

**RNN / Enc-Dec / Attention / Transformer**

## **Neural-LM [2003]**

1. Bengio et al., “A neural probabilistic language model“, NIPS 2000
2. Bengio et al., “A neural probabilistic language model“, JMLR 2003.

## **RNN-LM [2010]**

T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur. “Recurrent neural network based language model“, In INTERSPEECH, pages 1045–1048, 2010.

## **Seq2seq learning (Encoder-Decoder) [2014]**

Sutskever, I., Vinyals, O., & Le, Q. V., “Sequence to sequence learning with neural networks“, Advances in Neural Information Processing Systems (pp. 3104–3112), 2014

## **Attention [2015]**

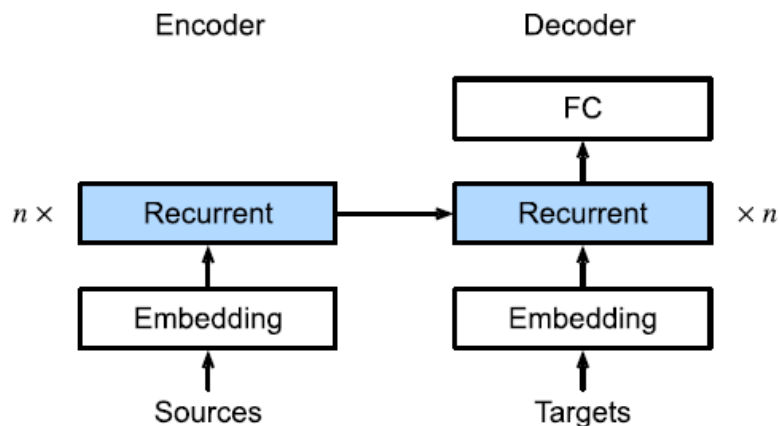
Bahdanau, D., Cho, K., & Bengio, Y., “Neural machine translation by jointly learning to align and translate“, Proc. ICLR 2015

## **Transformer [2017]**

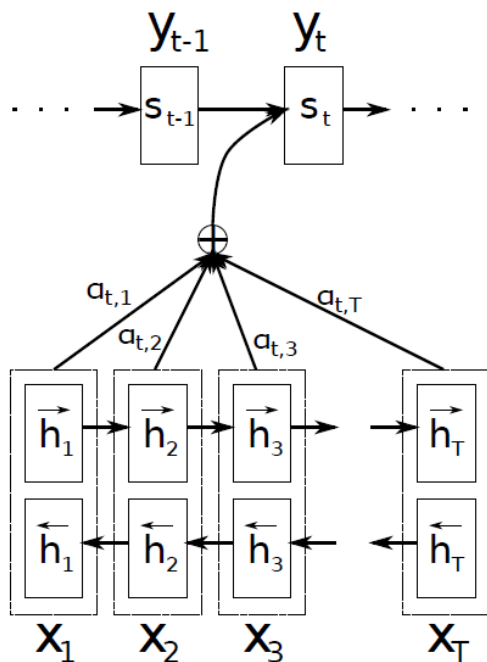
A. Vaswani et al., “Attention is all you need“, In NIPS, pages 6000–6010, 2017.

## **GPT [2018]**

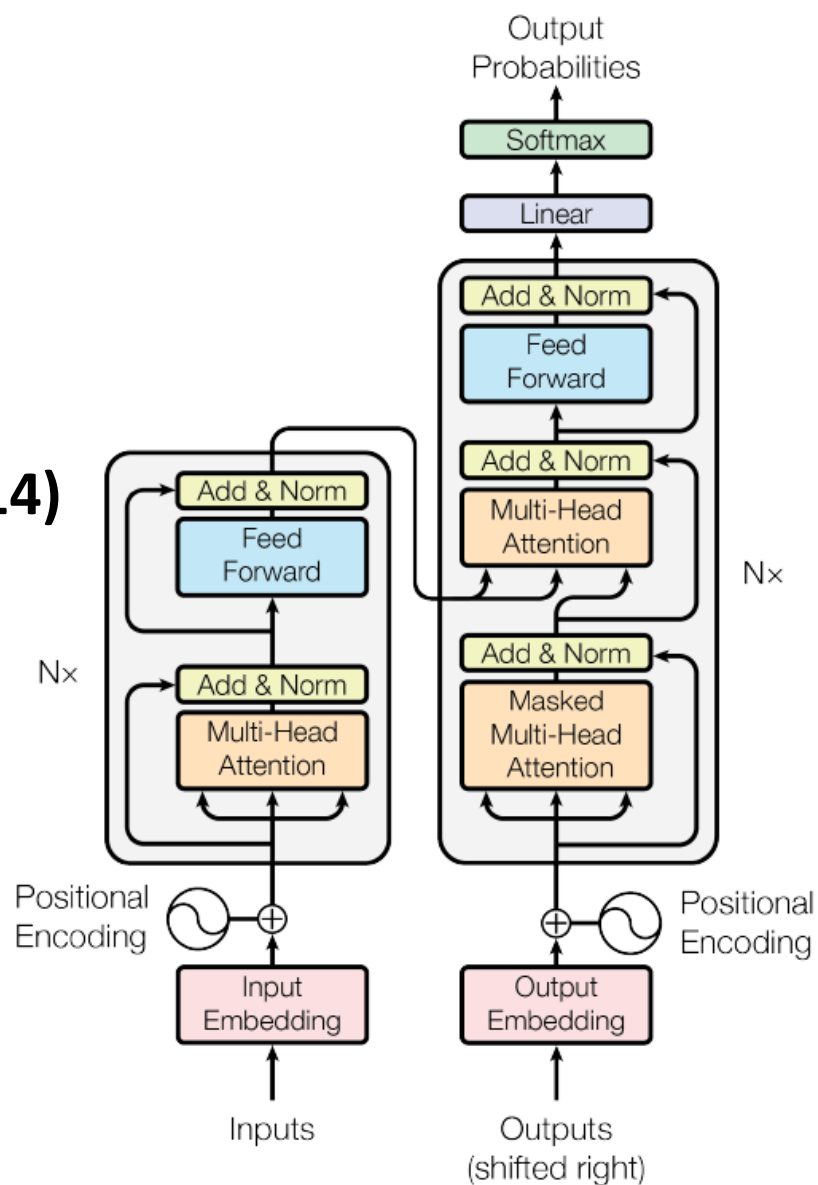
A. Radford, K. Narasimhan, T. Salimans and I. Sutskever, “Improving Language Understanding by Generative Pre-Training“, 2018 (arXiv)



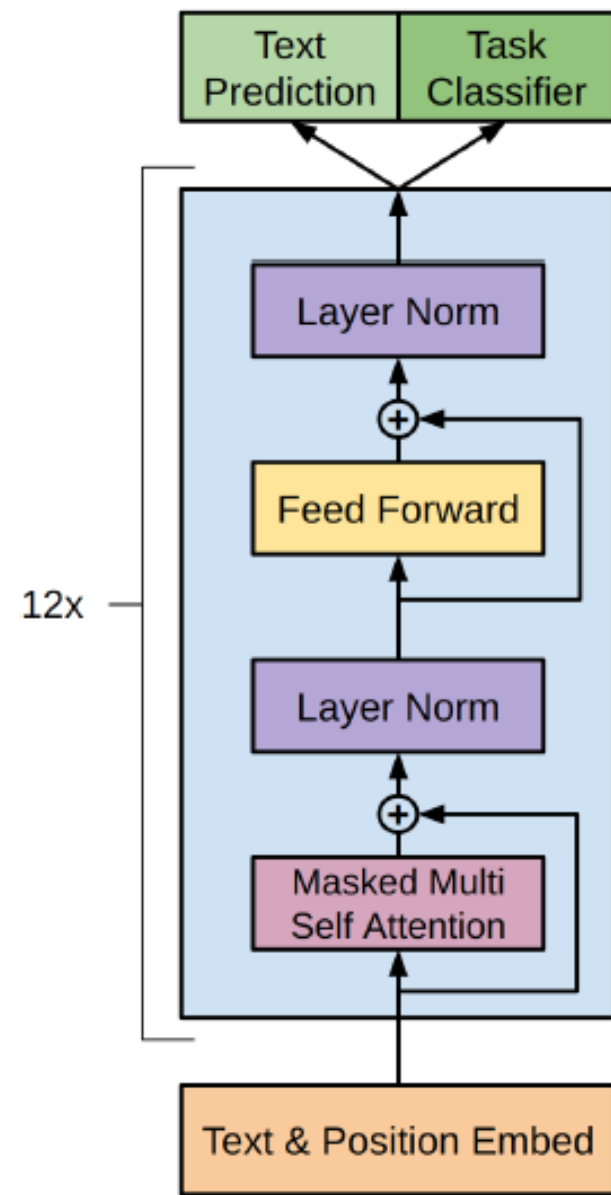
**RNN Encoder-Decoder Model (2014)**



**Bahdanau Attention (2015)**



**Transformer Model (2017)**

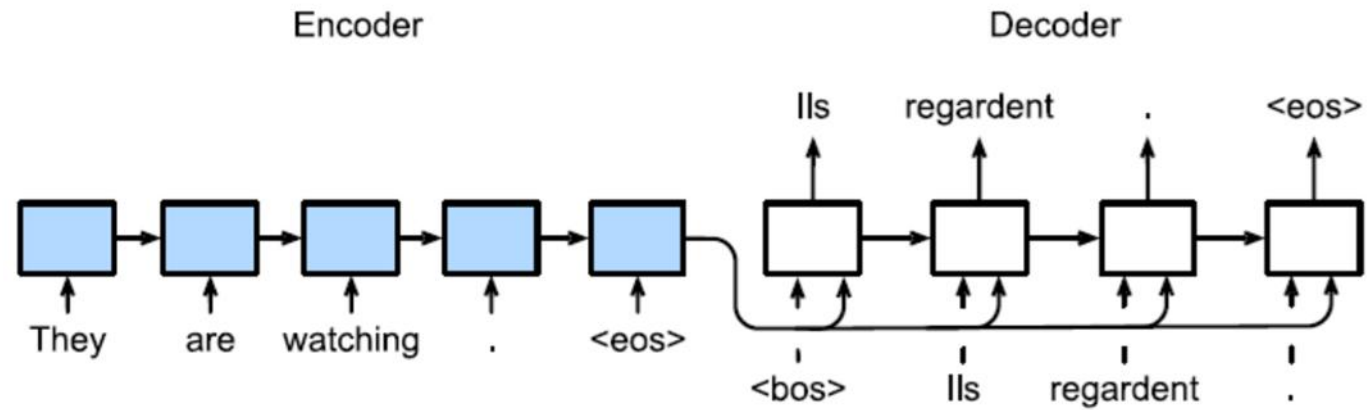


**Decoder-only GPT (2018)**

# Encoder-Decoder Architecture: Sequence-to-Sequence Learning

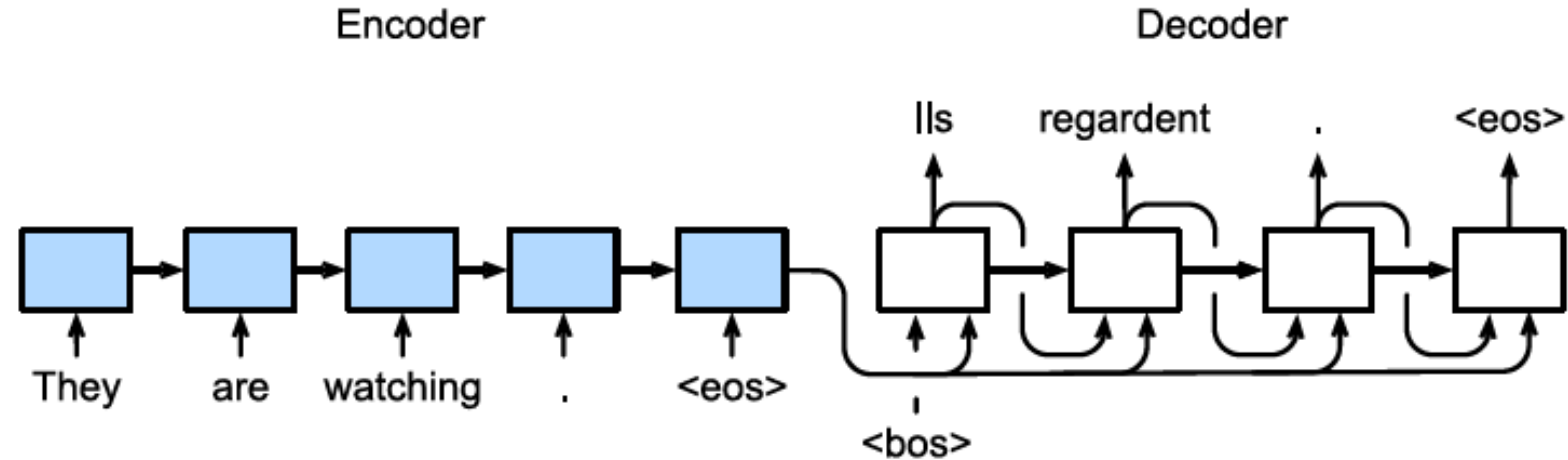
Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in Neural Information Processing Systems* (pp. 3104–3112).

# TRAINING: TEACHER FORCING



Sequence-to-sequence learning with an RNN encoder and an RNN decoder.

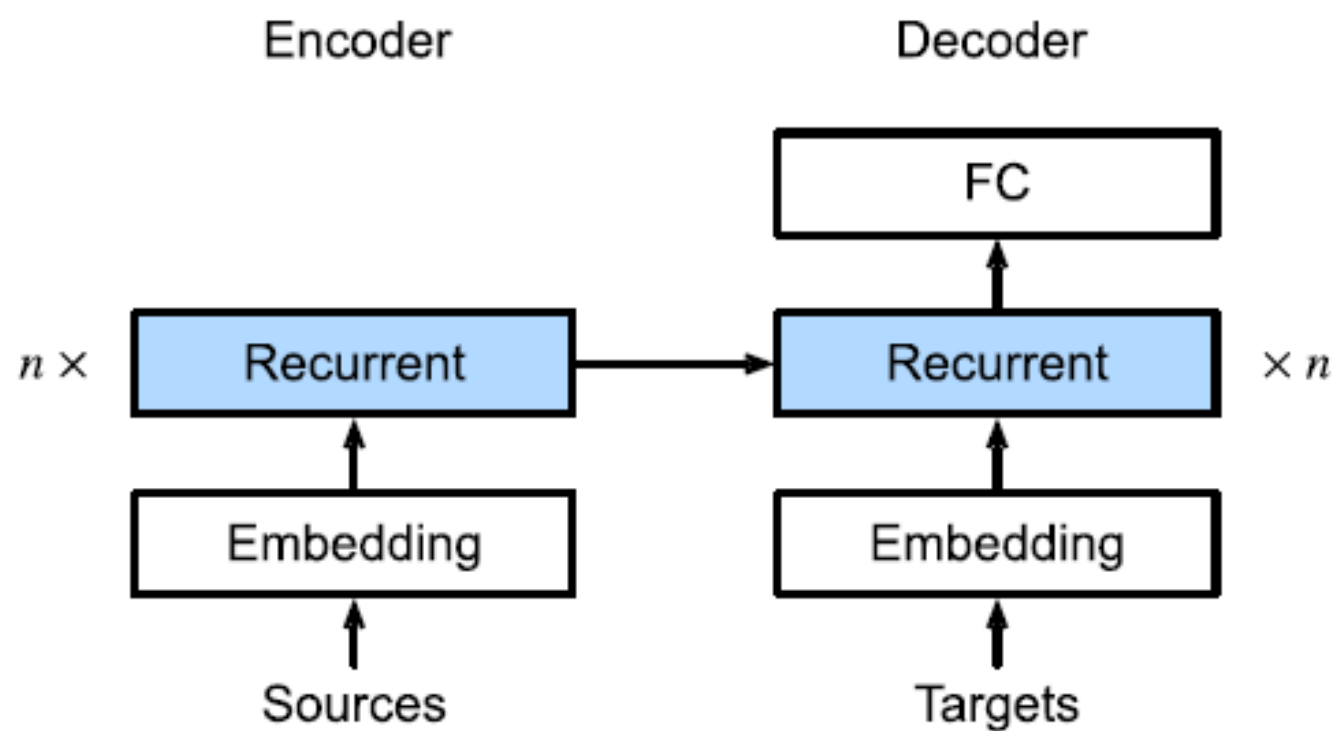
# TEST: PREDICTION



Predicting the output sequence token by token using an RNN encoder–decoder.

- ❑ Predict the output sequence at each step: Predicted token from the previous time step is fed into the decoder as an input.
- ❑ One simple strategy: Sample whichever token that has been assigned by the decoder the highest probability when predicting at each step.
- ❑ As in training, at the initial time step the beginning-of-sequence (“<bos>”) token is fed into the decoder.
- ❑ When the end-of-sequence (“<eos>”) token is predicted, the prediction of the output sequence is complete

# Encoder–Decoder for Sequence-to-Sequence Learning



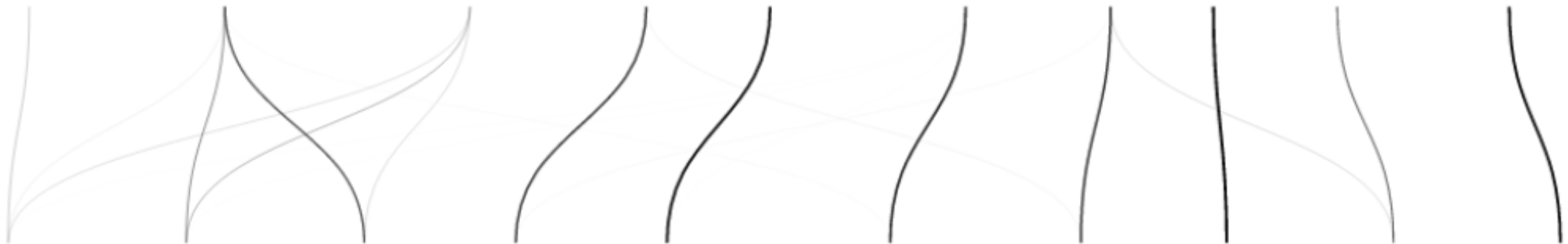
Layers in an RNN encoder–decoder model.

## Attention [2015]

Bahdanau, D., Cho, K., & Bengio, Y.,

“Neural machine translation by jointly learning to align and translate”,  
Proc. ICLR 2015

### Neural Machine Translation by Jointly Learning to Align and Translate



Neural Traduction Automatique par Conjointement Apprentissage Pour Aligner et Traduire



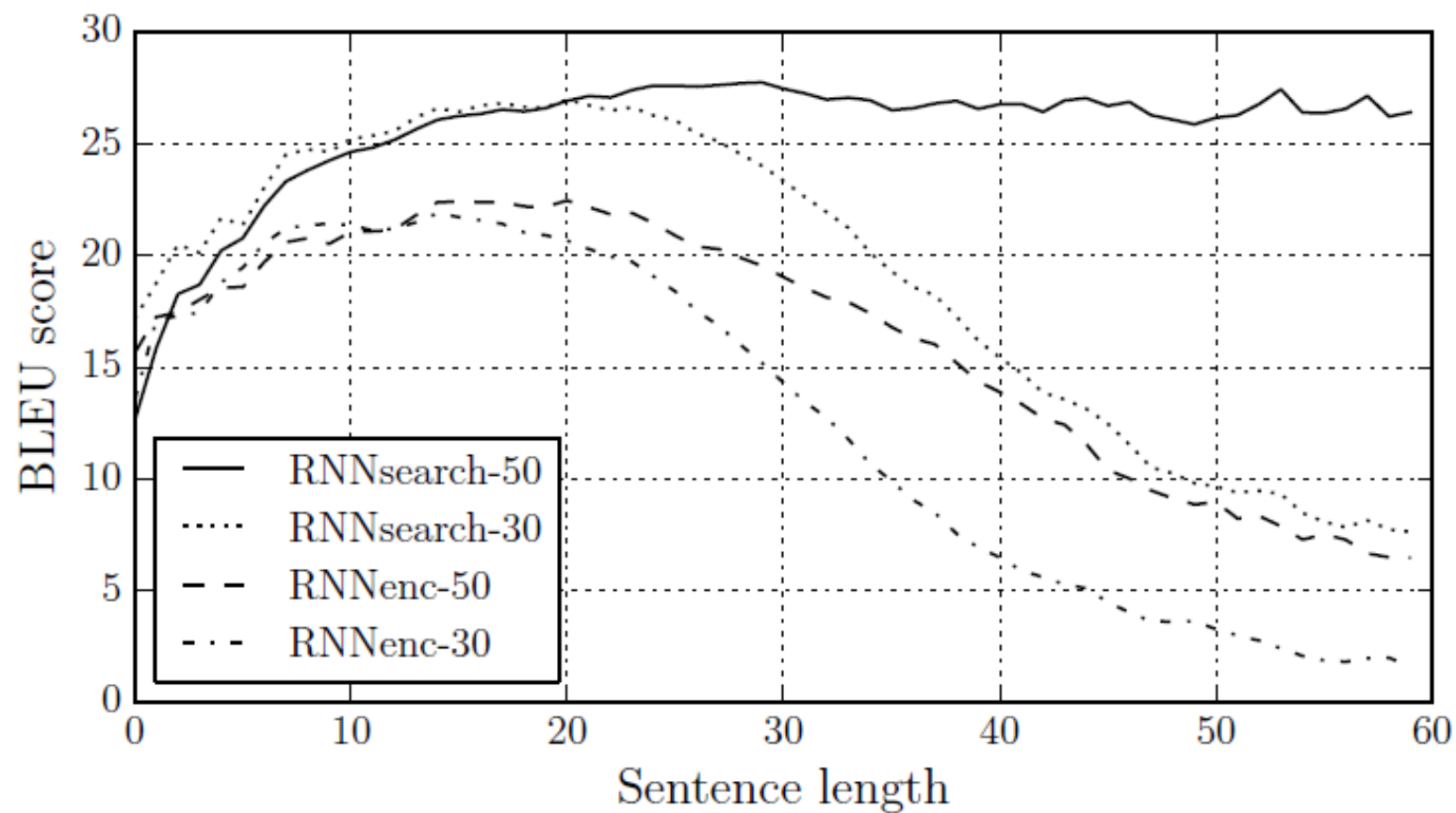
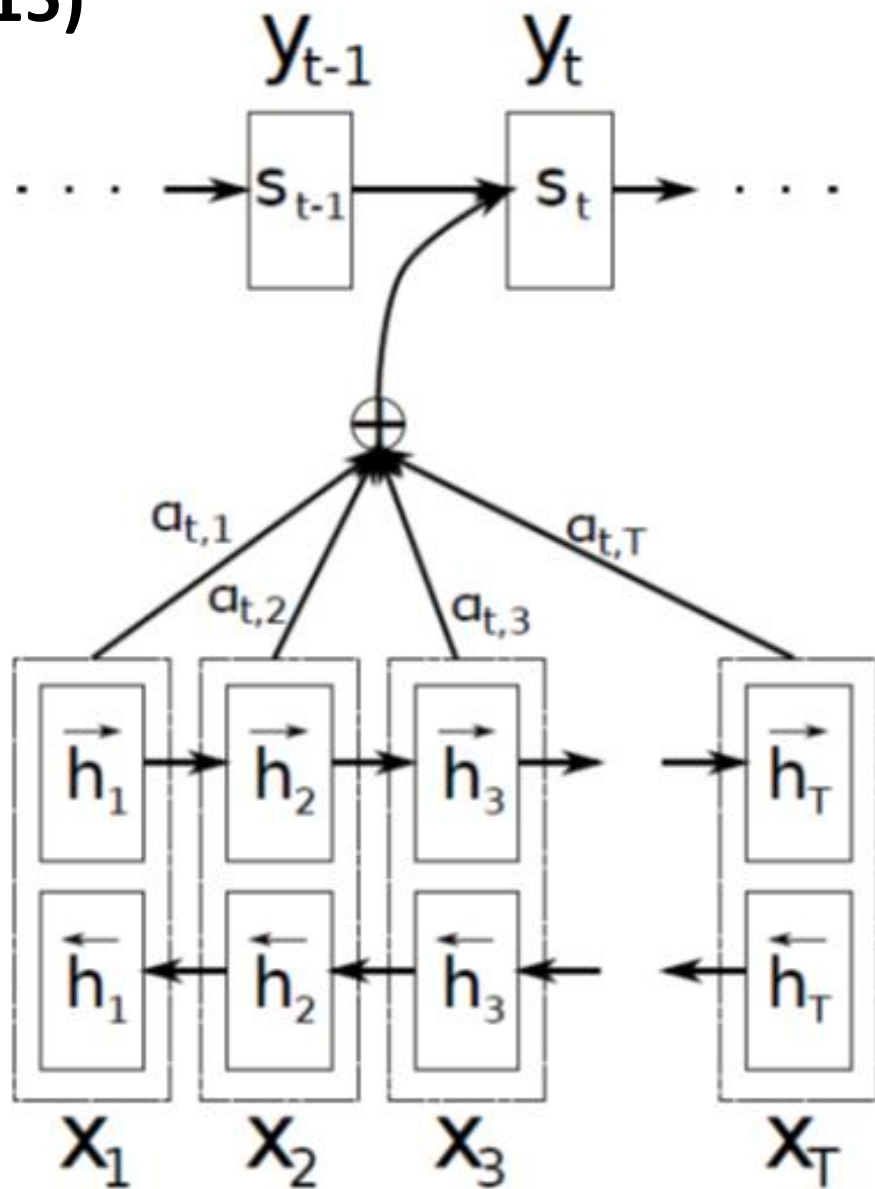


Figure 2: The BLEU scores of the generated translations on the test set with respect to the lengths of the sentences. The results are on the full test set which includes sentences having unknown words to the models.

## Bahdanau Attention (2015)

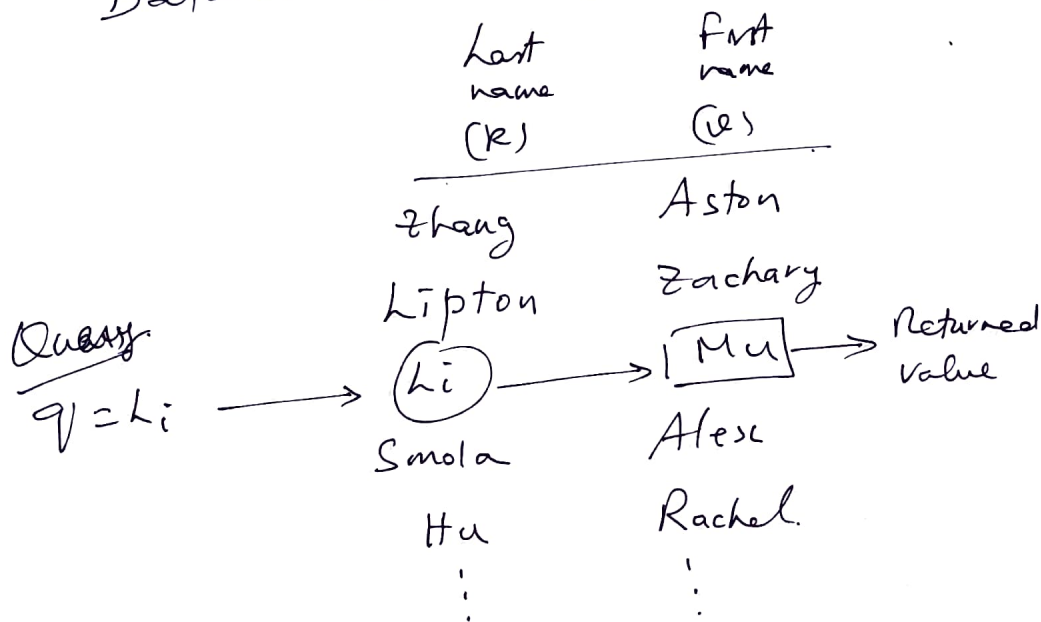


# ATTENTION

①

Query  $Q$ , Key  $K$ , Value  $V \Rightarrow [Q, K, V]$

Database  $D$ : Tuples  $[Keys(K), values(V)]$



Query ' $q$ ' operates on  $(K, V)$  pairs.

Exact Match

Query  $q = L_i$   
 $\Rightarrow (L_i, Mu)$

Approximate Match

$q = L_i$   
 $\Rightarrow (Lipton, Zachary)$

## Attention Mechanism [Bahdanau (2014)]

(2)

$$D = \{(k_1, v_1), (k_2, v_2), \dots, (k_m, v_m)\}$$

D database of  $m$  tuples of (keys, values).

Let  $q$  : query.

Attention over D

$$\text{Attention}(q, D) = \sum_{i=1}^m \alpha(q, k_i) v_i \quad \text{--- (1.)}$$

$\alpha \in \mathbb{R}$  are scalar attention weights

Eqn (1): Attention Pooling.

: Linear Combination of "values" in the db D.

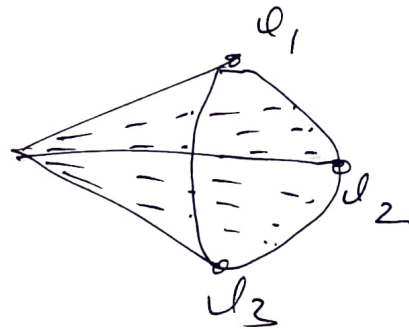
Attention: operation pays "particular attention" to the terms for which weight  $\alpha$  is large.

## Special Cases

$$\text{Attention}(q, D) = \sum_{i=1}^m \alpha(q, k_i) v_i$$

1.  $\alpha(q, k_i) \geq 0$  (non-negative)

o/p is contained in the "convex-cone" spanned by values of  $v_i$



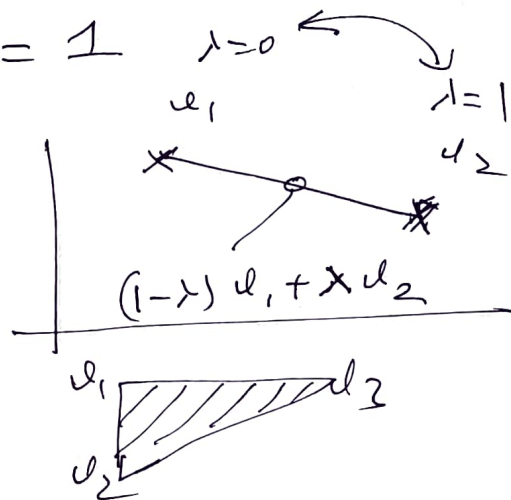
2.  $\alpha(q, k_i) \Rightarrow$  Convex Combination

i.e.  $\alpha(q, k_i) \geq 0$

$$\sum_{i=1}^m \alpha(q, k_i) = 1$$

Attention( $q, D$ ): Lies on Simplex Surface

MOST COMMON SETTING IN  
DEEP LEARNING



3. Exactly one  $\alpha(q, k_i) = 1$   
rest  $\dots = 0$

$\Rightarrow$  Traditional Database Querying

ie  $q \Rightarrow \underline{\underline{d_i \text{ for which } \alpha(q, k_i) = 1}}$

4. All weights are equal

$$\alpha(q, k_i) = \frac{1}{m} \quad \forall i$$

$\Rightarrow$  Averaging the entire Database  $D$ .

$\Rightarrow$  AVERAGE POOLING IN DEEP LEARNING

(4)

Common mechanism to ensure Case (2) | Convex Combination (5)

i.e.  $\sum_{i=1}^m \alpha(q, k_i) = 1$

is to define

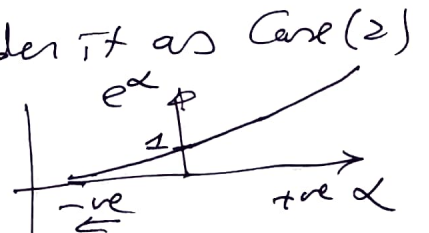
$$\underset{\substack{\uparrow \\ \text{normalized}}}{\alpha(q, k_i)} = \frac{\alpha(q, k_i)}{\sum_j \alpha(q, k_j)}$$

and to ensure also  $\alpha(q, k_i) \geq 0$

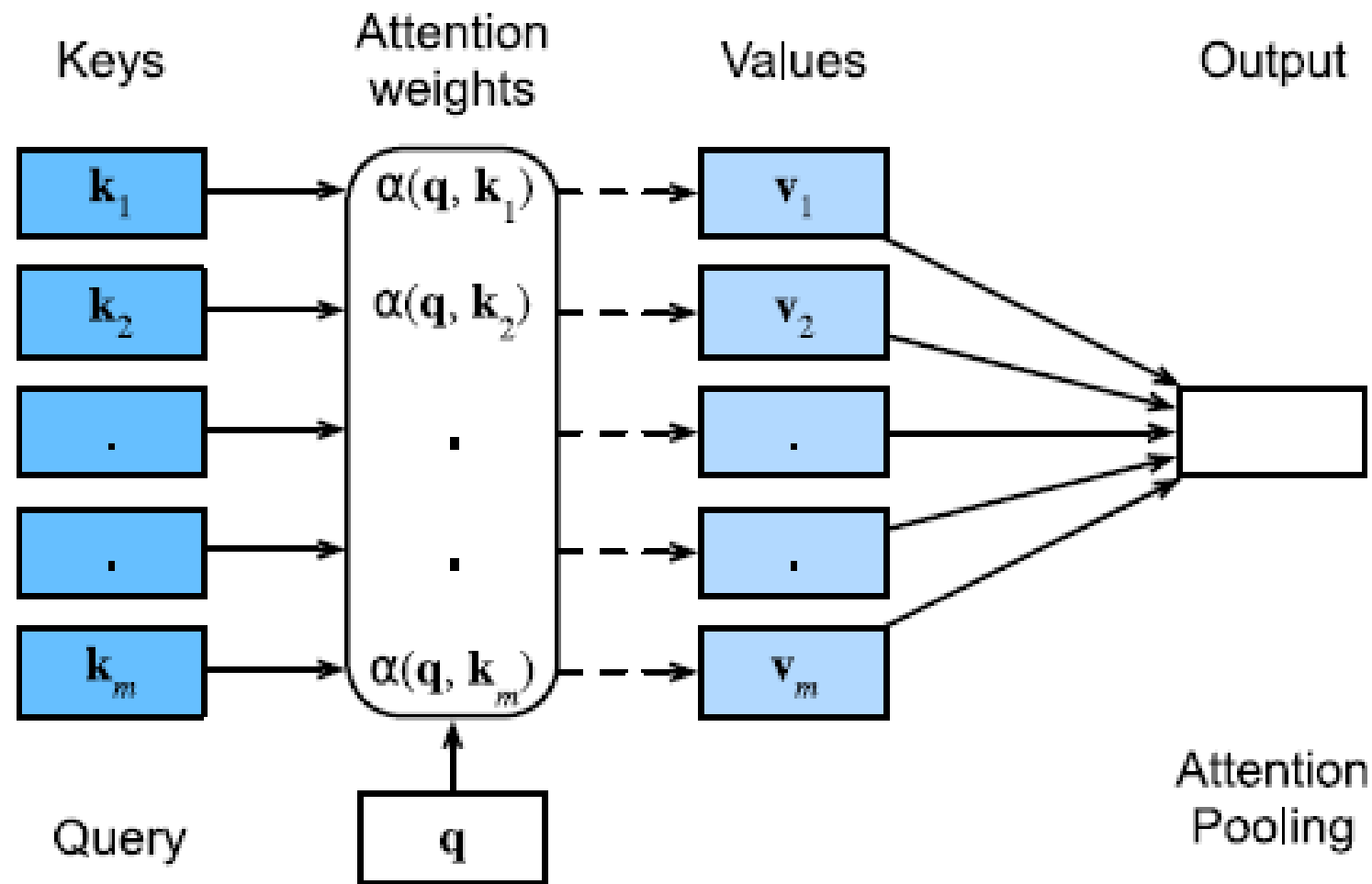
→ Let  $\alpha(q, k_i)$  be "any" function, but render it as Case (2)

$$\alpha(q, k_i) = \frac{\exp[\alpha(q, k_i)]}{\sum_j \exp[\alpha(q, k_j)]}$$

• SOFTMAX operation used for Multinomial models.



- Differentiable
- Gradient never vanishes.



The attention mechanism computes a linear combination over values  $v_i$  via attention pooling, where weights are derived according to the compatibility between a query  $q$  and keys  $k_i$ .



Interpretation of Core ②:  $\alpha$ : Non negative  
: Sums up to 1

⑥

Large weight  $\alpha_s$  as a way for the model to "select"  
 $U_i$  as of "RELEVANCE" to query  $q$ .

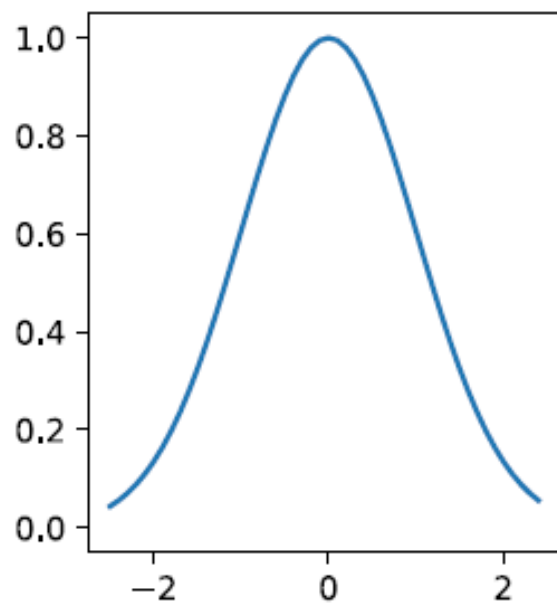
Some common kernels

$$\alpha(q, k) = \exp\left(-\frac{1}{2} \|q - k\|^2\right) \quad \text{Gaussian}$$

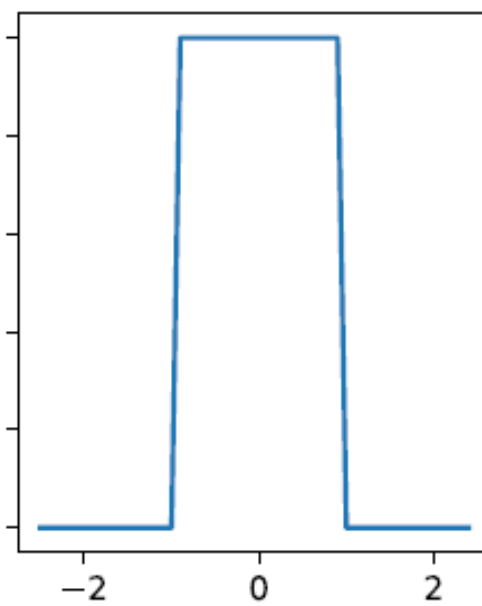
$$\alpha(q, k) = 1 \quad \text{if } \|q - k\| \leq 1 \quad \text{Box Car}$$

$$\alpha(q, k) = \max(0, 1 - \|q - k\|) \quad \text{Epanechnikov.}$$

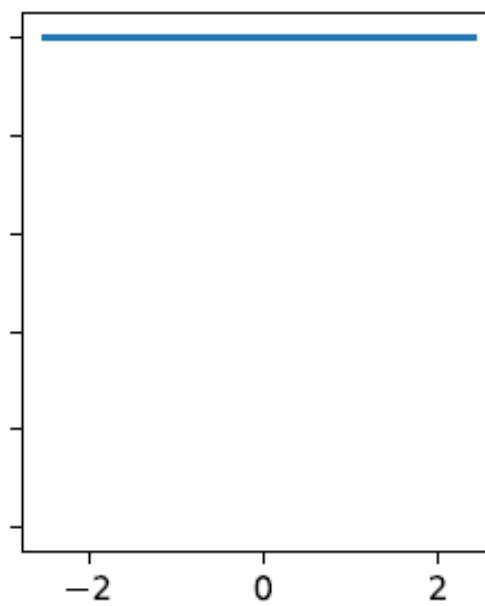
choiced kernels: Related to KDE / Parzen Windows



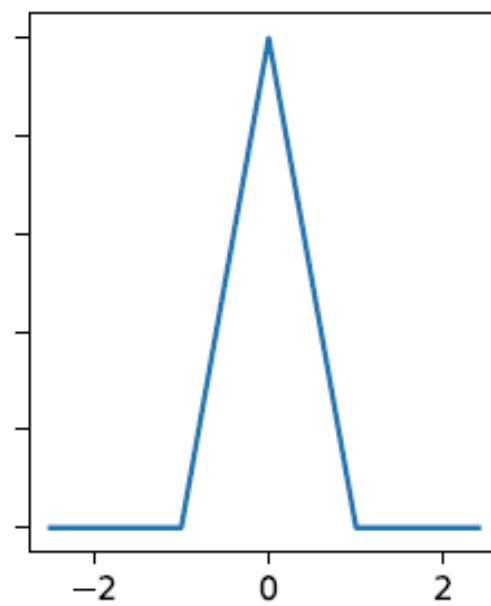
Gaussian



Boxcar

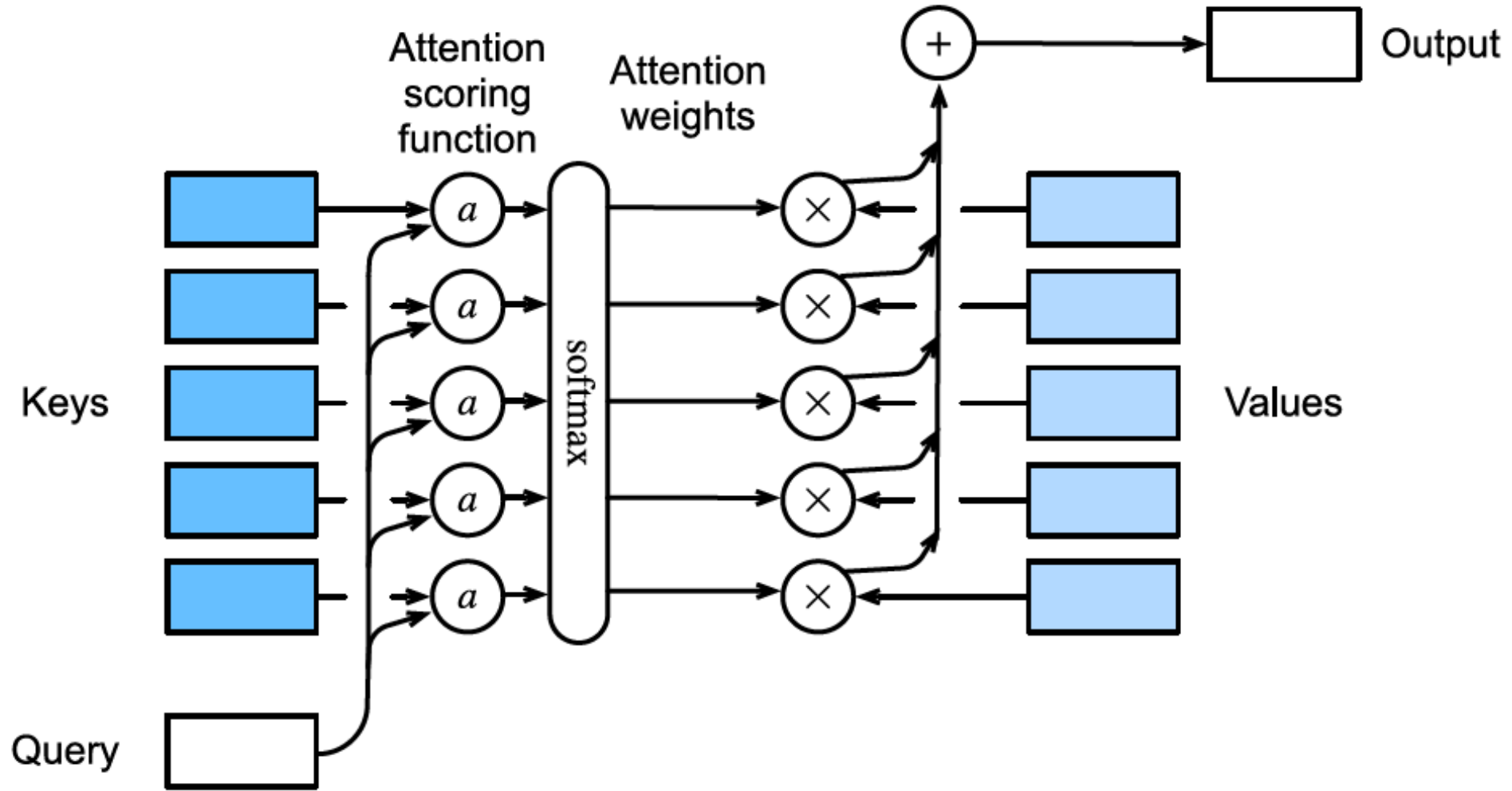


Constant



Epanechnikov

## ATTENTION SCORING FUNCTION



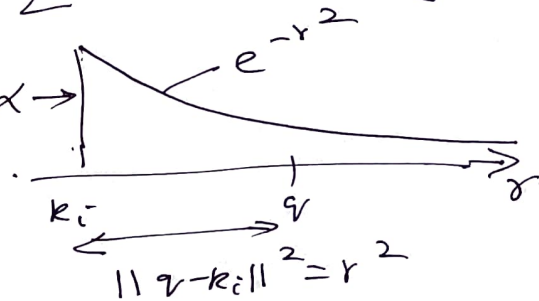
Computing the output of attention pooling as a weighted average of values, where weights are computed with the attention scoring function  $a$  and the softmax operation.

# DOT-PRODUCT ATTENTION

Gaussian kernel

$$\alpha(q, k_i) = \exp \left[ -\frac{1}{2} \|q - k_i\|^2 \right]$$

$$r = \|q - k_i\| \Rightarrow \alpha = e^{-r^2} \Rightarrow \alpha \rightarrow$$



w/o exponentiation

$$\alpha(q, k_i) = -\frac{1}{2} \|q - k_i\|^2$$

$$= q^T k_i - \frac{1}{2} \|k_i\|^2 - \frac{1}{2} \|q\|^2$$

from

$$\|q - k_i\|^2 = \|q\|^2 + \|k_i\|^2 - 2q^T k_i$$

$$\alpha(q, k_i) = q^T k_i - \frac{1}{2} \|k_i\|^2 - \frac{1}{2} \|q\|^2 \quad (9)$$

Only term  
that matters

$\Rightarrow$  Batch &  
Layer  
Normalization.

$\Rightarrow$  Lead to  
activations  
that have  
"constant" norms  
 $\|k_i\|$   
in DL Context

$\therefore$  Can be dropped

$$\therefore \boxed{\alpha(q, k_i) = q^T k_i}$$

- $\hookrightarrow$  Depends only on "q"
- Can be dropped  
(identical for all  $(q, k_i)$ )
- Show that "normalization"  
as in "softmax"  
cancels this term in  
Numerator & Denominator.

$$q \in \mathbb{R}^d$$

$$k_i \in \mathbb{R}^d$$

$$\alpha(q, k_i) = q^T k_i \quad \begin{array}{l} q \in \mathbb{R}^d \\ k_i \in \mathbb{R}^d \end{array} \quad (10)$$

To ensure Variance of the dot-product  
remains 1 regardless of vector lengths  
(dim  $d$ )

$\Rightarrow$  Use "Scaled Dot Product Attention" Scoring fn.

$$\text{ie } \alpha(q, k_i) = \frac{q^T k_i}{\sqrt{d}}$$

First Commonly used  
Attn. Fn. in  
TRANSFORMERS  
(Vaswani, 2017.)

$$\alpha(q, k_i) = \frac{q^T k_i}{\sqrt{d}}$$

(11)

2 Normalize attention weights by softmax operation.

$$\begin{aligned} \alpha(q, k_i) &= \text{Softmax} [\alpha(q, k_i)] \\ &= \frac{\exp(q^T k_i / \sqrt{d})}{\sum_{j=1}^m \exp(q^T k_j / \sqrt{d})} \end{aligned}$$

If  $q \in \mathbb{R}^d$  &  $k_i \in \mathbb{R}^{d'}$  i.e.  $q$  &  $k_i$  are not in same dim.

Then, use  $q^T M k$   $\rightarrow$  Matrix of suitable dim (e.g.  $d \times d'$ )  
 $\rightarrow$  Linear Projection by a Fully Connected (FC) layer.

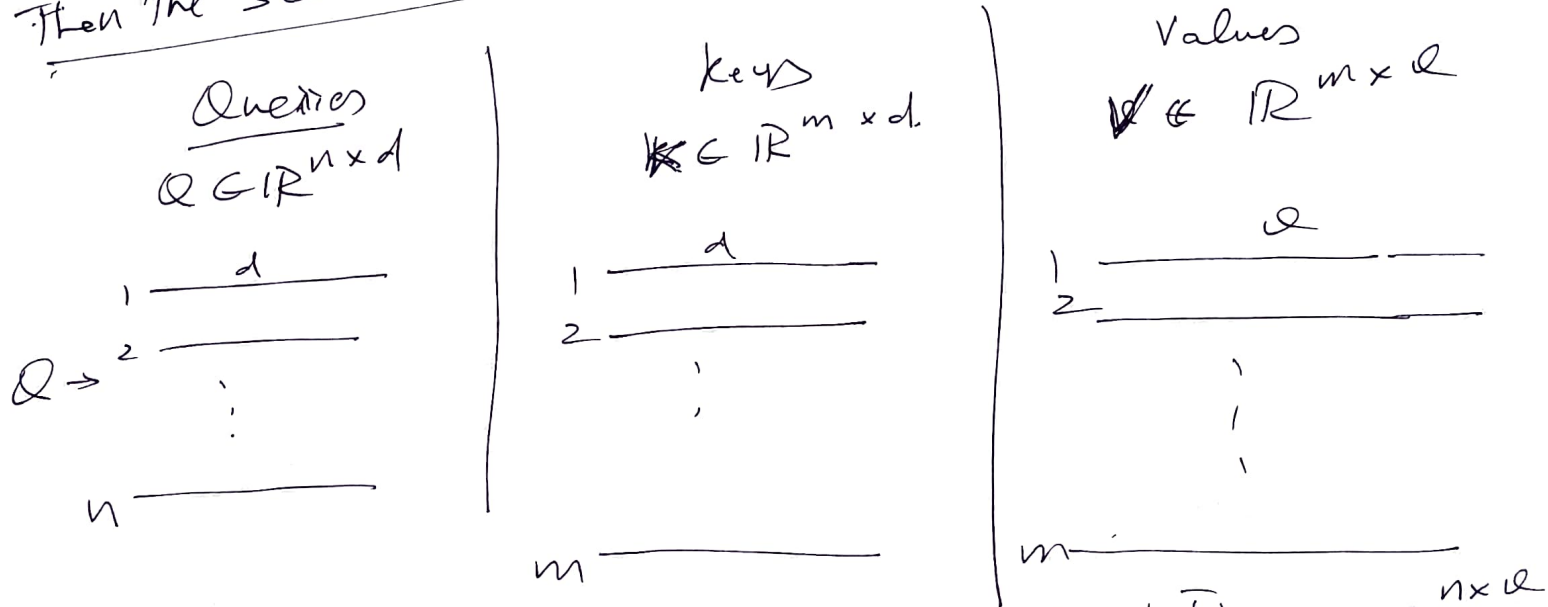
$n$  queries  
 $m$  key-value pairs

$q_i, k_i : \dim d \in \mathbb{R}^d$   
 values  $v_i : \dim v \in \mathbb{R}^v$

(12)

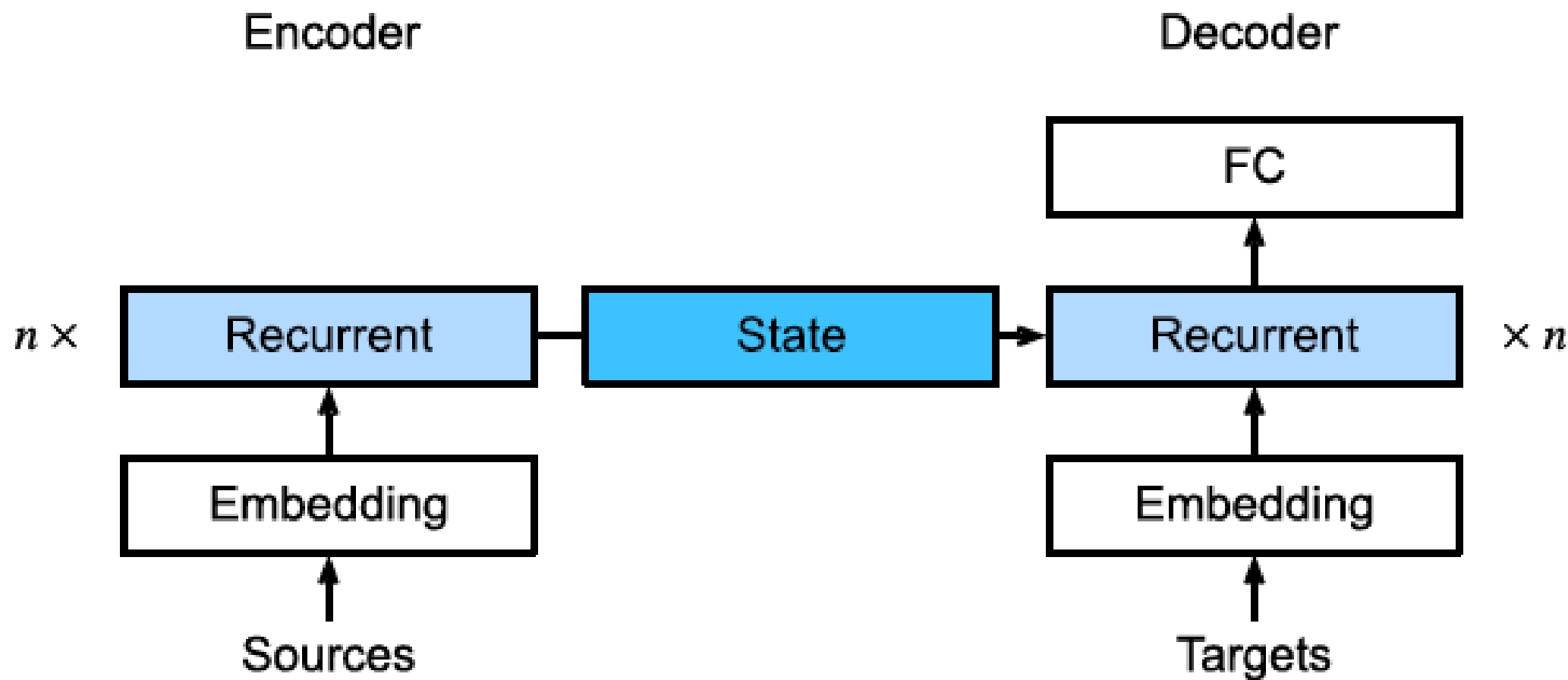
ie for ~~mini~~ batch of  $n$  queries

Then The Scaled Dot Prod. Attn of Queries



$$\text{Scaled DPA} \Rightarrow \alpha(Q, K) \Rightarrow \text{Softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{n \times v}$$





Sequence-to-sequence model. The state, as generated by the encoder, is the only piece of information shared between the encoder and the decoder.

# Bahdanau Attention (2014)

13

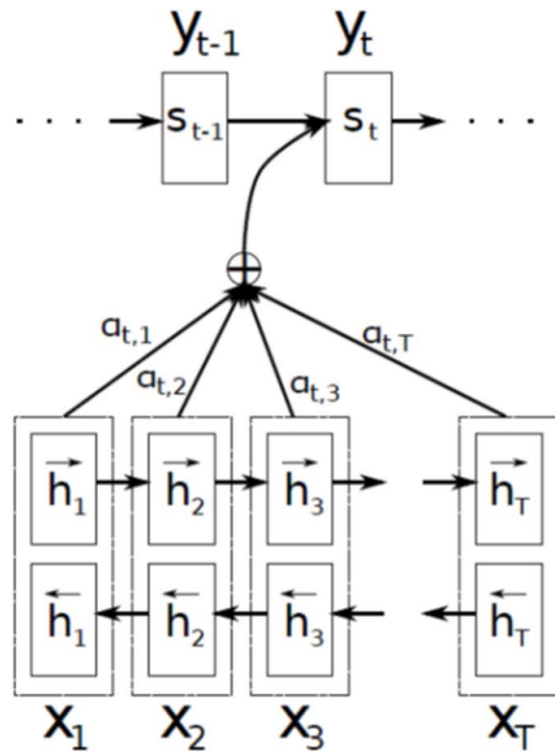
## ▷ "Learning to Attend"

- Predict a token in o/p @ decoder
- If not all the i/p tokens are relevant

→ The model always ["attends"]

only to parts of the i/p

sequence that are considered "relevant" to the current prediction.



Bahdanau A-M  $\Rightarrow$  Most influential idea in past decade in DL.

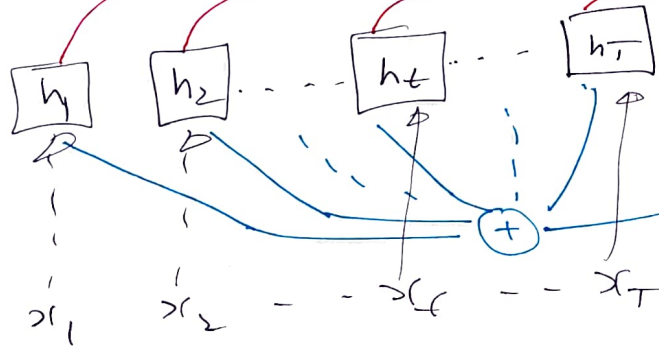
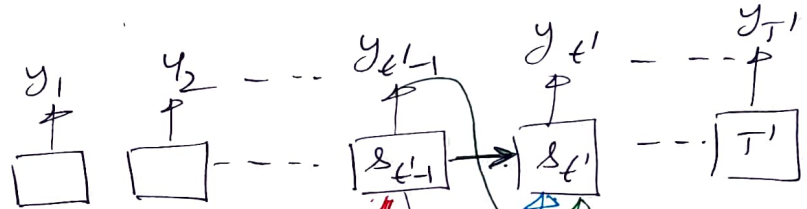
$\Rightarrow$  Giving rise to TRANSFORMERS [Vaswani, 2017]

RAM Model

Takes the place of a single context  $C = f(h_1, h_2, \dots, h_T)$

$$C_{t'} = \sum_{t=1}^T \alpha(\delta_{t'-1}, h_t) h_t$$

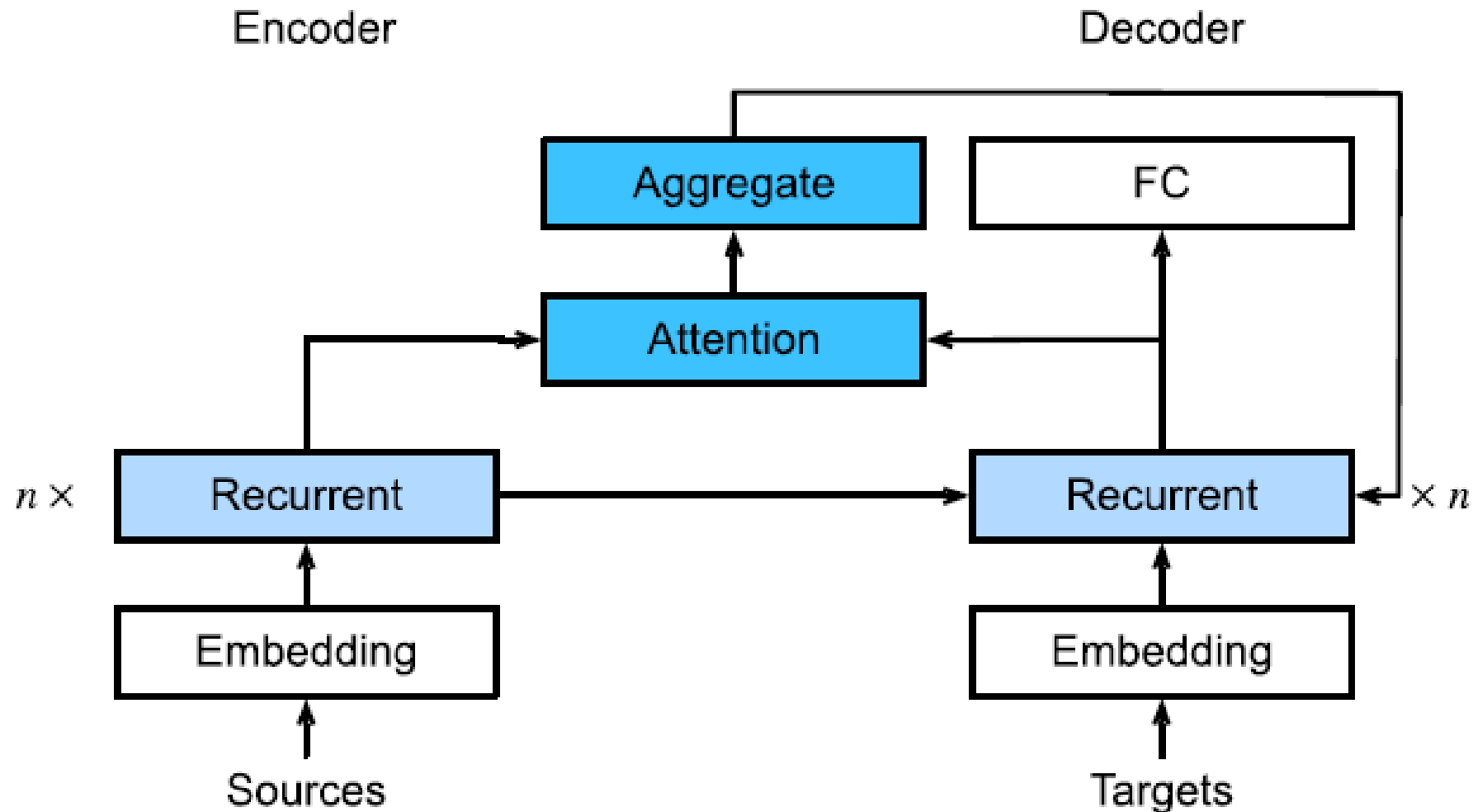
$\nearrow$  QUERY  $\nearrow$  KEY  $\nearrow$  VALUE  
 $\delta_{t'-1}$   $h_t$



Predict

$$\delta_{t'} = g(y_{t'-1}, C_{t'}, \delta_{t'-1})$$

## Bahdanau Attention (2015)

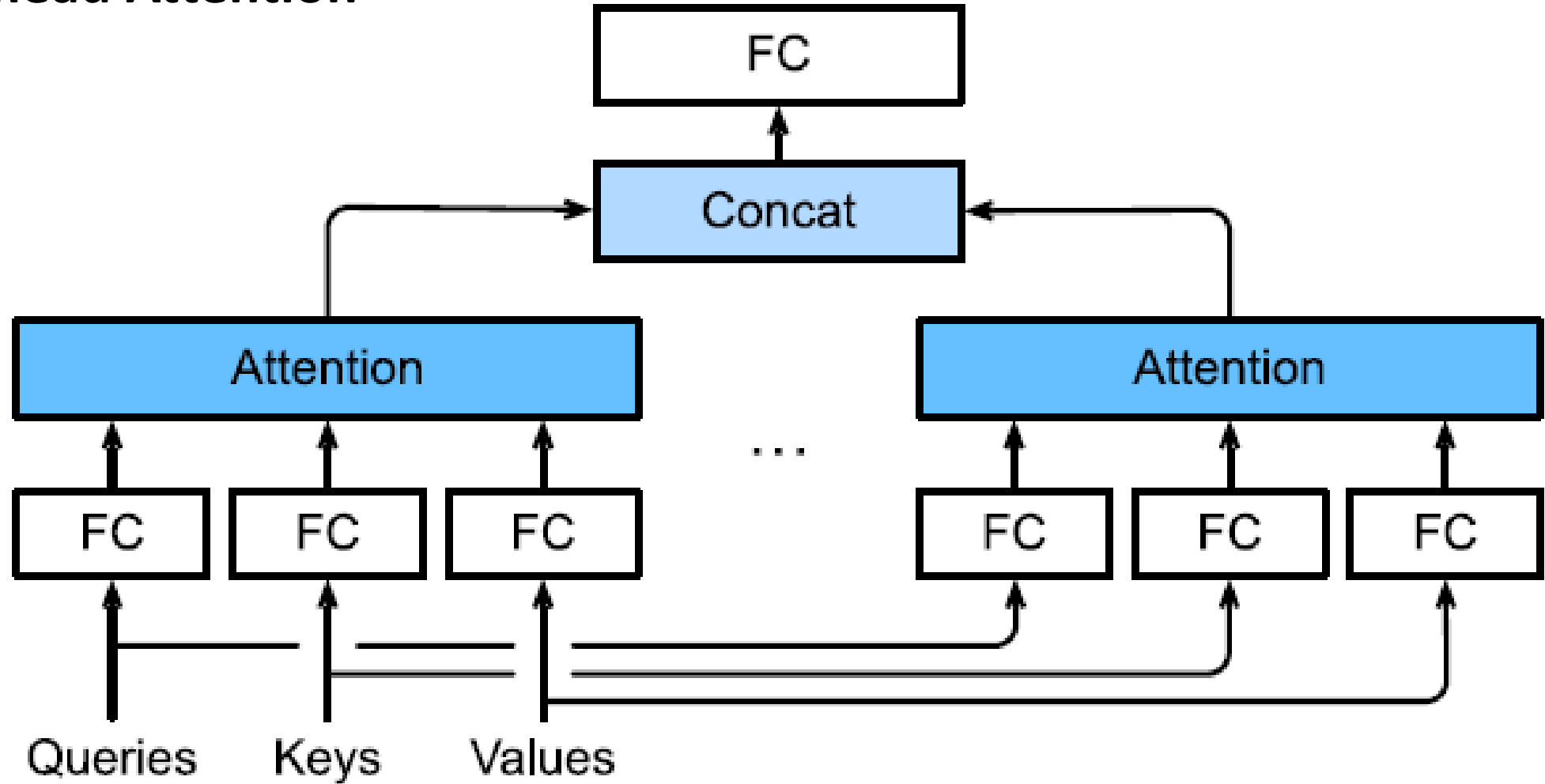


Layers in an RNN encoder–decoder model with the Bahdanau attention mechanism.

15.

- These "h" projected  $Q, K, V_k$  are fed into Attention Pooling in Parallel.
- Each of "h" attn pooling o/p = "head"

## Multi-head Attention



Multi-head attention, where multiple heads are concatenated then linearly transformed.

# SELF-ATTENTION (SA)

16

- Every token (in an i/p seq)  
attends to every other token.

[ Unlike where Decoder Steps attend to  
Encoder Steps ] as in B.A.M.

Thank you !!