# Lecture 13: Transformers

Cornell University · Founded A.D. 1865
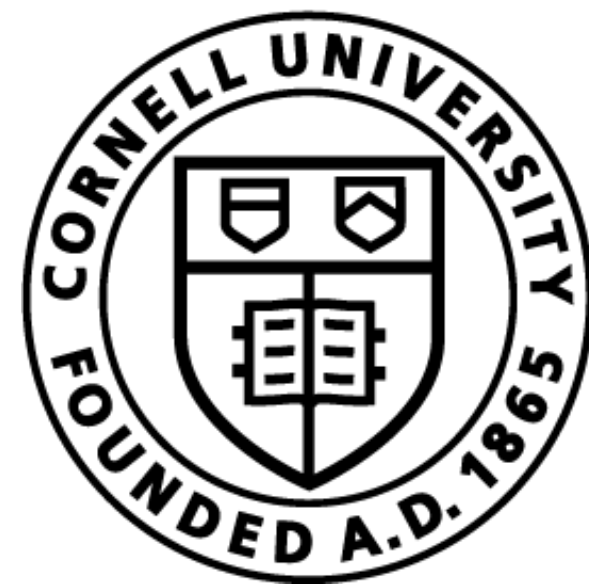
Cornell Bowers C·IS
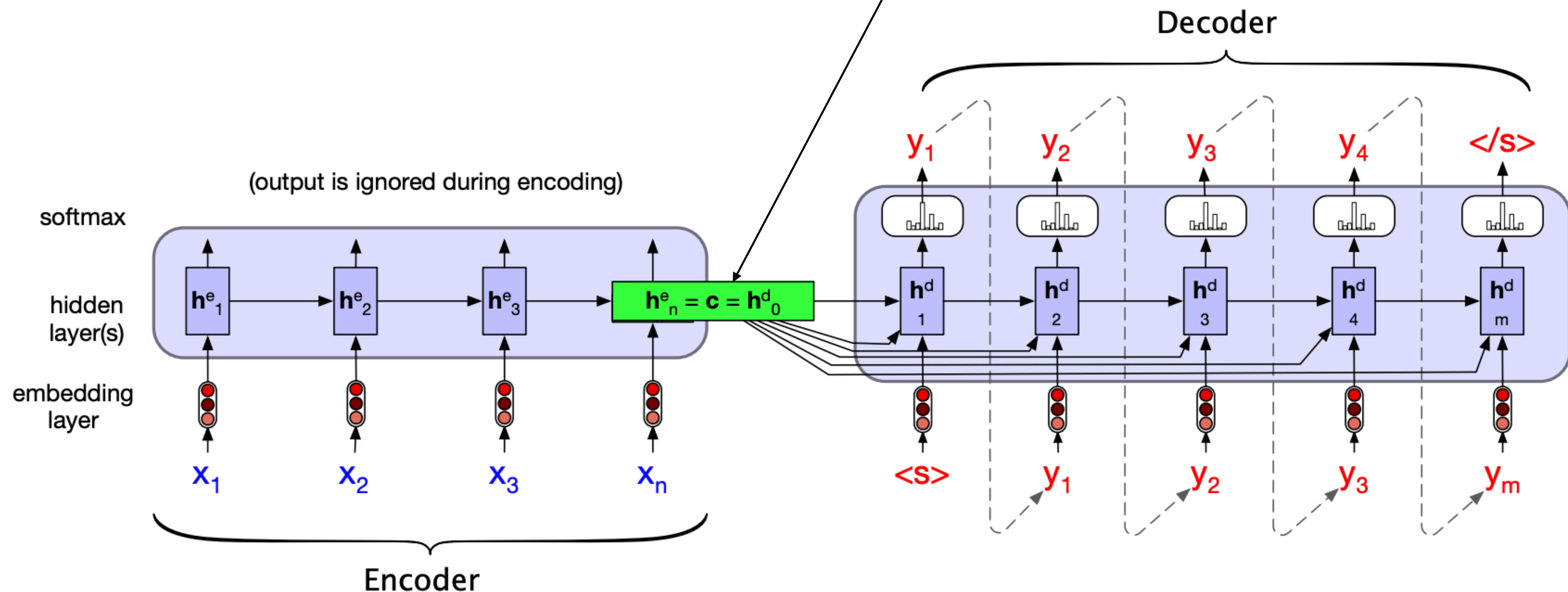**Computer Science**

Claire Cardie, Tanya Goyal

CS 4740 (and crosslists): Introduction to Natural Language Processing

# Today

- Recap: Attention in RNNs

- Transformers

  - Self-Attention

    - Single-head

    - Multi-head

  - Position Embedding

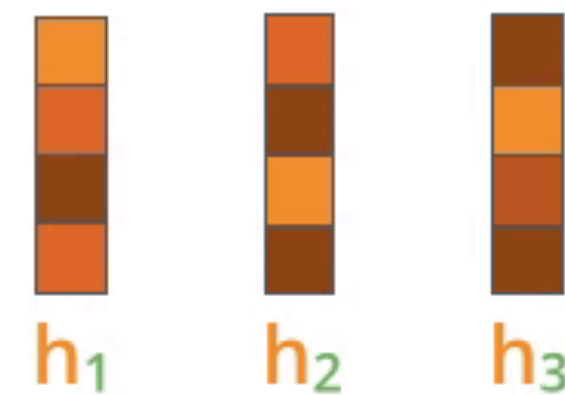# Recap: motivation for attention in the enc/dec framework

If we knew a single relevant encoder state $h_t^{enc}$ to use for our particular decoding step, we could use that instead of a fixed $c^{enc}$.

## Attention at time step 3

1. Prepare inputs

    Encoder hidden states — $h_1$ $h_2$ $h_3$

   Decoder hidden state at time step 2

2. Score each hidden state

   | 13 | 9 | 9 |

   scores
   Attention weights for decoder time step #2

3. Softmax the scores

   | 0.96 | 0.02 | 0.02 |

   softmax scores

4. Multiply each vector by its softmaxed score
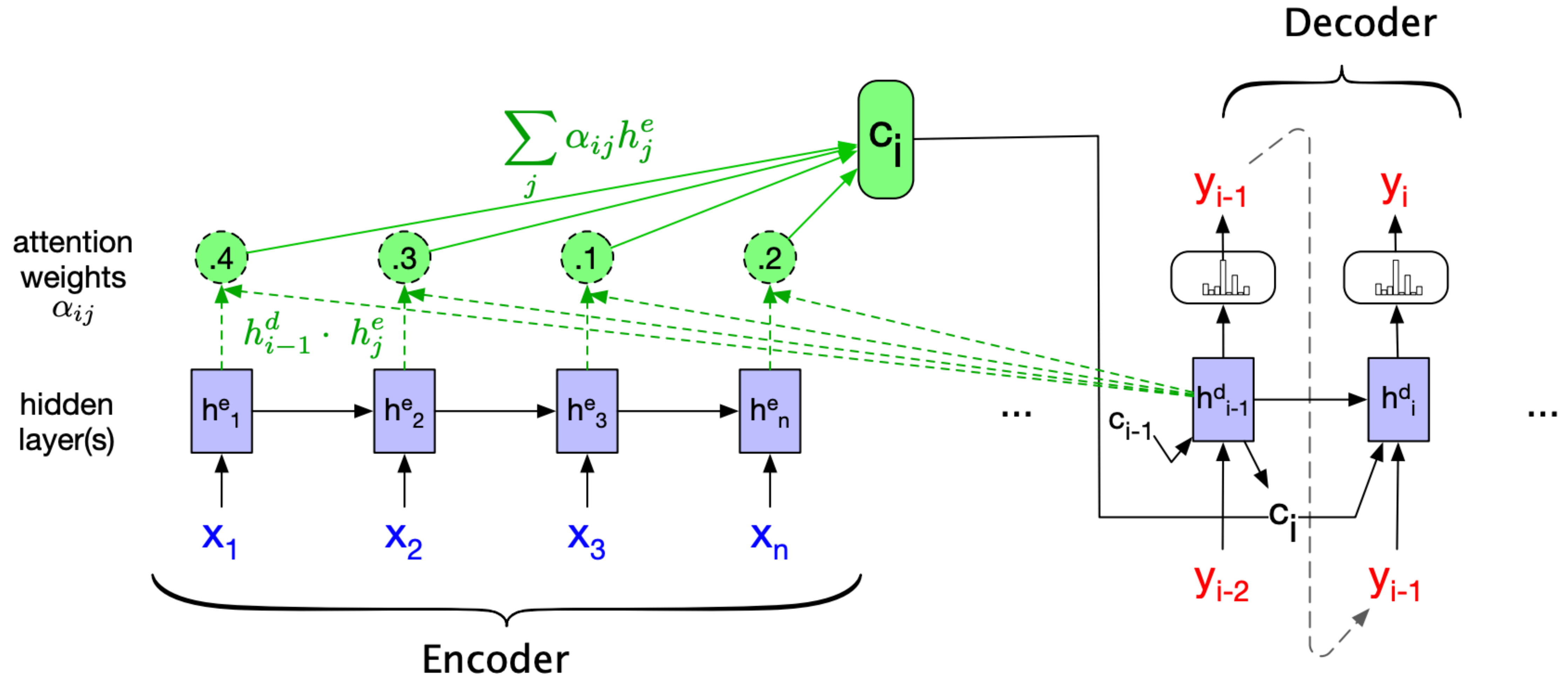
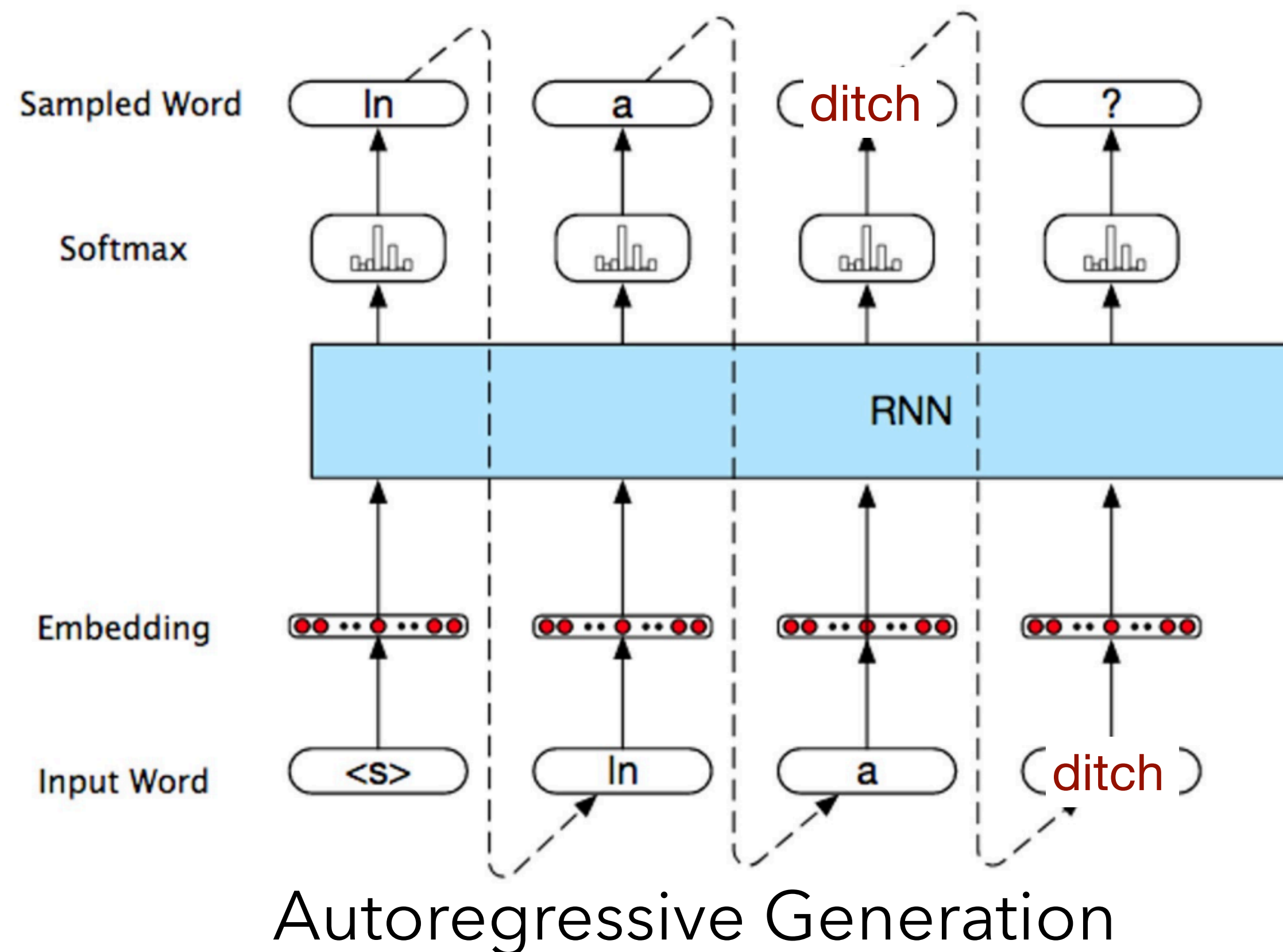5. Sum up the weighted vectors

   Context vector for decoder time step #3

# Recap: Attention allows <u>all</u> enc. hidden states to participate to a weighted degree

# Mini quiz

Q1: What is teacher forcing? Why do we use it?

Compute loss wrt to gold seq.

Autoregressive Generation
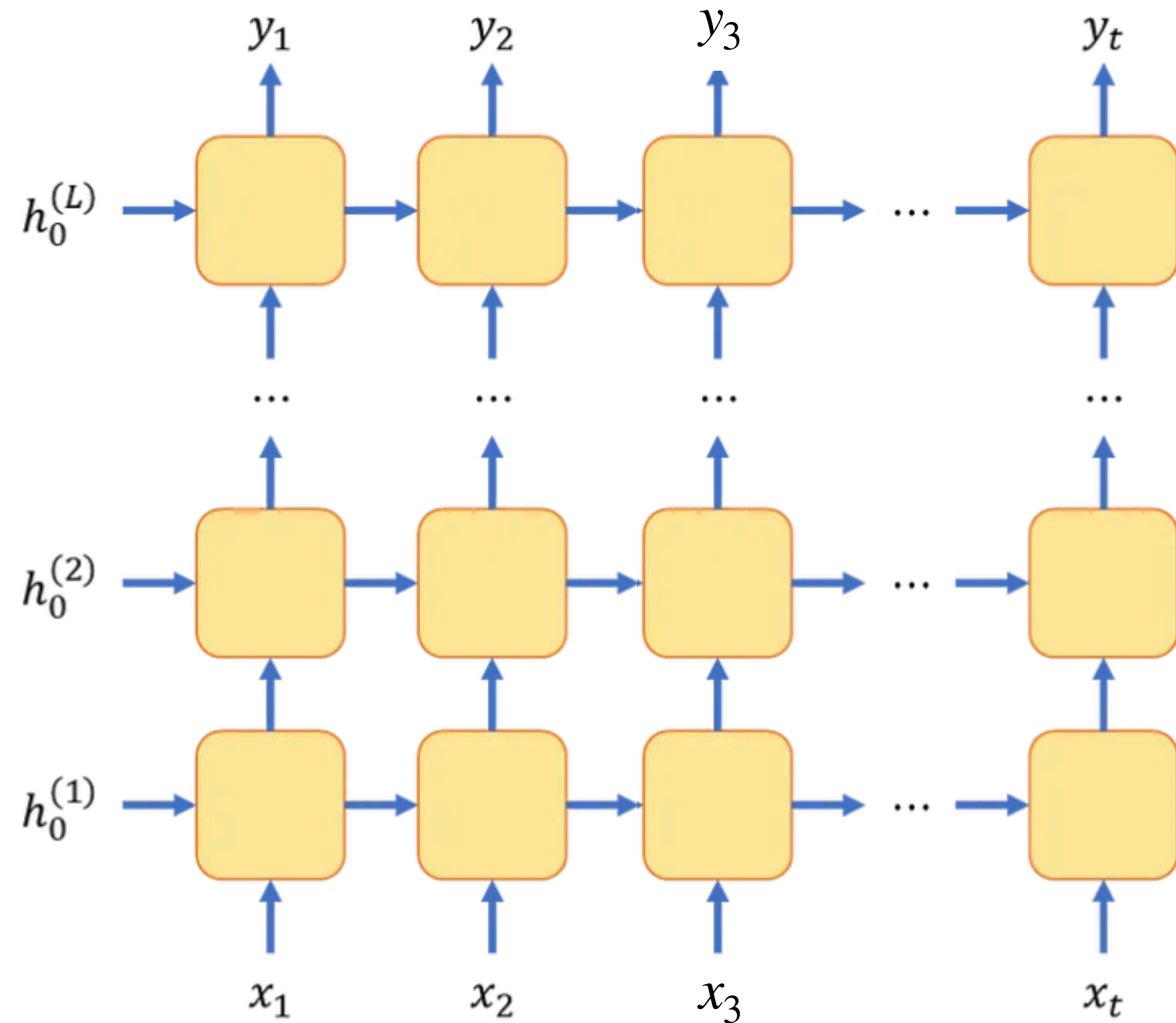
Gold seq.: <s> In a hole in the wall …

Input is always the gold sequence..

**Teacher Forcing during training**

Sampled Word | Softmax | RNN | Embedding | Input Word

In | a | ditch | ?
<s> | In | a | ditch

# With attention, do we need recurrence? Maybe not!



Multi-layer RNN

Computation at time *i* takes into account the computation (hidden values) from time *i-1*.

Above: the computation at time i just looks at the outputs from the previous layer.
Computations at the same layer are parallelizable!
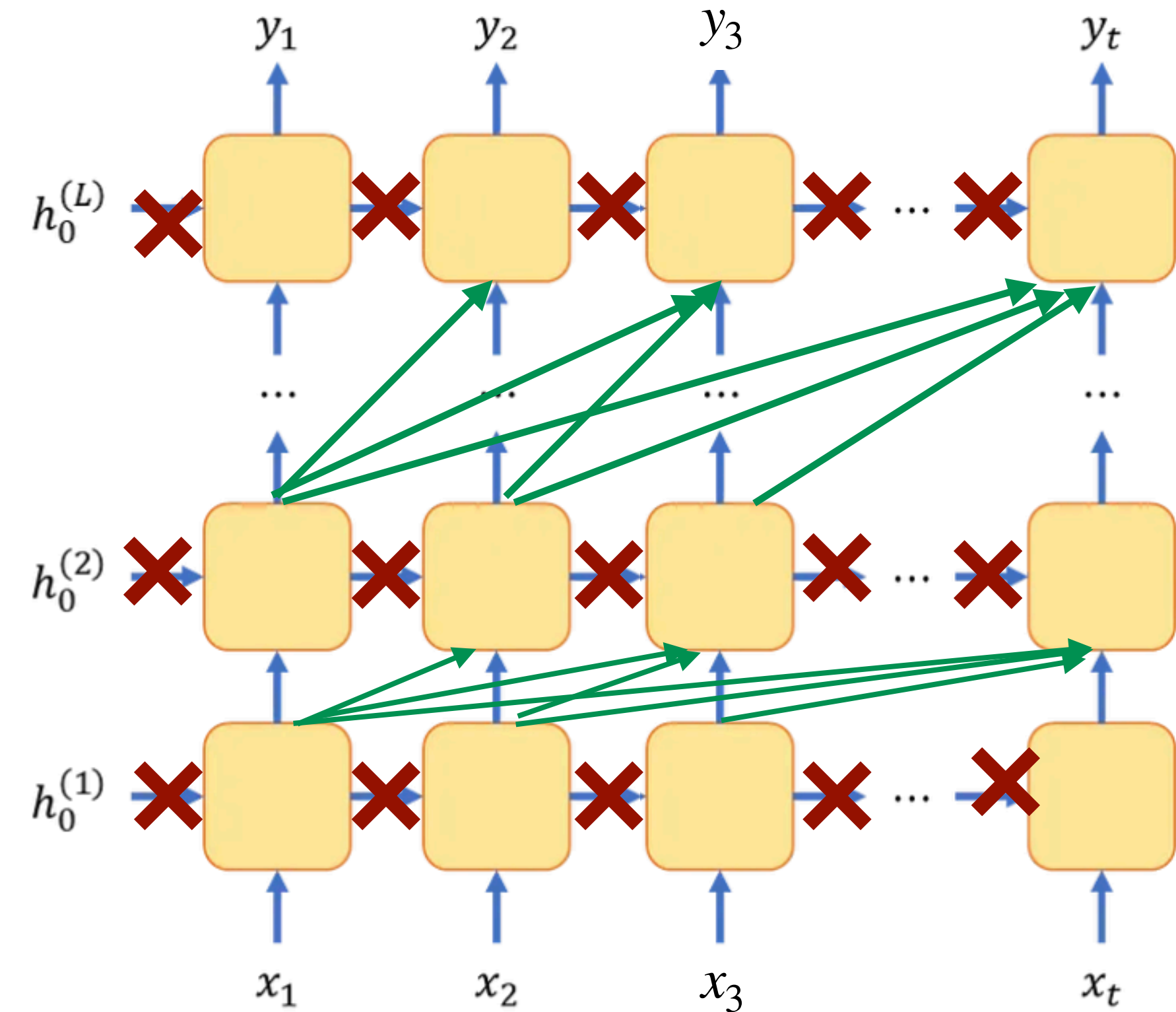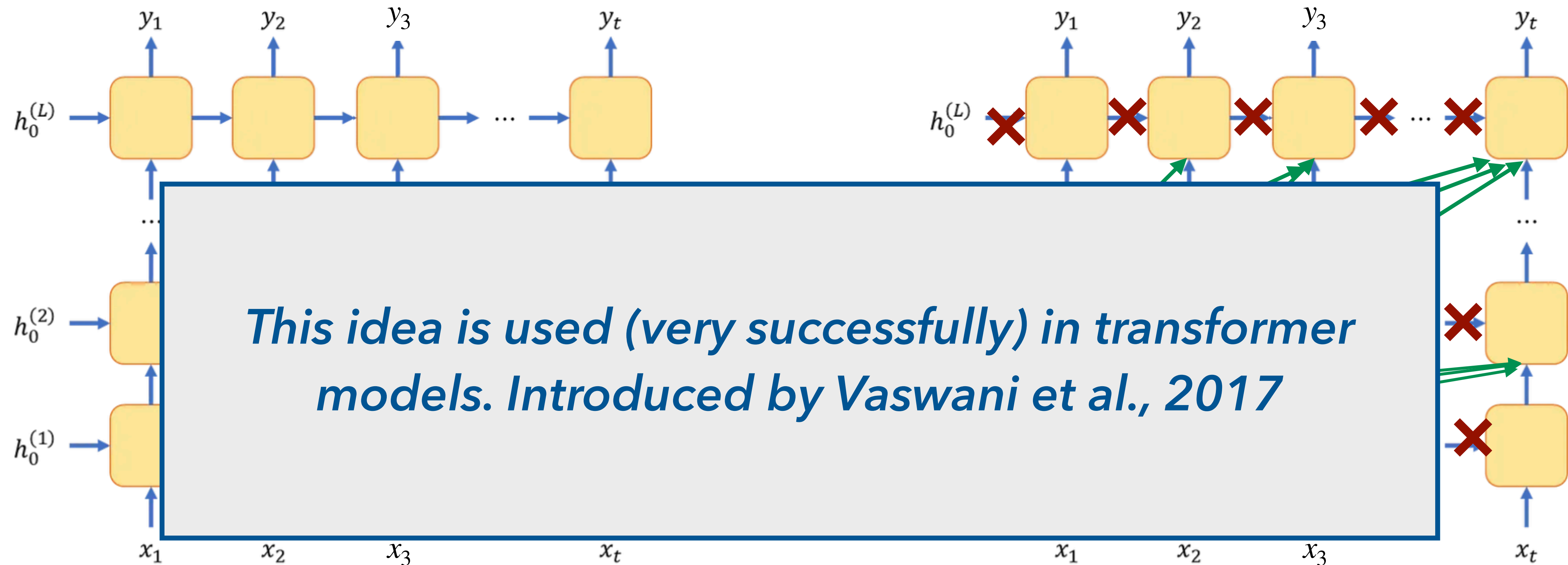
# With attention, do we need recurrence? Maybe not!



**Multi-layer RNN**

Computation at time *i* takes into account the computation (hidden values) from time *i-1*.

*This idea is used (very successfully) in transformer models. Introduced by Vaswani et al., 2017*

Above: the computation at time i just looks at the outputs from the previous layer.
Computations at the same layer are parallelizable!

# Why should we learn about transformers?

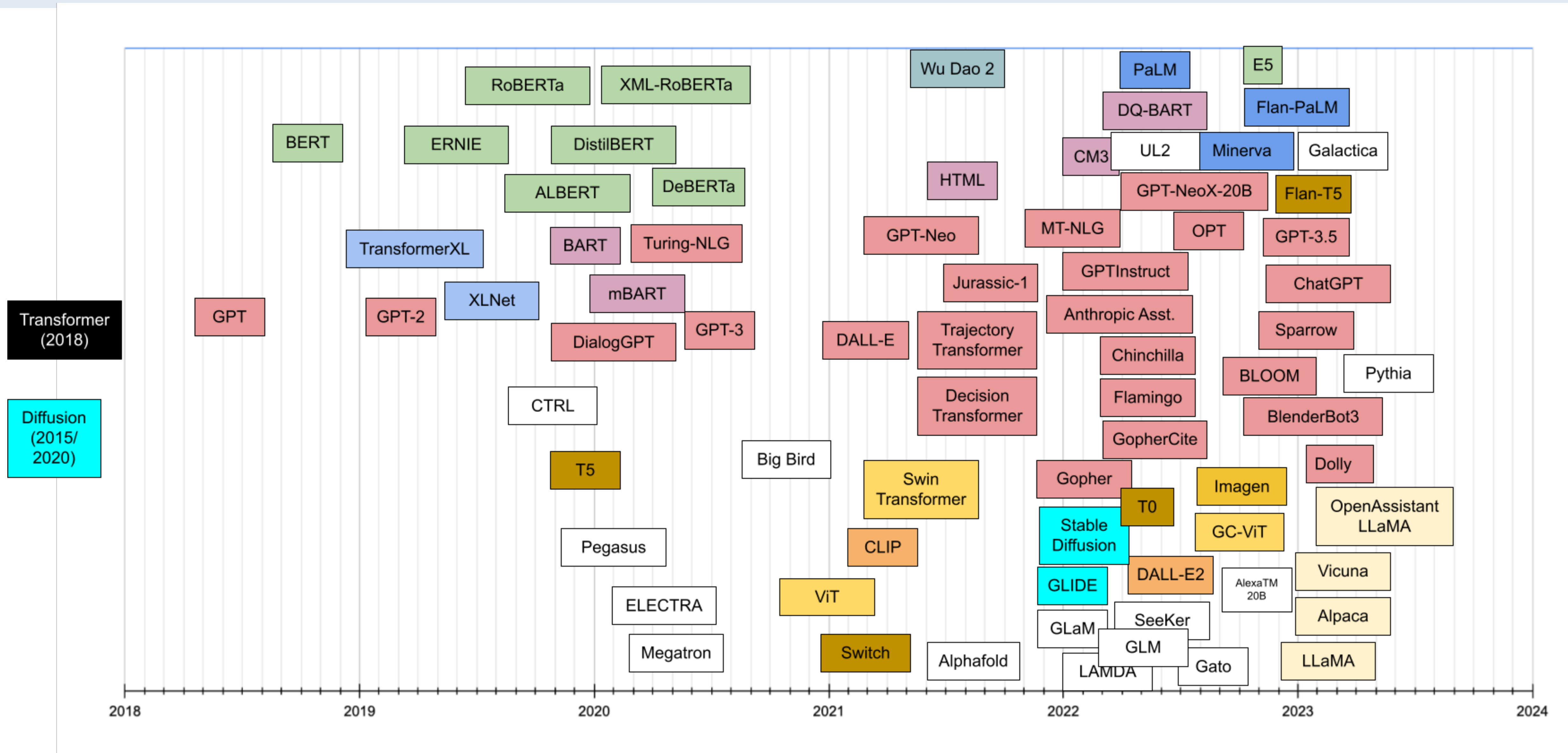- Transformer (variants) are the backbone of all powerful LLMs today!



- Tons of visualizations to trace influence of transformers architecture:
  - Amatriain's: https://amatriain.net/blog/transformer-models-an-introduction-and-catalog-2d1e9039f376/
  - Victor Gaske's: https://ai.v-gar.de/ml/transformer/timeline/

# Why should we learn about transformers?

# Why should we learn about transformers?

- Public Wager: https://www.isattentionallyouneed.com/

- Proposition: On January 1, 2027, a Transformer-like model will continue to hold the state-of-the-art position in most benchmarked tasks in natural language processing.

.

# Today: Transformer Models

- Introduced in Attention Is All You Need (Vaswani et al. NeurIPS 2017)

- A purely attention-based architecture (highly parallelizable), i.e. no recurrence

- Very deep model for NLP (12 layers)

- Originally envisioned for seq2seq tasks (encoder is 6 layers, decoder is 6 layers)

- The encoder and decoder are the same "architecture" applied differently

- **We will first look at the decoder-only transformer today**

.

# Transformer Architecture (Decoder-only)



- We will build up to this!!

- Main components of a transformer model
  - **(Multi-head) Attention**
  - Feed forward
  - Layer Norm

  - Position Encoding

# Simplified Attention

Computation at time step 3, ie. $\mathbf{a}_3$



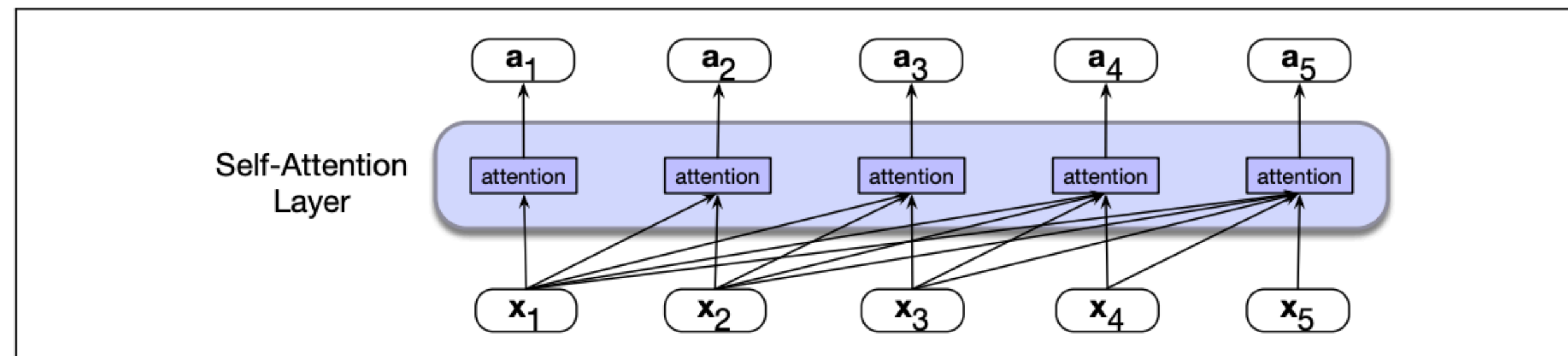**Figure 9.3** Information flow in causal self-attention. When processing each input $\mathbf{x}_i$, the model attends to all the inputs up to, and including $\mathbf{x}_i$.

Step1: prepare inputs

$$\mathbf{x}_1 \quad \mathbf{x}_2 \quad \mathbf{x}_3 \qquad \mathbf{x}_3$$

Step2: compute scores

| 9 | 9 | 13 |
|---|---|----|

- Simplified attention (Similar to RNNs)
  - $\alpha_{ij} = \text{softmax}\left(\text{score}\left(\mathbf{x_i}, \mathbf{x_j}\right)\right), \forall j \leq i$
  - $\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij}\mathbf{x}_i$

Step3: softmax scores
(Attention weights)

| 0.02 | 0.02 | 0.96 |
|------|------|------|

Step4: multiply each vector by softmax scores

$+ \qquad +$

$=$

Step5: sum up the weighted vectors

$\mathbf{a}_3$

# Attention in Transformer models

- Attention in Transformer architectures
  - For a given input $\mathbf{x_i}$ (could be the input at any layer of an encoder or decoder) create three different "roles" or "versions":

  **query** vector:  $\mathbf{q}_i = \mathbf{x}_i\, W^Q$

  **key** vector: $\mathbf{k}_i = \mathbf{x}_i\, W^K$

  **value** vector: $\mathbf{v}_i = \mathbf{x}_i\, W^V$

  $W^Q, W^K, W^V$ are learned matrices

Computation at time step 3, ie. $\mathbf{a}_3$

Step1: prepare inputs

$\mathbf{x_1}$  $\mathbf{x_2}$  $\mathbf{x_3}$      $\mathbf{x_3}$

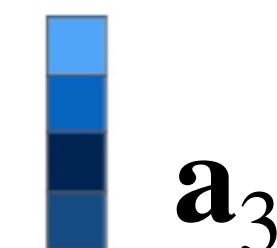| Step2: compute scores | 9 | 9 | 13 |
|---|---|---|---|

Step3: softmax scores
(Attention weights)

| 0.02 | 0.02 | 0.96 |
|---|---|---|

Step4: multiply each vector by softmax scores

Step5: sum up the weighted vectors  $\mathbf{a}_3$

# Attention in Transformer models

- Attention in Transformer architectures
  - For a given input $\mathbf{x_i}$ (could be the input at any layer of an encoder or decoder) create three different "roles" or "versions":

    **query** vector: $\mathbf{q}_i = \mathbf{x}_i \, W^Q$

    **key** vector: $\mathbf{k}_i = \mathbf{x}_i \, W^K$

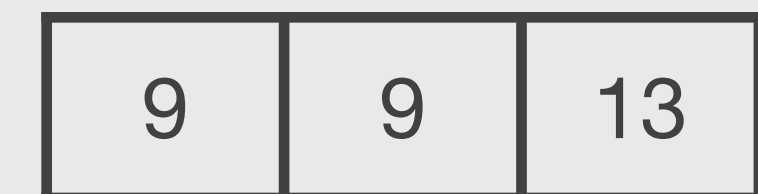    **value** vector: $\mathbf{v}_i = \mathbf{x}_i \, W^V$

    $W^Q, W^K, W^V$ are learned matrices

Step1: prepare inputs

$\mathbf{k_1}$  $\mathbf{k_2}$  $\mathbf{k_3}$  $\mathbf{q_3} = W^Q \mathbf{x_3}$

Step2: compute scores

| 9 | 9 | 13 |
|---|---|----|

Step3: softmax scores
(Attention weights)

| 0.02 | 0.02 | 0.96 |
|------|------|------|

Step4: multiply each vector by softmax scores

$+$  $+$

$=$

Step5: sum up the weighted vectors

$\mathbf{a}_3$

# Attention in Transformer models

- Attention in Transformer architectures
  - For a given input $\mathbf{x_i}$ (could be the input at any layer of an encoder or decoder) create three different "roles" or "versions":

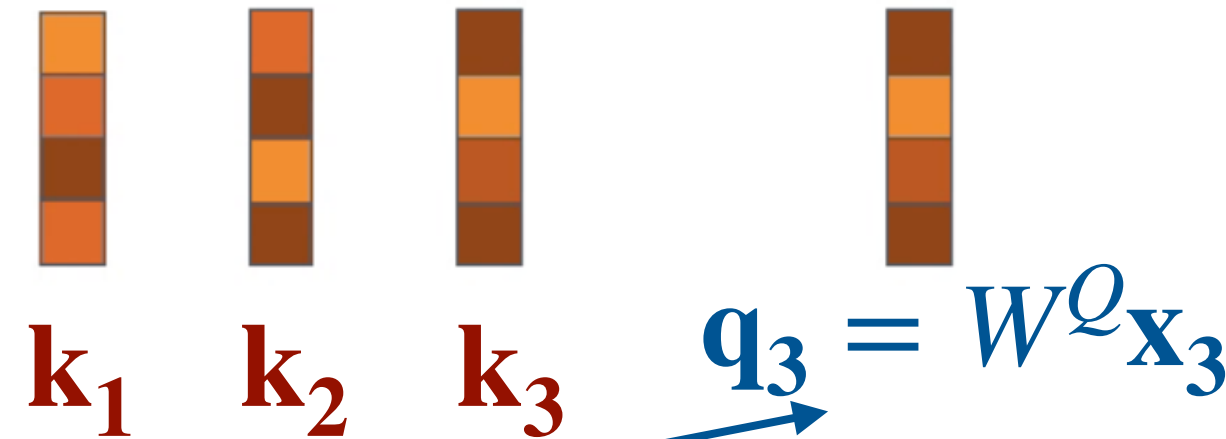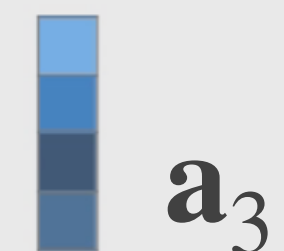    **query** vector: $\mathbf{q}_i = \mathbf{x}_i \, W^Q$

    **key** vector: $\mathbf{k}_i = \mathbf{x}_i \, W^K$

    **value** vector: $\mathbf{v}_i = \mathbf{x}_i \, W^V$

    $W^Q, W^K, W^V$ are learned matrices

Step1: prepare inputs

$\mathbf{k}_1 \quad \mathbf{k}_2 \quad \mathbf{k}_3 \qquad \mathbf{q_3} = W^Q\mathbf{x_3}$

Step2: compute scores

| 9 | 9 | 13 |

$\dfrac{\mathbf{q}_3 \cdot \mathbf{k}_1}{\sqrt{d_k}}$

Step3: softmax scores
(Attention weights)

| 0.02 | 0.02 | 0.96 |

$d_k$ is dim of query, key vectors

Step4: multiply each vector by softmax scores

+  +

=

Step5: sum up the weighted vectors

$\mathbf{a}_3$

# Attention in Transformer models

- Attention in Transformer architectures
  - For a given input $\mathbf{x_i}$ (could be the input at any layer of an encoder or decoder) create three different "roles" or "versions":
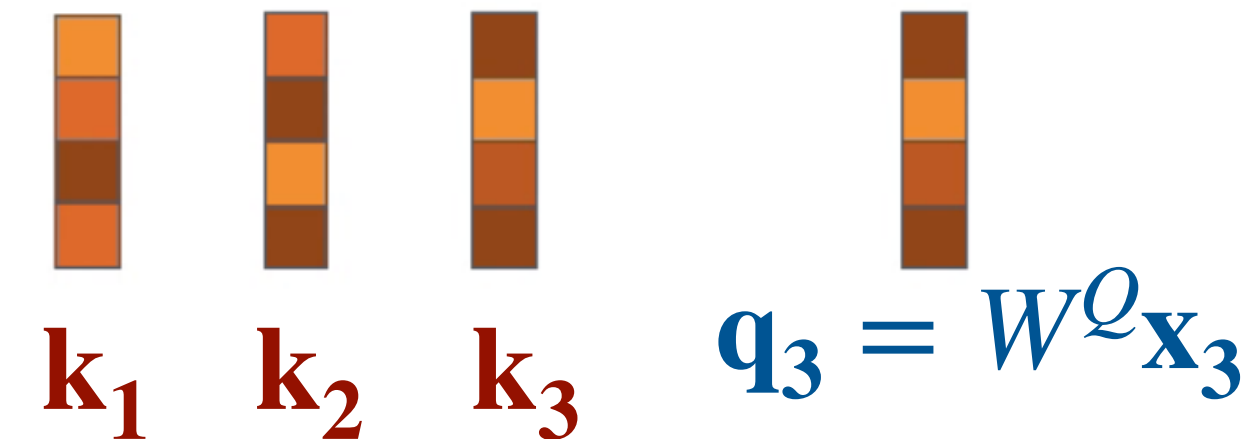
  **query** vector: $\mathbf{q}_i = \mathbf{x}_i\, W^Q$

  **key** vector: $\mathbf{k}_i = \mathbf{x}_i\, W^K$
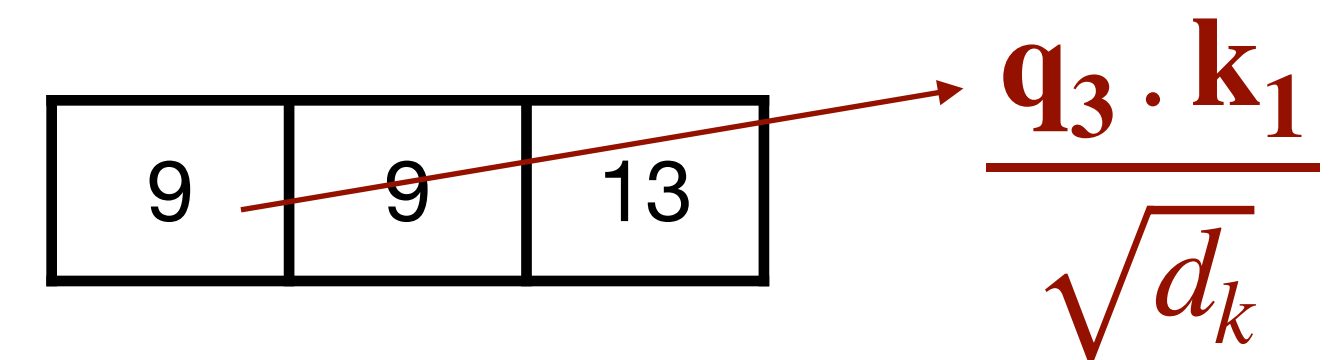
  **value** vector: $\mathbf{v}_i = \mathbf{x}_i\, W^V$

  $W^Q, W^K, W^V$ are learned matrices

Step1: prepare inputs

$\mathbf{k}_1$ $\quad$ $\mathbf{k}_2$ $\quad$ $\mathbf{k}_3$ $\qquad$ $\mathbf{q}_3 = W^Q\mathbf{x}_3$

Step2: compute scores

| 9 | 9 | 13 |

$\dfrac{\mathbf{q}_3 \cdot \mathbf{k}_1}{\sqrt{d_k}}$

Step3: softmax scores
(Attention weights)

| 0.02 | 0.02 | 0.96 |

$\alpha_{31}\mathbf{v}_1 \quad \alpha_{32}\mathbf{v}_2 \quad \alpha_{33}\mathbf{v}_3$

Step4: multiply each vector by softmax scores

$+ \qquad + \qquad =$

$=$

Step5: sum up the weighted vectors

$\mathbf{a}_3$

# Attention in Transformer models

- Attention in Transformer architectures

  - For a given input $\mathbf{x_i}$ (could be the input at any layer of an encoder or decoder) create three different "roles" or "versions":
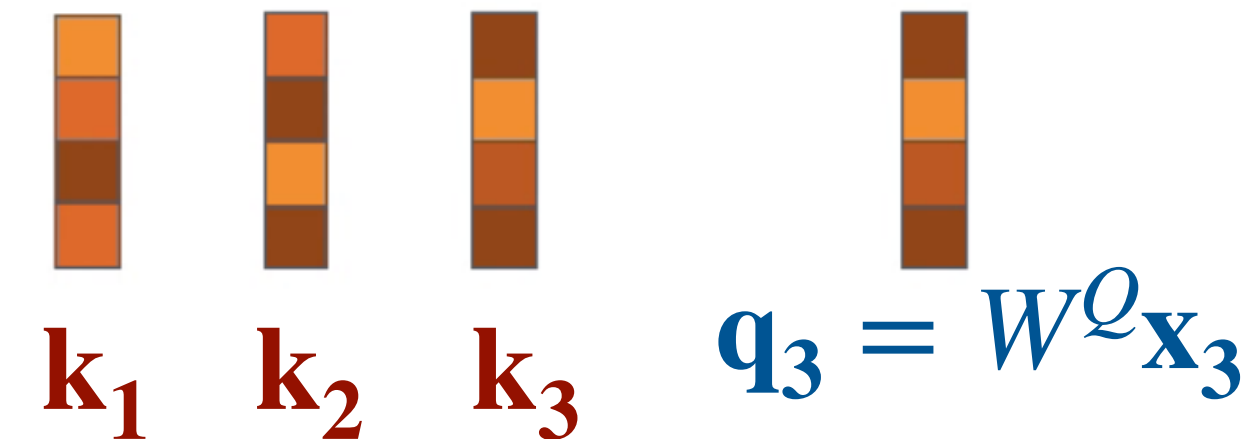
  **query** vector: $\mathbf{q}_i = \mathbf{x}_i\, W^Q$

  **key** vector: $\mathbf{k}_i = \mathbf{x}_i\, W^K$
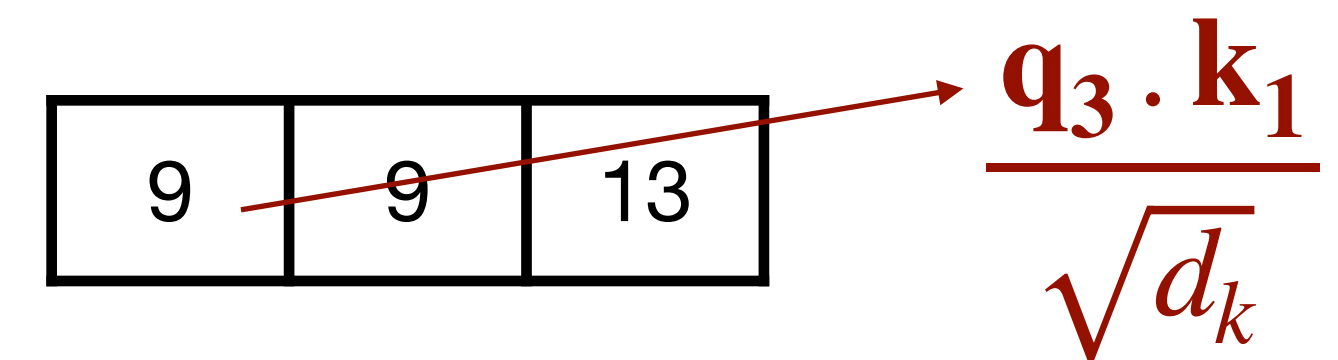
  **value** vector: $\mathbf{v}_i = \mathbf{x}_i\, W^V$

  $W^Q, W^K, W^V$ are learned matrices

Step1: prepare inputs

$\mathbf{k_1}$  $\mathbf{k_2}$  $\mathbf{k_3}$    $\mathbf{q_3} = W^Q\mathbf{x_3}$

Step2: compute scores

| 9 | 9 | 13 |

$\dfrac{\mathbf{q_3} \cdot \mathbf{k_1}}{\sqrt{d_k}}$

Step3: softmax scores
(Attention weights)

| 0.02 | 0.02 | 0.96 |

$\alpha_{31}\mathbf{v}_1$  $\alpha_{32}\mathbf{v}_2$  $\alpha_{33}\mathbf{v}_3$

Step4: multiply each vector by softmax scores

$+$   $+$   $=$

$=$

Step5: sum up the weighted vectors **and project**

$W^O = \mathbf{a}_3$

# Attention in Transformer models

Computation at time step 3, ie. $\mathbf{a_3}$



Step1: prepare inputs

$$\mathbf{k_1} \quad \mathbf{k_2} \quad \mathbf{k_3} \qquad \mathbf{q_3} = W^Q \mathbf{x_3}$$

Step2: compute scores

| 9 | 9 | 13 |
|---|---|---|

$$\frac{\mathbf{q_3} \cdot \mathbf{k_1}}{\sqrt{d_k}}$$

Step3: softmax scores
(Attention weights)

| 0.02 | 0.02 | 0.96 |
|------|------|------|

$$\alpha_{31}\mathbf{v}_1 \quad \alpha_{32}\mathbf{v}_2 \quad \alpha_{33}\mathbf{v}_3$$

Step4: multiply each vector by softmax scores

Step5: sum up the weighted vectors **and project**

$$W^O = \mathbf{a_3}$$

8. Output of self-attention $a_3$ [1 × d]

7. Reshape to [1 x d]  $W^O$ [$d_v$ × d]

[1 × $d_v$]

6. Sum the weighted value vectors  $\Sigma$

[1 × $d_v$]  [1 × $d_v$]  [1 × $d_v$]

5. Weigh each **value** vector
$\alpha_{3,1}$  $\alpha_{3,2}$  $\alpha_{3,3}$

4. Turn into $\alpha_{i,j}$ weights via softmax

3. Divide scalar score by √$d_k$   √$d_k$   √$d_k$   √$d_k$

2. Compare x3's **query** with the **keys** for x1, x2, and x3
[1 × $d_v$]  [1 × $d_v$]  [1 × $d_v$]

1. Generate **key, query, value** vectors
k q v  $W^K W^Q W^V$   k q v  $W^K W^Q W^V$   k q v  $W^K W^Q W^V$

$x_1$   $x_2$   $x_3$

[1 × d]   [1 × d]   [1 × d]

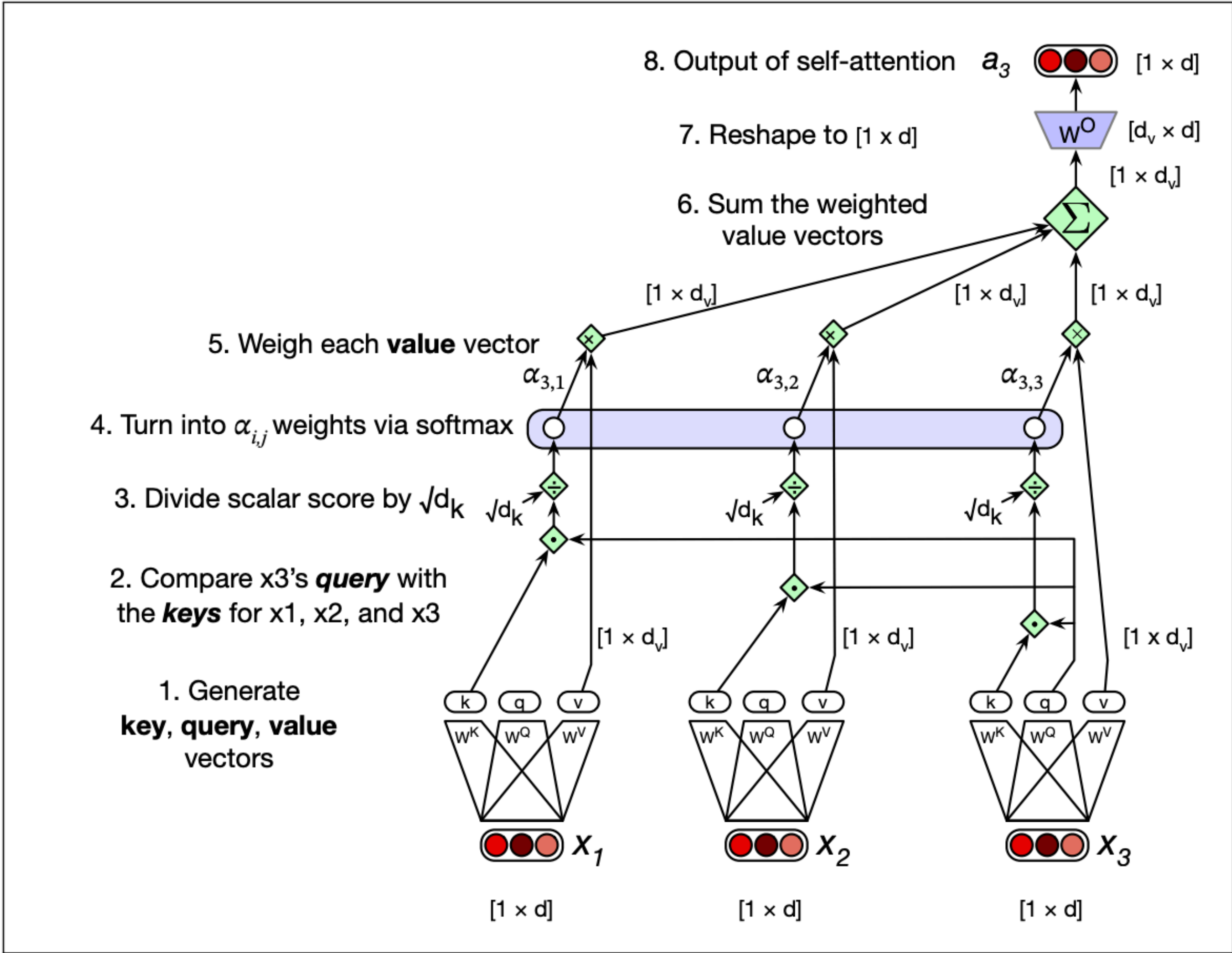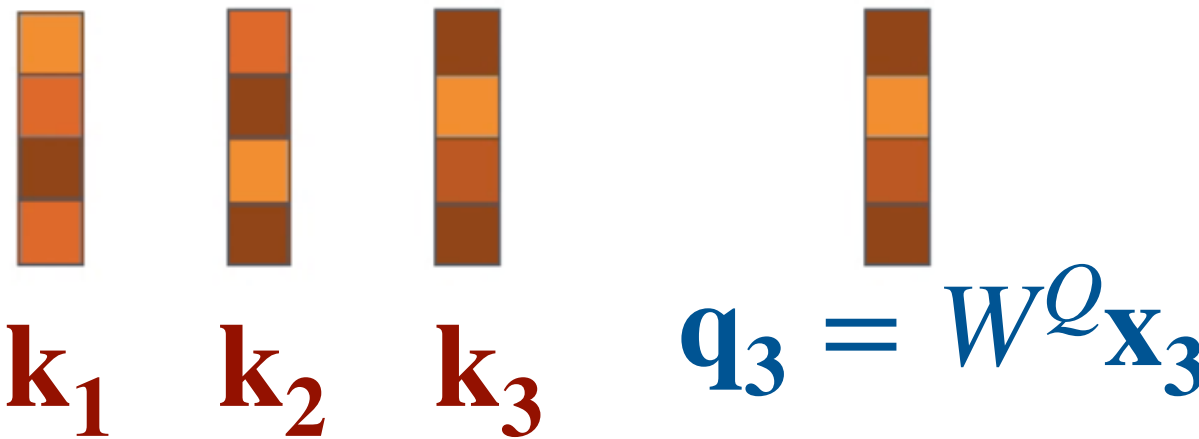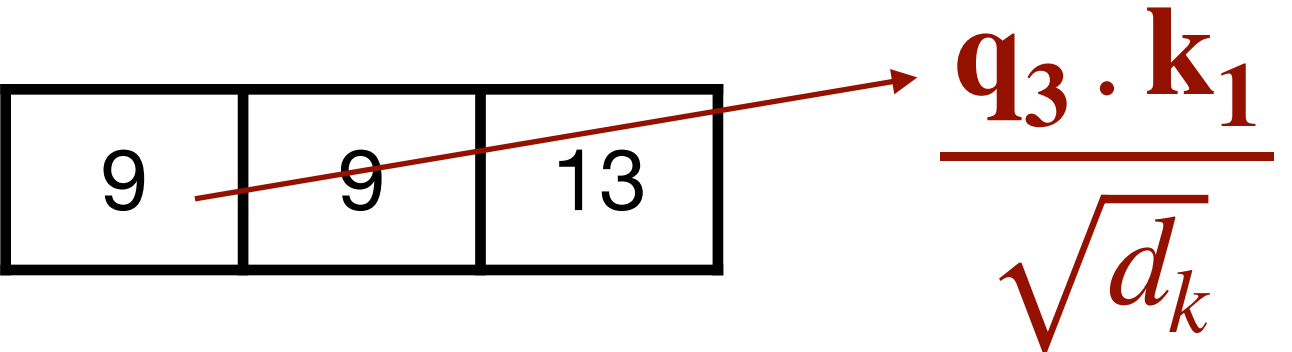**Figure 9.4** Calculating the value of $\mathbf{a}_3$, the third element of a sequence using causal (left-to-right) self-attention.

# Attention Computation (matrix form)
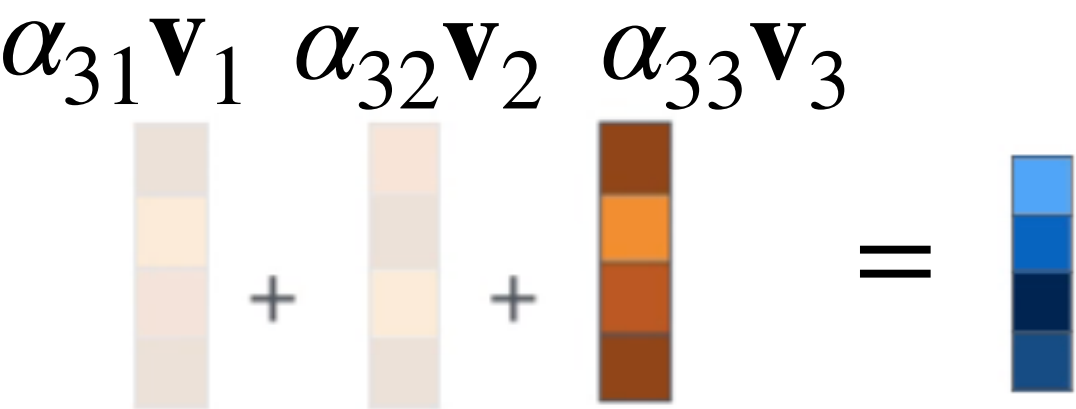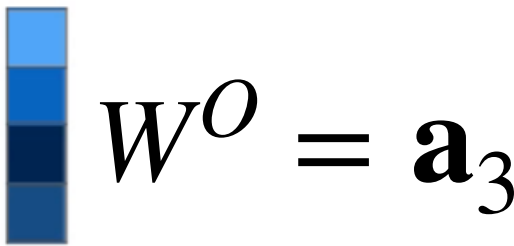


$Q$

| $q_1$ |
|---|
| $q_2$ |
| $q_3$ |
| $q_4$ |

$K$

| $k_1$ |
|---|
| $k_2$ |
| $k_3$ |
| $k_4$ |

$V$

| $v_1$ |
|---|
| $v_2$ |
| $v_3$ |
| $v_4$ |

$$\left( \dfrac{\overset{Q}{\boxed{\begin{array}{c} q_1 \\ q_2 \\ q_3 \\ q_4 \end{array}}} \; \overset{K^T}{\boxed{k_1 \, k_2 \, k_3 \, k_4}}}{\sqrt{d_k}} \right) \oplus$$

| 0 | -∞ | -∞ | -∞ |
|---|---|---|---|
| 0 | 0 | -∞ | -∞ |
| 0 | 0 | 0 | -∞ |
| 0 | 0 | 0 | 0 |

Causal mask so that inputs at time i only attend to previous inputs (Aside: in transformer encoders, inputs at time i attend to all other encoder inputs)

**H = Softmax operation (along dim = 0)**

$(HV) \, W^O =$

| $a_1$ |
|---|
| $a_2$ |
| $a_3$ |
| $a_4$ |

# Multi-headed Attention



**Figure 9.5** The multi-head attention computation for input $\mathbf{x}_i$, producing output $\mathbf{a}_i$. A multi-head attention layer has $A$ heads, each with its own key, query and value weight matrices. The outputs from each of the heads are concatenated and then projected down to $d$, thus producing an output of the same size as the input.

- Multiple heads –> multiple ``independent'' projections (keys, queries, values) for each input.
- Each head has different $W^Q, W^K, W^V$ matrices
- Different heads can potentially capture different phenomenon.
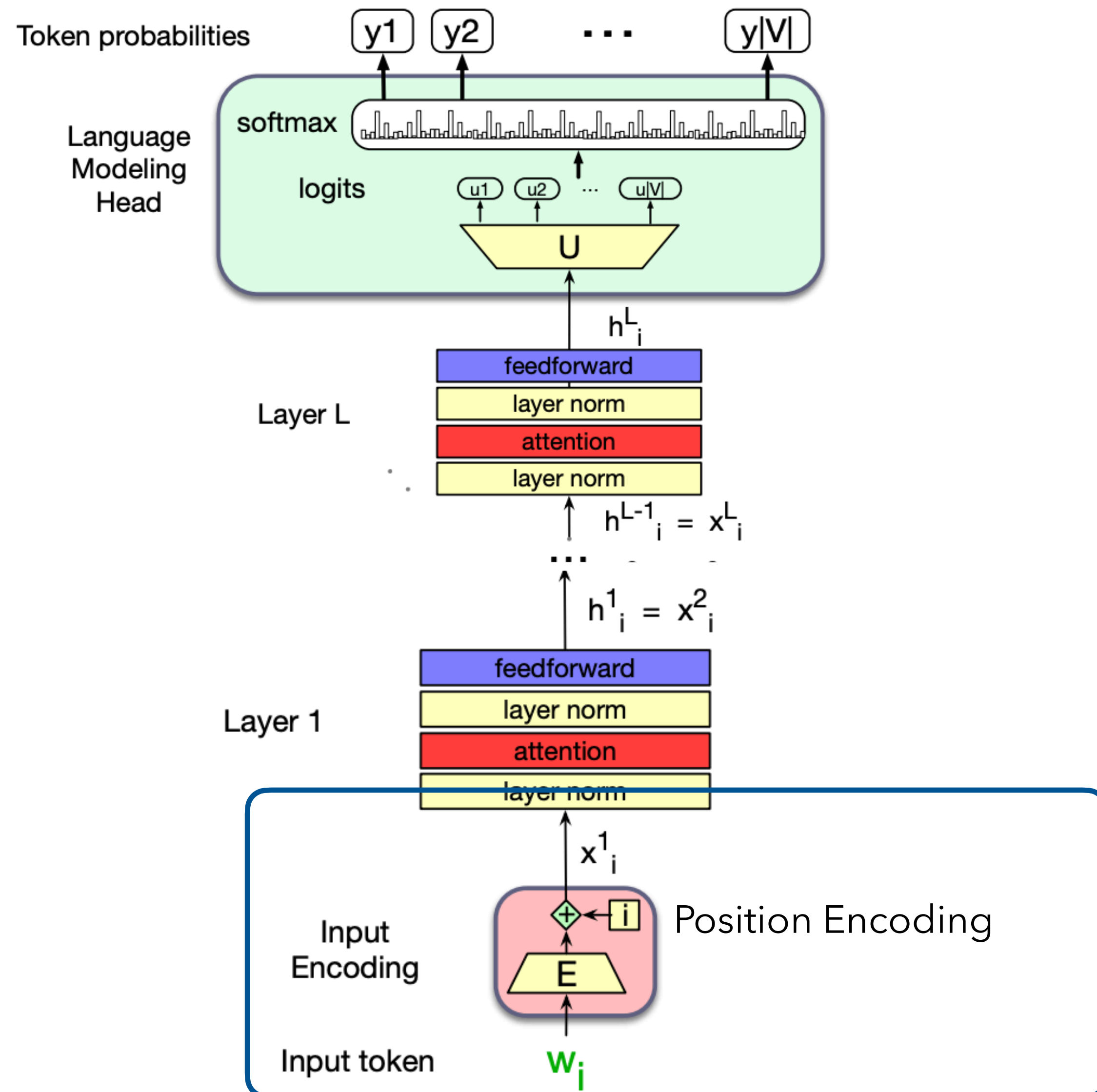
# Zooming out



- <u>Self-attention layer</u> transformed the input $\mathbf{x_i}$ to output $\mathbf{a_i}$
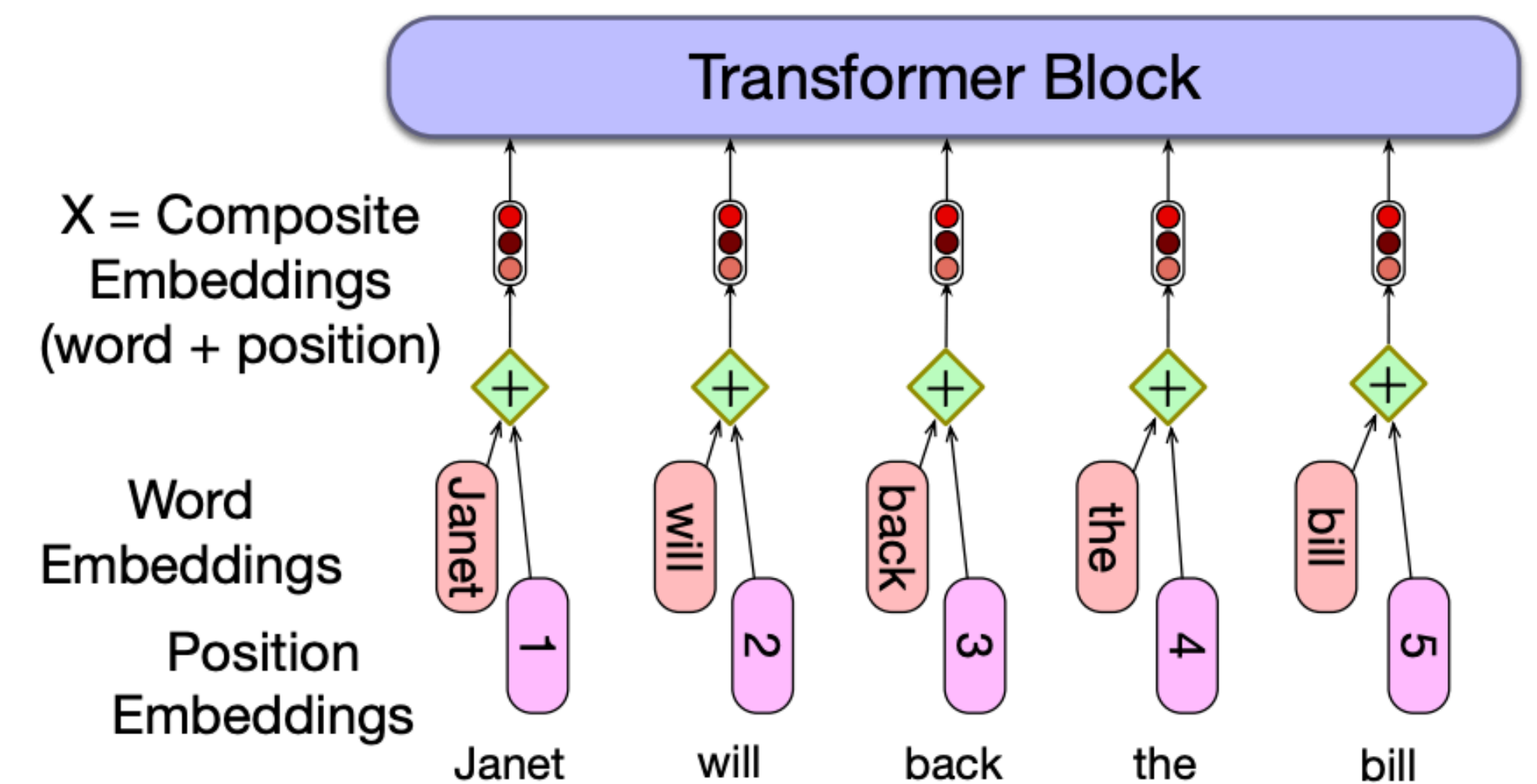- Word order information is lost!

*An <u>old</u> dog and a <u>young</u> boy ….*

- *boy* attends to both <u>old</u> and <u>young</u>. We <u>young</u> to have a higher influence on *boy*'s hidden representation than <u>old</u>. Attention does not ensure this.
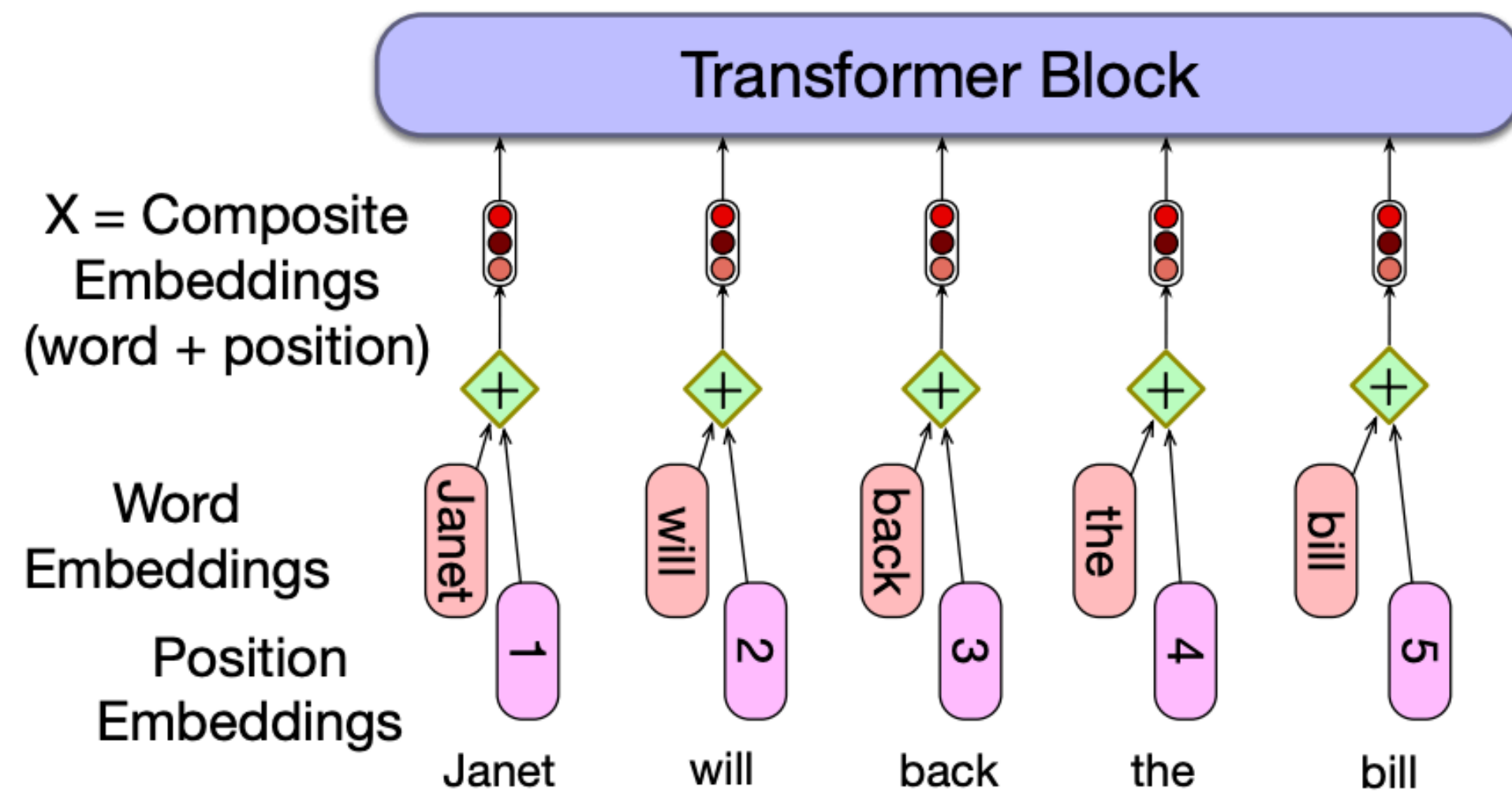- Q: How do RNNs include this information?

# Let's go back to our transformer arch



- <u>Solution:</u> Add a "position" embedding to the word embedding to produce a new embedding of the same dimension.

# Let's go back to our transformer arch
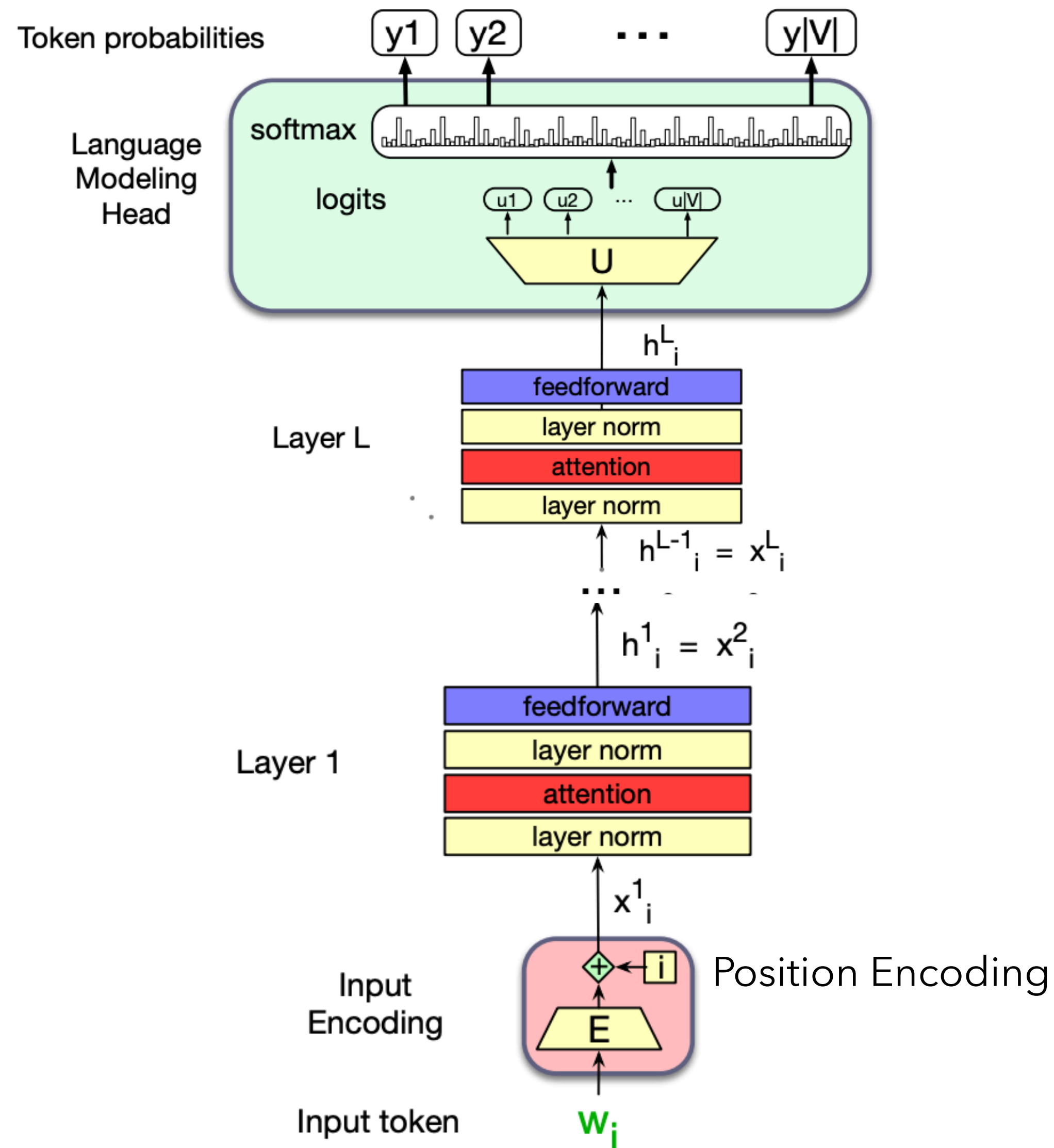


- <u>Solution:</u> Add a "position" embedding to the word embedding to produce a new embedding of the same dimension.

How do we get these positional embeddings?

- Assume all sequences will have length between 0 to N (say 512). Randomly initialize embeddings for each position.

- These will get trained with other transformer parameters.

# Let's go back to our transformer arch



- <u>Today</u>:
  - Multi-head self-attention
  - Position Embeddings

- <u>Next Class</u>:
  - Layer Norm
  - Feedforward layer
  - Putting it all together
  - Encoder Decoder

# Slide Acknowledgements

‣ Earlier versions of this course offerings including materials from Claire Cardie, Marten van Schijndel, Lillian Lee.

‣