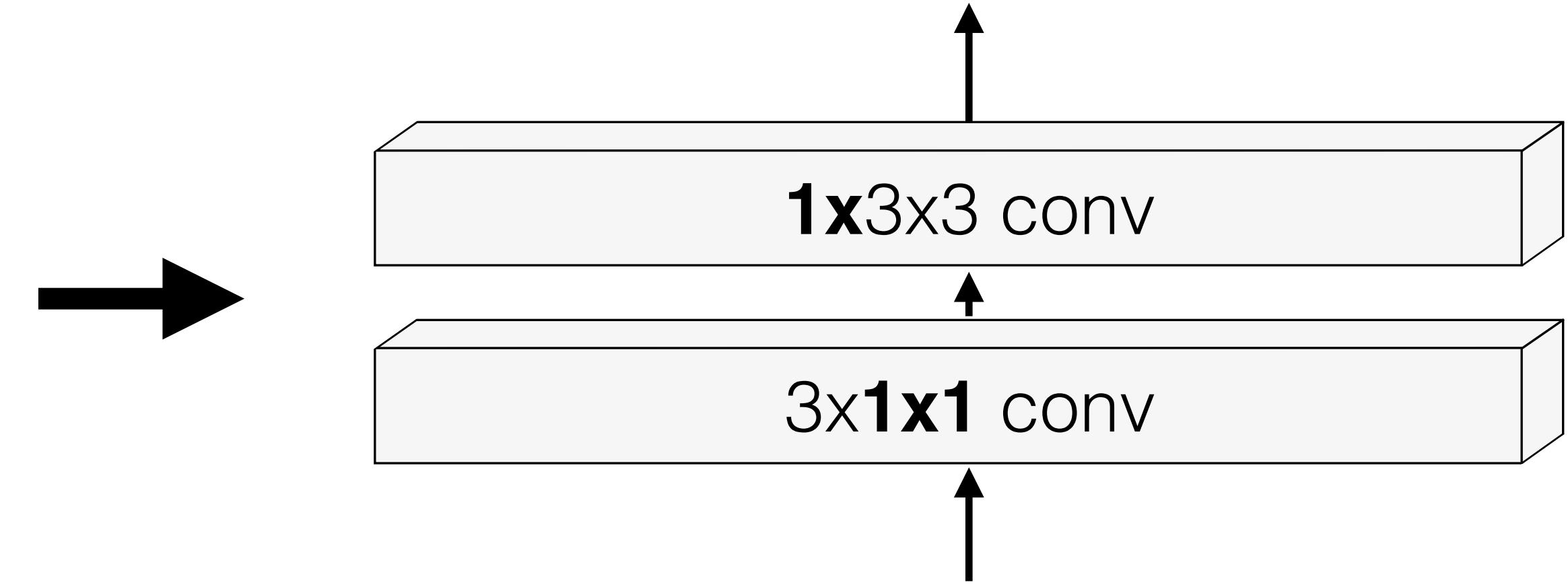
- Can reuse 2D architectures. [Carreira et al. 2017]
- Pretrain with 2D nets (“inflating” 2D filter to 3D)

Separable convolutions

3D convolution



Separate space/time



Faster and fewer parameters

Are we solving the right problem?



Photo by H. Edgerton

Are our models using motion?

- Single-frame models do badly.
- Shuffling frames after training hurts performance
- But 2D convolutions + temporal pooling often is competitive!

When do we actually need motion?



Making latte art



Jaywalking



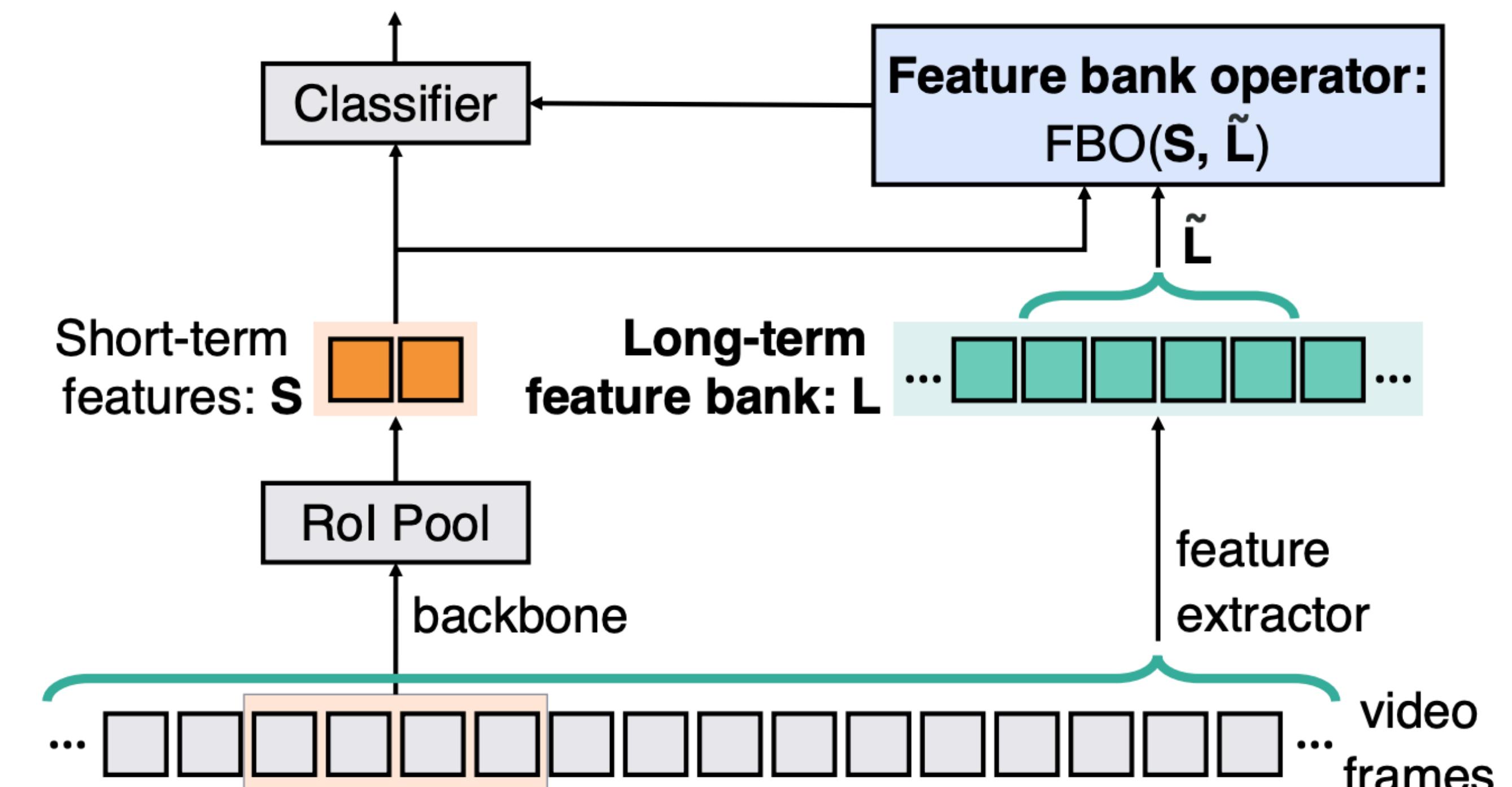
Grooming dog

Multiple frames often means multiple samples

Localizing objects in space and time



AVA dataset [Gu et al. 2018]

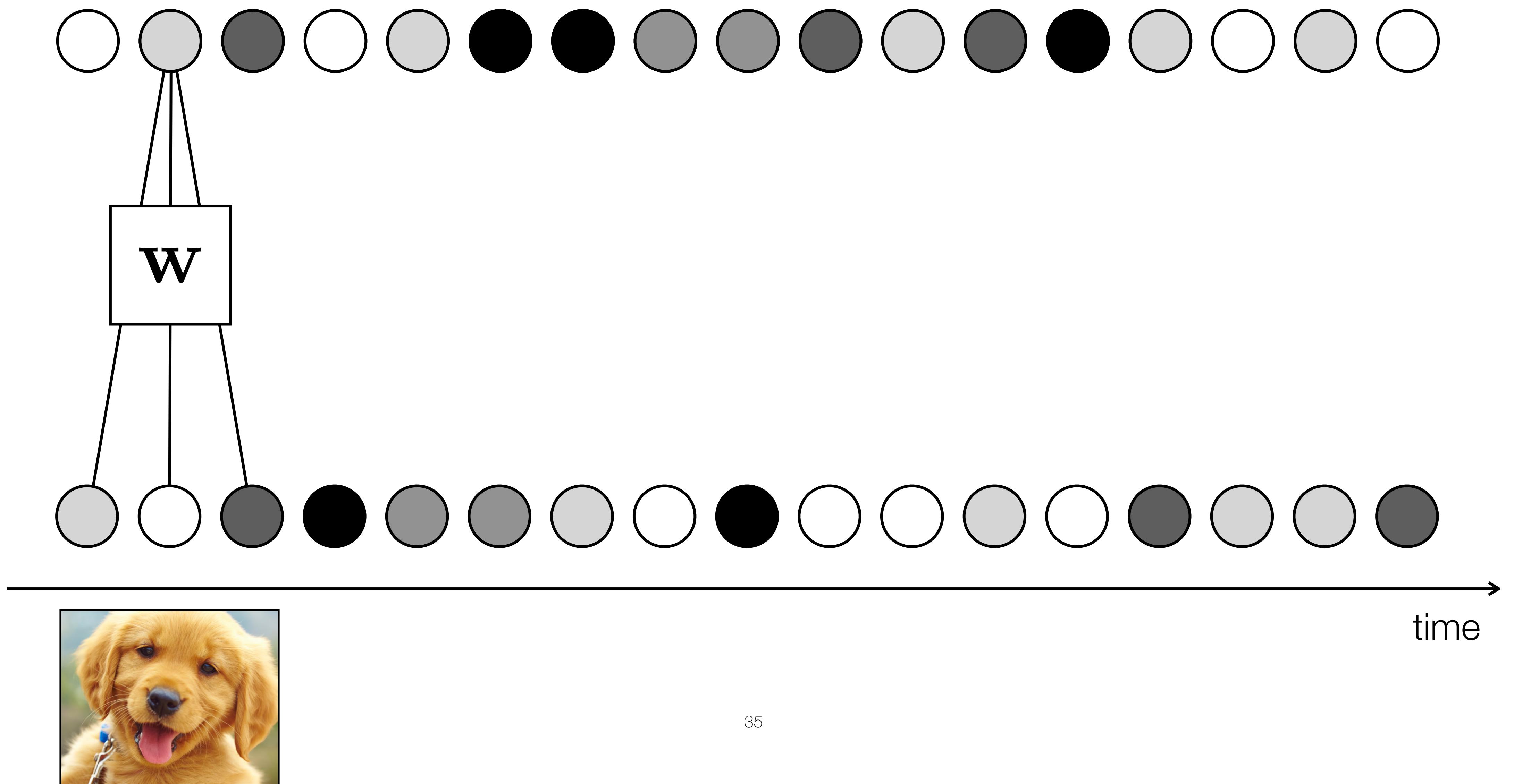


[Wu et al. 2019]

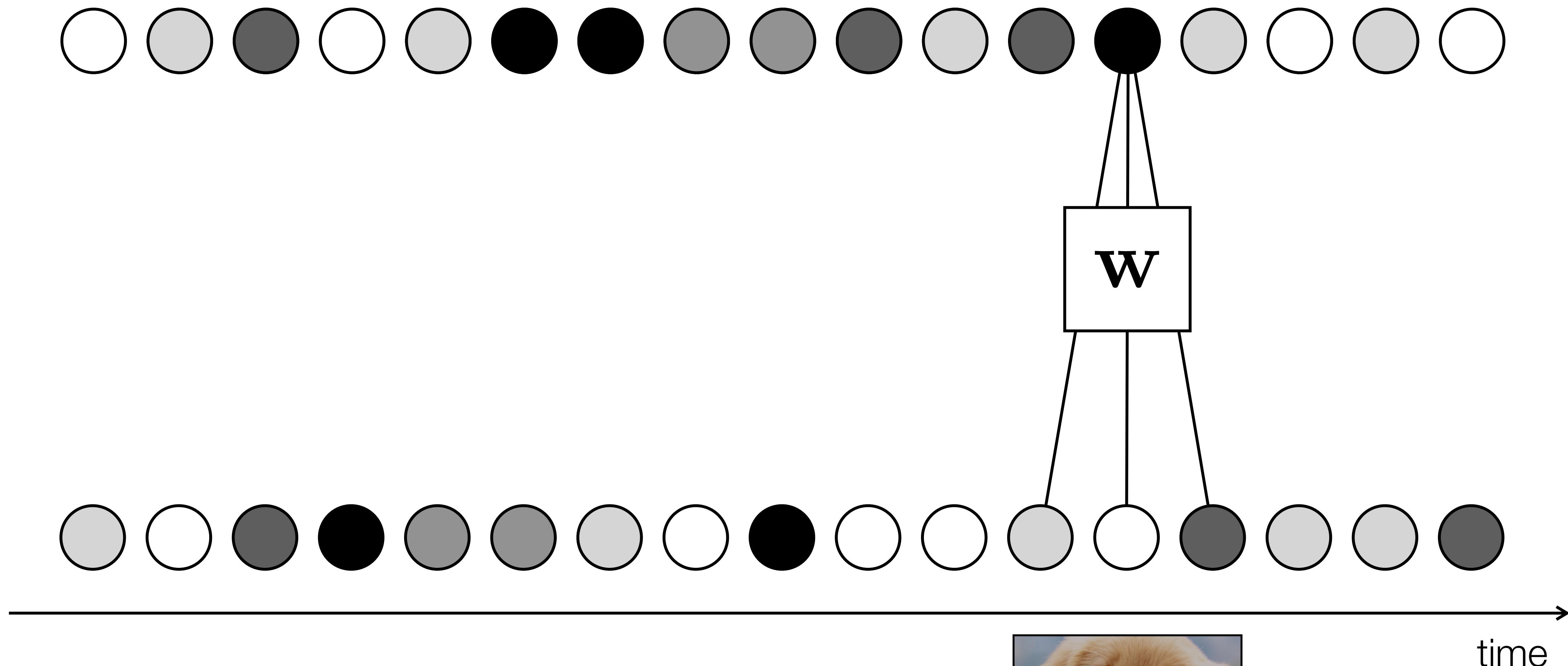
Recurrent networks



Rufus

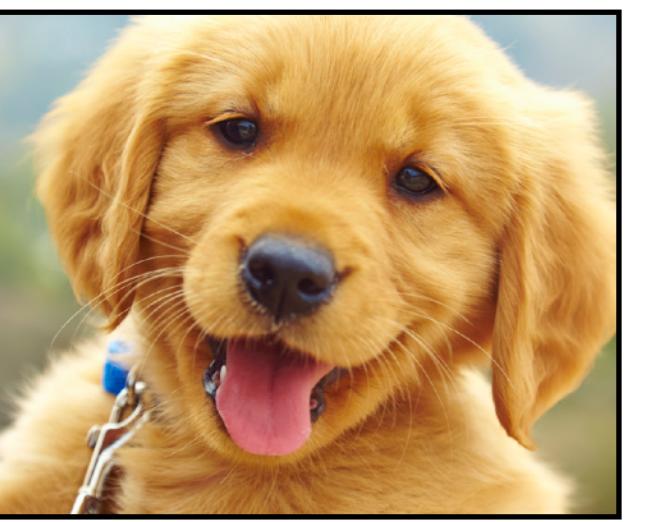
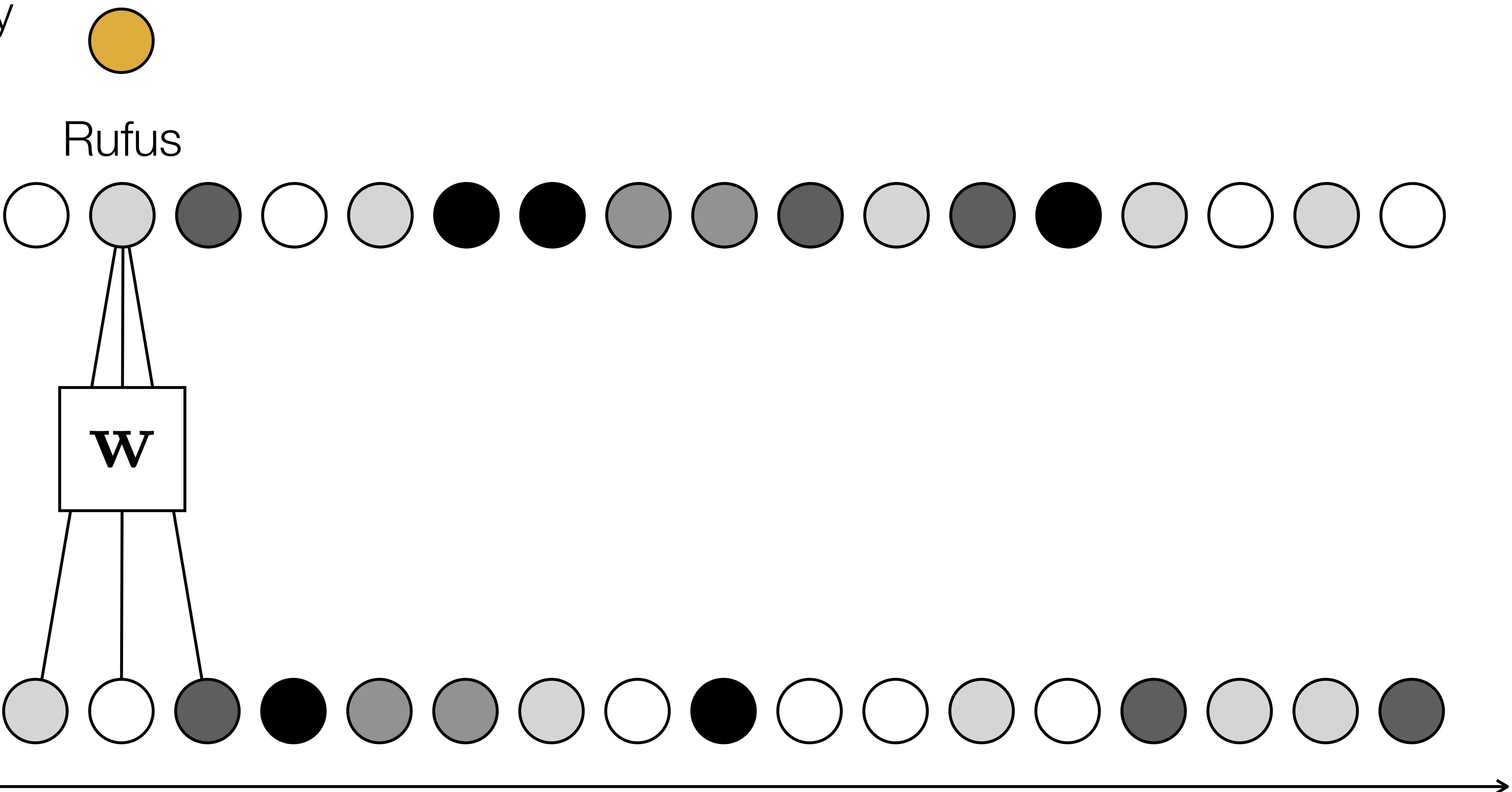


Douglas

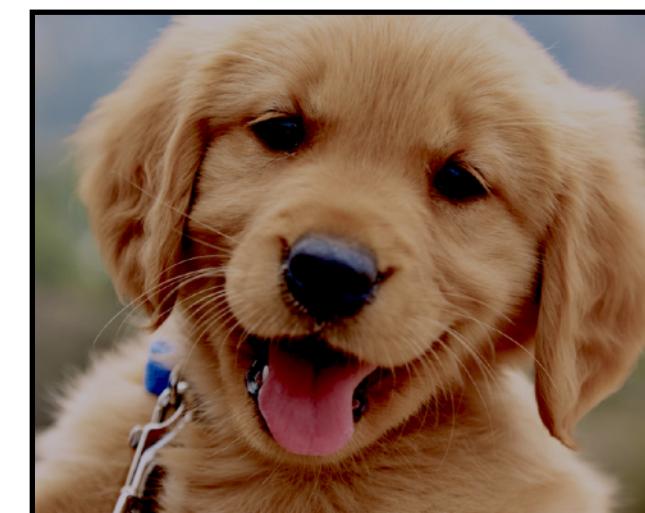
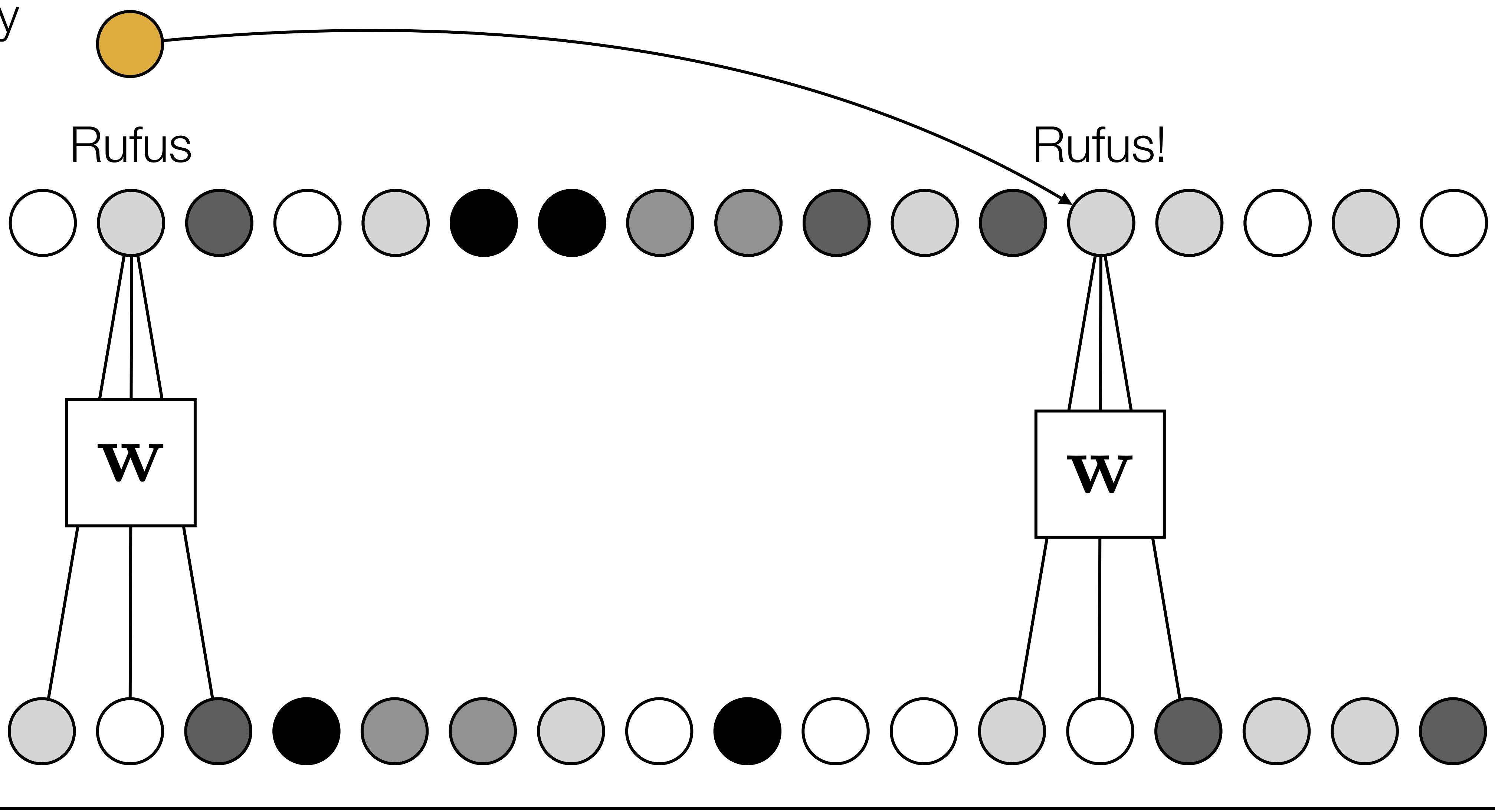


Memory
unit

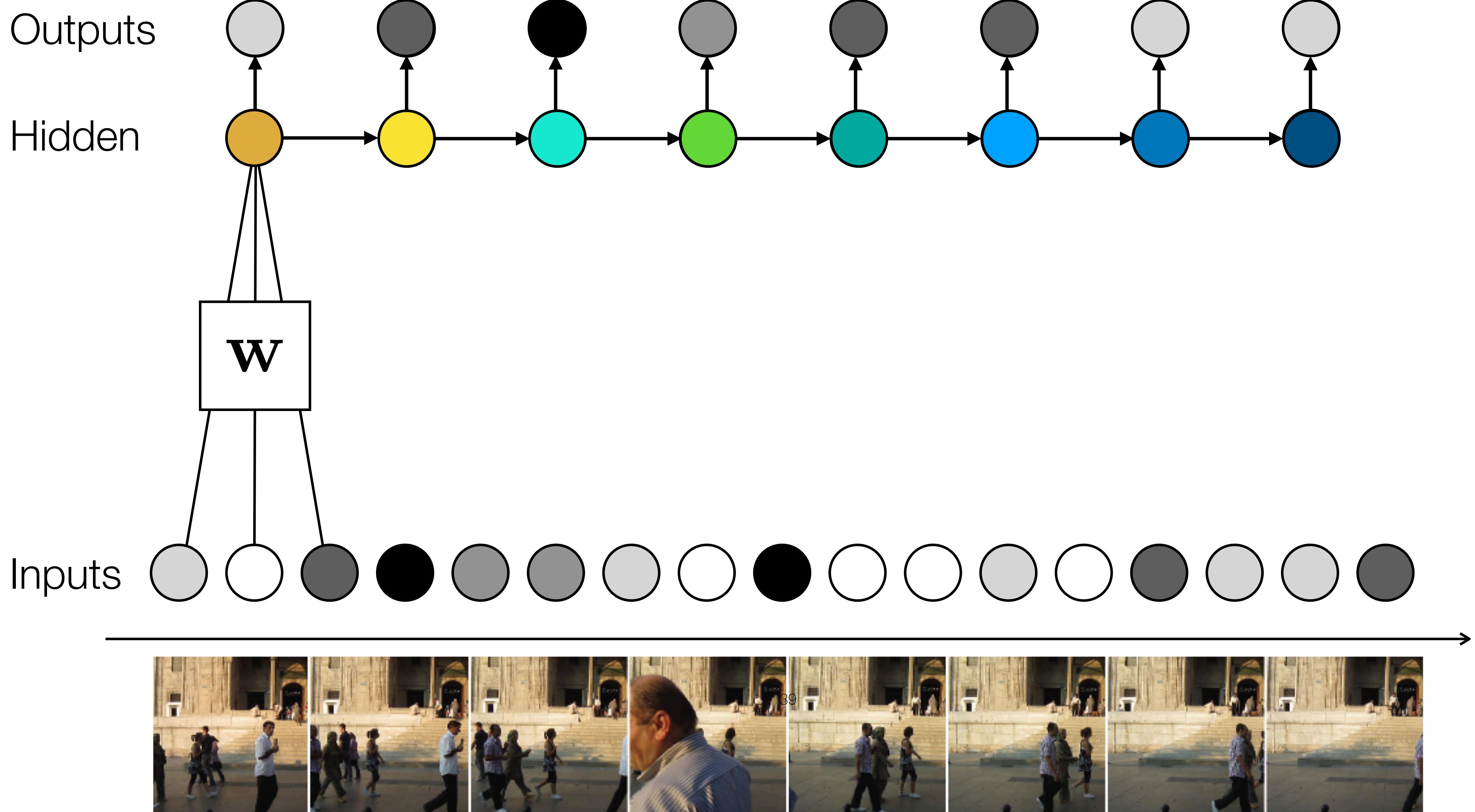
Rufus



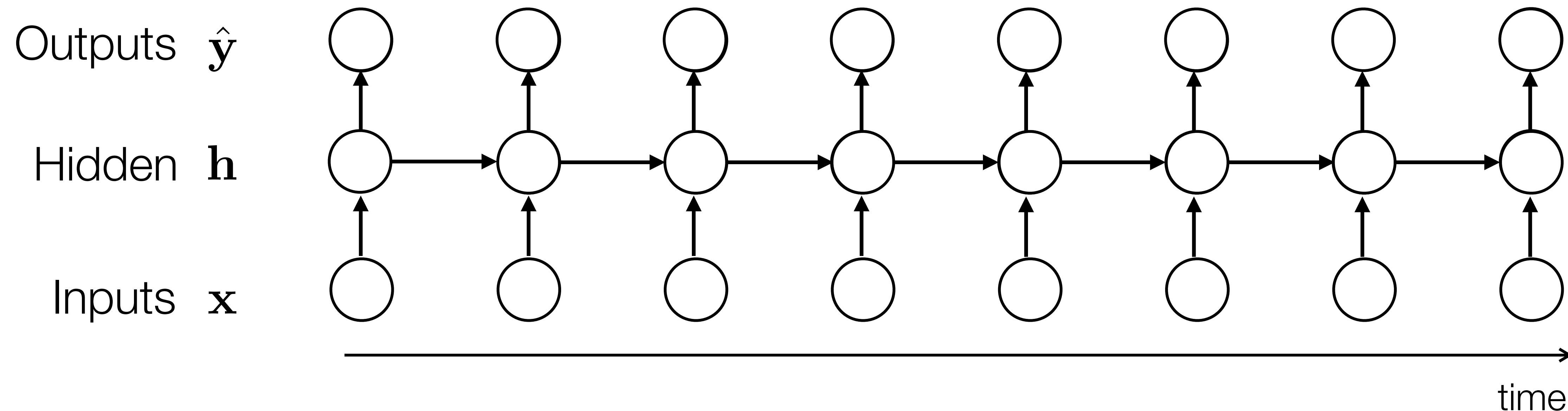
Memory
unit



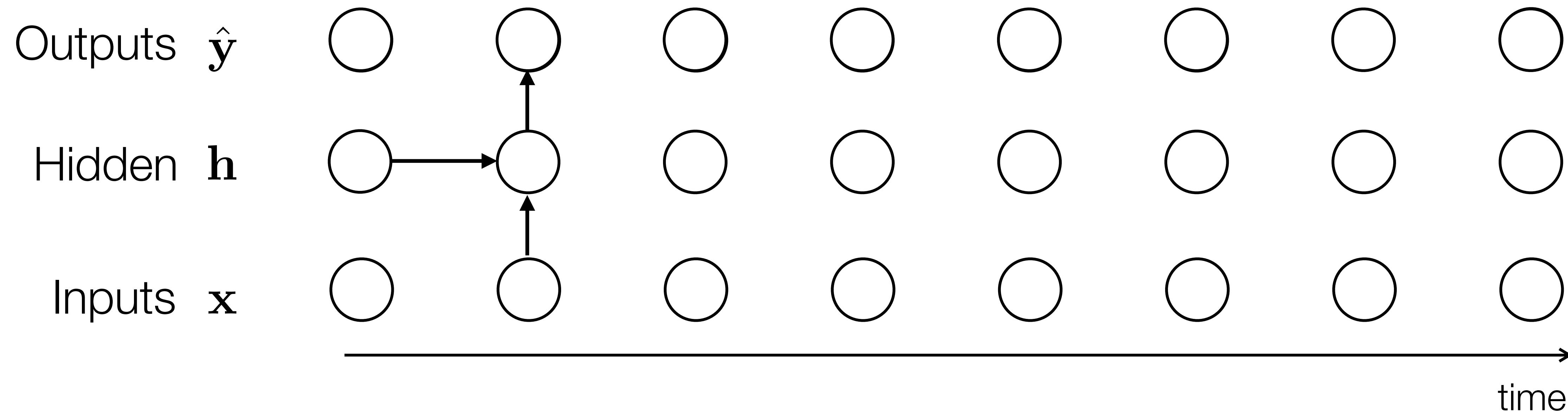
Recurrent Neural Networks (RNNs)



Recurrent Neural Networks (RNNs)



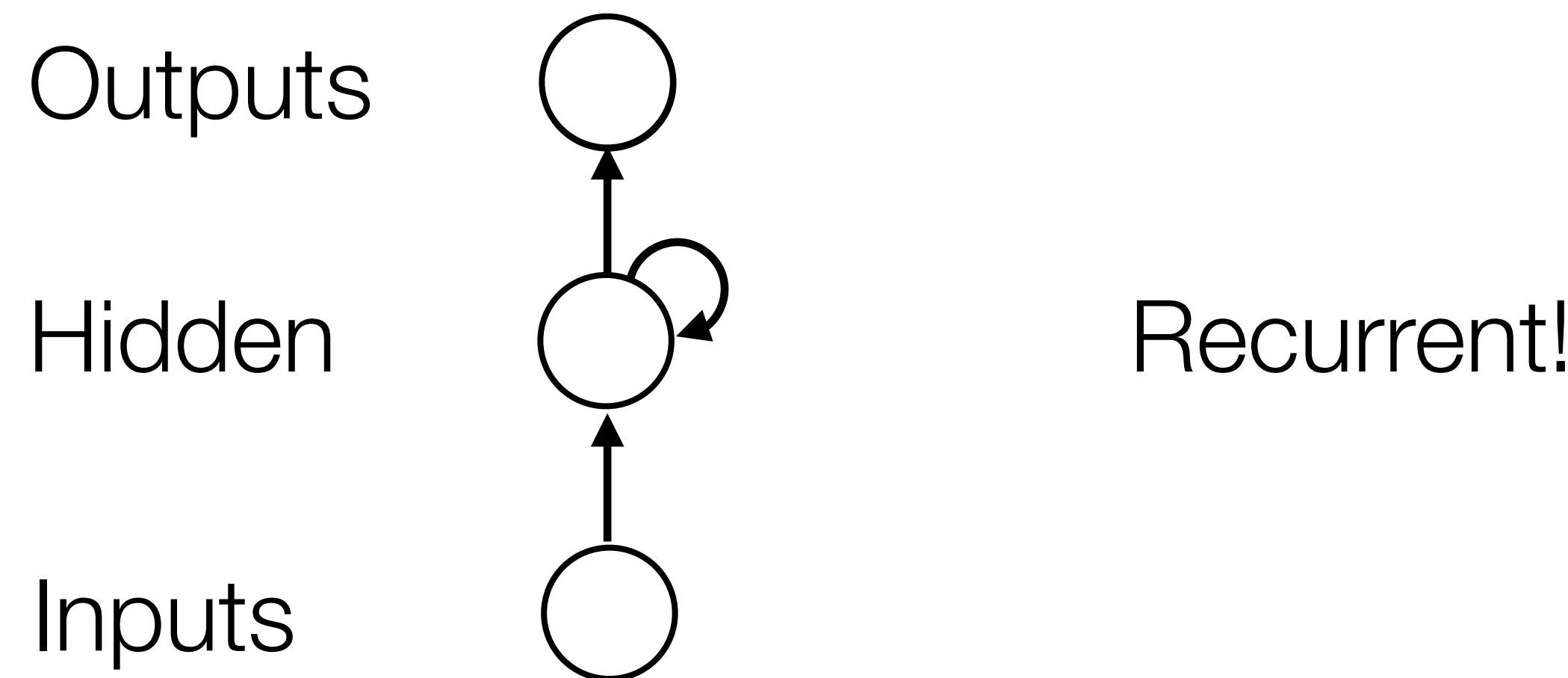
Recurrent Neural Networks (RNNs)



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

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

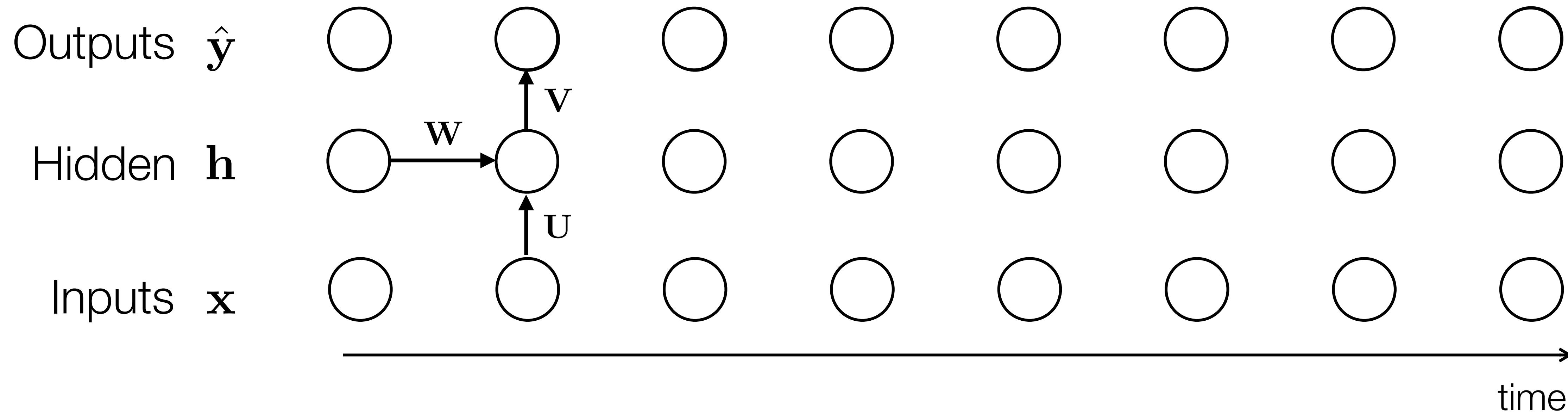
Recurrent Neural Networks (RNNs)



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

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

Recurrent Neural Networks (RNNs)



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

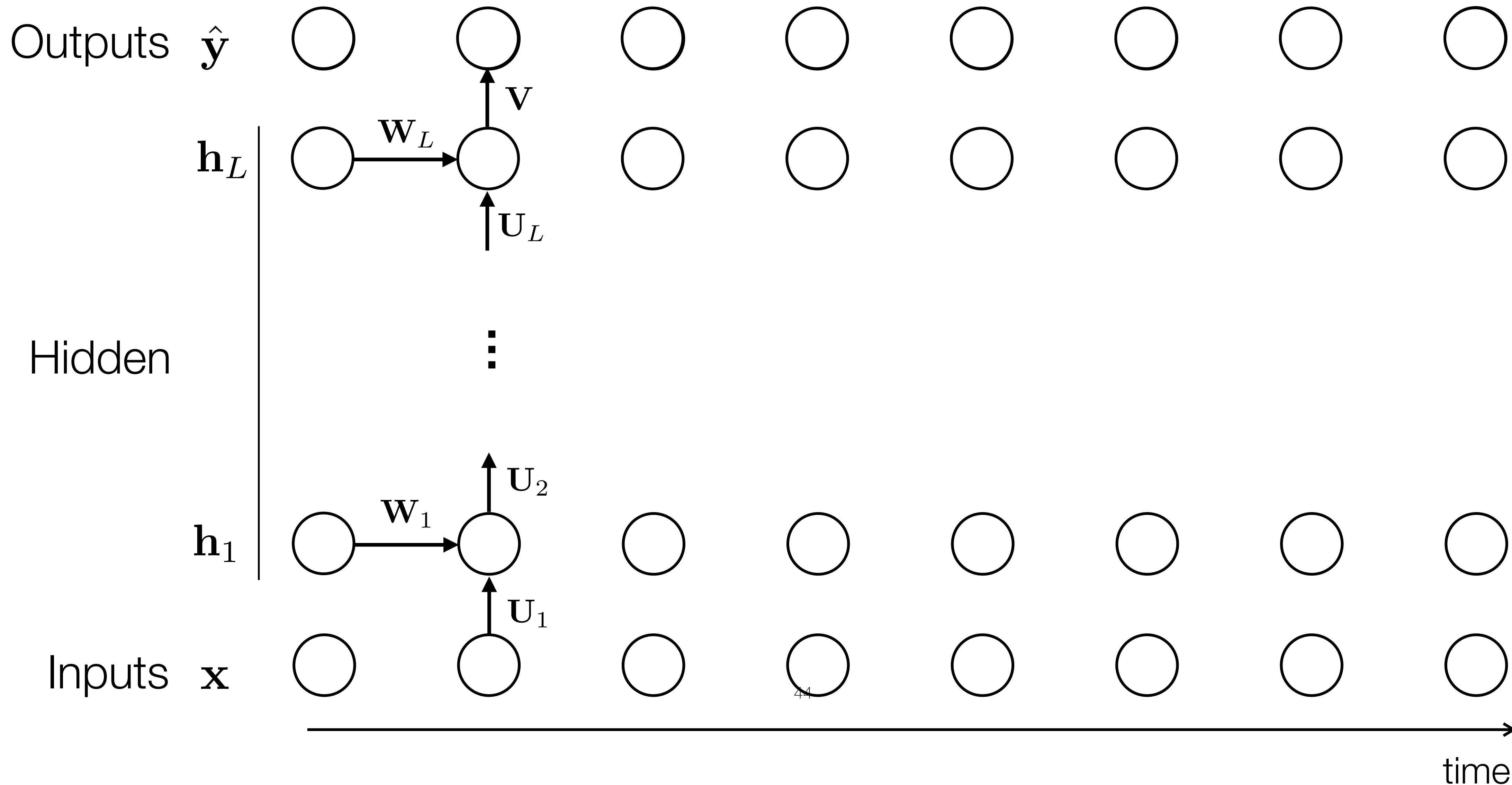
$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)})$$

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

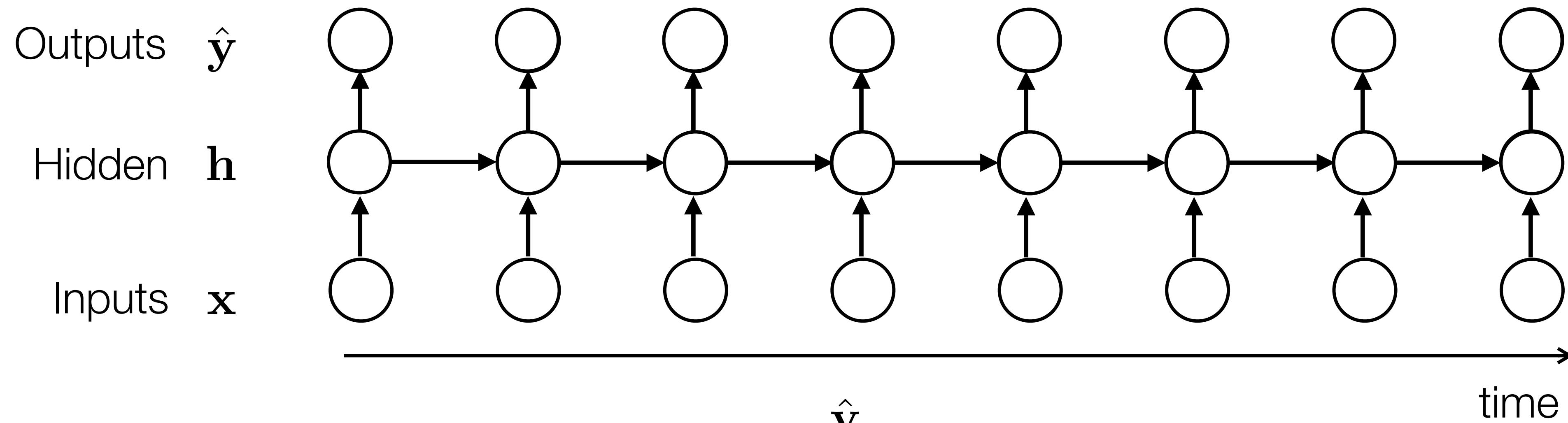
43

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

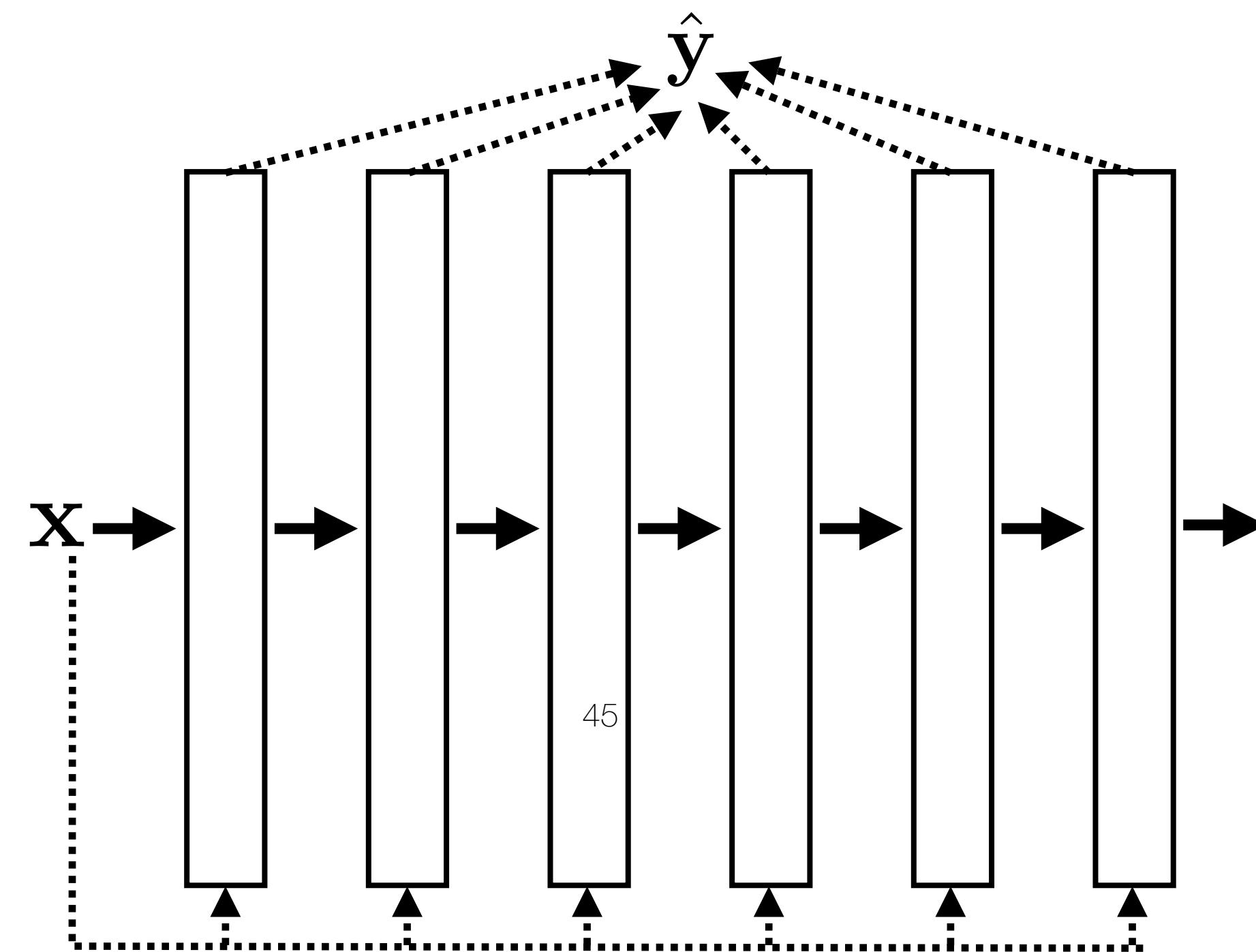
Deep Recurrent Neural Networks (RNNs)



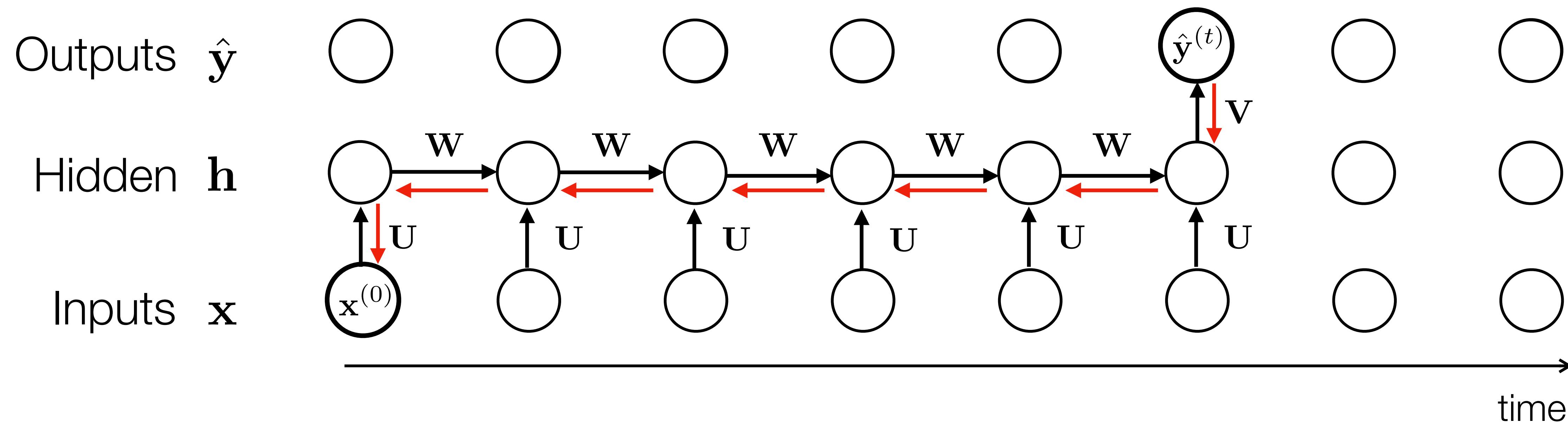
Unrolling an RNN



Equivalent to
"unrolled"
network with
lots of layers

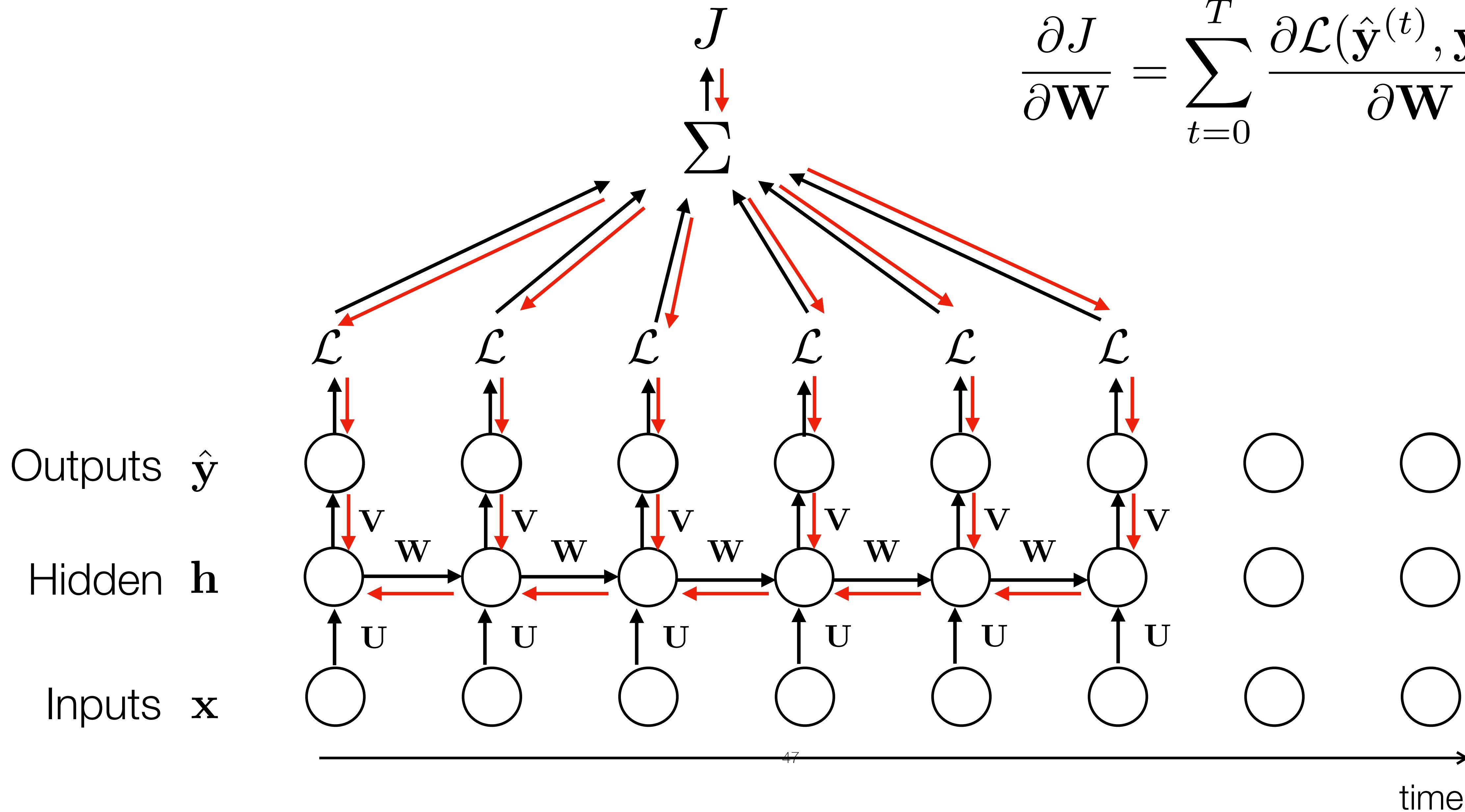


Backprop through time

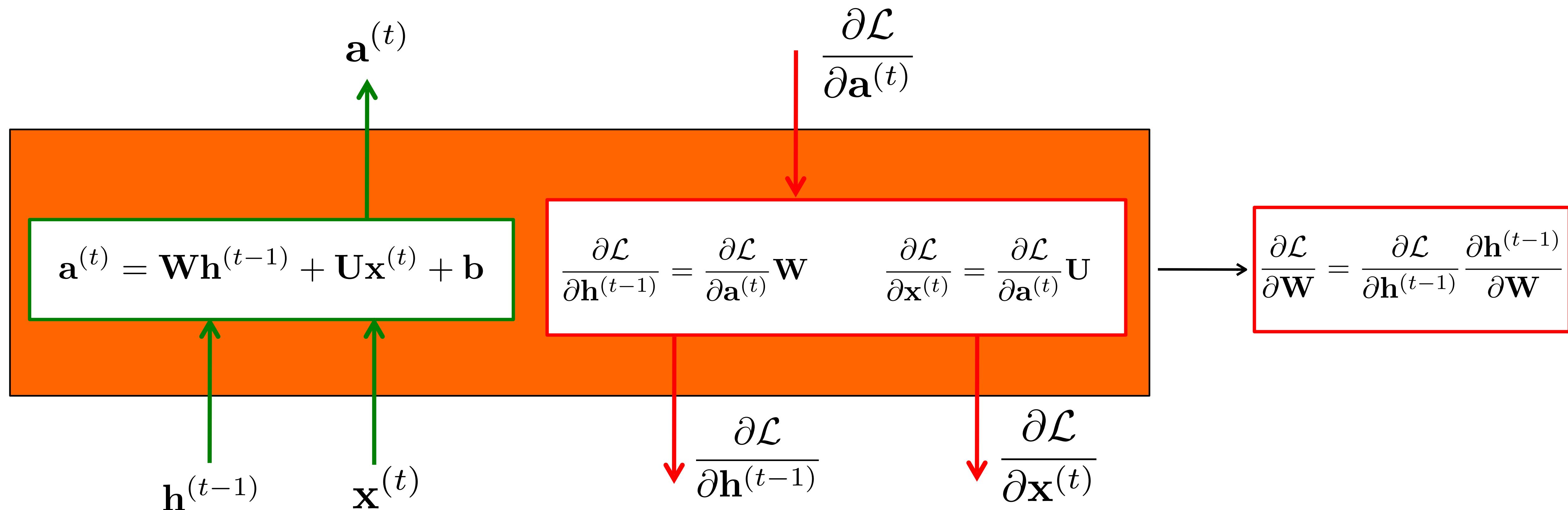


$$\frac{\partial \hat{y}^{(t)}}{\partial \mathbf{x}^{(0)}} = \frac{\partial \hat{y}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \cdots \frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{h}^{(0)}} \frac{\partial \mathbf{h}^{(0)}}{\partial \mathbf{x}^{(0)}}$$

$$\frac{\partial J}{\partial \mathbf{W}} = \sum_{t=0}^T \frac{\partial \mathcal{L}(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)})}{\partial \mathbf{W}}$$

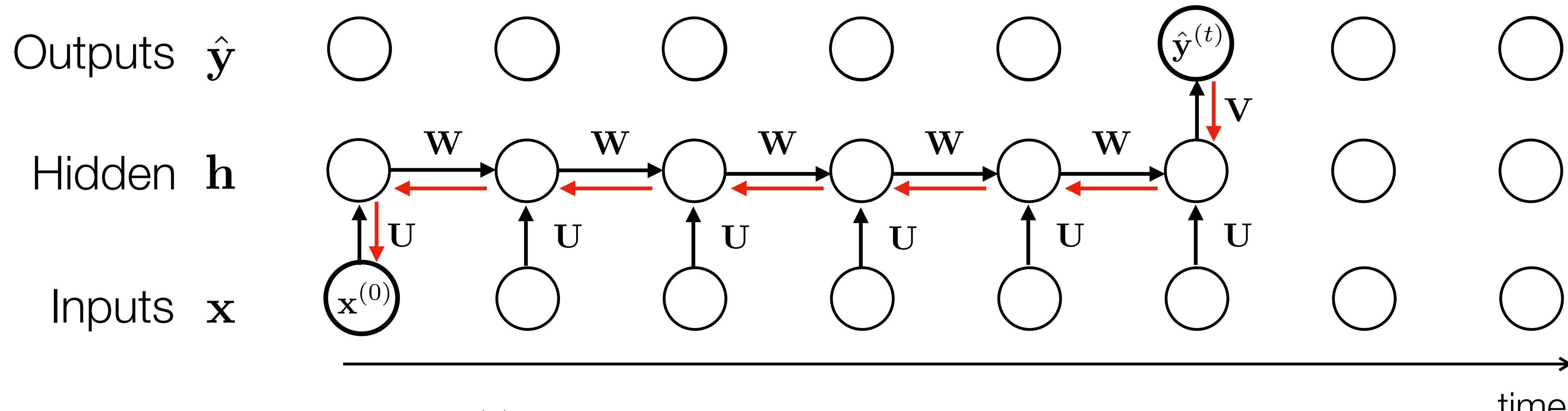


Recurrent linear layer



$$\frac{\partial J}{\partial \mathbf{W}} = \sum_{t=0}^T \frac{\partial \mathcal{L}(\hat{\mathbf{y}}^{(t)}, \mathbf{y}^{(t)})}{\partial \mathbf{W}}$$

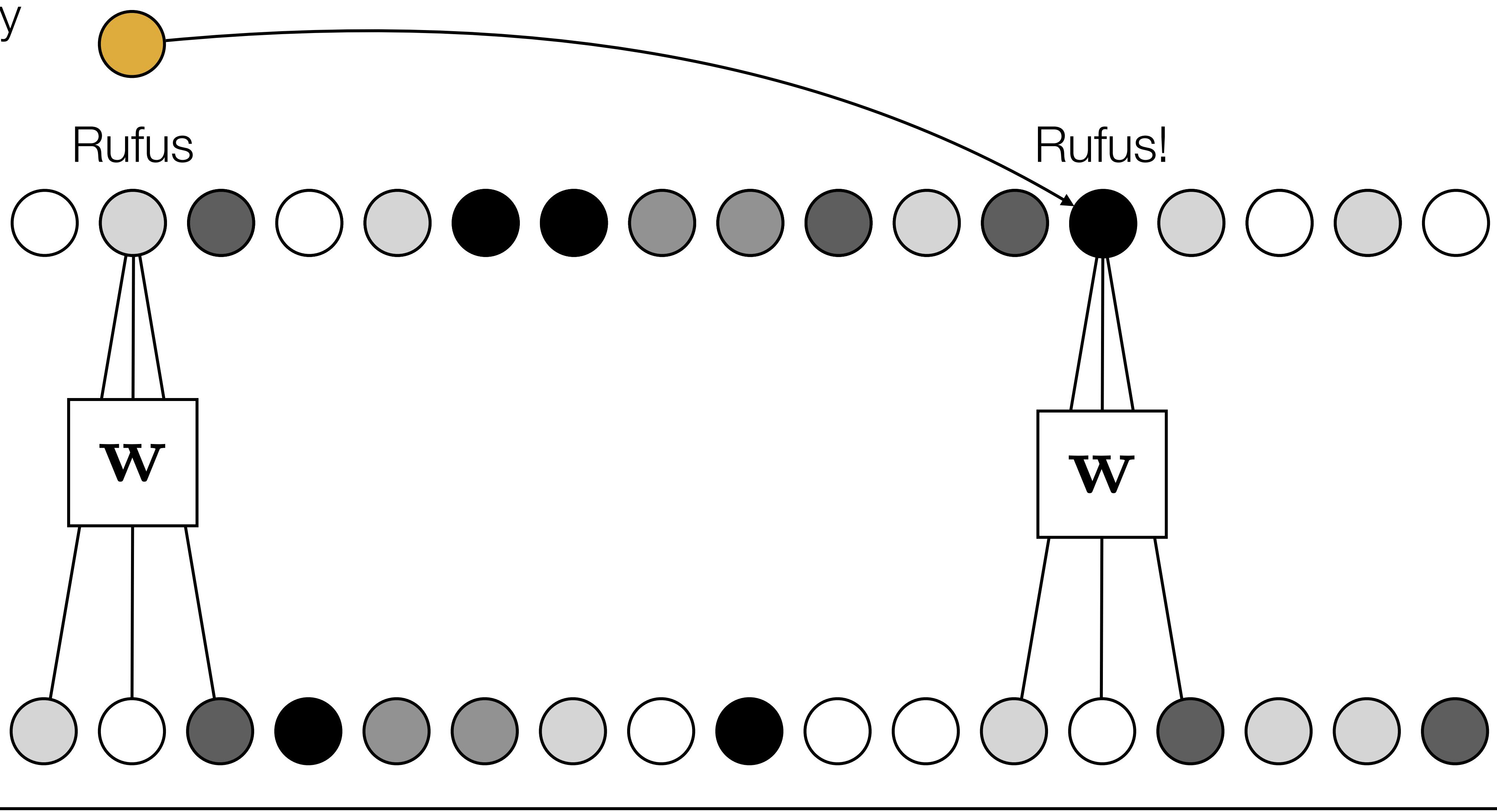
The problem of long-range dependences



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

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

Memory
unit



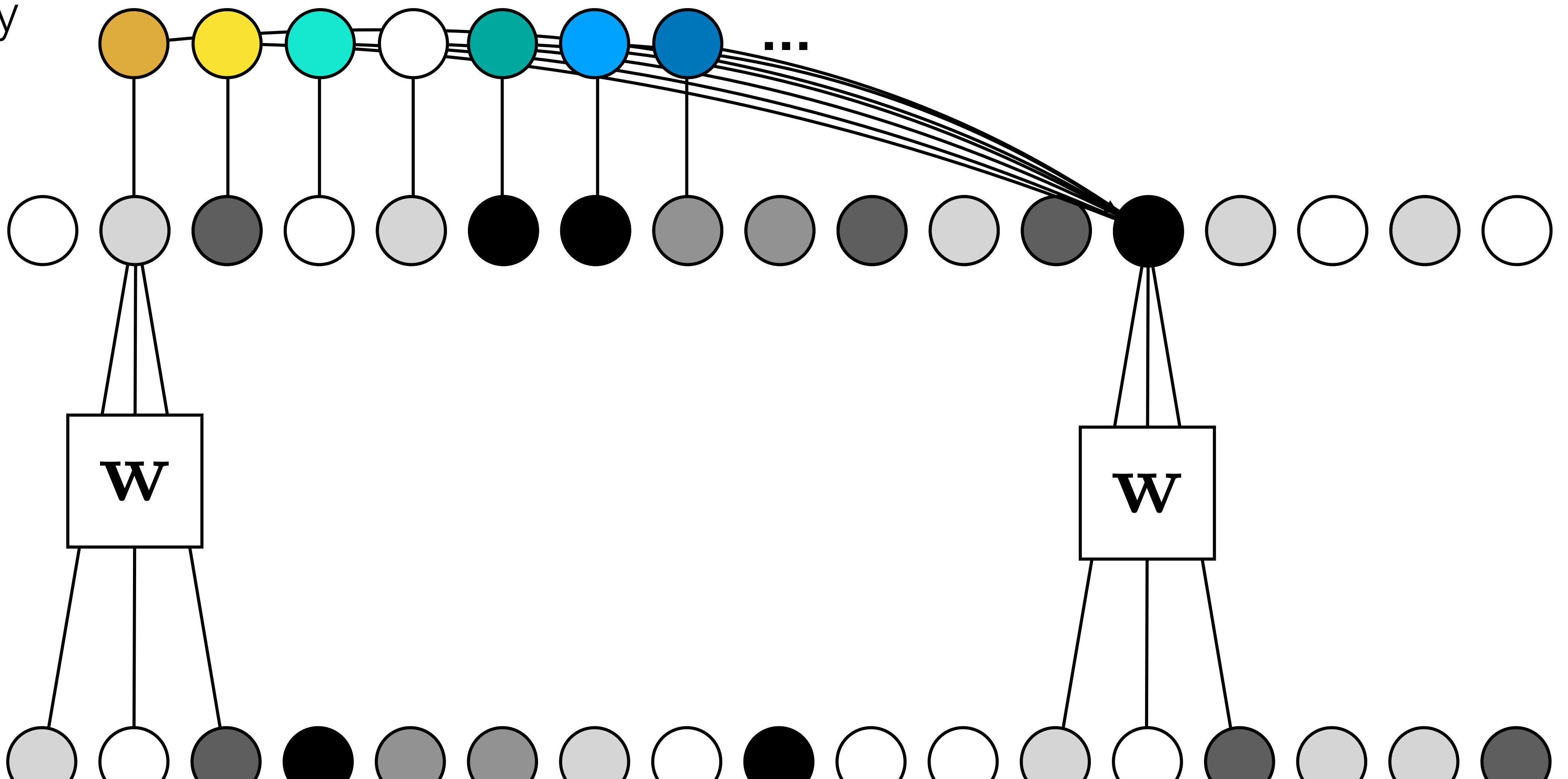
50



time

Source: Torralba, Freeman, Isola

Memory
units



51

time

Source: Torralba, Freeman, Isola

The problem of long-range dependences

Why not remember everything?

- Memory size grows with t
- This kind of memory is **nonparametric**: there is no finite set of parameters we can use to model it
- Markov-like assumption: future state only dependent on preceding hidden state. Very easy to forget things.
- By putting the right info in to the hidden state, RNNs can model dependencies that are arbitrarily⁶² far apart

LSTMs

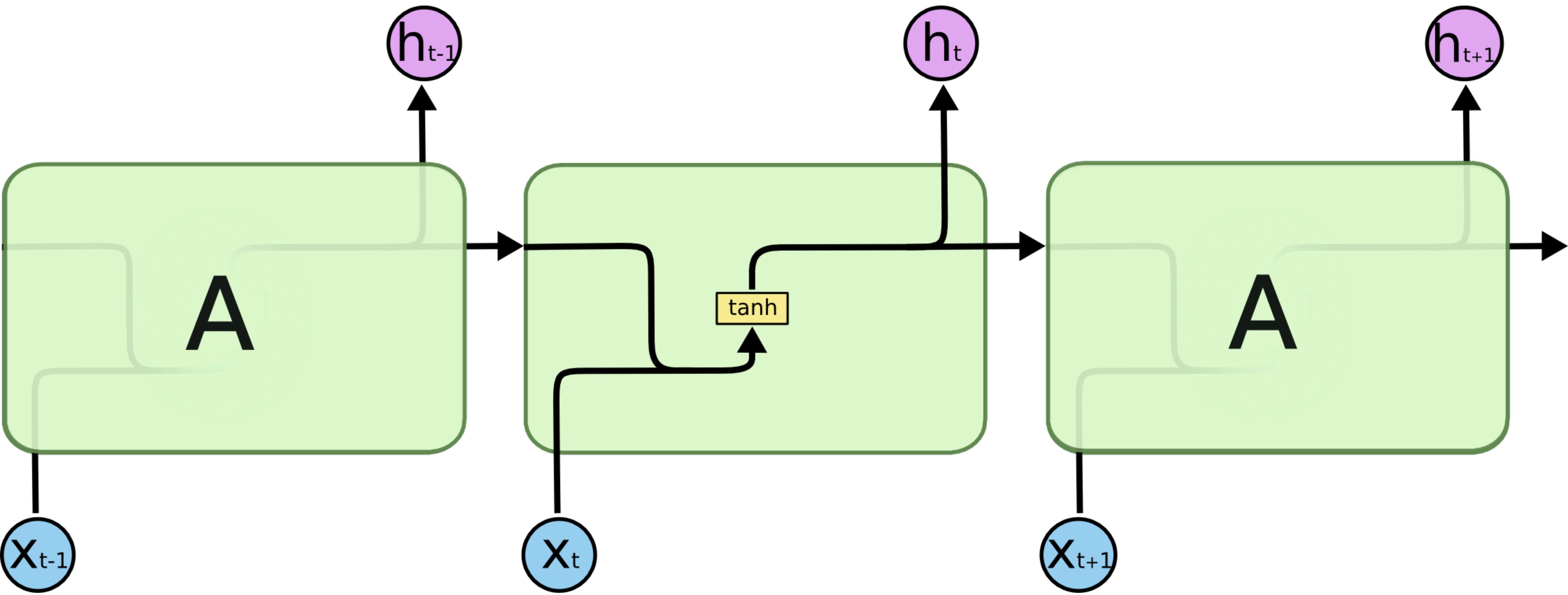
Long Short Term Memory

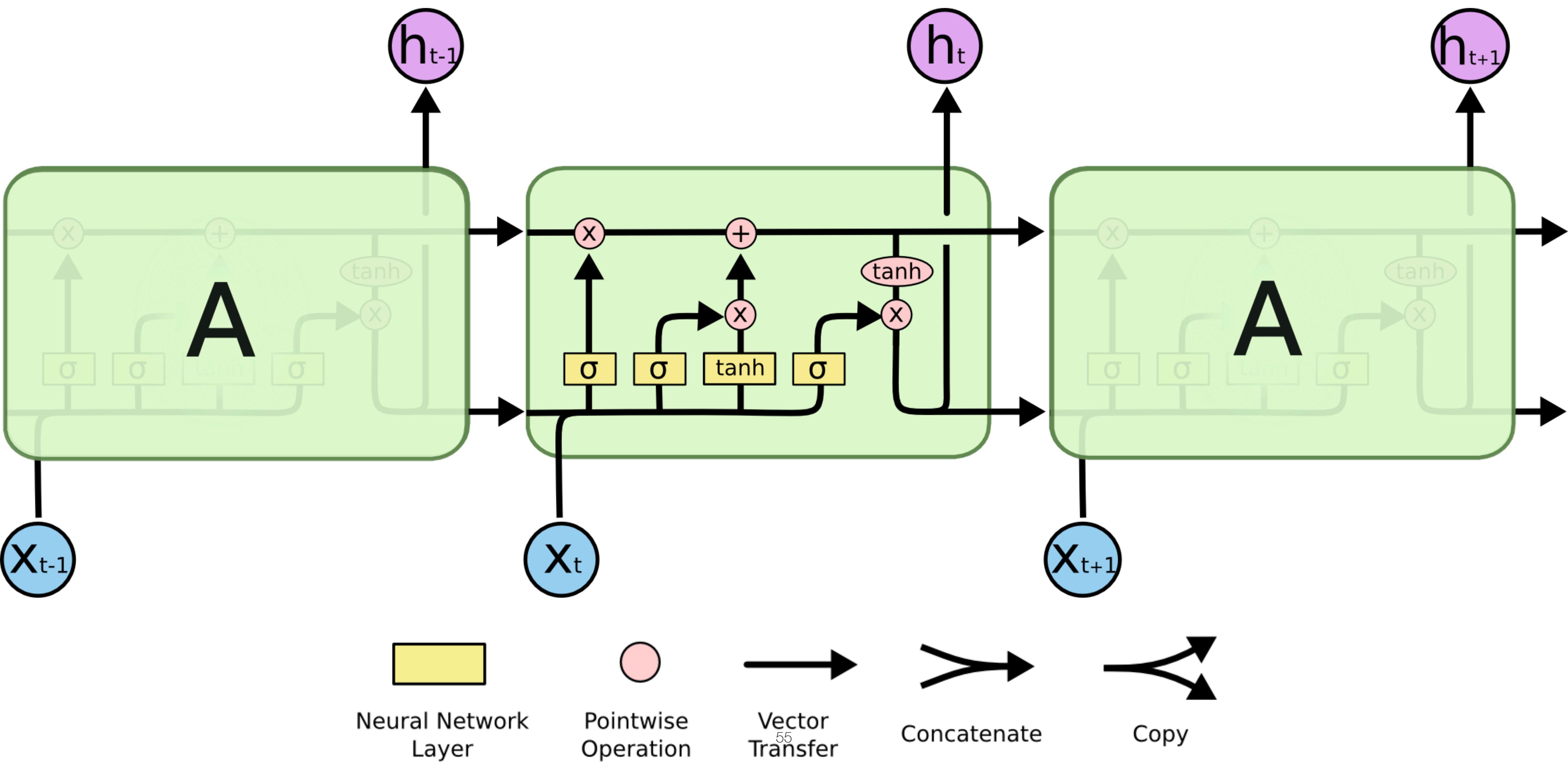
A special kind of RNN designed to avoid forgetting.

Related to ResNets inductive bias is that state transition is an identity function.

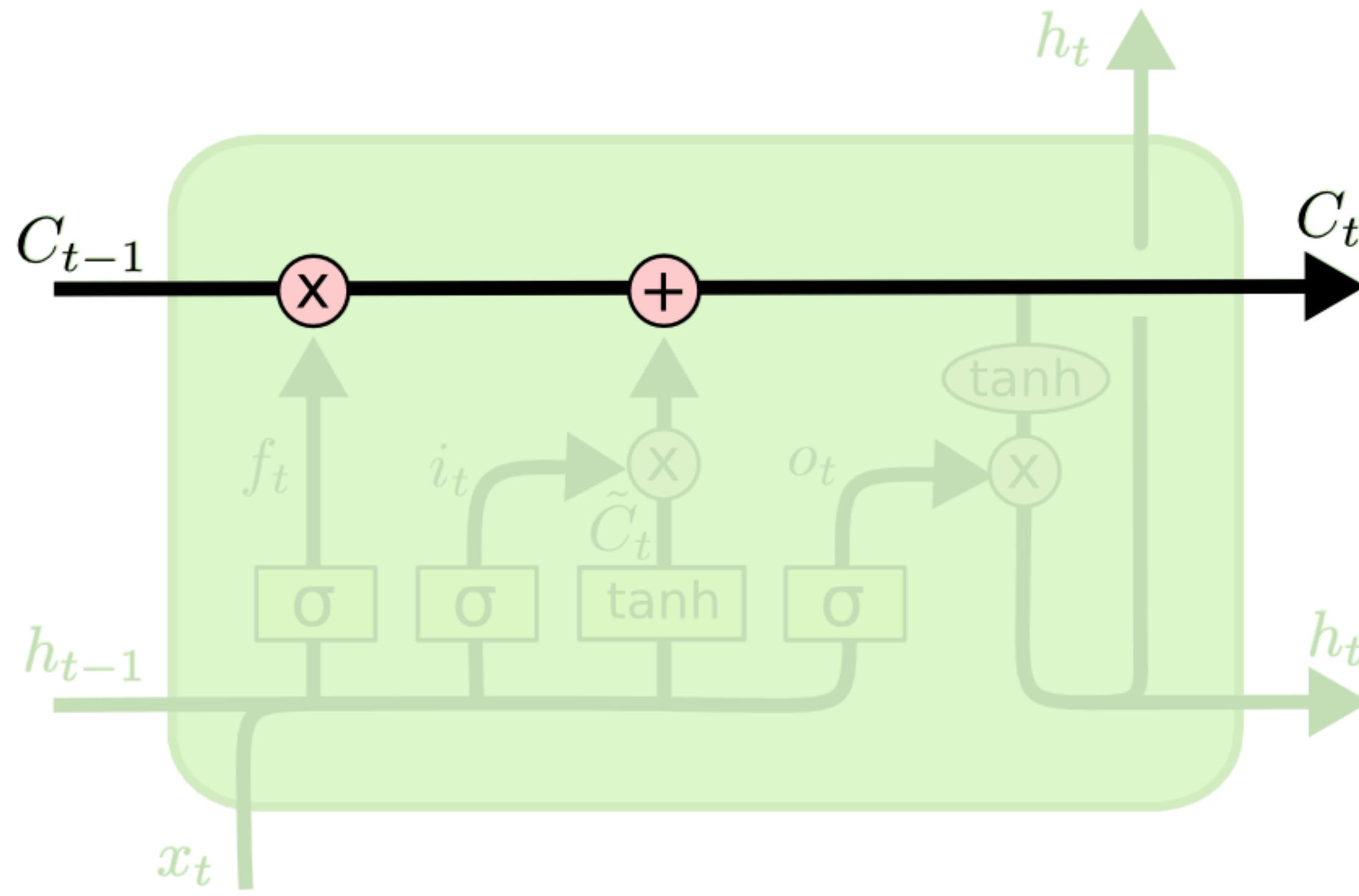
This way the default behavior is not to forget an old state. Instead of forgetting by default, the network has to *learn to forget*.

Bit of a complex design. Works well but simpler methods like Gated Recurrent Unit (GRU) are competitive [Jozefowicz et al. 2015].

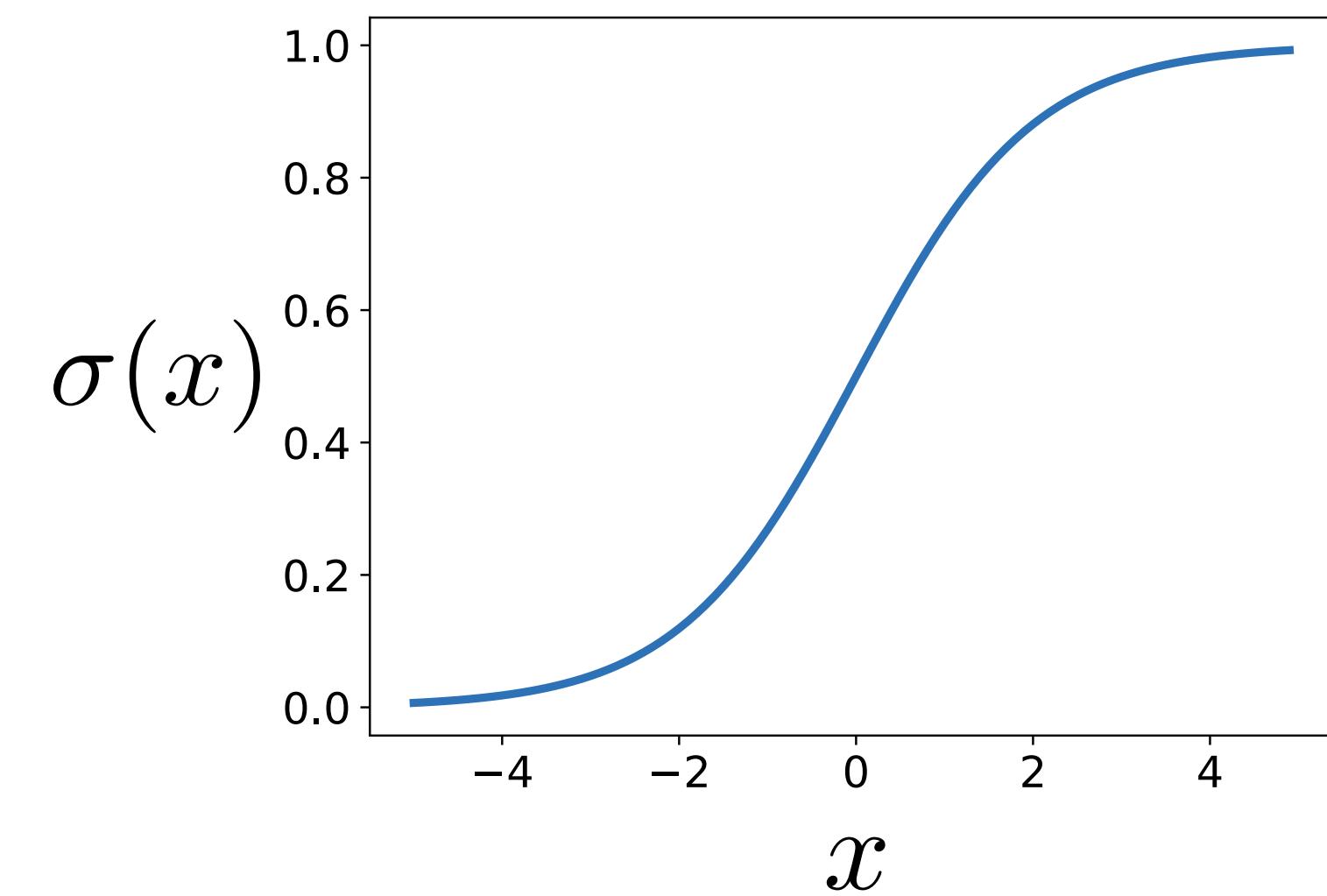
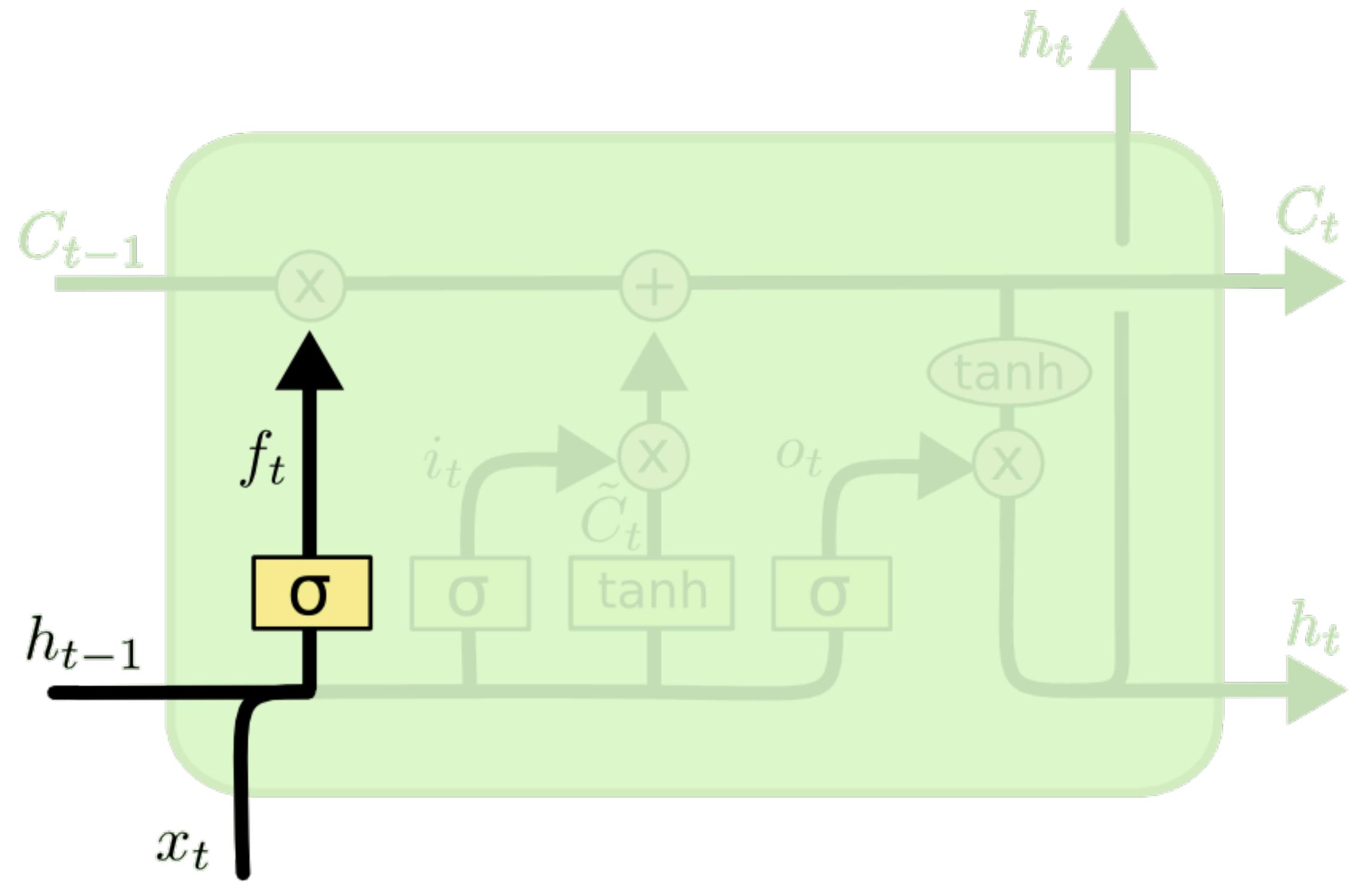




[Slide derived from Chris Olah: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>]



$C_t = \text{Cell state}$

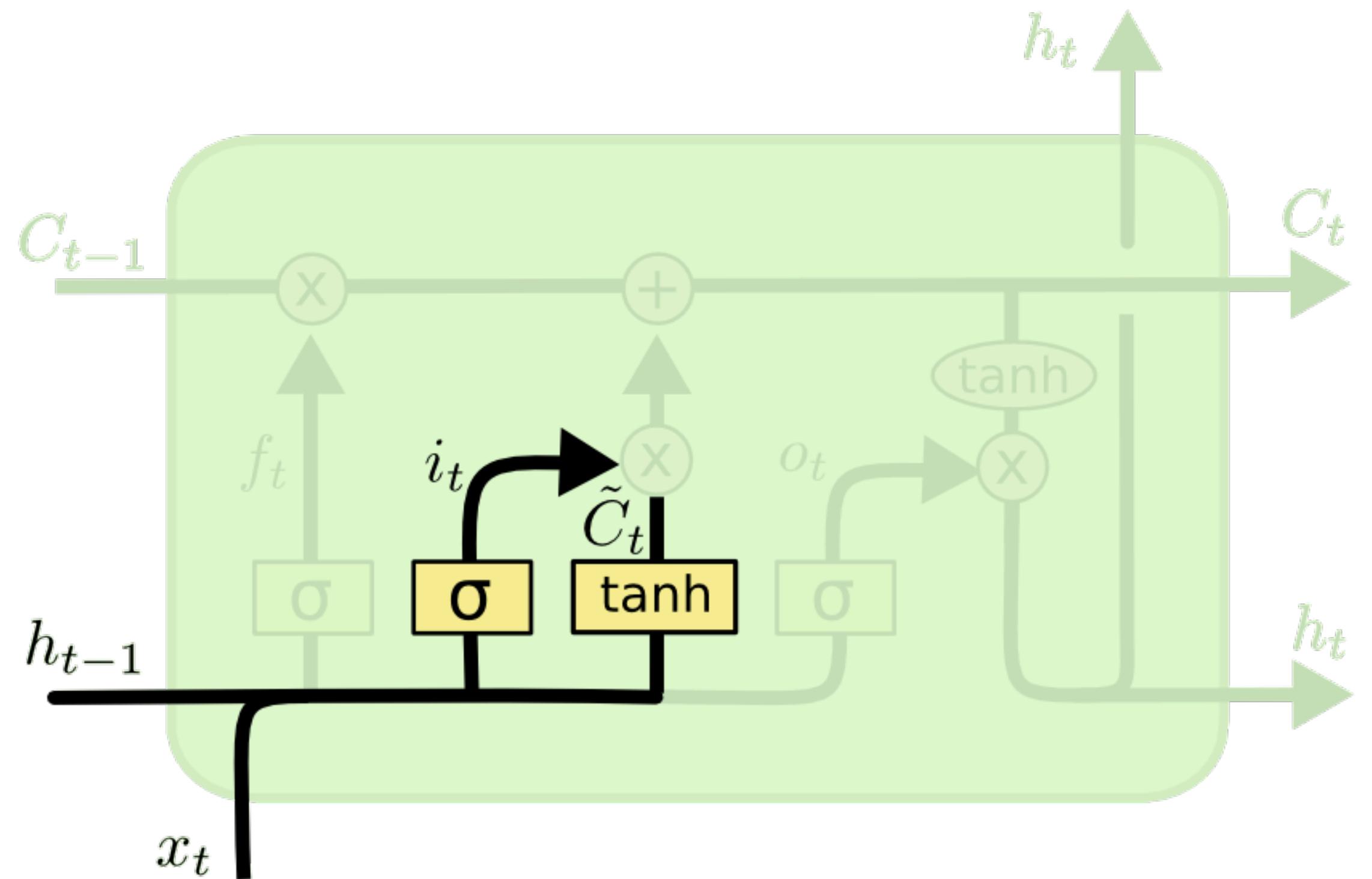


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

Decide what information to throw away from the cell state.

Each element of cell state is multiplied by ~ 1 (remember) or ~ 0 (forget).

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



which indices to write to

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

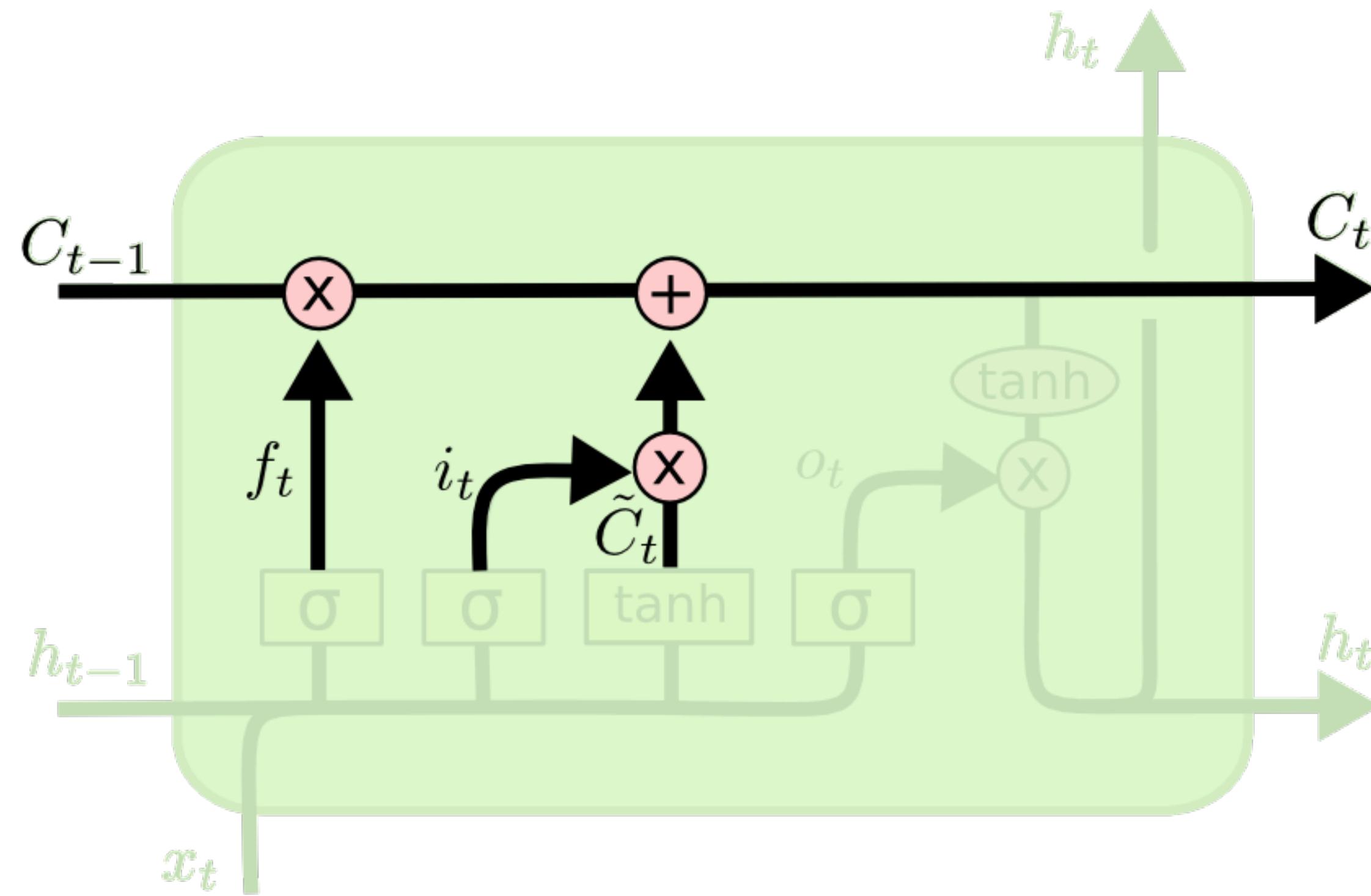
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

what to write to those indices

Decide what new information to add to the cell state.

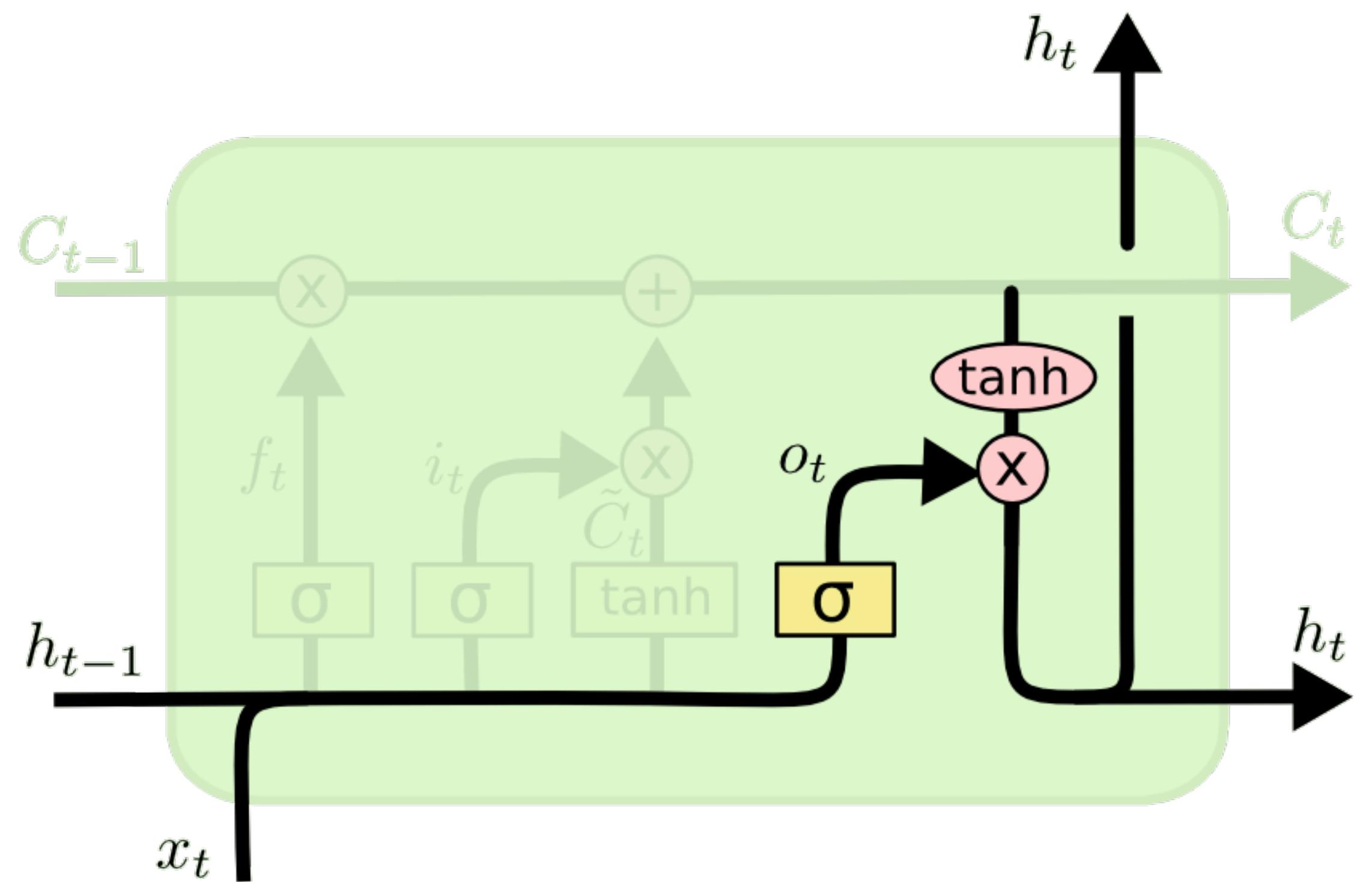
58

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



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

Forget selected old information, write selected new information.



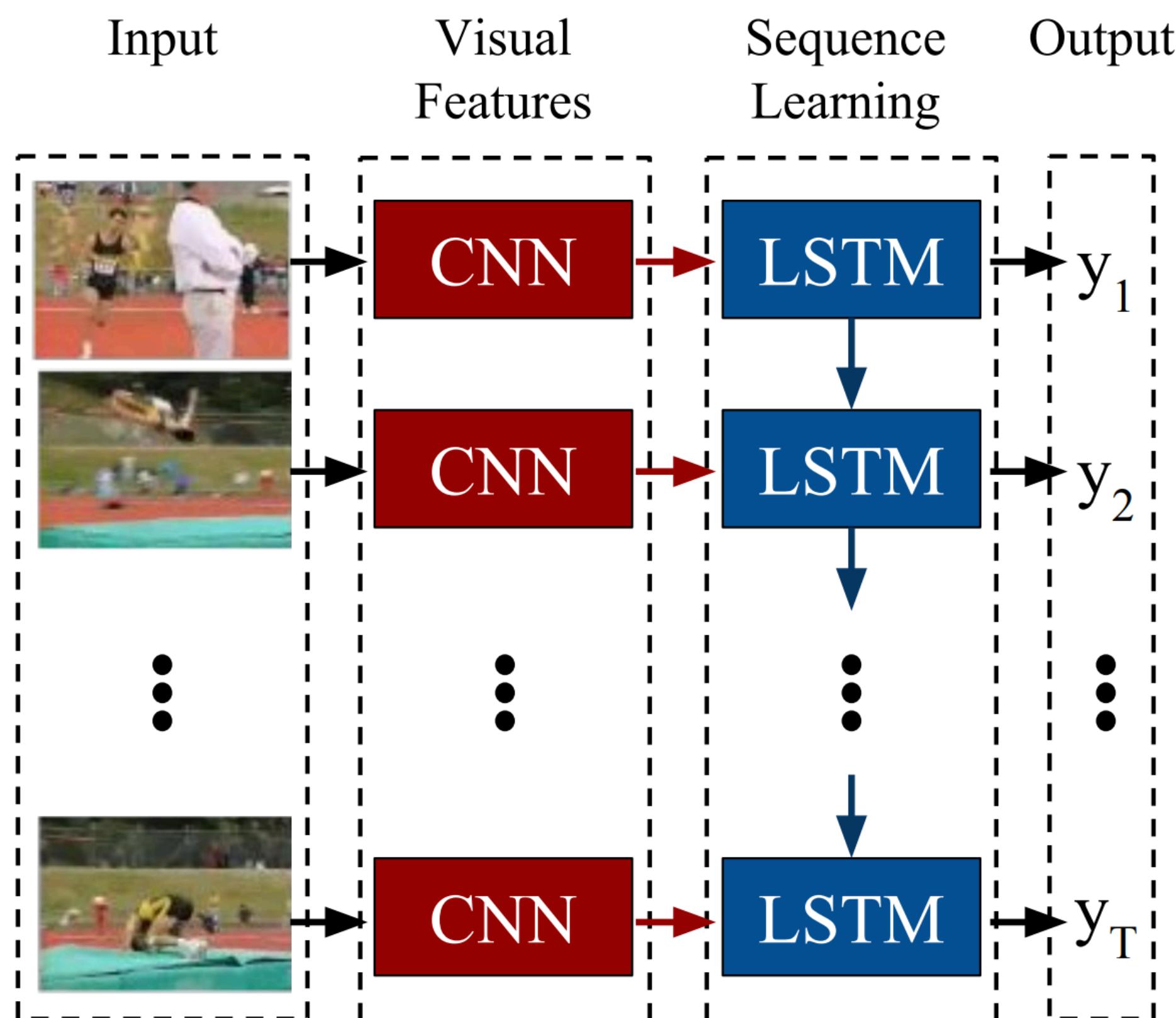
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

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

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 LSTMs



Activity Recognition
Sequences in the Input

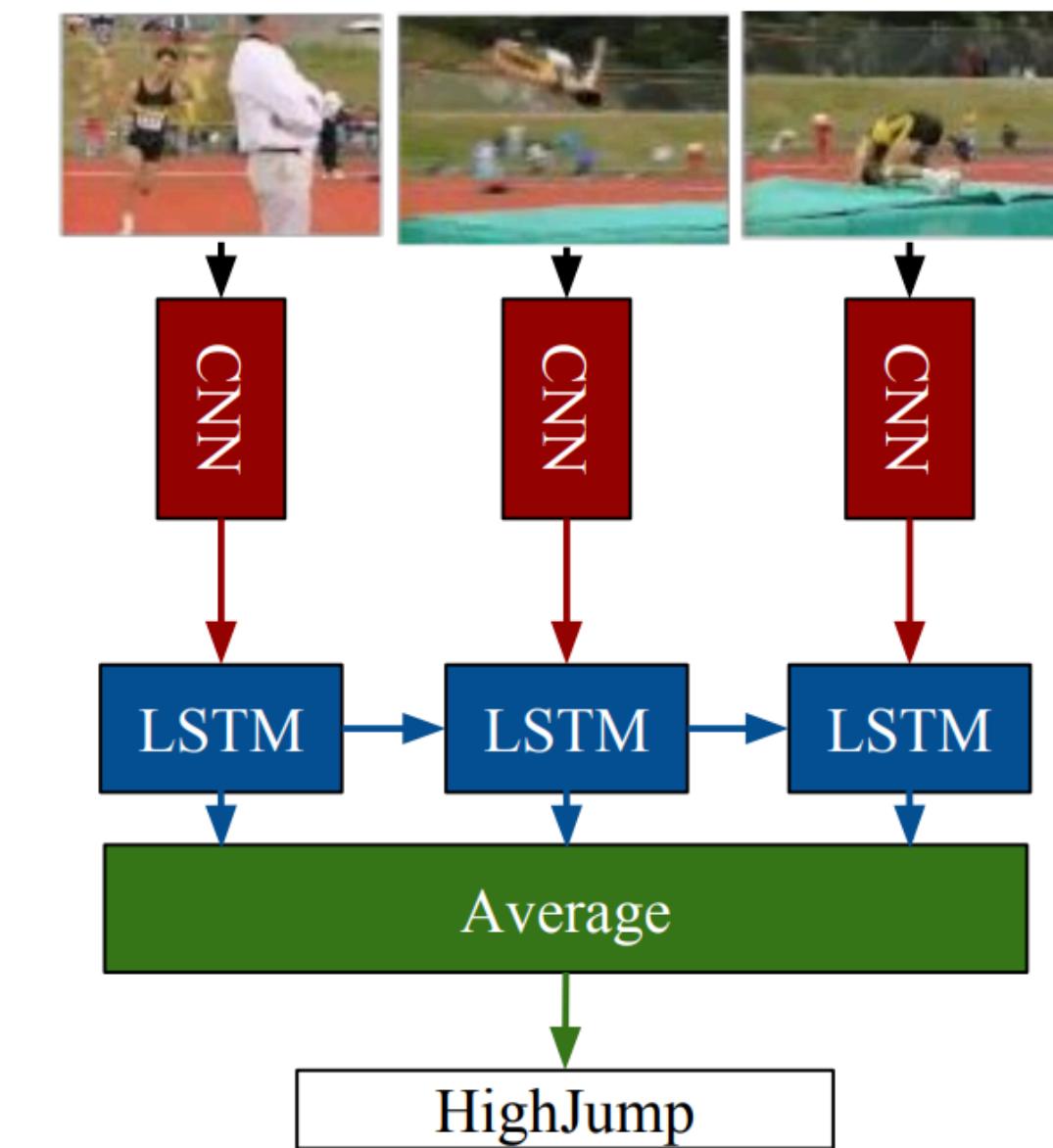
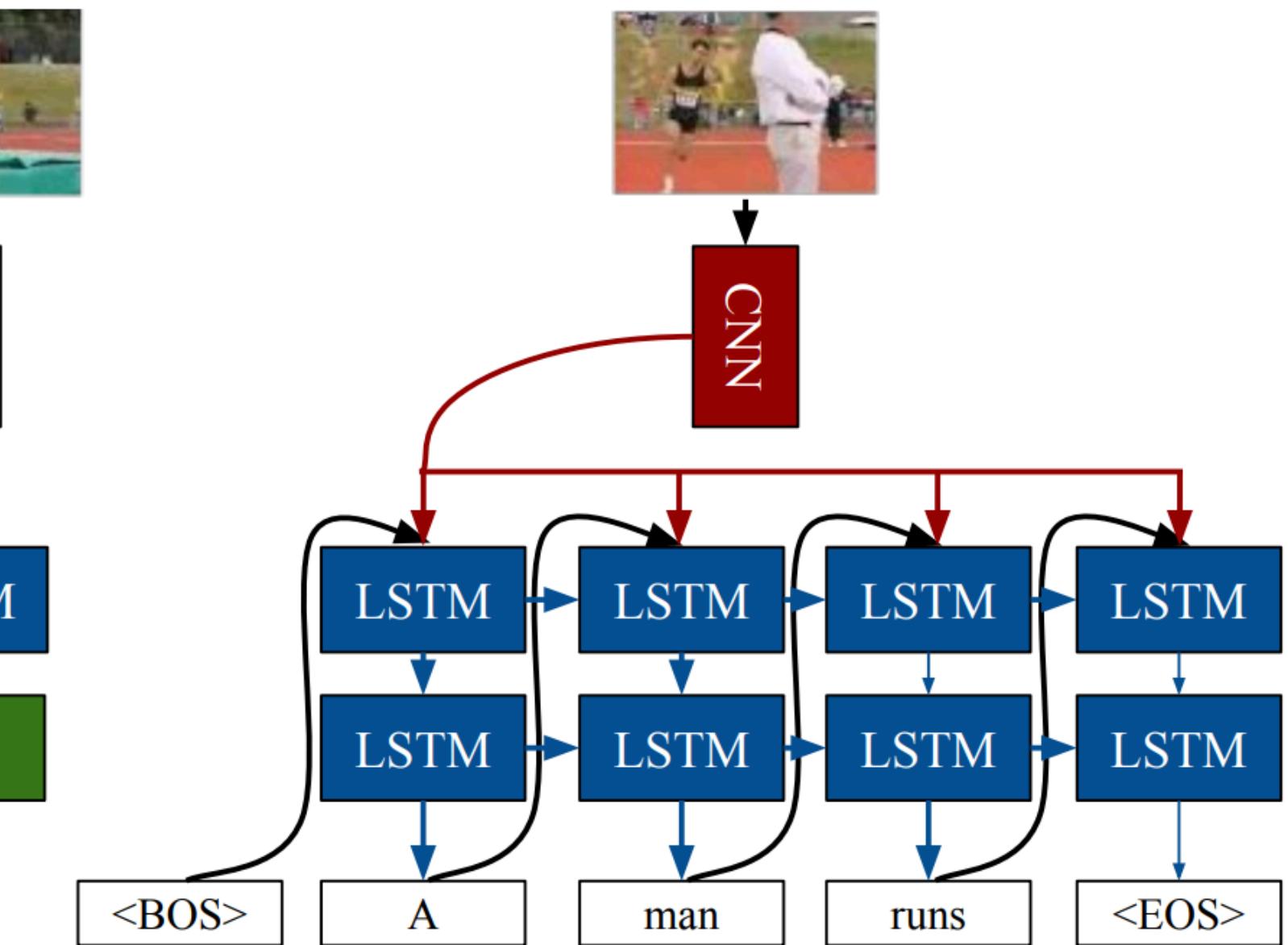


Image Captioning
Sequences in the Output



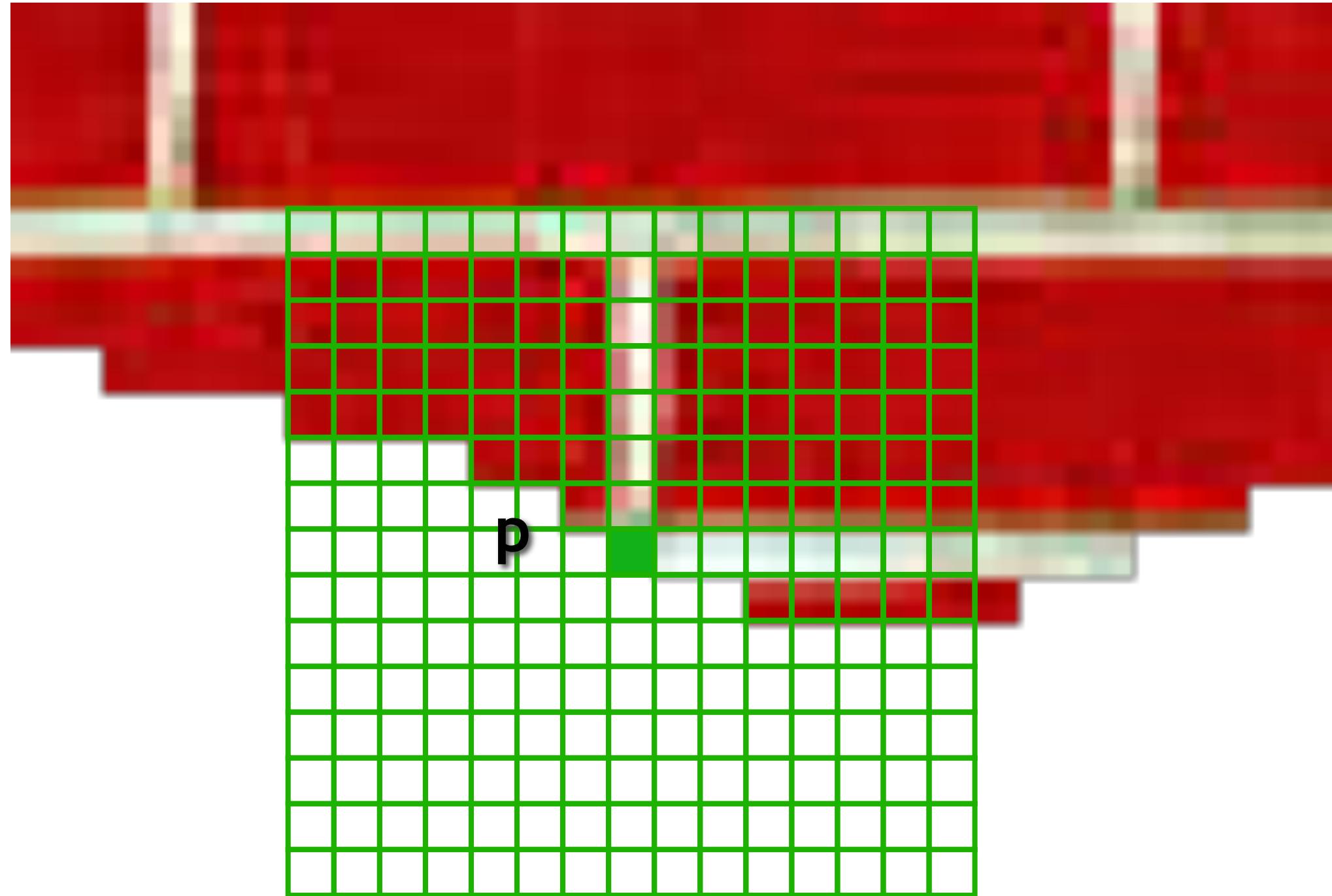
The problem of long-range dependences

Other methods exist that do directly link old “memories” (observations or hidden states) to future predictions:

- Temporal convolutions (but temporal horizon often short)
- Attention (see <https://arxiv.org/abs/1706.03762>)
- Memory networks (see <https://arxiv.org/abs/1410.3916>)

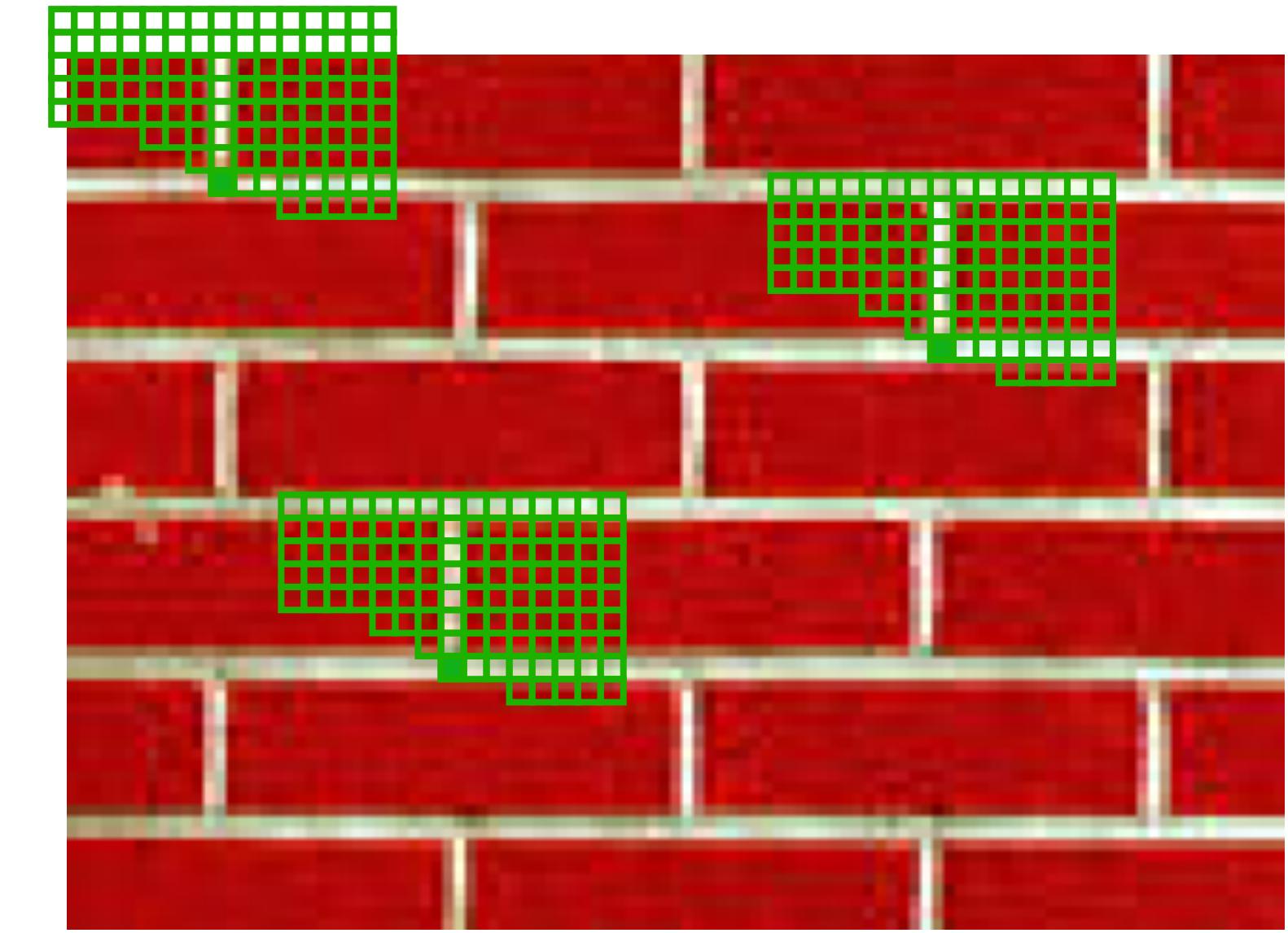
Other uses for sequence models

Texture synthesis by non-parametric sampling



Synthesizing a pixel

non-parametric
sampling



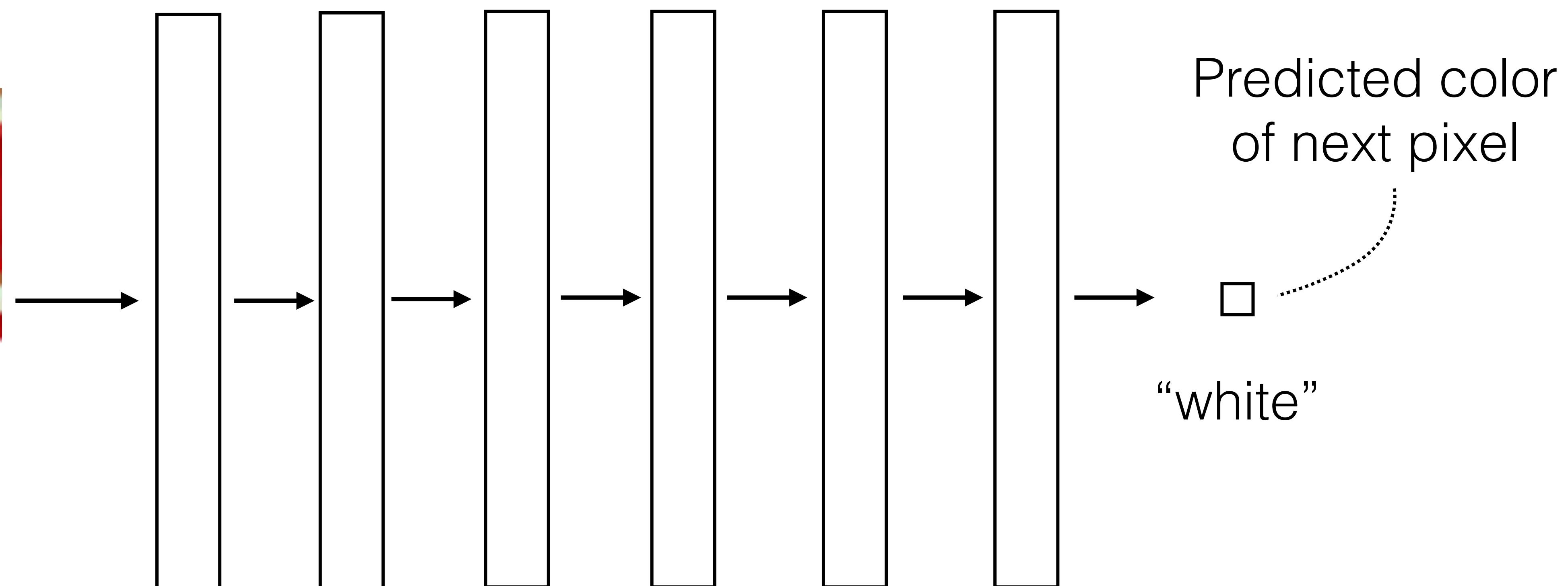
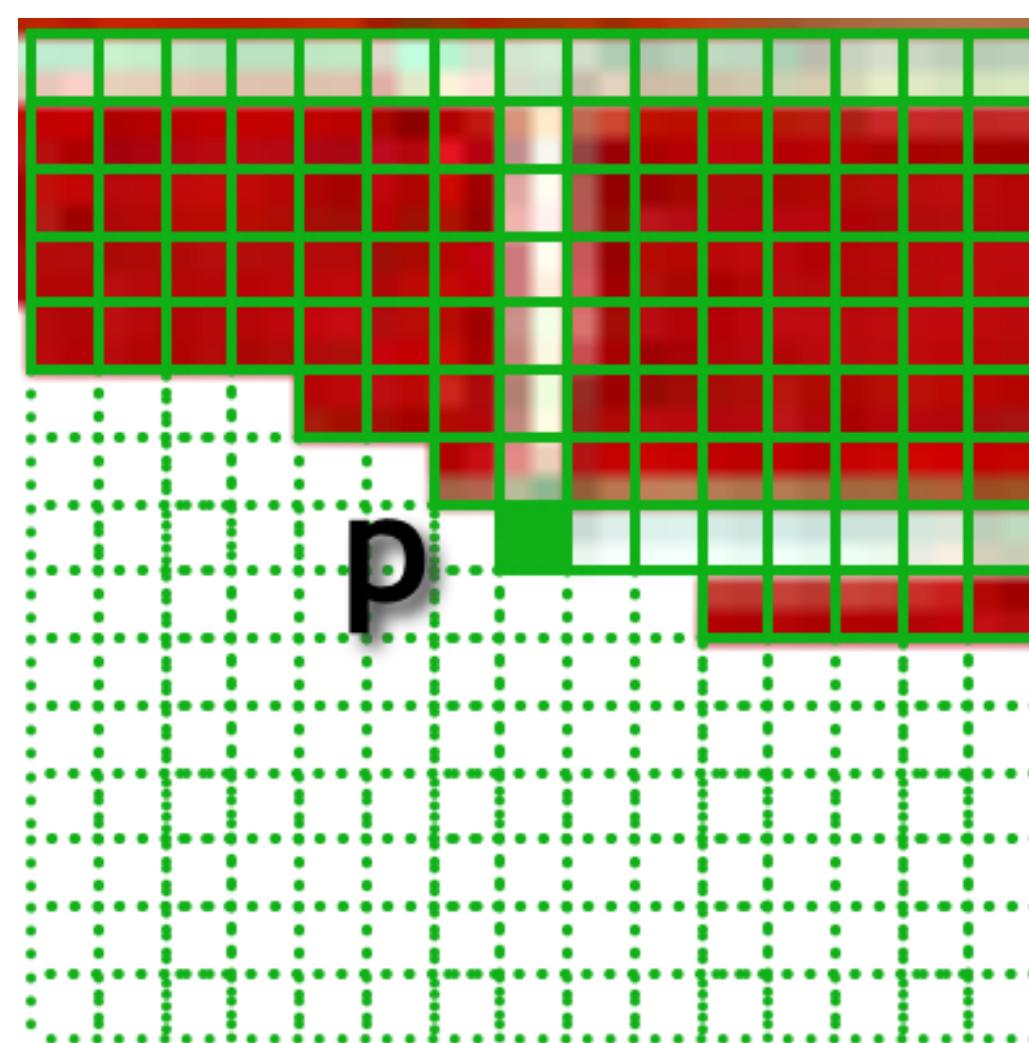
Input image

Models $P(p|N(p))$

[Efros & Leung 1999]

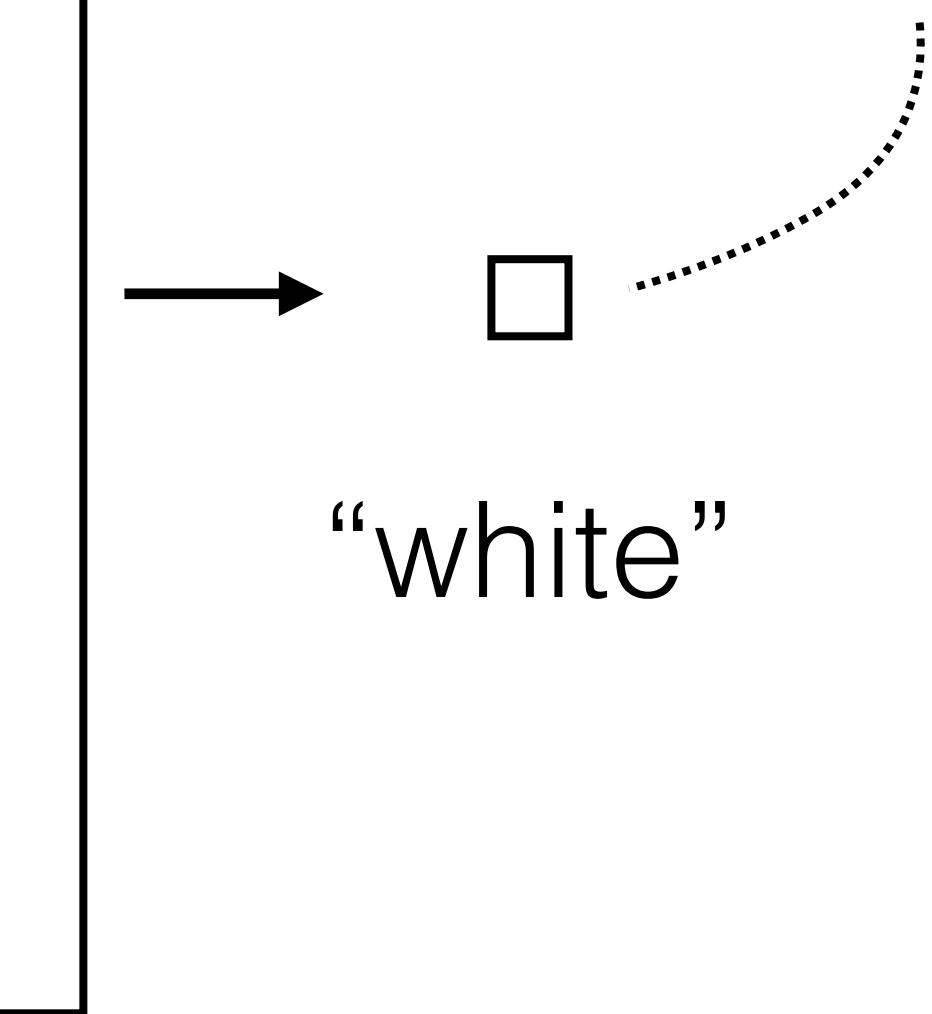
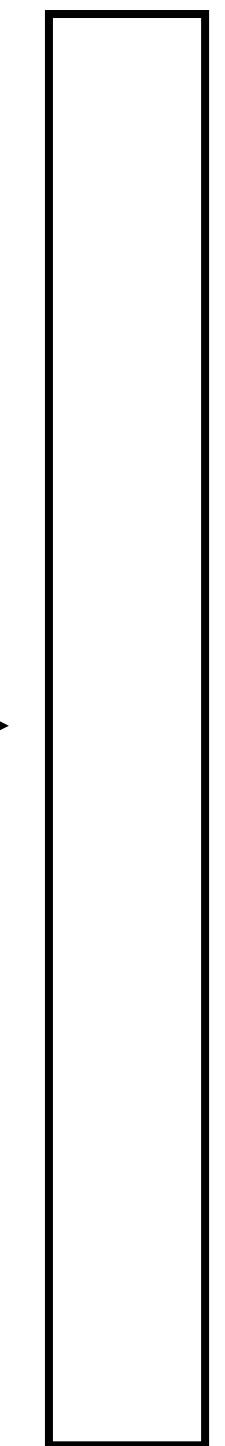
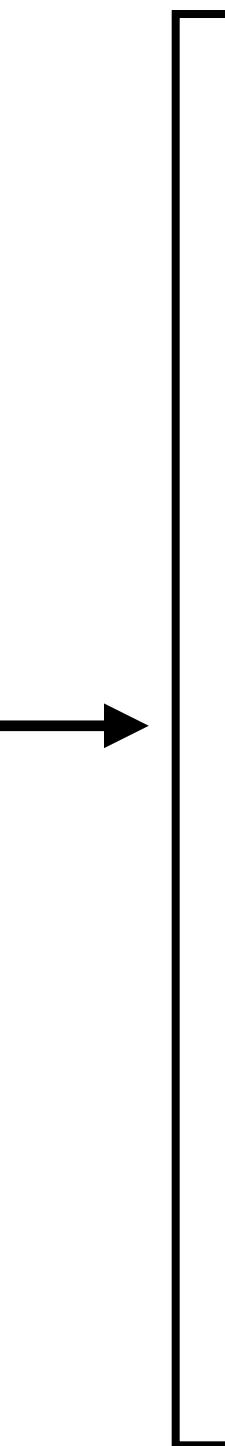
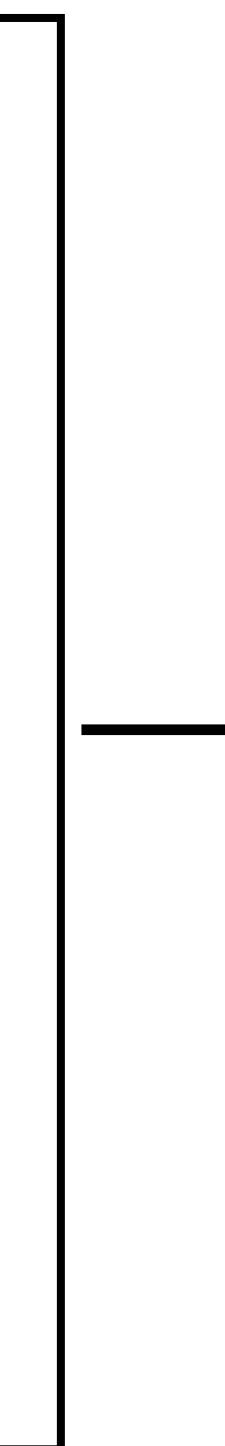
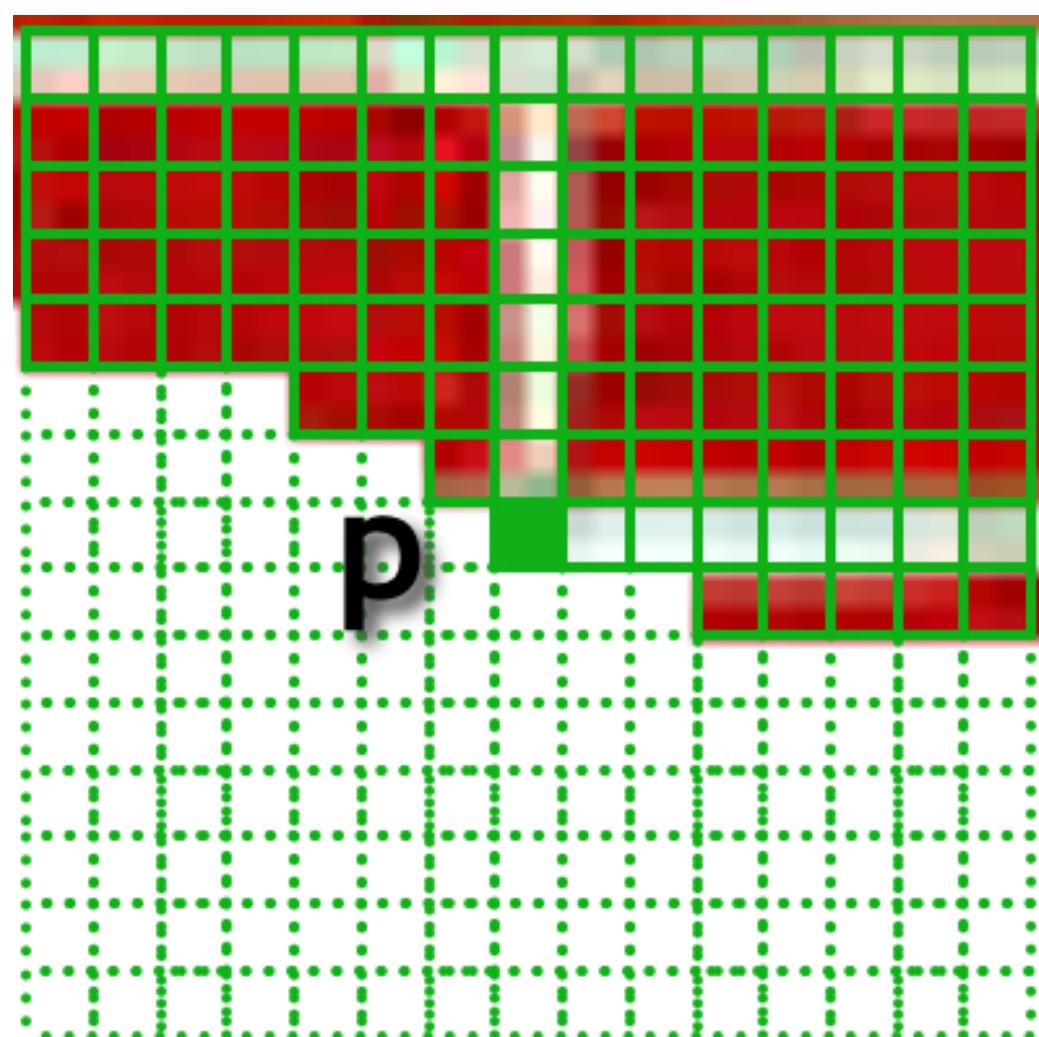
Texture synthesis with a deep net

Input partial
image



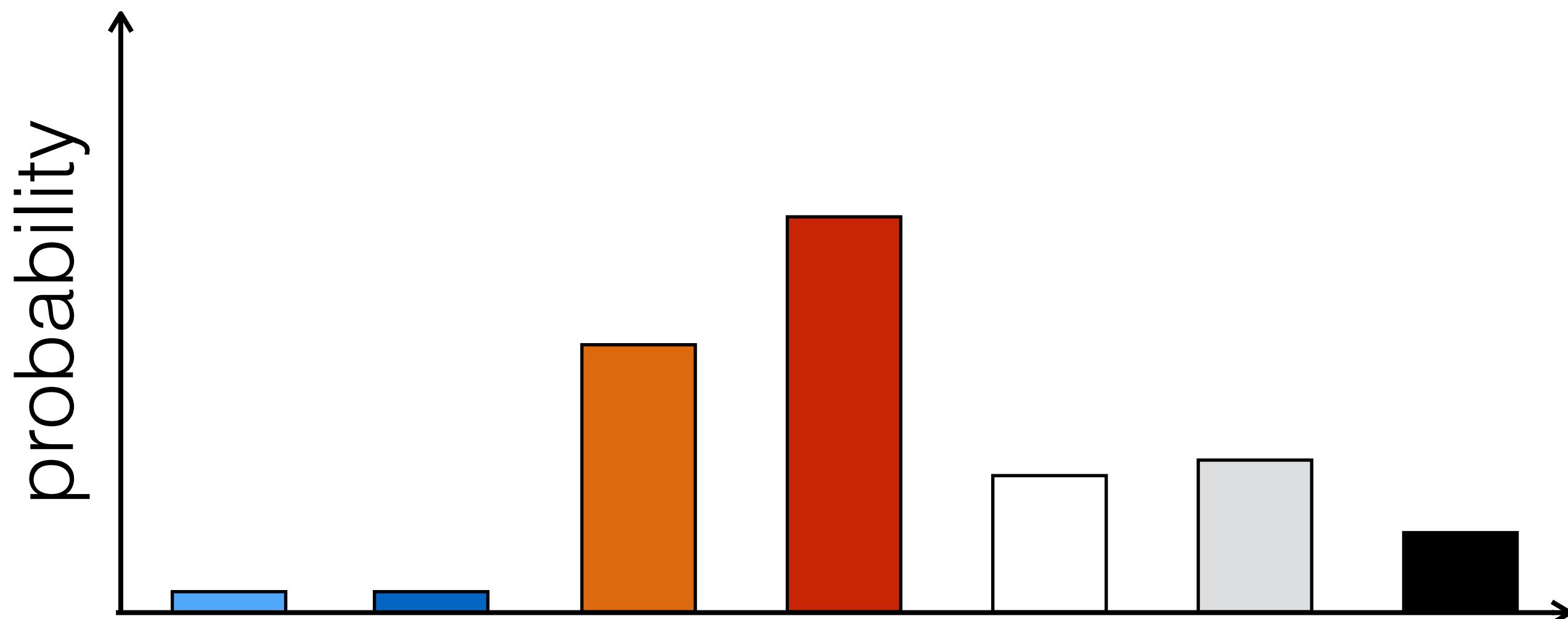
[PixelRNN, PixelCNN, van der Oord et al. 2016]

Input partial
image



Predicted color
of next pixel

Represent colors as discrete categories



- One class for each intensity value, e.g. 256 classes
- One label per color channel

And we can interpret the learner as modeling $P(\text{next pixel} \mid \text{previous pixels})$:

Softmax regression (a.k.a. multinomial logistic regression)

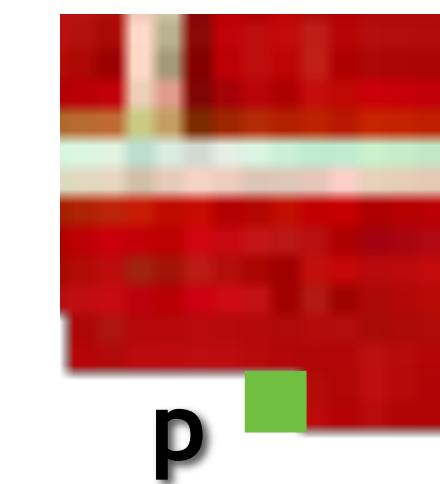
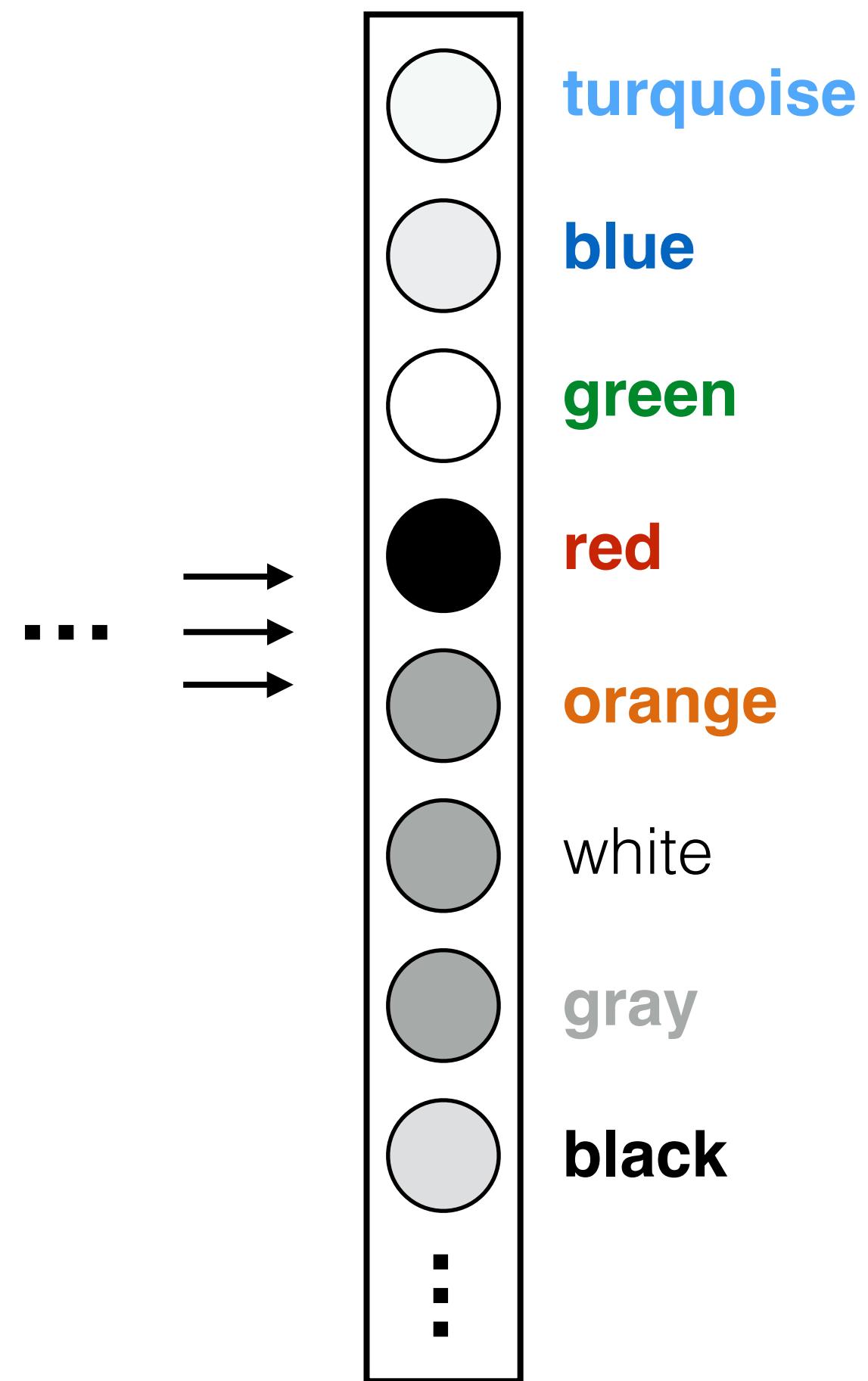
$$\hat{\mathbf{y}} \equiv [P_{\theta}(Y = 1 | X = \mathbf{x}), \dots, P_{\theta}(Y = K | X = \mathbf{x})] \leftarrow \text{predicted probability of each class given input } \mathbf{x}$$

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k \leftarrow \text{picks out the -log likelihood of the ground truth class } \mathbf{y \text{ under the model prediction } \hat{\mathbf{y}}}$$

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N H(\mathbf{y}_i, \hat{\mathbf{y}}_i) \leftarrow \text{max likelihood learner!}$$

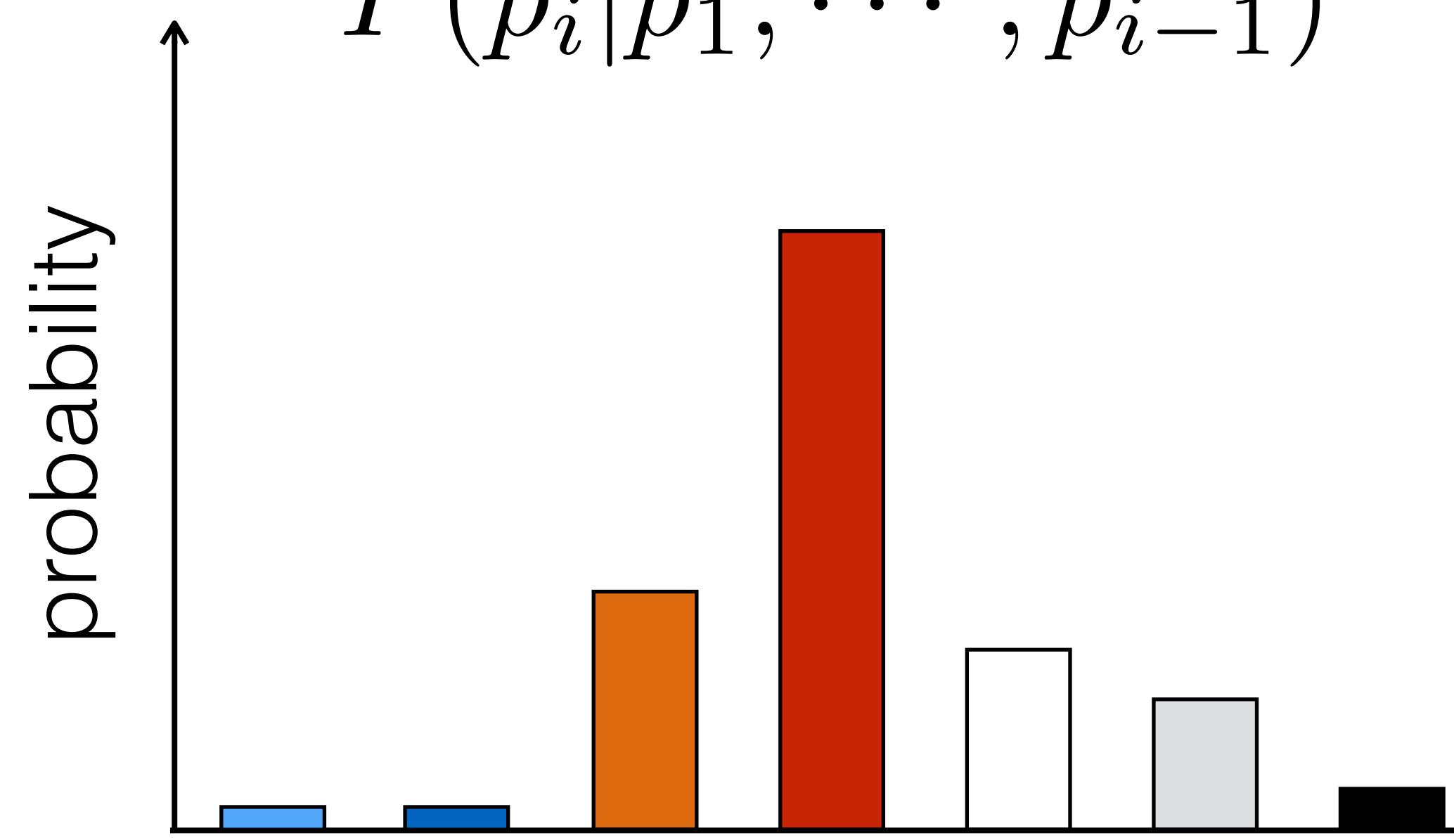
68

Network output

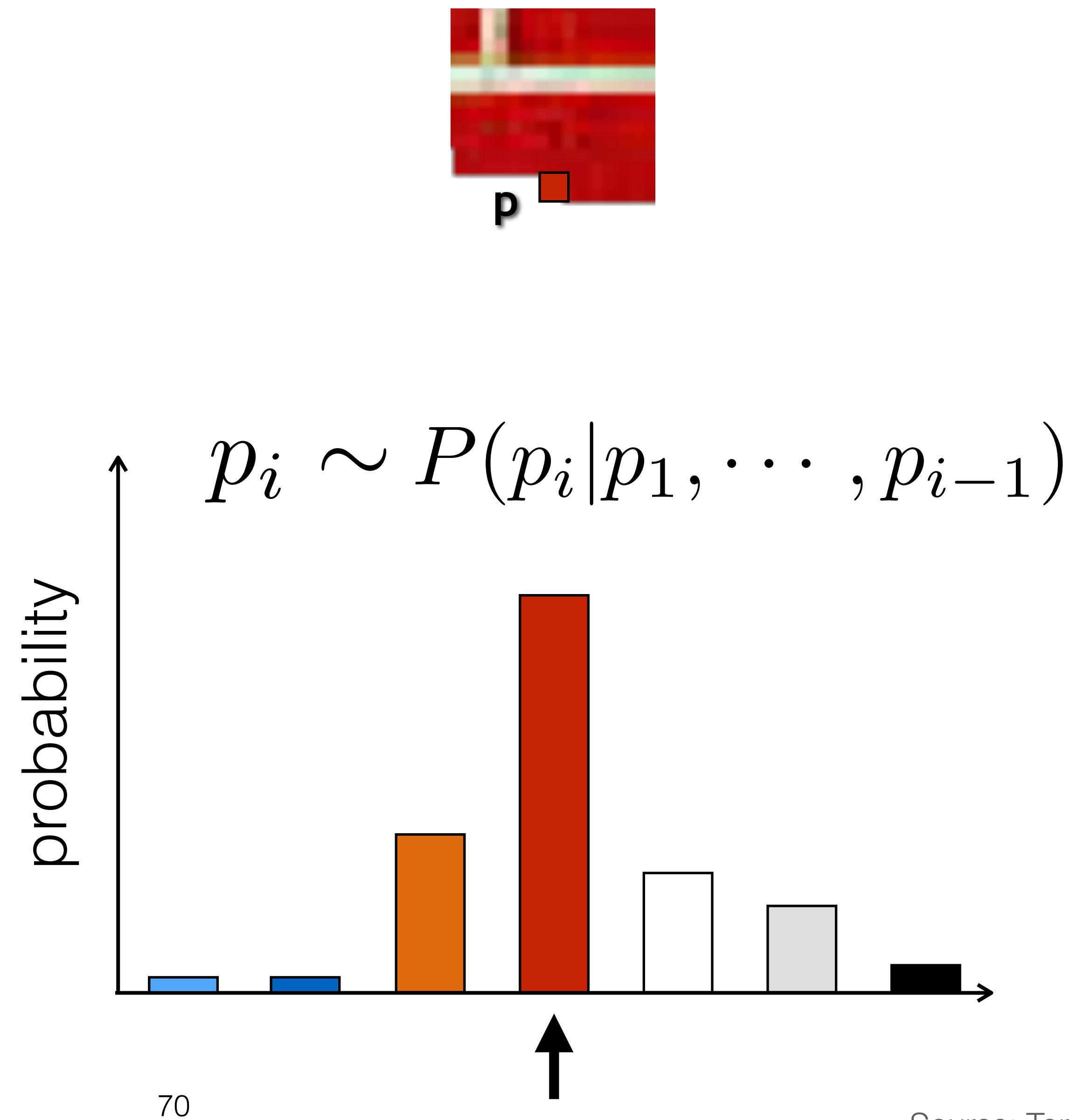
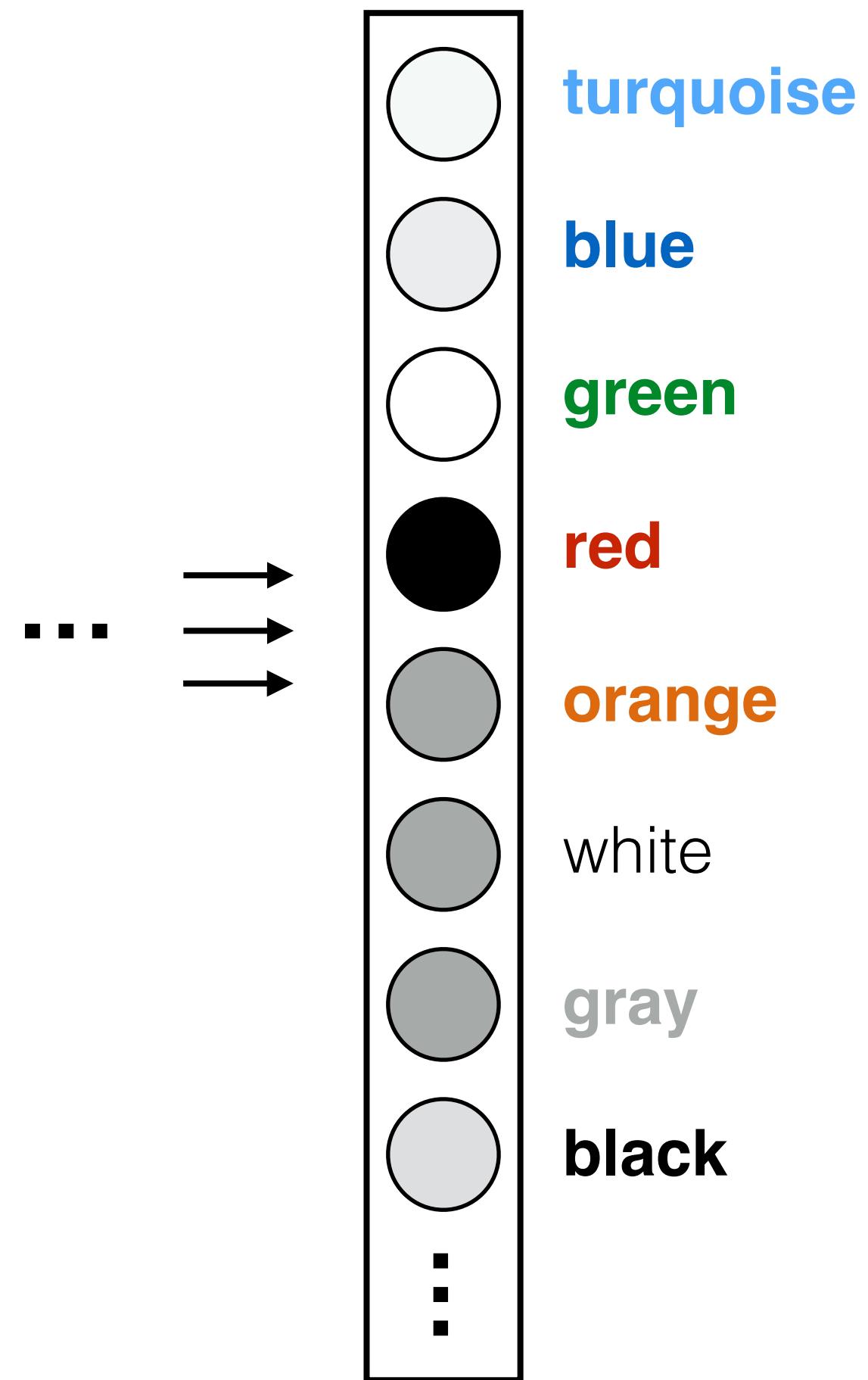


$P(\text{next pixel} \mid \text{previous pixels})$

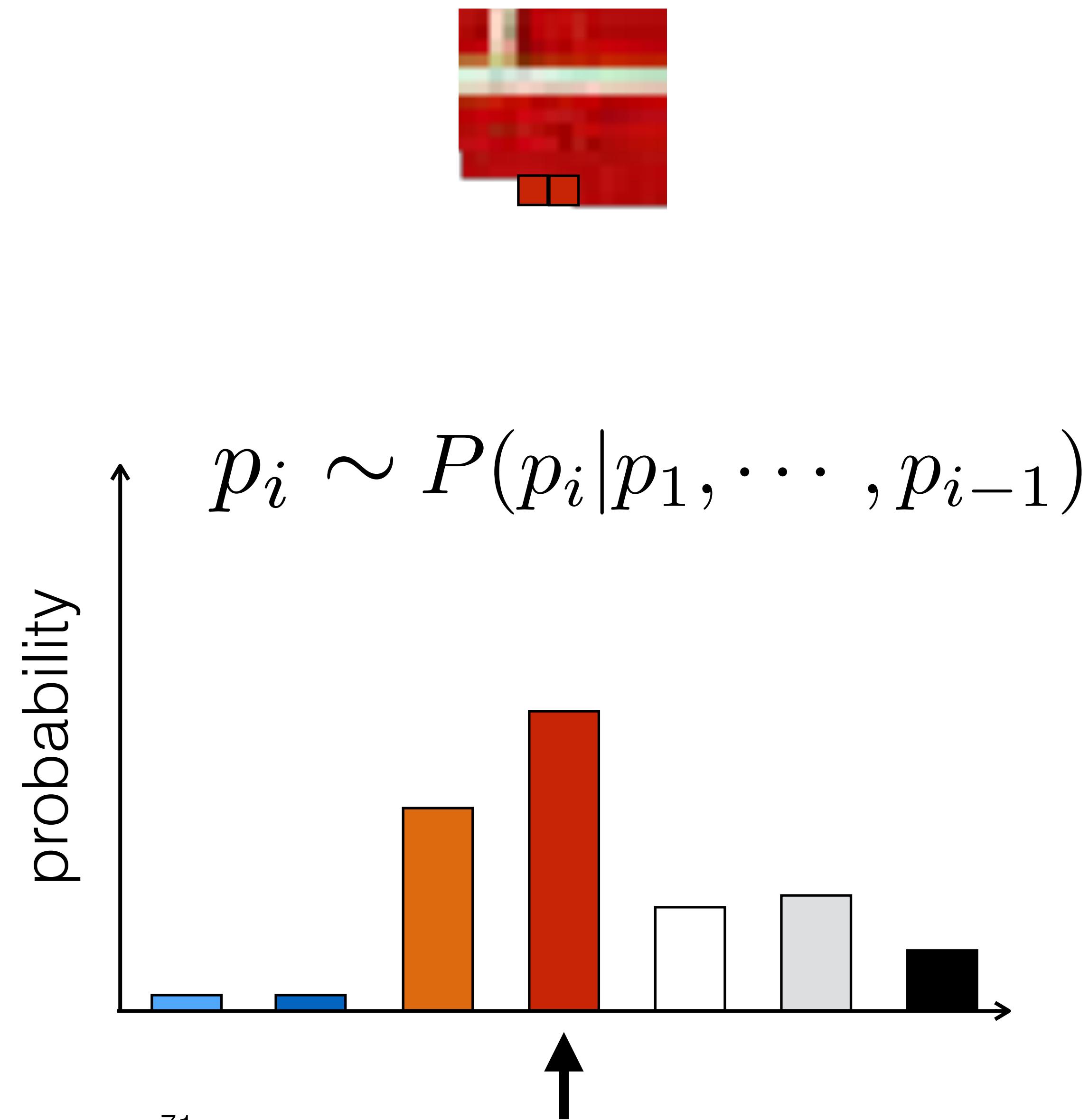
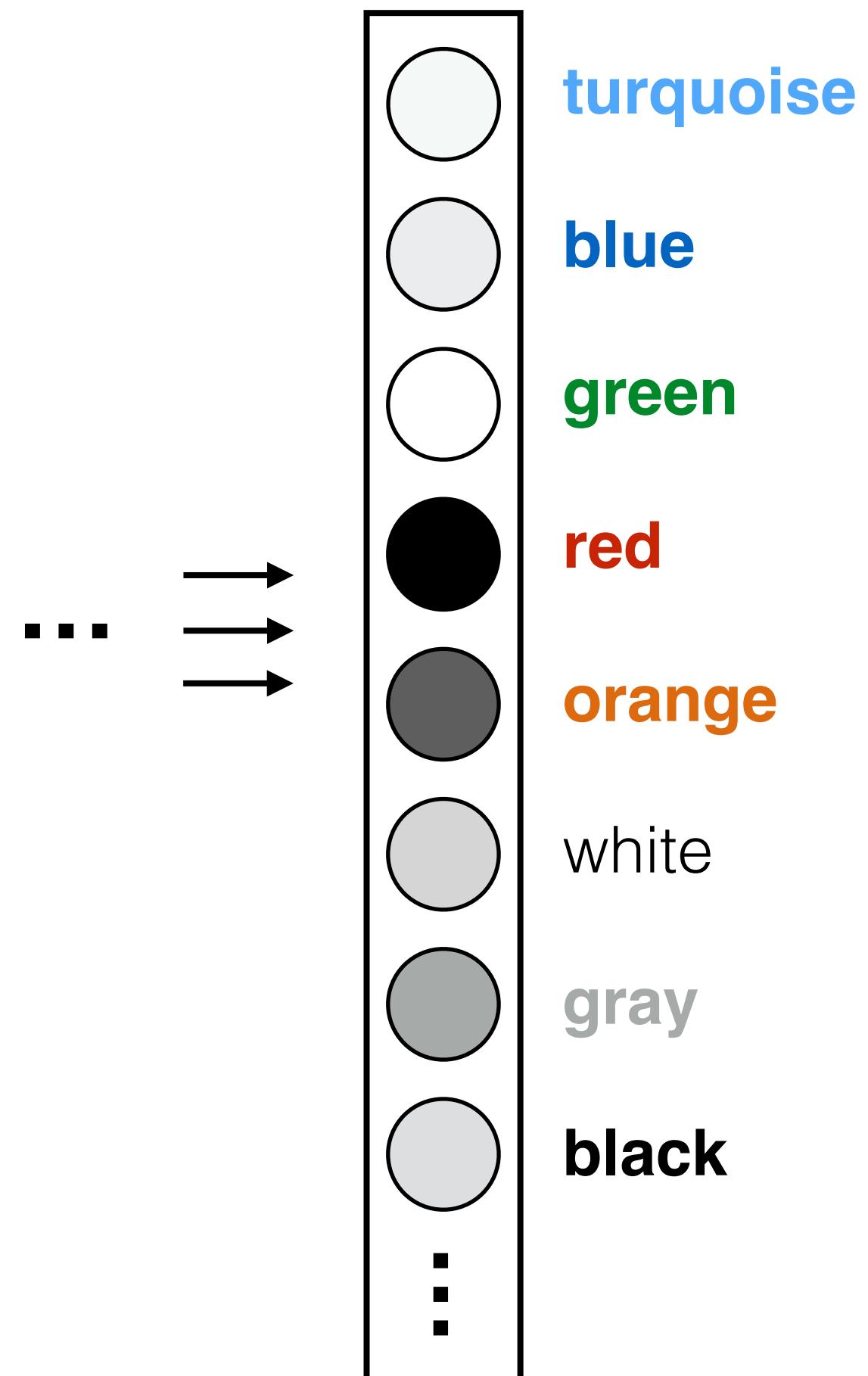
$$P(p_i \mid p_1, \dots, p_{i-1})$$



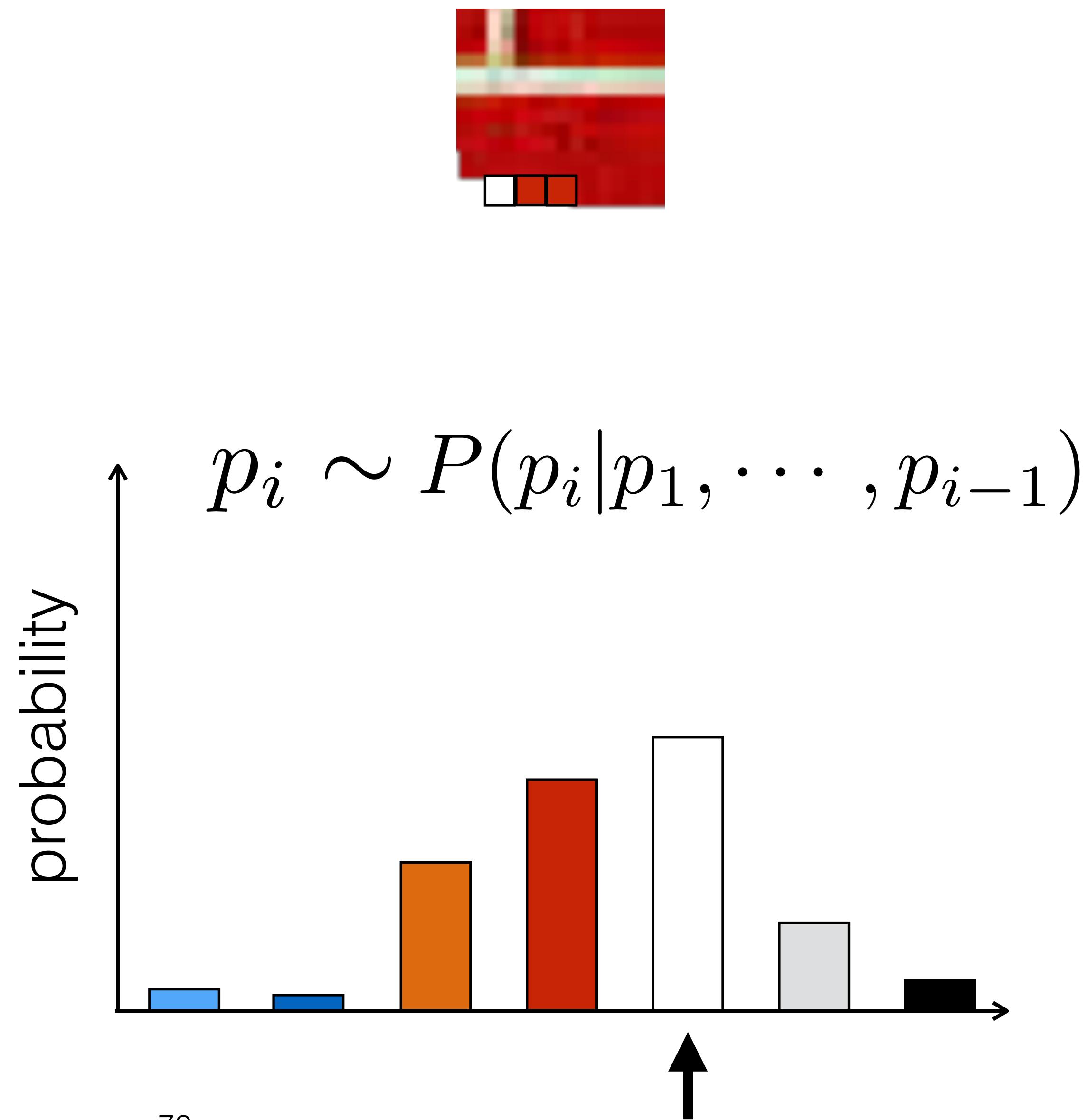
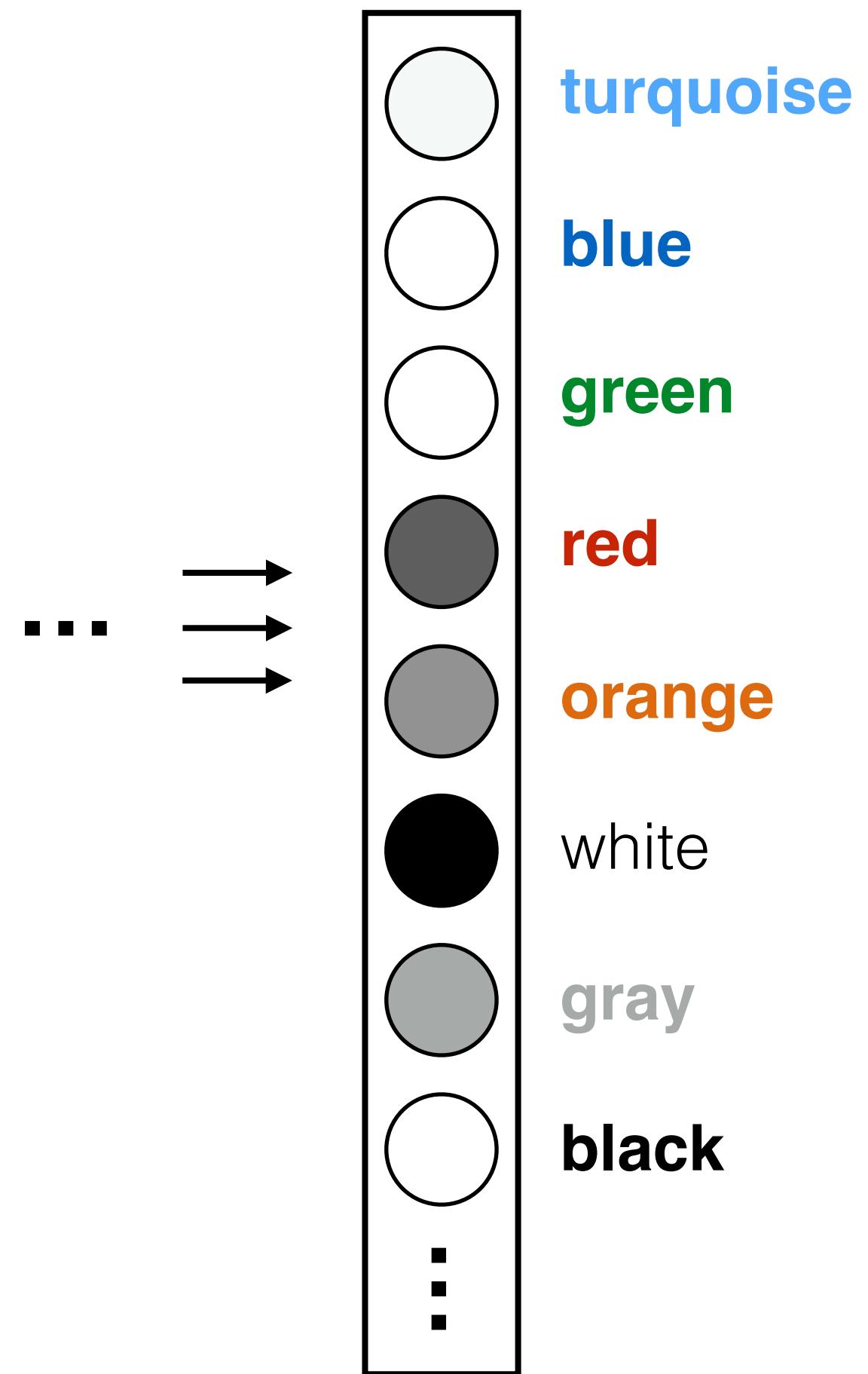
Network output



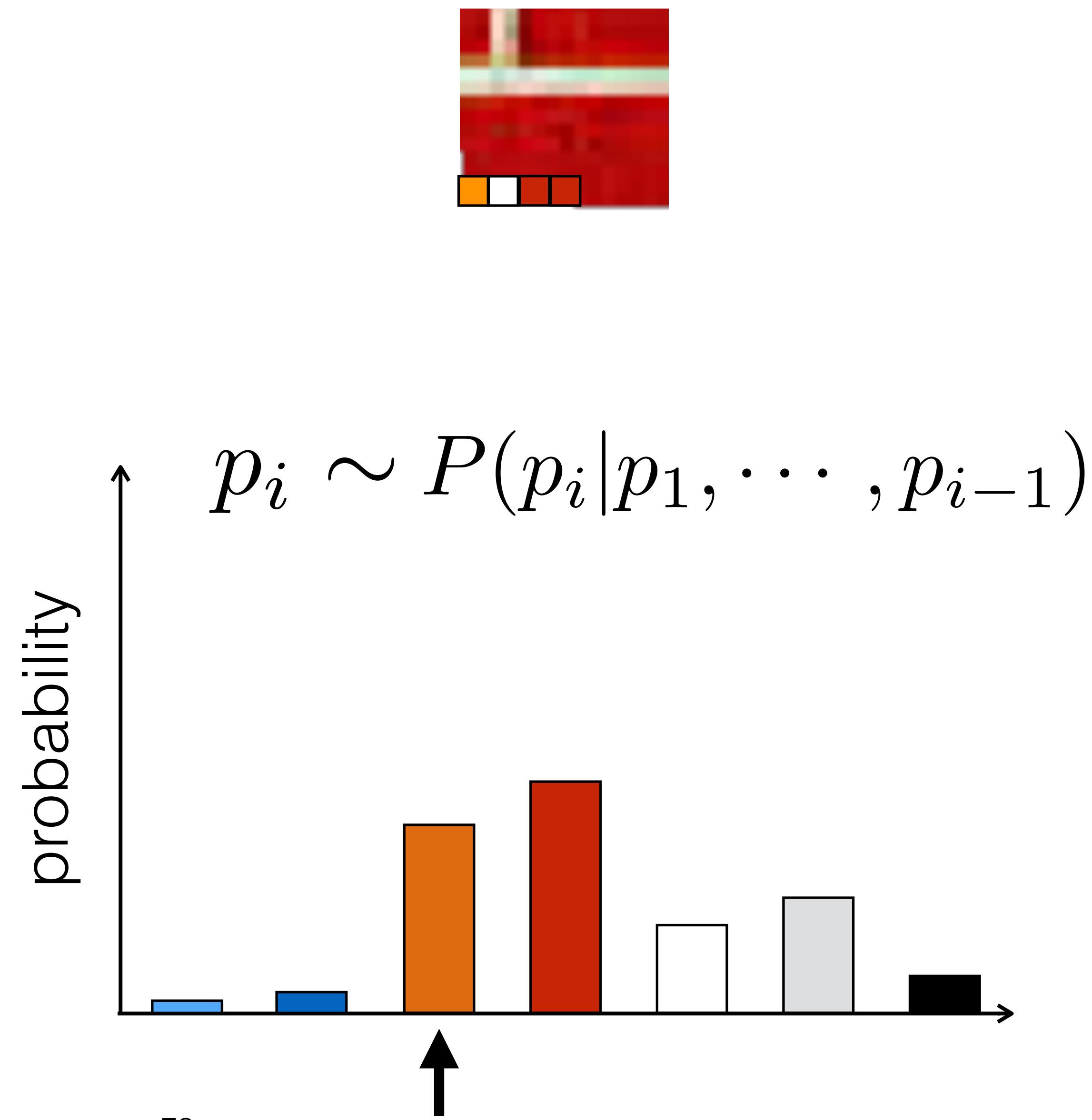
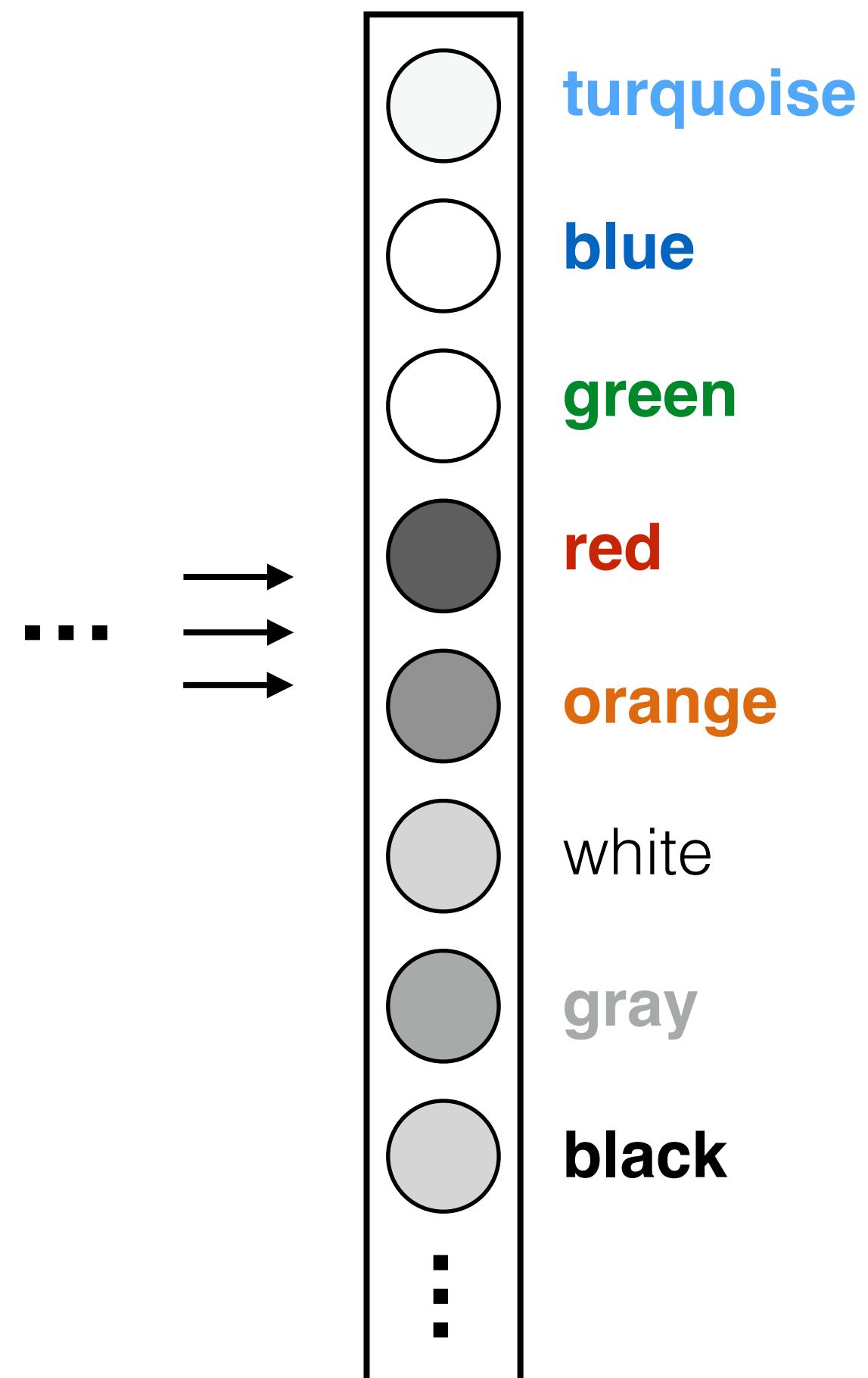
Network output



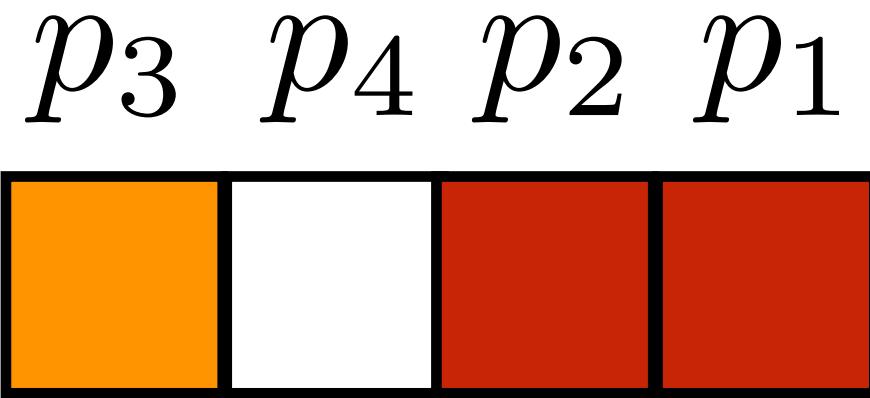
Network output



Network output



$$p_1 \sim P(p_1)$$



$$p_2 \sim P(p_2|p_1)$$

$$p_3 \sim P(p_3|p_1, p_2)$$

$$p_4 \sim P(p_4|p_1, p_2, p_3)$$

$$\{p_1, p_2, p_3, p_4\} \sim P(p_4|p_1, p_2, p_3)P(p_3|p_1, p_2)P(p_2|p_1)P(p_1)$$

$$p_i \sim P(p_i|p_1, \dots, p_{i-1})$$

$$\mathbf{p} \sim \prod_{i=1}^N P(p_i|p_1, \dots, p_{i-1})$$

Autoregressive probability model

$$\mathbf{p} \sim \prod_{i=1}^N P(p_i | p_1, \dots, p_{i-1})$$

$$P(\mathbf{p}) = \prod_{i=1}^N P(p_i | p_1, \dots, p_{i-1}) \quad \leftarrow \textbf{General product rule}$$

The sampling procedure we defined above takes exact samples from the learned probability distribution (pmf).

Multiplying all conditionals evaluates the probability of a full joint configuration of pixels.

Autoregressive probability model

$$\mathbf{p} \sim P(\mathbf{p})$$

Models that allow us to sample, i.e. *generate*, images from scratch are called **generative models**.

We will see more examples in a future lecture.

Samples from PixelRNN

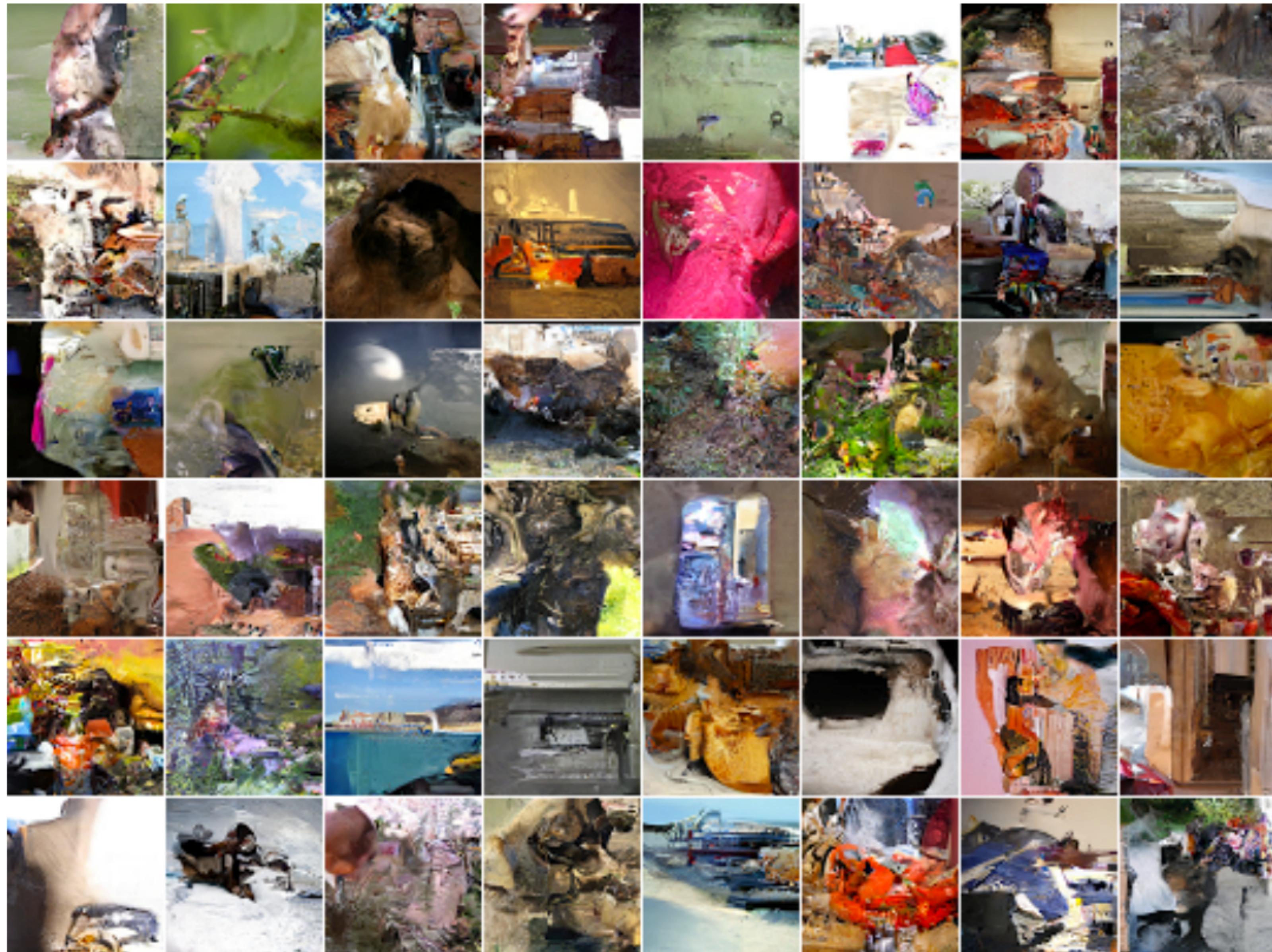


Image completions (conditional samples) from PixelRNN

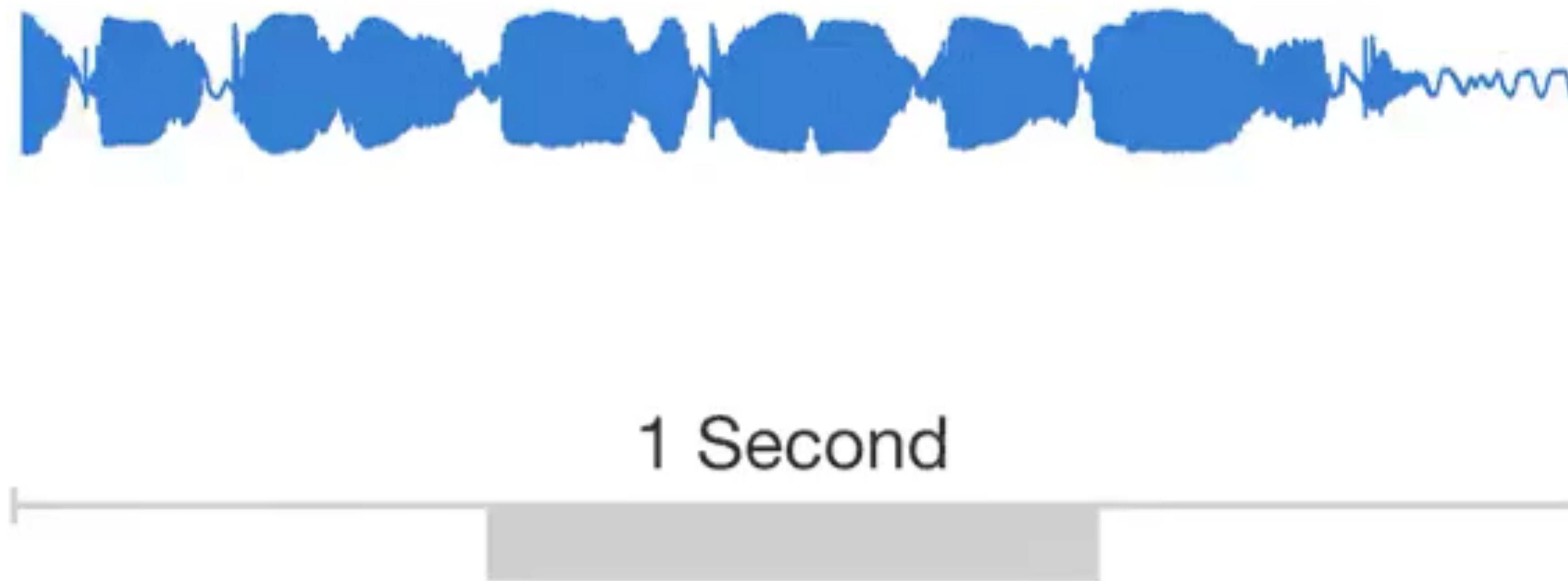
occluded

completions

original

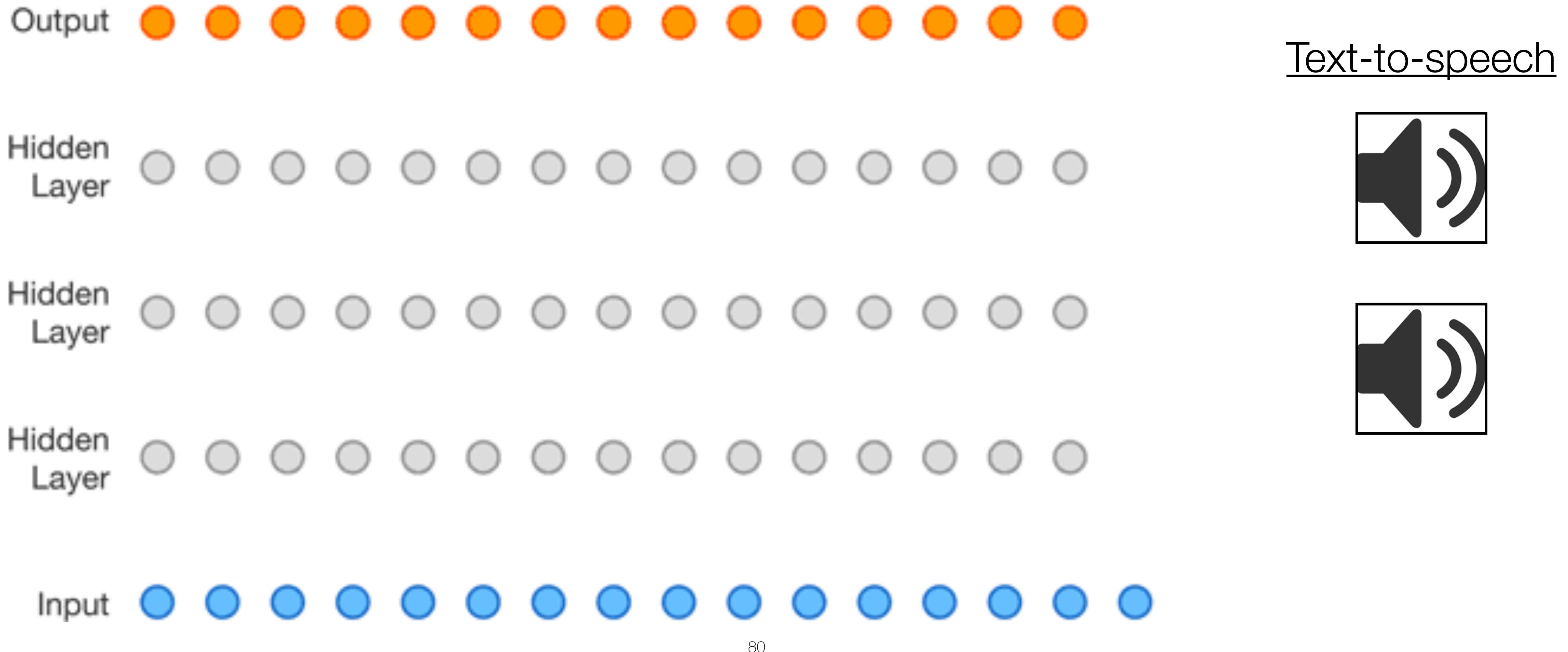


Synthesize audio with the same model



[**Wavenet**, <https://deepmind.com/blog/wavenet-generative-model-raw-audio/>]
79

Synthesize audio with the same model



[**Wavenet**, <https://deepmind.com/blog/wavenet-generative-model-raw-audio/>]

Next week: more image synthesis