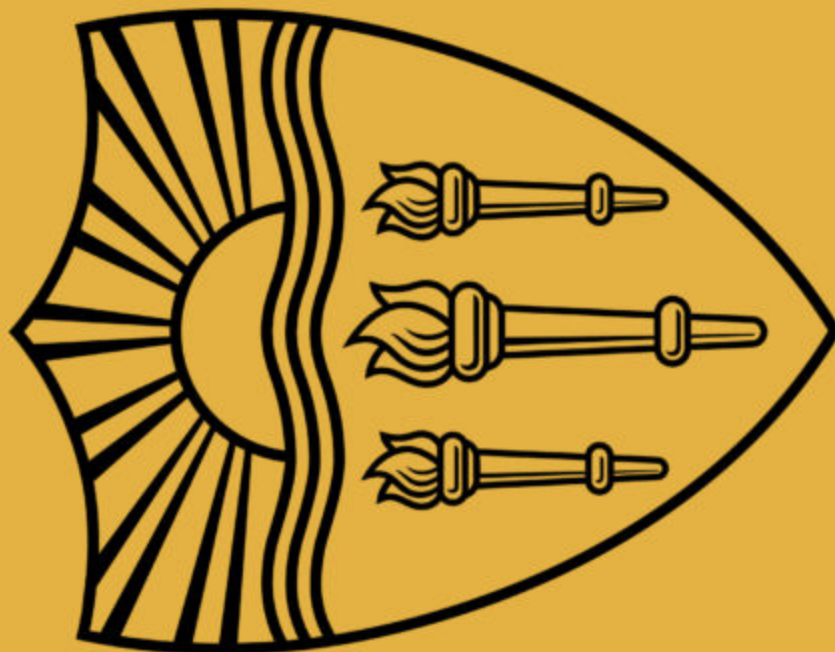


CSCI 499

Lecture 11: Transformers: Building Blocks

Instructor: Swabha Swayamdipta
USC CSCI 499 LMs in NLP
Feb 26, Spring 2024



Logistics / Announcements

- Today:
 - HW2
 - Quiz 3 in class
- This Wednesday:
 - In-class project discussions
 - Make sure the entire team is present!
- Upcoming guest / TA lectures
- No Office Hours in the Spring Break week

Feb 26:	Transformers - Building Blocks I	HW2 Due
Feb 28:	PROJECT DISCUSSIONS	
Mar 4:	Transformers - Building Blocks II	HW3 Released PROGRESS REPORT DUE
Mar 6:	TA Lecture: PyTorch for Transformers	
Mar 11:	No Class SPRING BREAK	
Mar 13:	No Class SPRING BREAK	

Large Language Models (LLMs)

Mar 18:	Pre-training Transformers I	HW3 Due
Mar 20:	Pre-training Transformers II	HW4 Released
Mar 25:	Guest Lecture: Limitations and Harms of LLMs	
Mar 27:	Generating from Language Models I	

Lecture Outline

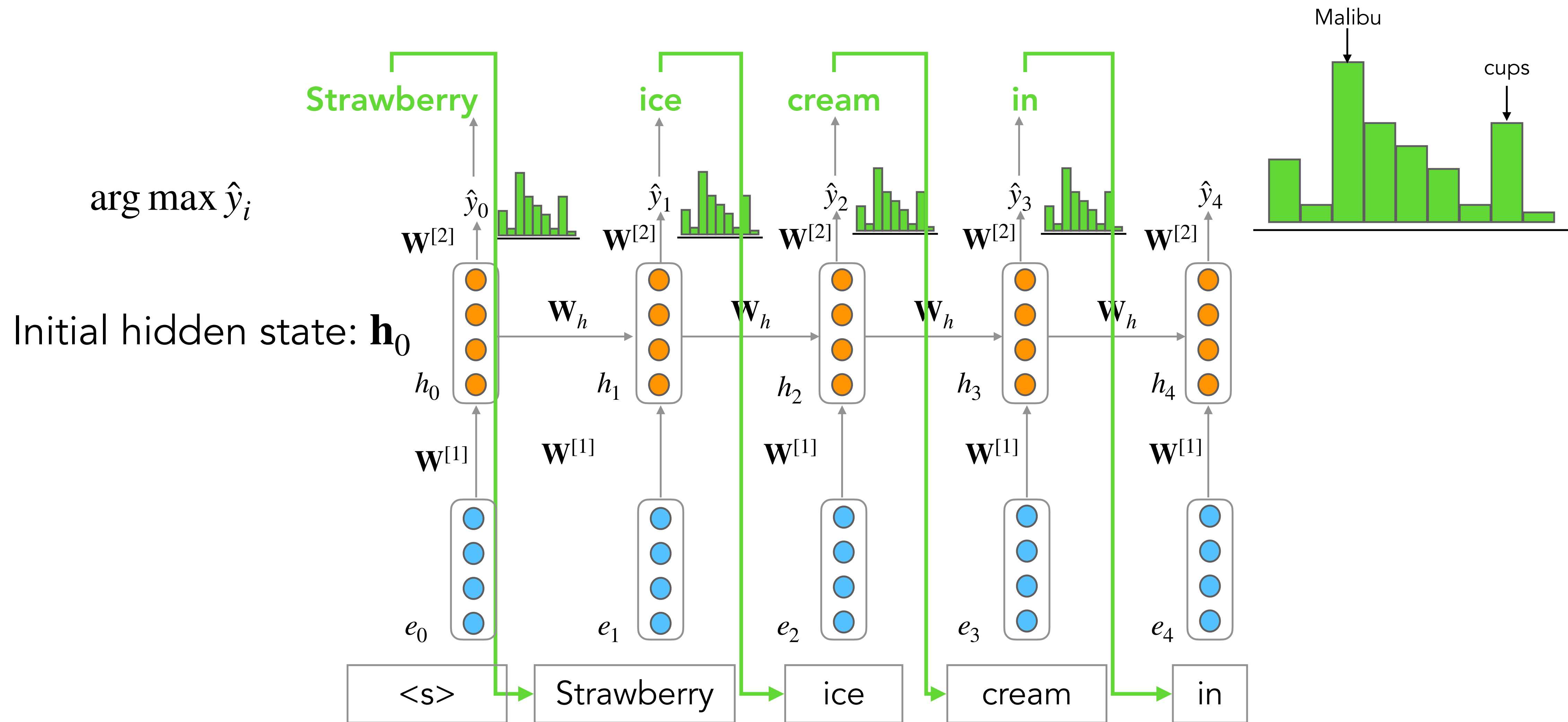
- Quiz 3: FFNNs and RNNs
- Recap: Seq2Seq and Attention
- More on Attention
- Transformers: Self-Attention
- Transformers: Multi-headed Attention
- Transformers: Positional Embeddings
- Putting it all together: Transformer Blocks

Quiz 3!

Recap: Seq2Seq and Attention

Generation with RNNLMs

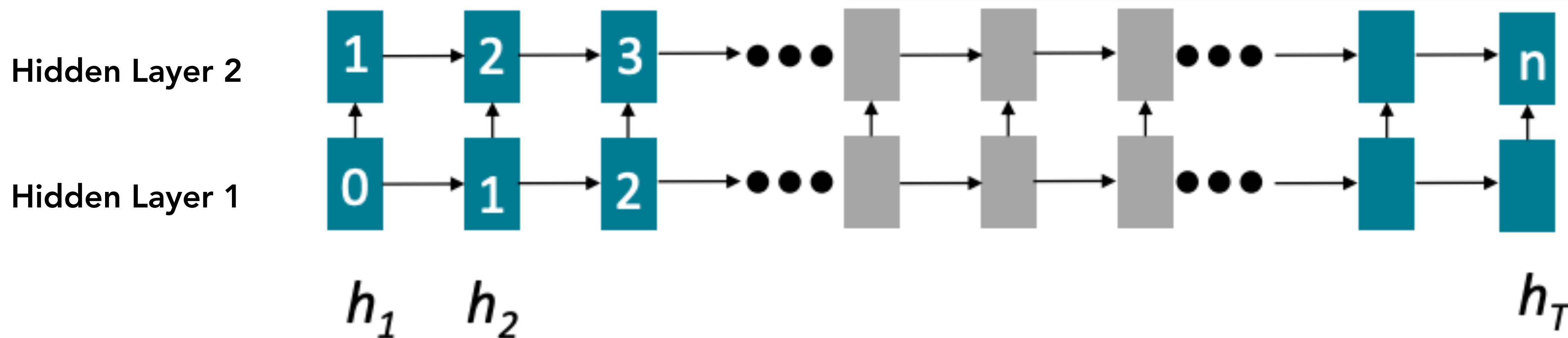
$$\hat{y}_4 = P(x_5 | \text{Strawberry ice cream in})$$



RNNs and parallelizability

- Forward and backward passes have **$O(\text{sequence length})$** unparallelizable operations!
 - Future RNN hidden states can't be computed in full before past RNN hidden states have been computed
 - But GPUs can perform a bunch of independent computations at once!

Inhibits training on very large datasets!

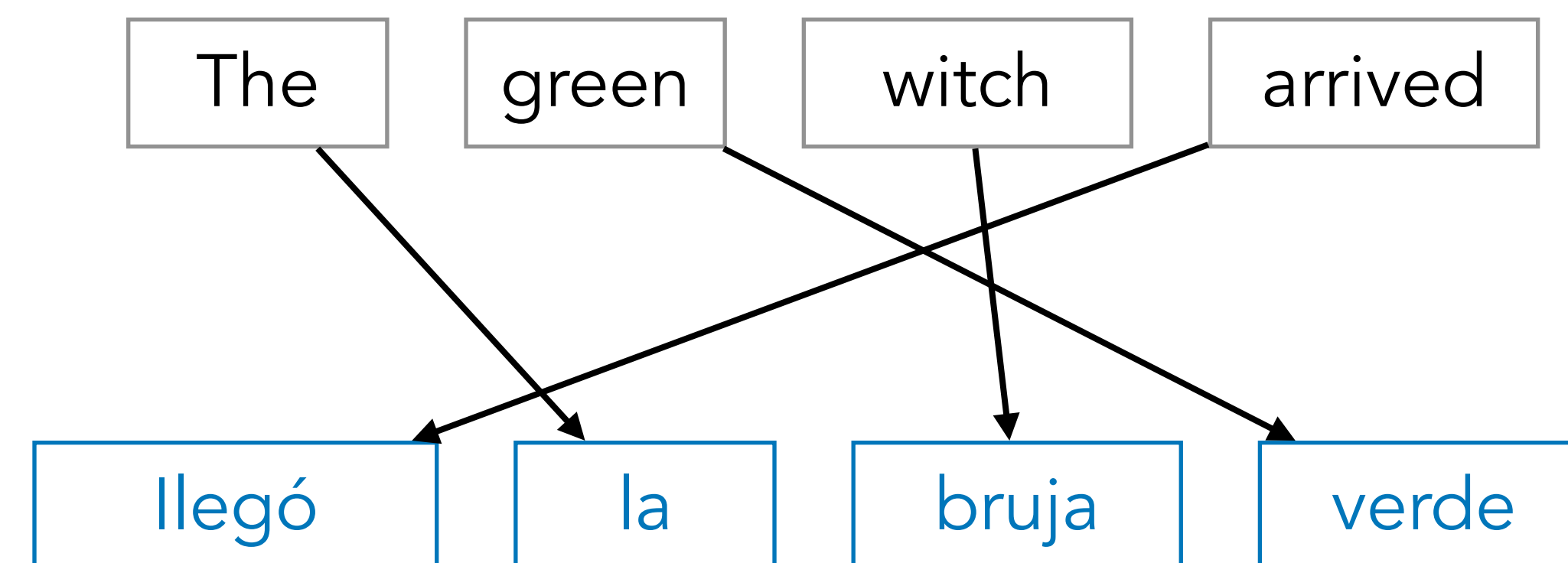


Numbers indicate min # of steps before a state can be computed

(Neural) Machine Translation

Seq2Seq uses rich, task-appropriate context!

- Sequence Generation Problem (as opposed to sequence classification)
 - \mathbf{x} = Source sequence of length n
 - \mathbf{y} = Target sequence of length m
- Different from regular generation from an LM
 - Since we expect the target sequence to serve a specific utility (translate the source)



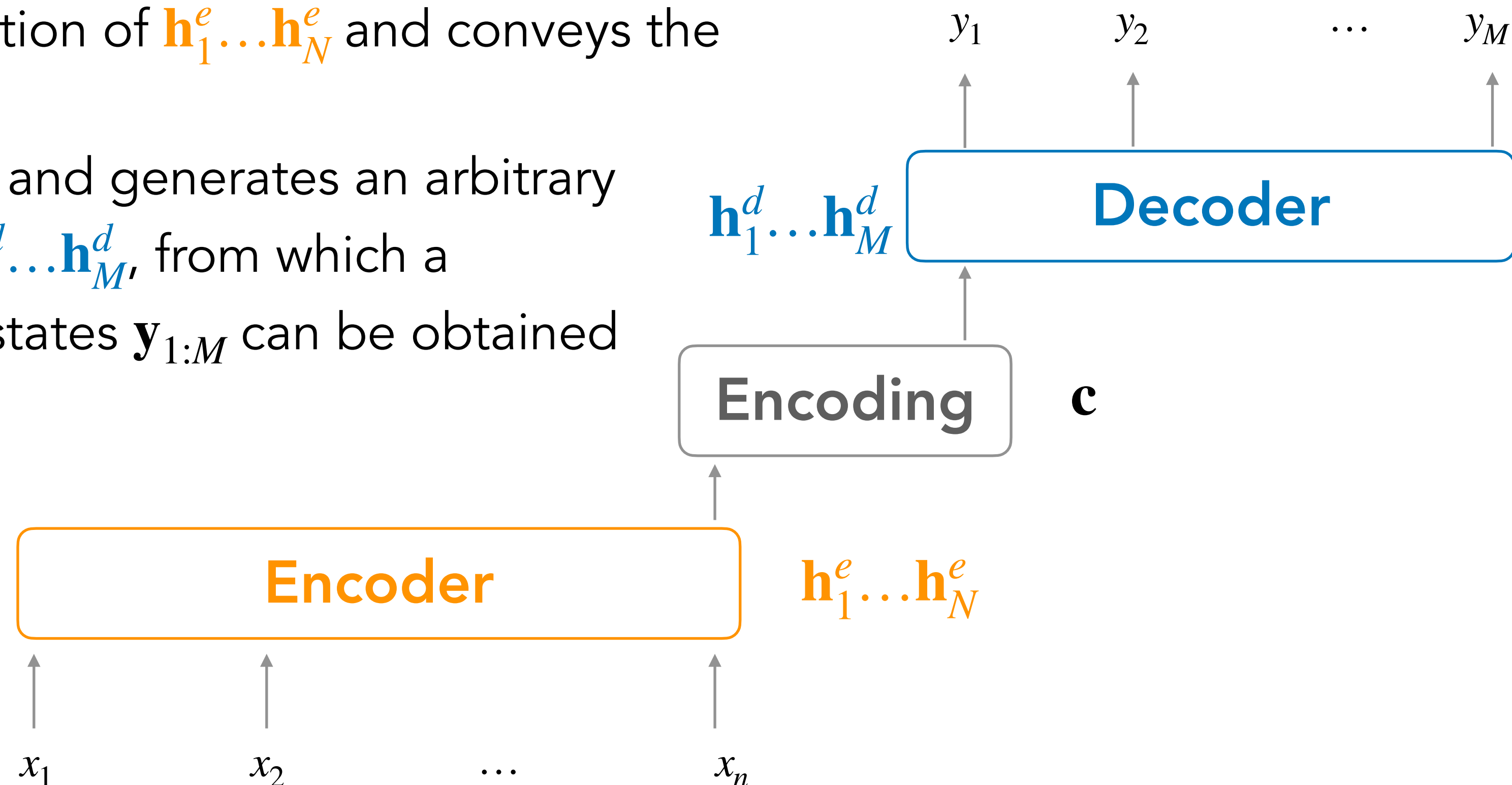
Sequence-to-Sequence (Seq2seq)

Encoder-Decoder Networks

Encoder-decoder networks consist of three components:

1. An **encoder** that accepts an input sequence, $\mathbf{x}_{1:N}$ and generates a corresponding sequence of contextualized representations, $\mathbf{h}_1^e \dots \mathbf{h}_N^e$
2. A **encoding** vector, \mathbf{c} which is a function of $\mathbf{h}_1^e \dots \mathbf{h}_N^e$ and conveys the essence of the input to the decoder
3. A **decoder** which accepts \mathbf{c} as input and generates an arbitrary length sequence of hidden states $\mathbf{h}_1^d \dots \mathbf{h}_M^d$, from which a corresponding sequence of output states $\mathbf{y}_{1:M}$ can be obtained

Encoders and decoders can be made of FFNNs, RNNs, or Transformers



Produces an encoding of the source sequence

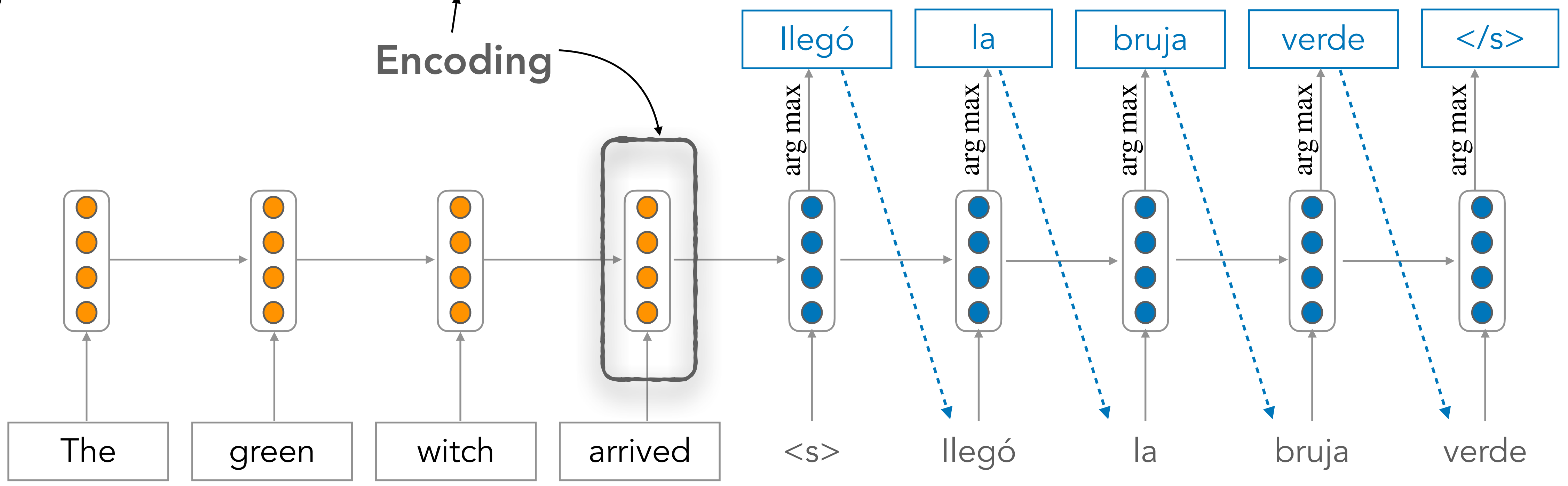
Represents input sequence. Provides initial hidden state for Decoder RNN

Encoding

Target Sentence y

Encoder RNN

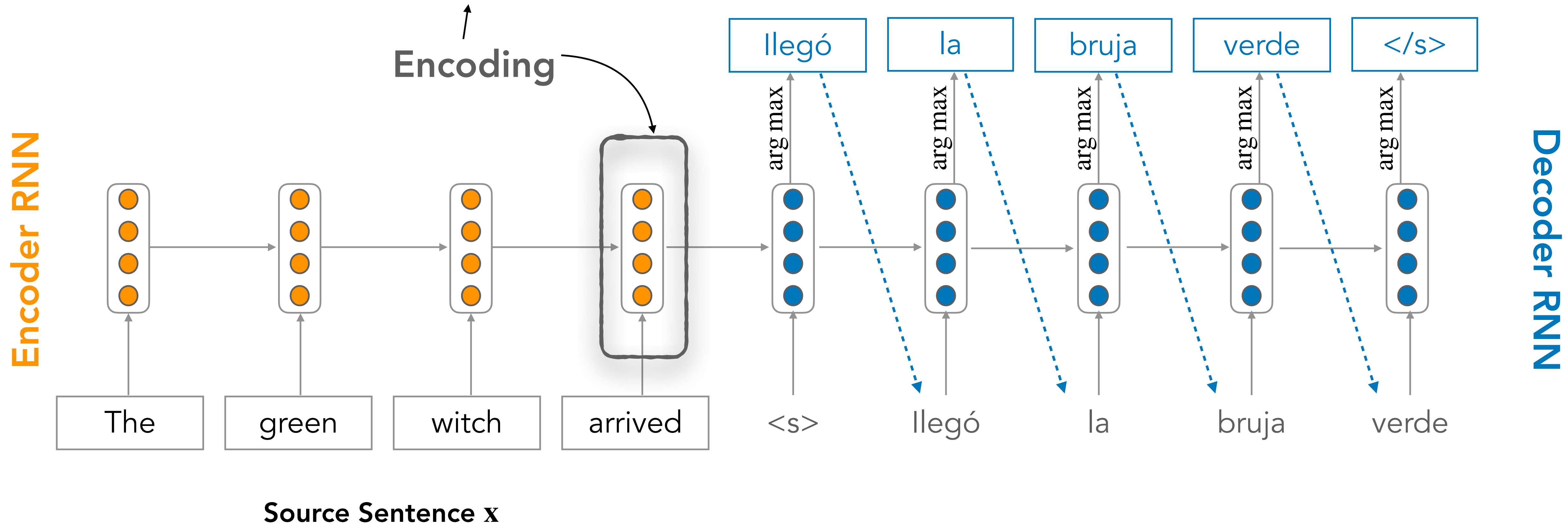
Decoder RNN



Source Sentence x

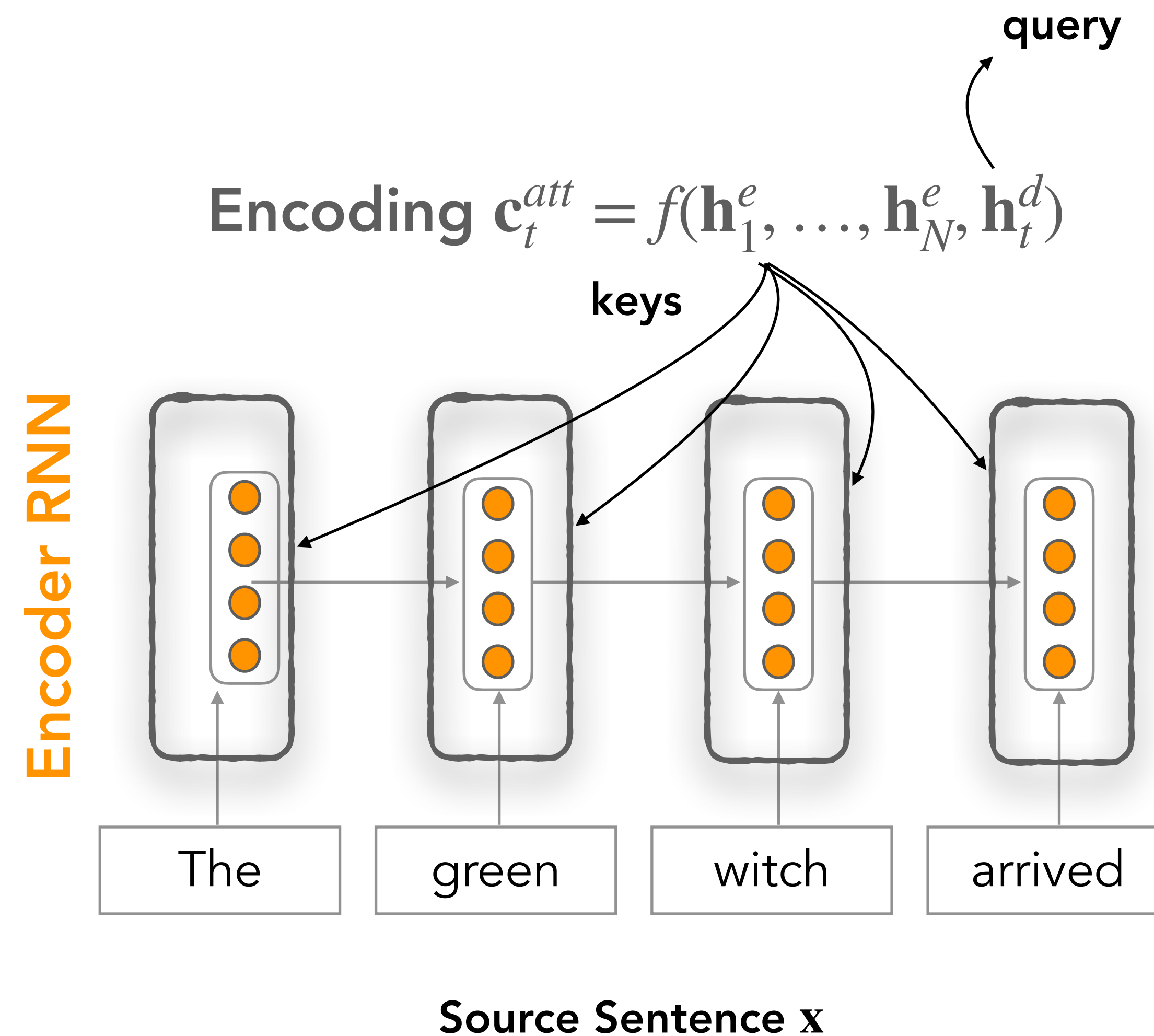
Language Model that produces the target sentence conditioned on the encoding

This needs to capture all information about the source sentence. Information bottleneck!



Attention Mechanism

- Attention mechanisms allow the decoder to focus on a particular part of the source sequence at each time step
- Fixed-length vector \mathbf{c}_t^{att} (attention context vector)
 - Take a weighted sum of all the encoder hidden states
 - One vector per time step *of the decoder!*
 - Weights *attend* to part of the source text relevant for the token the decoder is producing at step t
- In general, we have a single **query** vector and multiple **key** vectors.
 - We want to score each query-key pair



Bahdanau et al., 2015

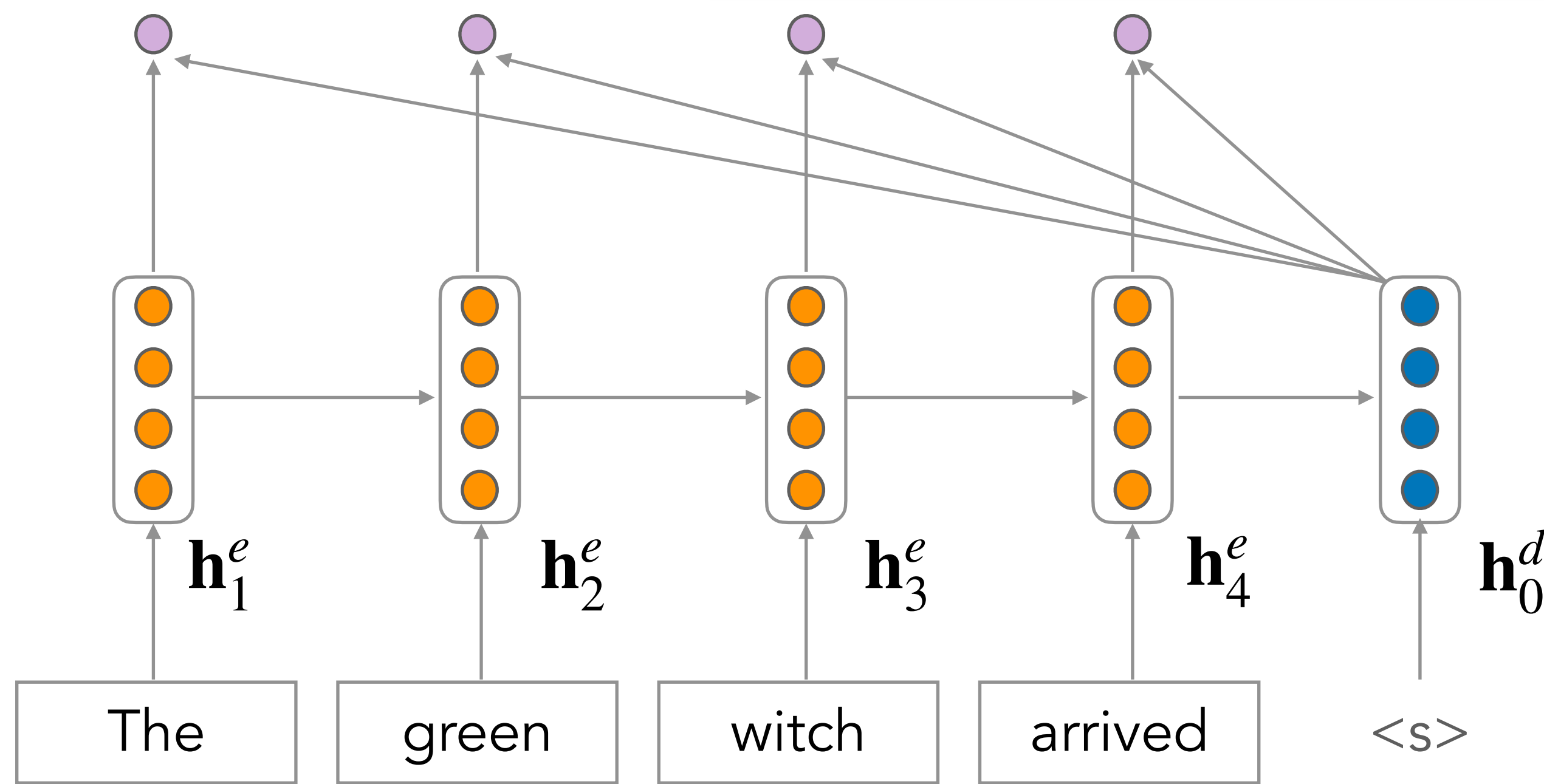
Note: Notation different from J&M

Seq2Seq with Attention

Encoder RNN
Attention Scores / Attention Logits

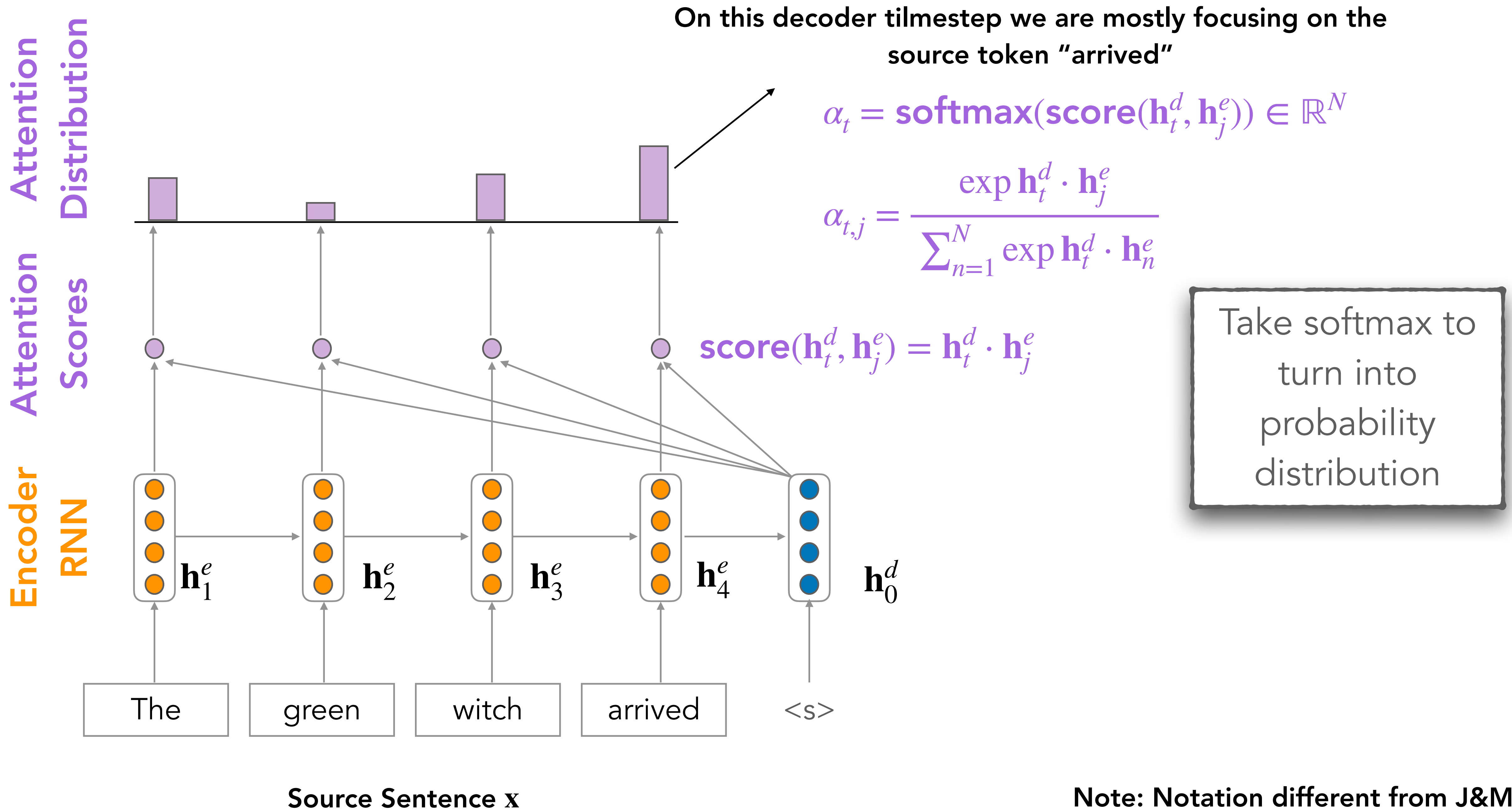
$$\text{score}(\mathbf{h}_t^d, \mathbf{h}_j^e) = \mathbf{h}_t^d \cdot \mathbf{h}_j^e$$

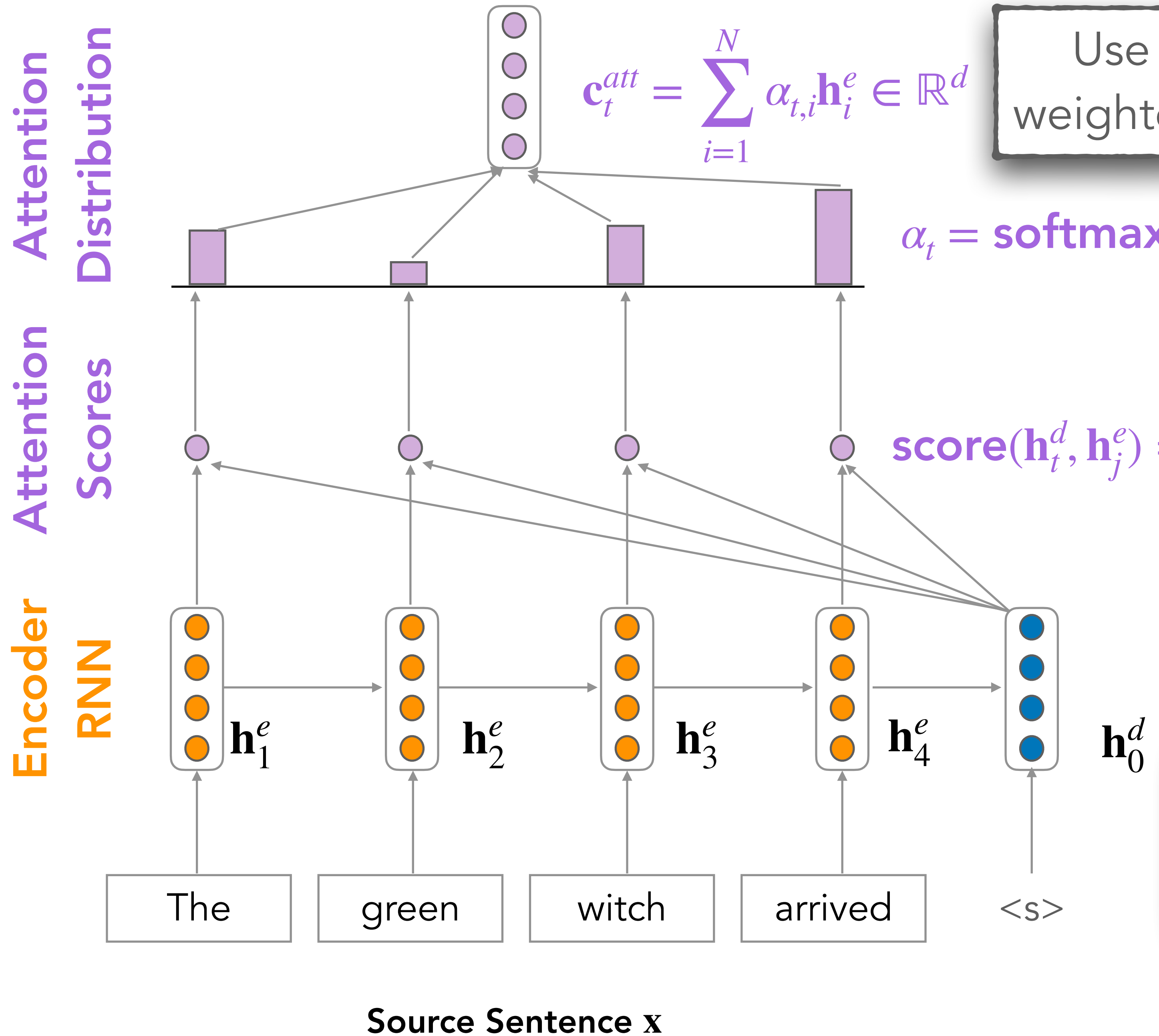
Dot product with keys (encoder hidden states) to encode similarity with what is decoded so far...



Query 1: Decoder, first time step

Dot product attention





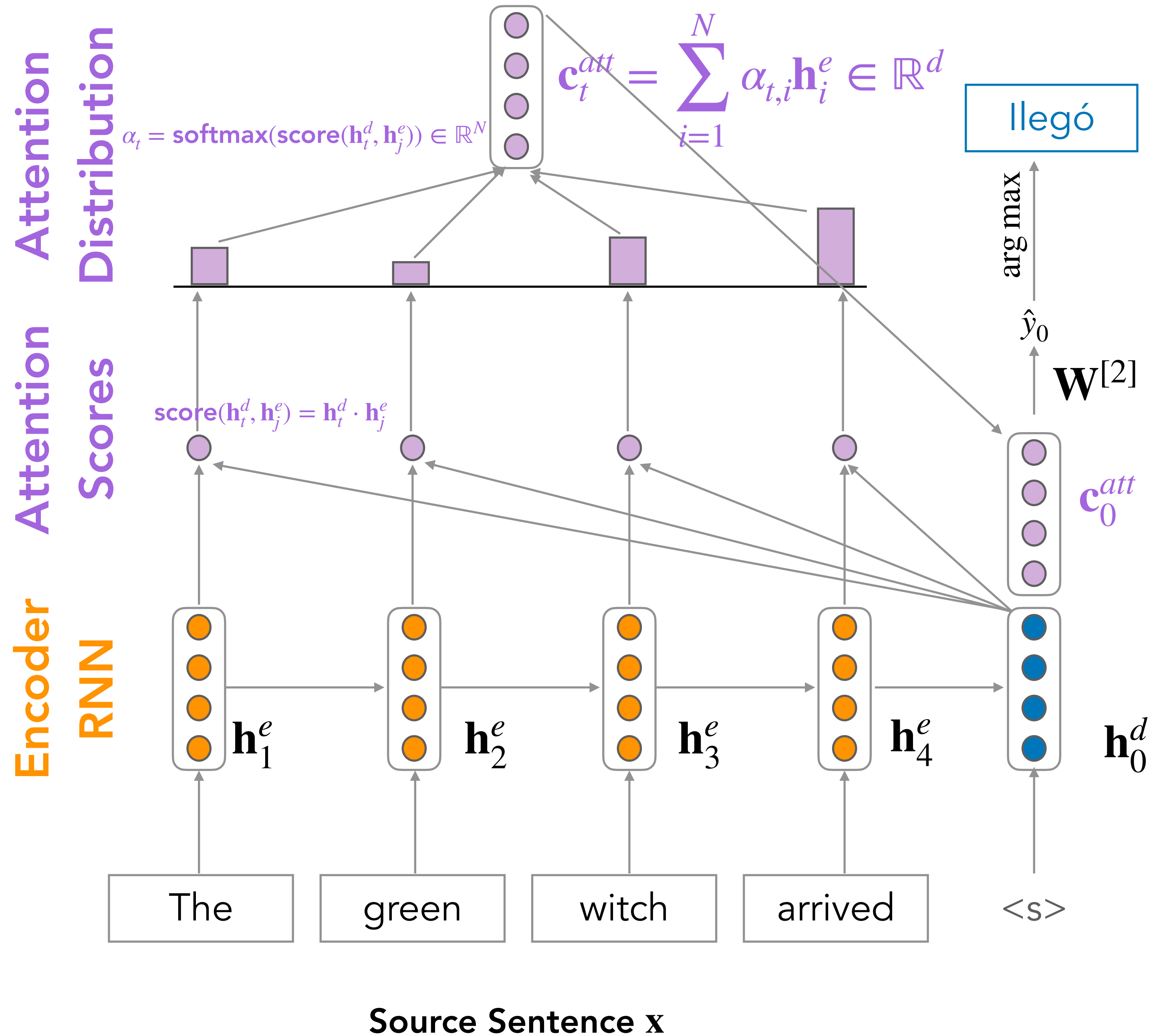
Use the attention distribution to take a weighted sum of the encoder hidden states.

$$\alpha_t = \text{softmax}(\text{score}(\mathbf{h}_t^d, \mathbf{h}_j^e)) \in \mathbb{R}^N$$

$$\text{score}(\mathbf{h}_t^d, \mathbf{h}_j^e) = \mathbf{h}_t^d \cdot \mathbf{h}_j^e$$

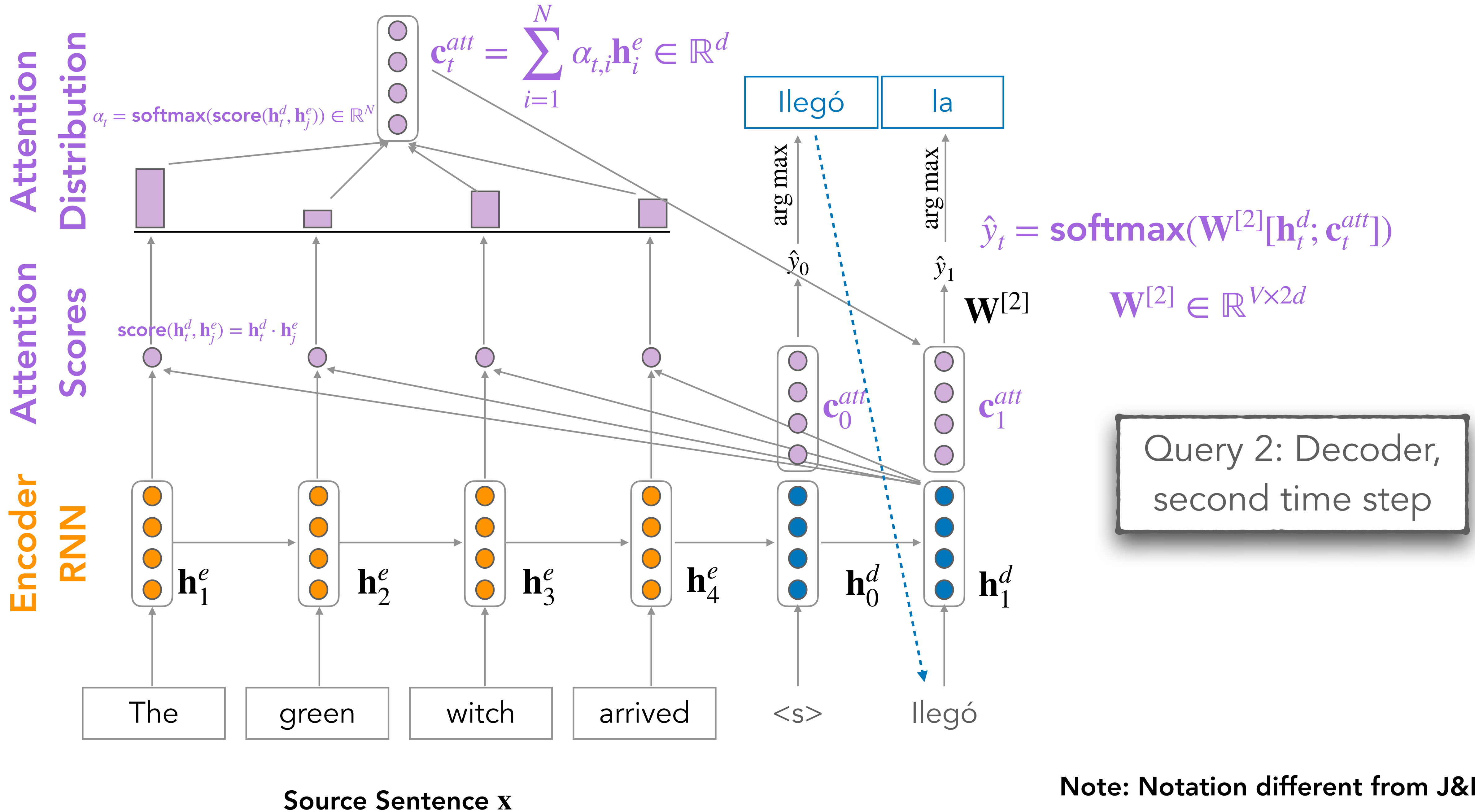
The attention output mostly contains information the hidden states that received high attention.

Note: Notation different from J&M



Concatenate attention output with decoder hidden state, then use to compute \hat{y}_0 as before

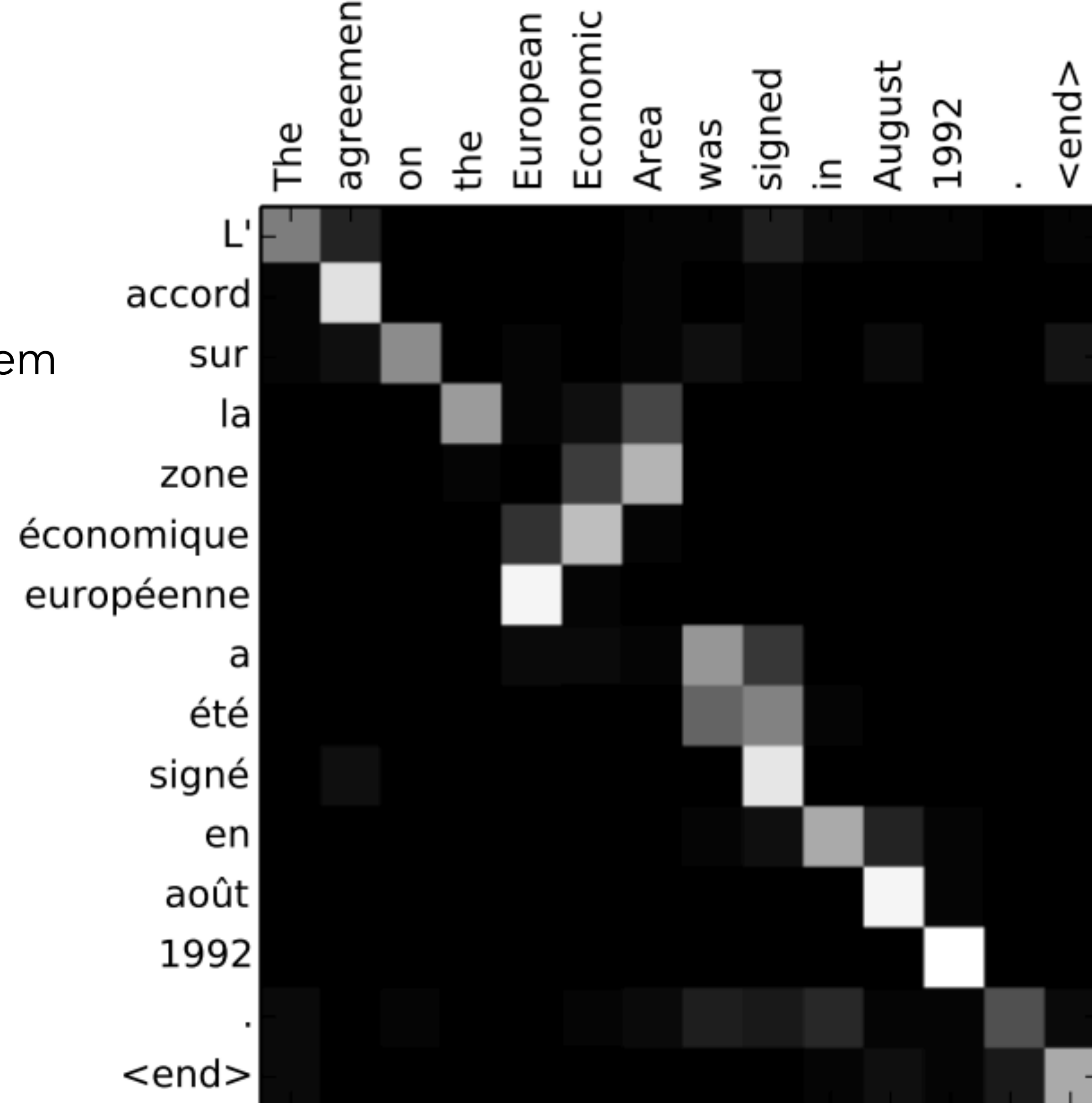
Note: Notation different from J&M



Note: Notation different from J&M

Why Attention?

- Attention significantly **improves** neural machine translation **performance**
 - Very useful to allow decoder to focus on certain parts of the source
- Attention **solves the information bottleneck** problem
 - Attention allows decoder to look directly at source; bypass bottleneck
- Attention **helps with vanishing gradient problem**
 - Provides shortcut to faraway states
- Attention provides some **interpretability**
 - By inspecting attention distribution, we can see what the decoder was focusing on →
 - We get alignment for free! We never explicitly trained an alignment system! The network just learned alignment by itself




More on Attention

Attention Variants

- In general, we have some keys $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$ and a query $\mathbf{q} \in \mathbb{R}^{d_2}$

- Attention always involves

1. Computing the attention scores, $e(\mathbf{q}, \mathbf{h}_{1:N}) \in \mathbb{R}^N$ 
2. Taking softmax to get attention distribution $\alpha_t = \text{softmax}(e(\mathbf{q}, \mathbf{h}_{1:N})) \in [0, 1]^N$
3. Using attention distribution to take weighted sum of values:

Can be done in multiple ways!

$$\mathbf{c}_t^{att} = \sum_{i=1}^N \alpha_{t,i} \mathbf{h}_i \in \mathbb{R}^{d_1}$$

This leads to the attention output \mathbf{c}_t^{att} (sometimes called the attention context vector)

Attention Variants

- There are several ways you can compute $e(\mathbf{q}, \mathbf{h}_{1:N}) \in \mathbb{R}^N$ from $\mathbf{h}_1 \dots \mathbf{h}_N \in \mathbb{R}^{d_1}$ and $\mathbf{q} \in \mathbb{R}^{d_2}$
- Basic dot-product attention: $e(\mathbf{q}, \mathbf{h}_{1:N}) = [\mathbf{q} \cdot \mathbf{h}_j]_{j=1:N}$
 - This assumes $d_1 = d_2$
 - We applied this in encoder-decoder RNNs
- Multiplicative attention: $e(\mathbf{q}, \mathbf{h}_{1:N}) = [\mathbf{q}^T \mathbf{W} \mathbf{h}_j]_{j=1:N}$
 - Where $\mathbf{W} \in \mathbb{R}^{d_2 \times d_1}$ is a learned weight matrix.
 - Also called “bilinear attention”

More on Attention

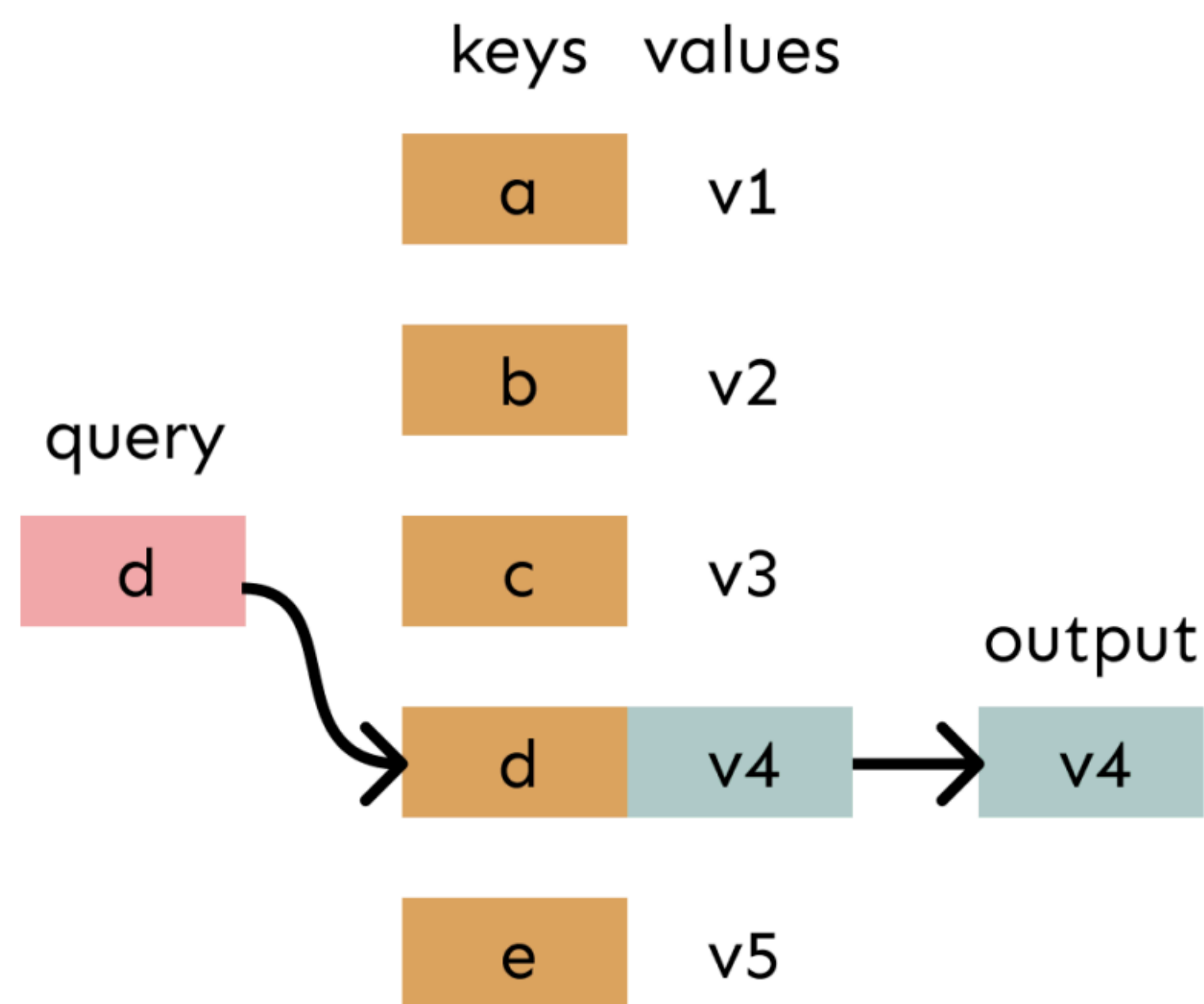
Given a set of vector values, and a vector query, attention is a technique to compute a weighted sum of the values, dependent on the query

- We sometimes say that the query attends to the values.
 - For example, in the seq2seq + attention model, each decoder hidden state (query) attends to all the encoder hidden states (values)
 - Here, keys and values are the same!
- The weighted sum is a **selective summary** of the information contained in the values, where the query determines which values to focus on.
- Attention is a way to obtain a **fixed-size representation** of an **arbitrary set of representations** (the values), dependent on some other representation (the query).
- Attention is a powerful, flexible, general deep learning technique in all deep learning models.
 - A new idea from after 2010! Originated in NMT

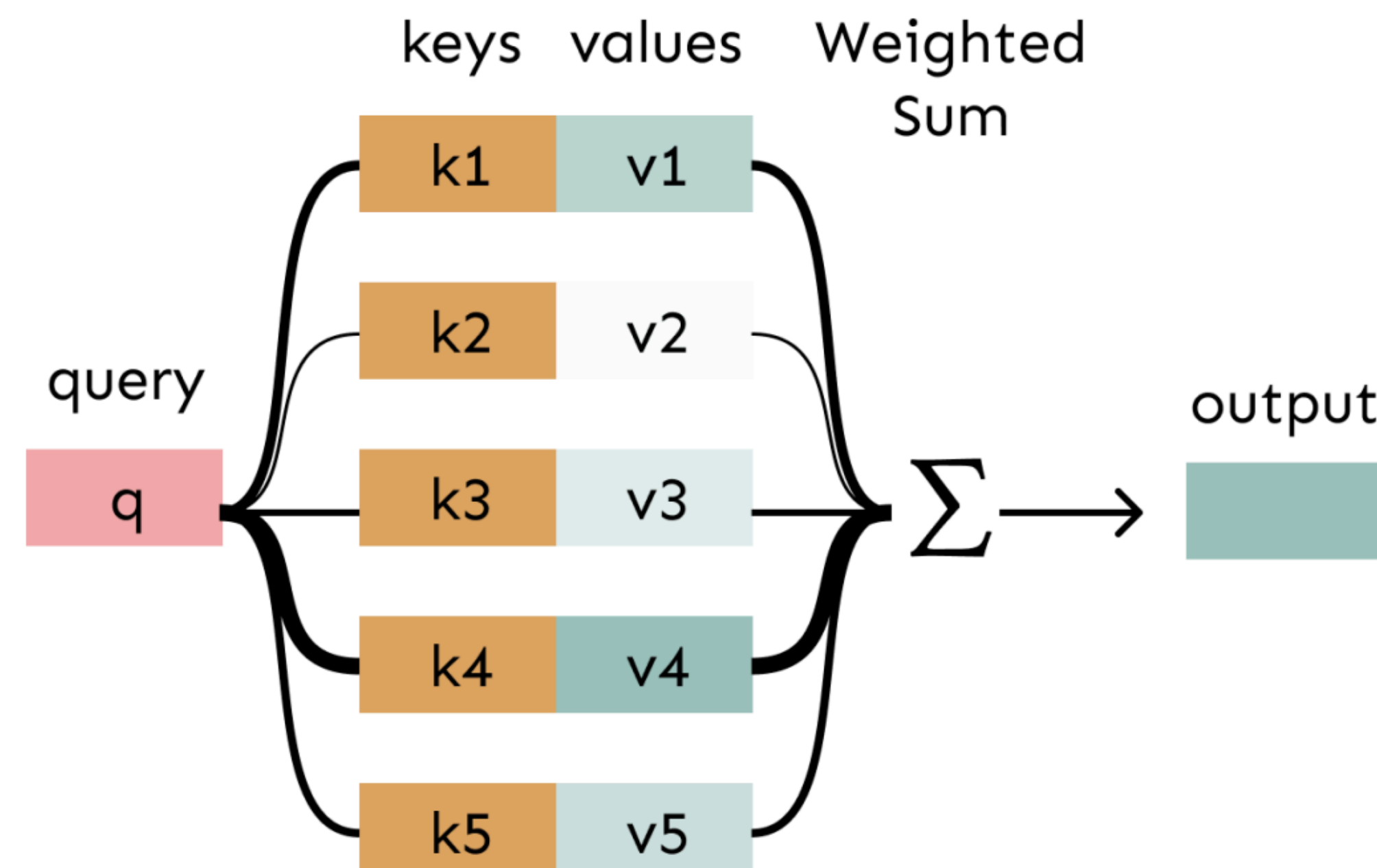
Attention and lookup tables

Attention performs fuzzy lookup in a key-value store

In a lookup table, we have a table of keys that map to values. The query matches one of the keys, returning its value.

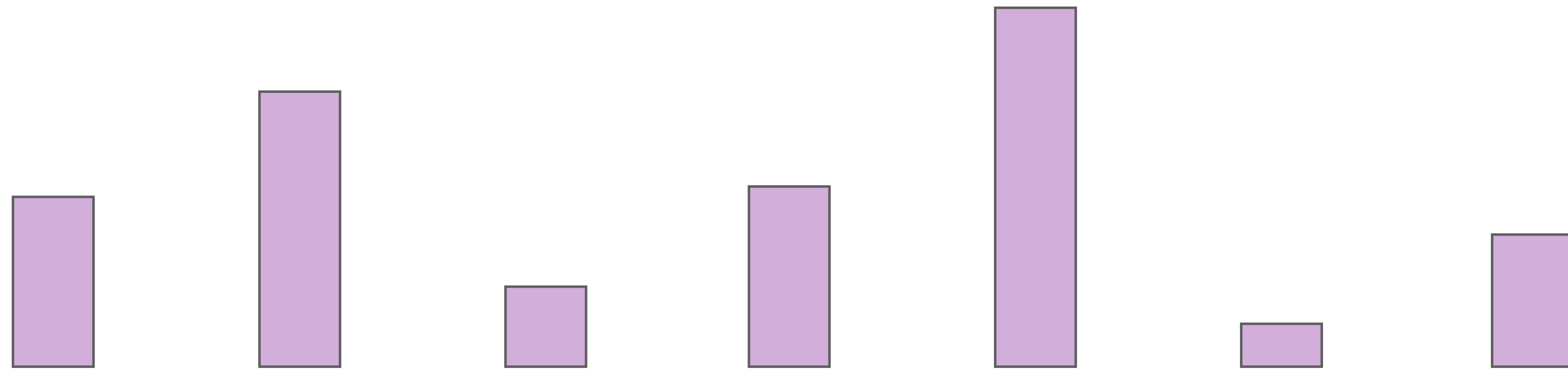


In attention, the query matches all keys softly, to a weight between 0 and 1. The keys' values are multiplied by the weights and summed.

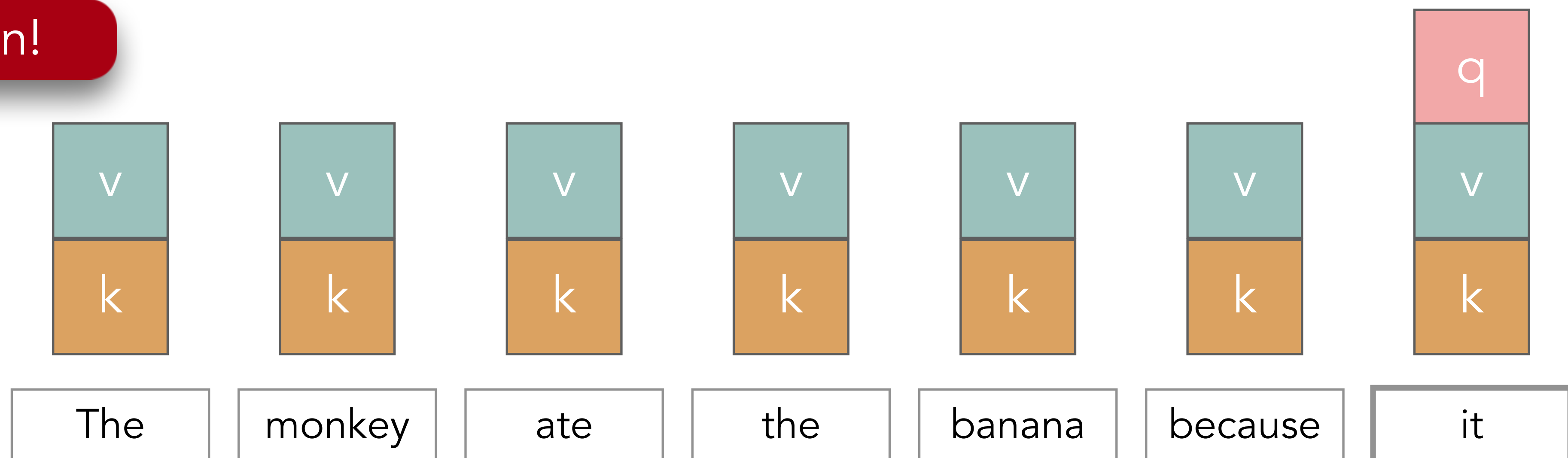


Attention in the decoder

Attention
Distribution



Self-Attention!

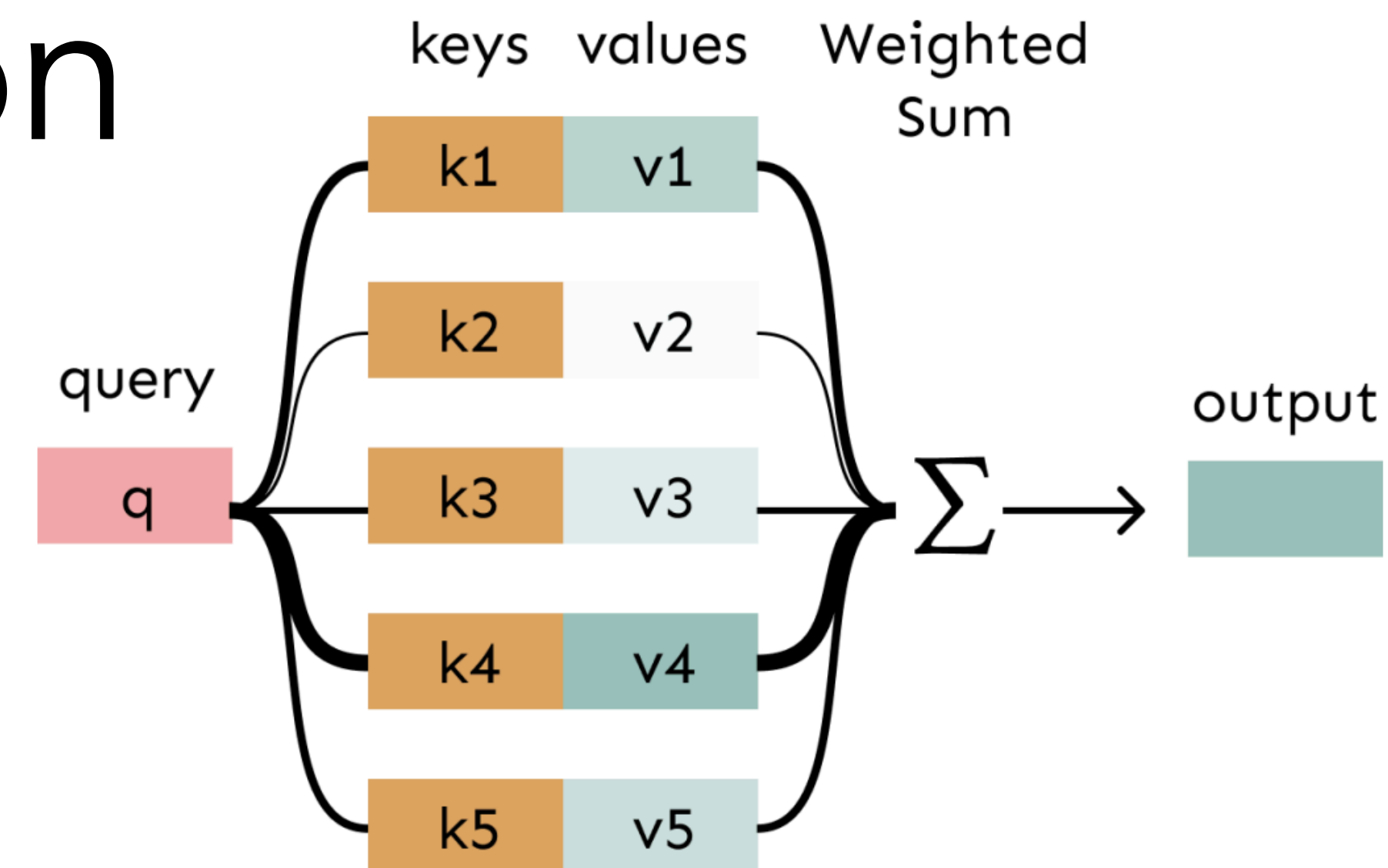


Transformers: Self-Attention Networks

Self-Attention

Keys, Queries, Values from the same sequence

Let $\mathbf{w}_{1:N}$ be a sequence of words in vocabulary V
 For each \mathbf{w}_i , let $\mathbf{x}_i = \mathbf{E}_{\mathbf{w}_i}$ where $\mathbf{E} \in \mathbb{R}^{d \times V}$ is an embedding matrix.



1. Transform each word embedding with weight matrices \mathbf{Q} , \mathbf{K} , \mathbf{V} , each in $\mathbb{R}^{d \times d}$

$$\mathbf{q}_i = \mathbf{Q}\mathbf{x}_i \text{ (queries)}$$

$$\mathbf{k}_i = \mathbf{K}\mathbf{x}_i \text{ (keys)}$$

$$\mathbf{v}_i = \mathbf{V}\mathbf{x}_i \text{ (values)}$$

2. Compute pairwise similarities between keys and queries; normalize with softmax

$$\mathbf{e}_{ij} = \mathbf{q}_i^\top \mathbf{k}_j \quad \alpha_{ij} = \frac{\exp(\mathbf{e}_{ij})}{\sum_{j'} \exp(\mathbf{e}_{ij'})}$$

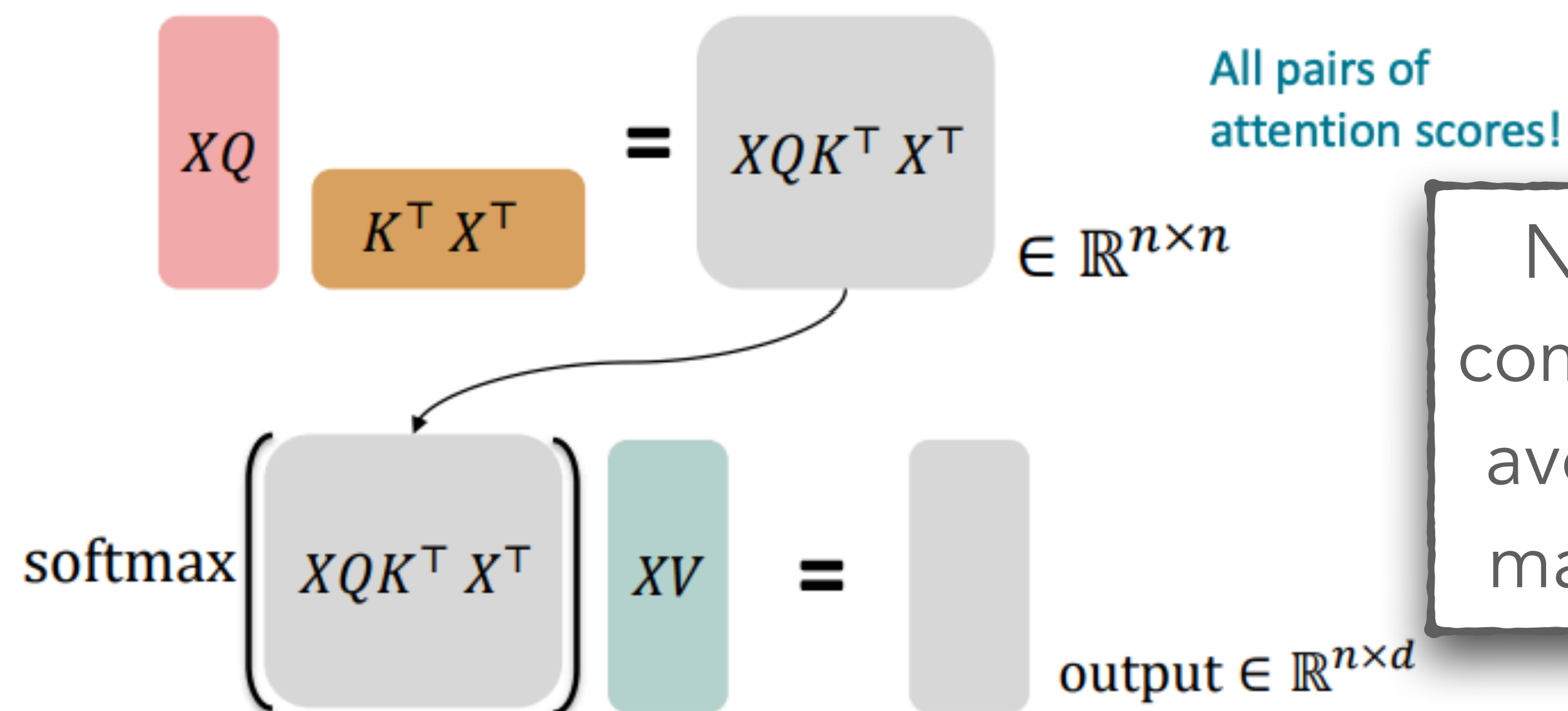
3. Compute output for each word as weighted sum of values

$$\mathbf{o}_i = \sum_j \alpha_{ij} \mathbf{v}_j$$

Self-Attention as Matrix Multiplications

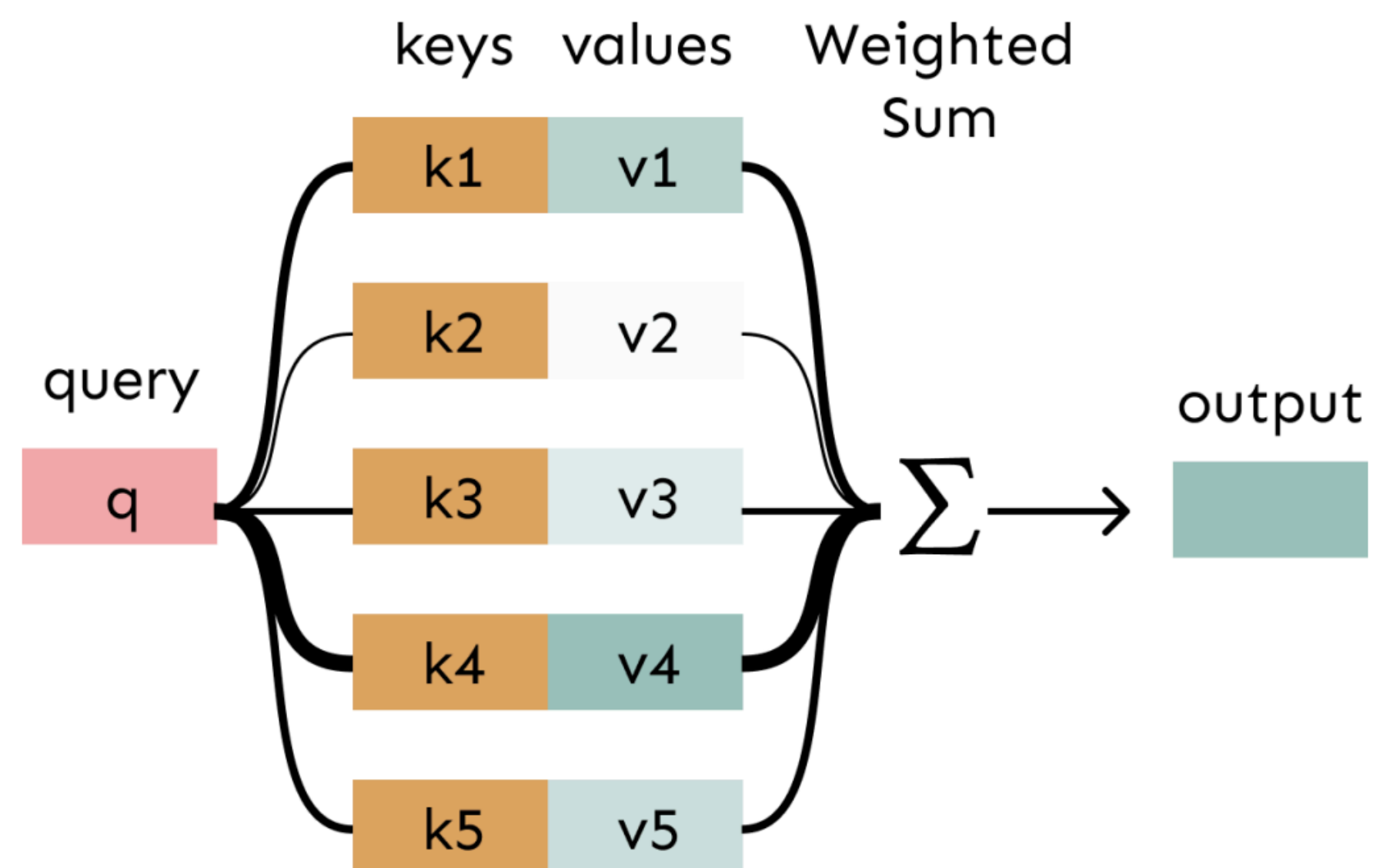
- Key-query-value attention is typically computed as matrices.
 - Let $\mathbf{X} = [\mathbf{x}_1; \dots; \mathbf{x}_n] \in \mathbb{R}^{n \times d}$ be the concatenation of input vectors
 - First, note that $\mathbf{X}\mathbf{K} \in \mathbb{R}^{n \times d}$, $\mathbf{X}\mathbf{Q} \in \mathbb{R}^{n \times d}$, and $\mathbf{X}\mathbf{V} \in \mathbb{R}^{n \times d}$
 - The output is defined as $\text{softmax}(\mathbf{X}\mathbf{Q}(\mathbf{X}\mathbf{K})^T)\mathbf{X}\mathbf{V} \in \mathbb{R}^{n \times d}$

First, take the query-key dot products in one matrix multiplication:

$$\mathbf{X}\mathbf{Q}(\mathbf{X}\mathbf{K})^T$$


Next, softmax, and compute the weighted average with another matrix multiplication.

Why Self-Attention?



- Self-attention allows a network to directly extract and use information from arbitrarily large contexts without the need to pass it through intermediate recurrent connections as in RNNs
- Used often with feedforward networks!

Transformers are Self-Attention Networks

- Self-Attention is the key innovation behind Transformers!
- Transformers map sequences of input vectors $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ to sequences of output vectors $(\mathbf{y}_1, \dots, \mathbf{y}_n)$ of the same length.
- Made up of stacks of Transformer blocks
 - each of which is a multilayer network made by combining
 - simple linear layers,
 - feedforward networks, and
 - self-attention layers



Self-Attention and Weighted Averages

- **Problem:** there are no element-wise nonlinearities in self-attention; stacking more self-attention layers just re-averages value vectors
- **Solution:** add a feed-forward network to post-process each output vector.

