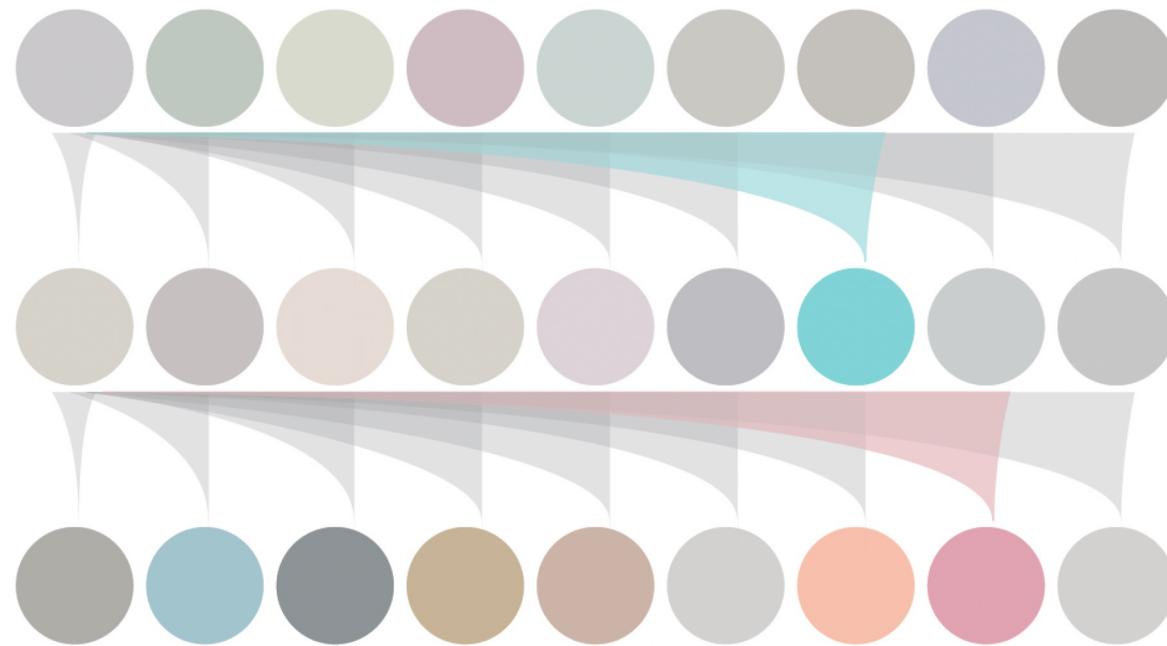


Lecture 11

Sequence Modeling



Overview

- Sequence Problems in Vision
- ConvNets for Sequence Problems
- Recurrent Neural Networks
- Attention and Transformers

Sequence Problems in Vision

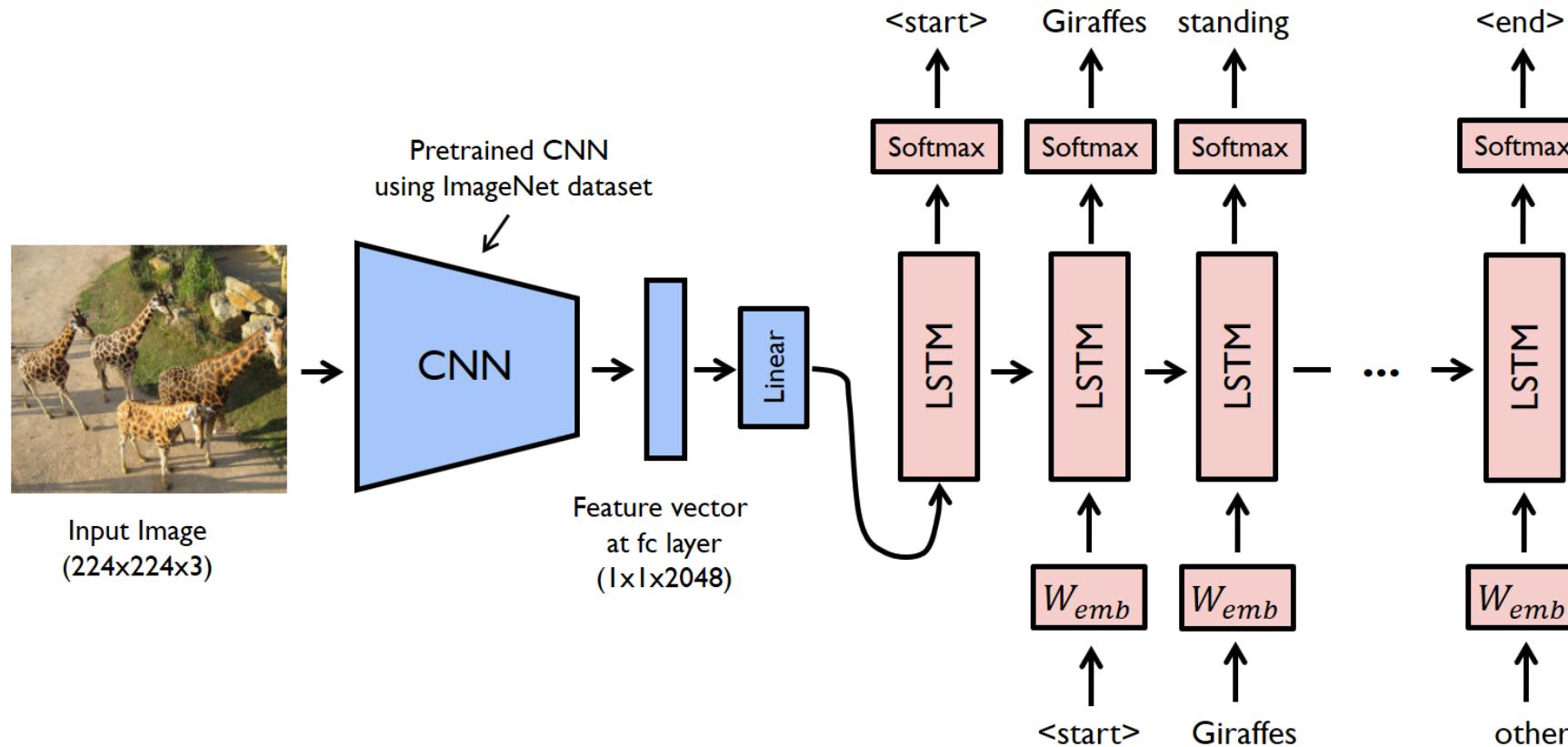
Video: Sequence of Images

- Video counterparts of 2D problems
 - classification, detection, segmentation, ...
- Video-specific problems
 - tracking, temporal reasoning, structure-from-motion, ...



Vision-Language Problems

- pixels in, language out
 - captioning, visual question answering (VQA), reasoning, referring...



Vision-Language Problems

- language in, pixels out
 - text-to-image, text-to-video, text-to-3D generation

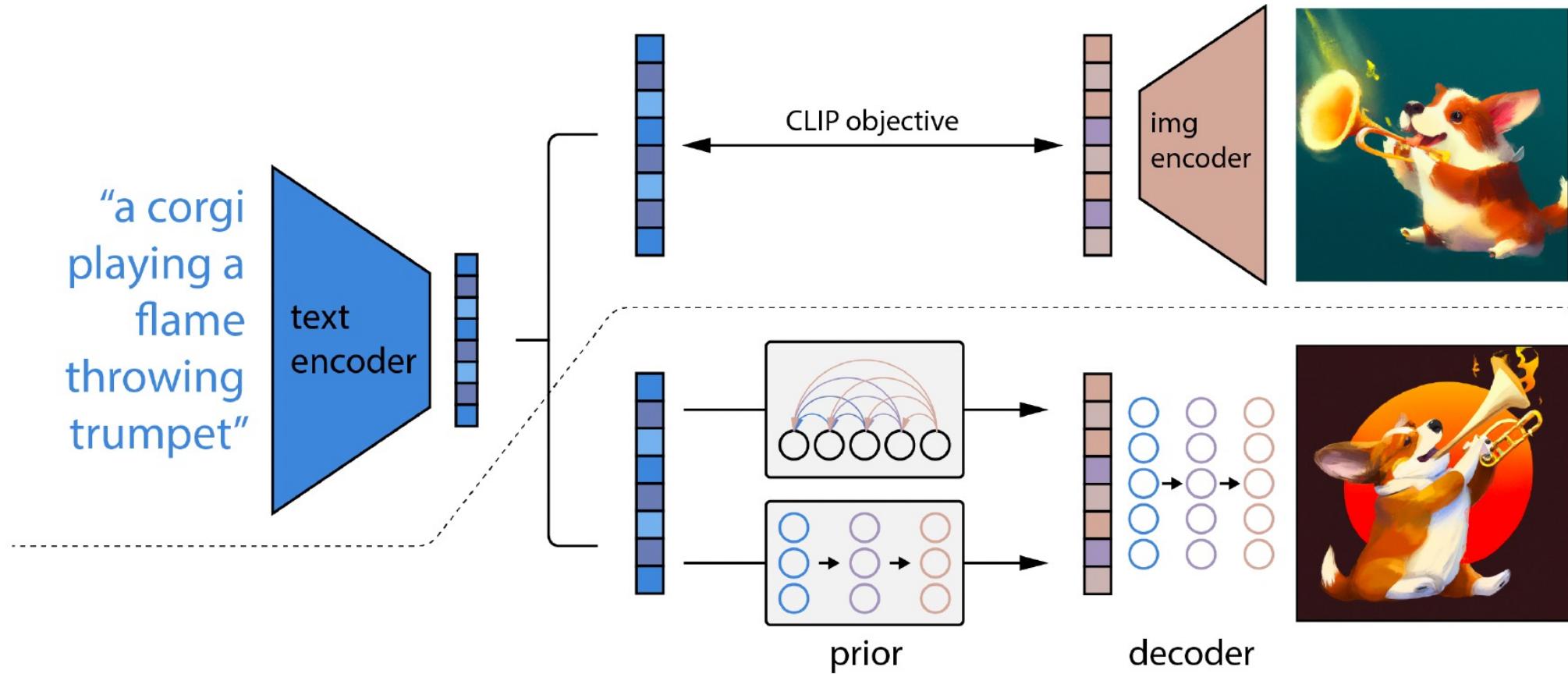


figure credit: DALL-E

Image as a Sequence

- Sequence of pixels (iGPT)

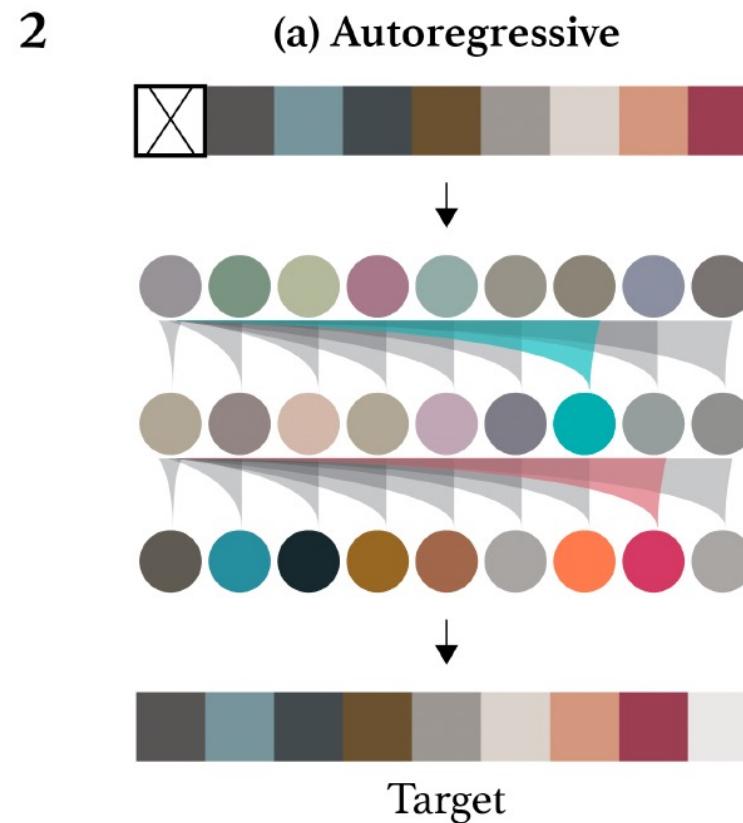
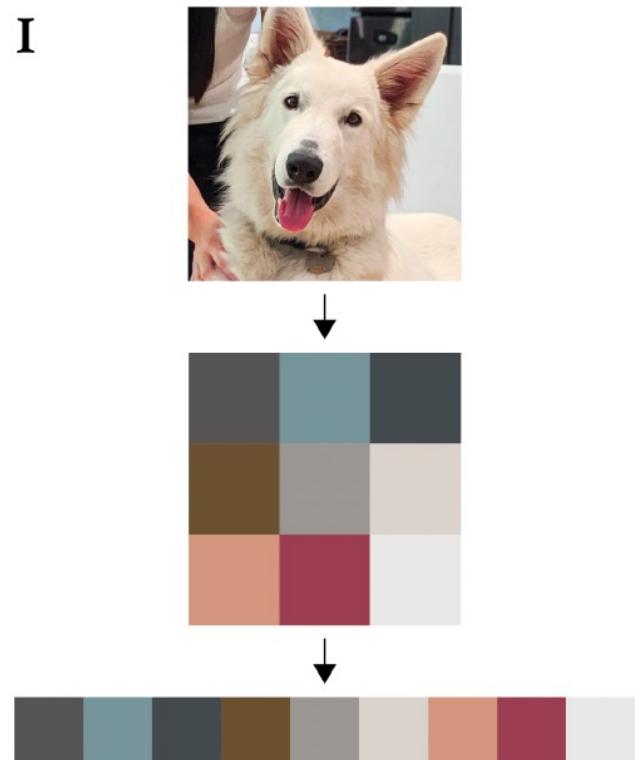
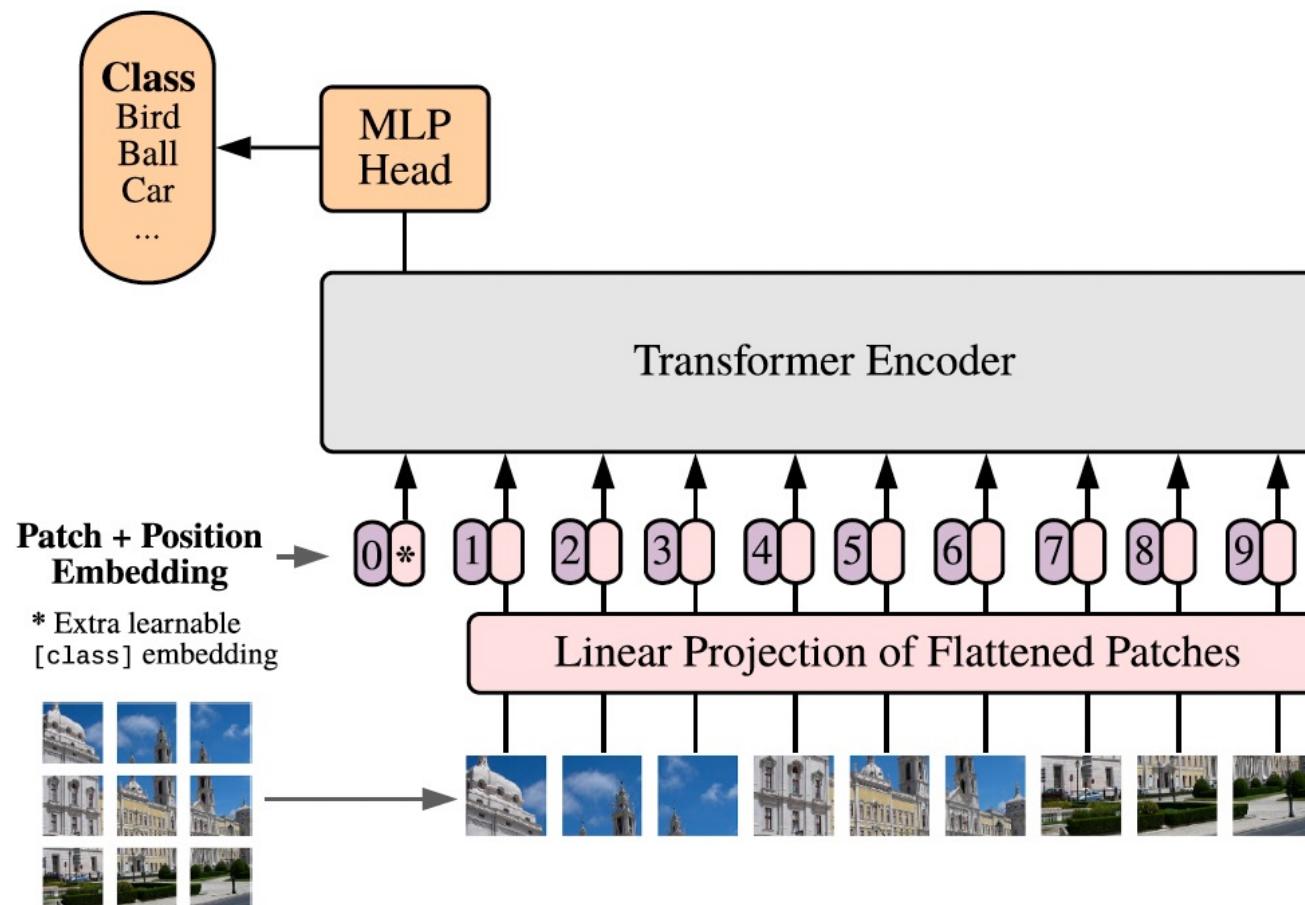


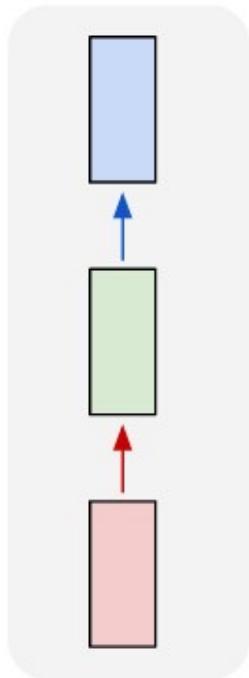
Image as a Sequence

- Sequence of patches (ViT)

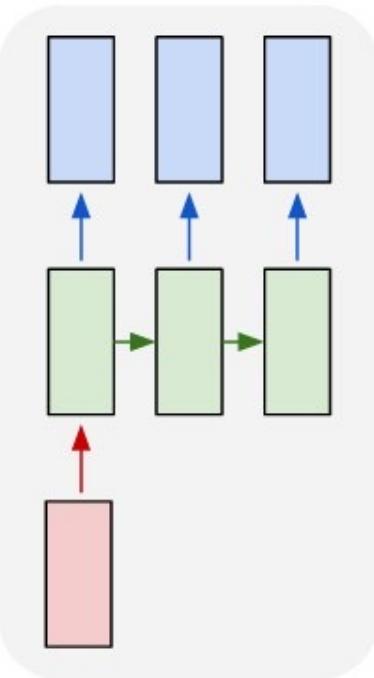


Seq2Seq Everywhere

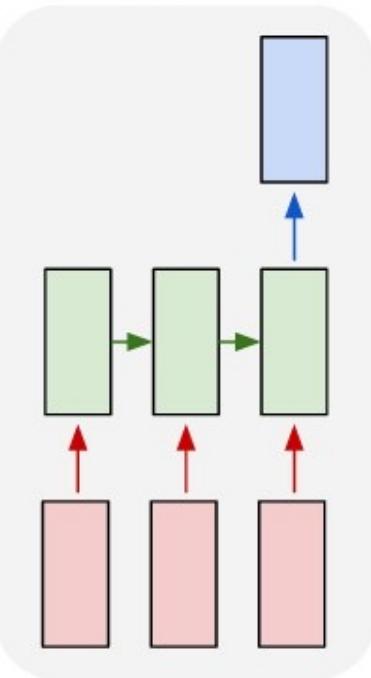
one to one



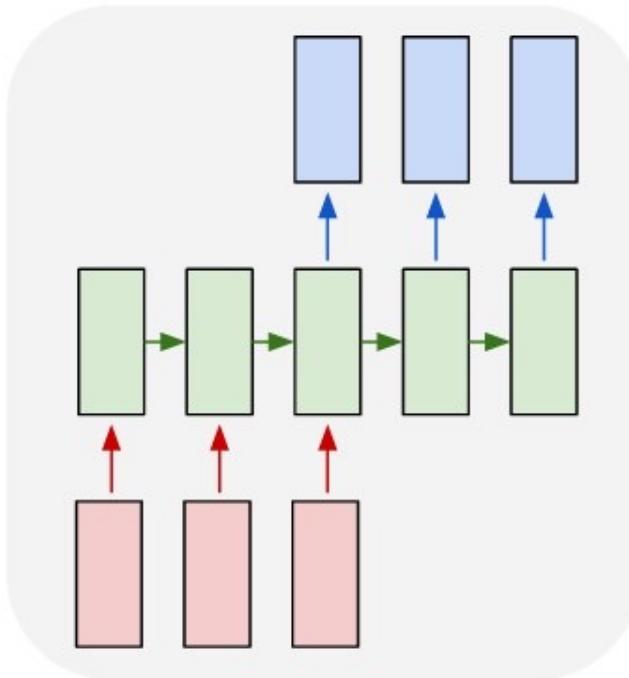
one to many



many to one



many to many



many to many

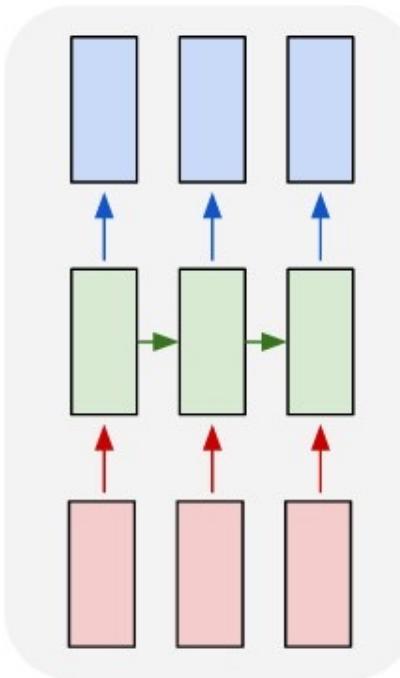


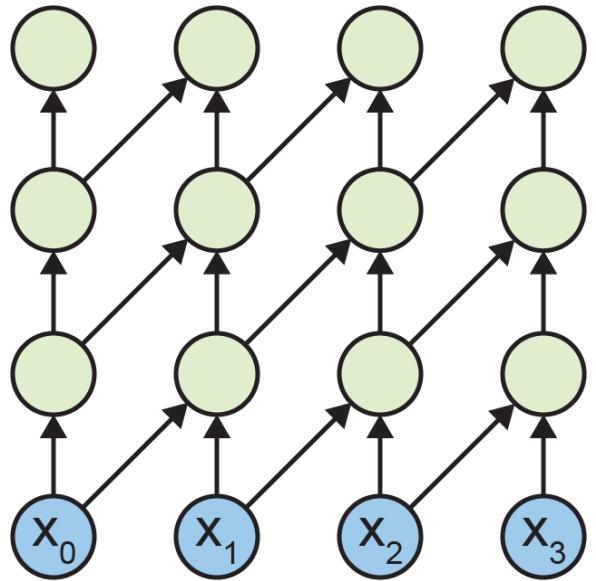
image classification

image captioning
object detection

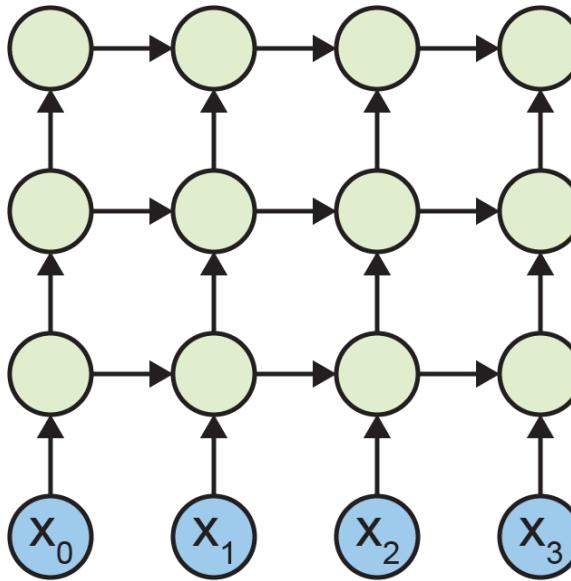
video classification

visual question answering, visual dialog
text-to-image, text-to-video generation, ...

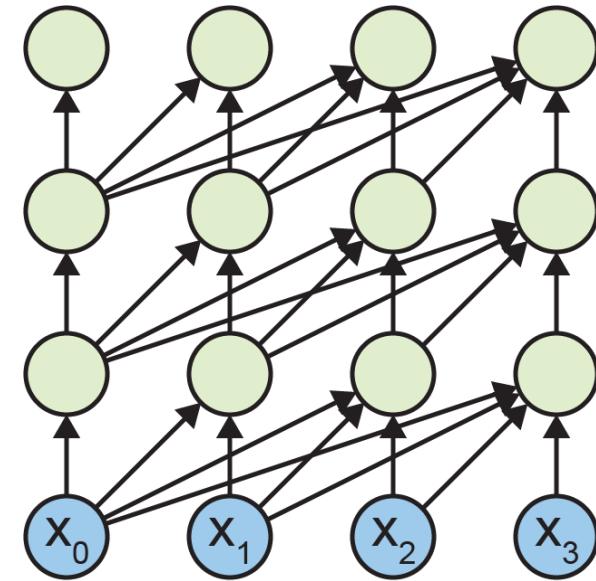
Paradigms of Sequence Modeling



CNN



RNN

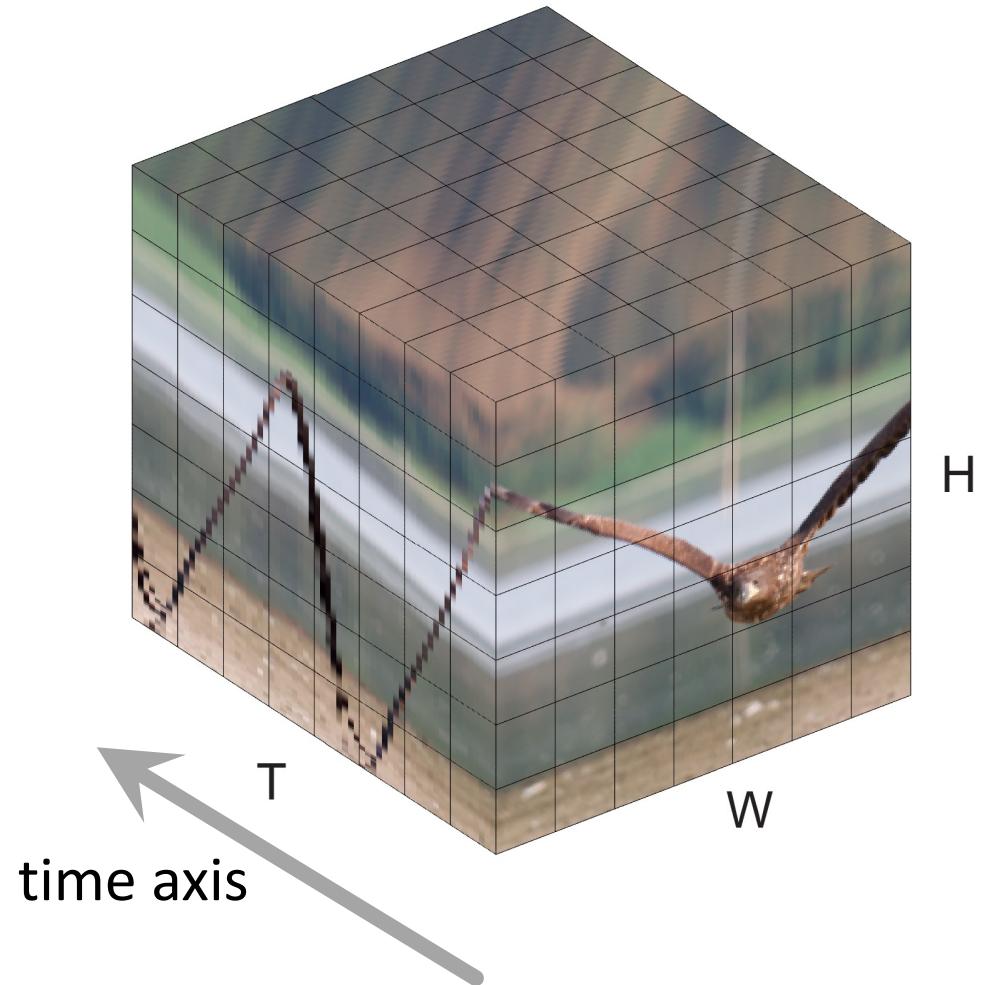


Attention

ConvNets for Sequence Problems

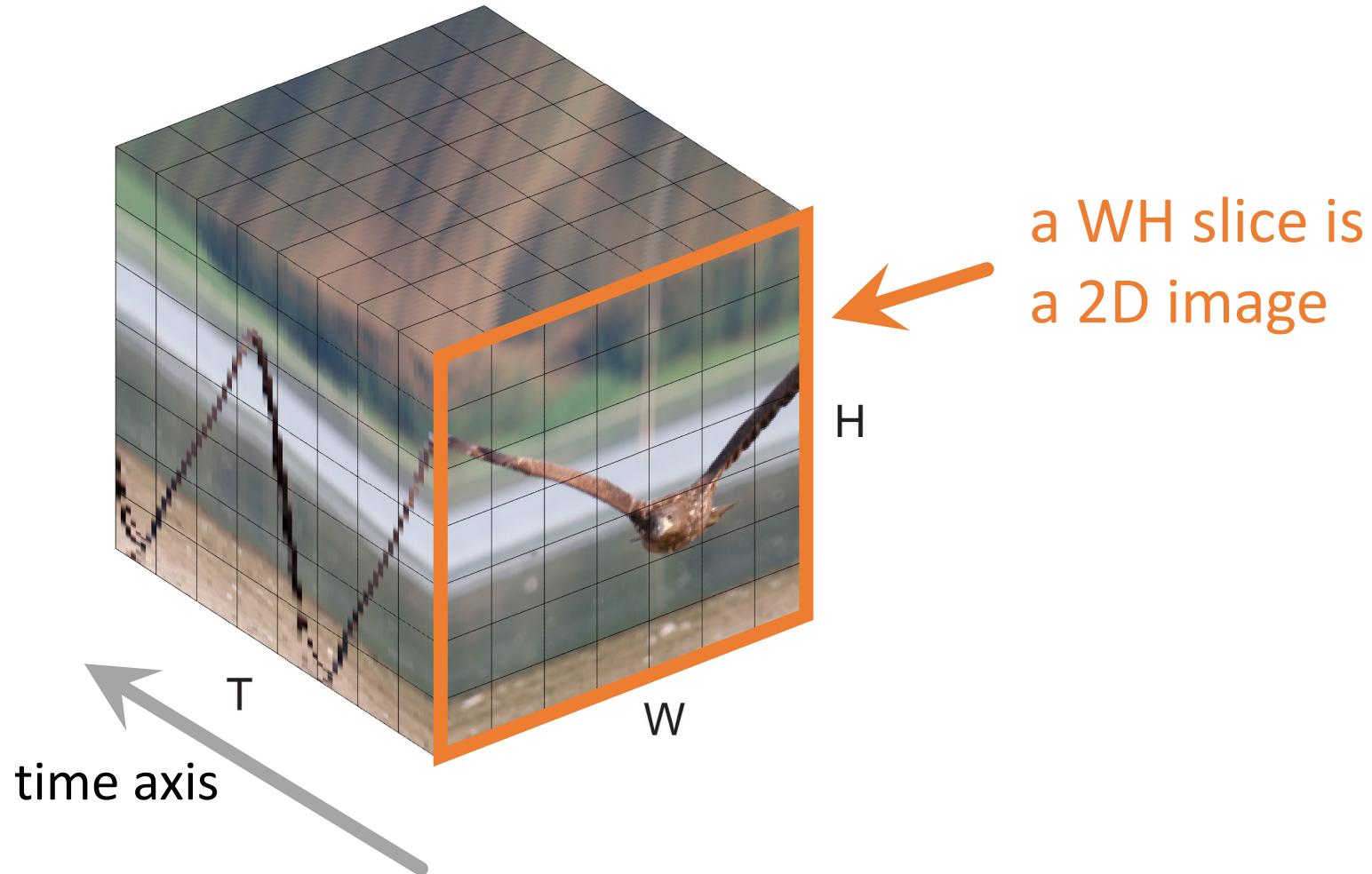
3D ConvNet for Videos

- Videos: 3D counterparts of images



3D ConvNet for Videos

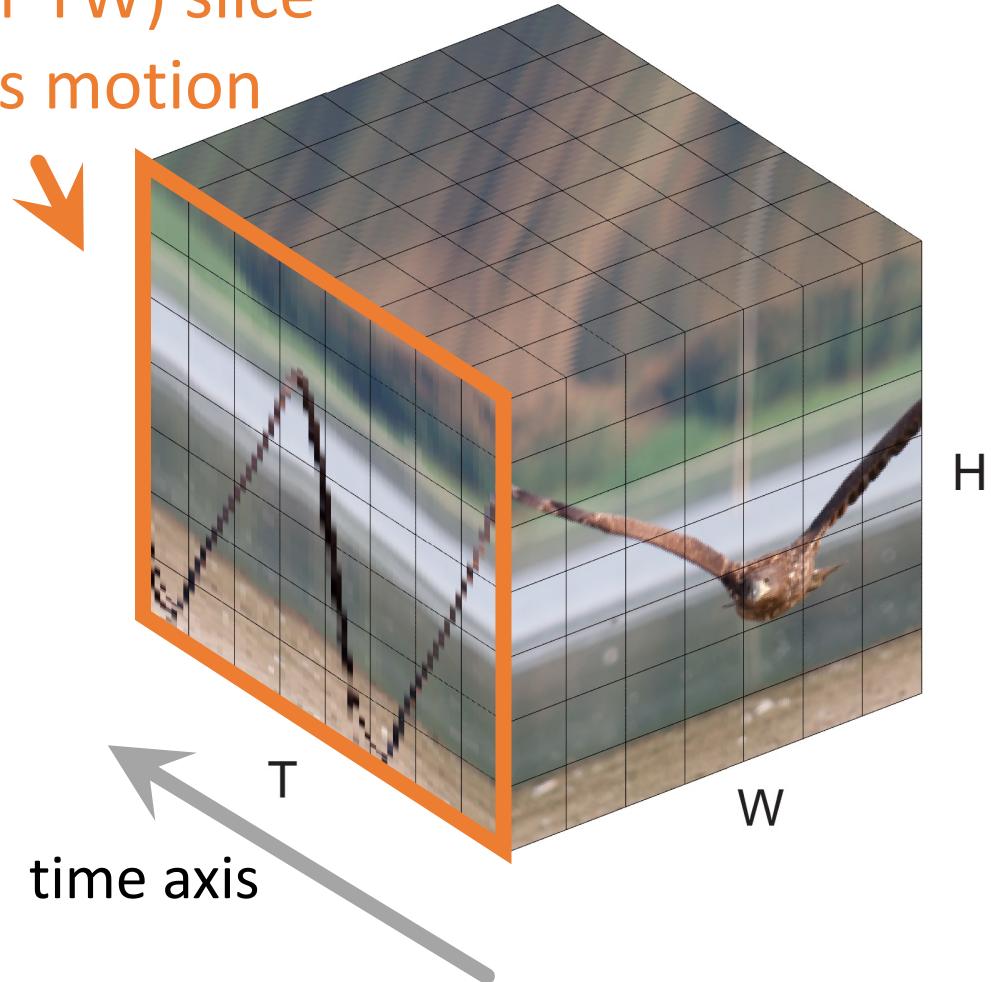
- Videos: 3D counterparts of images



3D ConvNet for Videos

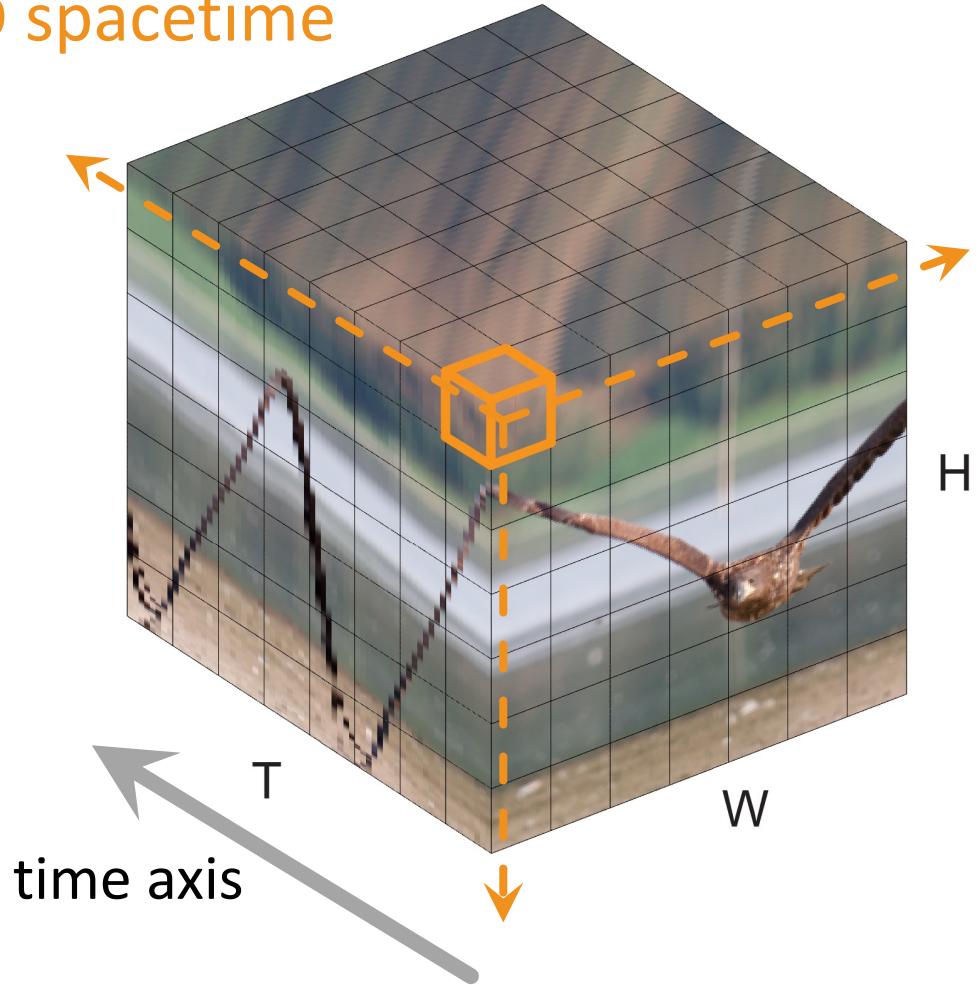
- Videos: 3D counterparts of images

a TH (or TW) slice
shows motion



3D ConvNet for Videos

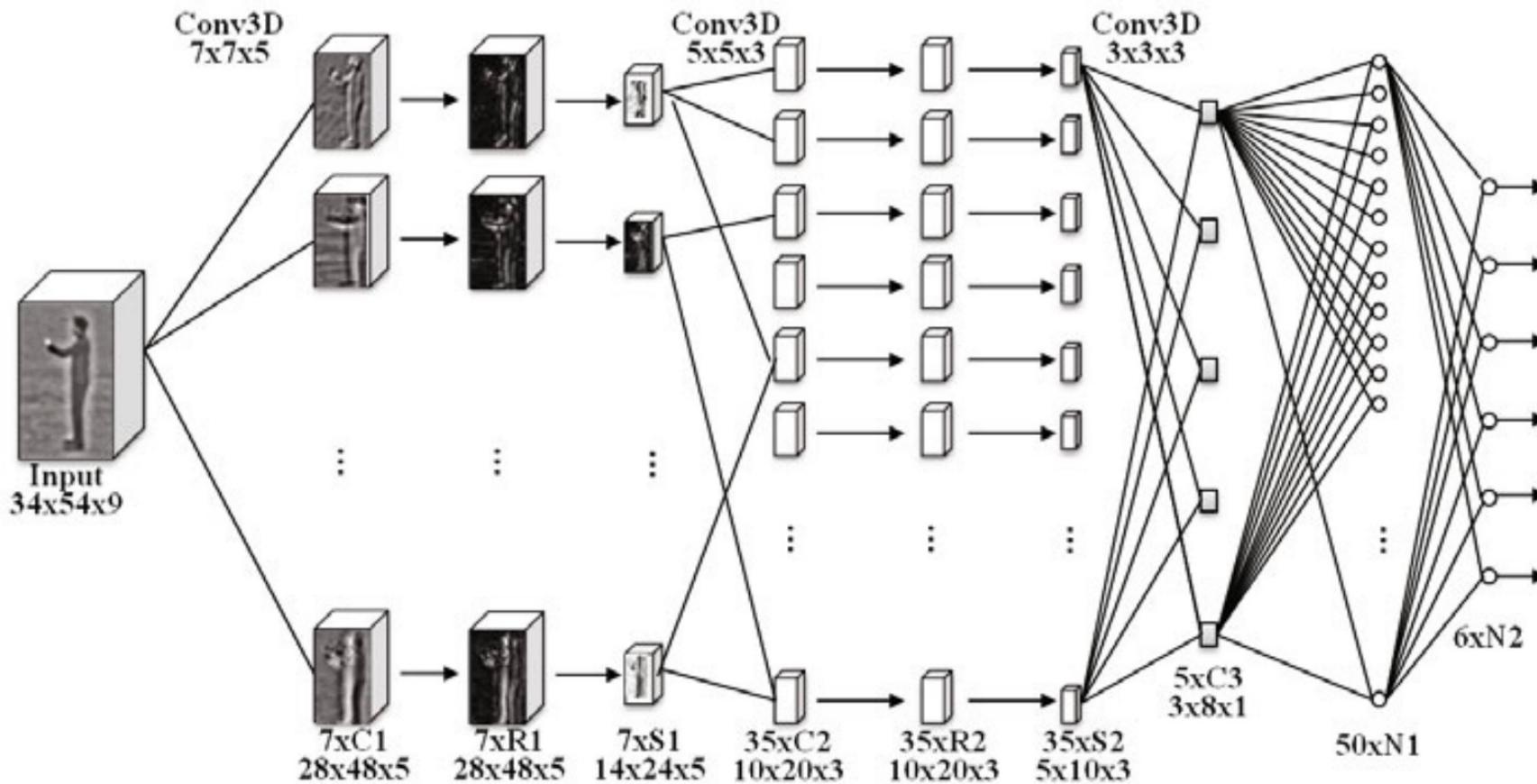
sliding a 3D window
in 3D spacetime



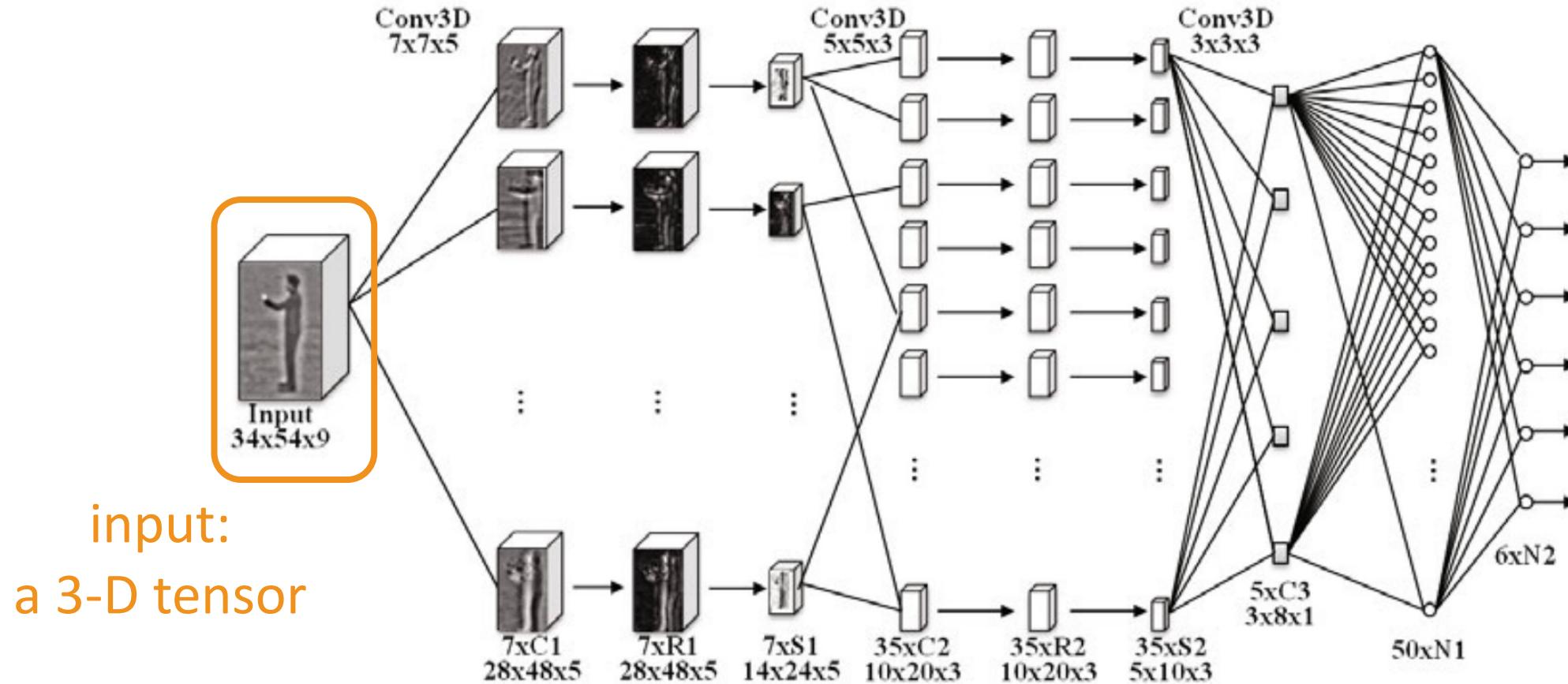
3D convolution

- features
 - **4-D tensor:** $C \times H \times W \times T$
 - C : channels
 - H : height
 - W : width
 - T : time
- filters
 - **5-D tensor:** $C_o \times C_i \times K_h \times K_w \times K_T$
 - C_o : output channels
 - C_i : input channels
 - K_h, K_w, K_T : filter size

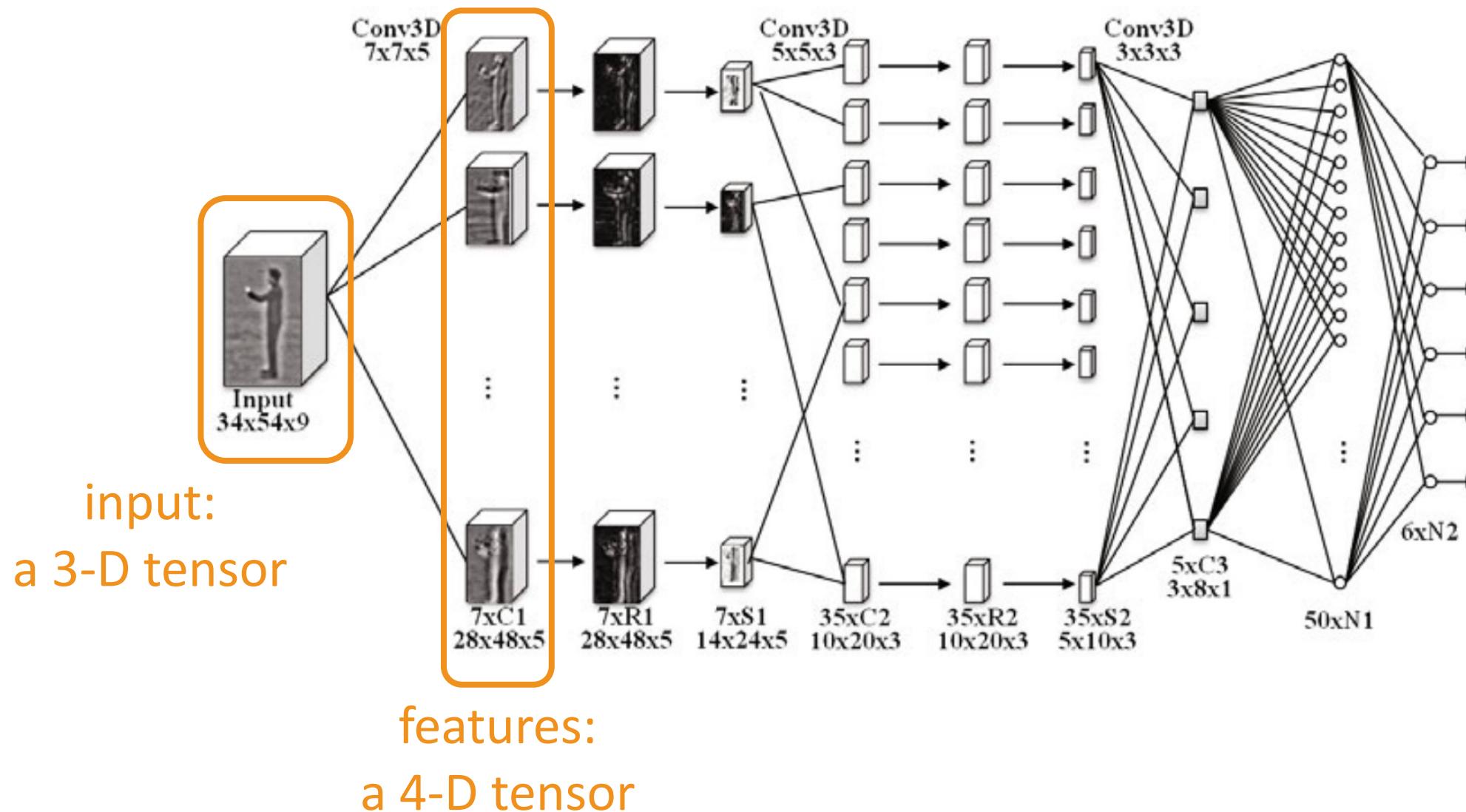
3D ConvNet for Videos



3D ConvNet for Videos



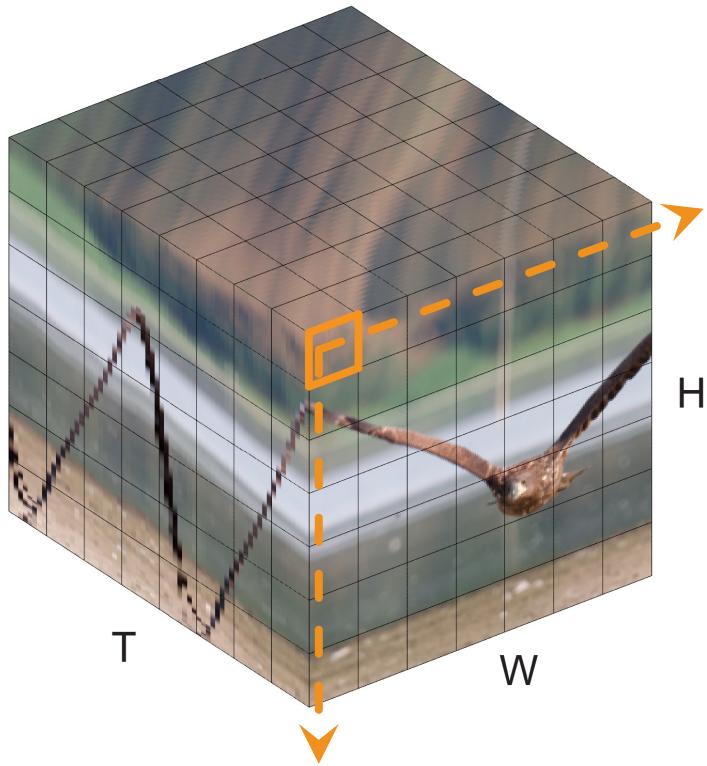
3D ConvNet for Videos



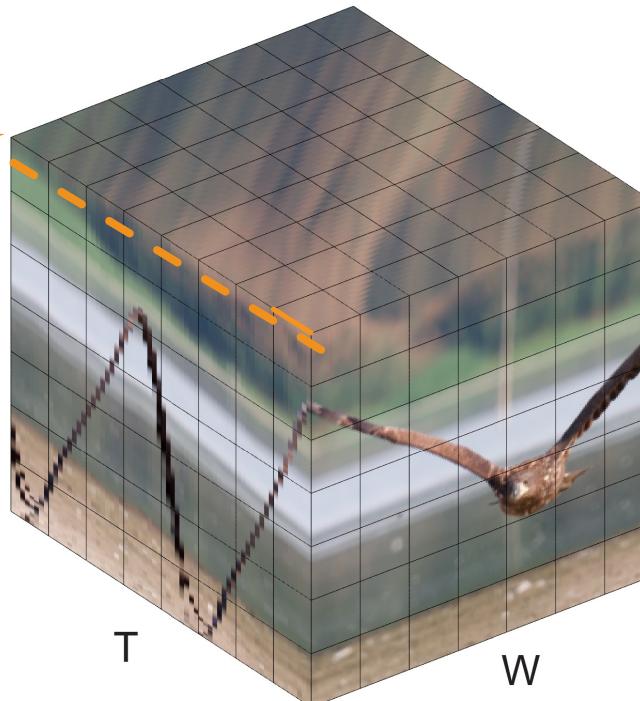
Separating Space and Time

- Separable convolution: 2D+1D

2D conv in space,
for each frame



1D conv in time,
for each pixel



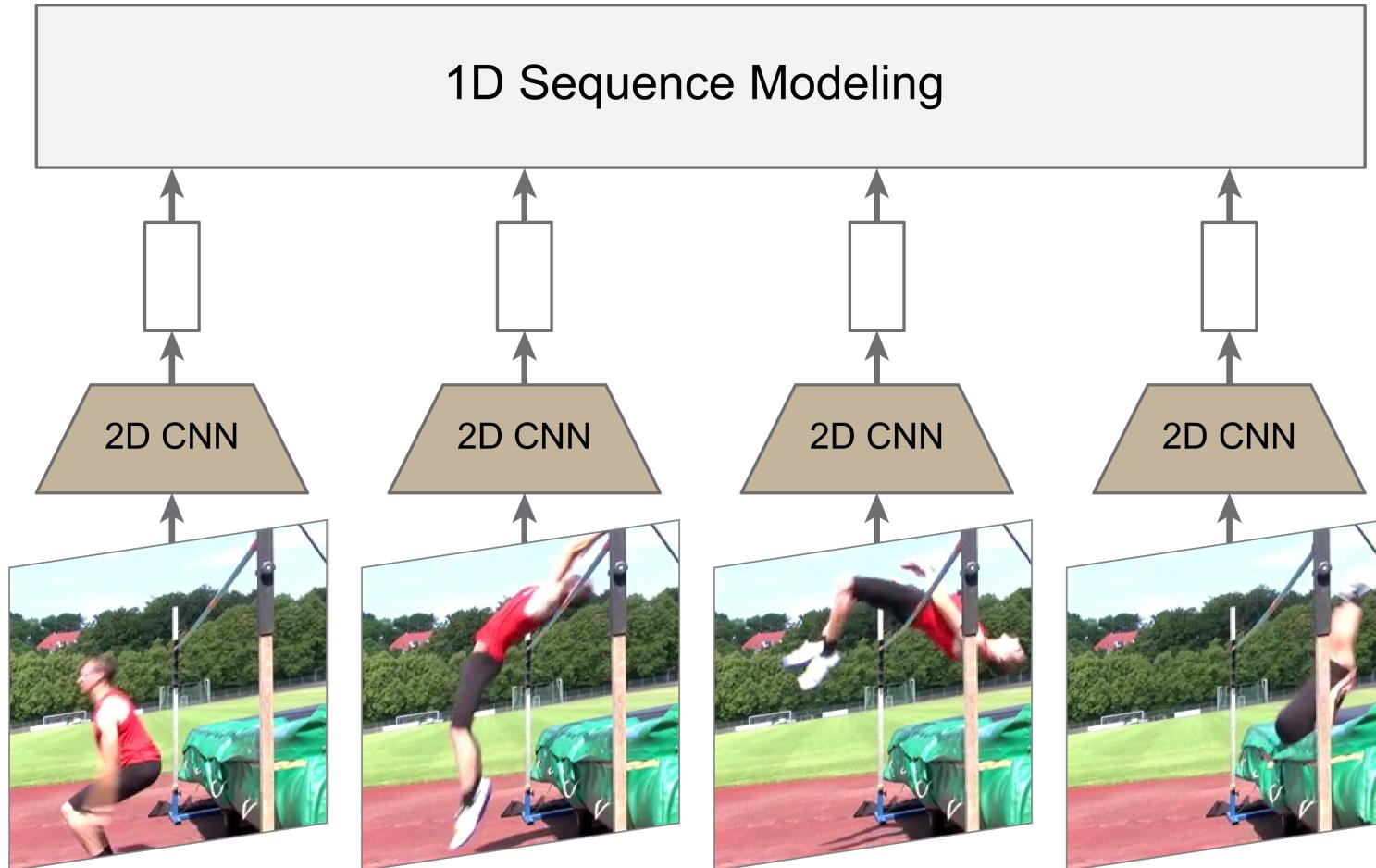
parameters

- **3D conv:** $C_o \times C_i \times K_h \times K_w \times K_T$
- **separable conv:** $C_{o'} \times C_i \times K_h \times K_w \times 1 + C_o \times C_{o'} \times 1 \times 1 \times K_T$

operations = # params \times H \times W \times T

Separating Space and Time

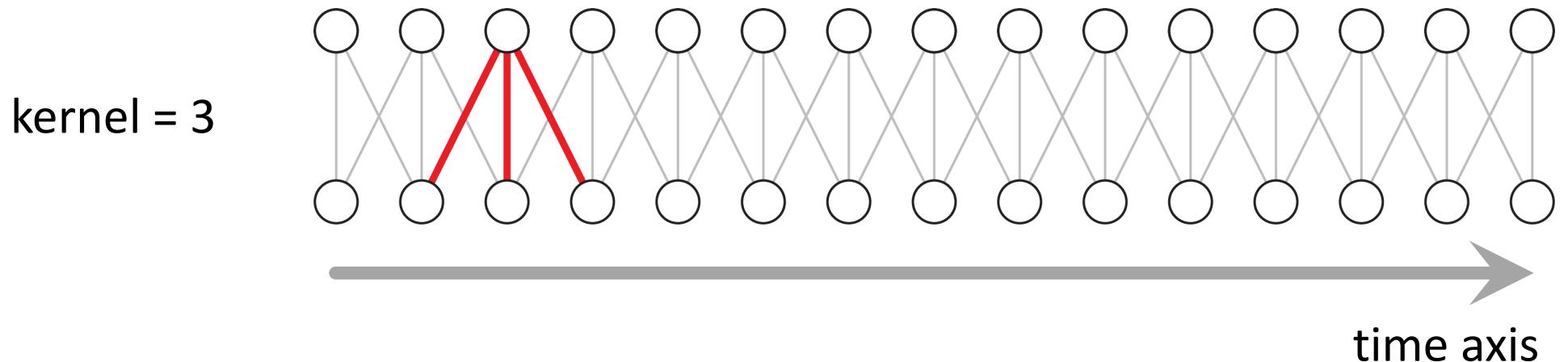
- 2D ConvNet per frame + 1D Sequence Modeling



1D ConvNet for Sequence Modeling

Convolution (in time)

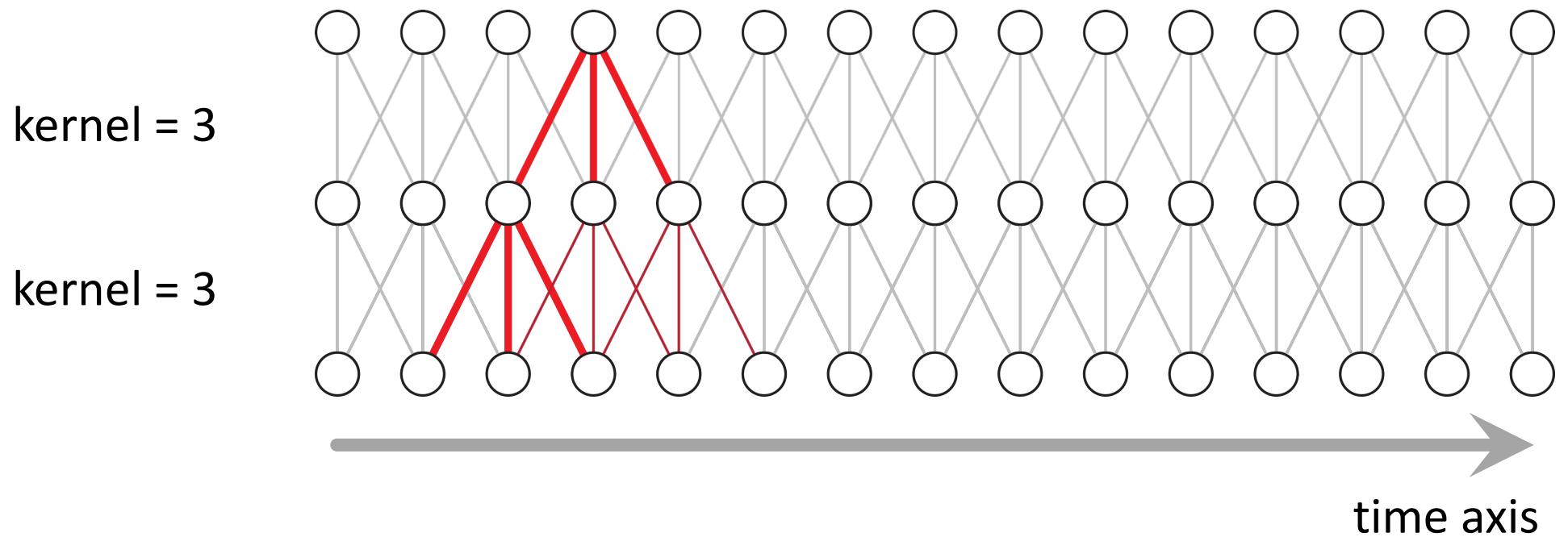
- sliding window (in time)
- local connection (in time)
- weight-sharing (in time)



1D ConvNet for Sequence Modeling

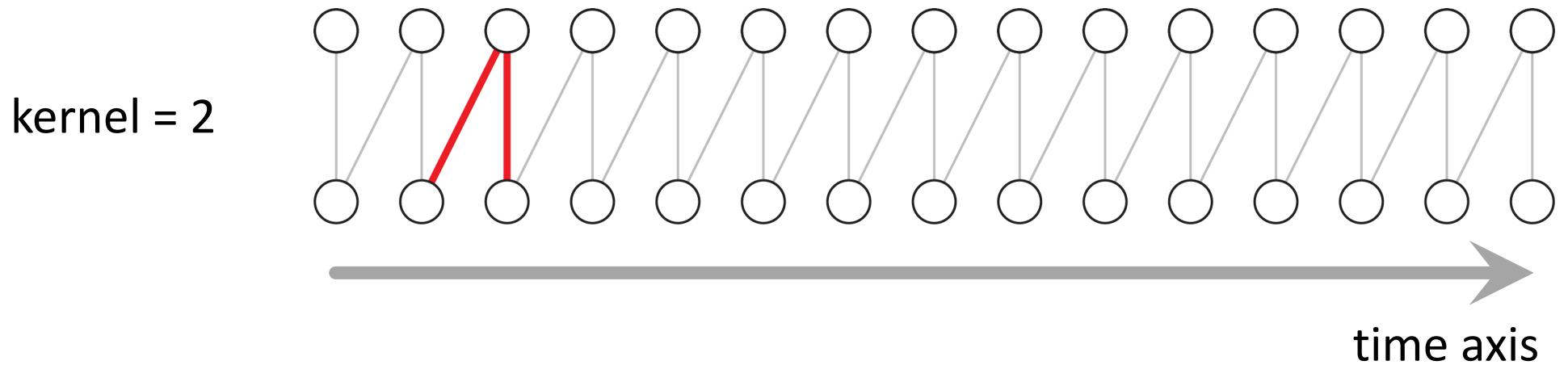
Convolution (in time)

- sliding window (in time)
- local connection (in time)
- weight-sharing (in time)



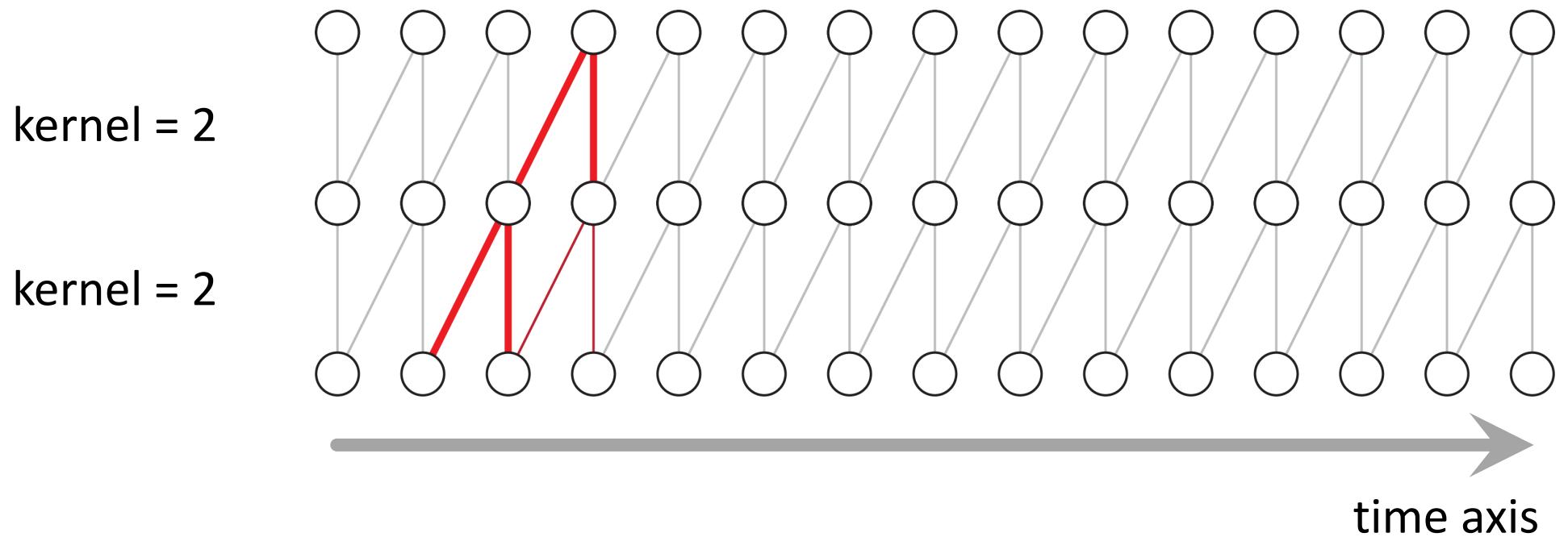
1D ConvNet: Causal convolution

- **Causal:** output at a time step does not depend on future time steps



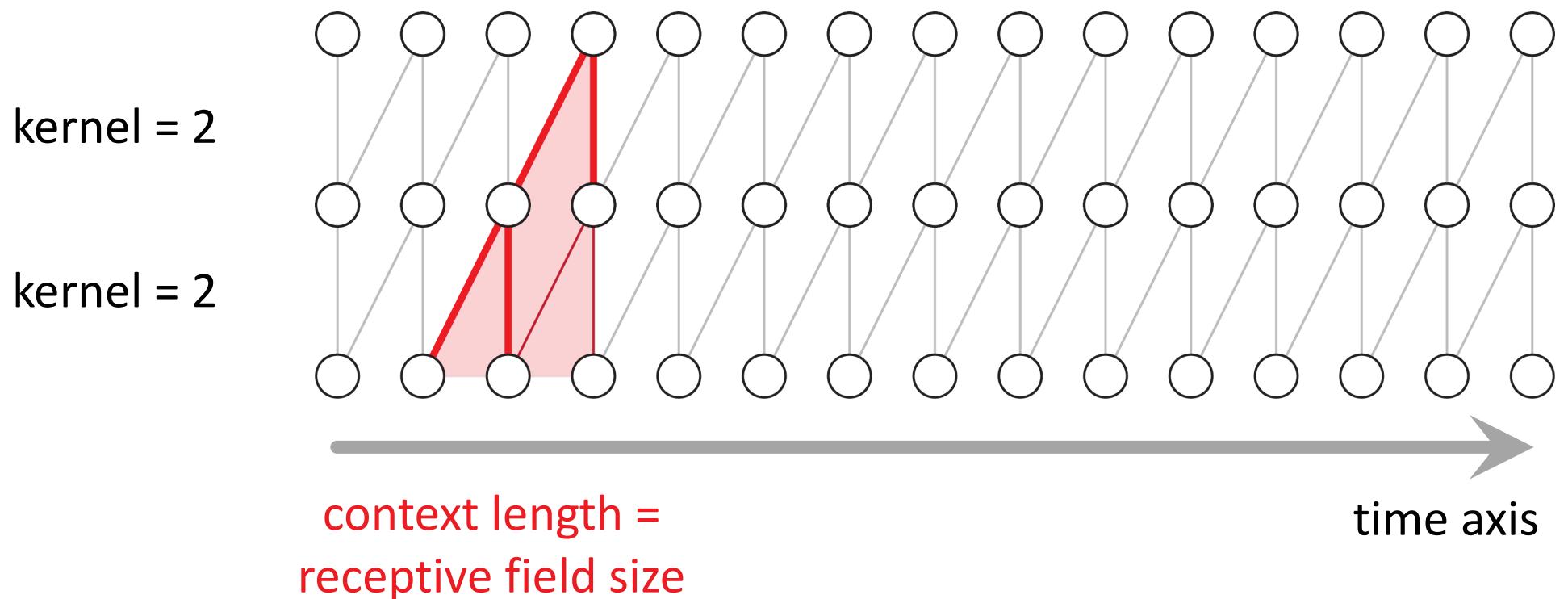
1D ConvNet: Causal convolution

- **Causal:** output at a time step does not depend on future time steps



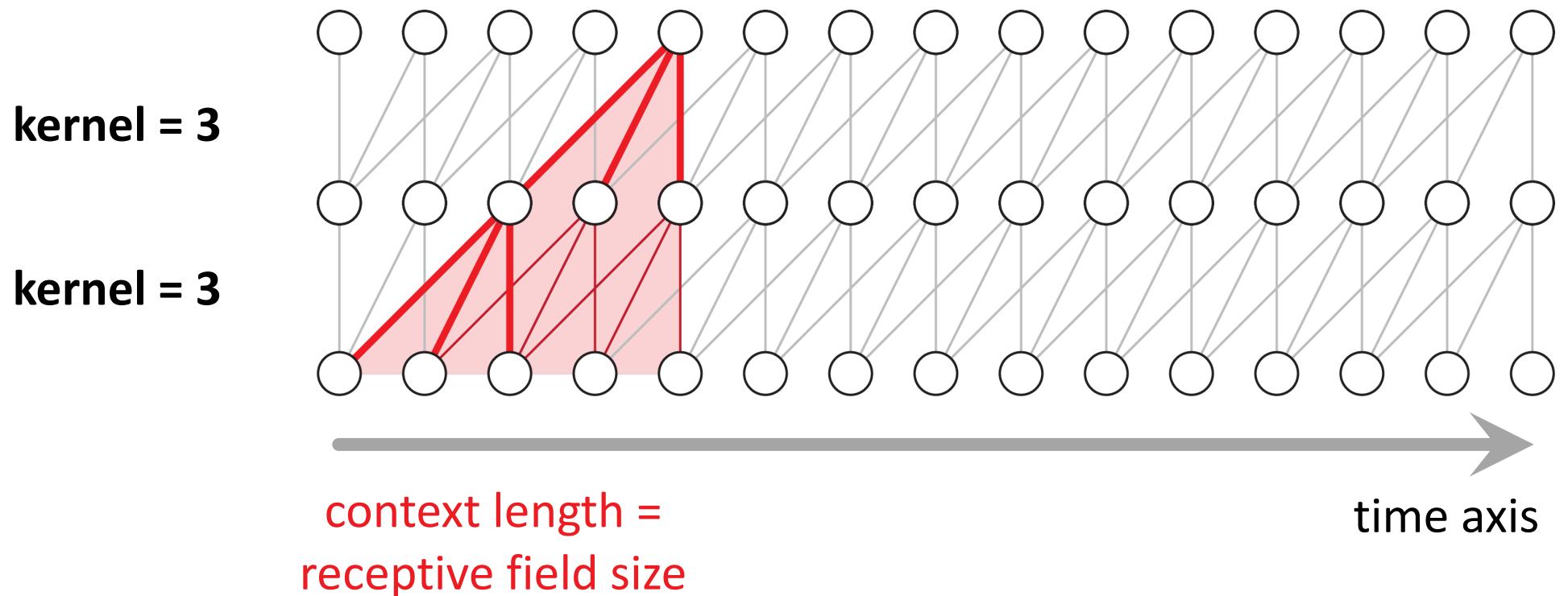
1D ConvNet: Causal convolution

- **Causal:** output at a time step does not depend on future time steps



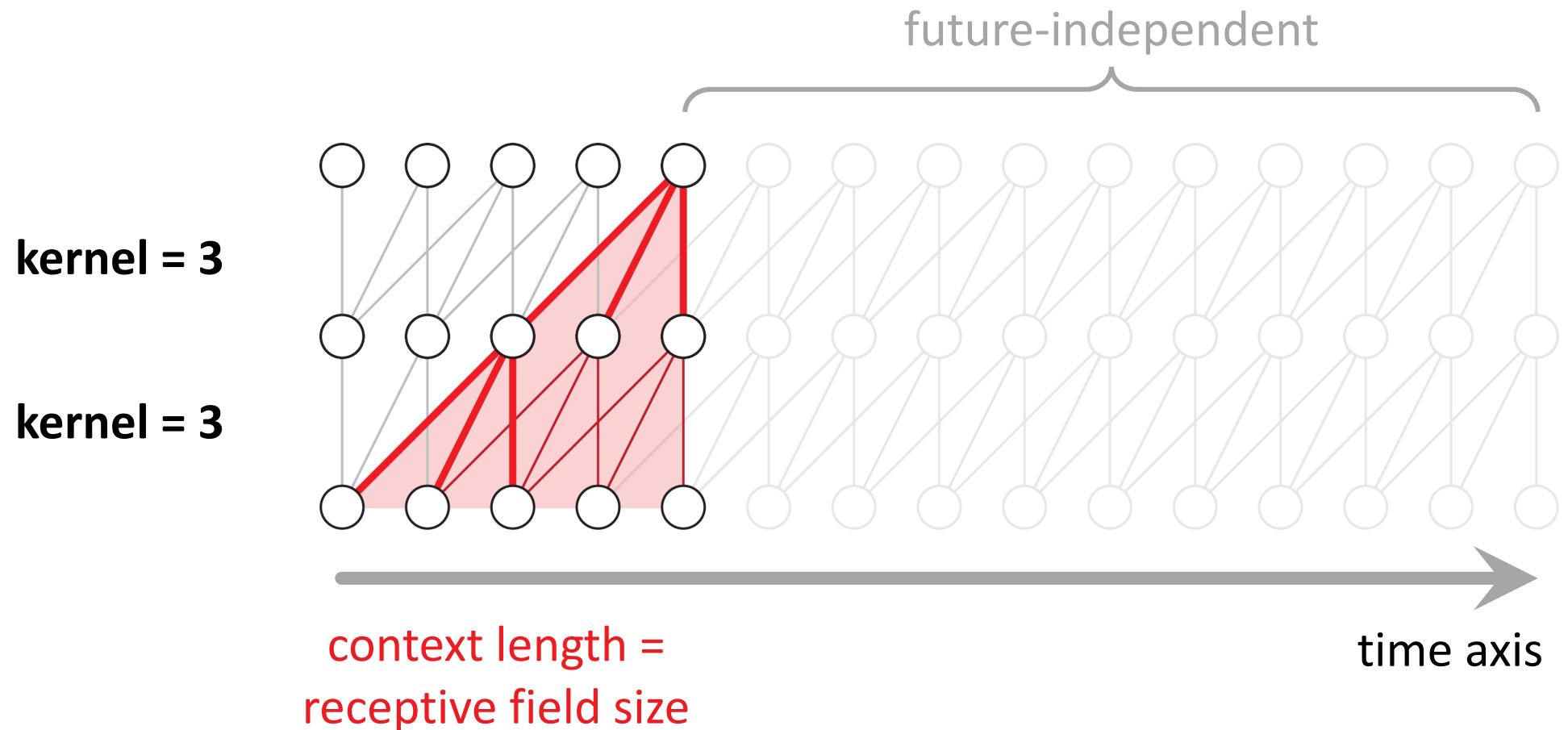
1D ConvNet: Causal convolution

- **Causal:** output at a time step does not depend on future time steps



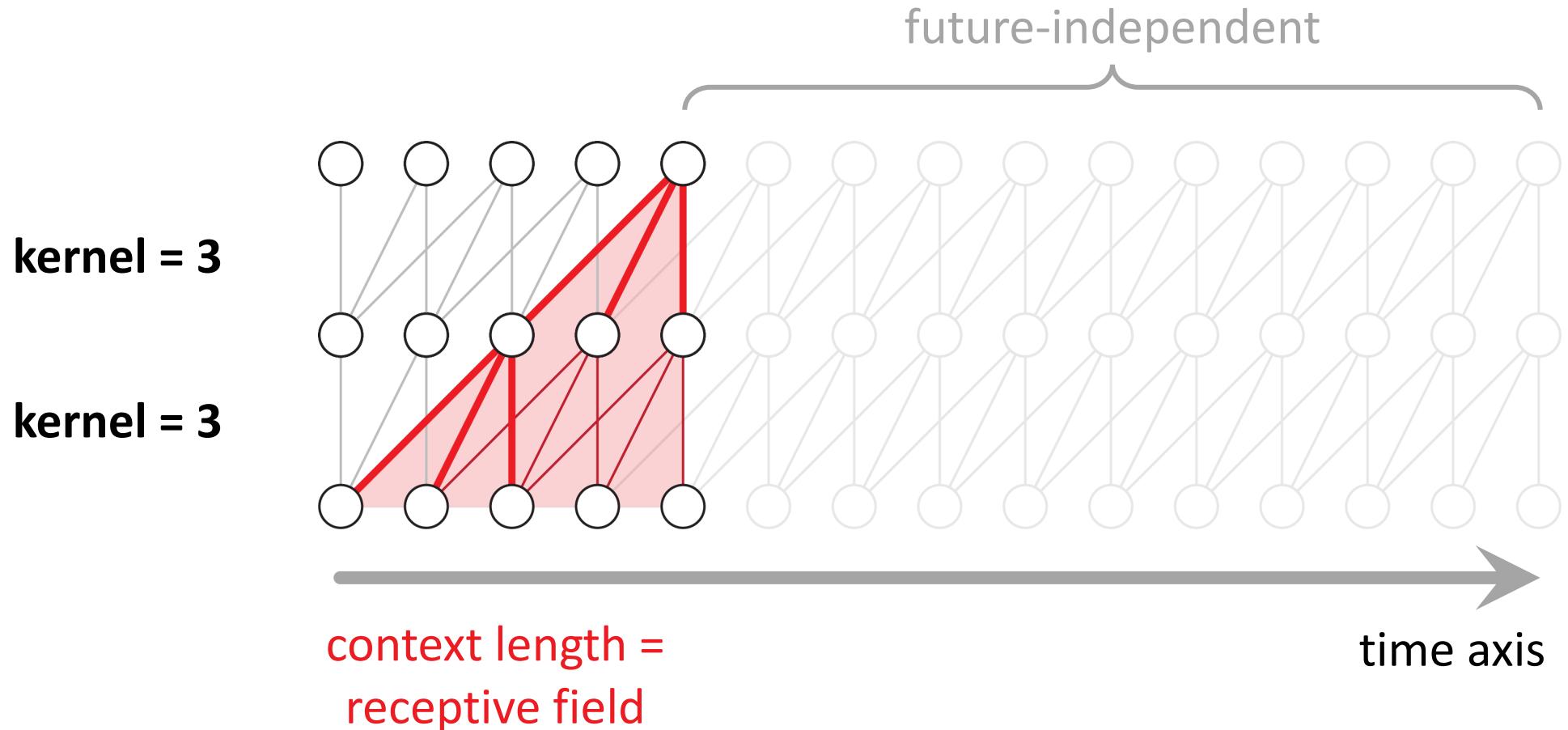
1D ConvNet: Causal convolution

- **Causal:** output at a time step does not depend on future time steps



1D ConvNet: Causal convolution

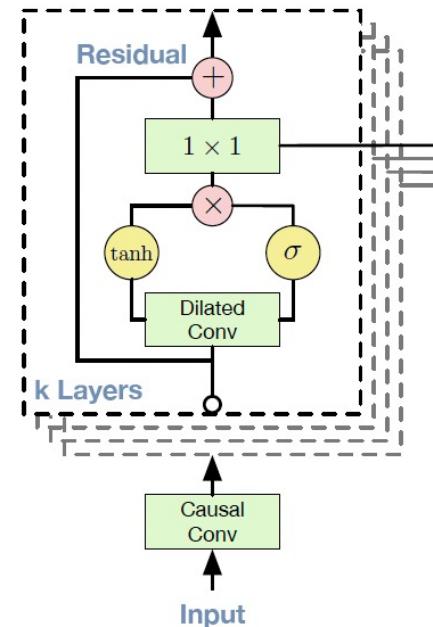
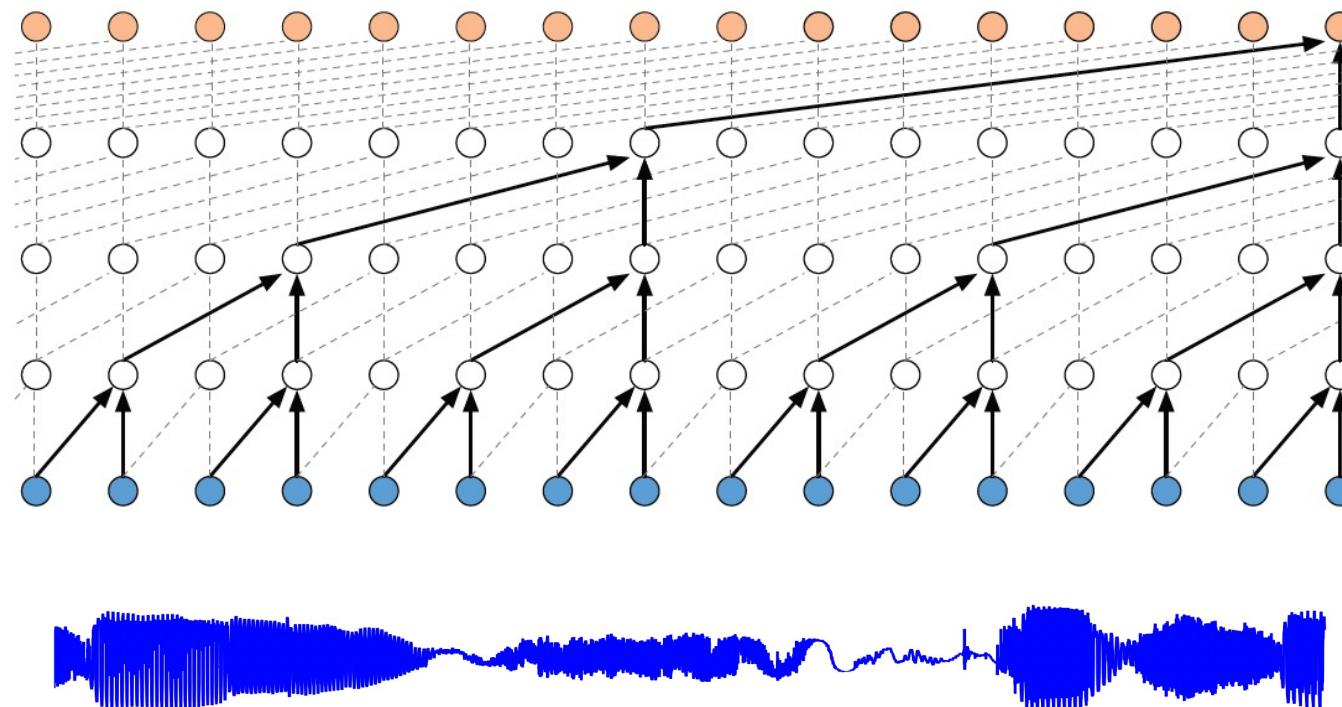
- **Causal:** output at a time step does not depend on future time steps
- conv w/ **pad = kernel – 1**, only on left (no pad on right)



Example: DeepMind's WaveNet for Audio Generation

Causal ConvNet with long context

- going deep with residual connections
- “dilated” convolution



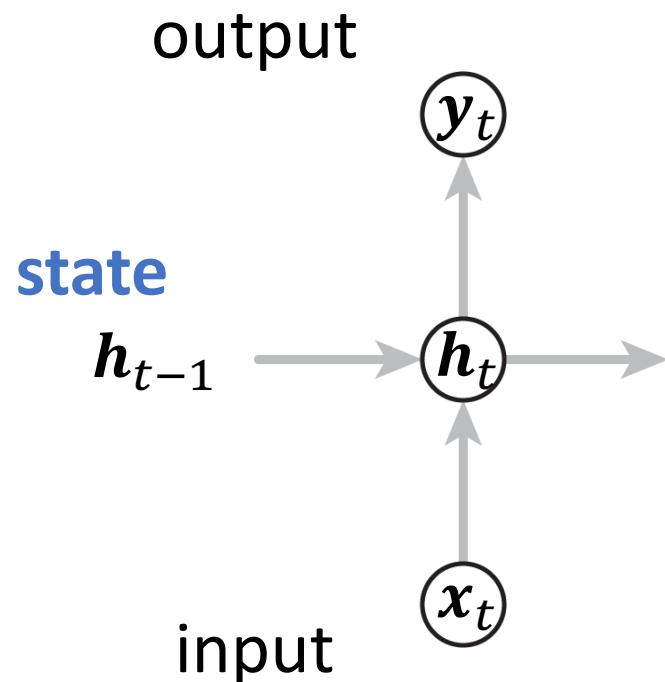
Recurrent Neural Networks (RNNs)

Recurrent Neural Networks (RNNs)

- Hidden state: summary of the past
- Weight-sharing in time
- Ordered sequences
- BackProp Through Time (BPTT)

Recurrent NN: Hidden State

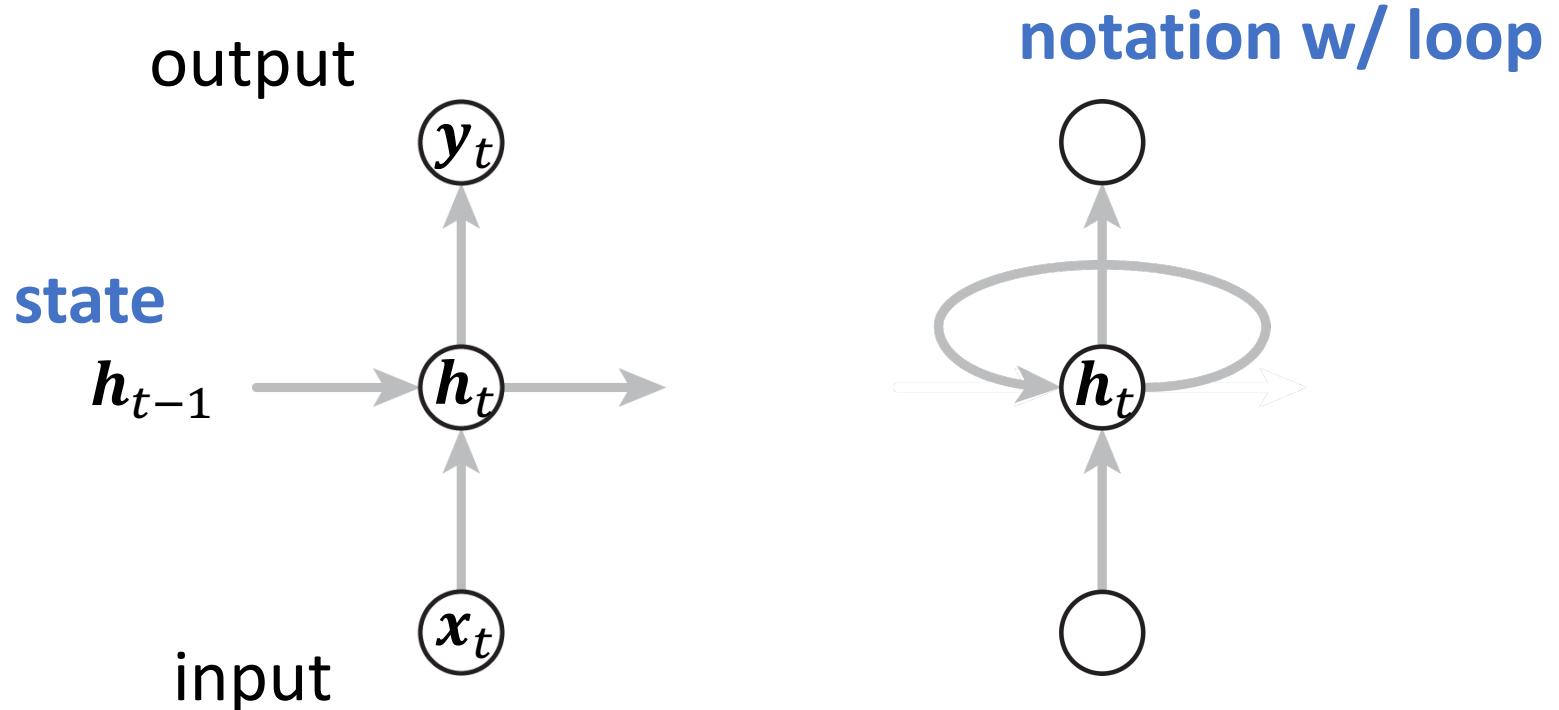
- Hidden state: summary/memory of the sequence seen so far



$$\begin{aligned} h_t &= f(h_{t-1}, x_t) && \bullet \text{ recurrent} \\ y_t &= g(h_t) && \bullet \text{ feedforward} \end{aligned}$$

Recurrent NN: Hidden State

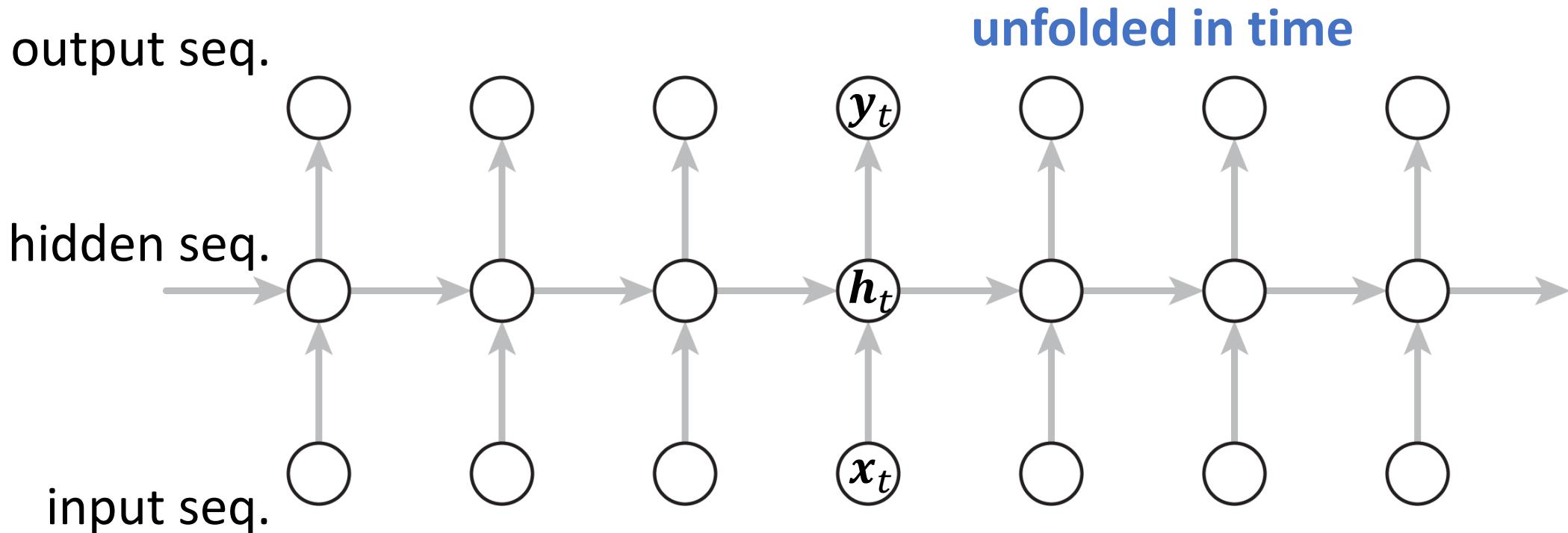
- Hidden state: summary/memory of the sequence seen so far



$$\begin{aligned} \mathbf{h}_t &= f(\mathbf{h}_{t-1}, \mathbf{x}_t) & \bullet \text{ recurrent} \\ \mathbf{y}_t &= g(\mathbf{h}_t) & \bullet \text{ feedforward} \end{aligned}$$

Recurrent NN: Hidden State

- Hidden state: summary/memory of the sequence seen so far

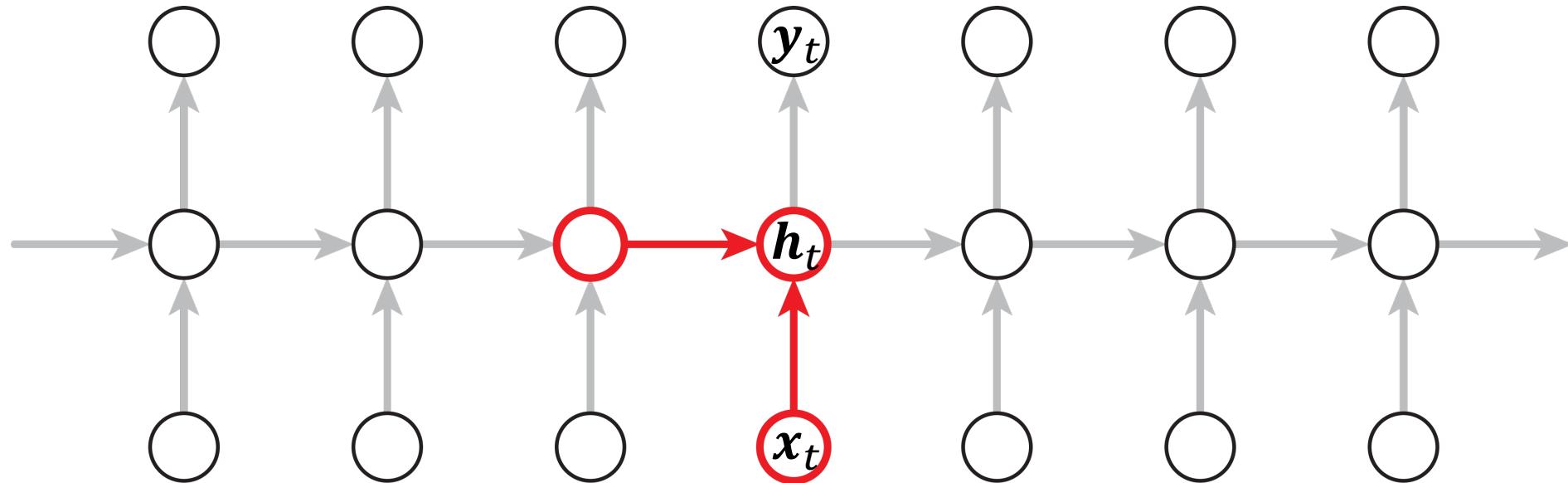


$$h_t = f(h_{t-1}, x_t)$$
$$y_t = g(h_t)$$

- recurrent
- feedforward

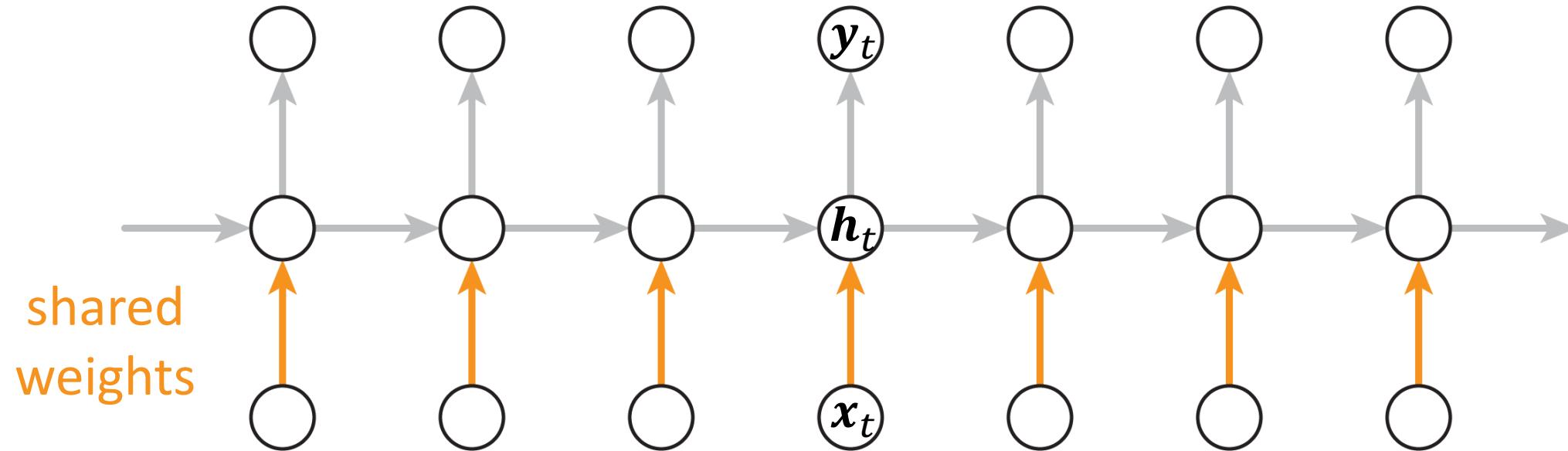
RNN: local connection + weight-sharing

- local connection in time



RNN: local connection + weight-sharing

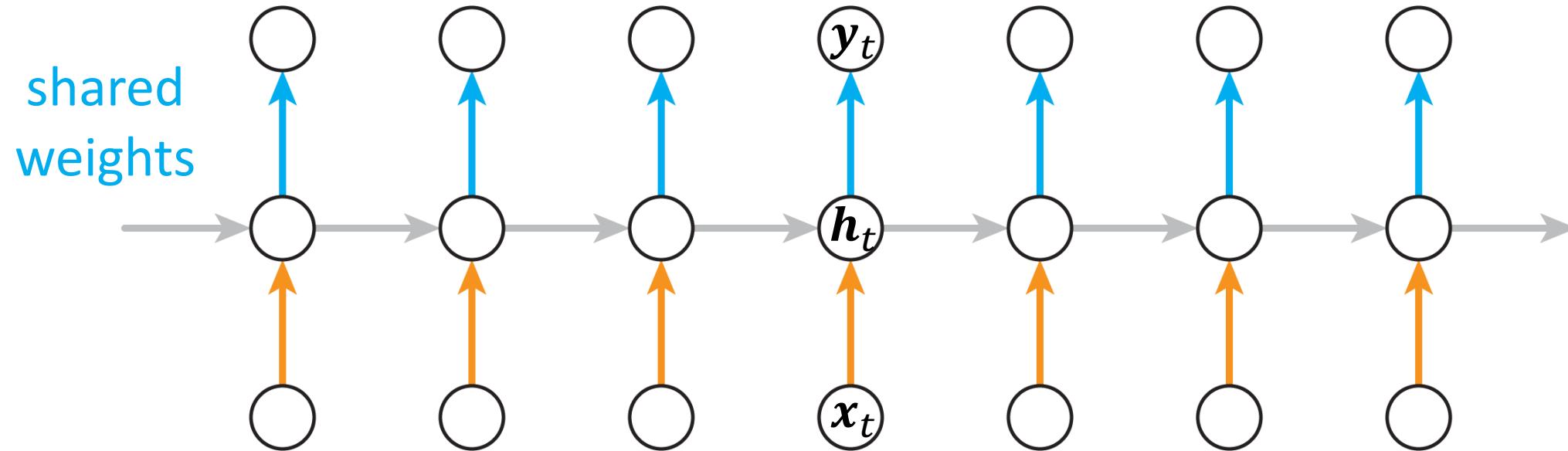
- weight-sharing in time



similar to conv
with kernel size 1

RNN: local connection + weight-sharing

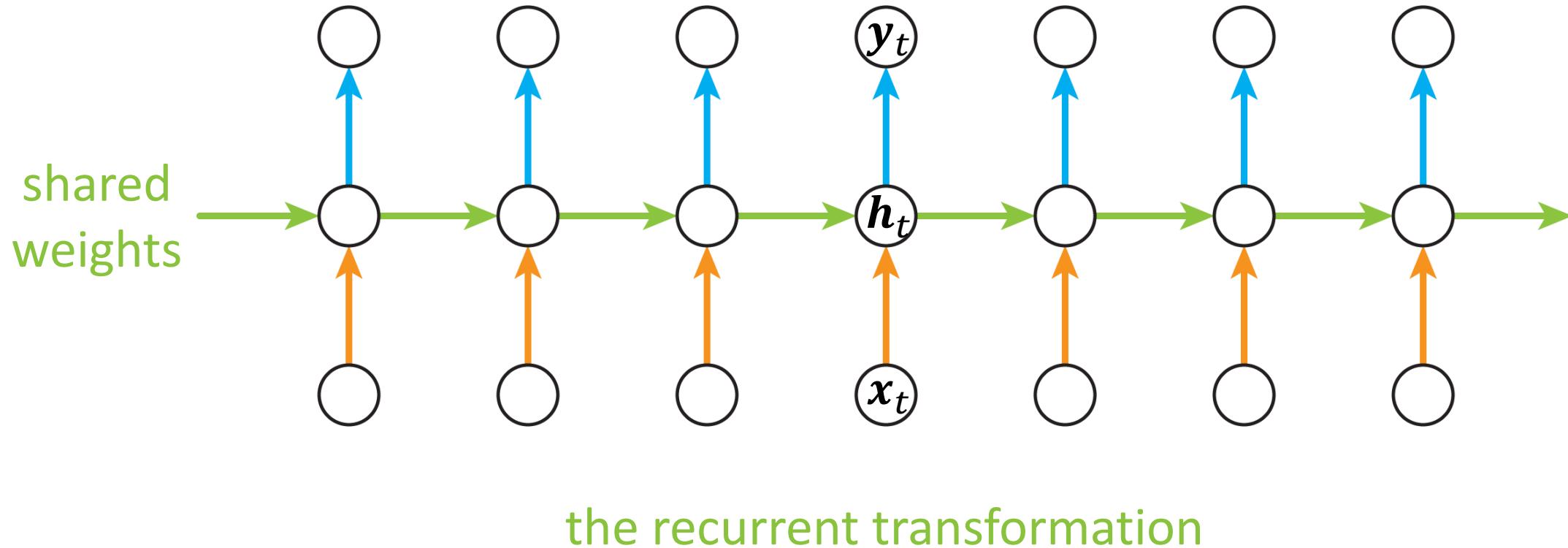
- weight-sharing in time



similar to conv
with kernel size 1

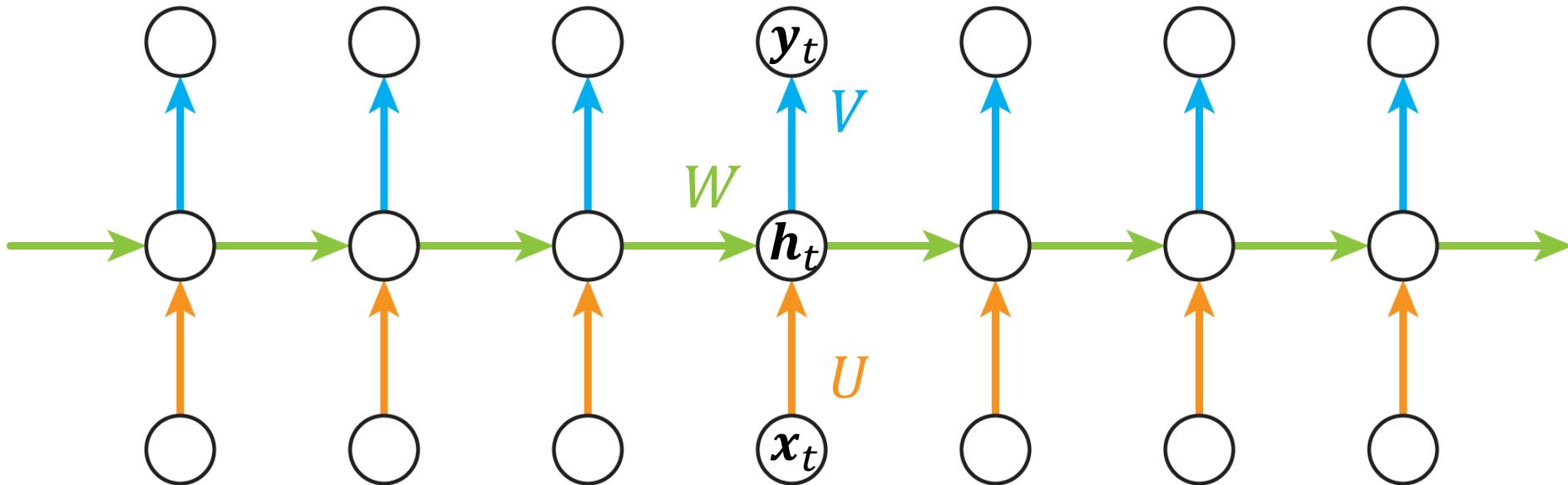
RNN: local connection + weight-sharing

- weight-sharing in time



Example: Vanilla RNN

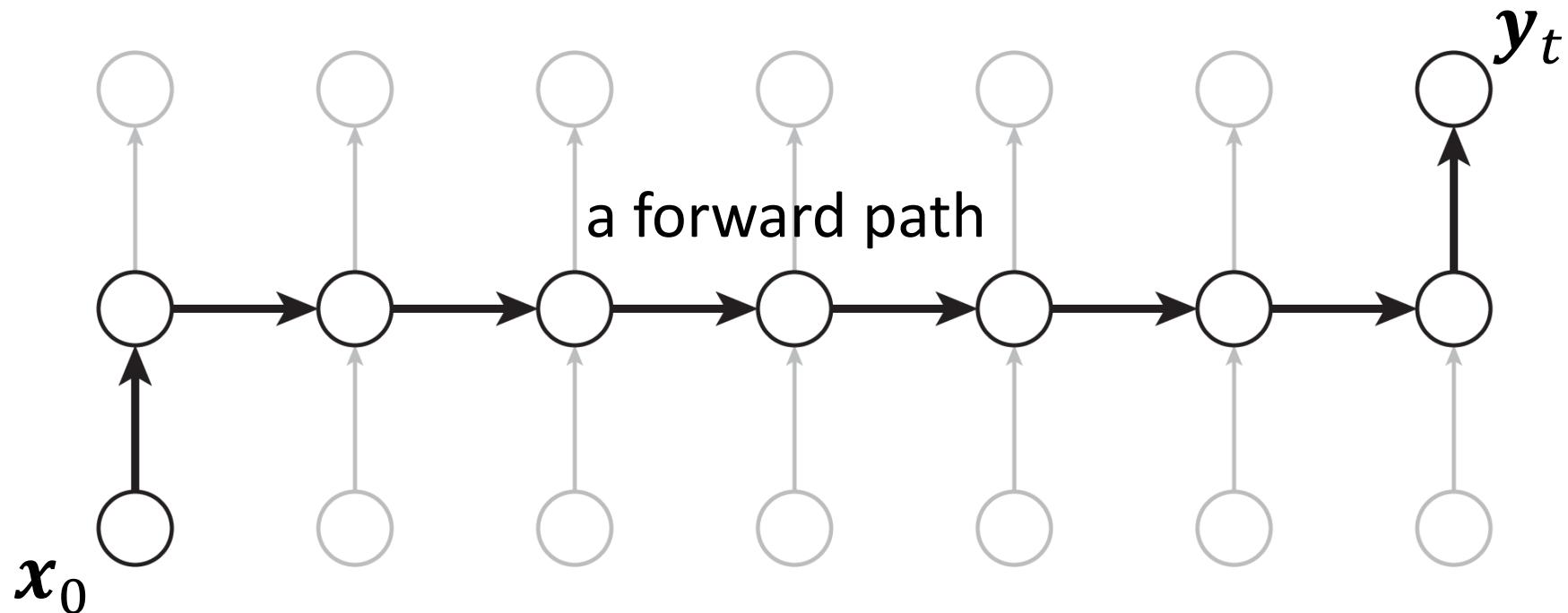
$$\begin{aligned}\mathbf{h}_t &= \tanh(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t) \\ \mathbf{y}_t &= \mathbf{V}\mathbf{h}_t\end{aligned}$$



*bias omitted for simplicity

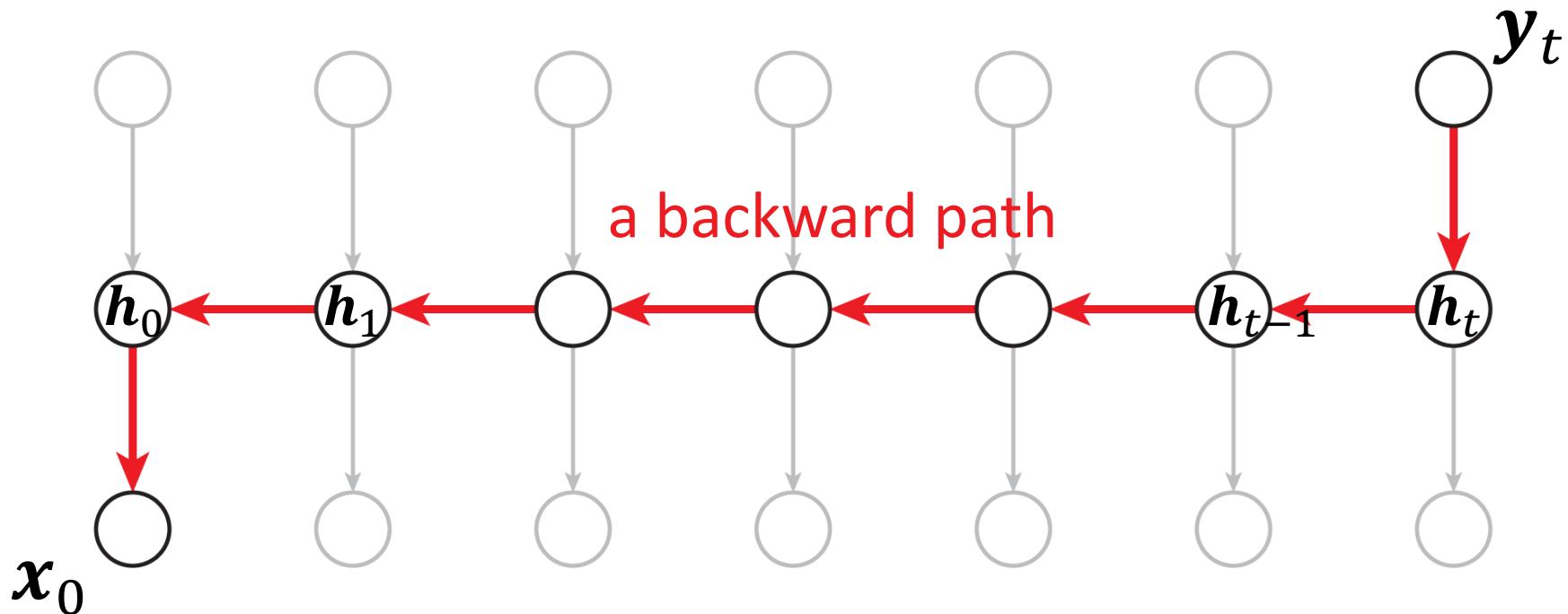
RNN: BackProp Through Time (BPTT)

- output at time t is like a t -layer deep net



RNN: BackProp Through Time (BPTT)

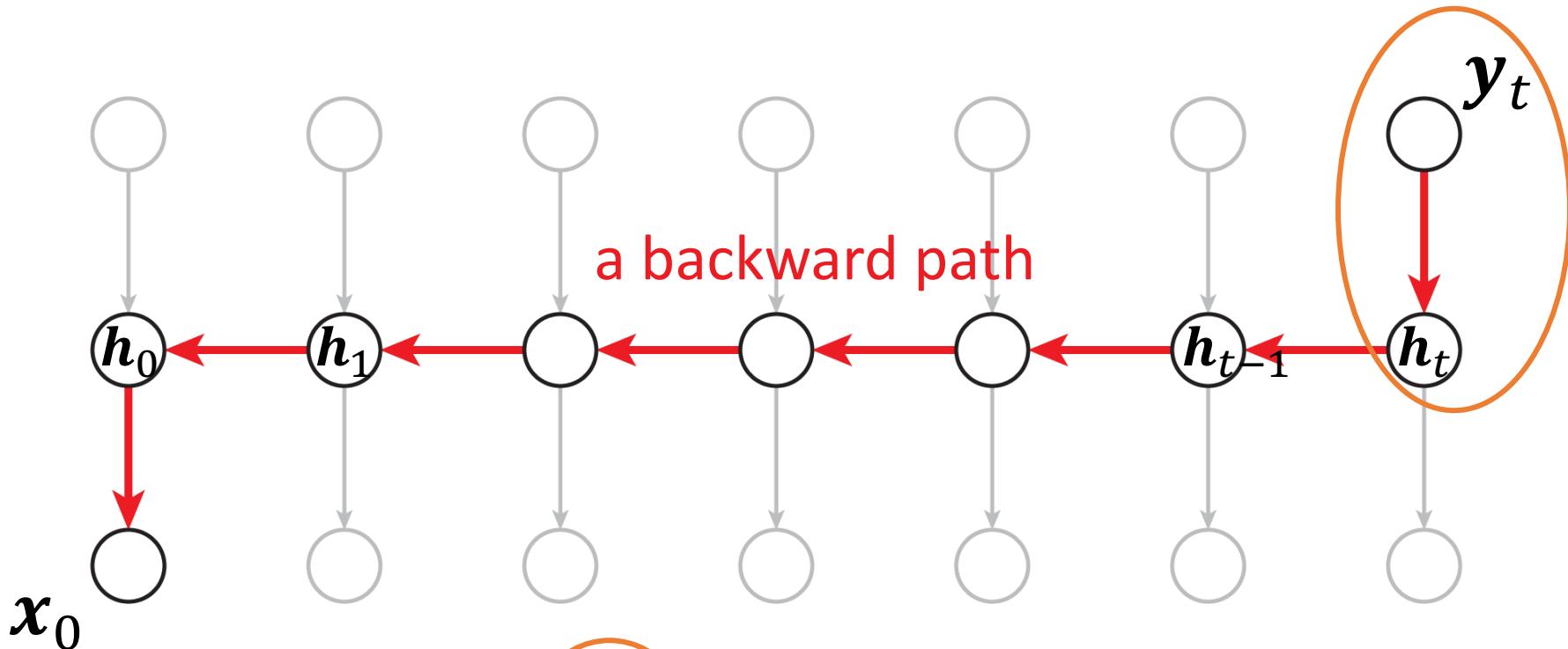
- output at time t is like a t -layer deep net



$$\frac{\partial y_t}{\partial x_0} = \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial h_0} \frac{\partial h_0}{\partial x_0}$$

RNN: BackProp Through Time (BPTT)

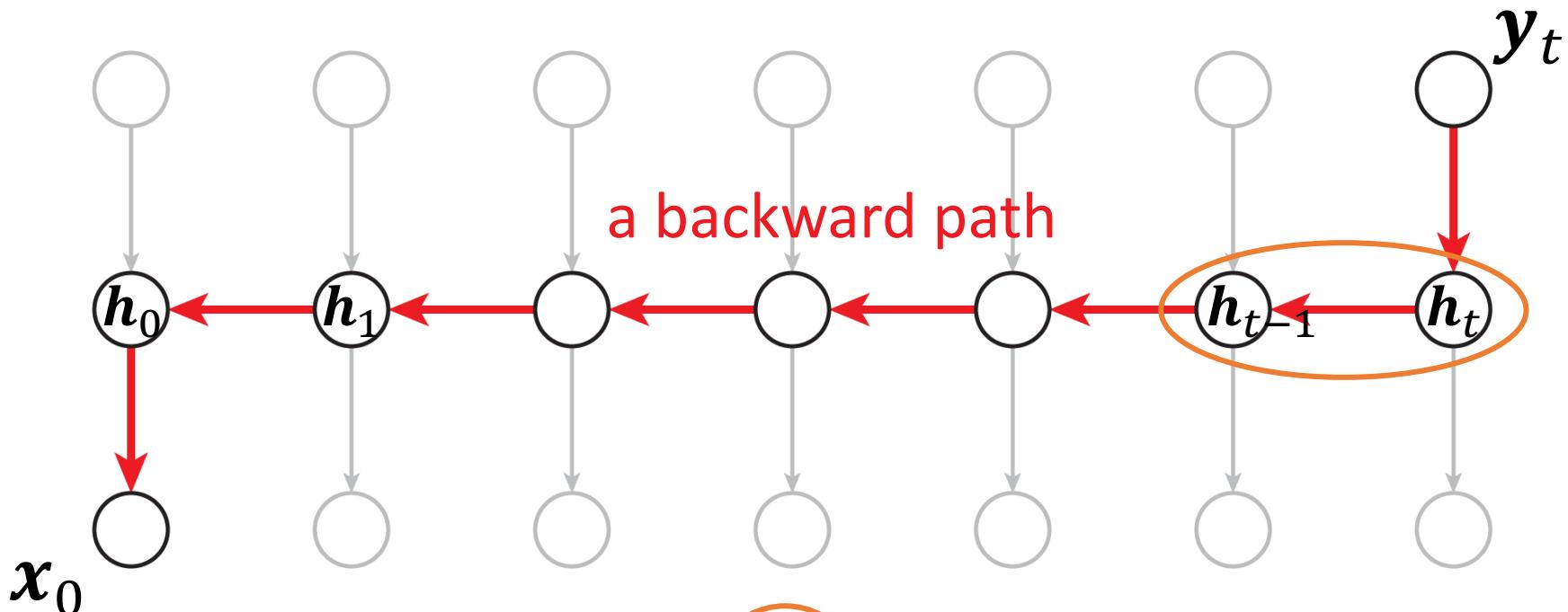
- output at time t is like a t -layer deep net



$$\frac{\partial y_t}{\partial x_0} = \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial h_0} \frac{\partial h_0}{\partial x_0}$$

RNN: BackProp Through Time (BPTT)

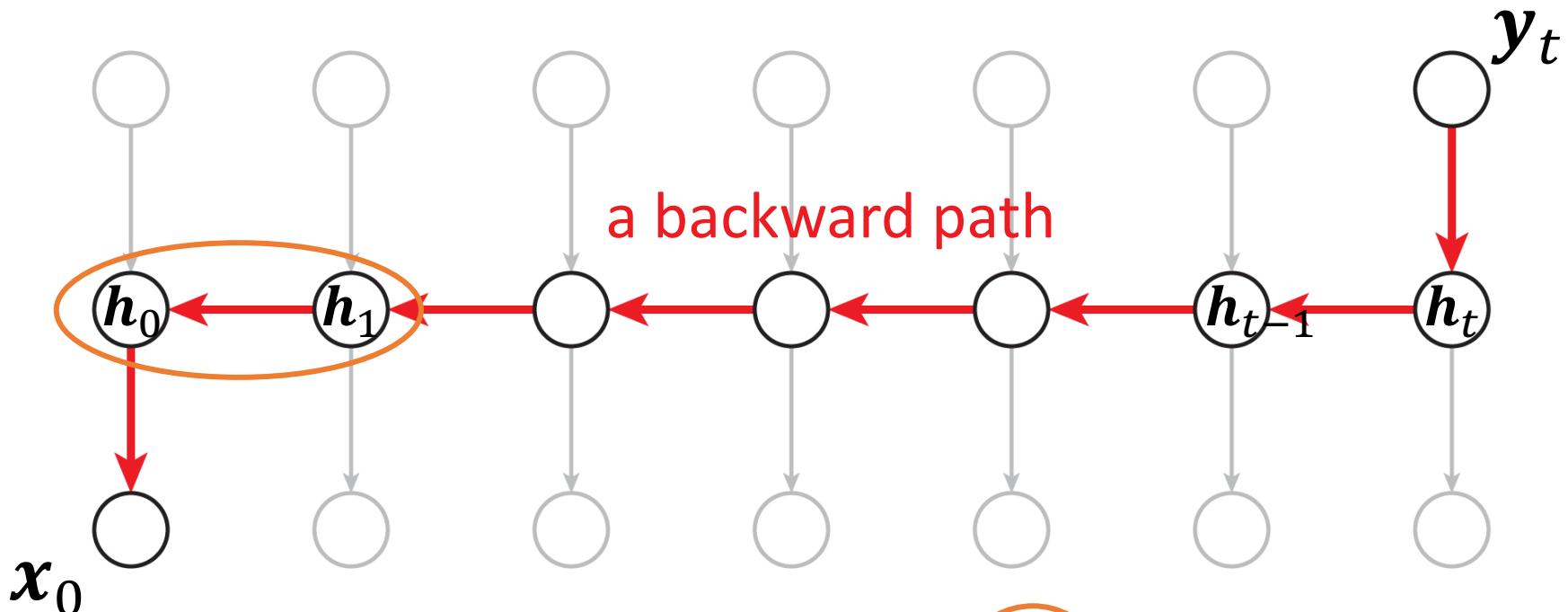
- output at time t is like a t -layer deep net



$$\frac{\partial y_t}{\partial x_0} = \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial h_0} \frac{\partial h_0}{\partial x_0}$$

RNN: BackProp Through Time (BPTT)

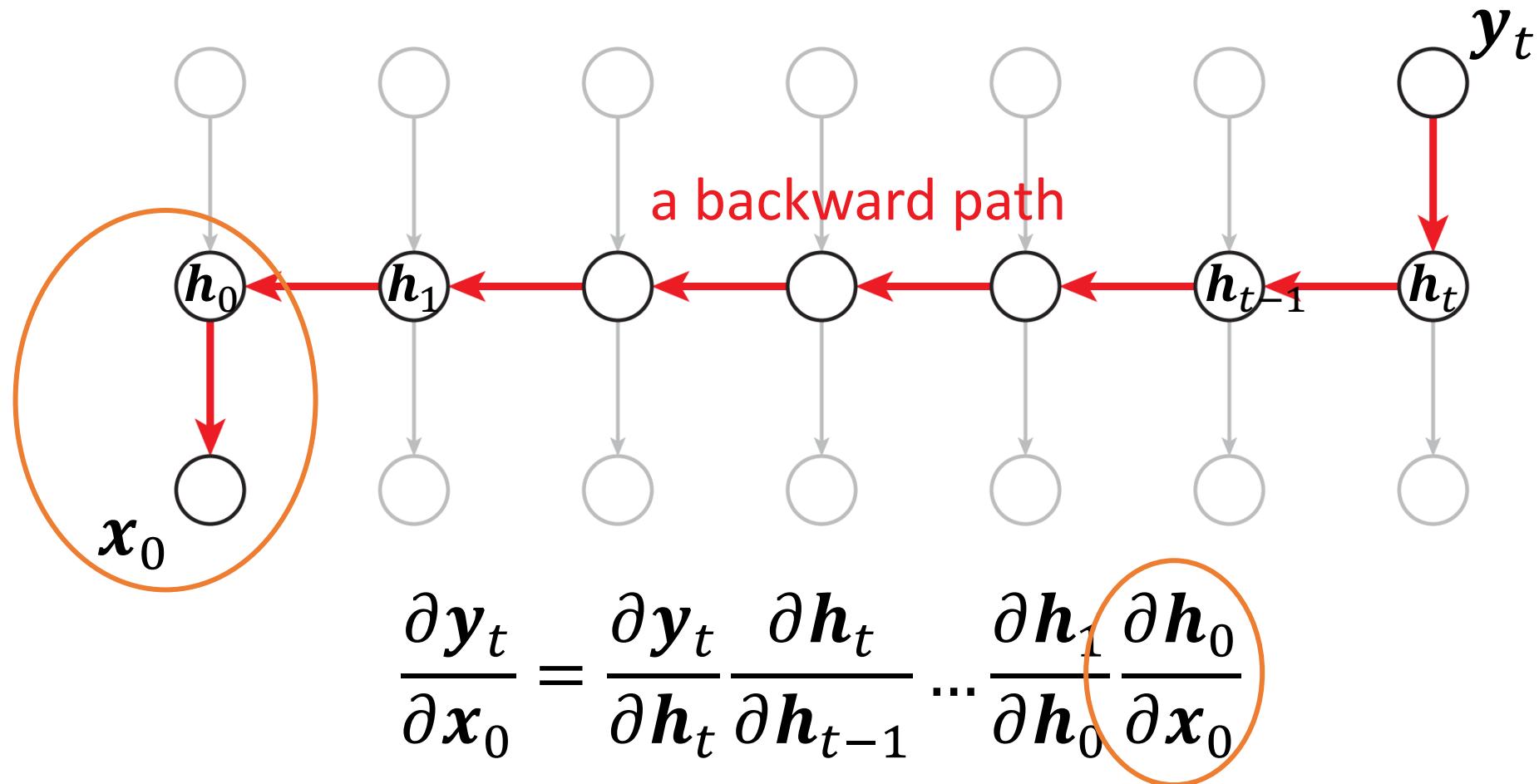
- output at time t is like a t -layer deep net



$$\frac{\partial y_t}{\partial x_0} = \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial h_0} \frac{\partial h_0}{\partial x_0}$$

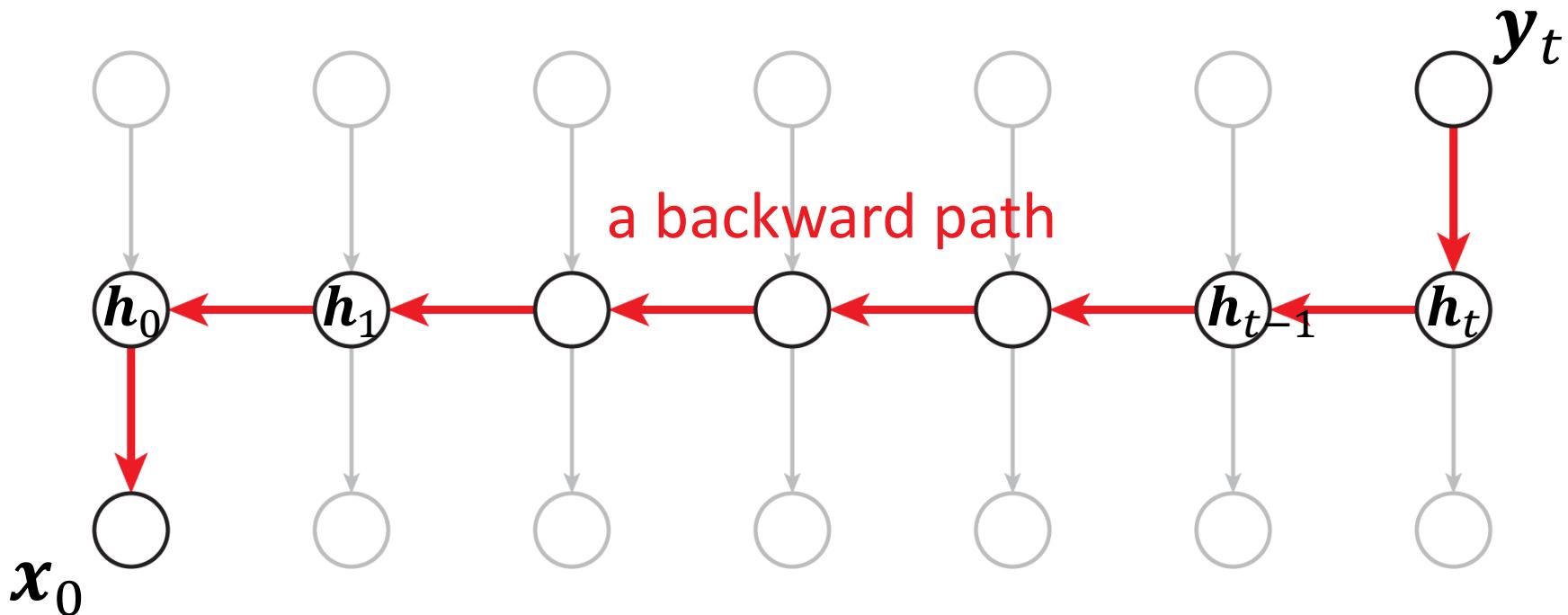
RNN: BackProp Through Time (BPTT)

- output at time t is like a t -layer deep net



RNN: BackProp Through Time (BPTT)

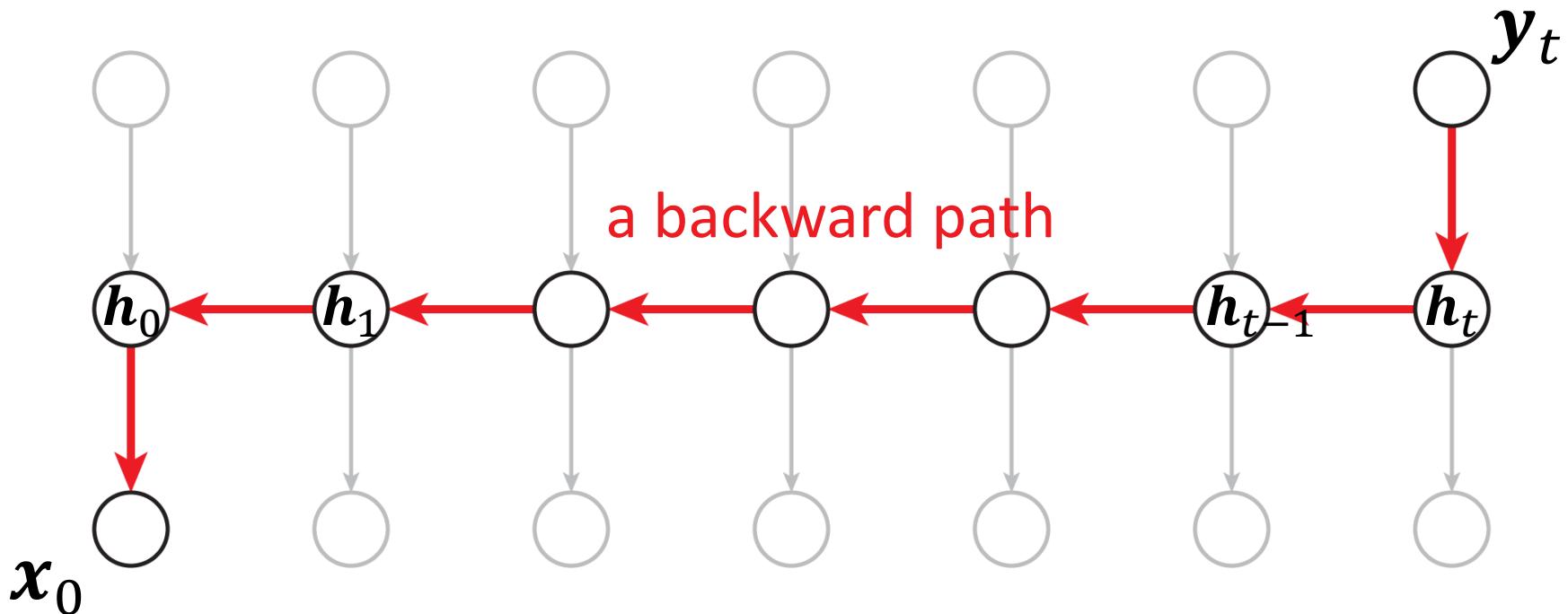
Example (linear vanilla RNN): $\mathbf{h}_t = W\mathbf{h}_{t-1} + U\mathbf{x}_t$



$$\frac{\partial \mathbf{y}_t}{\partial \mathbf{x}_0} = \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \cdots \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} \frac{\partial \mathbf{h}_0}{\partial \mathbf{x}_0}$$

RNN: BackProp Through Time (BPTT)

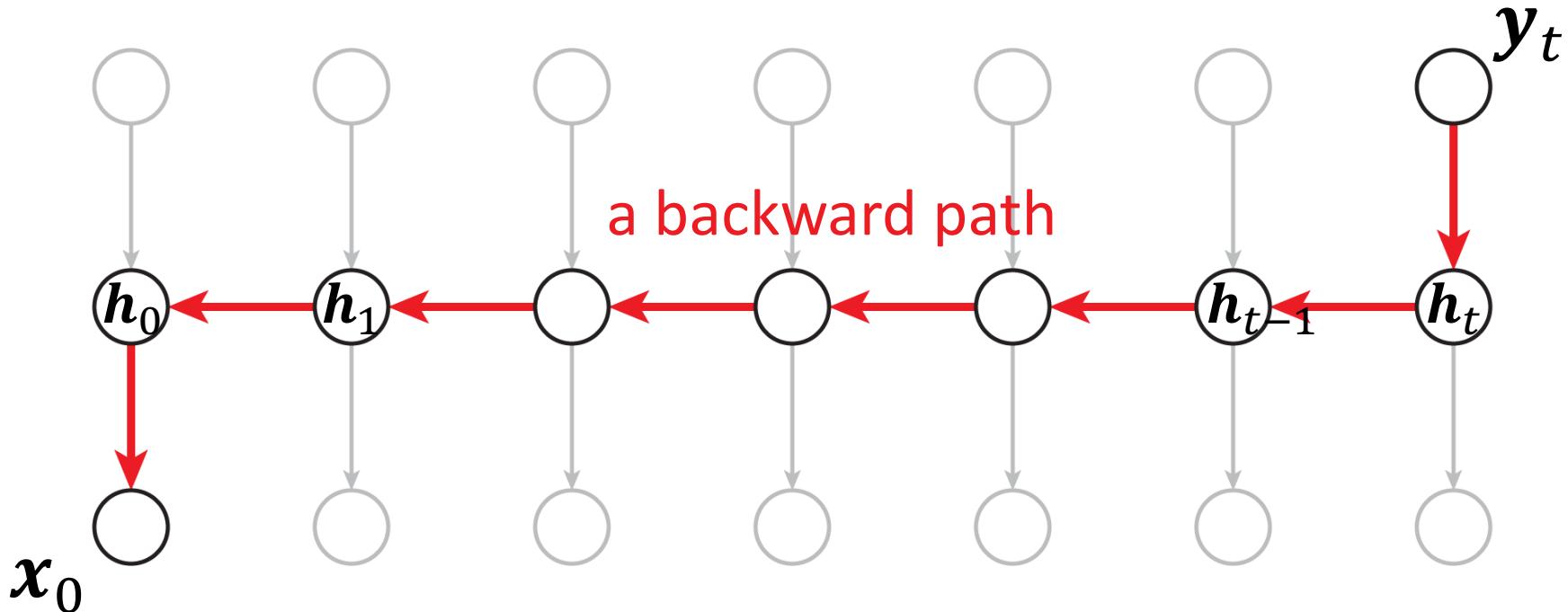
Example (linear vanilla RNN): $\mathbf{h}_t = W\mathbf{h}_{t-1} + U\mathbf{x}_t$



$$\frac{\partial \mathbf{y}_t}{\partial \mathbf{x}_0} = \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \cdot W \quad \dots \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} \frac{\partial \mathbf{h}_0}{\partial \mathbf{x}_0}$$

RNN: BackProp Through Time (BPTT)

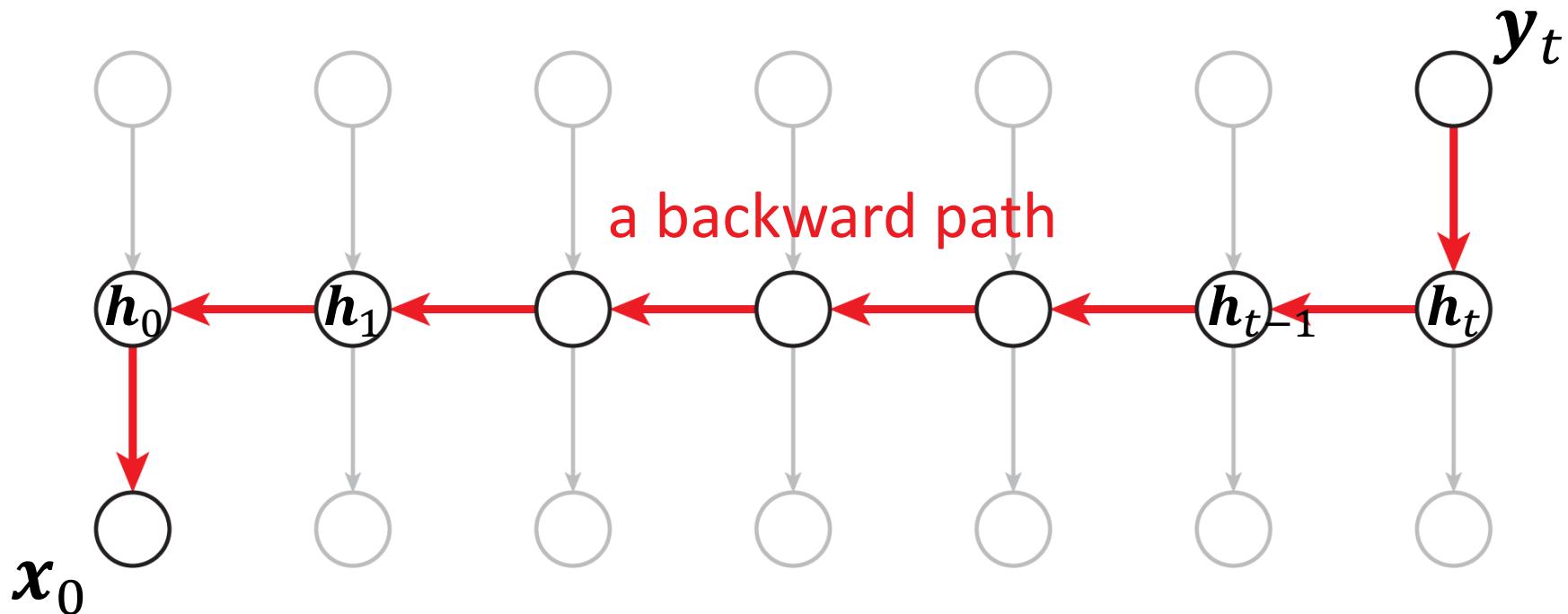
Example (linear vanilla RNN): $\mathbf{h}_t = W\mathbf{h}_{t-1} + U\mathbf{x}_t$



$$\frac{\partial \mathbf{y}_t}{\partial \mathbf{x}_0} = \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \underbrace{W \dots W}_{t \times} \frac{\partial \mathbf{h}_0}{\partial \mathbf{x}_0}$$

RNN: BackProp Through Time (BPTT)

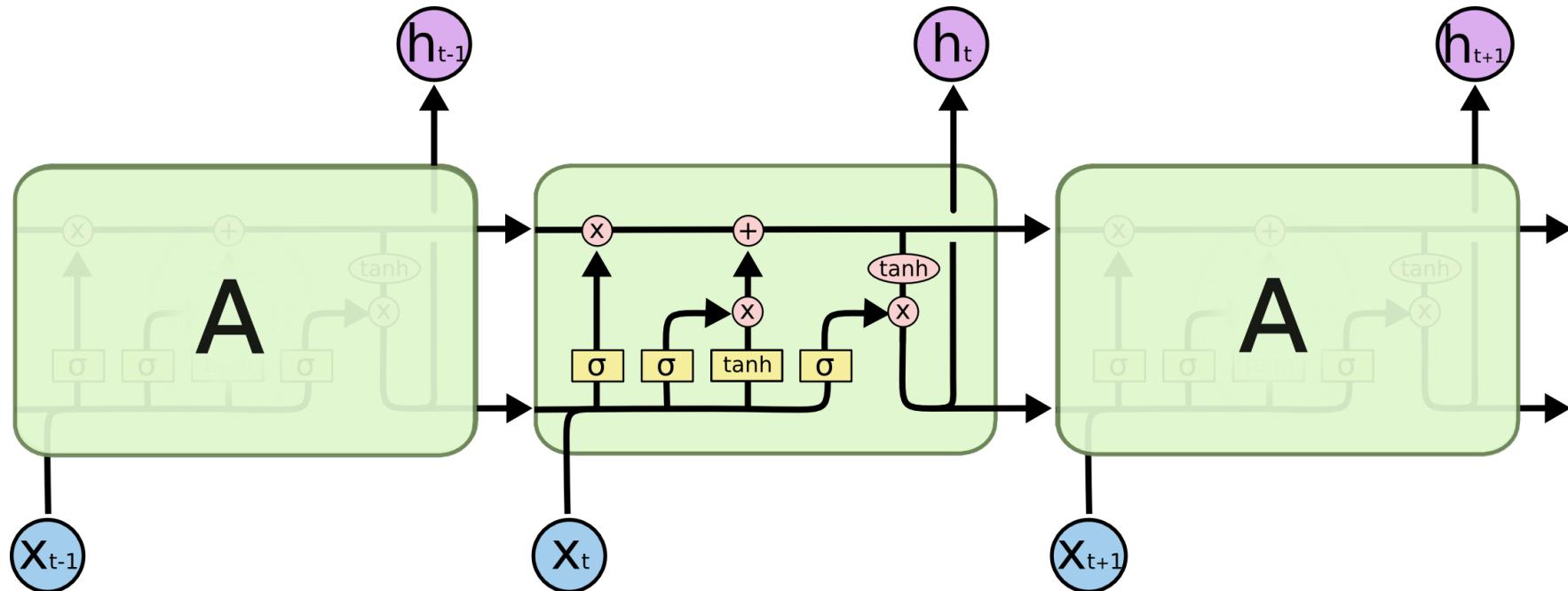
Example (linear vanilla RNN): $\mathbf{h}_t = W\mathbf{h}_{t-1} + U\mathbf{x}_t$



$$\frac{\partial \mathbf{y}_t}{\partial \mathbf{x}_0} = \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \cdot W^t \frac{\partial \mathbf{h}_0}{\partial \mathbf{x}_0}$$

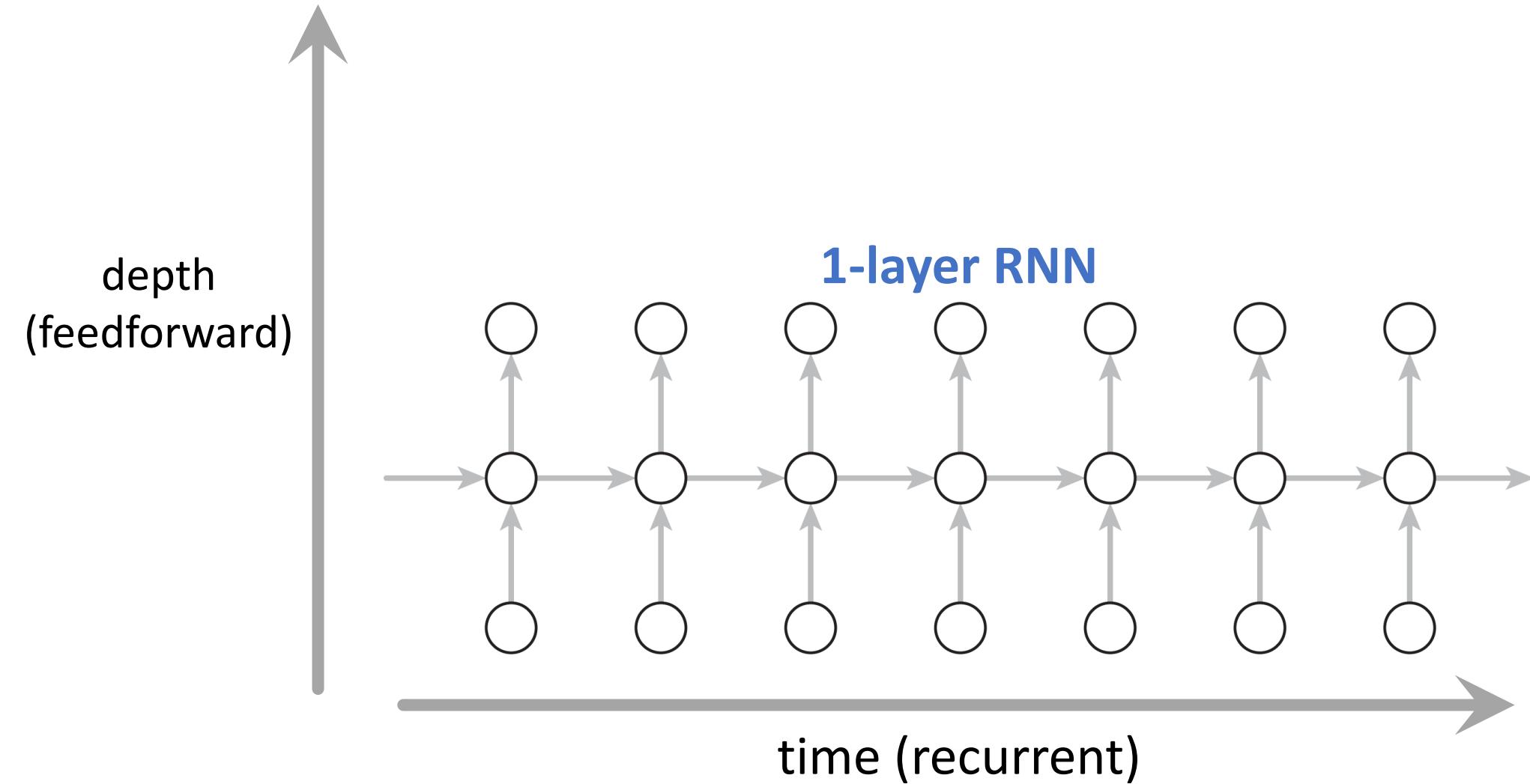
Long Short-Term Memory (LSTM) unit

- Recurrent NN w/ gating

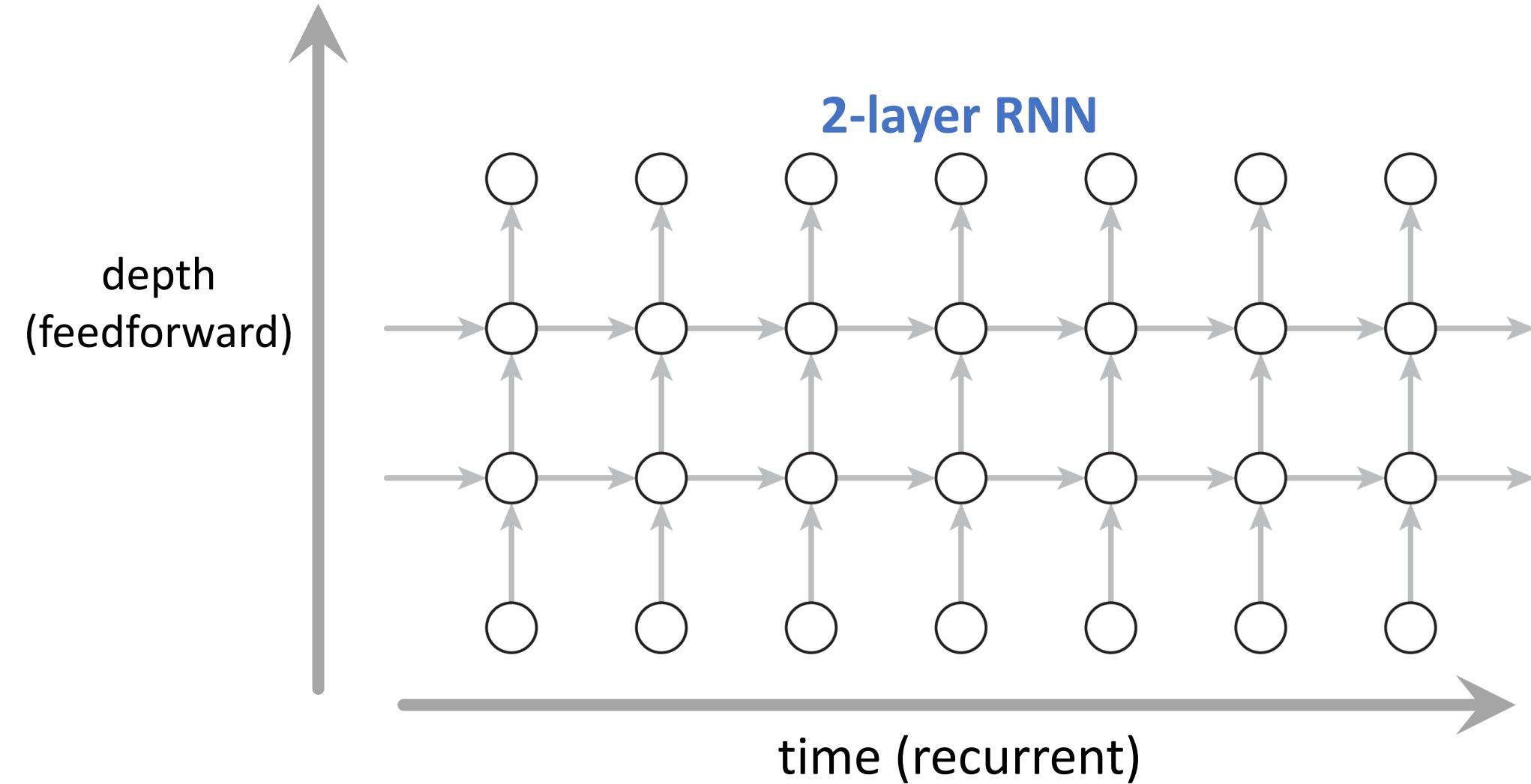


See Christopher Olah's blog (highly recommended!)
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Deep RNN

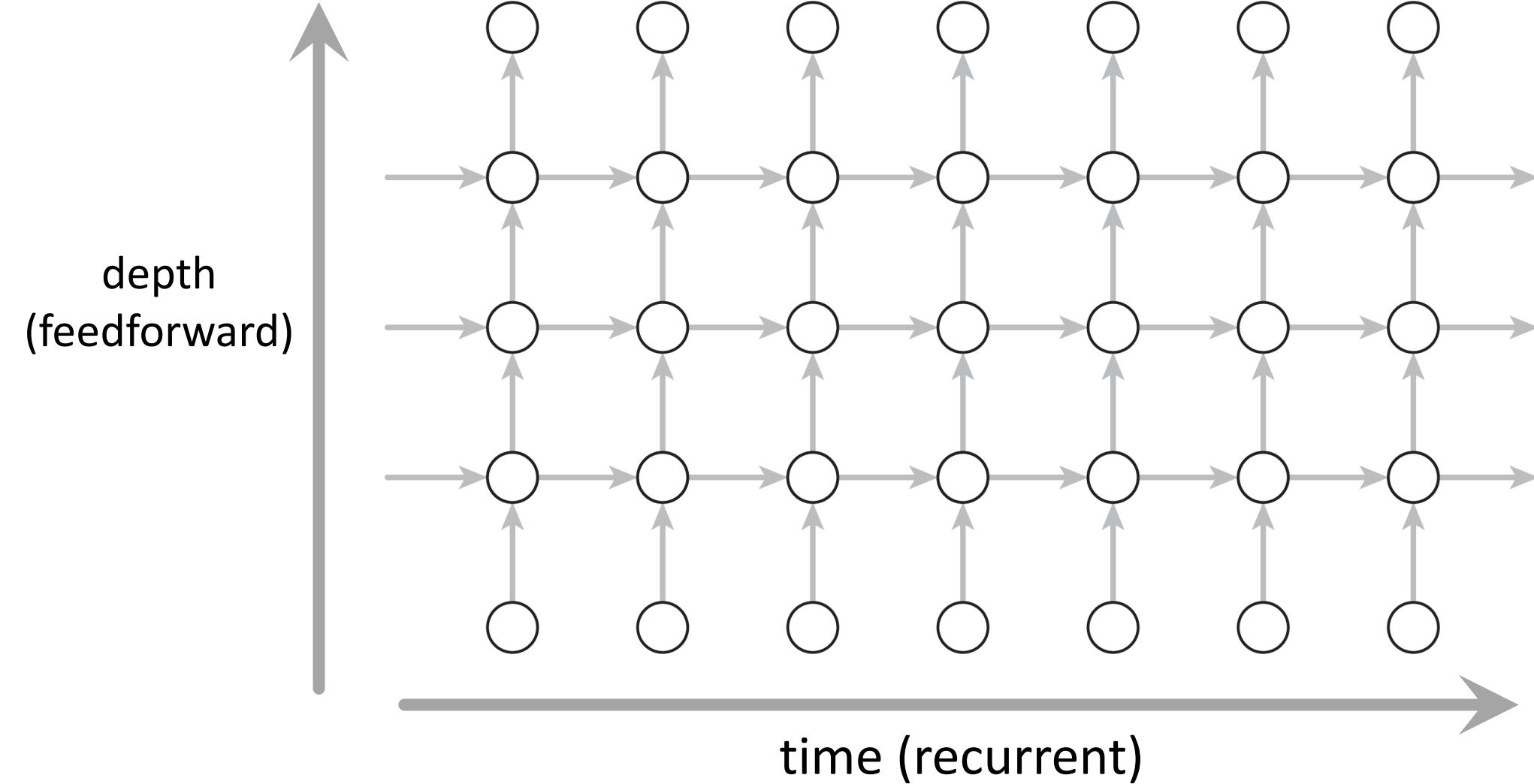


Deep RNN



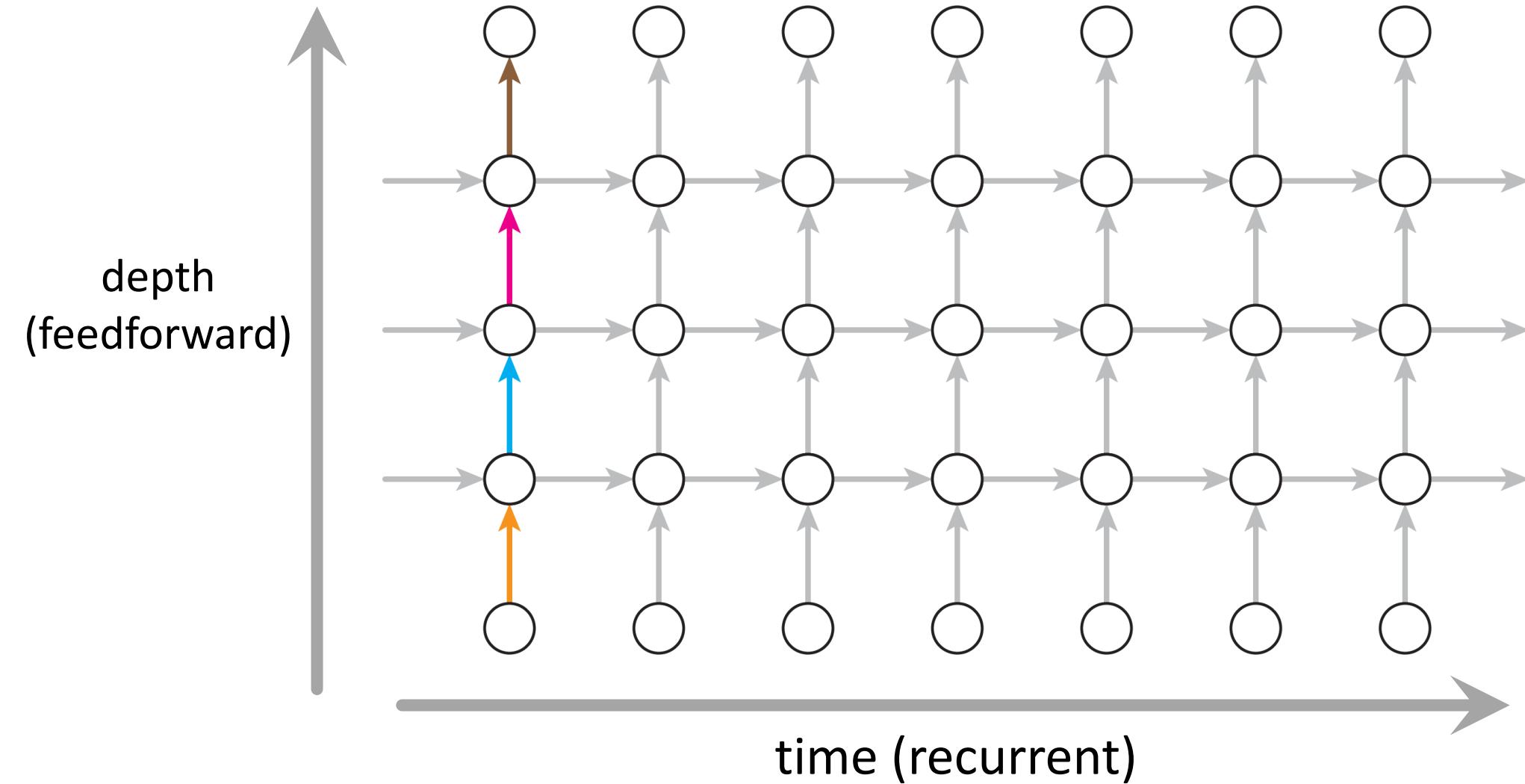
Deep RNN

3-layer RNN



Deep RNN

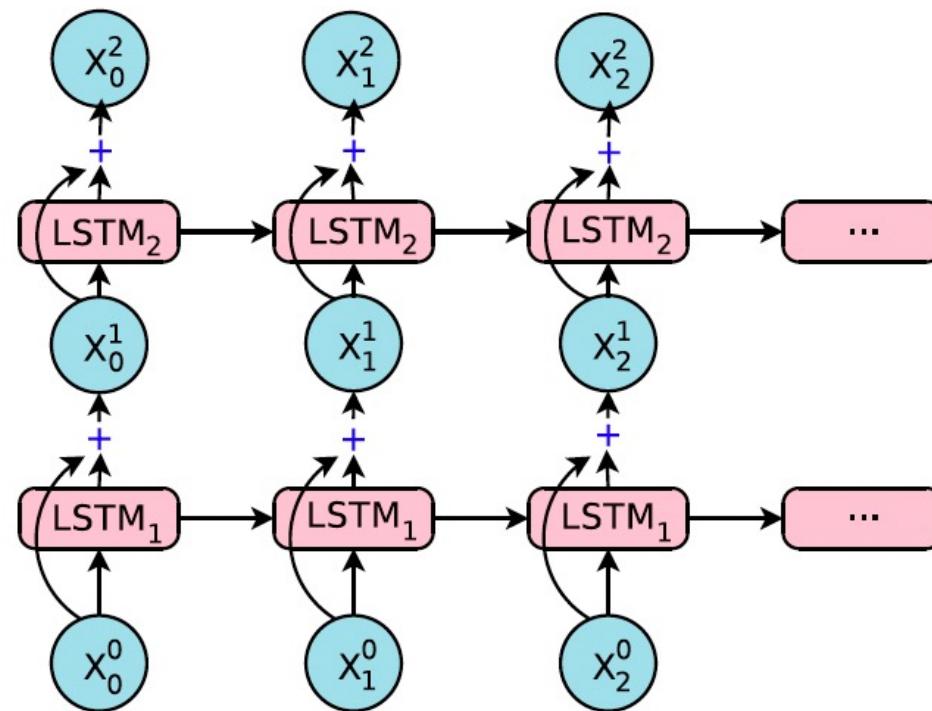
going **deep** in this direction



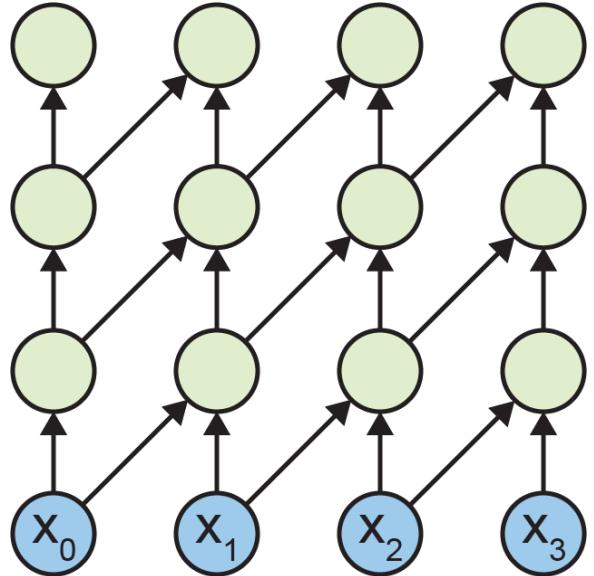
Example: Google's Neural Machine Translation System 2016

Deep RNN with LSTM units

- 16 layers enabled w/ residual connections
- degrade >4 layers w/o residual connections

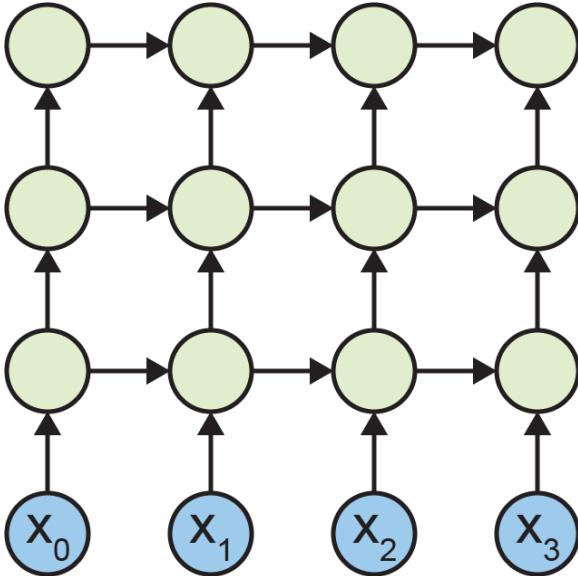


Sequence Modeling: CNN vs. RNN



CNN

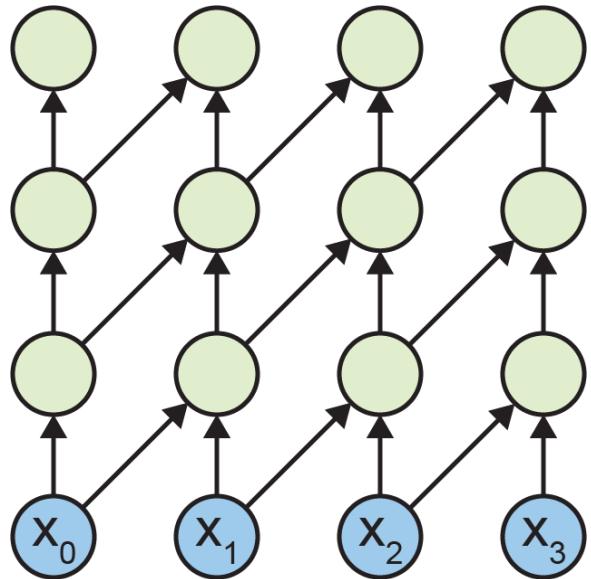
- limited context
- feedforward



RNN

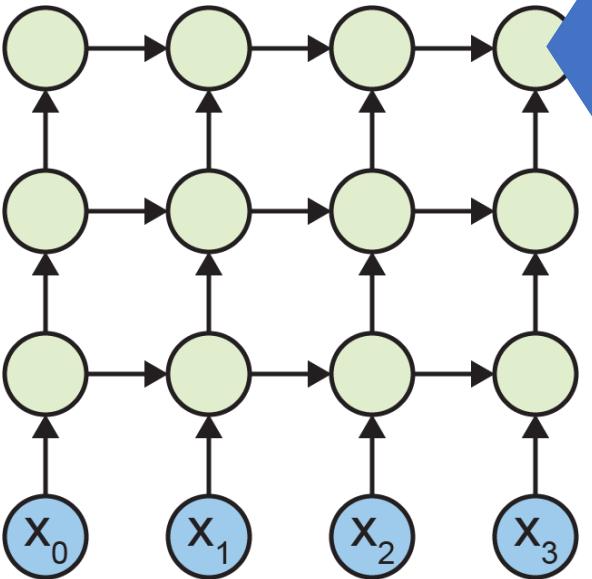
- long context
- not feedforward

Sequence Modeling: CNN vs. RNN



CNN

- limited context
- feedforward

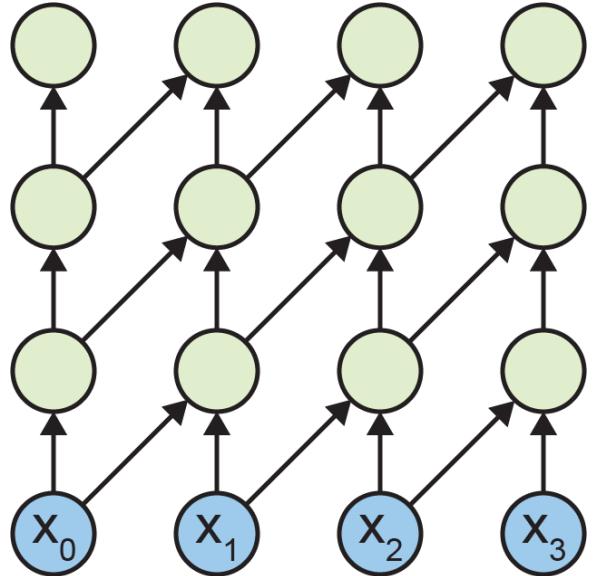


RNN

- long context
- not feedforward

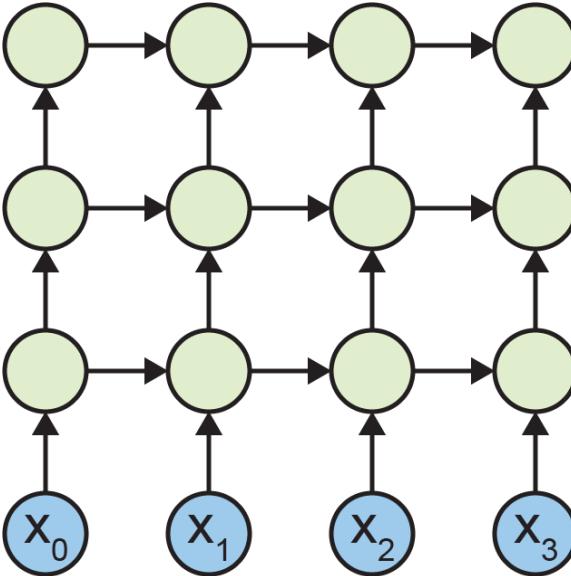
Not feedforward:
step at time t must wait for all
steps at time $< t$ to complete.
(not parallelism-friendly)

Sequence Modeling: CNN vs. RNN vs. ?



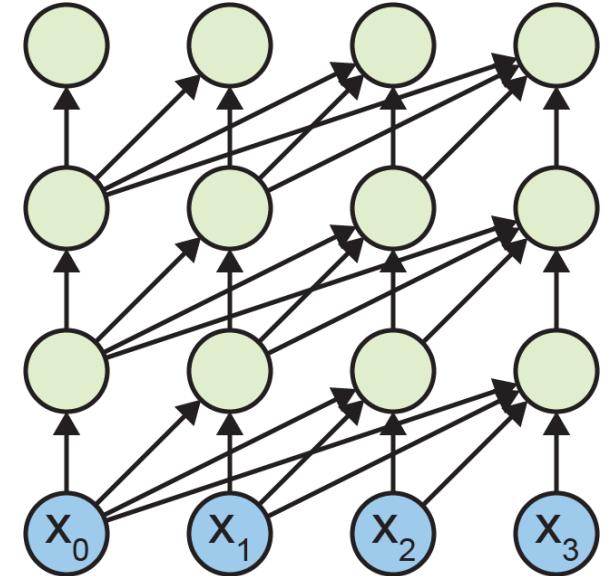
CNN

- limited context
- feedforward



RNN

- long context
- not feedforward



Attention

best of both world?

Attention and Transformers

Dictionary look-up

Think of a Python dict:

```
dict_fr2en = {  
    "pomme": "apple",  
    "banane": "banana",  
    "citron": "lemon"  
}
```

keys **values**

pomme : apple

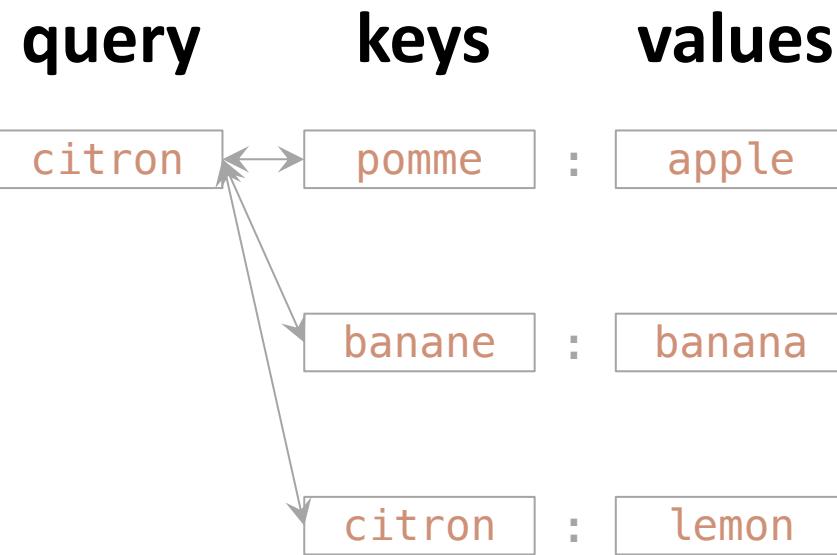
banane : banana

citron : lemon

Dictionary look-up

Think of a Python dict:

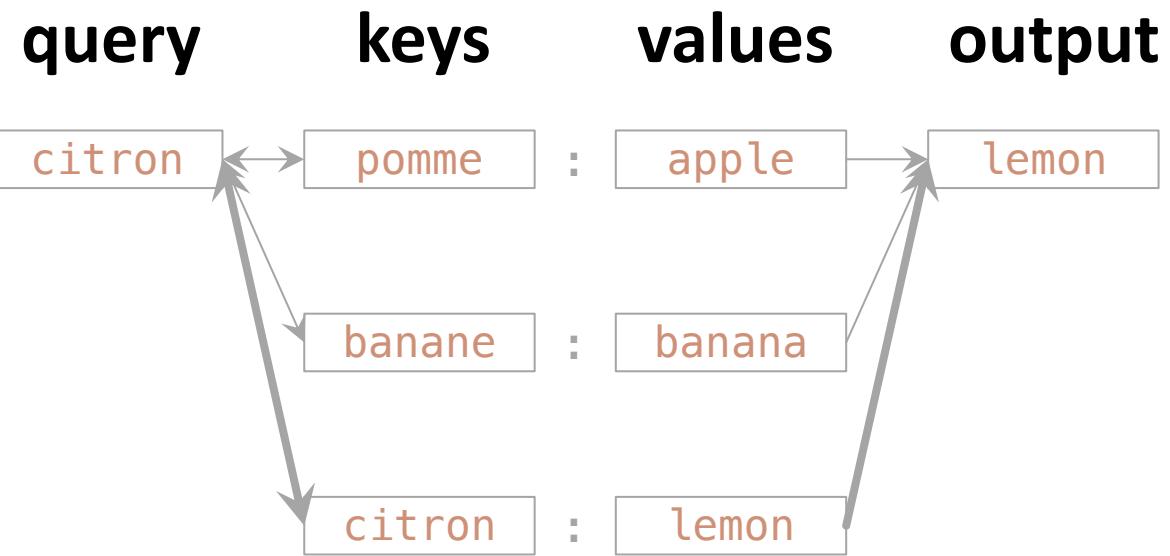
```
dict_fr2en = {  
    "pomme": "apple",  
    "banane": "banana",  
    "citron": "lemon"  
}  
  
query = "citron"  
output = dict_fr2en[query]
```



Dictionary look-up

Think of a Python dict:

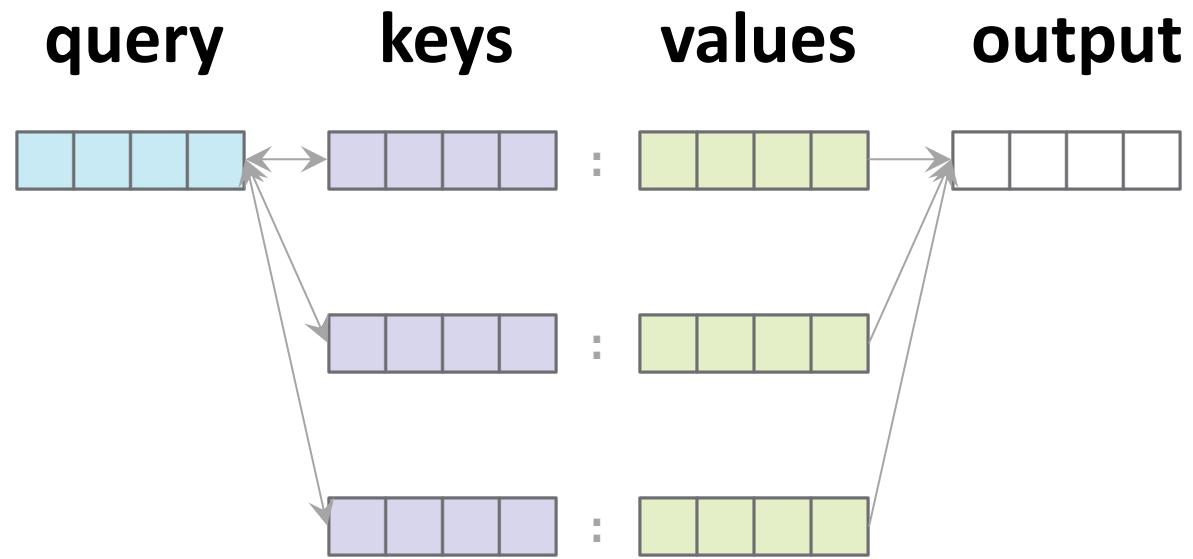
```
dict_fr2en = {  
    "pomme": "apple",  
    "banane": "banana",  
    "citron": "lemon"  
}  
  
query = "citron"  
output = dict_fr2en[query]
```



Dictionary look-up

Think of a Python dict:

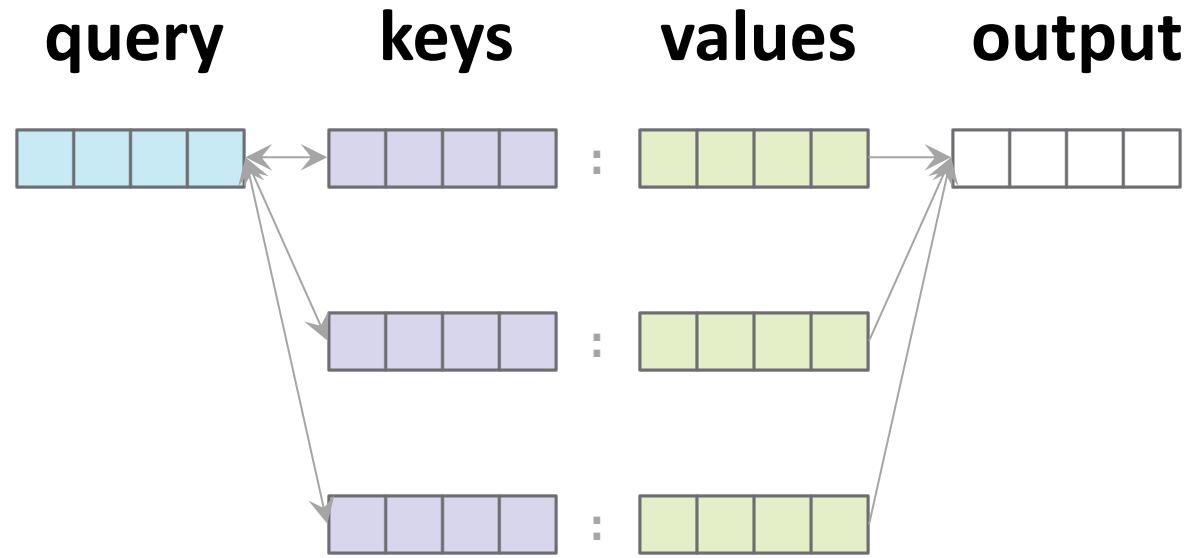
```
dict_fr2en = {  
    "pomme": "apple",  
    "banane": "banana",  
    "citron": "lemon"  
}  
  
query = "citron"  
output = dict_fr2en[query]
```



What if everything is an abstract representation (“embedding”)?

Attention

Attention is **soft** dict look-up in the representation space.

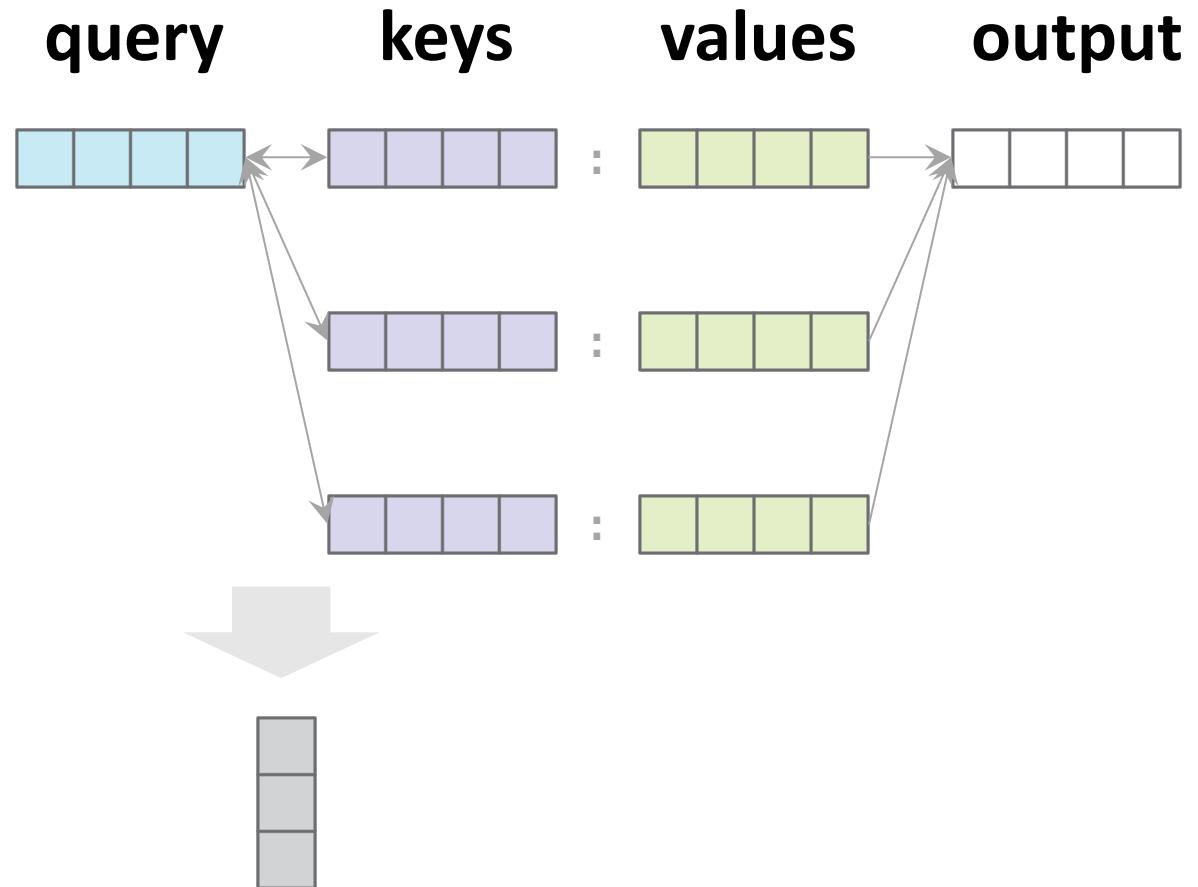


Attention

Attention is **soft** dict look-up in the representation space.

1. query-key similarity:

$$s_i = \mathbf{q} \cdot \mathbf{k}_i$$



Attention

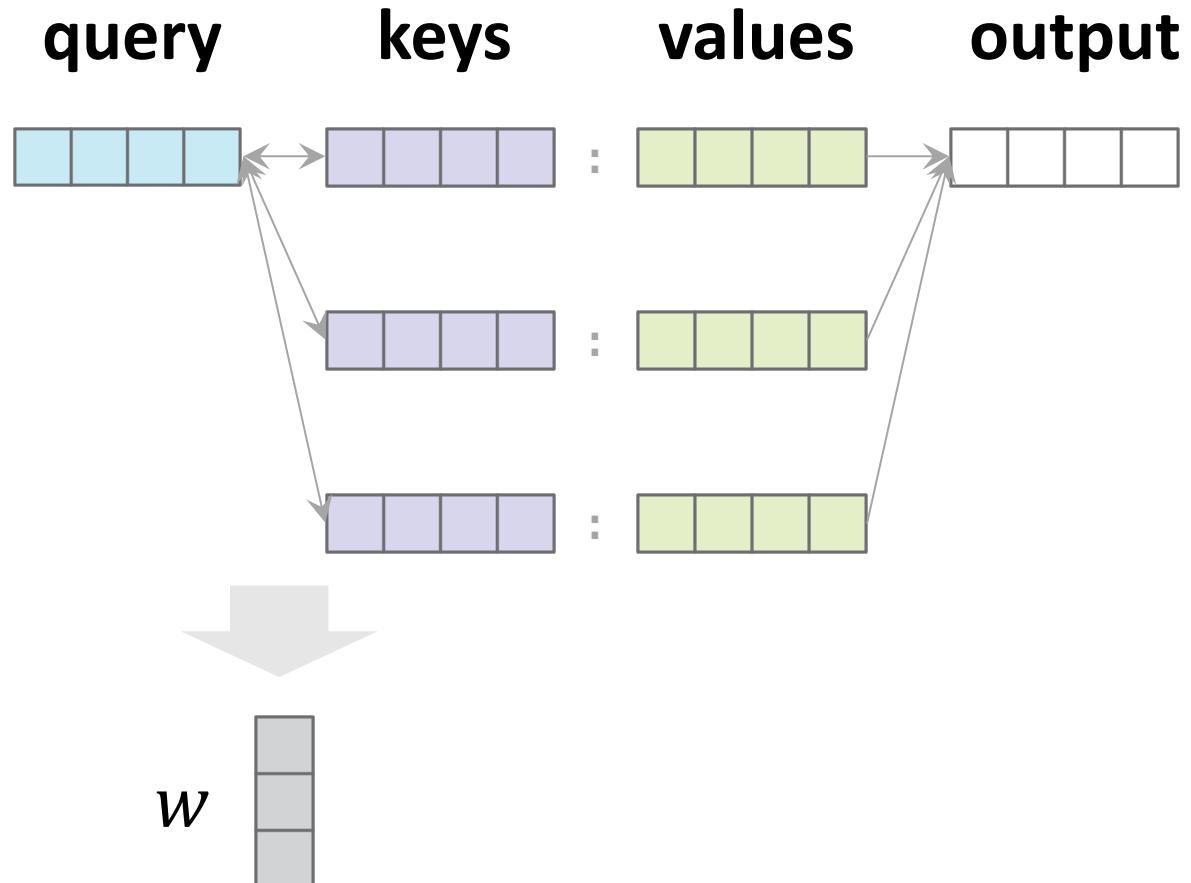
Attention is **soft** dict look-up in the representation space.

1. query-key similarity:

$$s_i = \mathbf{q} \cdot \mathbf{k}_i$$

2. attention weight (softmax):

$$w_i = e^{s_i} / \sum_j e^{s_j}$$



*In practice, the logits s_i are scaled by $1/\sqrt{\text{dim}}$ before softmax.

Attention

Attention is **soft** dict look-up in the representation space.

1. query-key similarity:

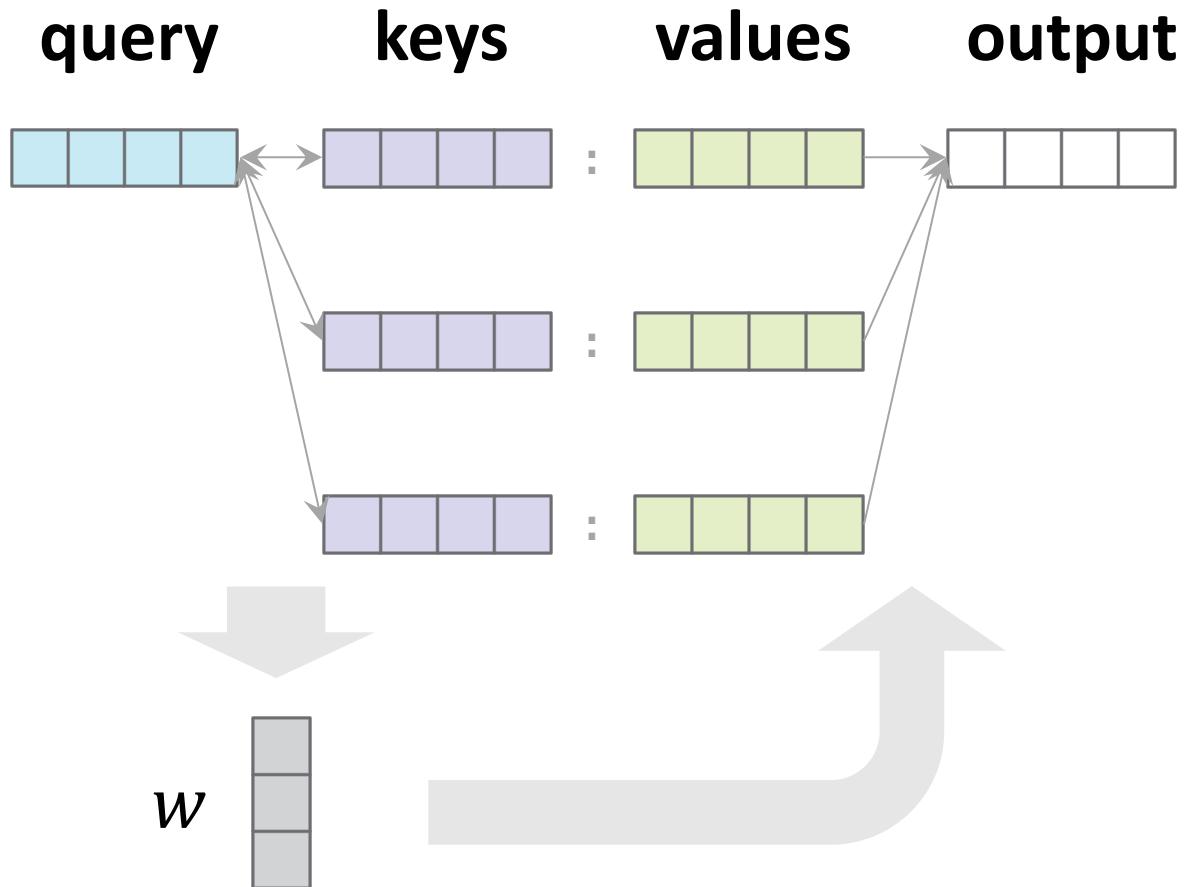
$$s_i = \mathbf{q} \cdot \mathbf{k}_i$$

2. attention weight (softmax):

$$w_i = e^{s_i} / \sum_j e^{s_j}$$

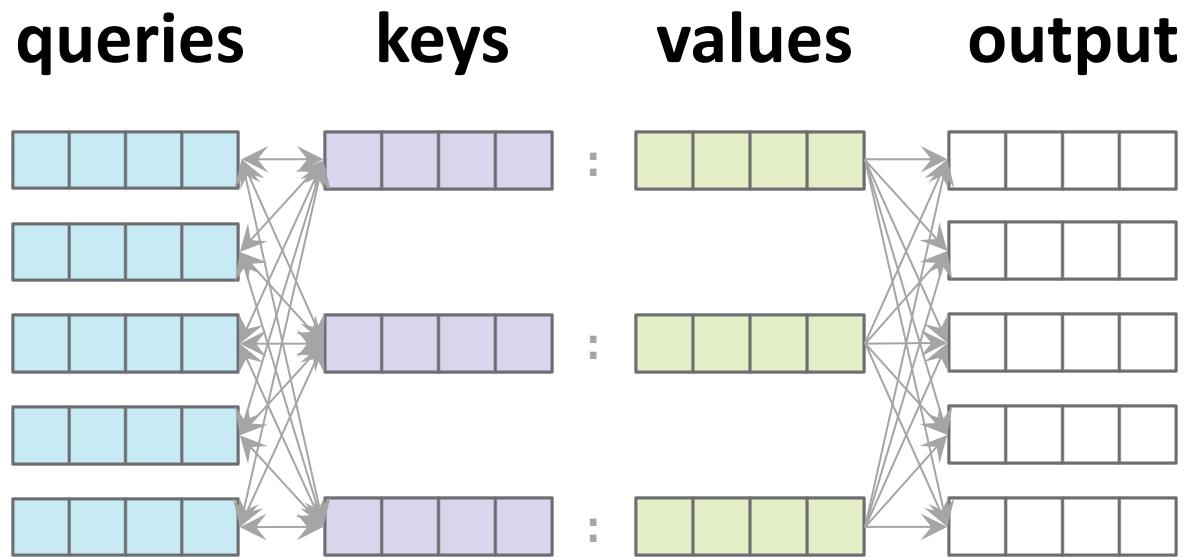
3. weighted average:

$$\mathbf{z} = \sum_i w_i \mathbf{v}_i$$



Attention

Multiple queries

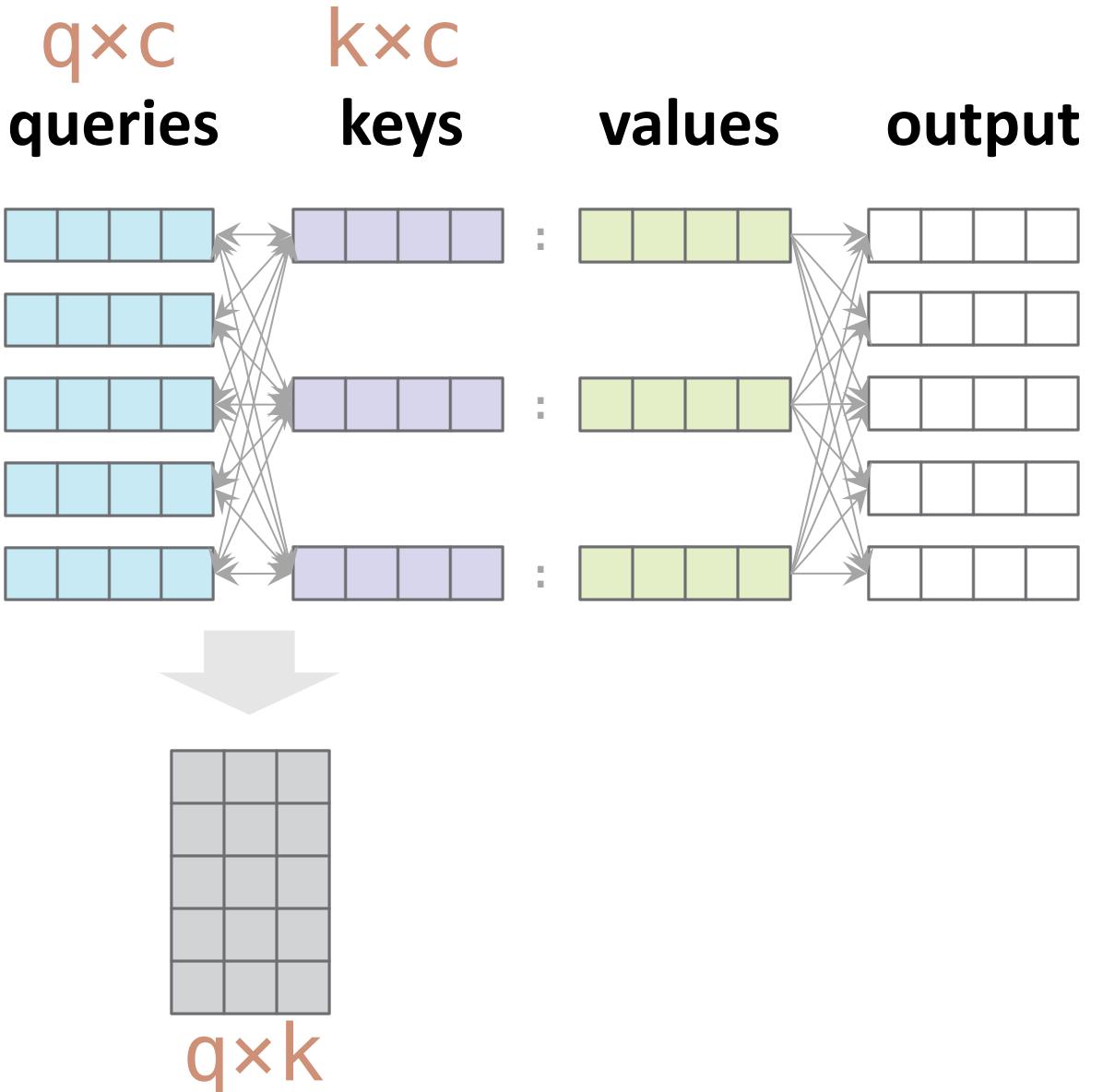


Attention

Multiple queries

1. query-key similarity:

$$S = QK^T$$



Attention

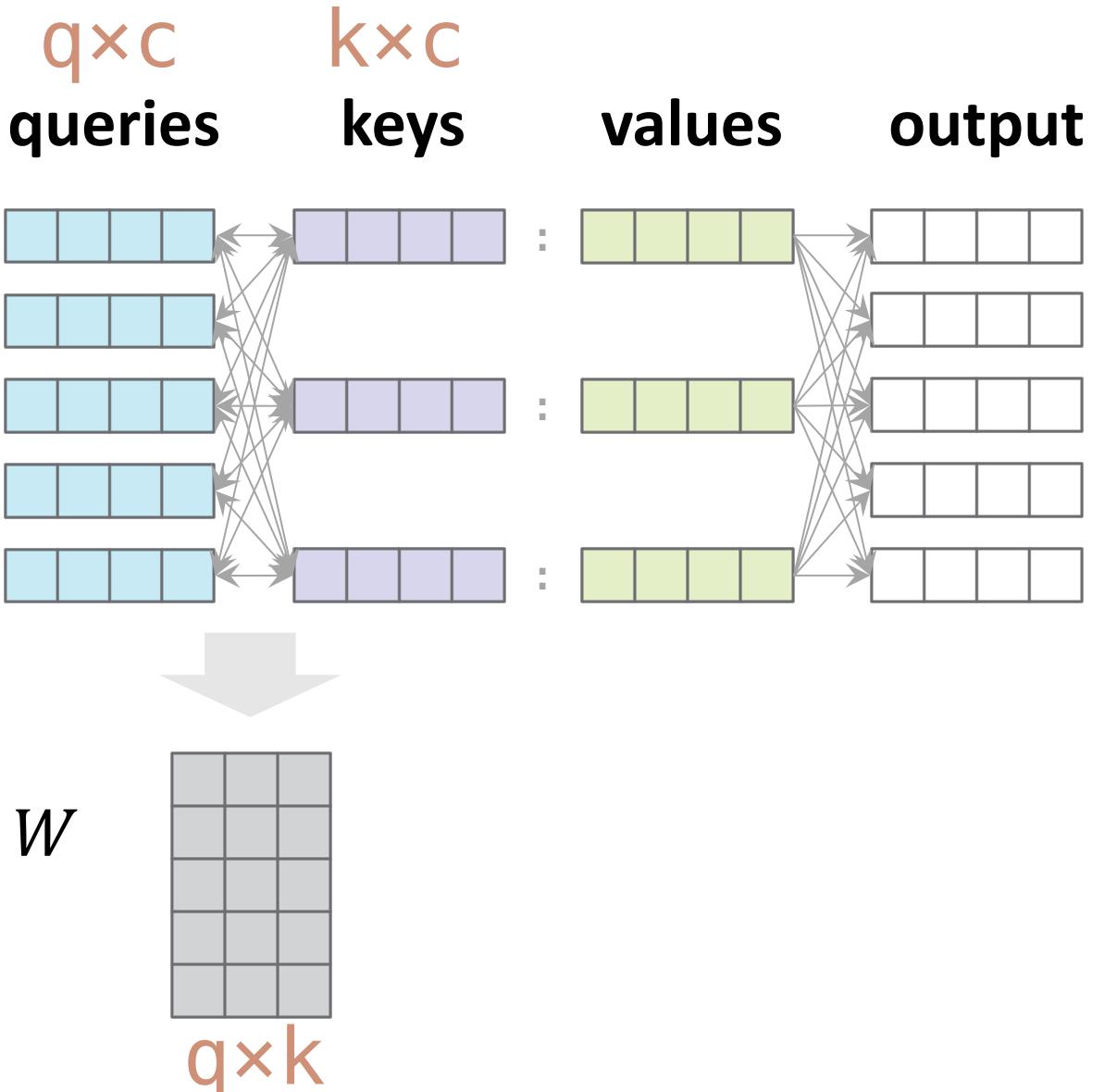
Multiple queries

1. query-key similarity:

$$S = QK^T$$

2. attention weight (softmax):

$$W = \text{softmax}(S)$$



*In practice, the logits s_i are scaled by $1/\sqrt{\text{dim}}$ before softmax.

Attention

Multiple queries

1. query-key similarity:

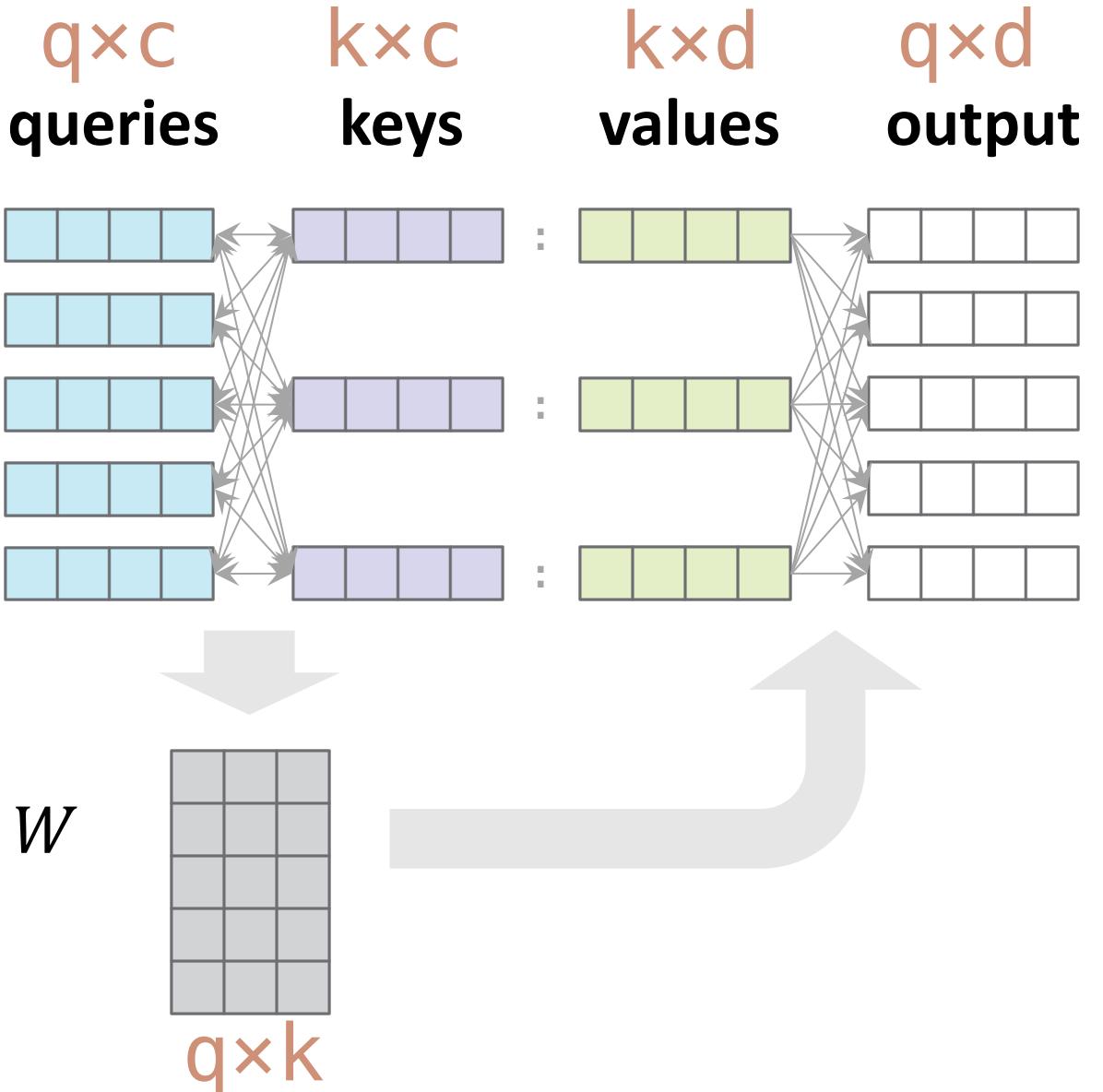
$$S = QK^T$$

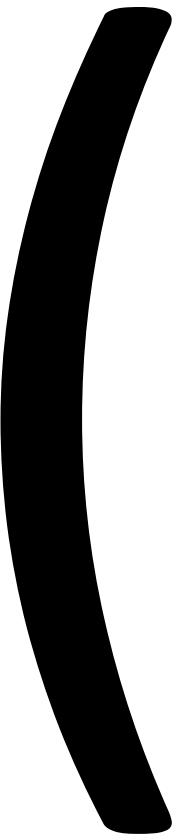
2. attention weight (softmax):

$$W = \text{softmax}(S)$$

3. weighted average:

$$Z = WV$$





Einstein Sum: notations of brevity

Einstein Sum: notations of brevity

Traditional notation (e.g., matrix-matrix product):

let $A \in \mathbb{R}^{I \times K}$, $B \in \mathbb{R}^{K \times J}$,

$$C = AB$$

where $C \in \mathbb{R}^{I \times J}$ with:

$$C_{ij} = \sum_k A_{ik}B_{kj}$$

Einstein Sum: notations of brevity

Traditional notation (e.g., matrix-matrix product):

let $A \in \mathbb{R}^{I \times K}$, $B \in \mathbb{R}^{K \times J}$,

$$C = AB$$

where $C \in \mathbb{R}^{I \times J}$ with:

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

most important
messages are in
subscripts

```
C = einsum('ik,kj->ij', A, B)
```

(An index that is summed over won't show in the output)

Examples

```
# column sum  
einsum('ij->j', X) # shape: (10, 5) -> (5,)  
  
# row sum  
einsum('ij->i', X) # shape: (10, 5) -> (10,)  
  
# sum  
einsum('ij->', X) # shape: (10, 5) -> scalar
```

(An index that is summed over won't show in the output)

Examples

```
# transpose  
einsum('ij->ji', X) # (10, 5) -> (5, 10)  
  
# permute  
einsum('nchw->nhwc', X) # (n, 3, 256, 256) -> (n, 256, 256, 3)
```

more intuitive than calling “permute” directly:

```
permute(X, [0, 2, 3, 1]) # (n, 3, 256, 256) -> (n, 256, 256, 3)
```

Examples

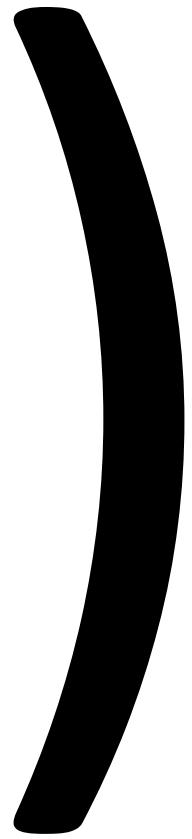
```
# inner prod  
einsum('i,i->', A, B) # (10,) x (10,) -> scalar  
  
# outer prod  
einsum('i,j->ij', A, B) # (10,) x (5,) -> (10, 5)  
  
# element-wise prod  
einsum('i,i->i', A, B) # (10,) x (10,) -> (10,)
```

Examples

```
# matrix-vector prod
einsum('ik,k->i', A, B) # (10, 4) x (4,) -> (10,)

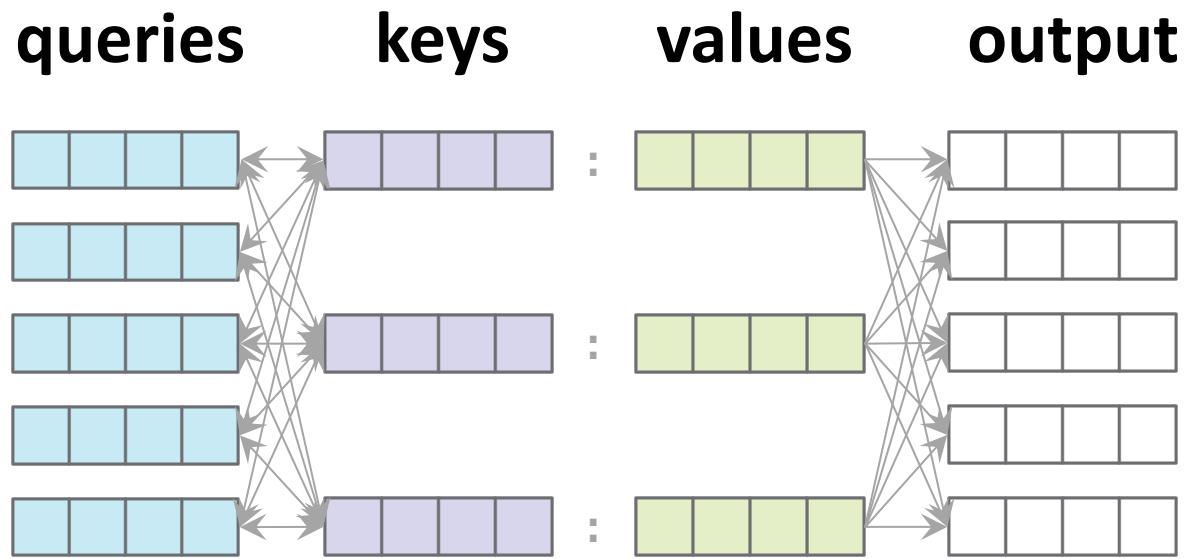
# matrix-matrix prod
einsum('ik,kj->ij', A, B) # (10, 4) x (4, 5) -> (10, 5)

# batch matrix-matrix prod
einsum('nik,nkj->nij', A, B) # (n, 10, 4) x (n, 4, 5) -> (n, 10, 5)
```



Attention

w/ Einsum notation

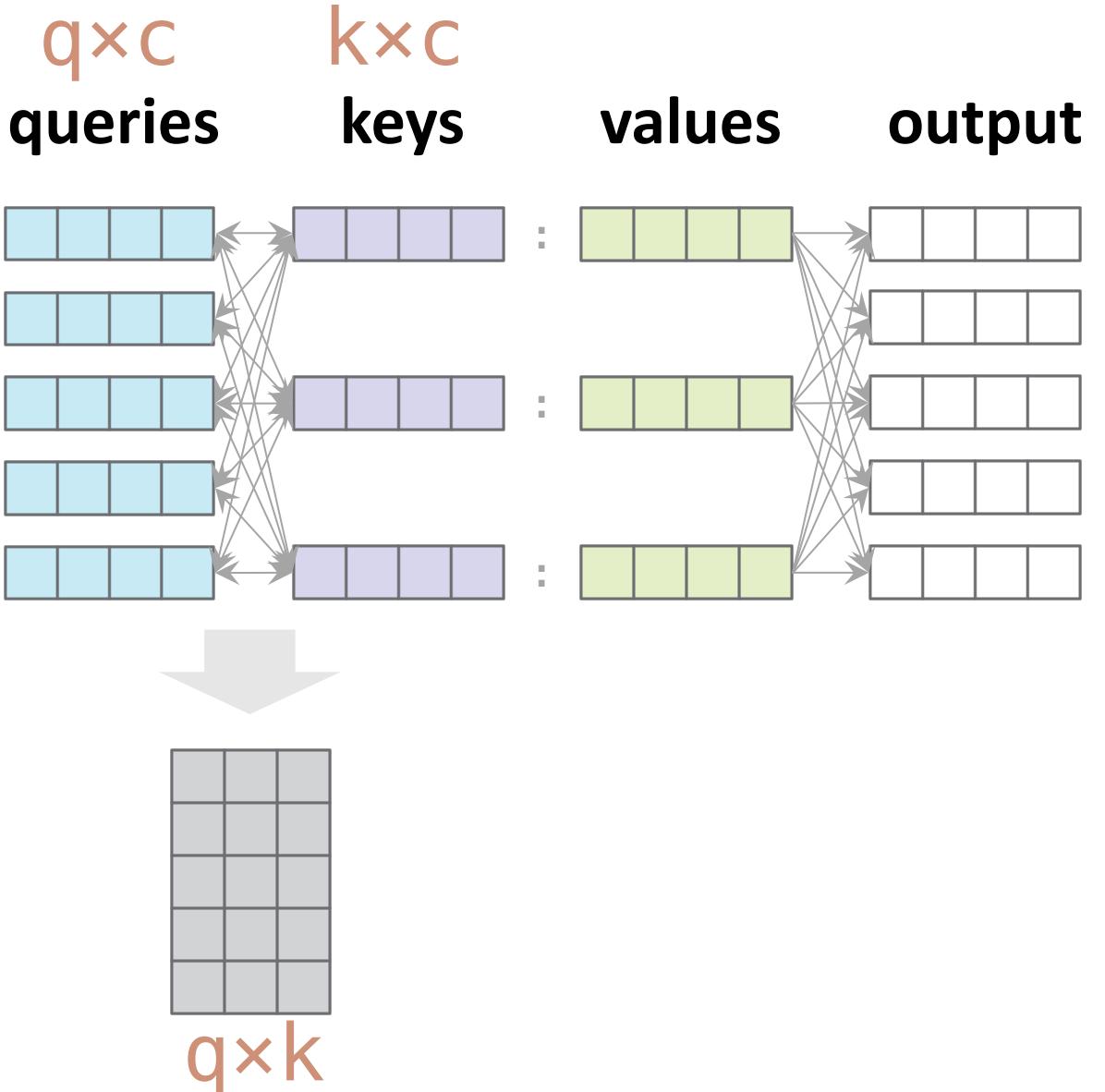


Attention

w/ Einsum notation

1. query-key similarity:

```
S = einsum('qc,kc->qk', Q, K)
```



Attention

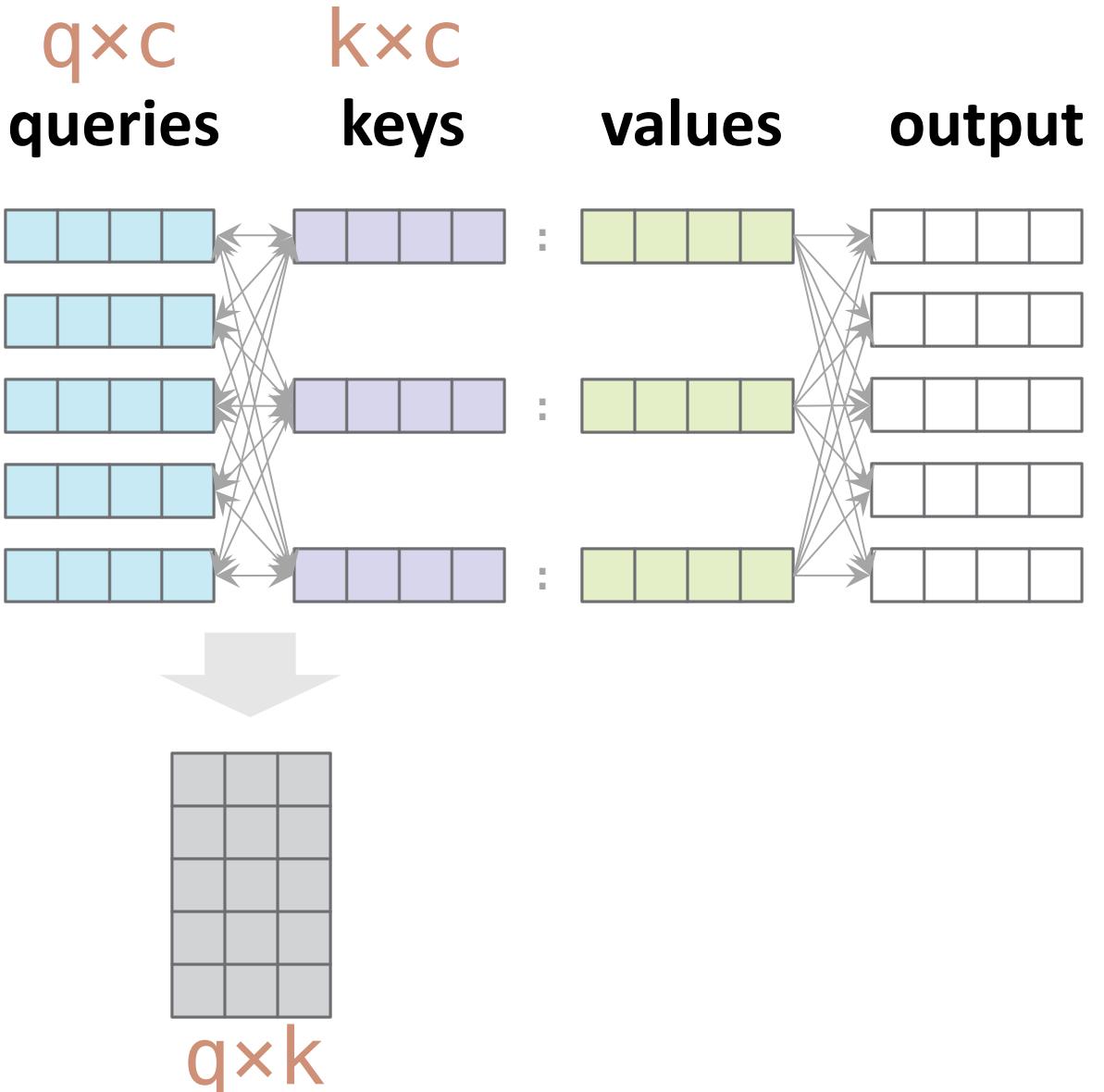
w/ Einsum notation

1. query-key similarity:

```
S = einsum('qc,kc->qk', Q, K)
```

2. attention weight (softmax):

```
W = softmax(S, dim=-1)
```



*In practice, the logits s_i are scaled by $1/\sqrt{\text{dim}}$ before softmax.

Attention

w/ Einsum notation

1. query-key similarity:

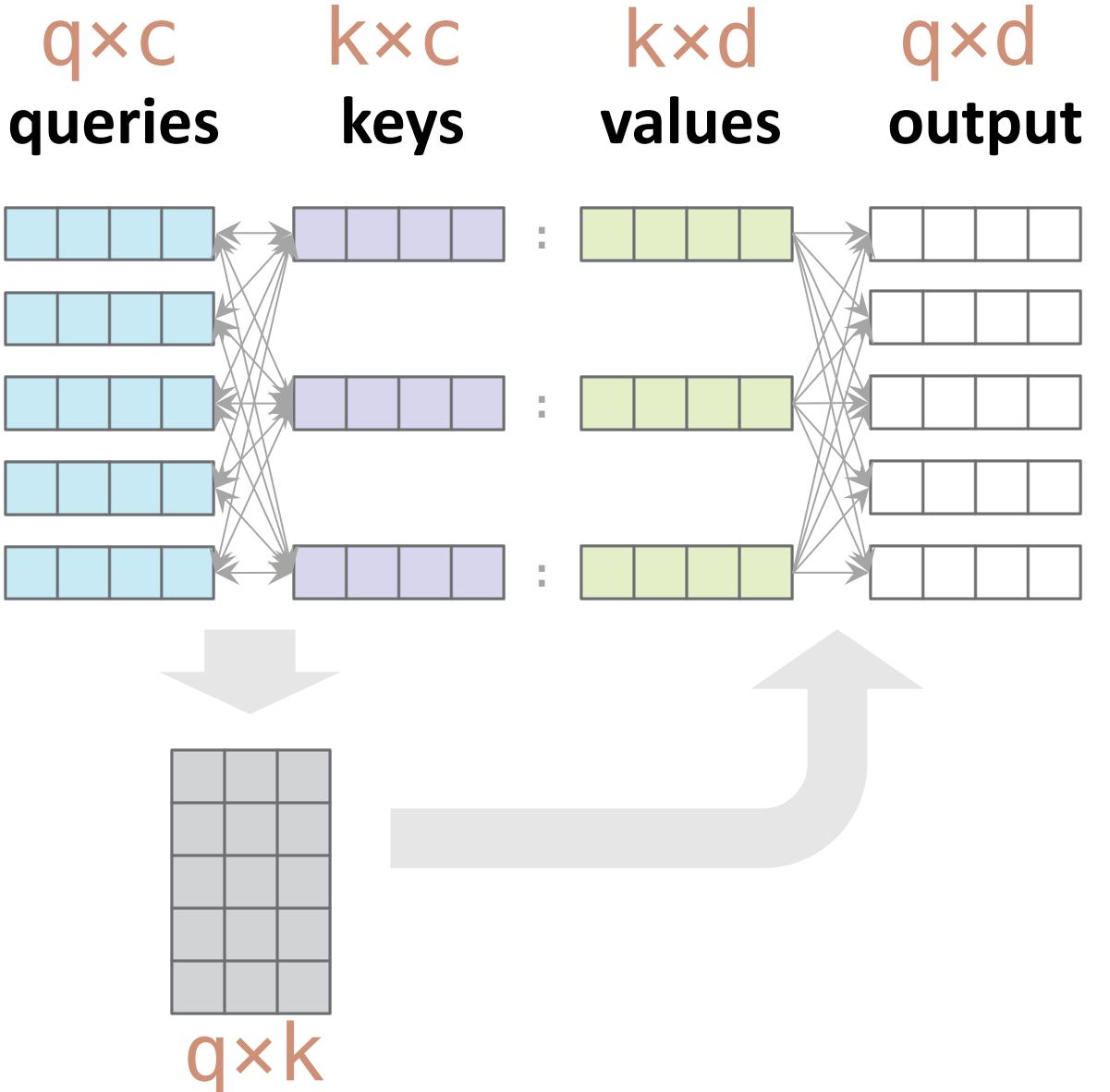
```
S = einsum('qc,kc->qk', Q, K)
```

2. attention weight (softmax):

```
W = softmax(S, dim=-1)
```

3. weighted average:

```
Z = einsum('qk,kd->qd', W, V)
```



Attention

Multi-Head: repeat multiple times

1. query-key similarity:

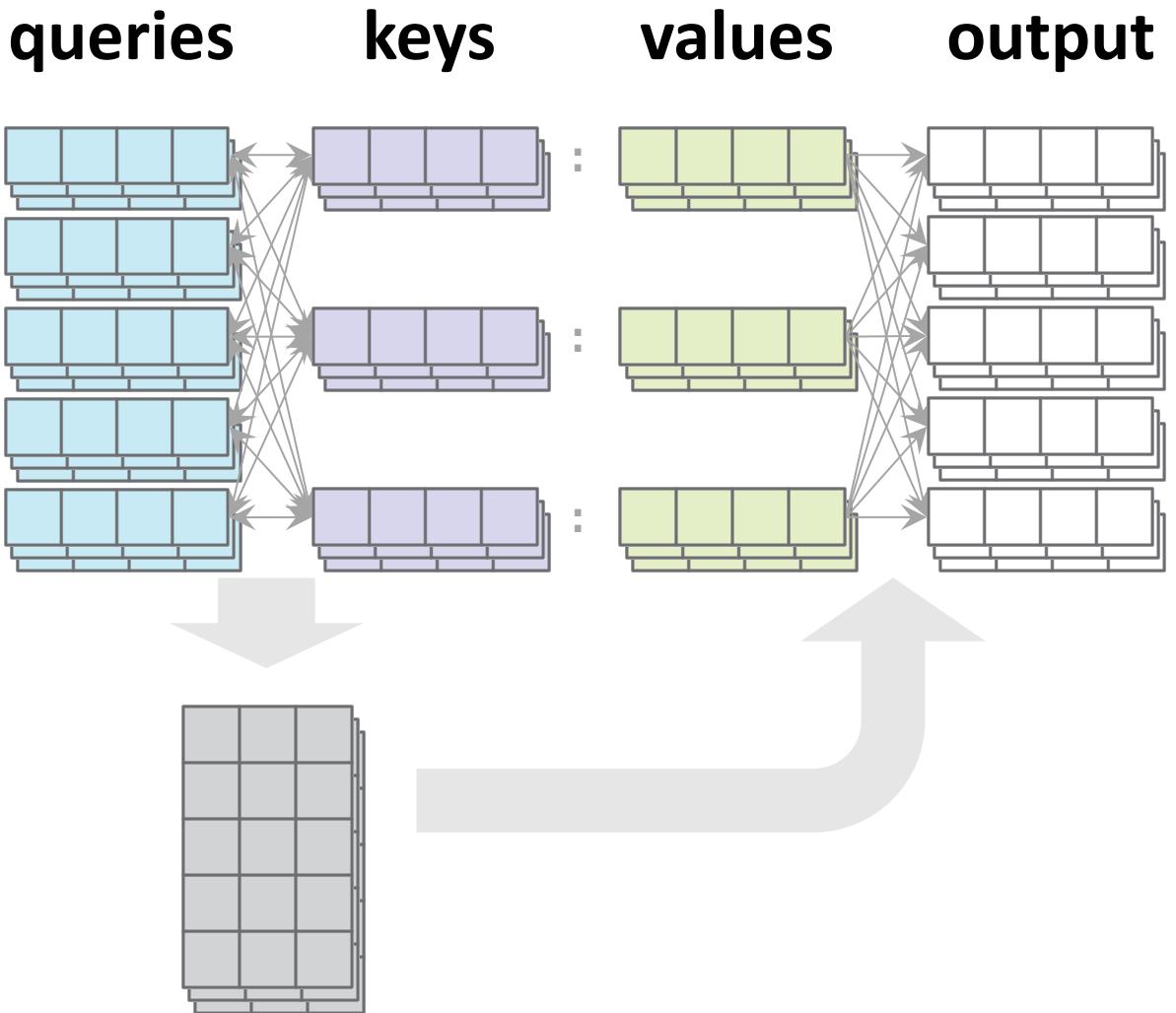
```
S = einsum('hqc,hkc->hqk', Q, K)
```

2. attention weight (softmax):

```
W = softmax(S, dim=-1)
```

3. weighted average:

```
Z = einsum('hqk,hkd->hqd', W, V)
```



Attention

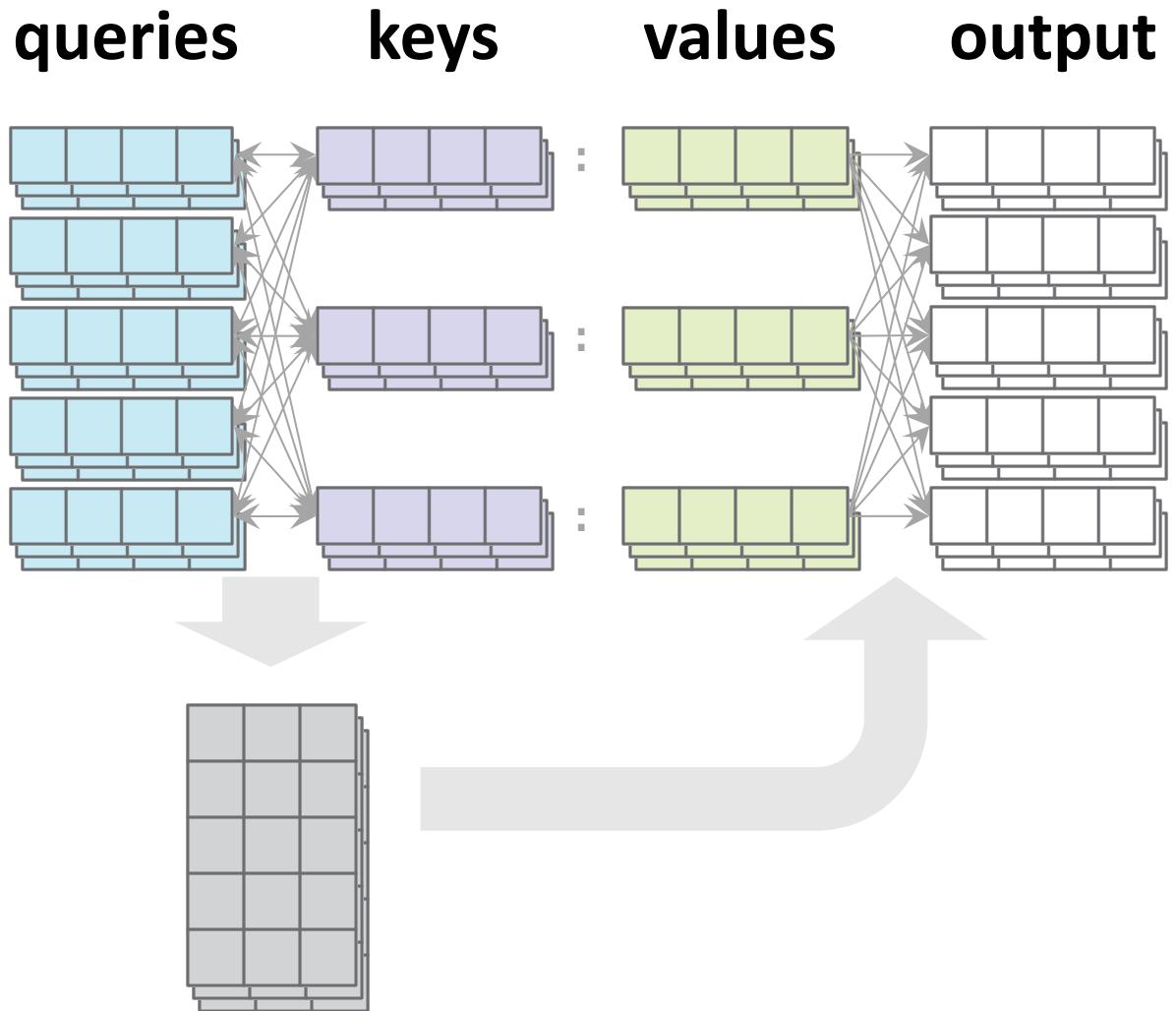
Multi-Head: repeat multiple times

h : # heads

1. query-key similarity:
`S = einsum('hqc,hkc->hqk', Q, K)`

2. attention weight (softmax):
`W = softmax(S, dim=-1)`

3. weighted average:
`Z = einsum('hqk,hkd->hqd', W, V)`



Attention

Multi-Head: repeat multiple times
(w/ batching)

n : batch size

1. query-key similarity:

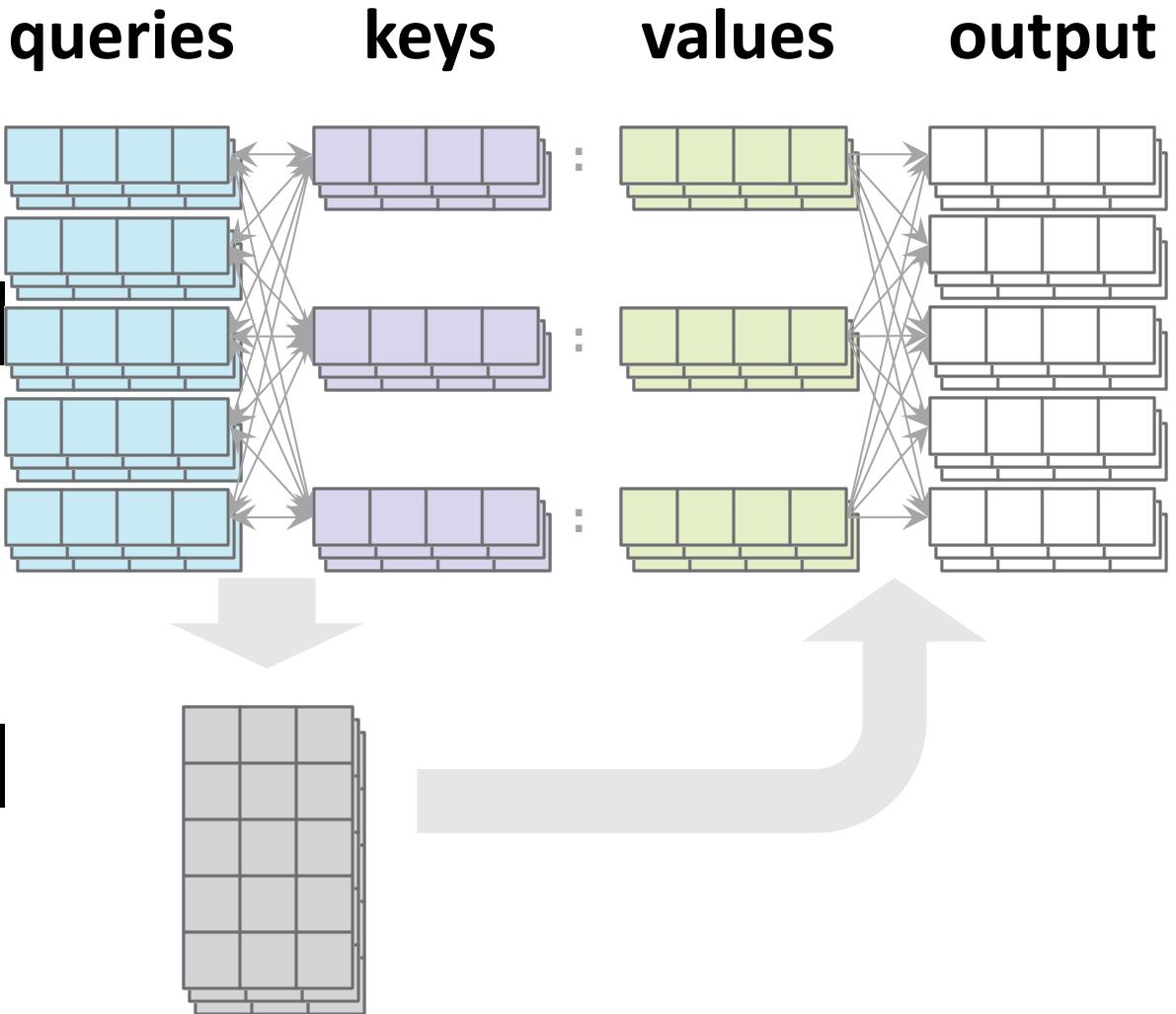
```
S = einsum('nhqc, nhkc->nhqk', Q, K)
```

2. attention weight (softmax):

```
W = softmax(S, dim=-1)
```

3. weighted average:

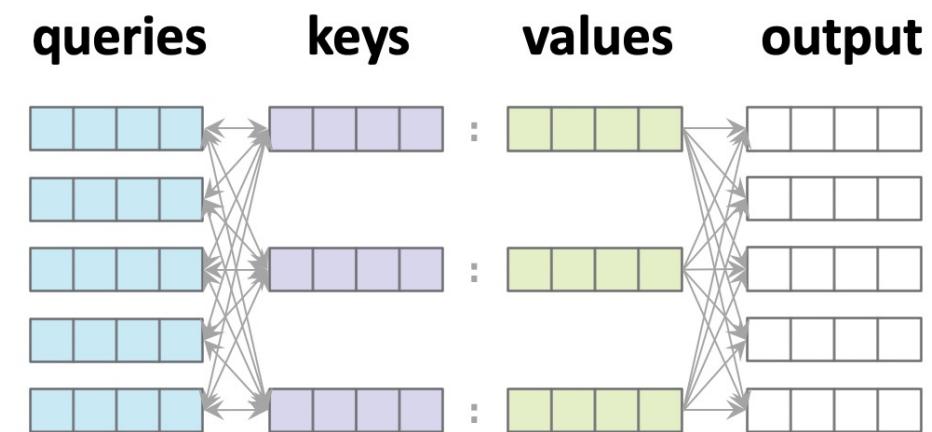
```
Z = einsum('nhqk, nhkd->nhqd', W, V)
```



Attention: Notes

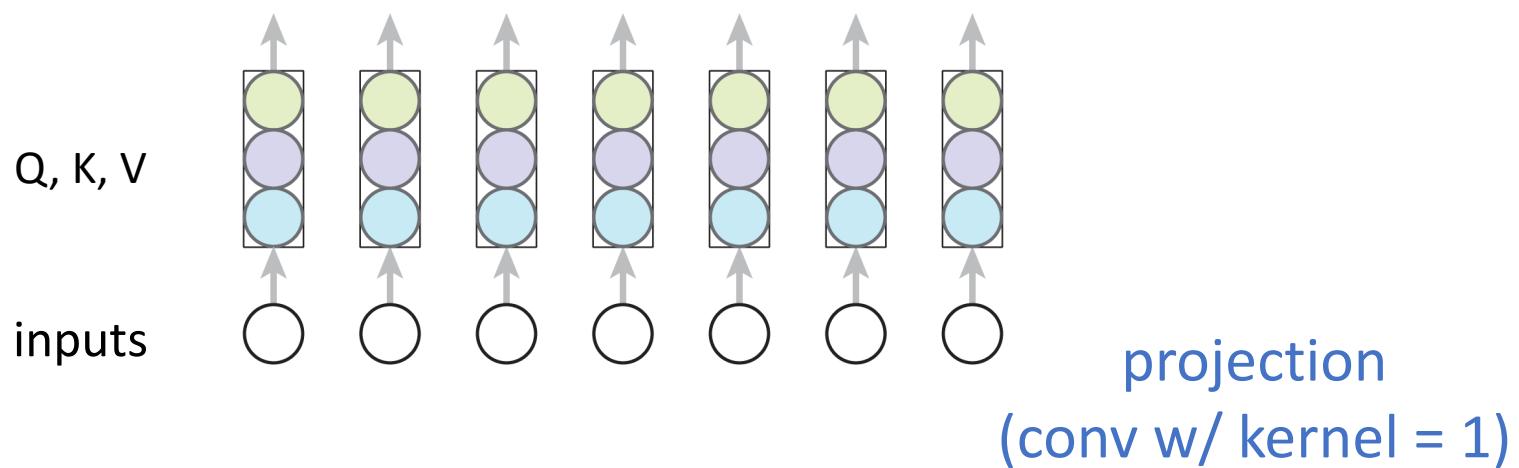
In its general form ...

- **Attention says nothing about how queries, keys, values are provided.**
 - # queries may not equal # keys
 - # queries = # outputs
 - # key channels may not equal # value channels
 - # key channels = # query channels
 - # value channels = # output channels
- **Attention is parameter-free!**
- **Every output is connected to every value and every key.**



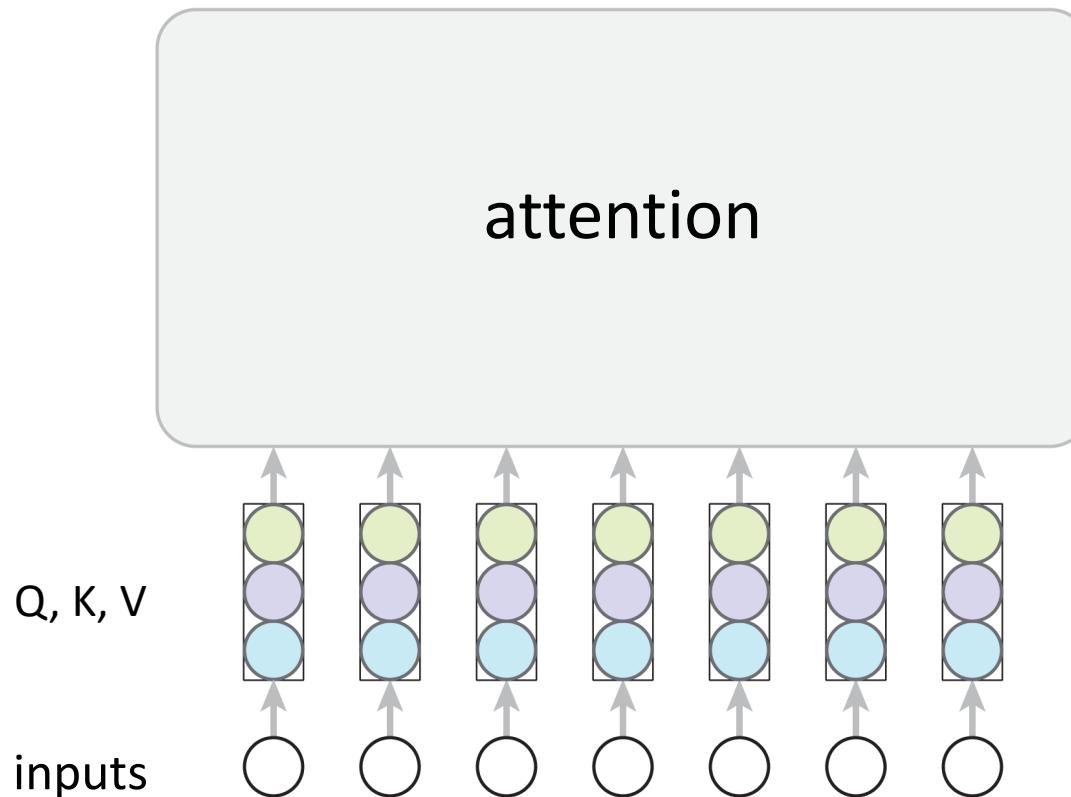
Self-attention

- Q, K, V sequences: produced by the input sequence **itself**.



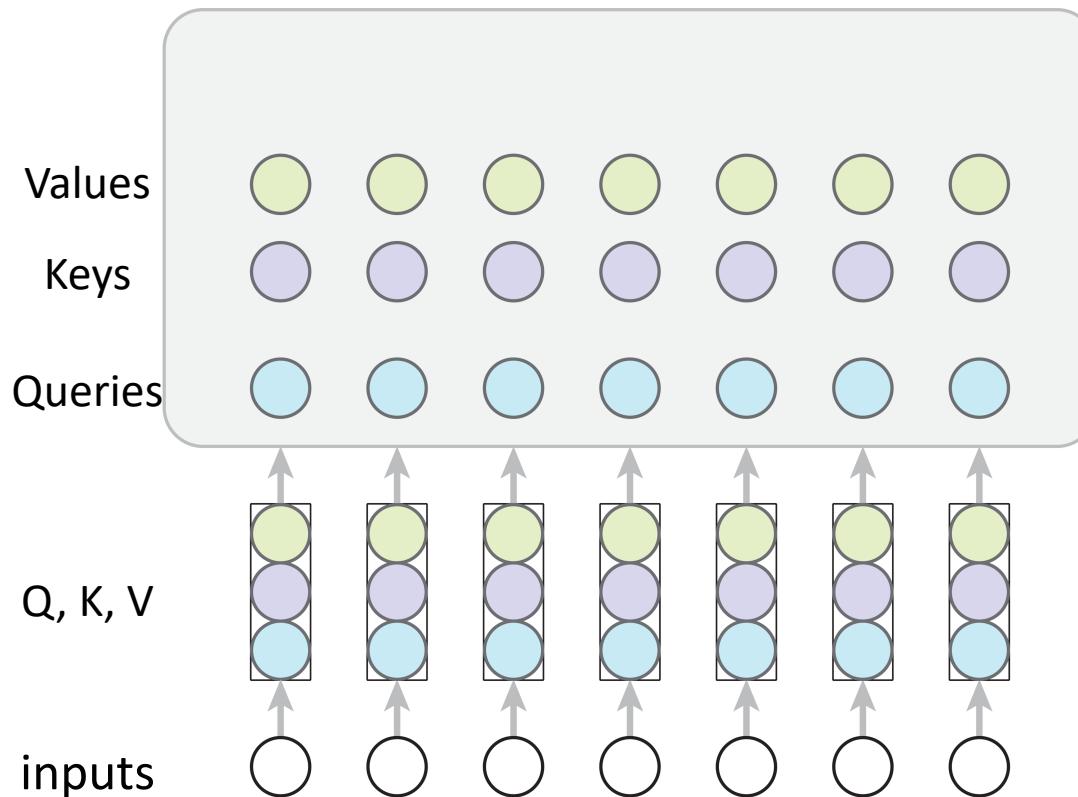
Self-attention

- Q, K, V sequences: produced by the input sequence **itself**.



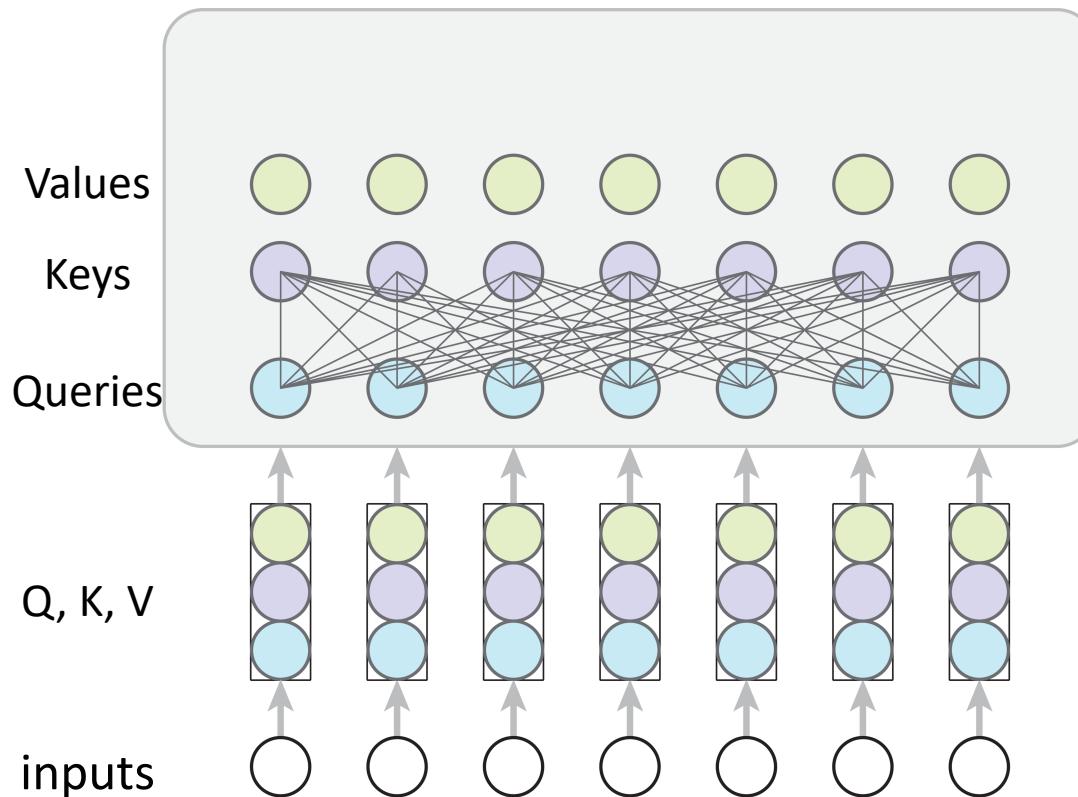
Self-attention

- Q, K, V sequences: produced by the input sequence **itself**.



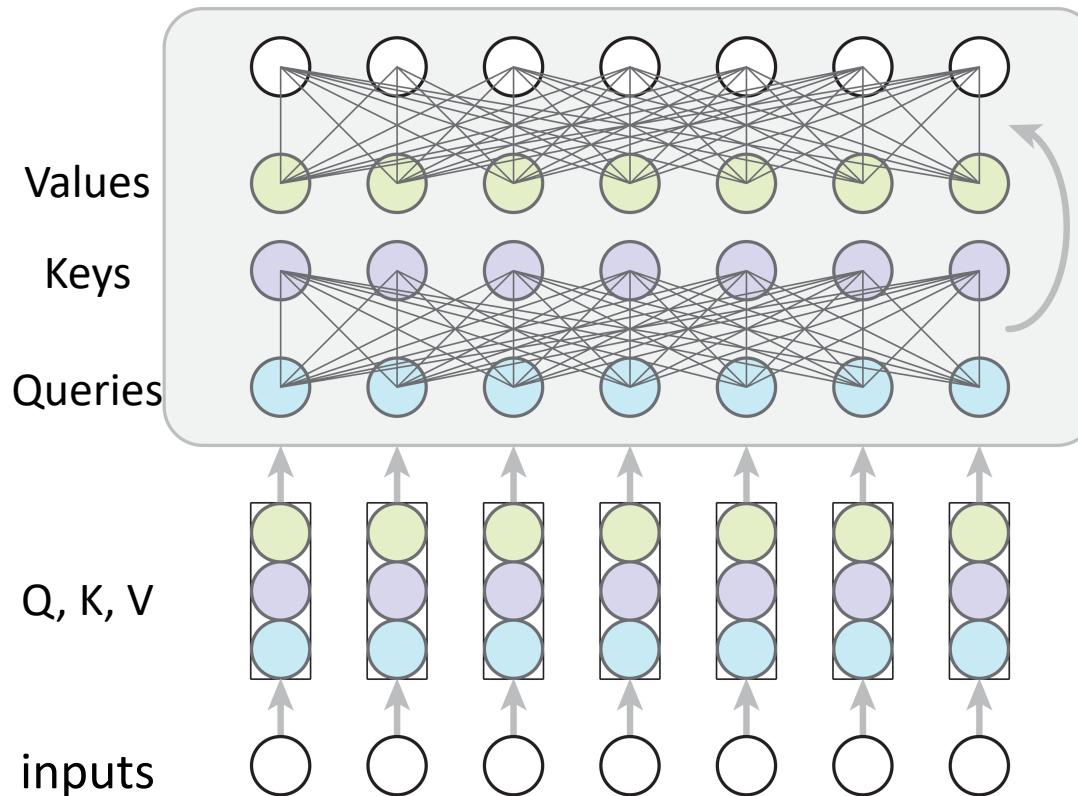
Self-attention

- Q, K, V sequences: produced by the input sequence **itself**.



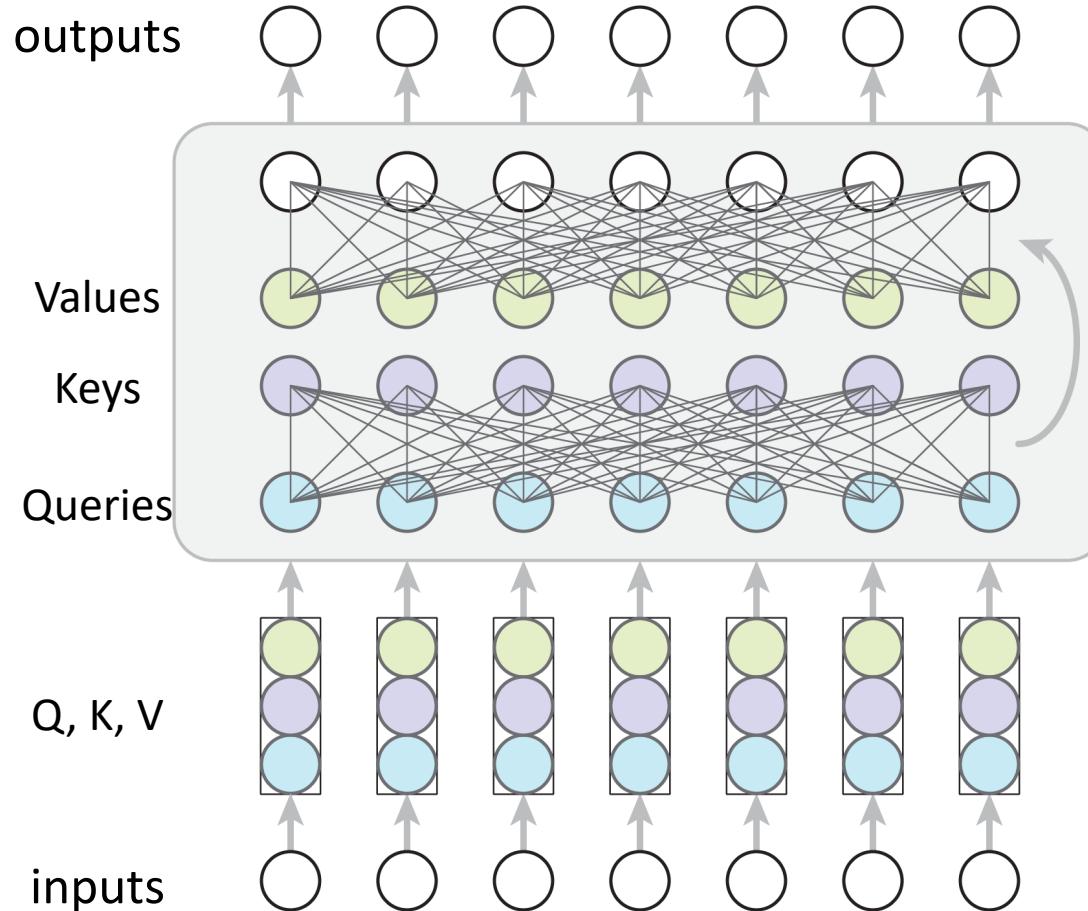
Self-attention

- Q, K, V sequences: produced by the input sequence **itself**.



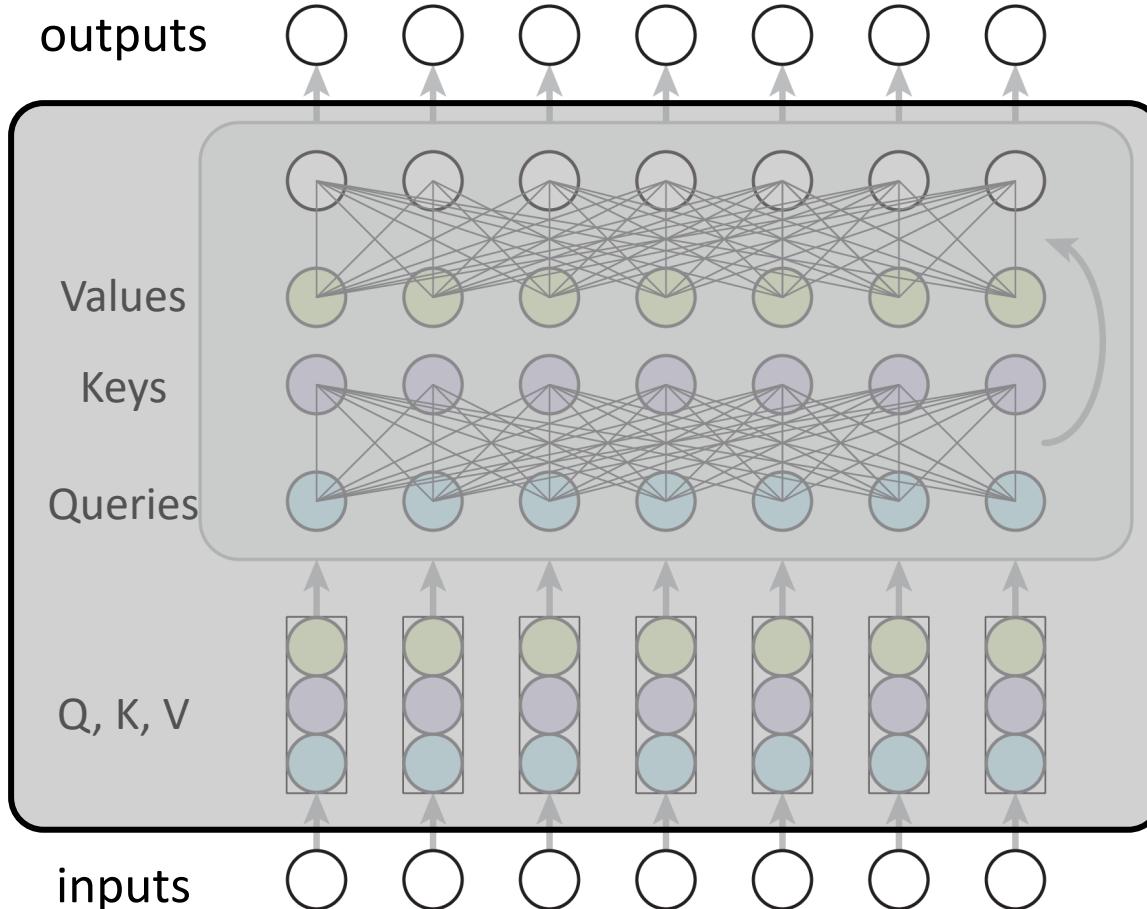
Self-attention

- Q, K, V sequences: produced by the input sequence **itself**.

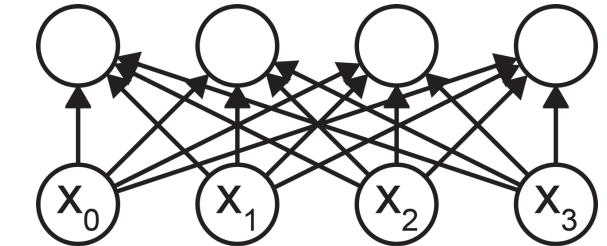


Self-attention

- Q, K, V sequences: produced by the input sequence **itself**.



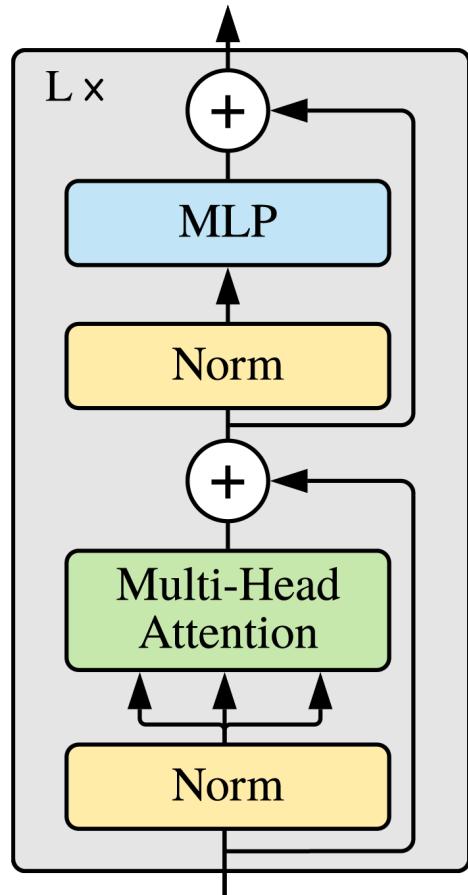
**Every output is connected
to every input!**



(Each arrow indicates
computational dependence)

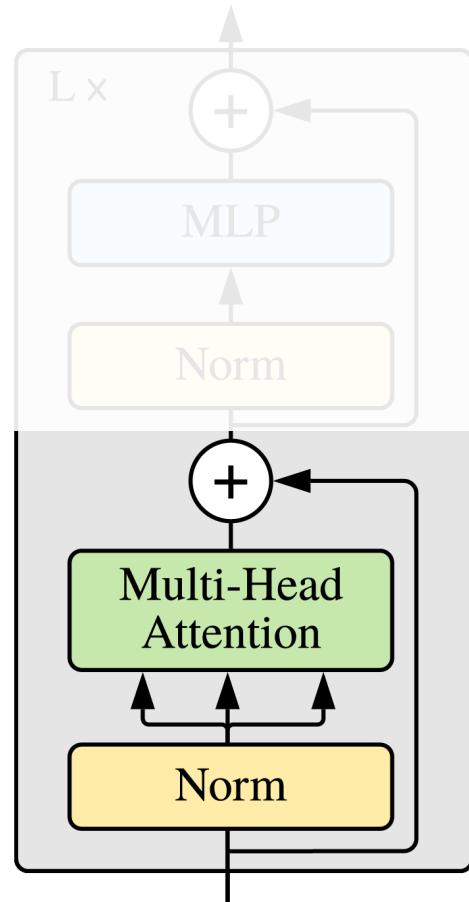
Transformer

- Transformer: a deep neural net based on the attention mechanism.

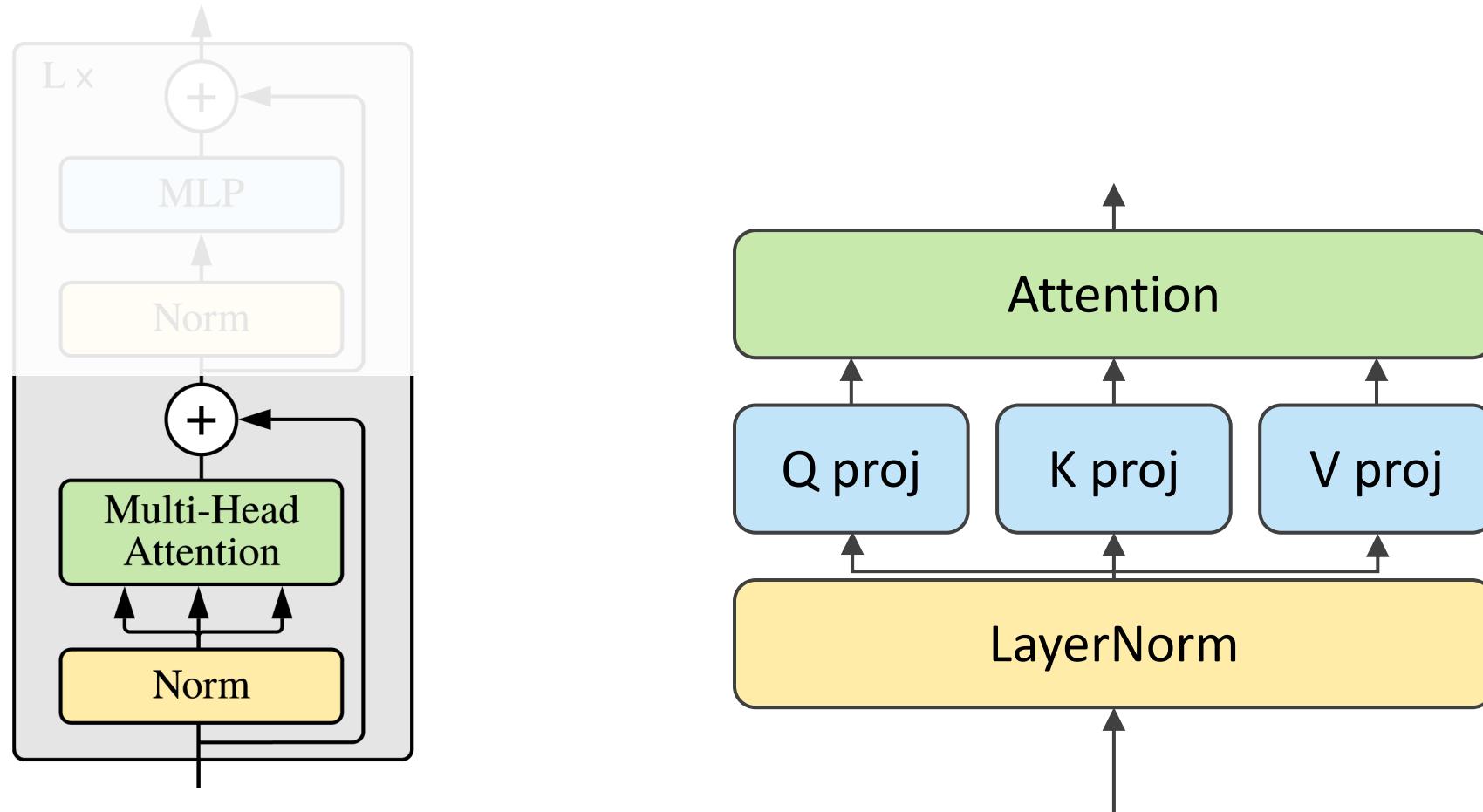


A Transformer block
(pre-norm variant)

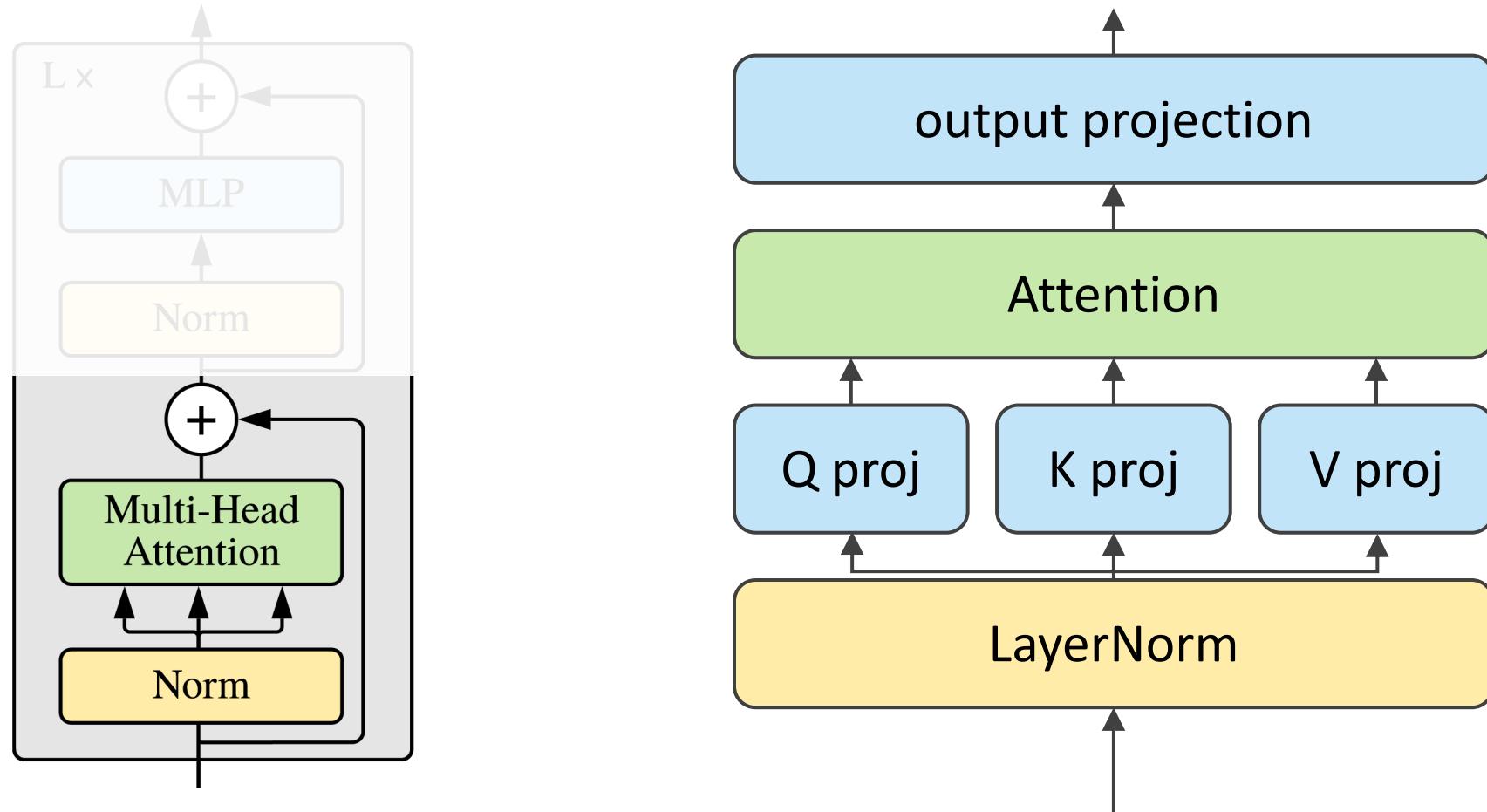
Transformer: Multi-Head Attention (MHA) Block



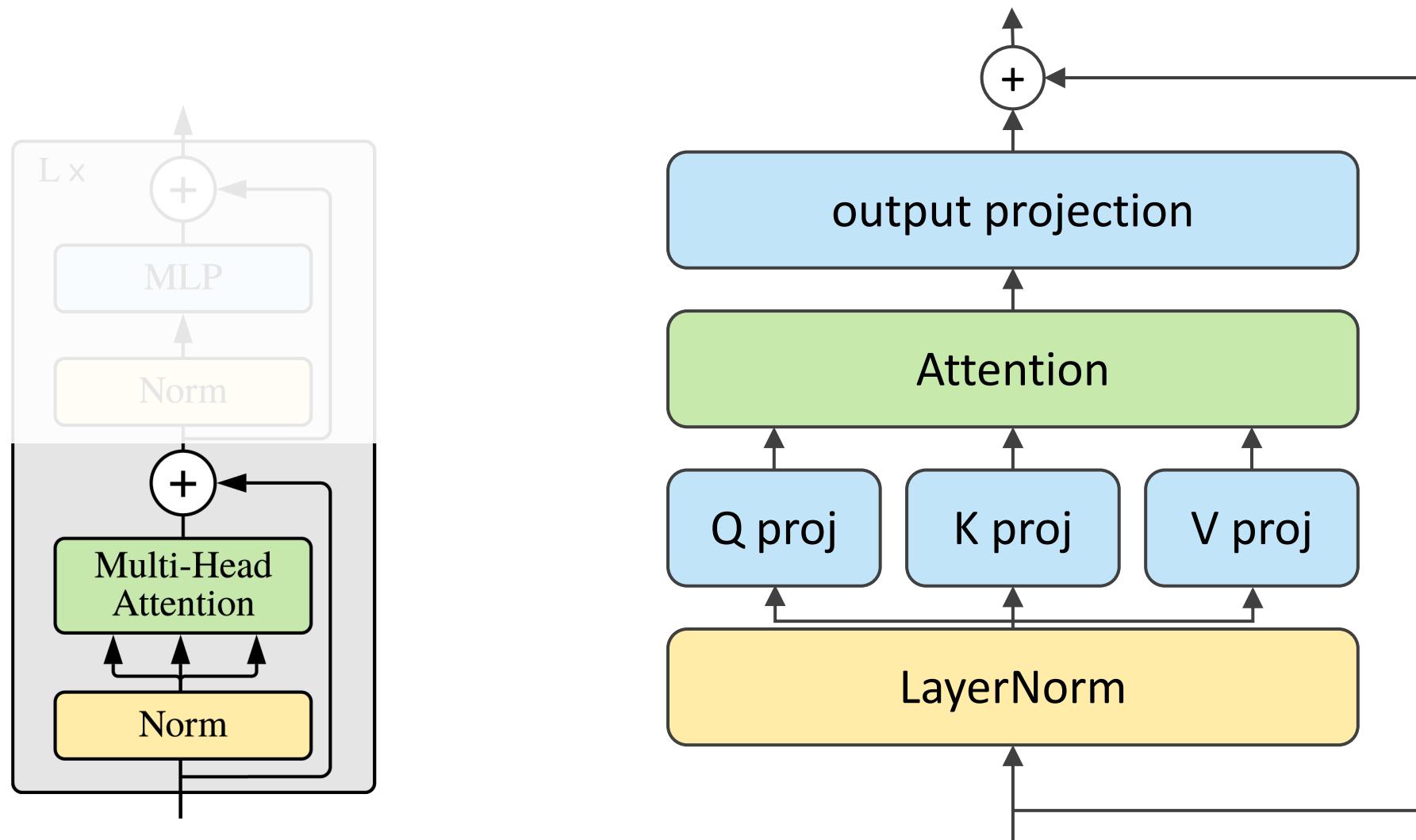
Transformer: Multi-Head Attention (MHA) Block



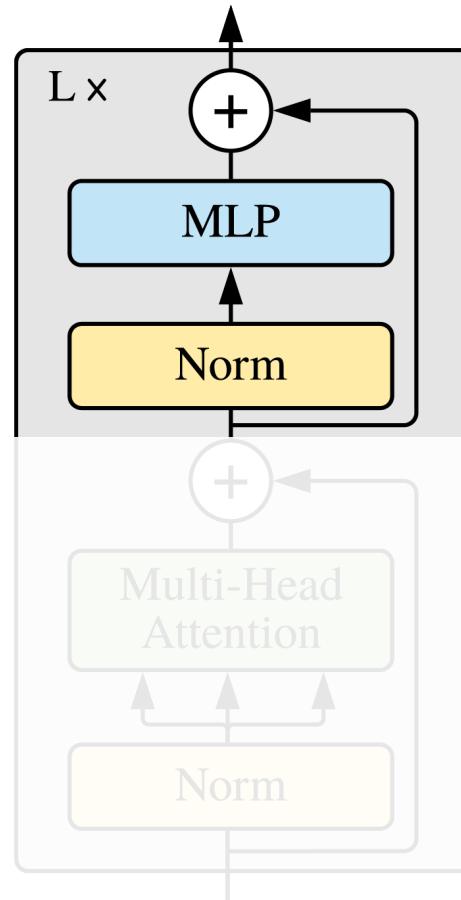
Transformer: Multi-Head Attention (MHA) Block



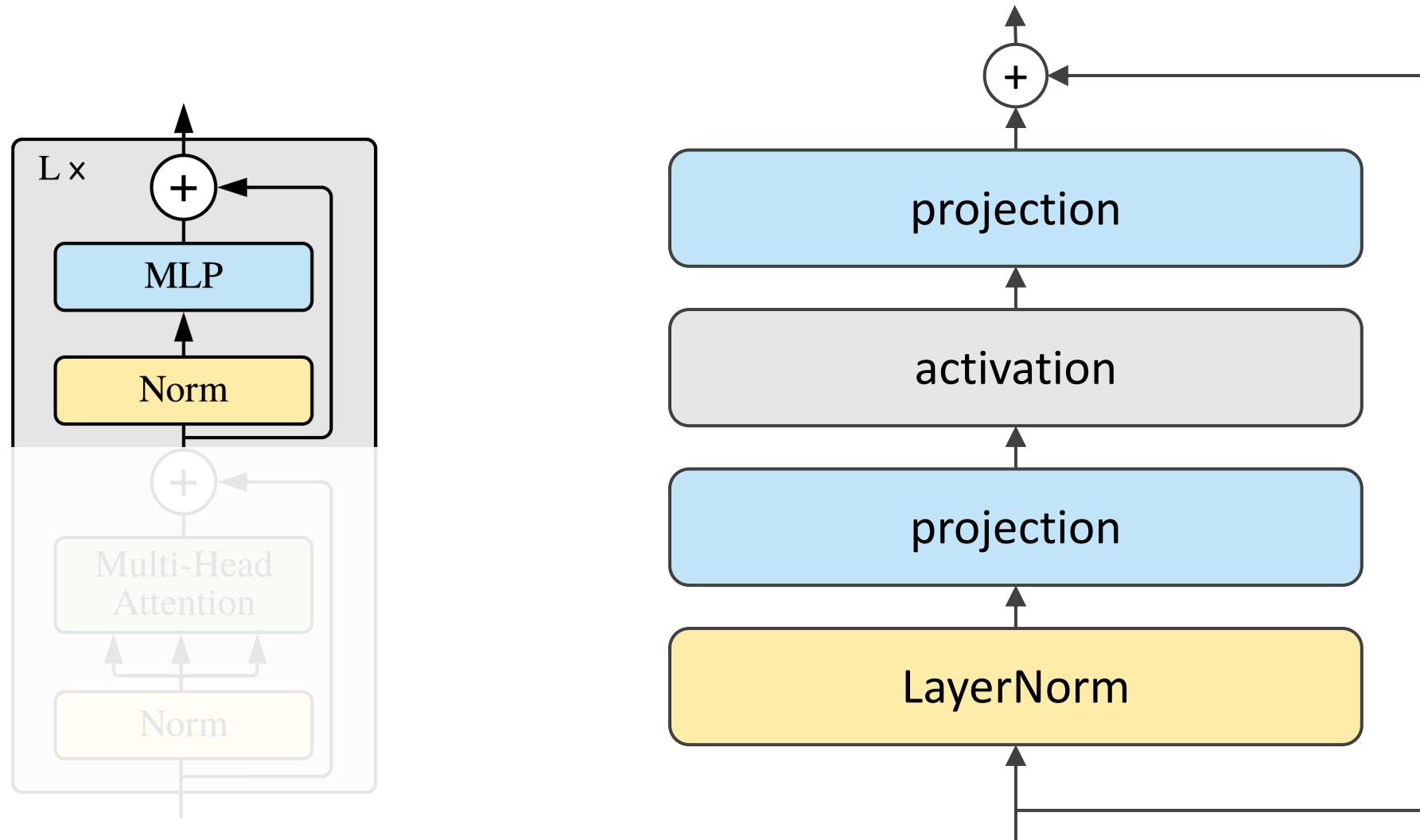
Transformer: Multi-Head Attention (MHA) Block



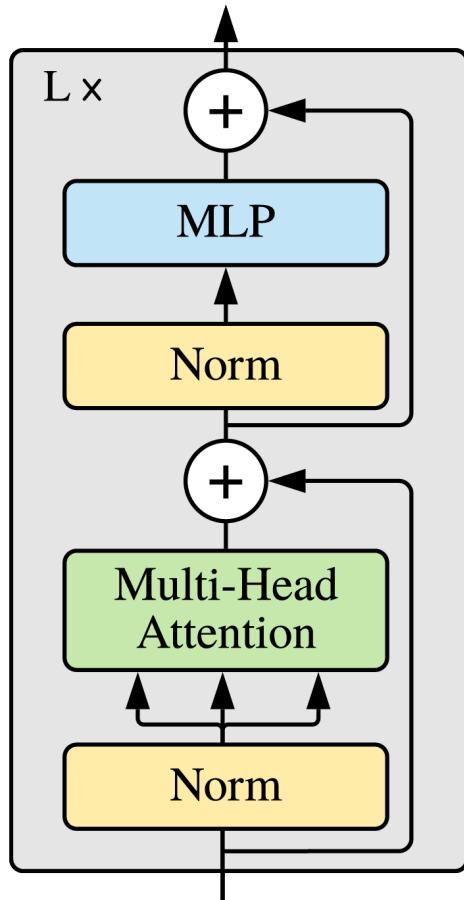
Transformer: Multi-Layer Perception (MLP) Block



Transformer: Multi-Layer Perception (MLP) Block

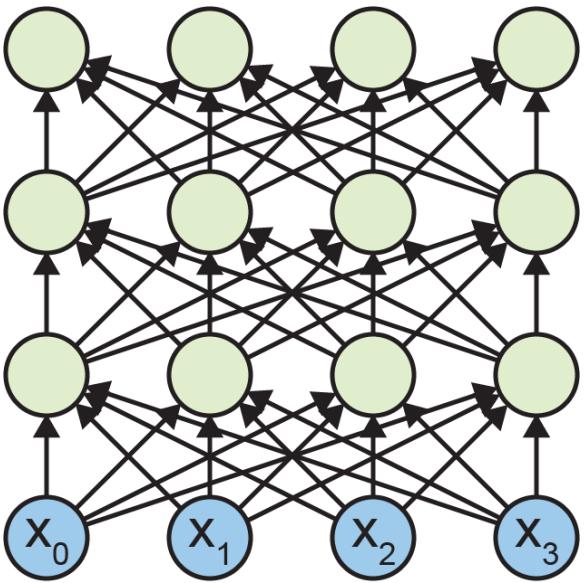


Transformer

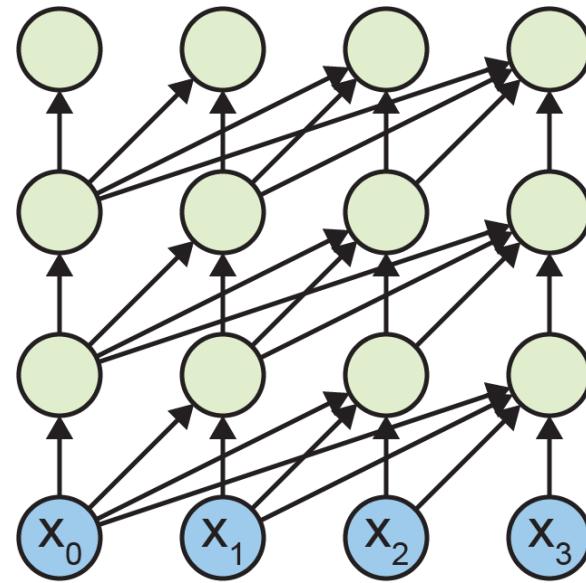


- MHA block
 - QKV **projection**
 - attention
 - output **projection**
- MLP block
 - two **projections**
- All parameters are in **projections**
 - conceptually, all conv with kernel size 1

Transformer: Ordered Sequences



bidirectional
(Transformer “Encoder”)



causal
(Transformer “Decoder”)

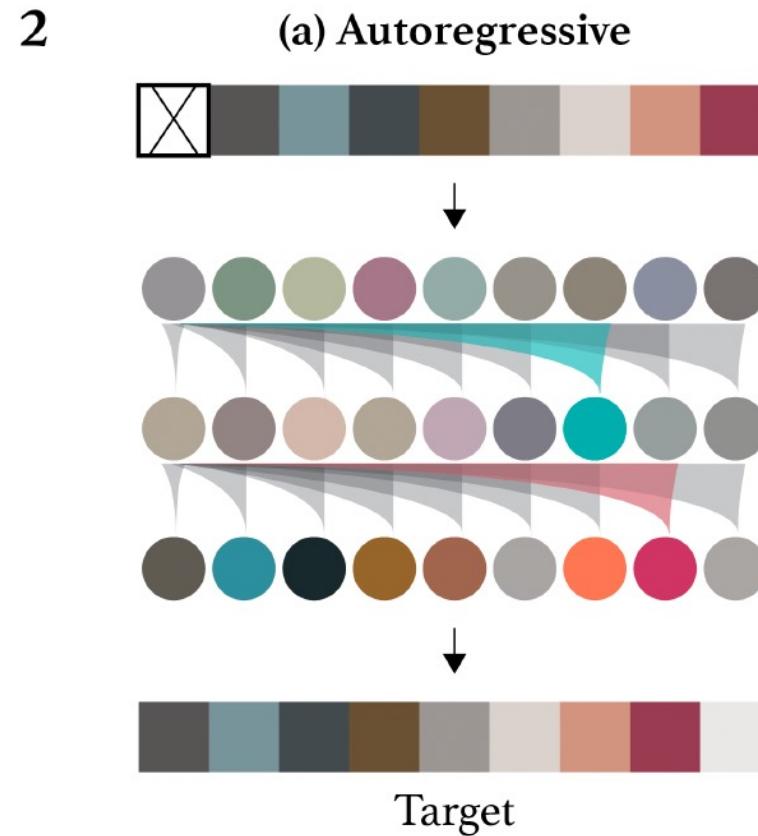
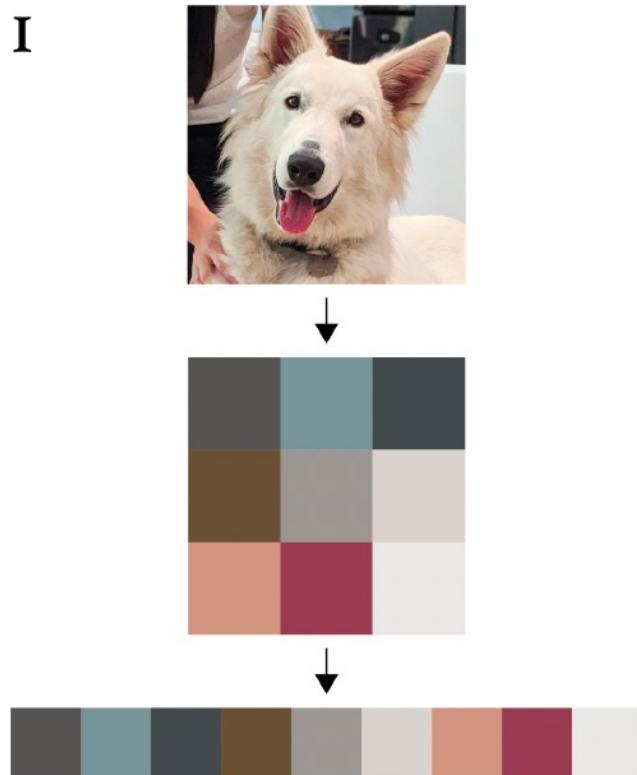
Notes:

- Every “layer” denotes a Transformer block
- Every arrow indicates computational dependence
- Causality is done by masking (before softmax)

Example: image GPT (iGPT)

iGPT: Transformer on sequences of pixels

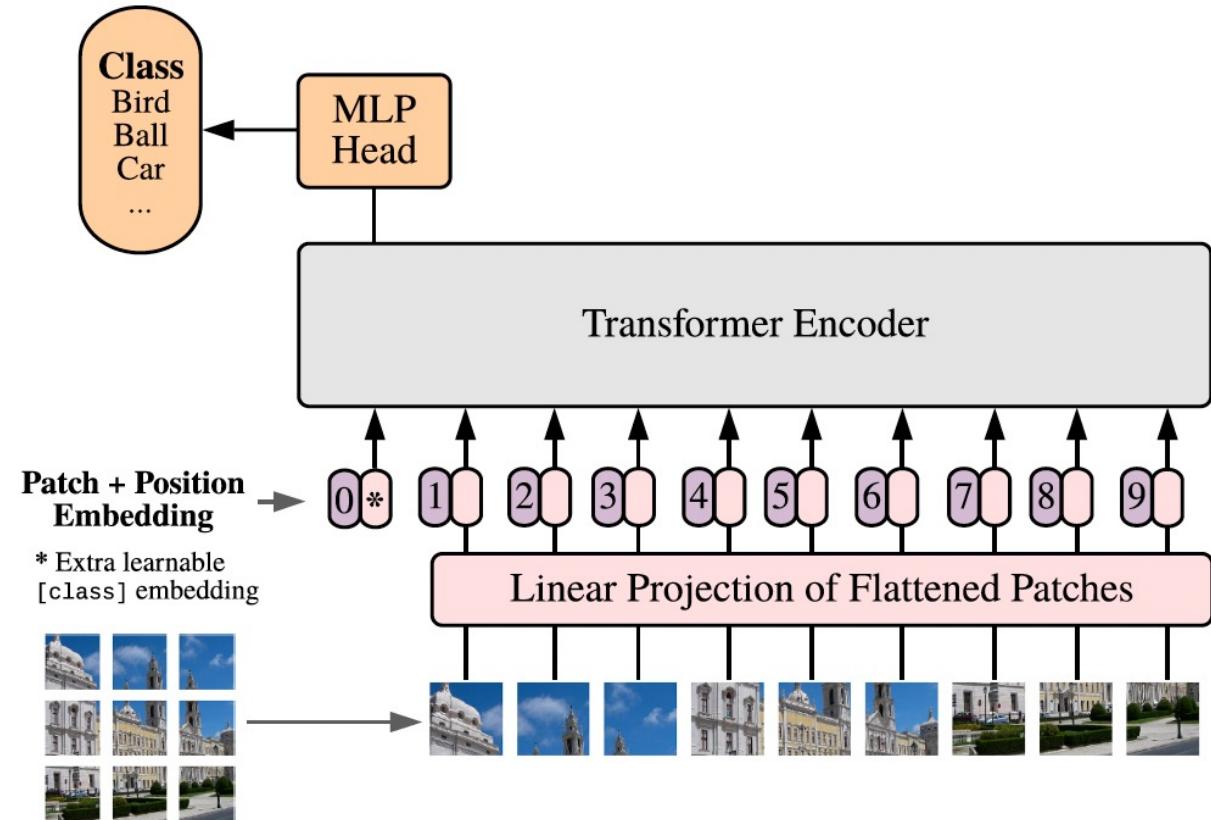
- Autoregressive models (GPT) on pixels for image generation



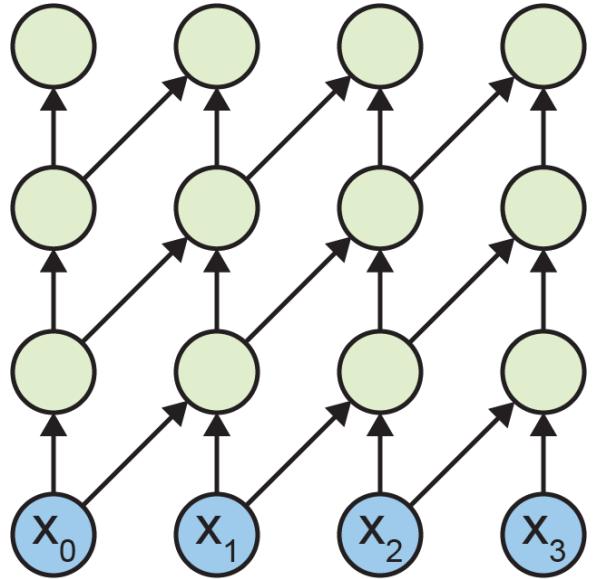
Example: Vision Transformer (ViT)

ViT: Transformer on sequences of patches

- Image: sequence of patches
- Patch embedding
 - 16×16 conv w/ stride 16
- MLP and QKV projection
 - 1×1 conv
- “Transforming” computer vision!

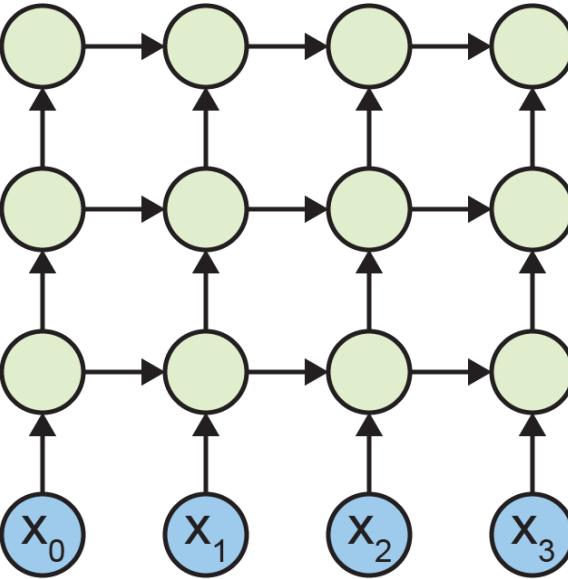


Sequence Modeling



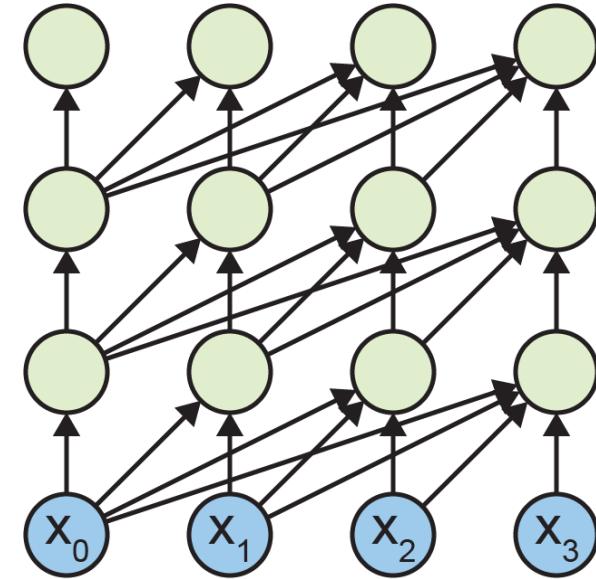
CNN

- limited context
- feedforward



RNN

- long context
- not feedforward



Attention

- long context
- feedforward

Lecture 11: Sequence Modeling

- Sequence Problems in Vision
- ConvNets for Sequence Problems
- Recurrent Neural Networks
- Attention and Transformers