# Attention-based architectures

INTELLIGENT SYSTEMS FOR PATTERN RECOGNITION (ISPR)

DAVIDE BACCIU – DIPARTIMENTO DI INFORMATICA - UNIVERSITA' DI PISA
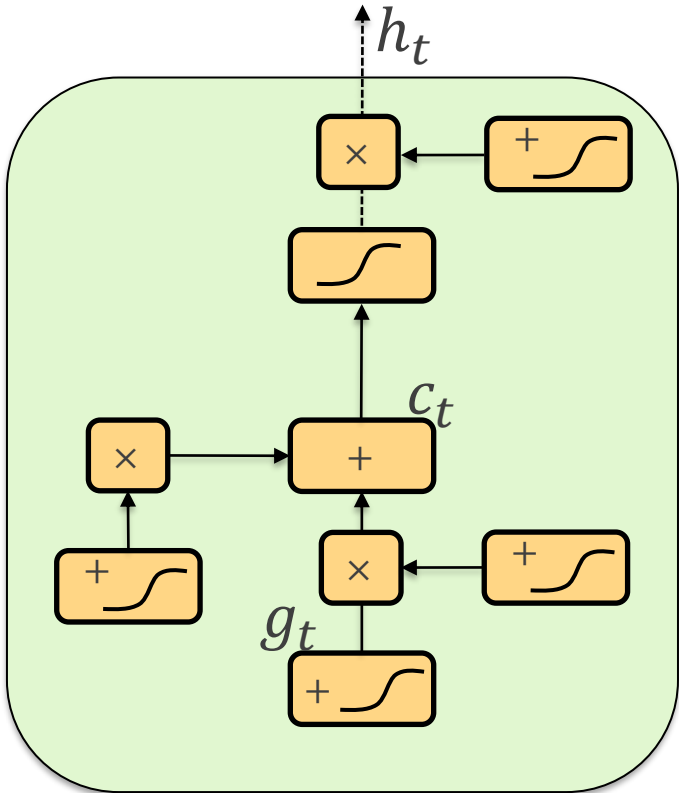
DAVIDE.BACCIU@UNIPI.IT

# A 2 Lectures Outline

- L22 - Neural attention for structured/compound data

  - Sequence-to-sequence
  - Attention models

- L23 - Dealing with very long-term dependencies

  - Multiscale networks
  - Neural memories (more attention)
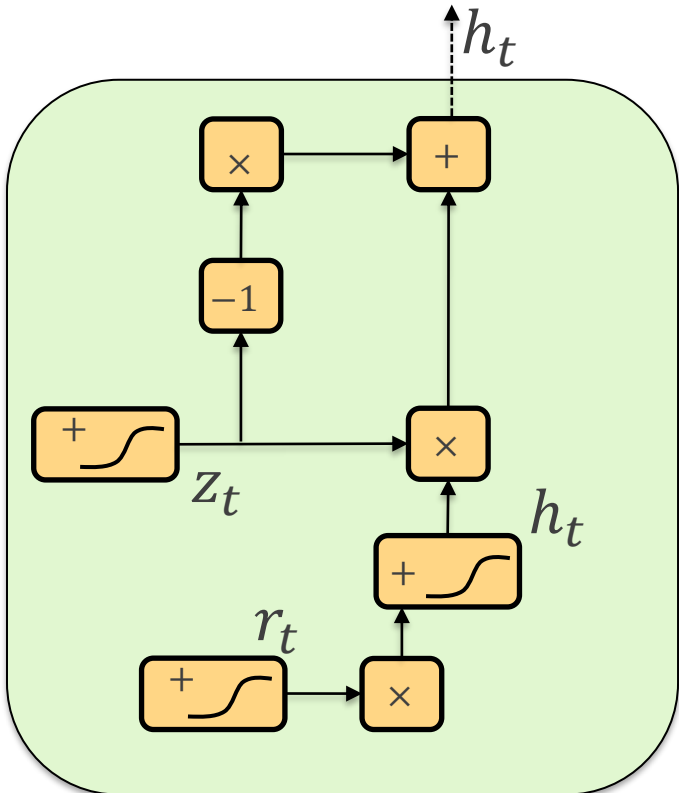  - Differentiable memory read, write, indexing

Extra Lecture
Tomorrow 12/04/2024
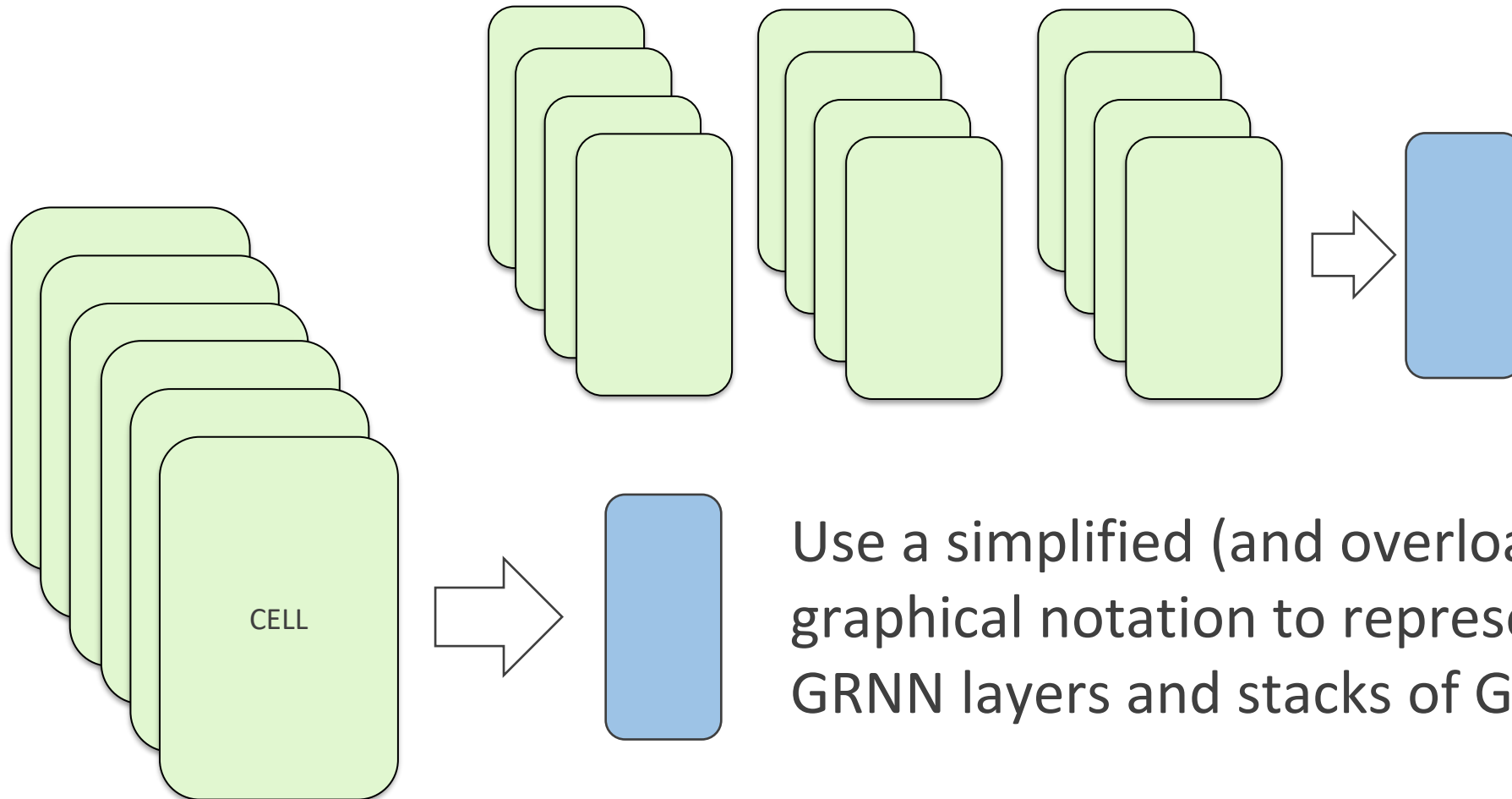– h16 – Aula E

# Gated RNN Refresher



LSTM Cell

GRU Cell

# Graphical Notation for Compositionality

CELL

Use a simplified (and overloaded) graphical notation to represent GRNN layers and stacks of GRNN

# Dealining with Compound Data

○ GRNN are excellent to handle size/topology varying data in input

- How can we handle size/topology varying outputs?

- Sequence-to-sequence

○ Structured data is compound information

- Efficient processing needs the ability to focus on certain parts of such information

- Attention mechanism
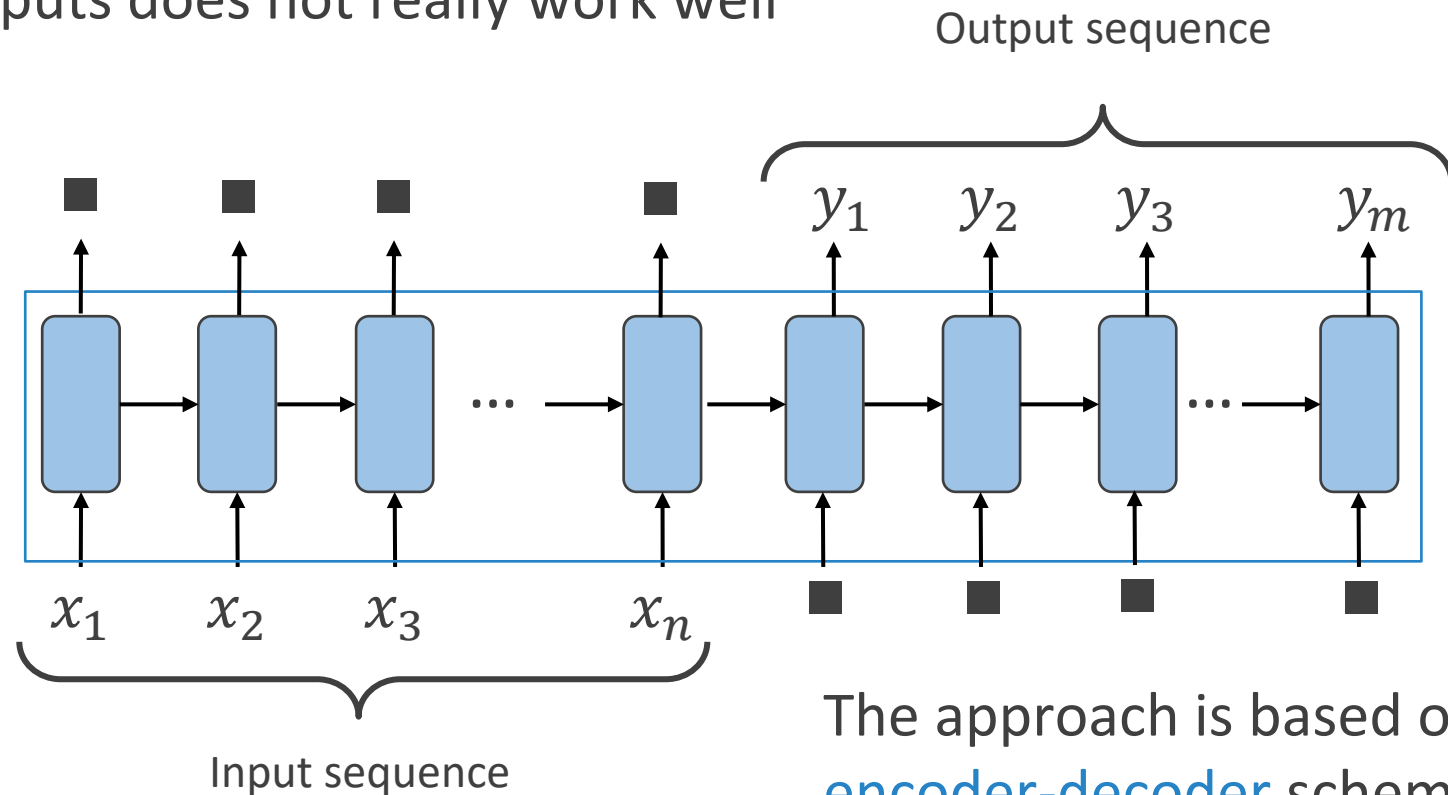
# Sequence-to-sequence

# Sequence Transduction

o Input and output are both sequences

o They may have different lengths

o Example: machine translation

**The cat is on the table** ⟶ **Il gatto è sul tavolo**

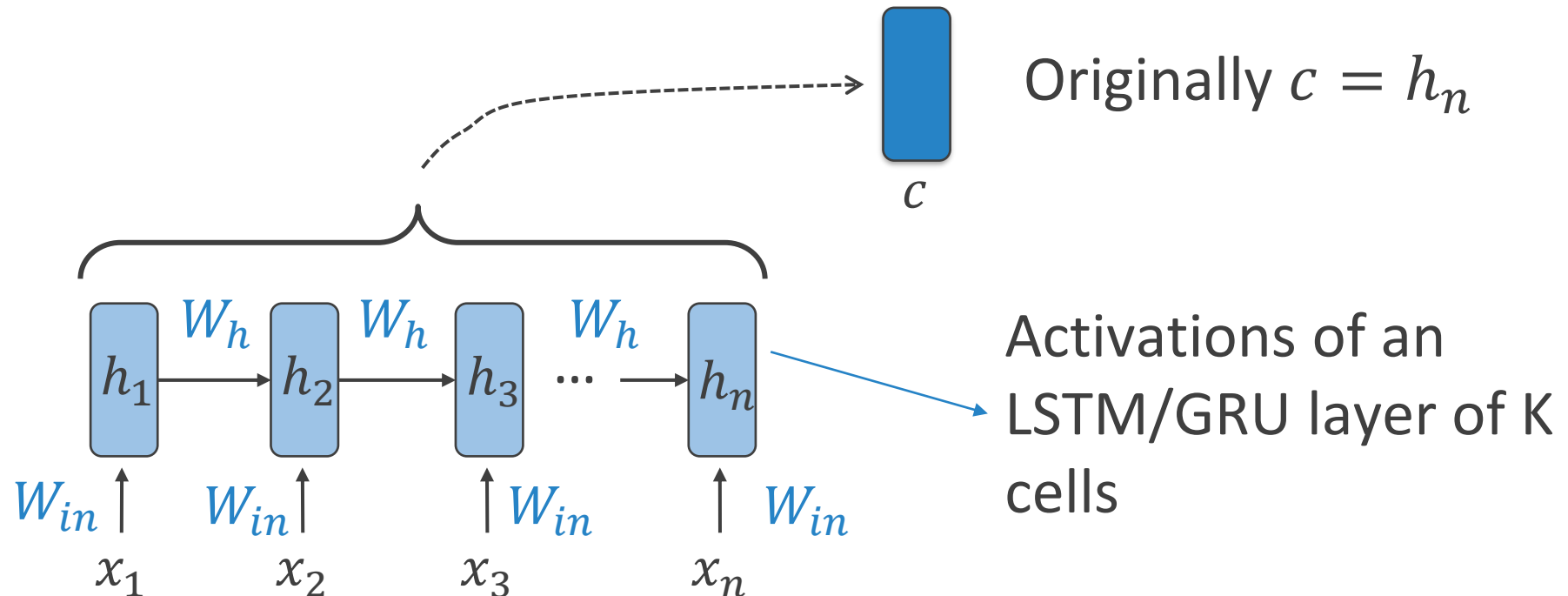How do we model the context here?

# Learning to Output Variable Length Sequences

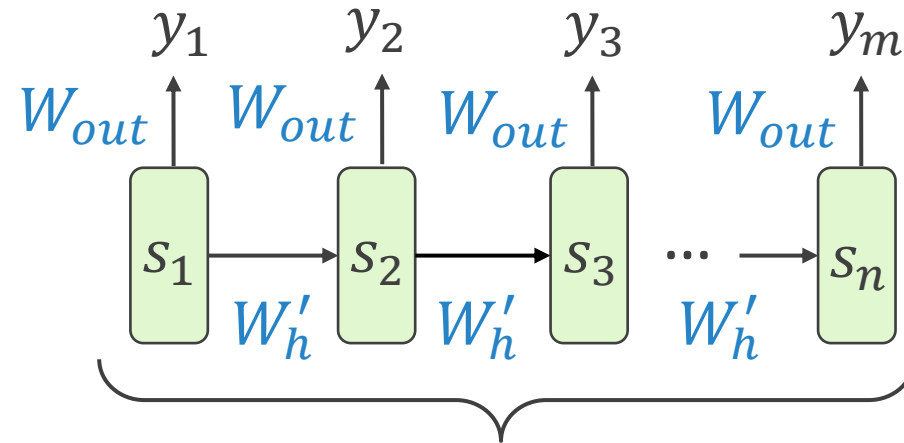The idea of an unfolded RNN with blank inputs-outputs does not really work well



Output sequence

Input sequence

The approach is based on an encoder-decoder scheme

# Encoder

Produce a compressed and fixed length representation $c$ of all the input sequence $x_1, \ldots, x_n$



Originally $c = h_n$

Activations of an LSTM/GRU layer of K cells
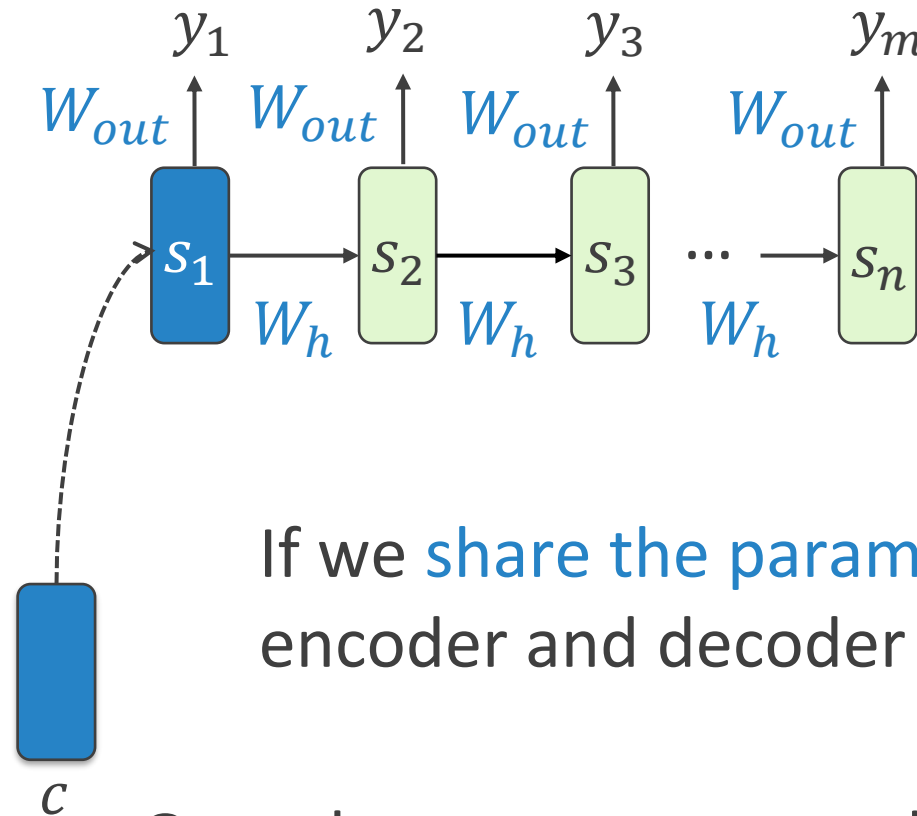
# Decoder



A LSTM/GRU layer of K cells seeded by the context vector $c$

Different approaches to realize this in practice
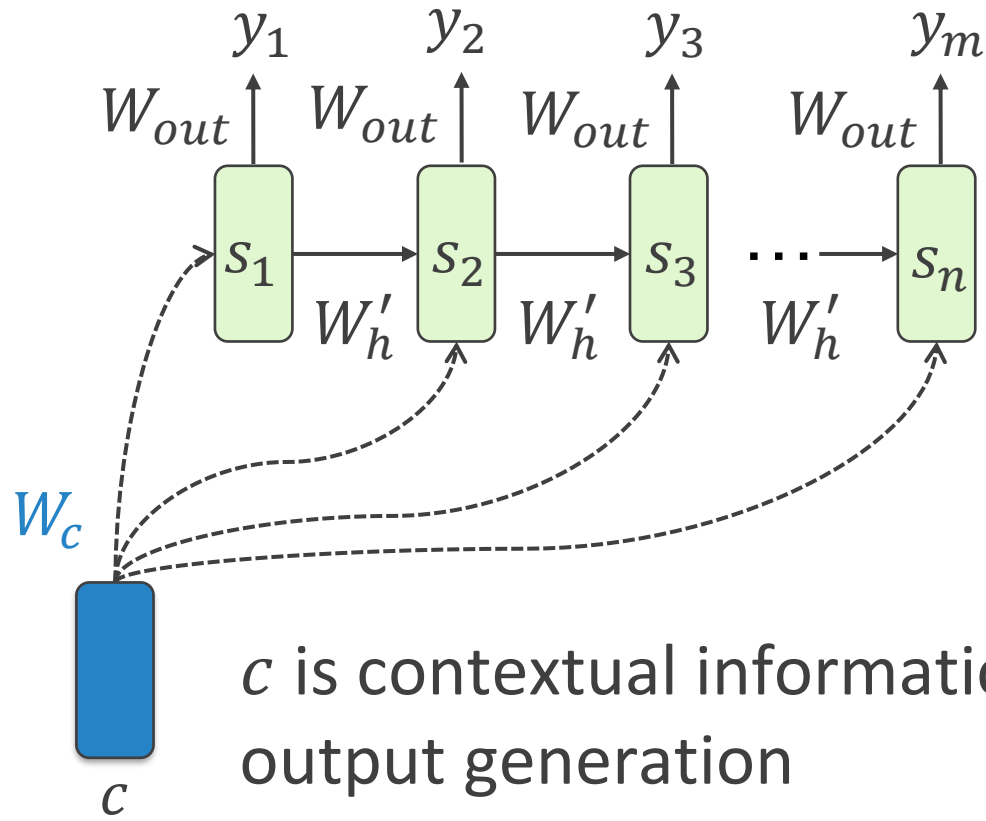
# Decoder



We risk to lose memory of $c$ soon

If we share the parameters between encoder and decoder we can take $s_1 = c$

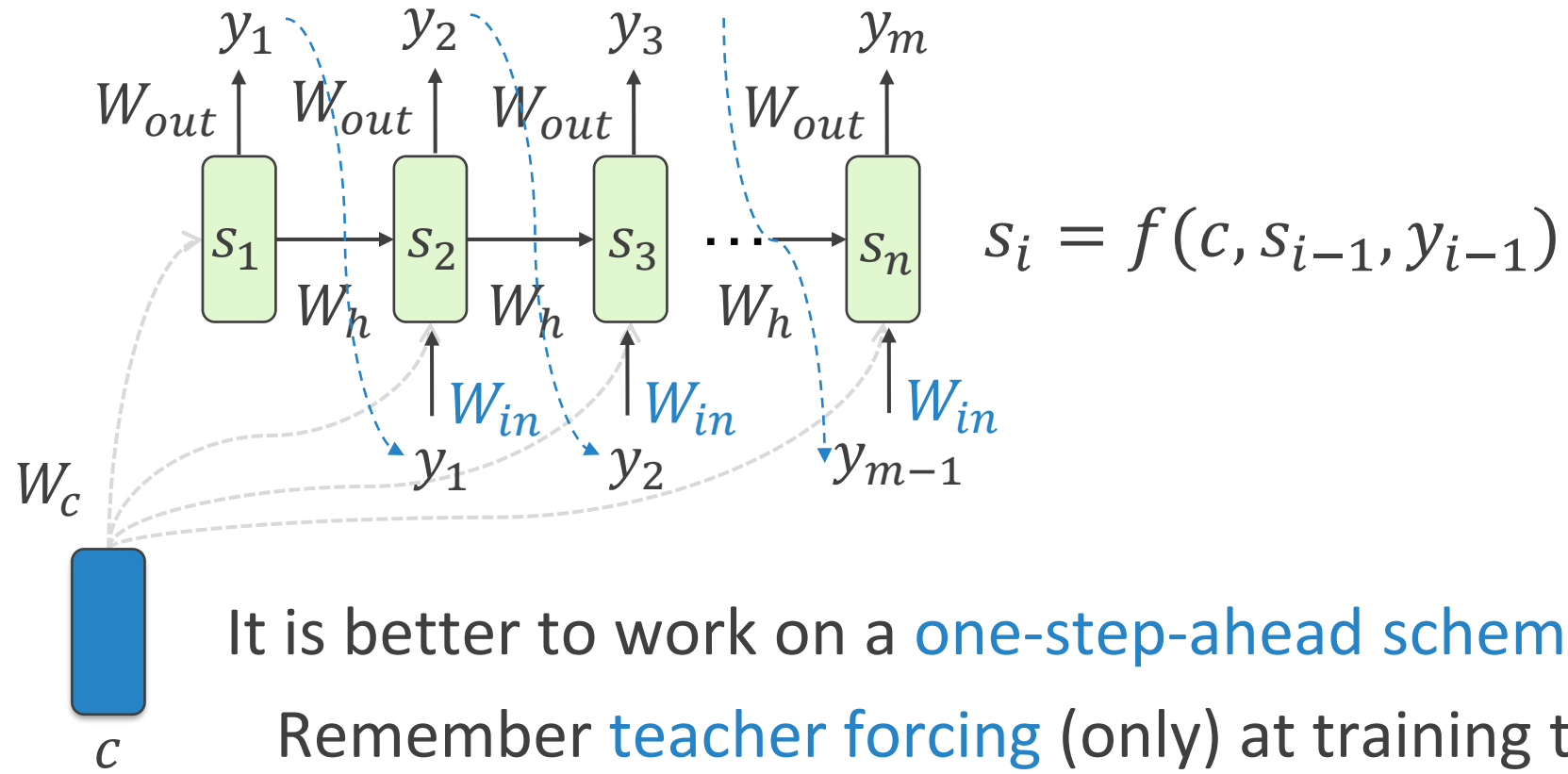Or, at least, assume $c$ and $s_1$ have compatible size

# Decoder



$c$ is contextual information kept throughout output generation

# Decoder



$$s_i = f(c, s_{i-1}, y_{i-1})$$

It is better to work on a one-step-ahead scheme

Remember teacher forcing (only) at training time
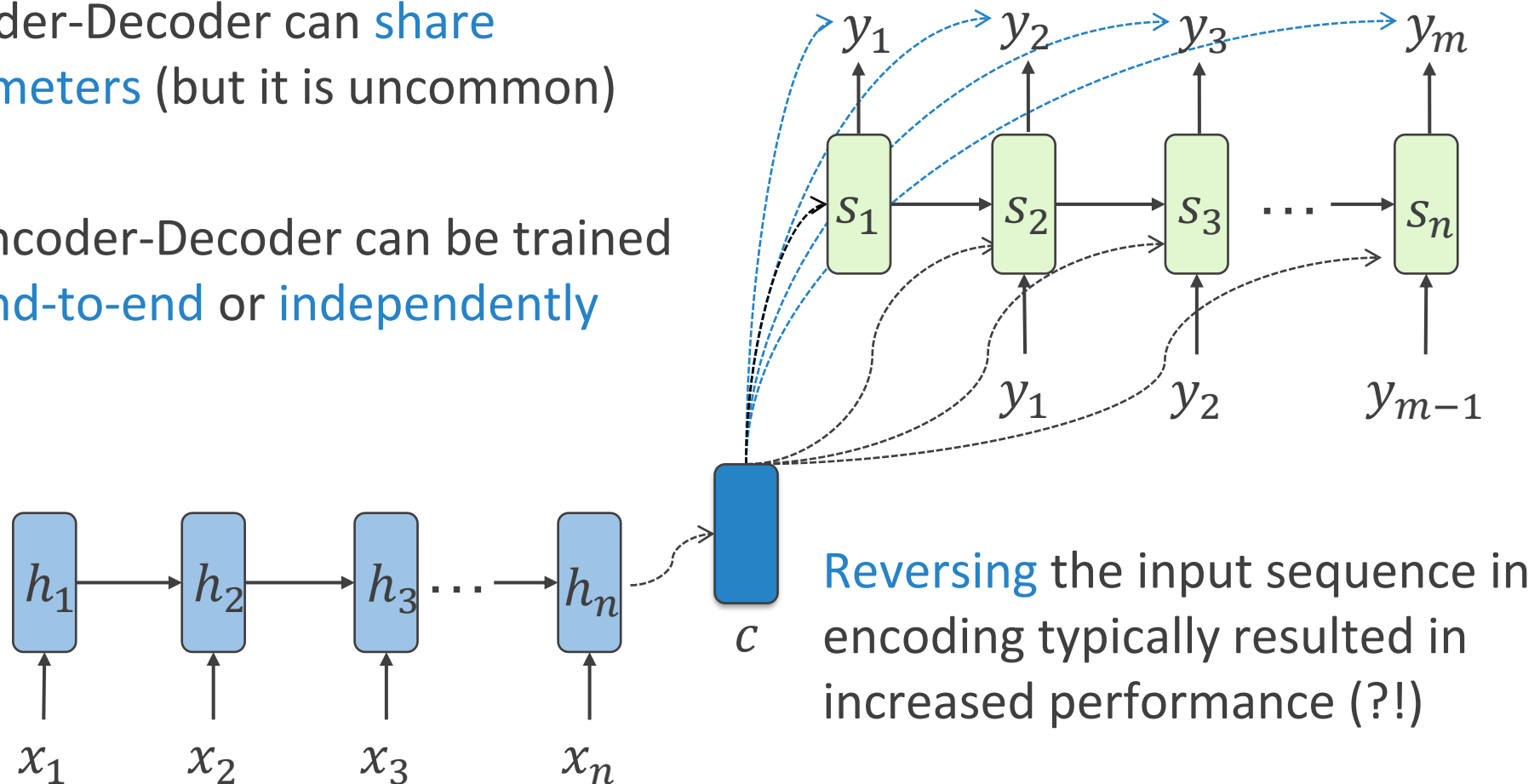
# Sequence-To-Sequence Learning

Encoder-Decoder can share parameters (but it is uncommon)

Encoder-Decoder can be trained end-to-end or independently



Reversing the input sequence in encoding typically resulted in increased performance (?!)

# A Motivating Example
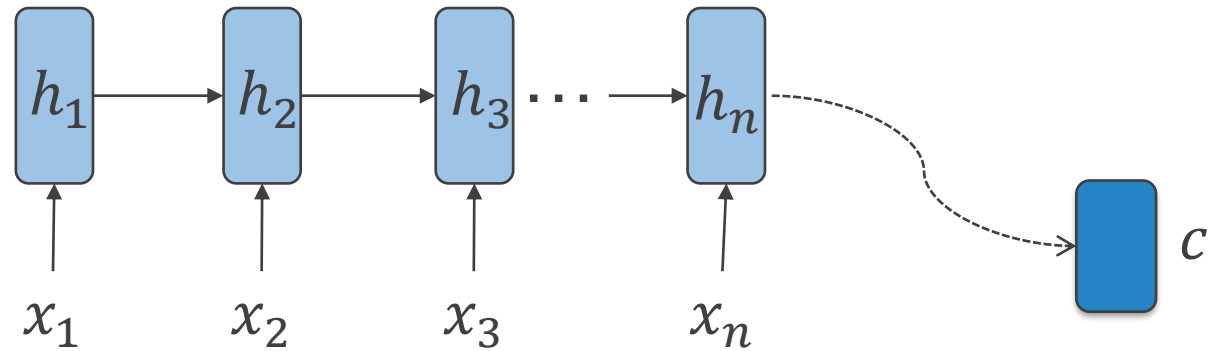
**The cat is on the table**
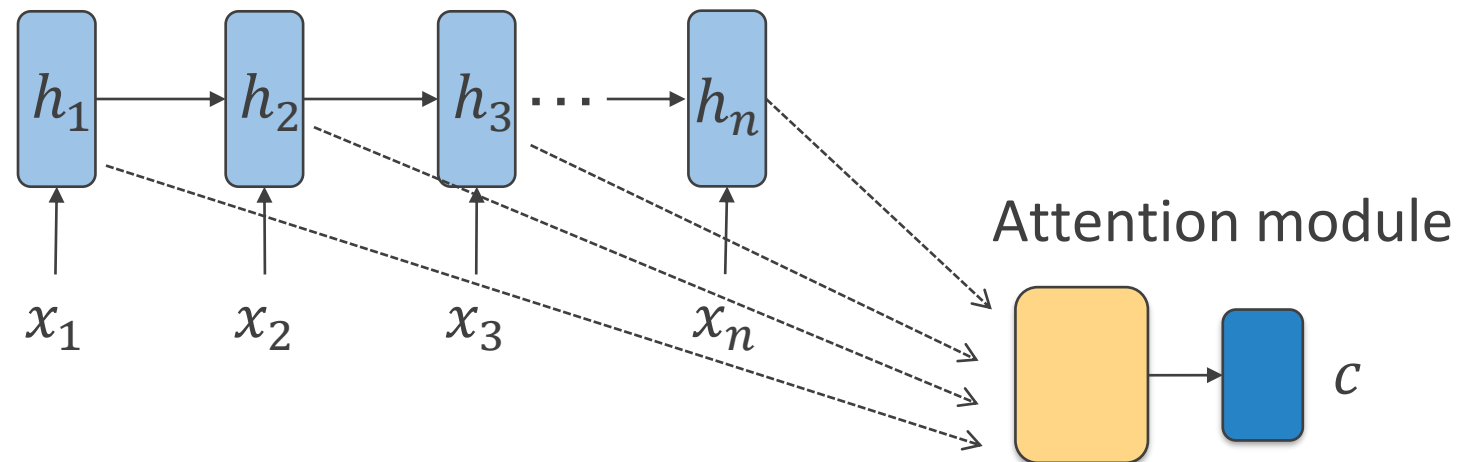
**Il gatto è sul tavolo**
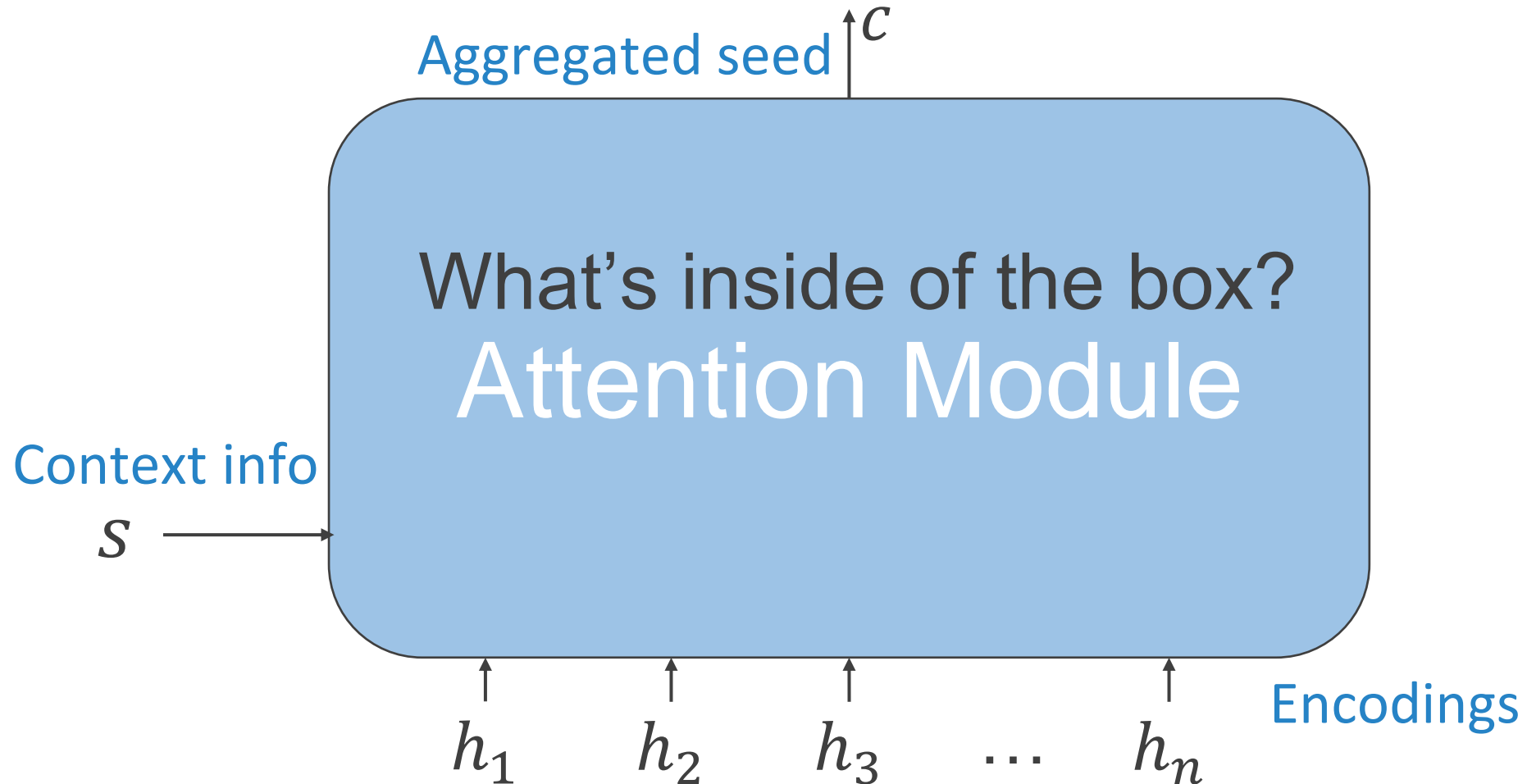
# Attention

# On the Need of Paying Attention



o Encoder-Decoder scheme assumes the hidden activation of the last input element summarizes sufficient information to generate the output
- Bias toward most recent past

o Other parts of the input sequence might be very informative for the task
- Possibly elements appearing very far from sequence end

# On the Need of Paying Attention



- Attention mechanisms select which part of the sequence to focus on to obtain a good $c$

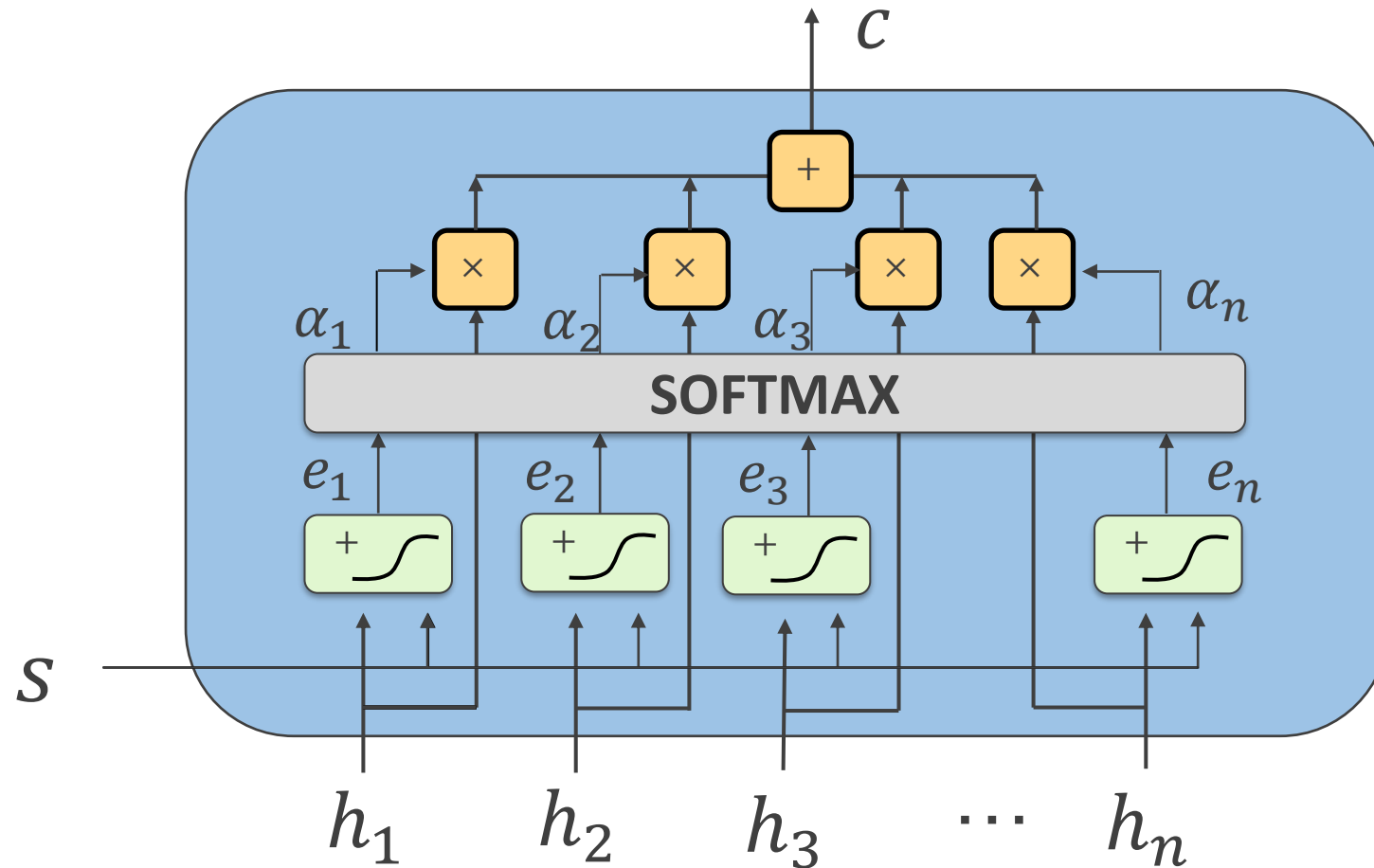# Attention Mechanisms – Blackbox View



Aggregated seed $\uparrow c$

What's inside of the box?
Attention Module

Context info
$s$

Encodings

$h_1 \quad h_2 \quad h_3 \quad \dots \quad h_n$

# What's inside of the box?

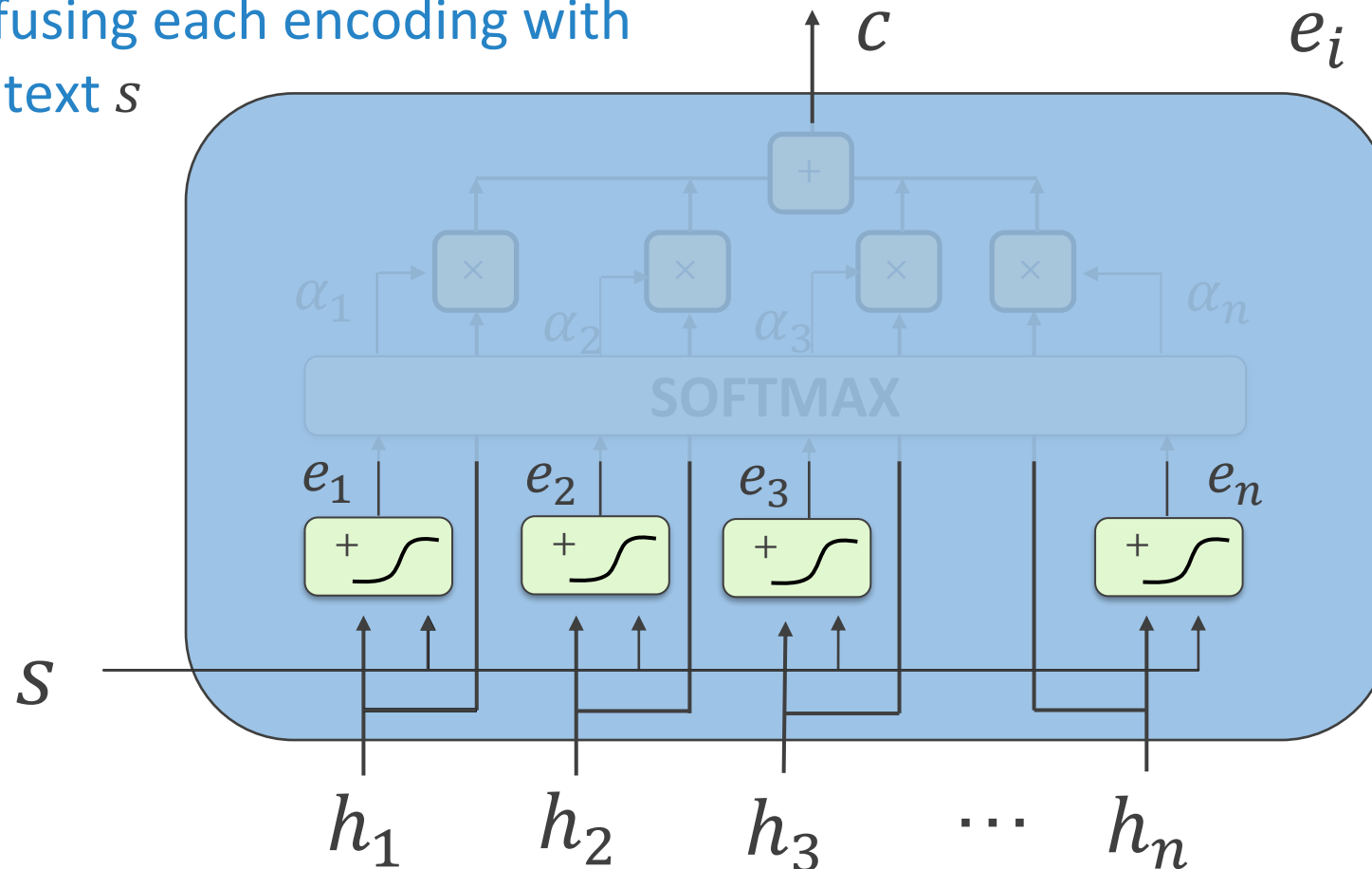## The Revenge of the Gates!

# Opening the Box

# Opening the Box – Relevance

Tanh layer fusing each encoding with current context $s$

$$e_i = a(s, h_i)$$

# Opening the Box – Softmax

# Opening the Box – Voting

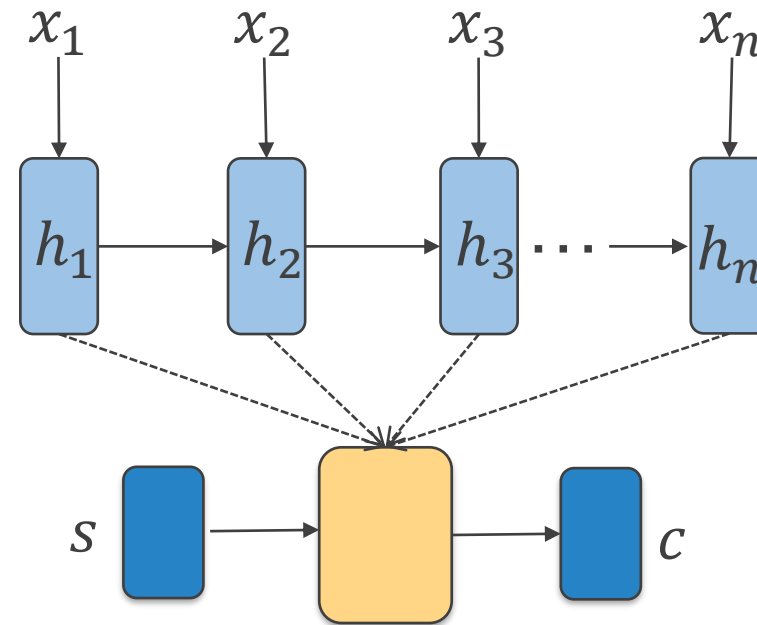Aggregated seed by (soft) attention voting



$$c = \sum_i \alpha_i h_i$$

# Attention - Equations

- Relevance: $e_i = a(s, h_i)$

- Normalization: $\alpha_i = \dfrac{\exp(e_i)}{\sum_j \exp(e_j)}$
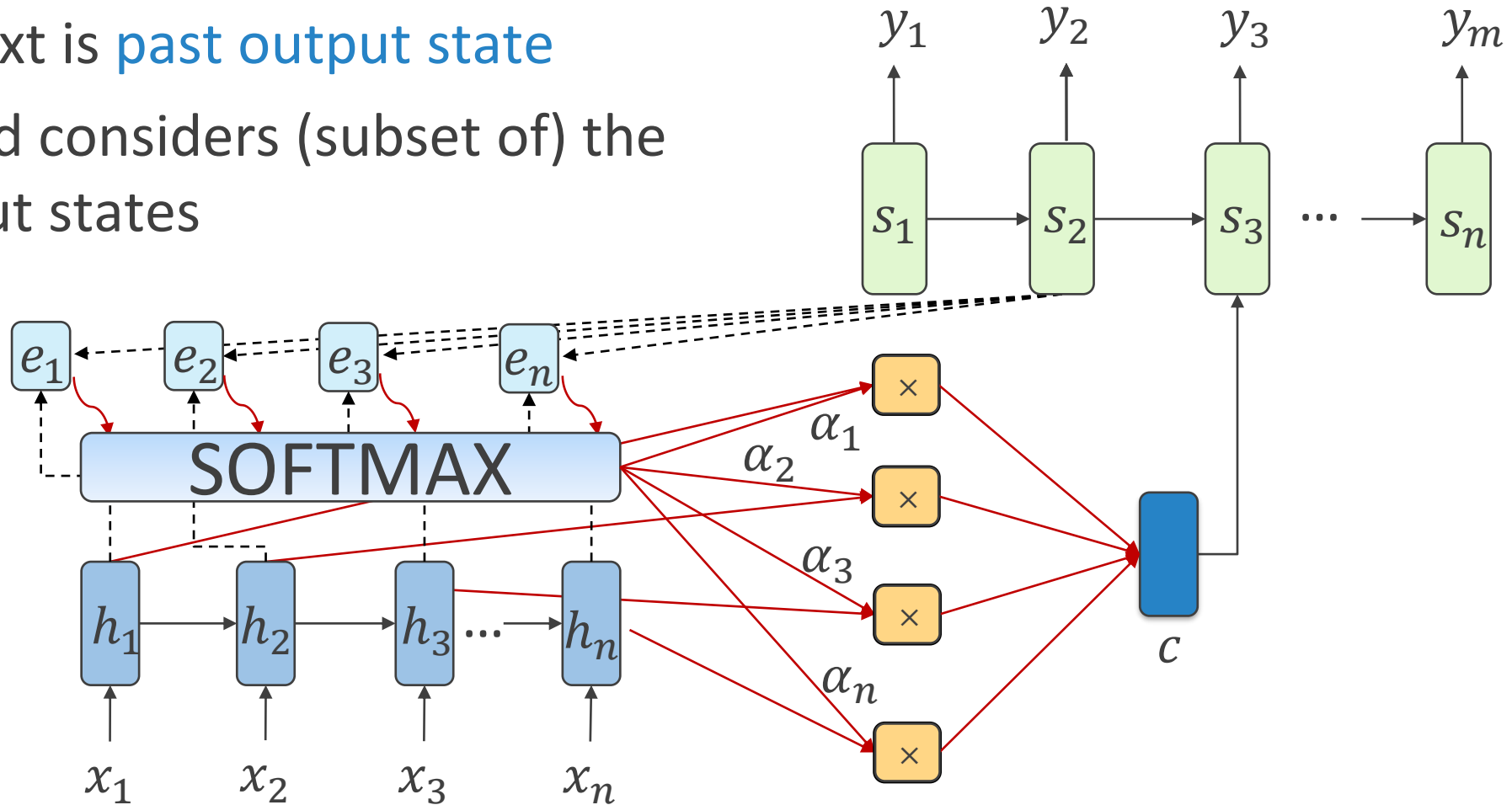
- Aggregation: $c = \sum_i \alpha_i h_i$
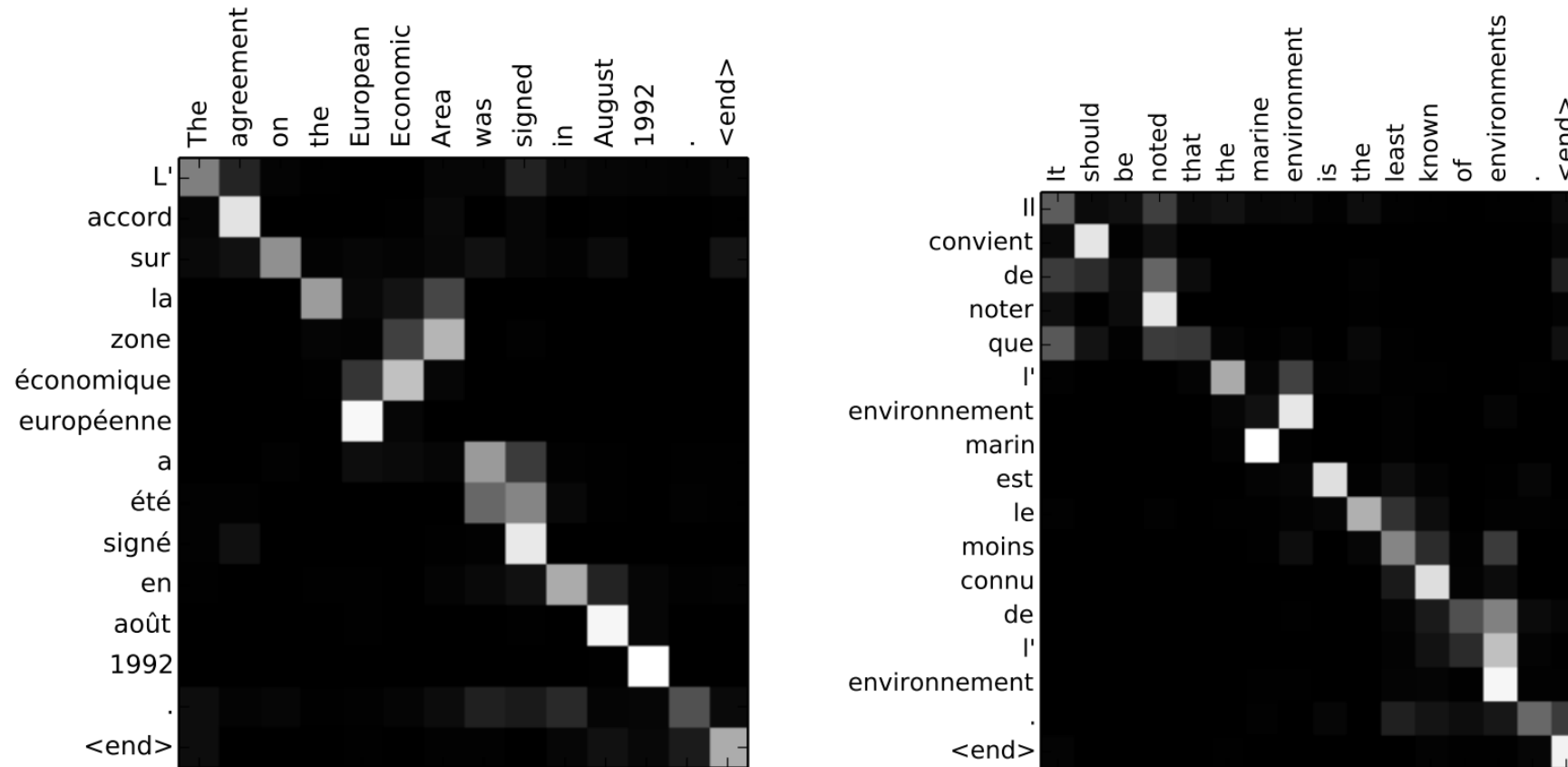


Attention module

# Attention in Seq2Seq

Context is past output state

Seed considers (subset of) the input states
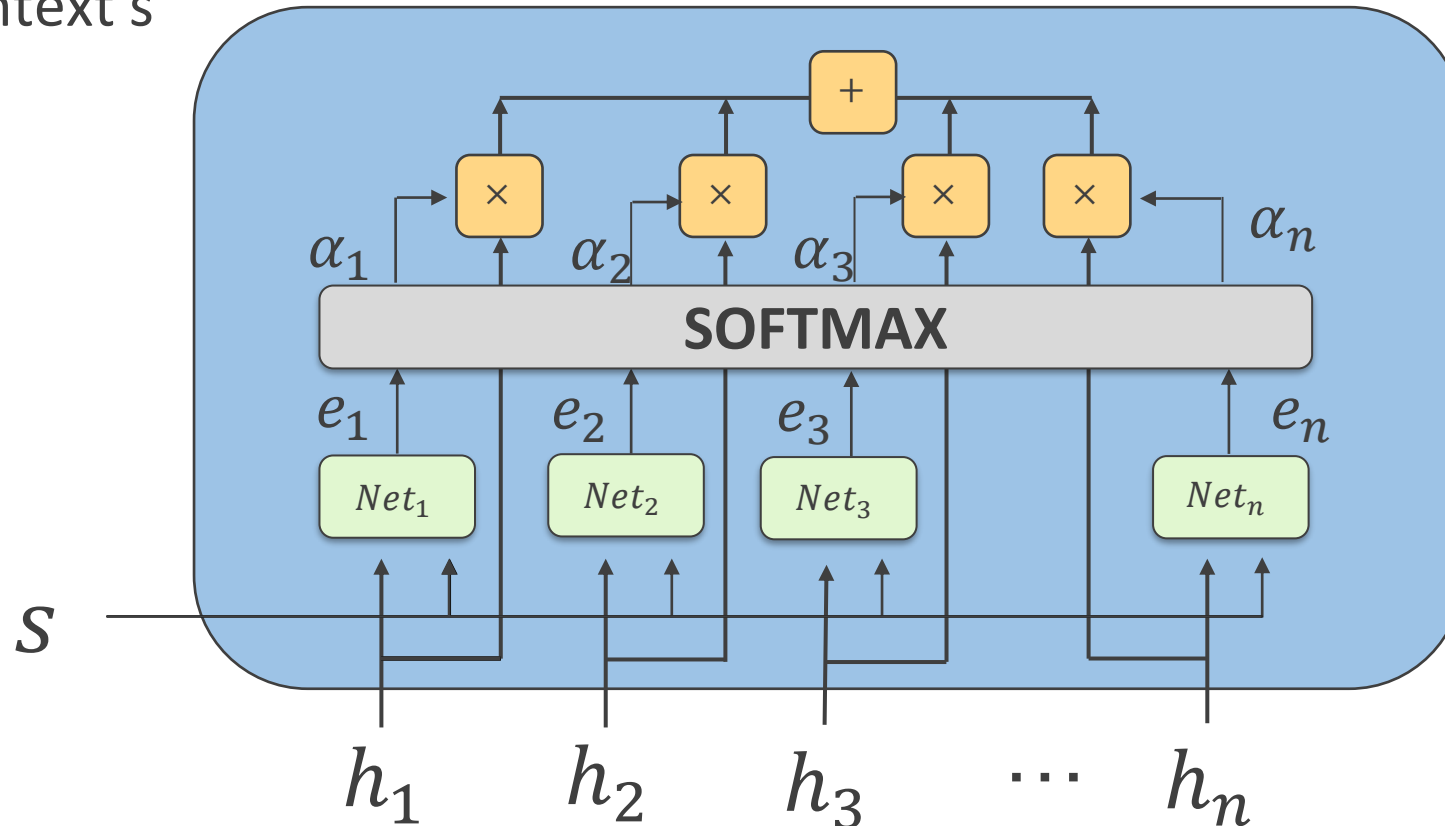
# Learning to Translate with Attention



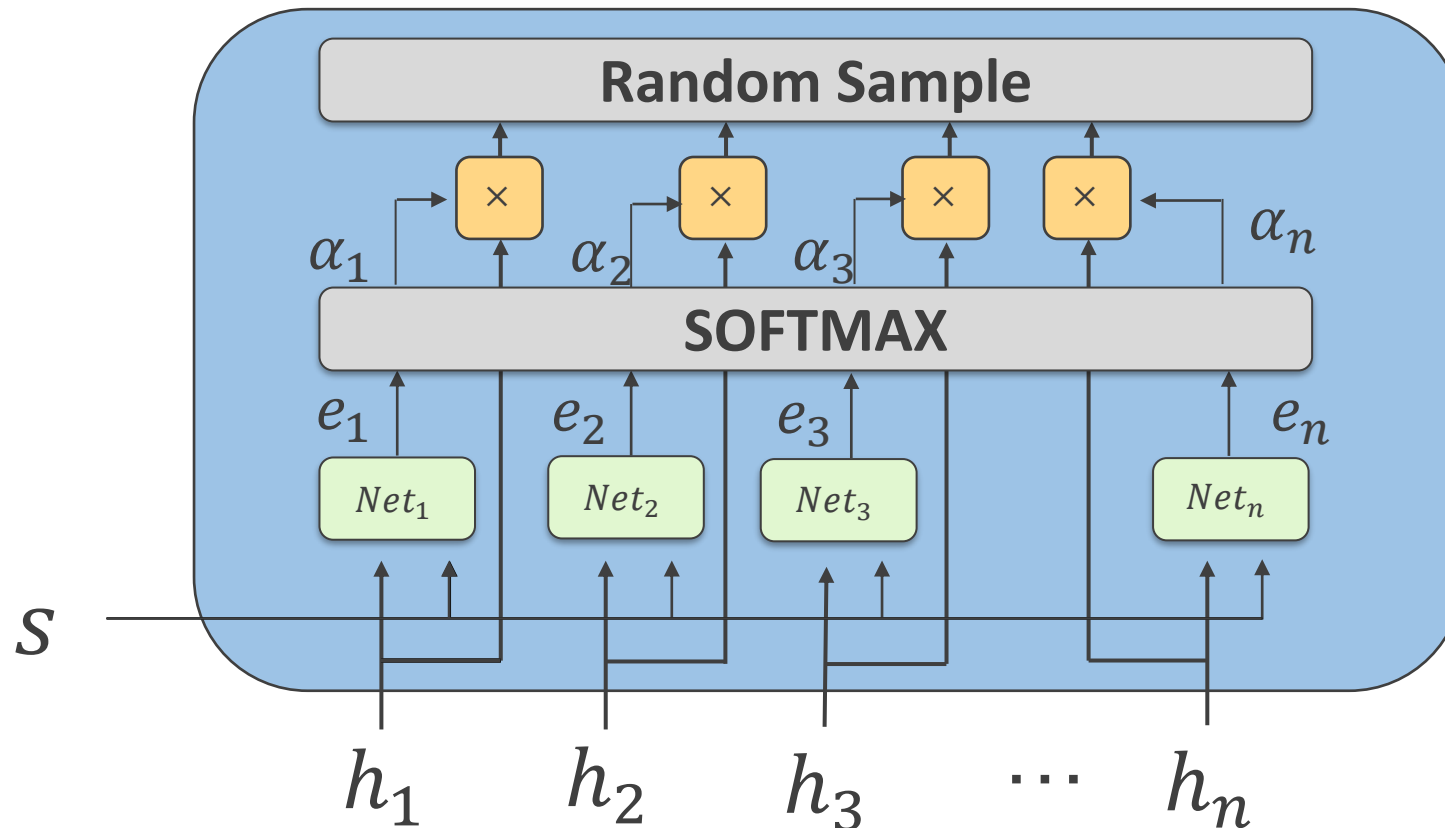Bahdanau et al, Show, Neural machine translation by jointly learning to align and translate, ICLR 2015

# Advanced Attention – Generalize Relevance

This component determines how much each h is correlated/associated with current context s
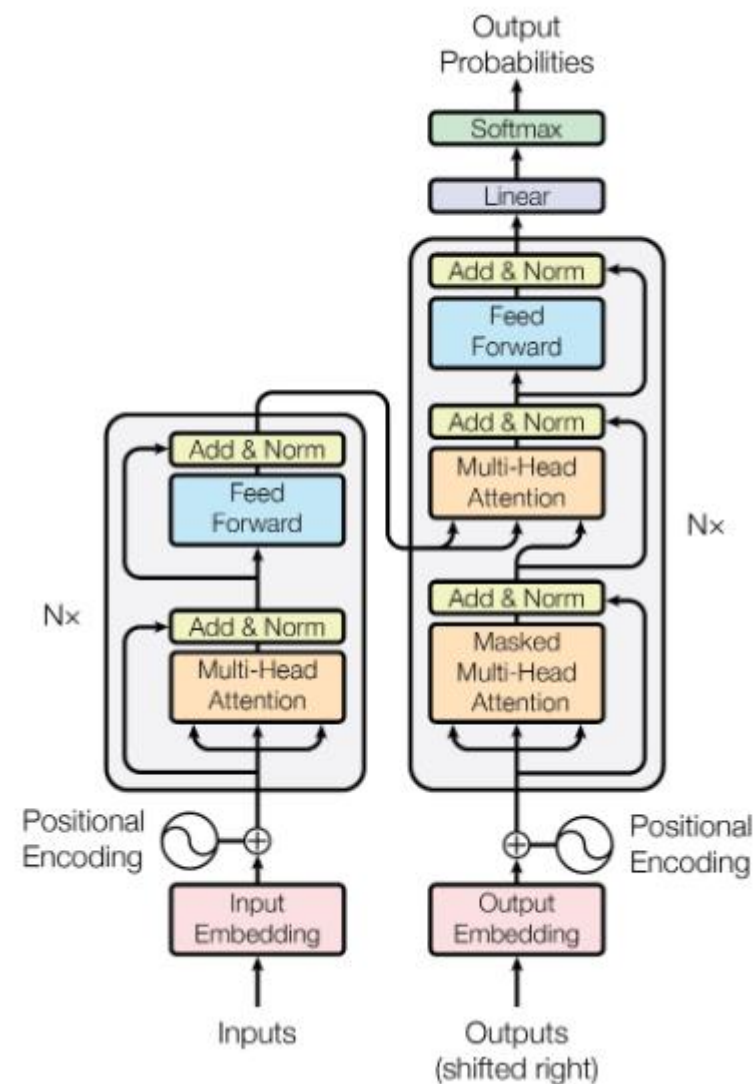
# Advanced Attention – Hard Attention

Sample a single encoding using probability $\alpha_i$

# Transformers

o First pure attention-based model

o Self-attention

o No recurrence
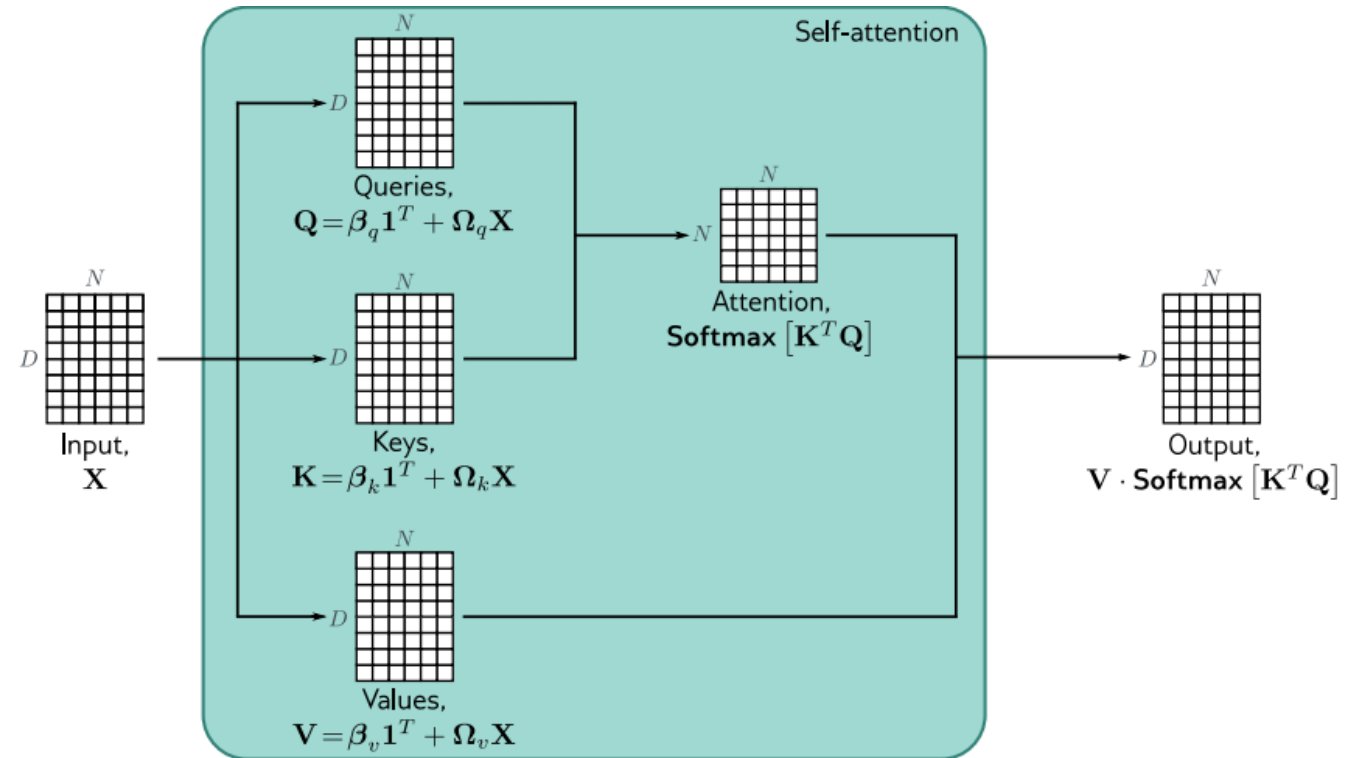
o Encoder-decoder architecture

# Self Attention

Each element of an input sequence $X_i$ projects into 3 vectors: **query**, **key** and **value**
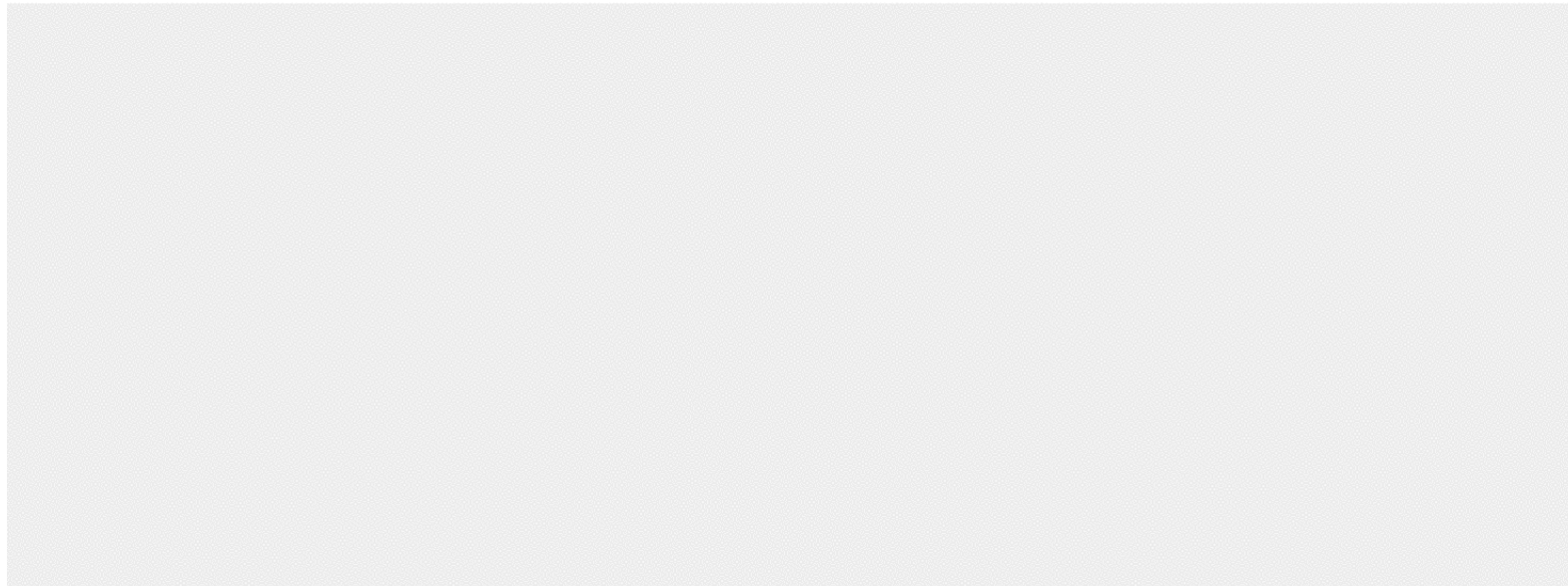
**Scaled self-attention**

$$\sum_j softmax_j \left( \frac{Q_i \cdot \boldsymbol{K}^T}{\sqrt{d_k}} \right) V_j$$

# Self Attention – K,V,Q Generation



Self-attention

input #1

| 1 | 0 | 1 | 0 |

input #2

| 0 | 2 | 0 | 2 |

input #3

| 1 | 1 | 1 | 1 |

Figure credit to this article

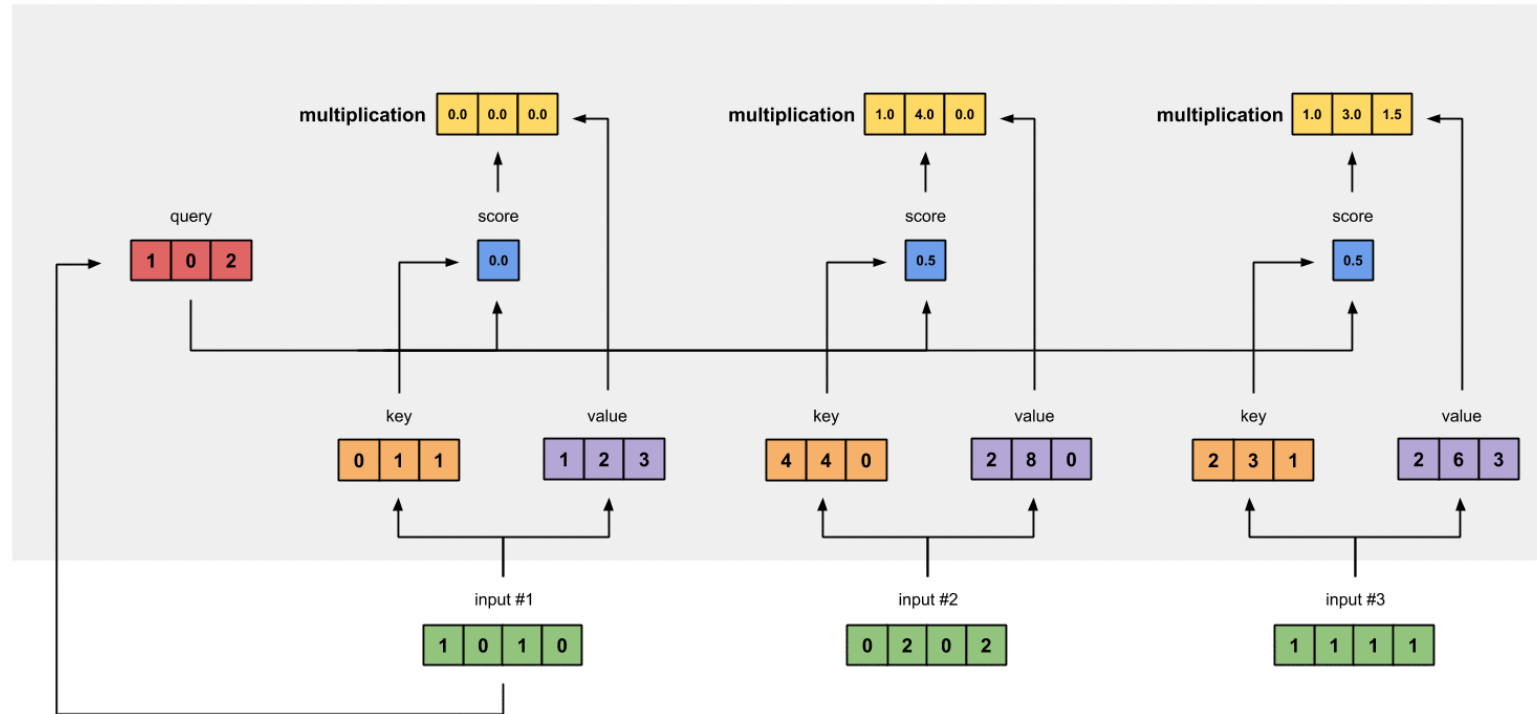# Self Attention – Compute Attention Score

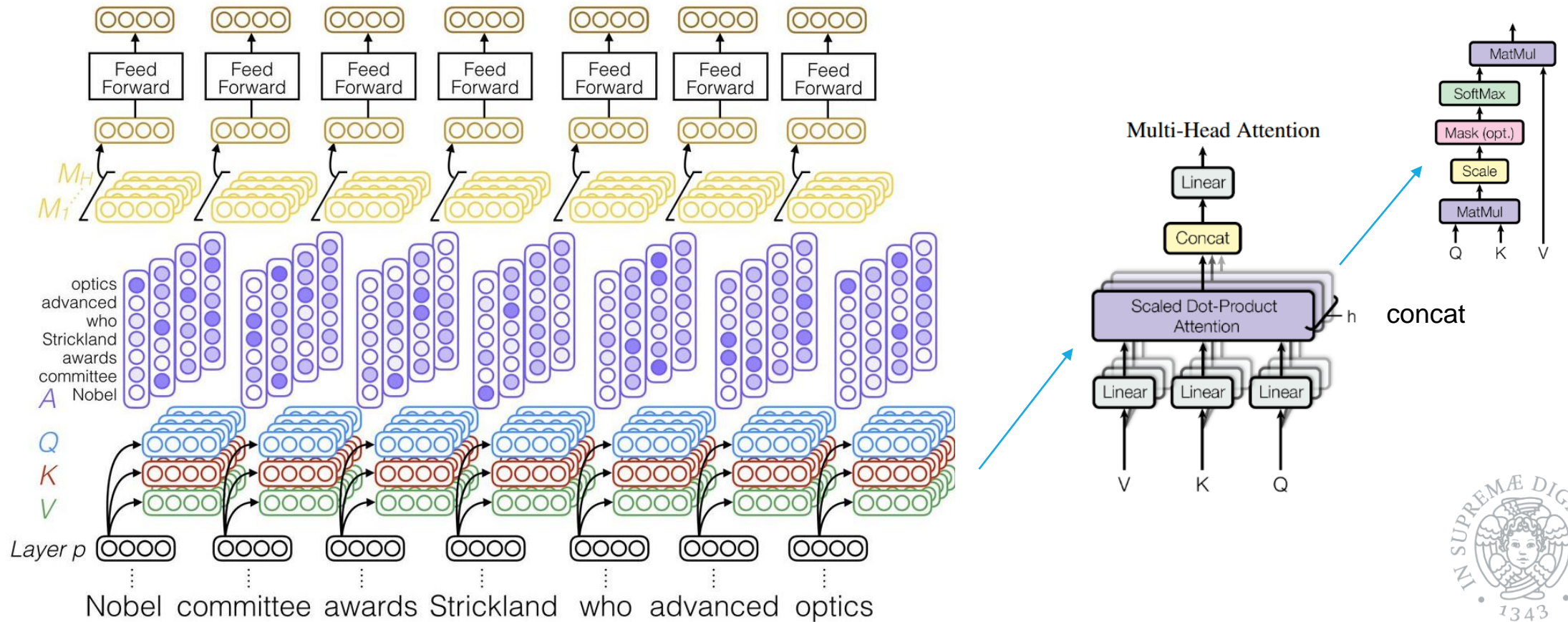# Self Attention – Produce Output

# Self Attention – MultiHead



Strubell et al, Linguistically-Informed Self-Attention for Semantic Role Labeling, EMNLP 2018
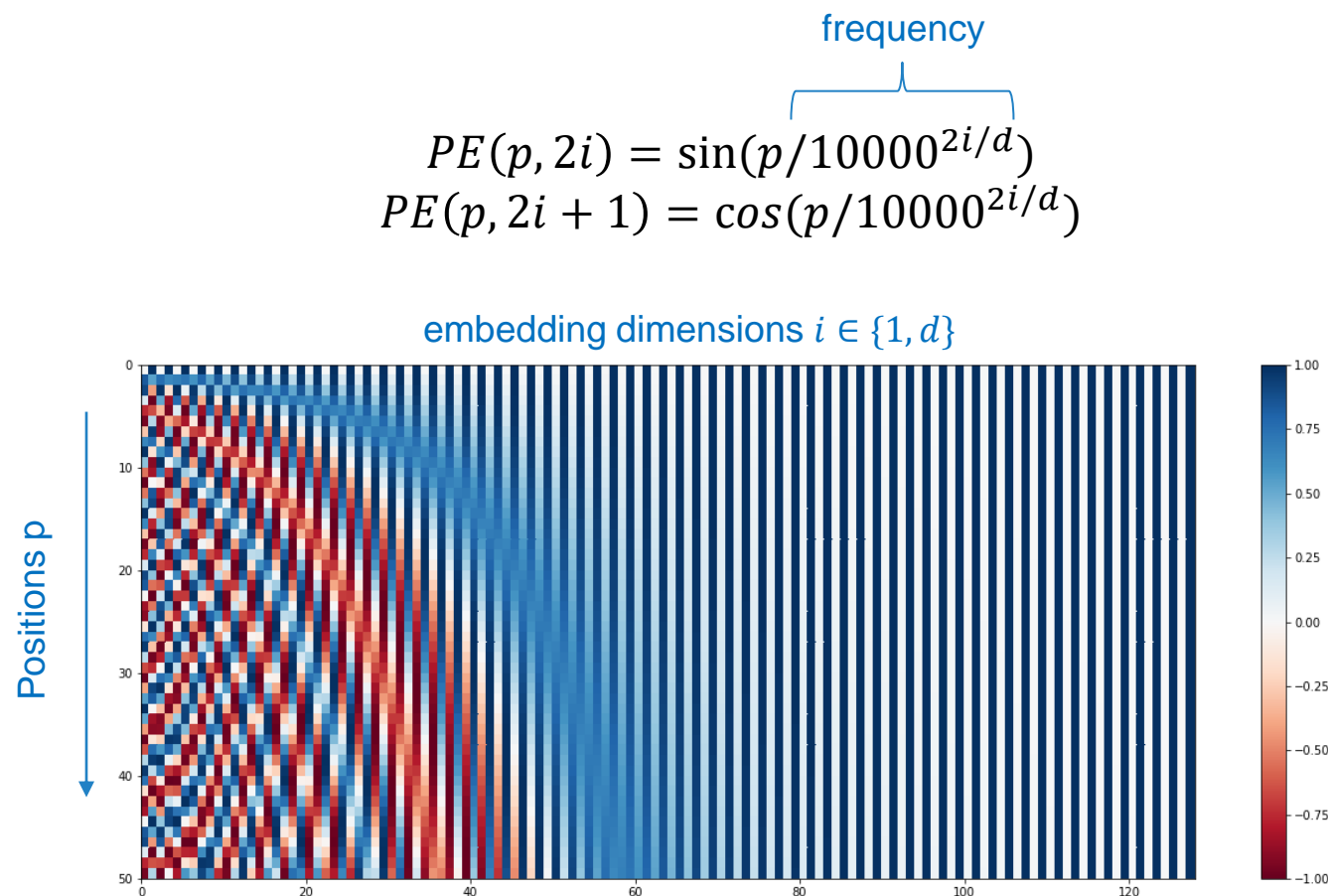
concat

Is self-attention a good mechanism to model temporal dependencies?

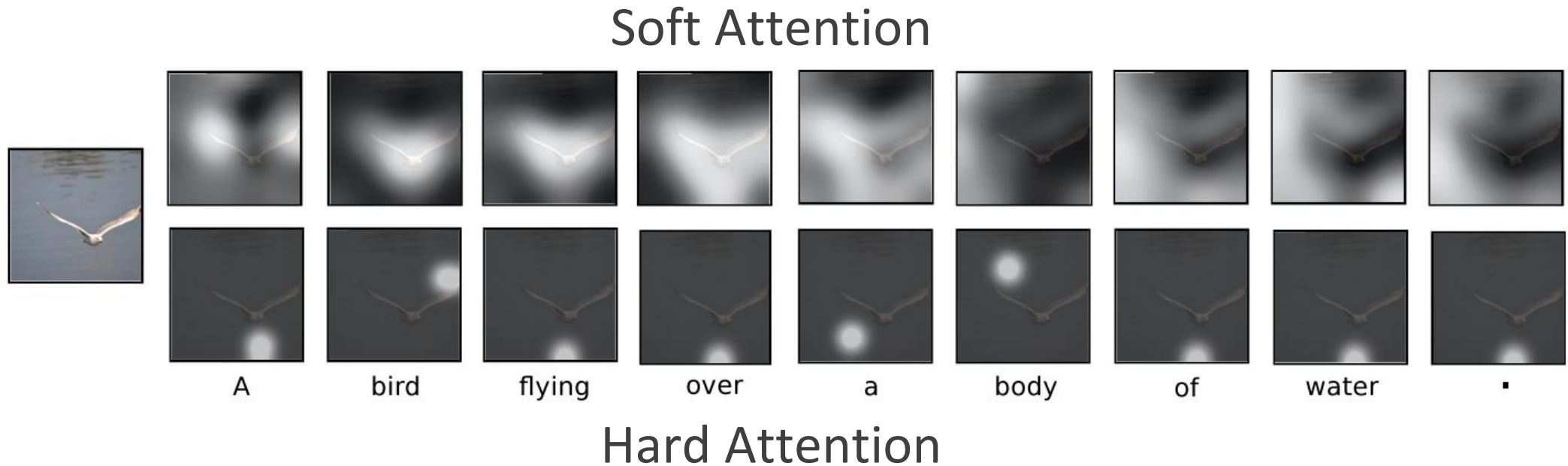What happens if I randomly shuffle some tokens?

# (Absolute) Positional Encoding

o Self-attention is order-independent

o But in sequences we need ordering information

o word embedding + positional embedding

frequency

$$PE(p, 2i) = \sin(p/10000^{2i/d})$$
$$PE(p, 2i + 1) = \cos(p/10000^{2i/d})$$

embedding dimensions $i \in \{1, d\}$

Positions p

# Attention in Vision

# Attention-Based Captioning – Focus Shifting

## Soft Attention



## Hard Attention

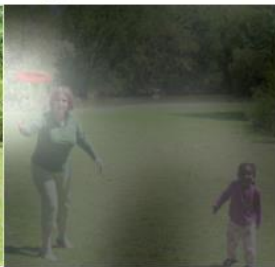A   bird   flying   over   a   body   of   water   .

Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

# Attention-Based Captioning - Generation

Learns to correlate textual and visual concepts



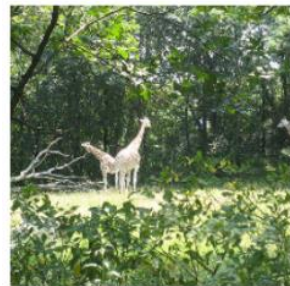A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

Helps understanding why the model fails



A large white <u>bird</u> standing in a forest.

A woman holding a <u>clock</u> in her hand.

Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015
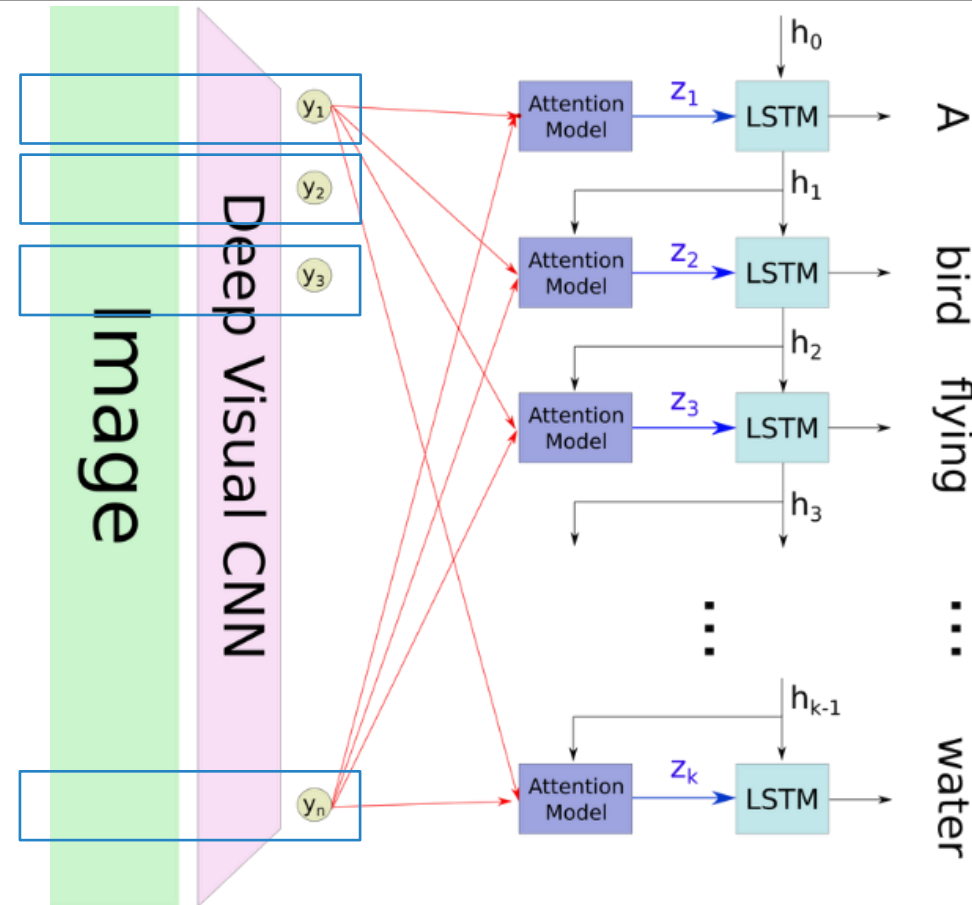
# Attention-Based Captioning – The Model

Encodings associated to *n* image regions

From convolutional layers rather than from fully connected



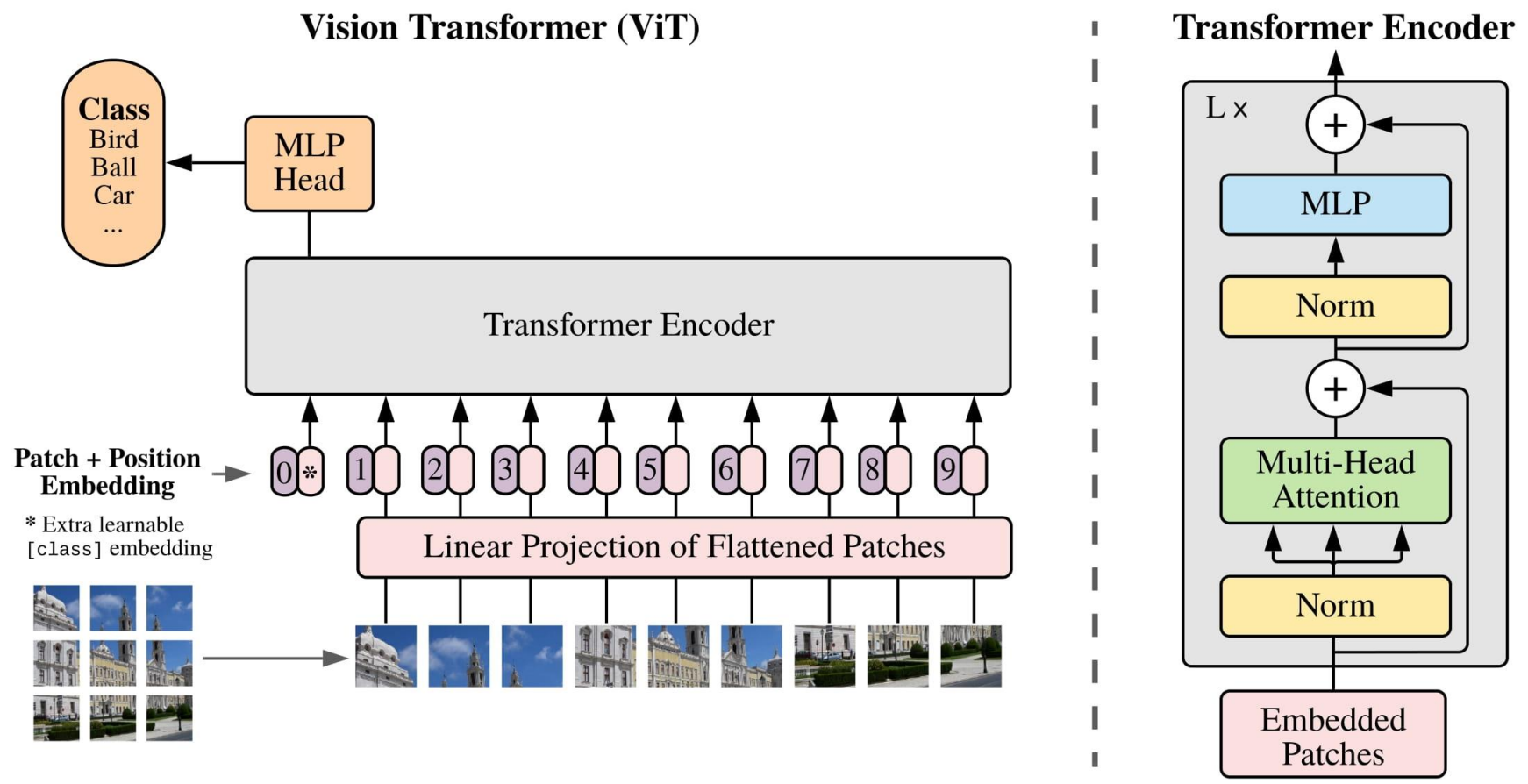Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

# The Vision Transformer (ViT)

# Take Home Messages

○ Attention.. Attention.. and, again, attention
- Soft attention is nice because makes everything fully differentiable
- Hard attention is stochastic hence cannot Backprop
- Empirical evidences of them being sensitive to different things

○ Encoder-Decoder scheme
- A general architecture to compose heterogeneous models and data
- Decoding allows sampling complex predictions from an encoding conditioned distribution

○ Transformers as low-inductive bias architectures
- Need huge amounts of data to generalize