

# TEK 5040/9040

## Memory and Attention in RNNs

Narada Warakagoda

background

Attention

Attn. Examples

MANN

Self-attention

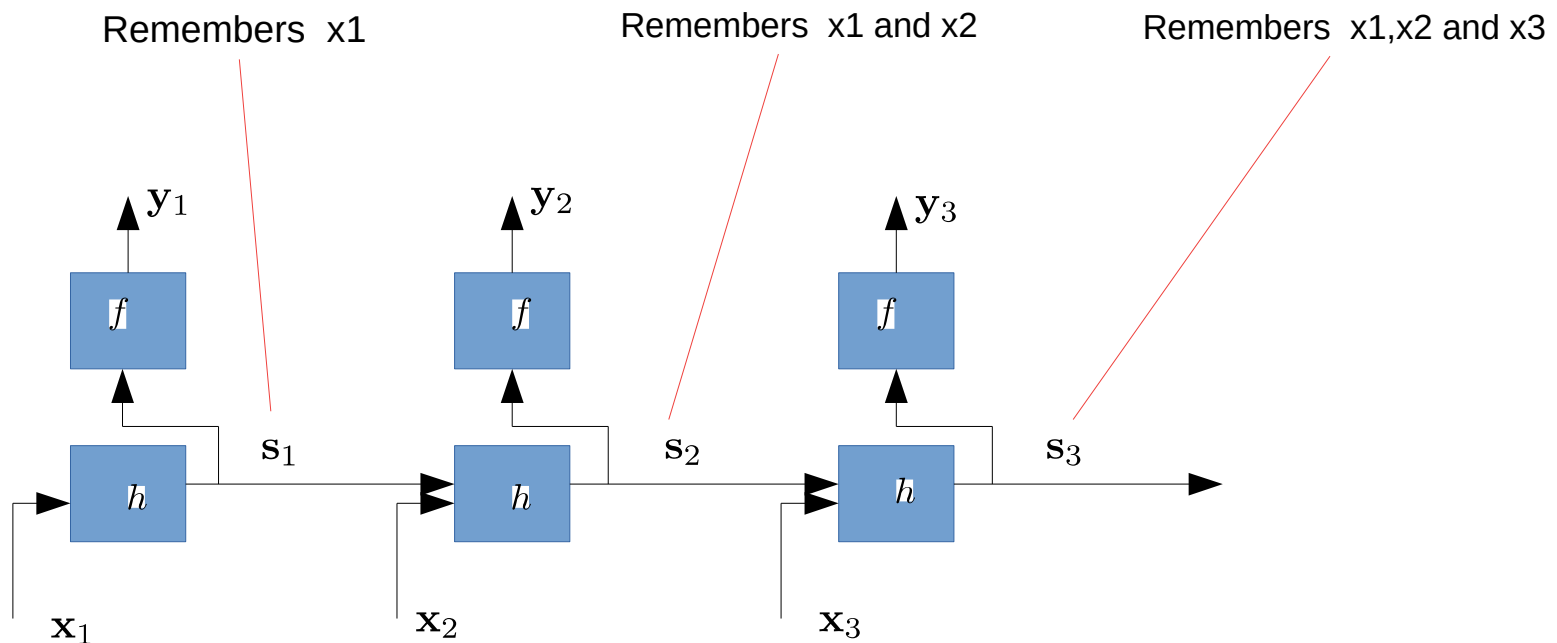
1 / 28

# Sequence processing and “memory”

- Processing inputs such as images does not depend on memory
  - All inputs (pixels) are fed at the same time
- However, sequence processing inherently needs “memory”
  - Need to remember the previous inputs
- Plain RNN provides a simple memory solution!

# The memory of plain RNN

- Cell states are a kind of memory of the previous inputs
- But these memories are highly restricted!



background

Attention

Attn. Examples

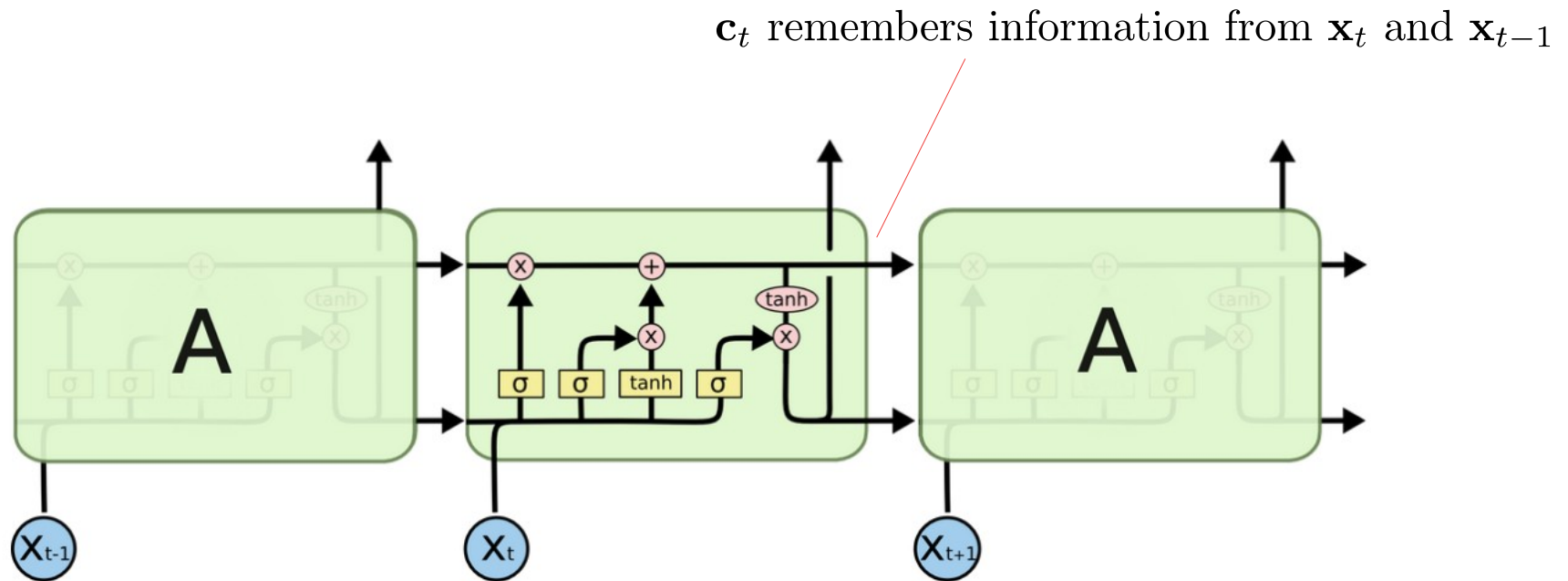
MANN

Self-attention

3 / 28

# Memory of the LSTM

- Control state has a better memory
  - Forgets irrelevant information
  - Remembers important information



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

background

Attention

Attn. Examples

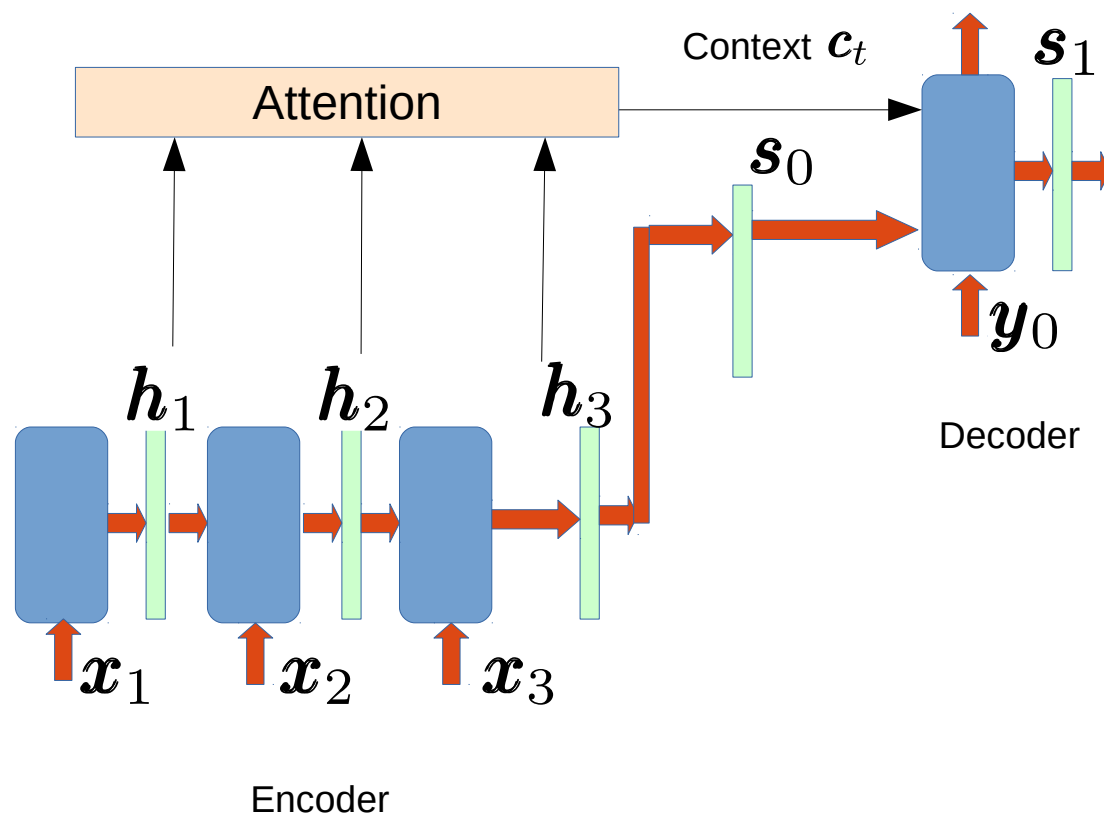
MANN

Self-attention

4 / 28

# Attention

- Direct access to all the previous states
- But not all previous states are not important
  - Pick the states which are important



background

Attention

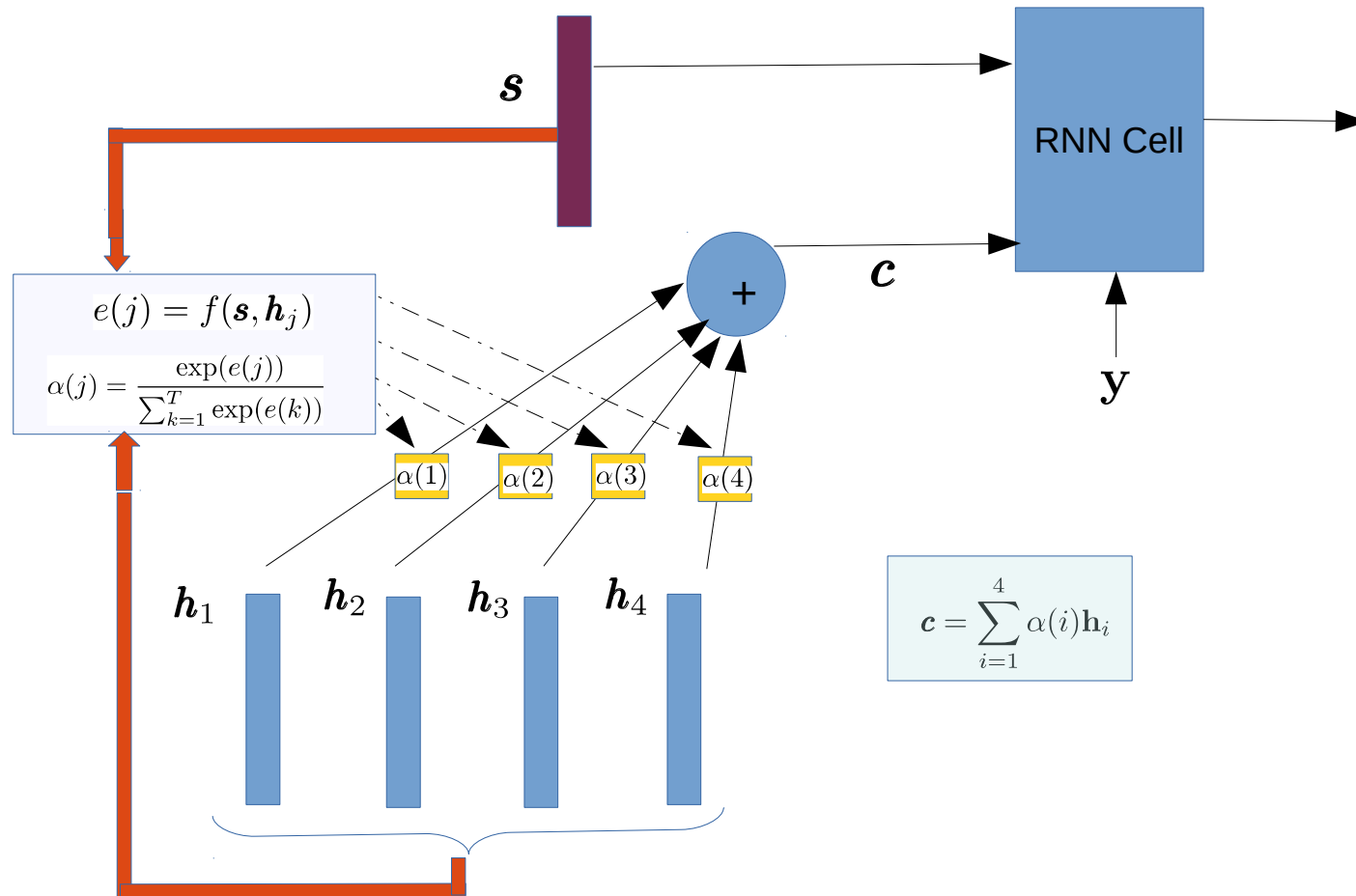
Attn. Examples

MANN

Self-attention

5 / 28

# Soft attention



background

Attention

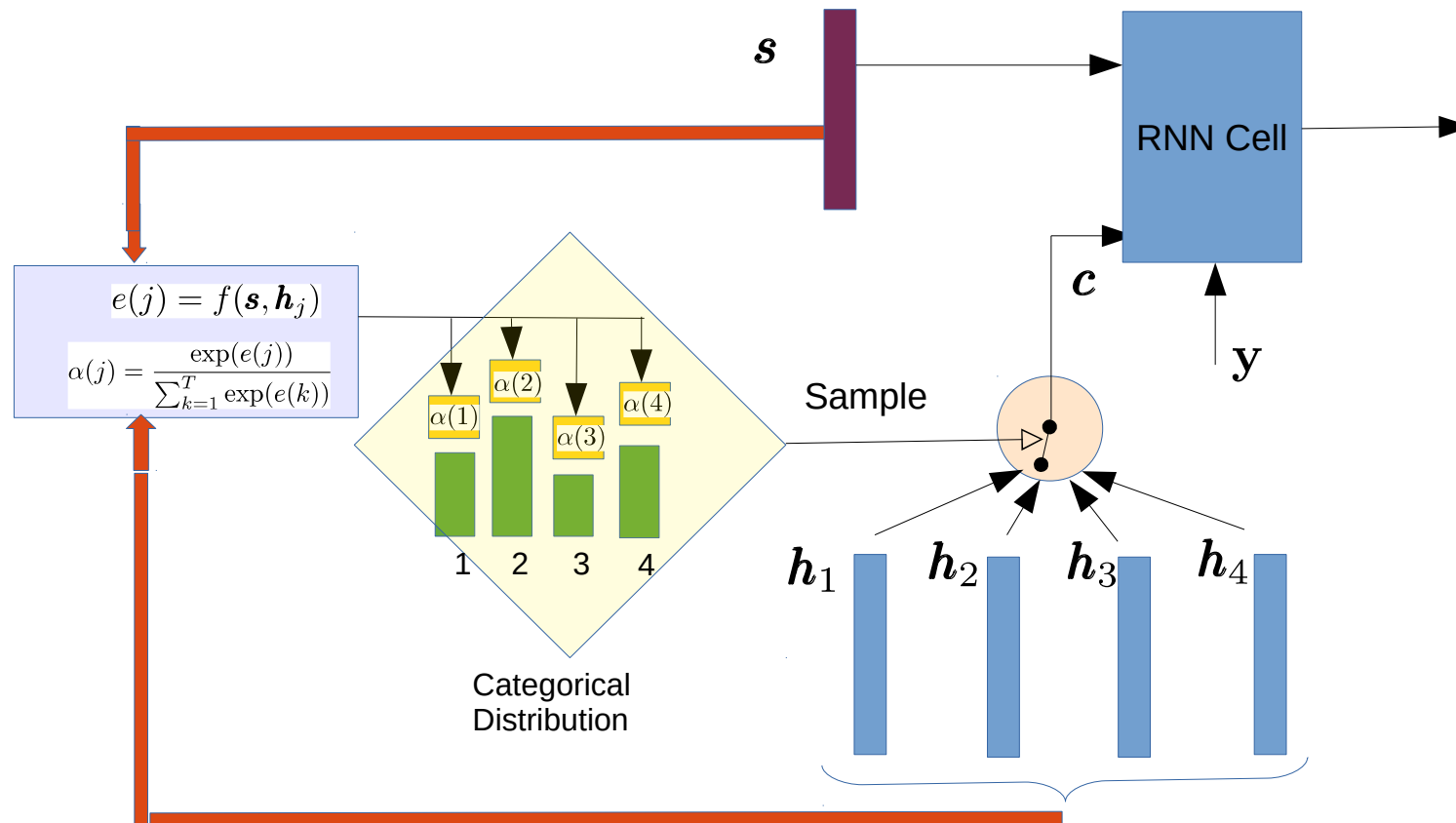
Attn. Examples

MANN

Self-attention

6 / 28

# Hard Attention



background

Attention

Attn. Examples

MANN

Self-attention

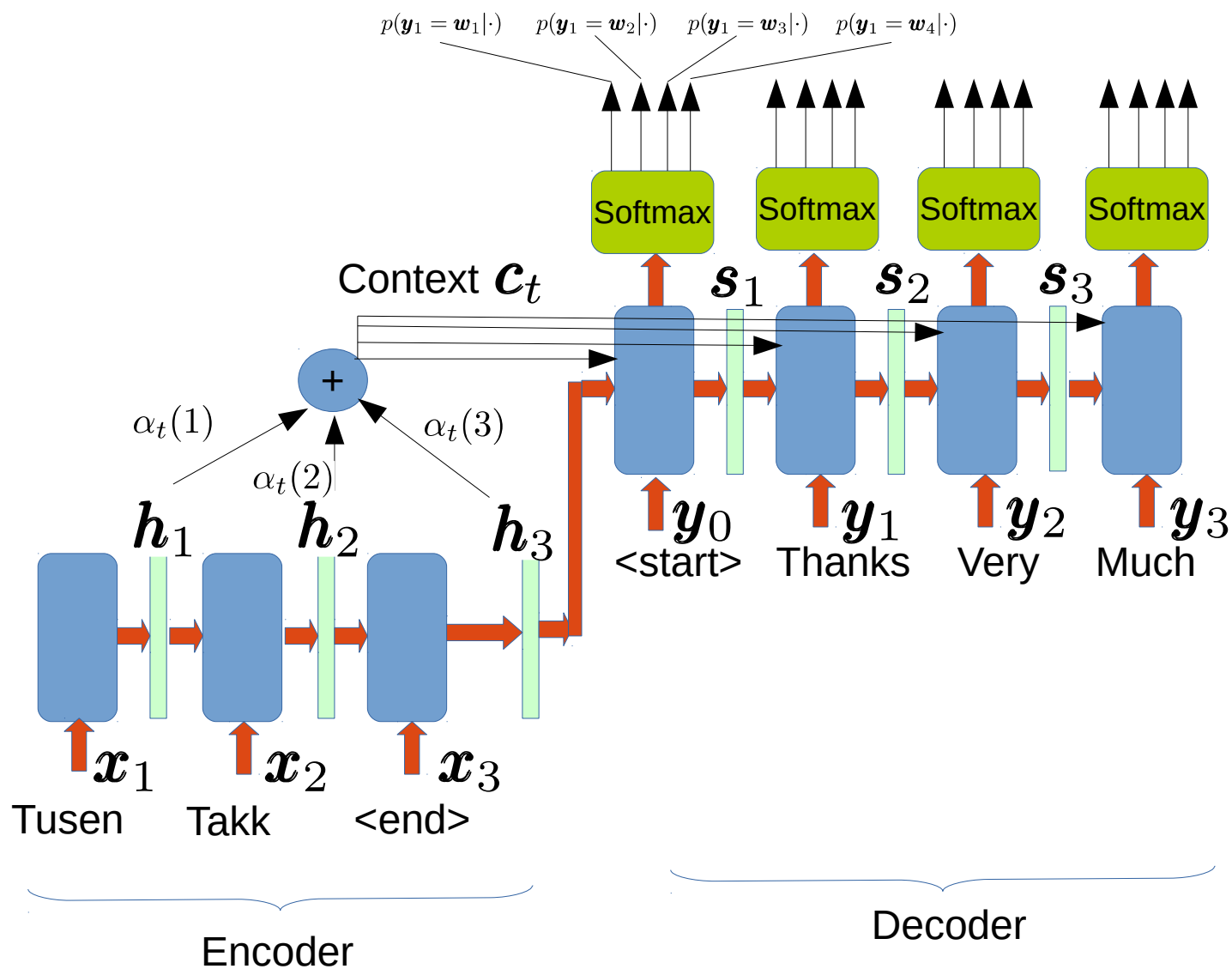
7 / 28

# Examples of attention

- Machine translation
  - Sequence to sequence
- Image captioning
  - Vector to sequence



# Machine translation



background

Attention

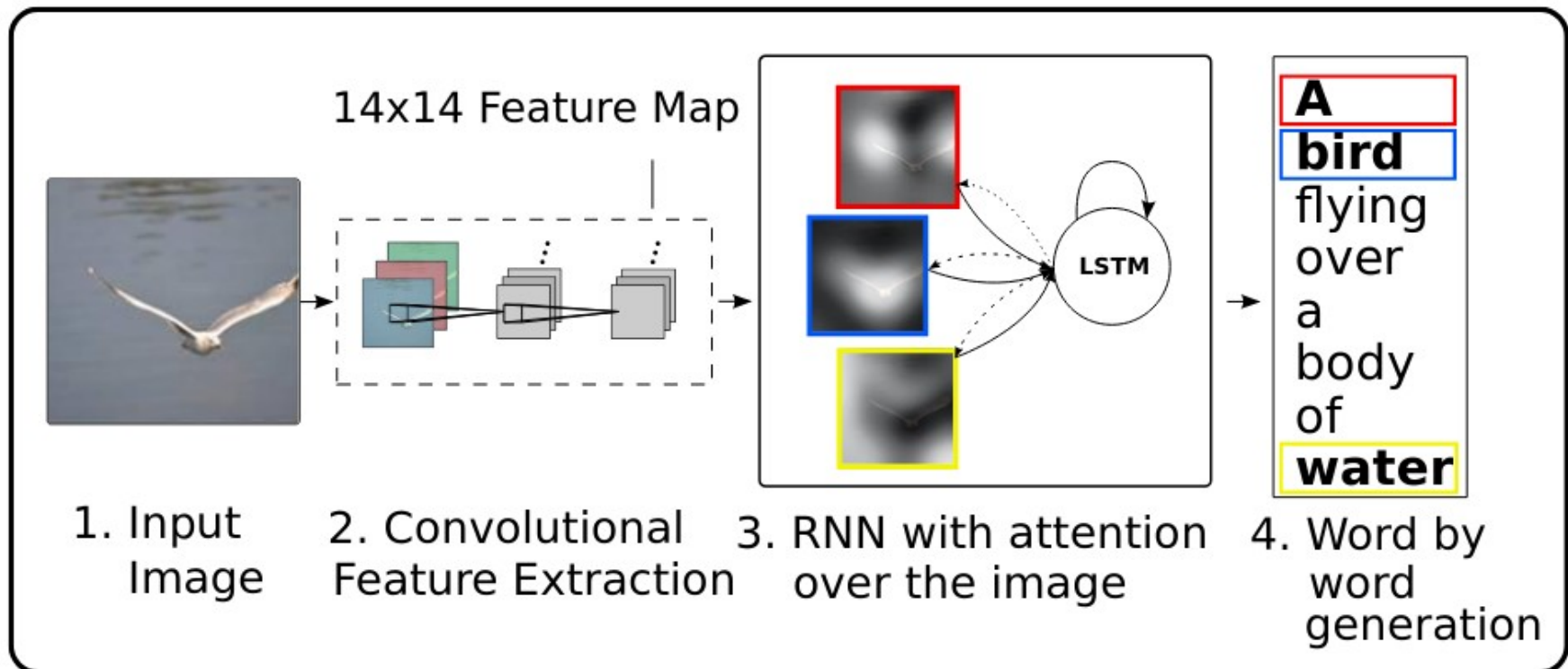
Attn. Examples

MANN

Self-attention

9 / 28

# Image captioning



background

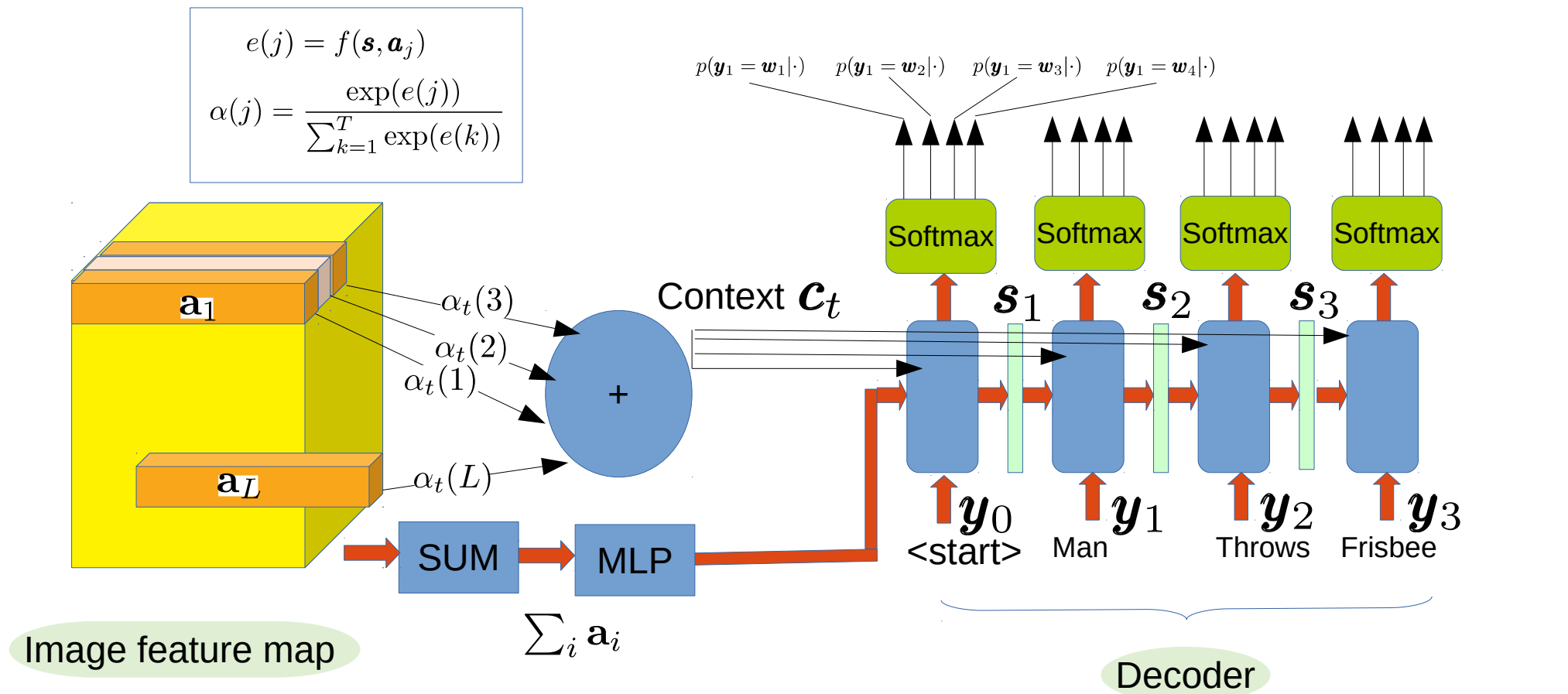
Attention

Attn. Examples

MANN

Self-attention

# Attend to locations



background

Attention

Attn. Examples

MANN

Self-attention

11 / 28

# Attention visualization



background

Attention

Attn. Examples

MANN

Self-attention

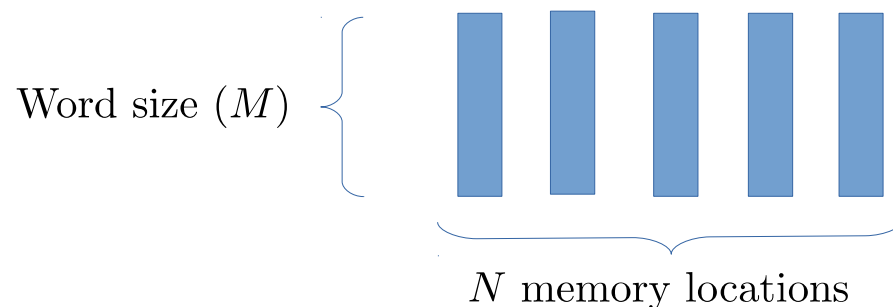
12 / 28

# Pros & cons of attention

- Pros:
  - Direct access to previous states (memory of previous inputs)
- Cons:
  - Memory = states (Not a general layout)
  - Writing to memory = generating states one by one

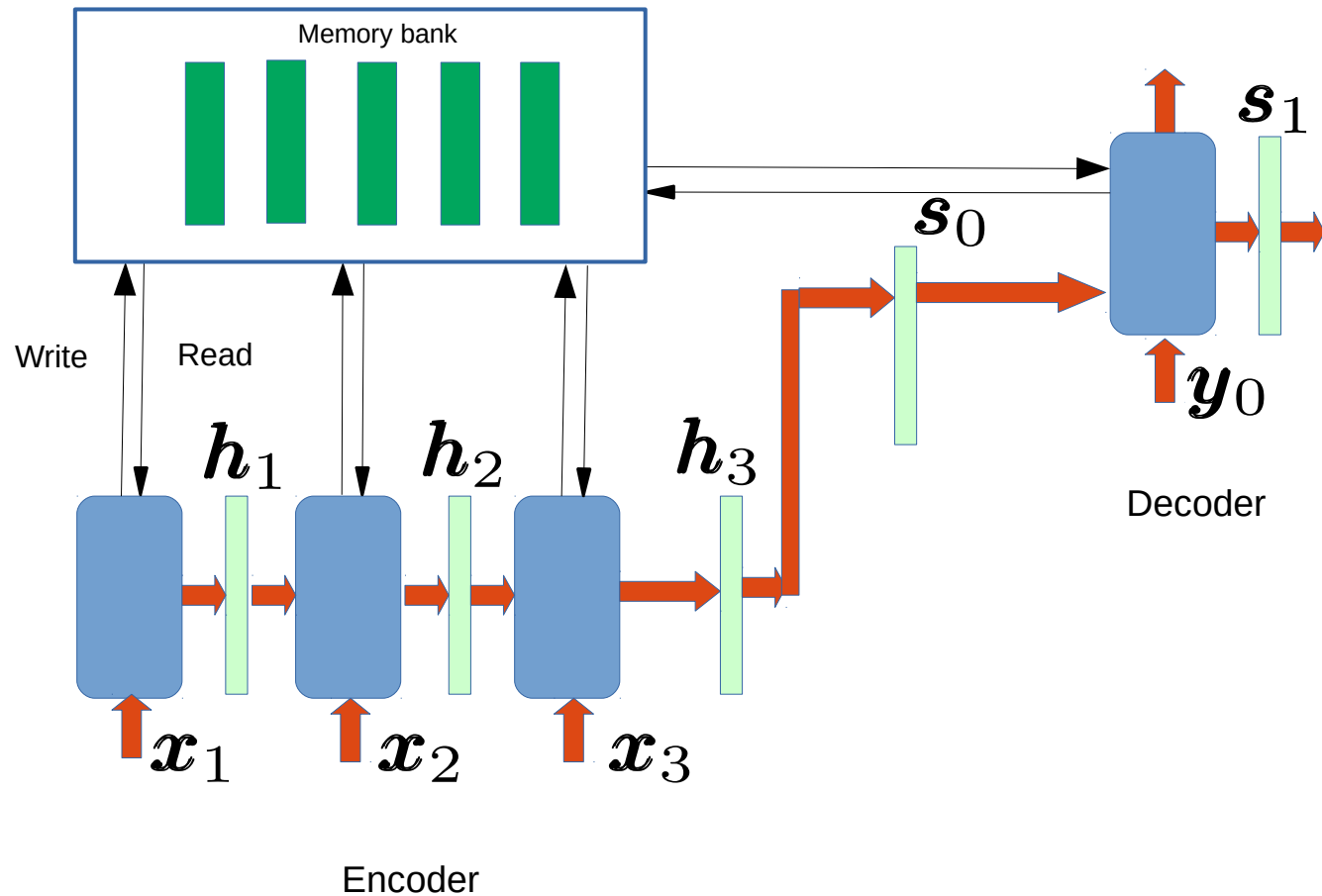
# Memory Augmented Neural Network (MANN)

- Provides a general memory layout
  - (independent of the RNN architecture and time-steps)



- Provides a more general reading/writing mechanism
- Choice in reading/writing memory
  - Content based addressing
  - Location based addressing

# MANN architecture



background

Attention

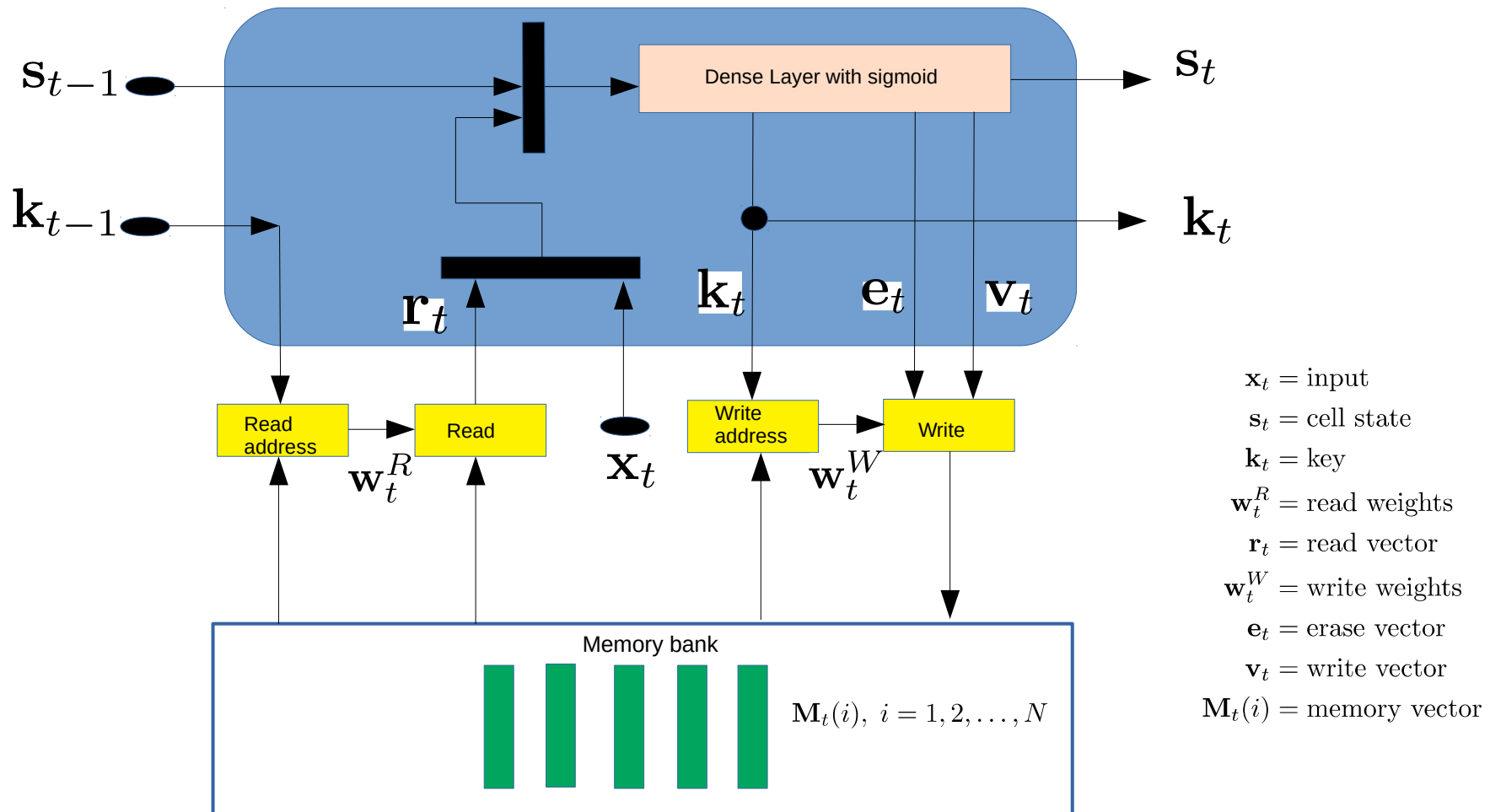
Attn. Examples

MANN

Self-attention

15 / 28

# Cell architecture



background

Attention

Attn. Examples

**MANN**

Self-attention



# Reading Memory

- Generate read weights

$$\alpha_j = \mathbf{k} \cdot \mathbf{M}(j) \quad (\text{Similarity measure})$$

$$w^R(j) = \text{softmax}(\alpha_j) = \frac{\exp(\alpha_j)}{\sum_k \exp(\alpha_k)}$$

- Generate the read vector

$$\mathbf{r} = \sum_j w_j^R \mathbf{M}(j)$$

- “Soft” operation and hence differentiable
- Interpreted as “attention” to memory

background

Attention

Attn. Examples

MANN

Self-attention

17 / 28

# Writing to memory

- Generate write weights

$$\alpha_j = \mathbf{k} \cdot \mathbf{M}^t(j) \quad (\text{Similarity measure})$$

$$w^W(j) = \text{softmax}(\alpha_j) = \frac{\exp(\alpha_j)}{\sum_k \exp(\alpha_k)}$$

- Perform the write operation (update memory)

$$\mathbf{M}^{t+1}(j) = \mathbf{M}^t(j) \circ (\mathbf{1} - w^W(j)\mathbf{e}) + w^W(j)\mathbf{v}$$

Elementwise multiplication      Erase vector      Write vector

# Addressing

- Content based addressing (previous slides)
  - Find the similarity between memory contents and the key
  - Locations with high similarity
    - contributes more to the read vector (in reading)
    - are written with “more” information (in writing)
  - Locations with low similarity
    - Do the opposite
- Location based addressing
  - Read from (write to) a specified location
  - Common in regular computer systems
  - Differentiable location based addressing is useful in some cases.

background

Attention

Attn. Examples

MANN

Self-attention

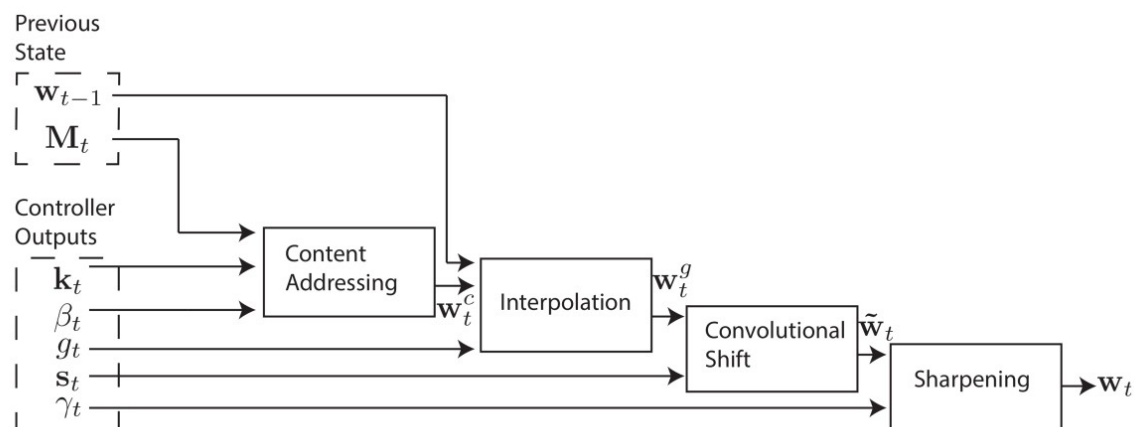
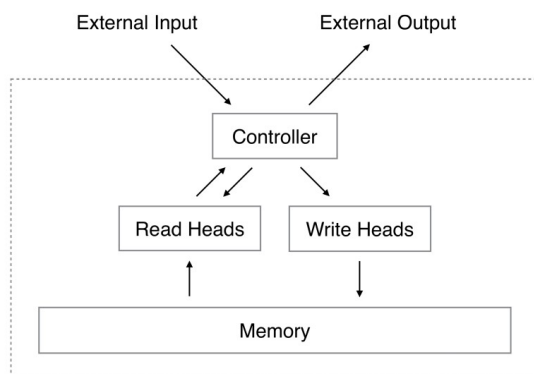
19 / 28

# MANN examples

- Neural Turing Machine
- Differentiable Neural Computer

# Neural Turing Machine

- Sophisticated weight generation mechanism
- Supports
  - Content based addressing
  - Location based addressing ( through address shifting distribution  $s_t$  )
  - Both addressing mechanisms are blended



background

Attention

Attn. Examples

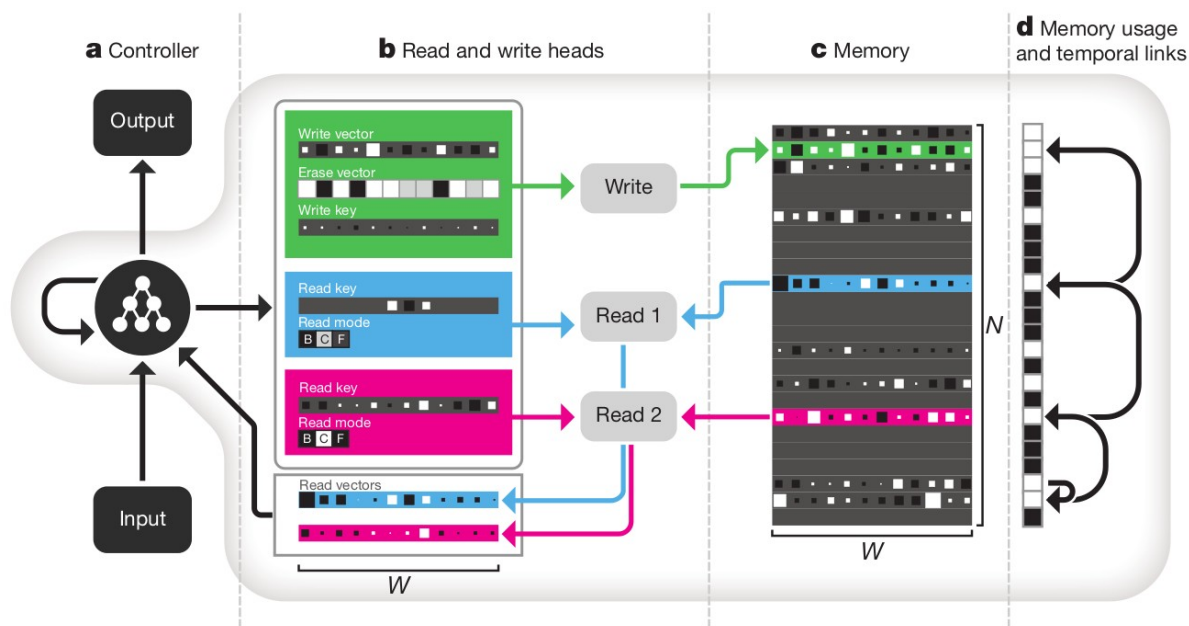
**MANN**

Self-attention

21 / 28

# Differentiable Neural Computer (DNC)

- Newer variant of NTM
- Main difference: Memory usage and temporal links.



background

Attention

Attn. Examples

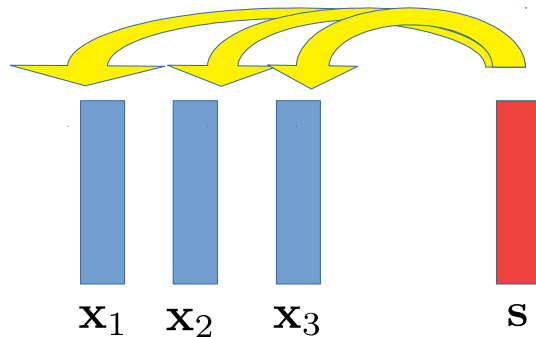
MANN

Self-attention

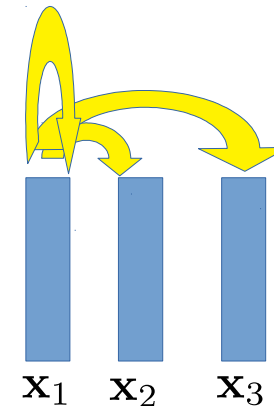
22 / 28

# Self attention

- An attention mechanism which considers its own elements of a sequence
- Different to the attention mechanisms related to RNNs
- Forms the basis of Transformers



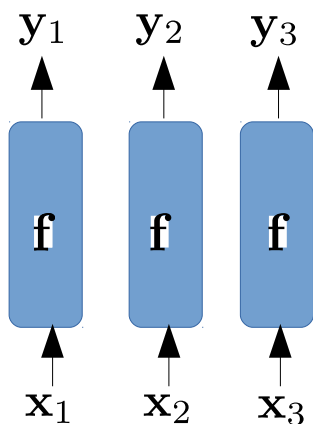
cross-attention



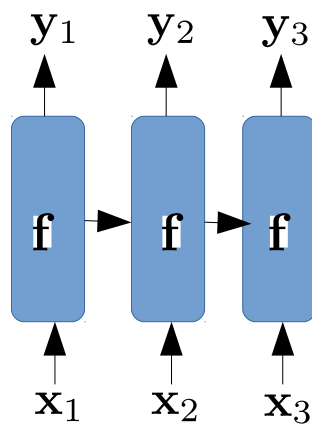
Self-attention

# The intuition

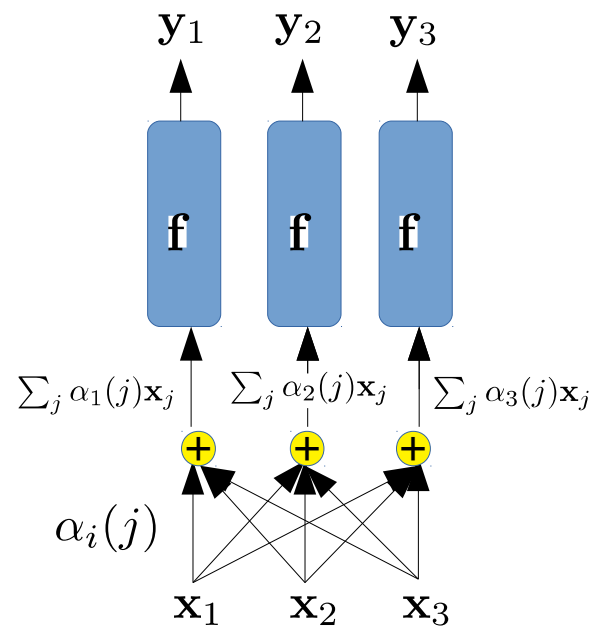
- Make the output a function of variable size sequence
- The number of learnable parameters should be a constant.



Naive



RNN



Self-attention

background

Attention

Attn. Examples

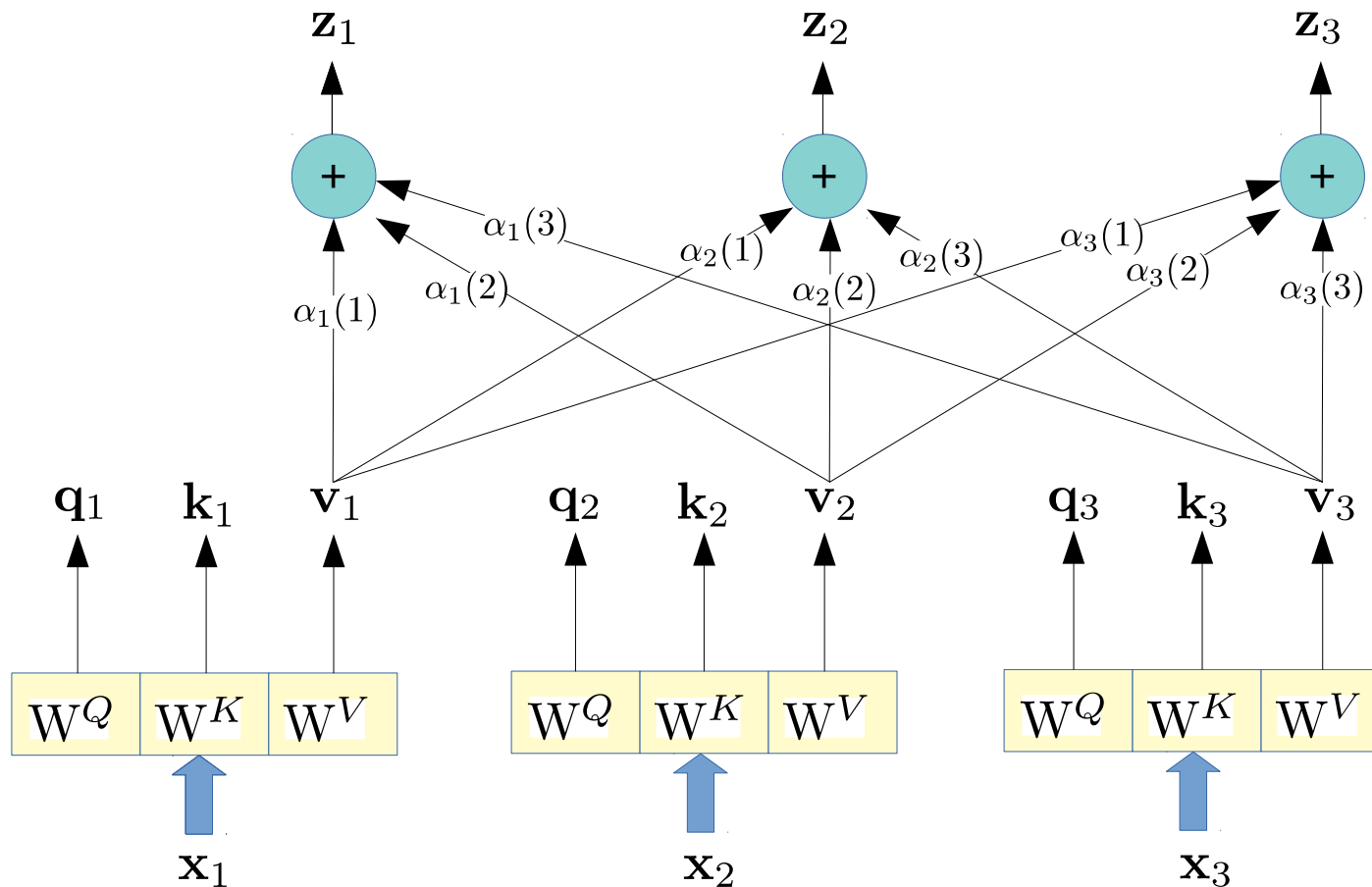
MANN

Self-attention

24 / 28



# Self-attention mechanics (I)



$$\begin{aligned} \mathbf{q}_i &= \mathbf{x}_i W^Q \\ \mathbf{k}_i &= \mathbf{x}_i W^K \\ \mathbf{v}_i &= \mathbf{x}_i W^V \end{aligned}$$

background

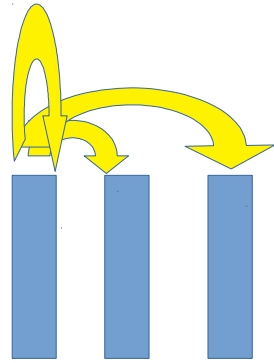
Attention

Attn. Examples

MANN

Self-attention 25 / 28

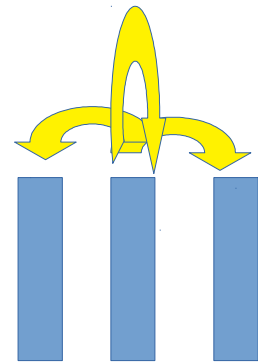
# Self-attention mechanics (II)



$$\alpha_1(1) = \frac{\exp(\mathbf{q}_1 \mathbf{k}_1^T)}{\sum_j \exp(\mathbf{q}_1 \mathbf{k}_j^T)}$$

$$\alpha_1(2) = \frac{\exp(\mathbf{q}_1 \mathbf{k}_2^T)}{\sum_j \exp(\mathbf{q}_1 \mathbf{k}_j^T)}$$

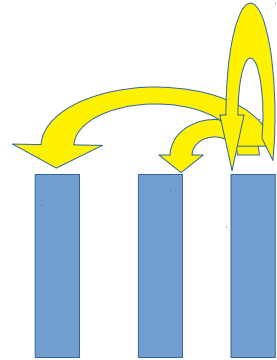
$$\alpha_1(3) = \frac{\exp(\mathbf{q}_1 \mathbf{k}_3^T)}{\sum_j \exp(\mathbf{q}_1 \mathbf{k}_j^T)}$$



$$\alpha_2(1) = \frac{\exp(\mathbf{q}_2 \mathbf{k}_1^T)}{\sum_j \exp(\mathbf{q}_2 \mathbf{k}_j^T)}$$

$$\alpha_2(2) = \frac{\exp(\mathbf{q}_2 \mathbf{k}_2^T)}{\sum_j \exp(\mathbf{q}_2 \mathbf{k}_j^T)}$$

$$\alpha_2(3) = \frac{\exp(\mathbf{q}_2 \mathbf{k}_3^T)}{\sum_j \exp(\mathbf{q}_2 \mathbf{k}_j^T)}$$



$$\alpha_3(1) = \frac{\exp(\mathbf{q}_3 \mathbf{k}_1^T)}{\sum_j \exp(\mathbf{q}_3 \mathbf{k}_j^T)}$$

$$\alpha_3(2) = \frac{\exp(\mathbf{q}_3 \mathbf{k}_2^T)}{\sum_j \exp(\mathbf{q}_3 \mathbf{k}_j^T)}$$

$$\alpha_3(3) = \frac{\exp(\mathbf{q}_3 \mathbf{k}_3^T)}{\sum_j \exp(\mathbf{q}_3 \mathbf{k}_j^T)}$$

$\mathbf{q}_1$	$\mathbf{q}_2$	$\mathbf{q}_3$
$\mathbf{k}_1$	$\mathbf{k}_2$	$\mathbf{k}_3$
$\mathbf{v}_1$	$\mathbf{v}_2$	$\mathbf{v}_3$
$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$

background

Attention

Attn. Examples

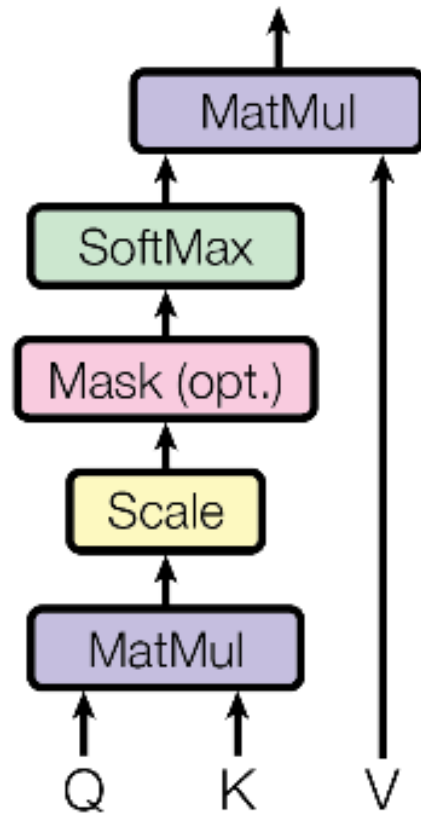
MANN

Self-attention

26 / 28



# Scaled dot product attention



Input word vectors  $X = [x_1, x_2, \dots, x_n]^T$

Query  $Q = [q_1, q_2, \dots, q_n]^T$

Keys  $K = [k_1, k_2, \dots, k_n]^T$

Values  $V = [v_1, v_2, \dots, v_n]^T$

$$Q = XW^Q$$

$$K = XW^K$$

$$V = XW^V$$

$W^Q, W^K, W^V$  Trainable weight vectors

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

# Self-attention Pros and Cons

- Advantages
  - Parallelizable
  - All elements gets equal treatment
  -
- Disadvantages
  - High computational complexity for long sequences
  - Position information gets lost