

# Modelos de Linguagem Neurais Com Transformers

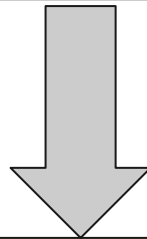
*"The true art of memory is the art of attention "* Samuel Johnson, Idler #74,  
September 1759



**Profa Aline Paes**  
**alinepaes@ic.uff.br**

# Tradução

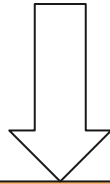
E a vida o que é, diga lá, meu irmão?



¿Y la vida, qué es? Dime tú, hermano.

# Encoder-Decoder

E a vida o que é, diga lá, meu irmão?



¿Y la vida, qué es? Dime tú, hermano.

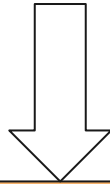
Encoder



Decoder

# Encoder-Decoder

E a vida o que é, diga lá, meu irmão?



¿Y la vida, qué es? Dime tú, hermano.

Constrói uma representação da entrada e passa para o decoder

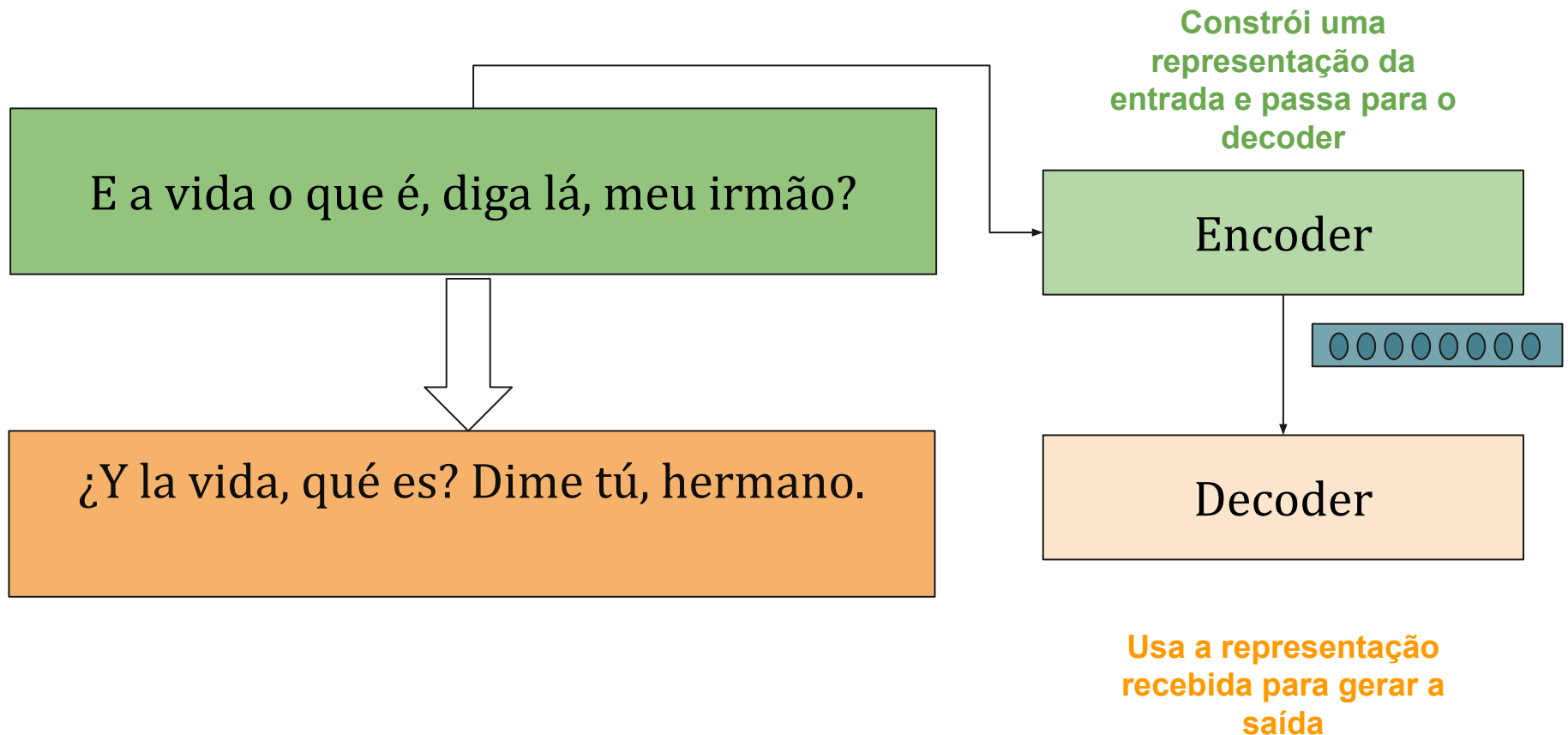
Encoder



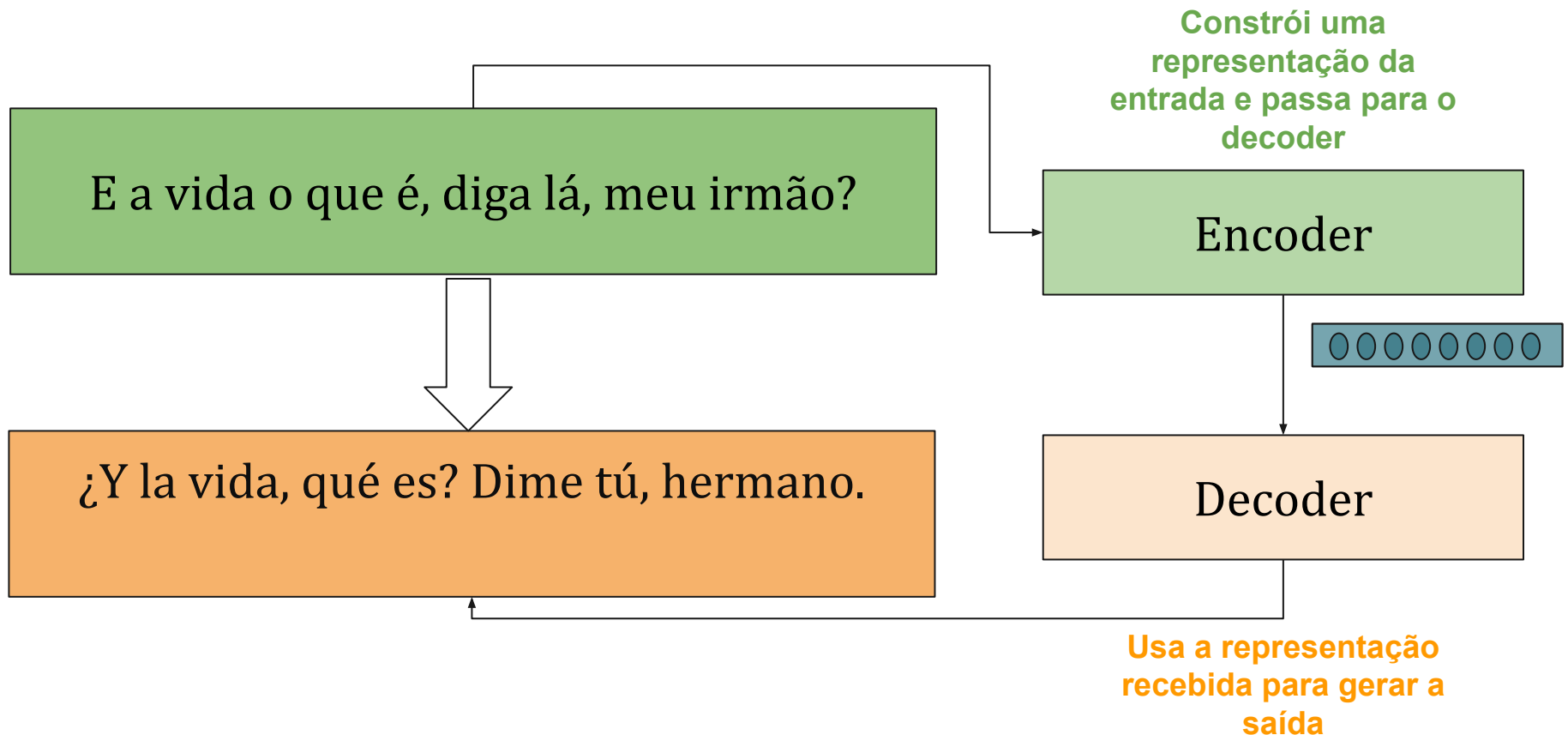
Decoder

Usa a representação recebida para gerar a saída

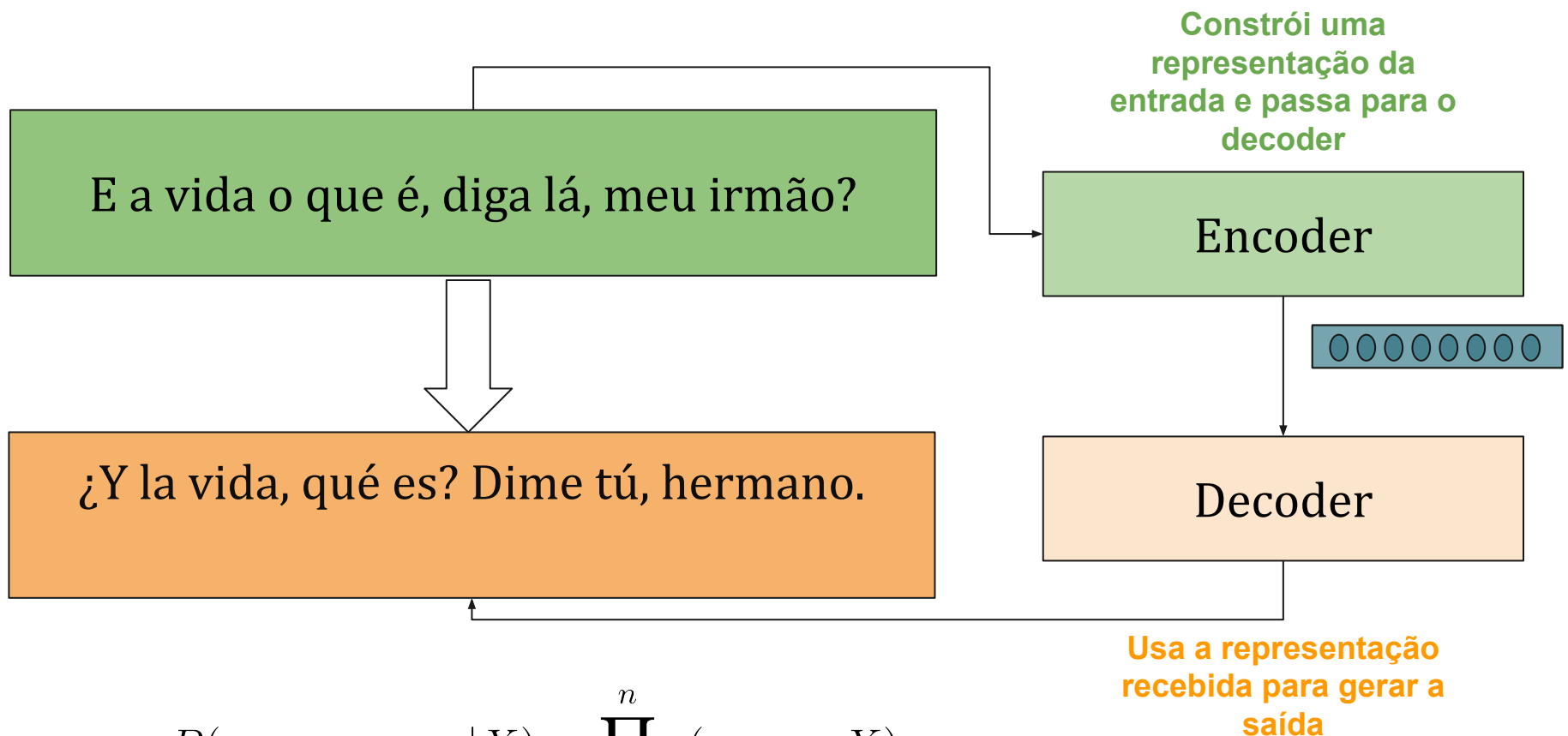
# Encoder-Decoder



# Encoder-Decoder

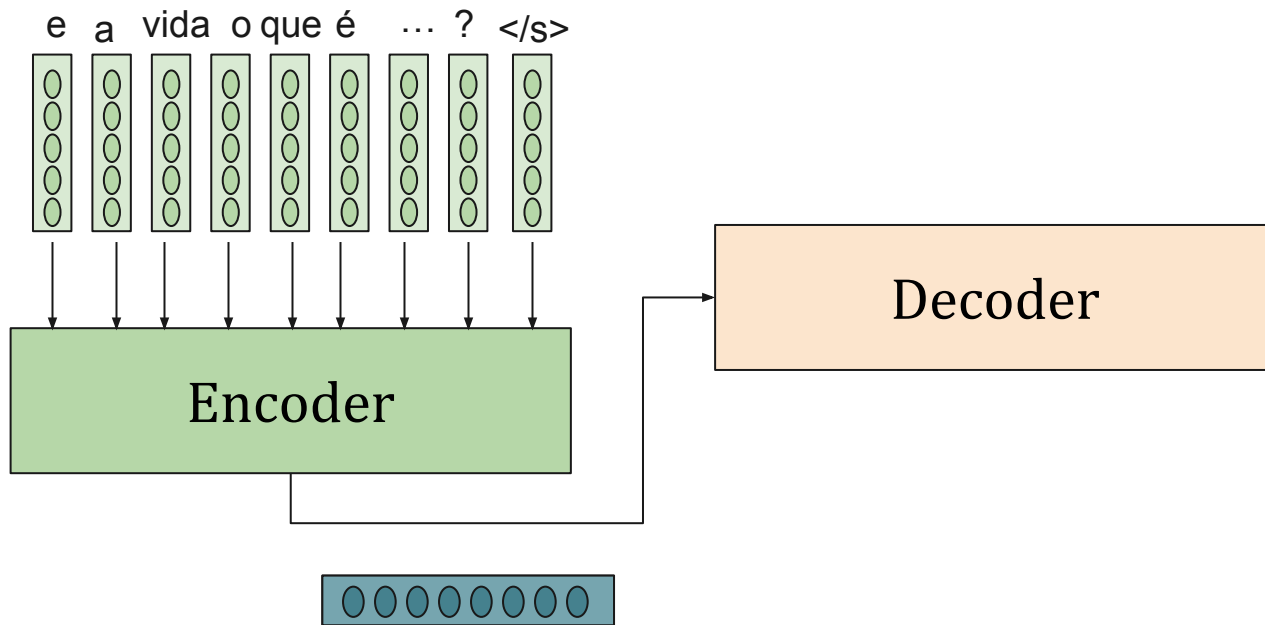


# Encoder-Decoder



$$P(y_1, y_2, \dots, y_n | X) = \prod_{t=1}^n p(y_t, y_{<t}, X)$$

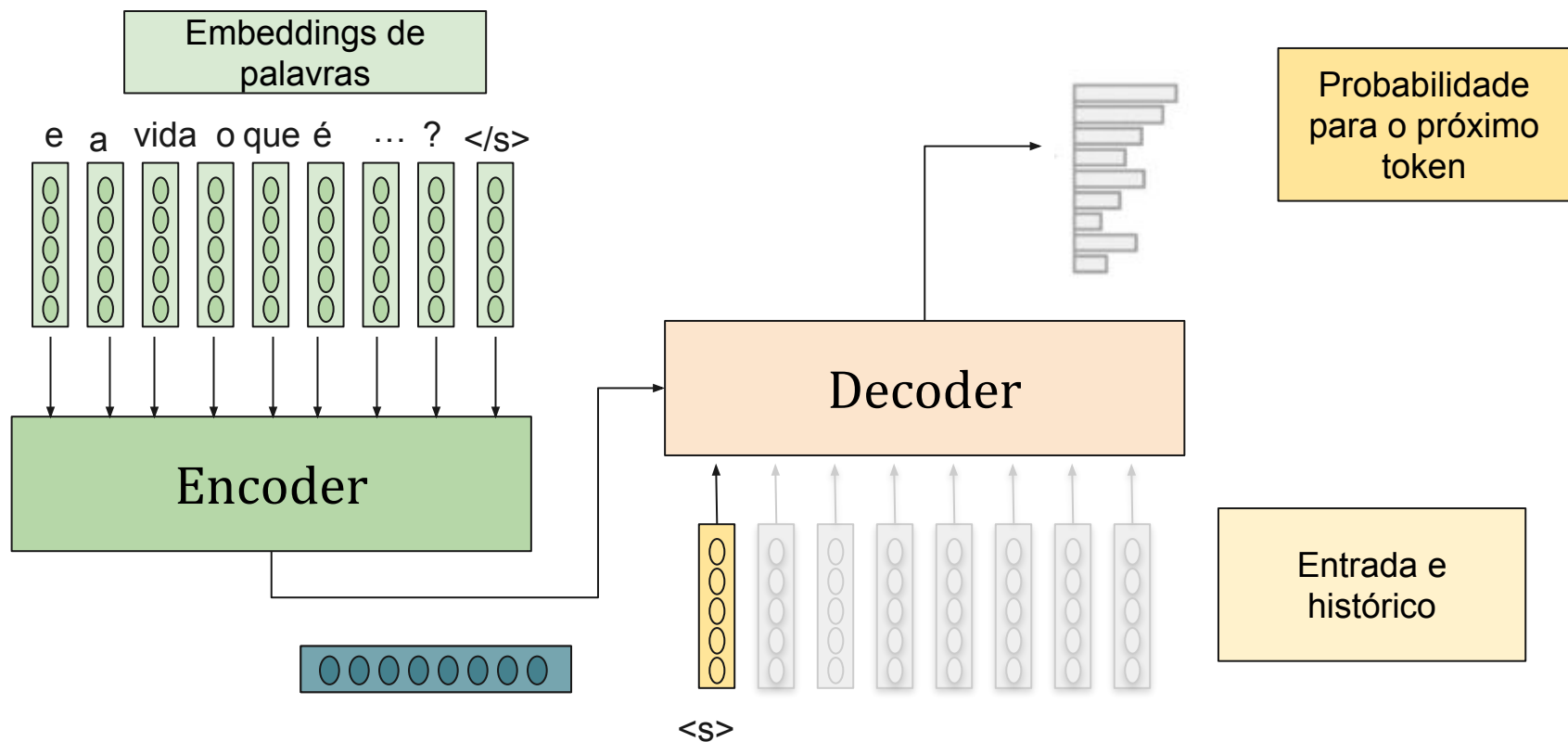
# Encoder-Decoder



$$P(y_1, y_2, \dots, y_n | X) = \prod_{t=1}^n p(y_t, y_{<t}, X)$$

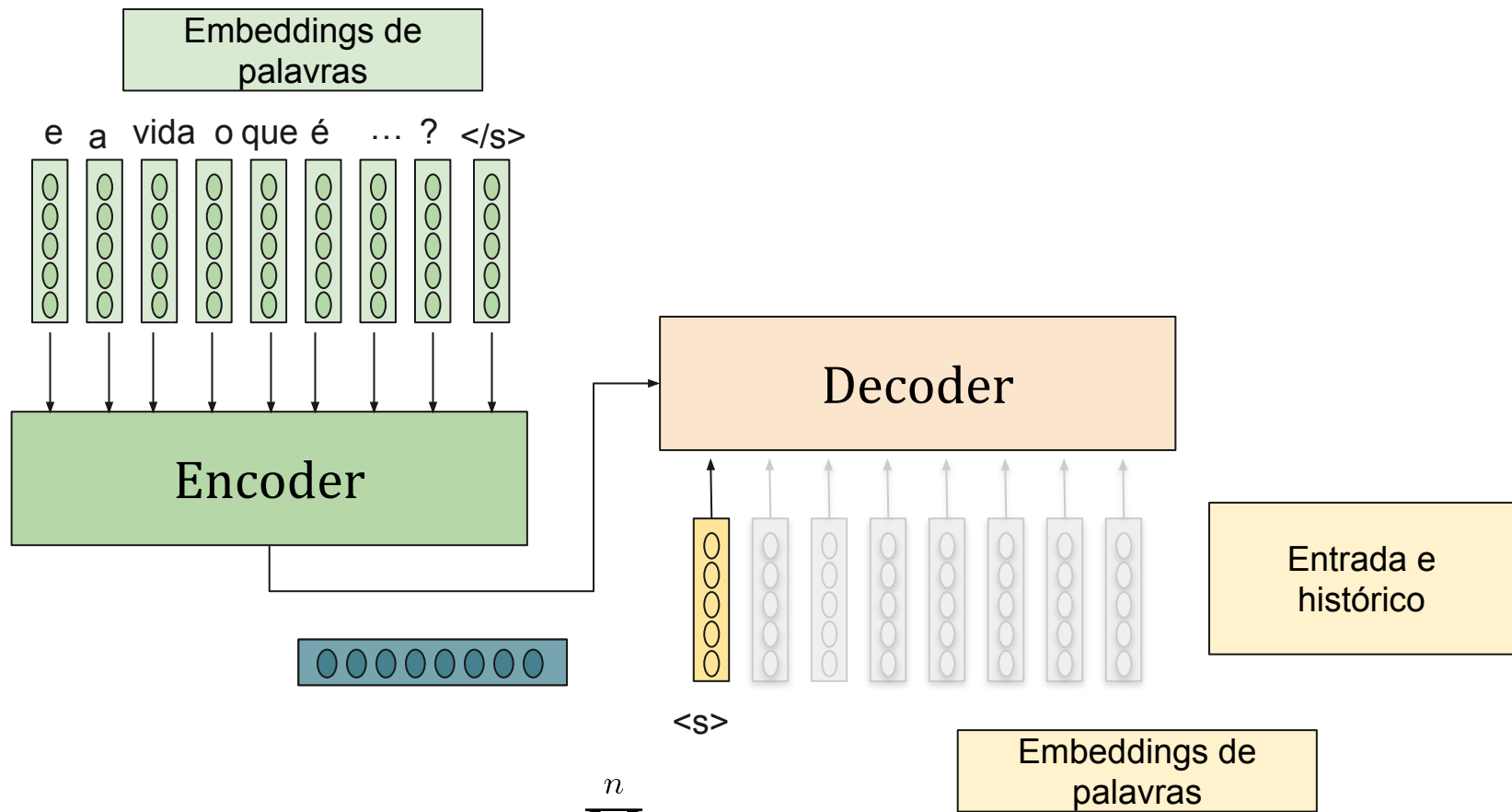


# Encoder-Decoder



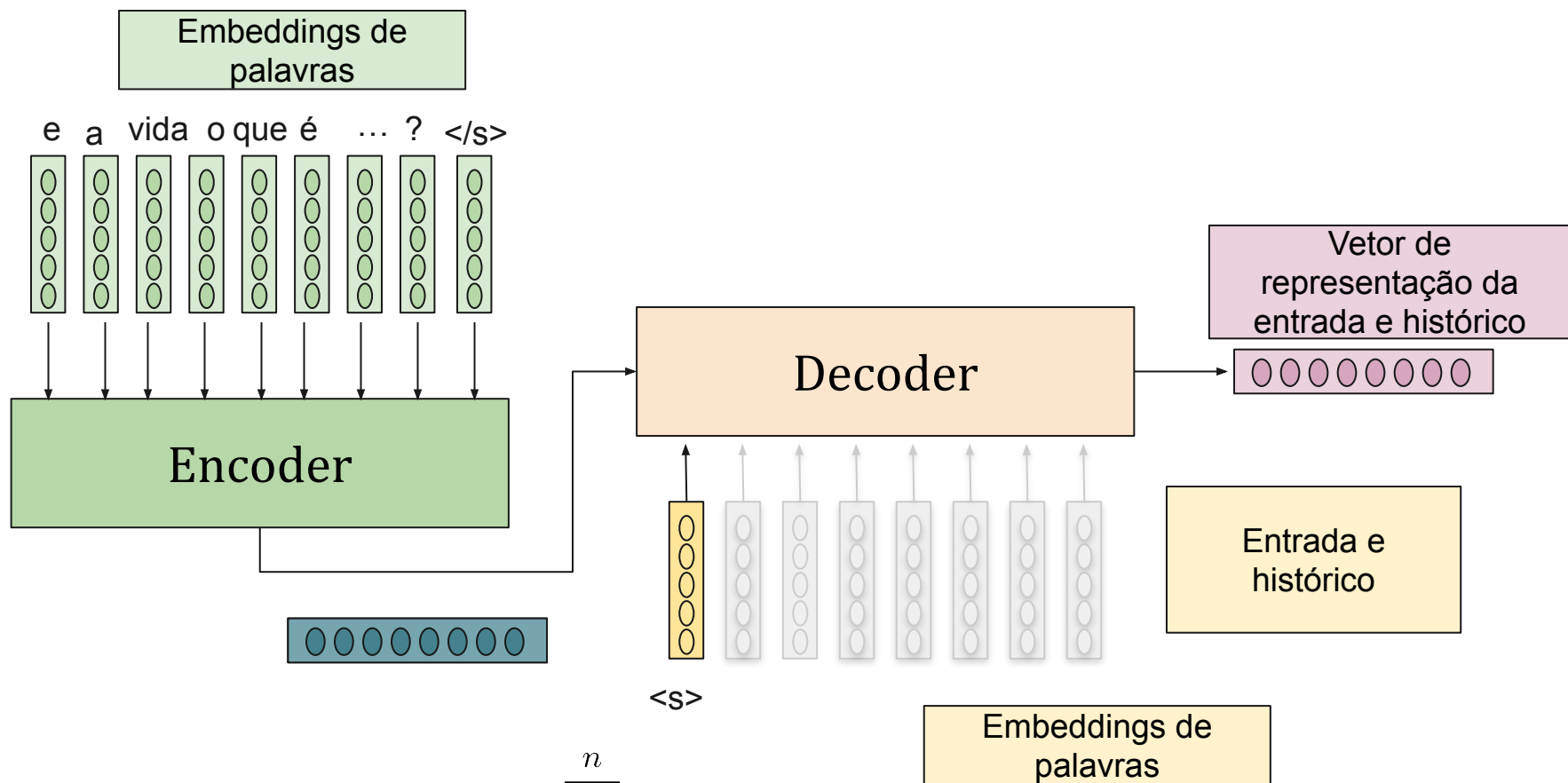
$$P(y_1, y_2, \dots, y_n | X) = \prod_{t=1}^n p(y_t, y_{<t}, X)$$

# Encoder-Decoder



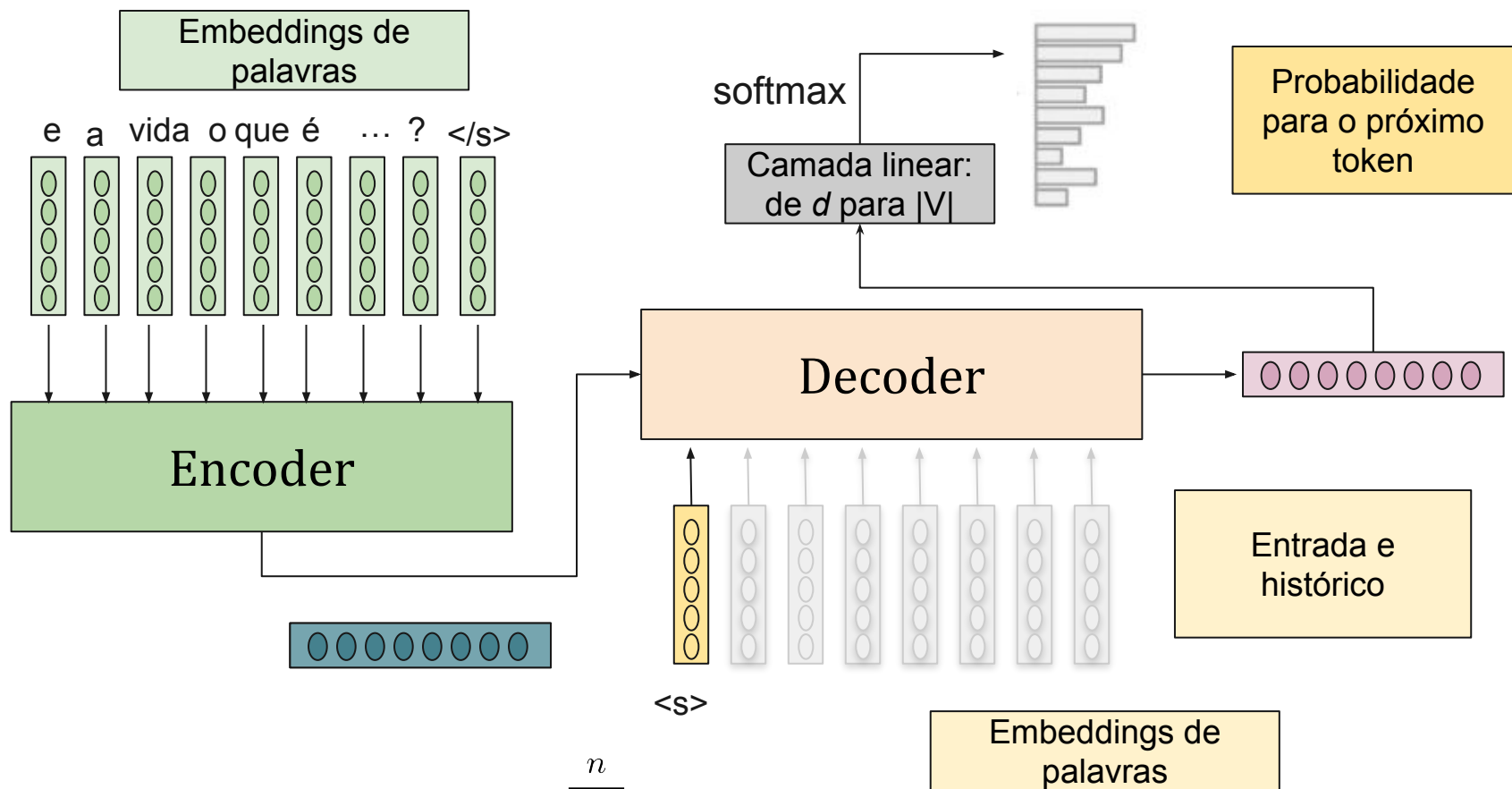
$$P(y_1, y_2, \dots, y_n | X) = \prod_{t=1}^n p(y_t, y_{<t}, X)$$

# Encoder-Decoder



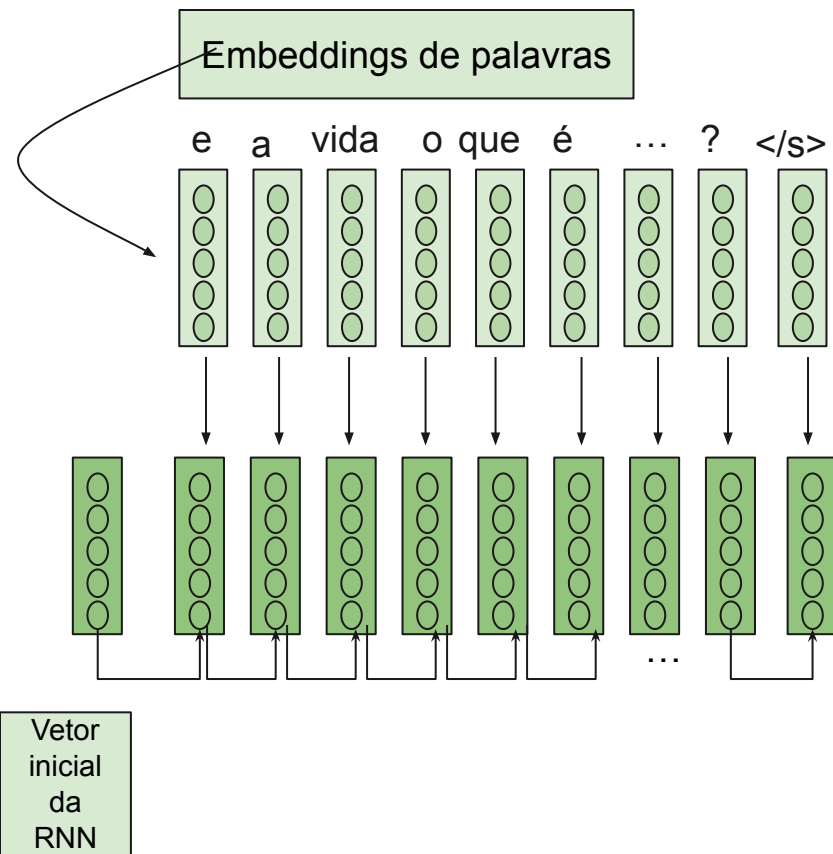
$$P(y_1, y_2, \dots, y_n | X) = \prod_{t=1}^n p(y_t, y_{<t}, X)$$

# Encoder-Decoder

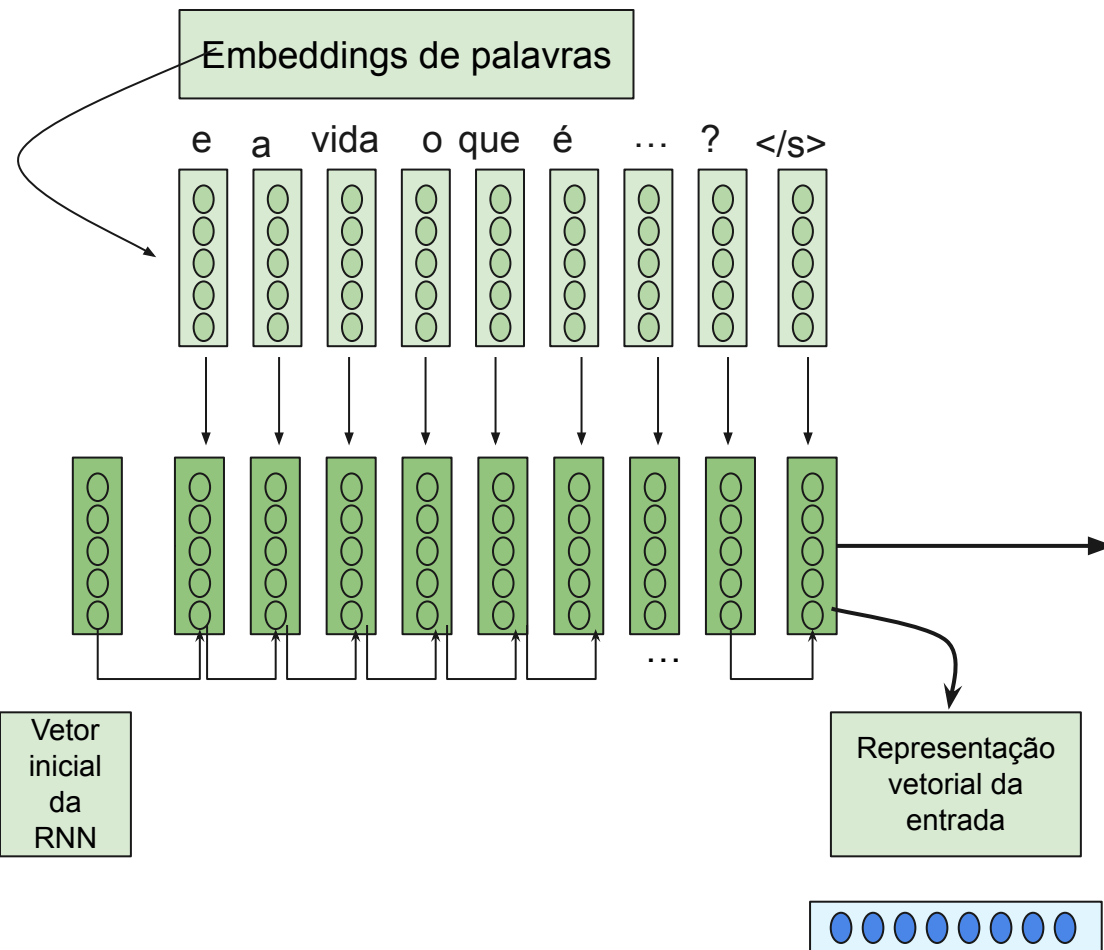


$$P(y_1, y_2, \dots, y_n | X) = \prod_{t=1}^n p(y_t, y_{<t}, X)$$

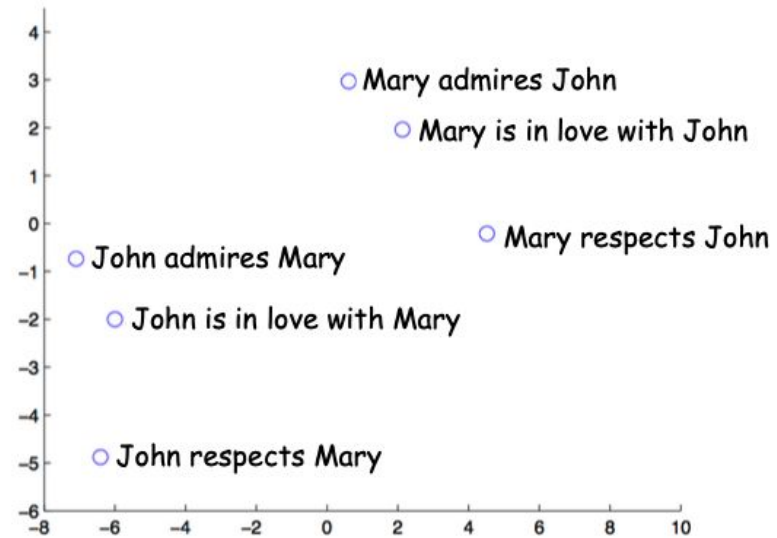
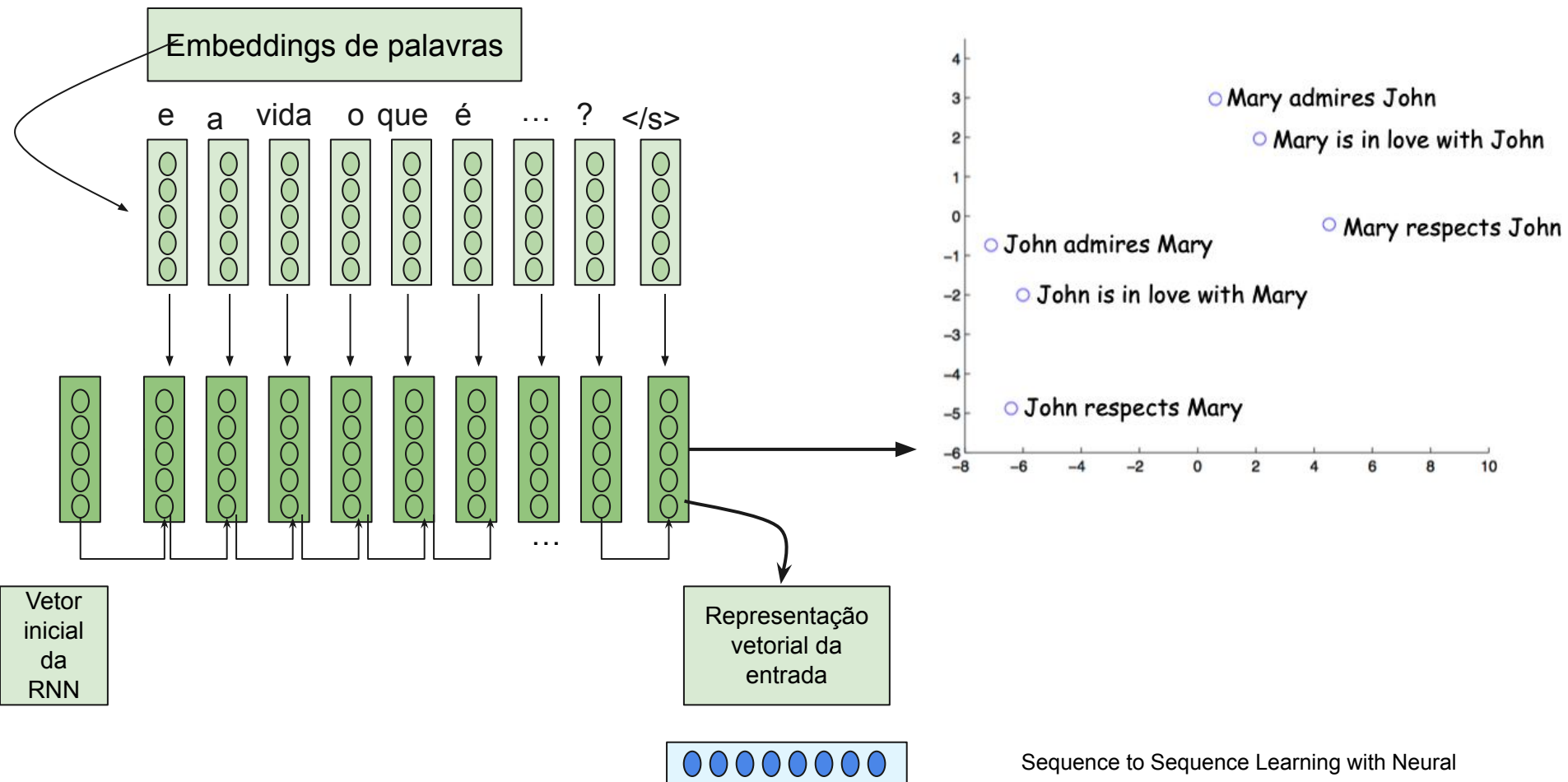
# Encoder-Decoder com RNN



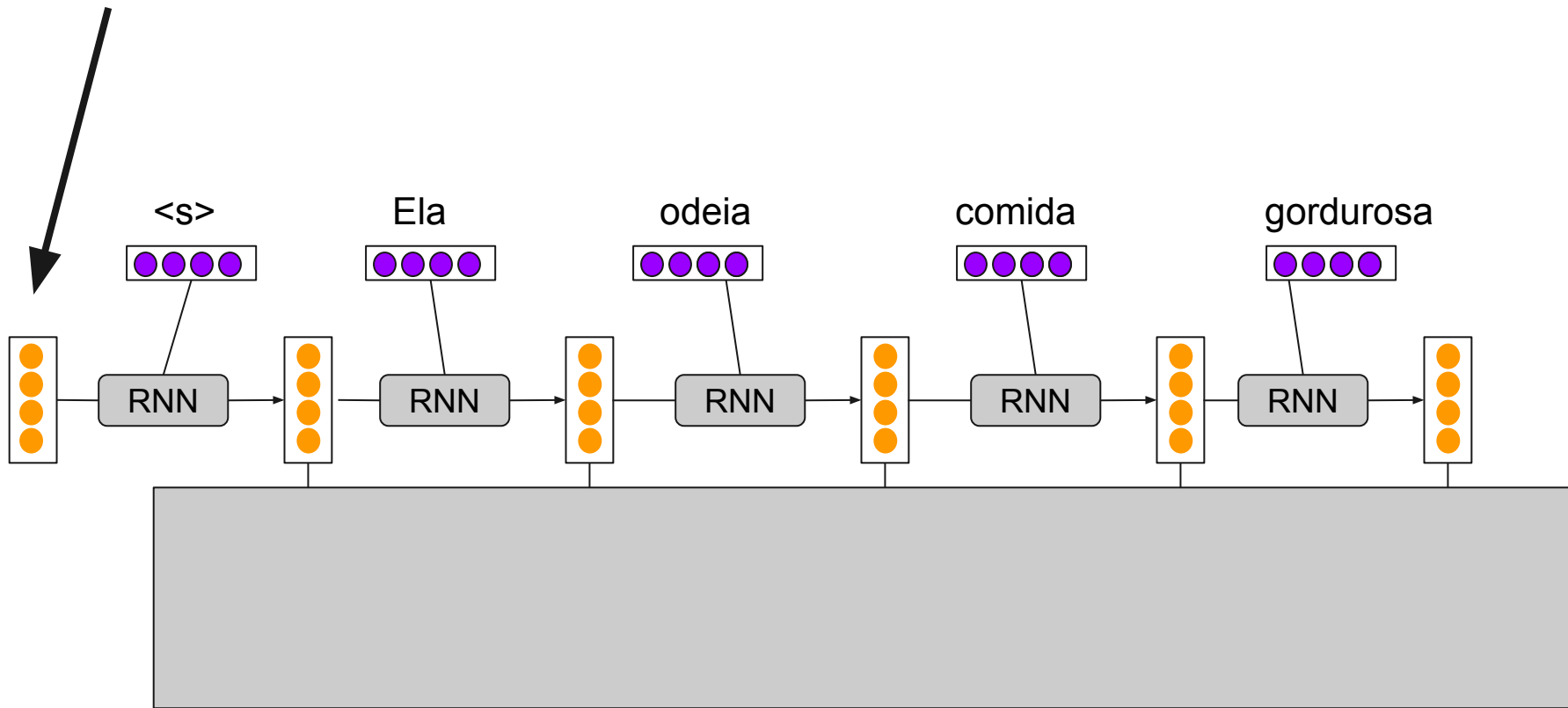
# Encoder-Decoder com RNN



# Encoder-Decoder com RNN



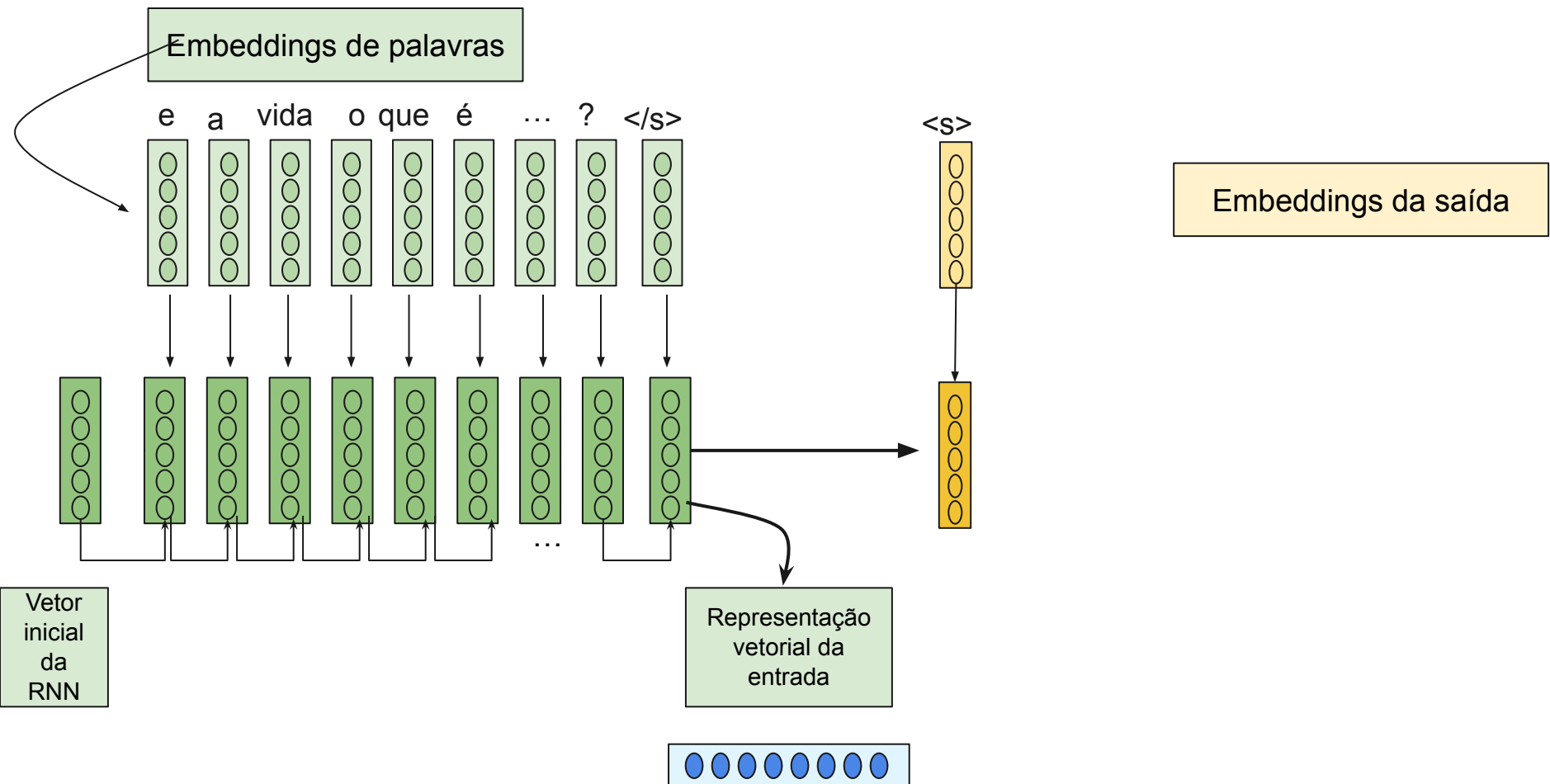
# Inferência - geração de texto



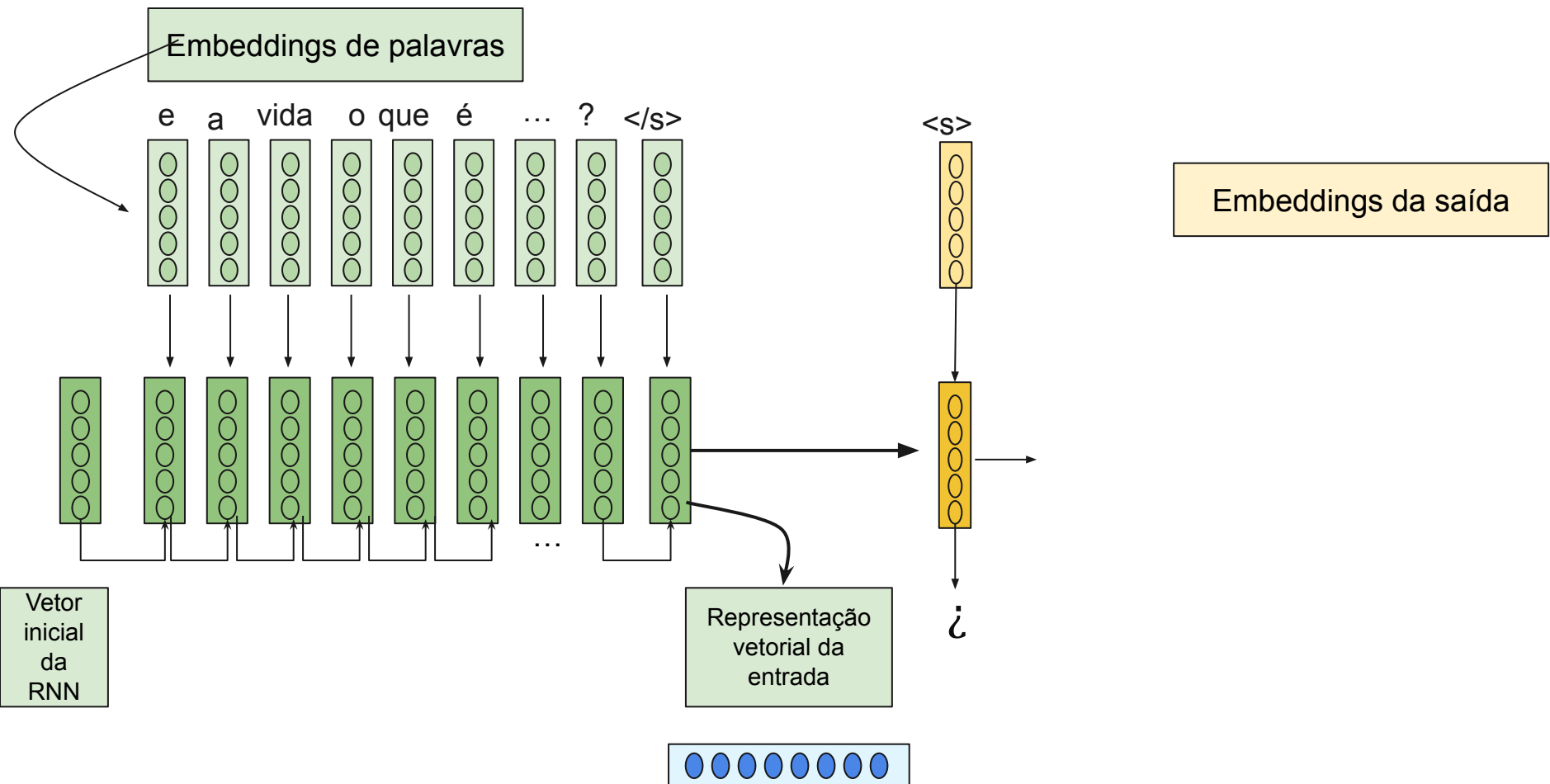
Geração autorregressiva



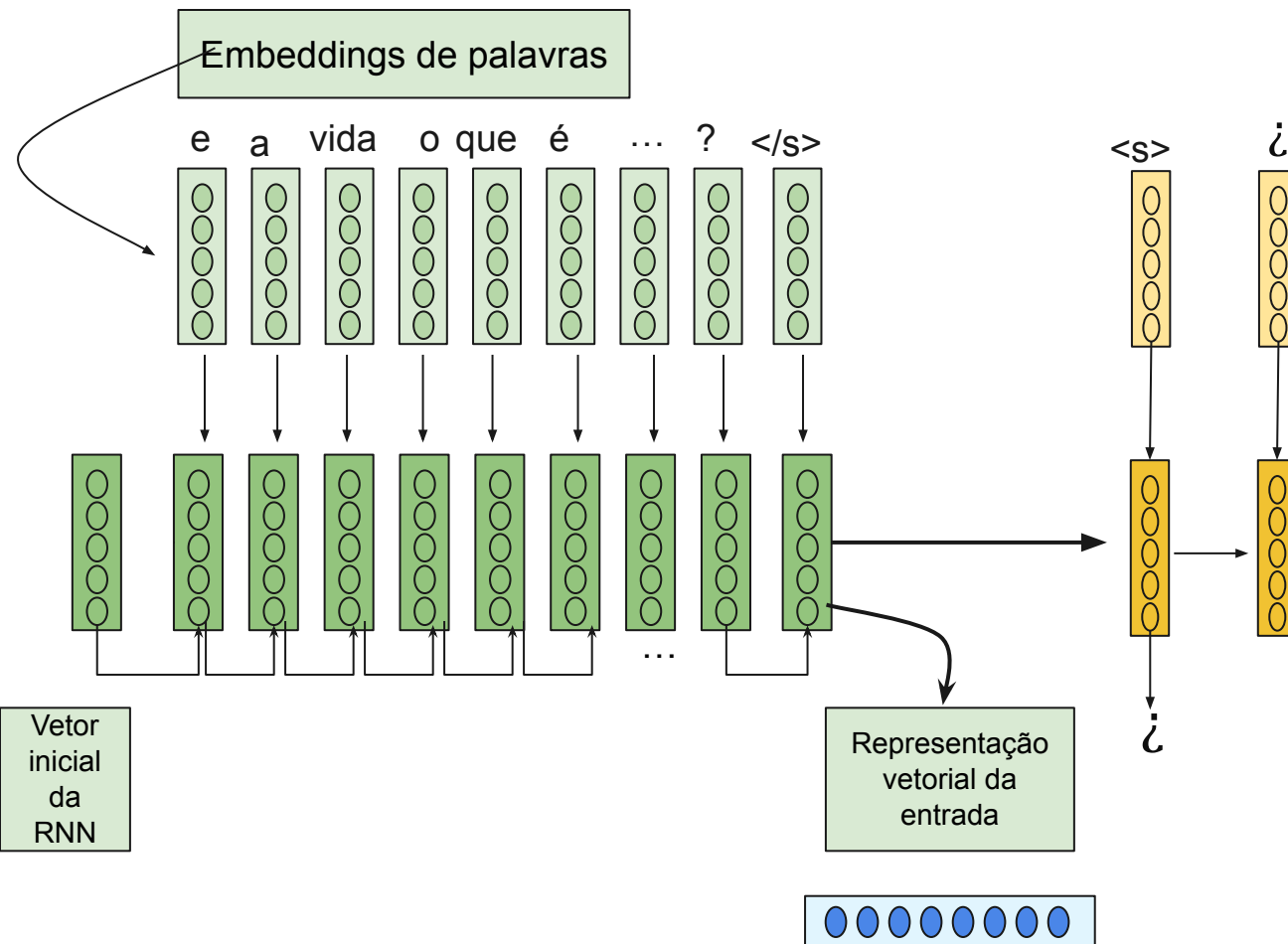
# Encoder-Decoder com RNN



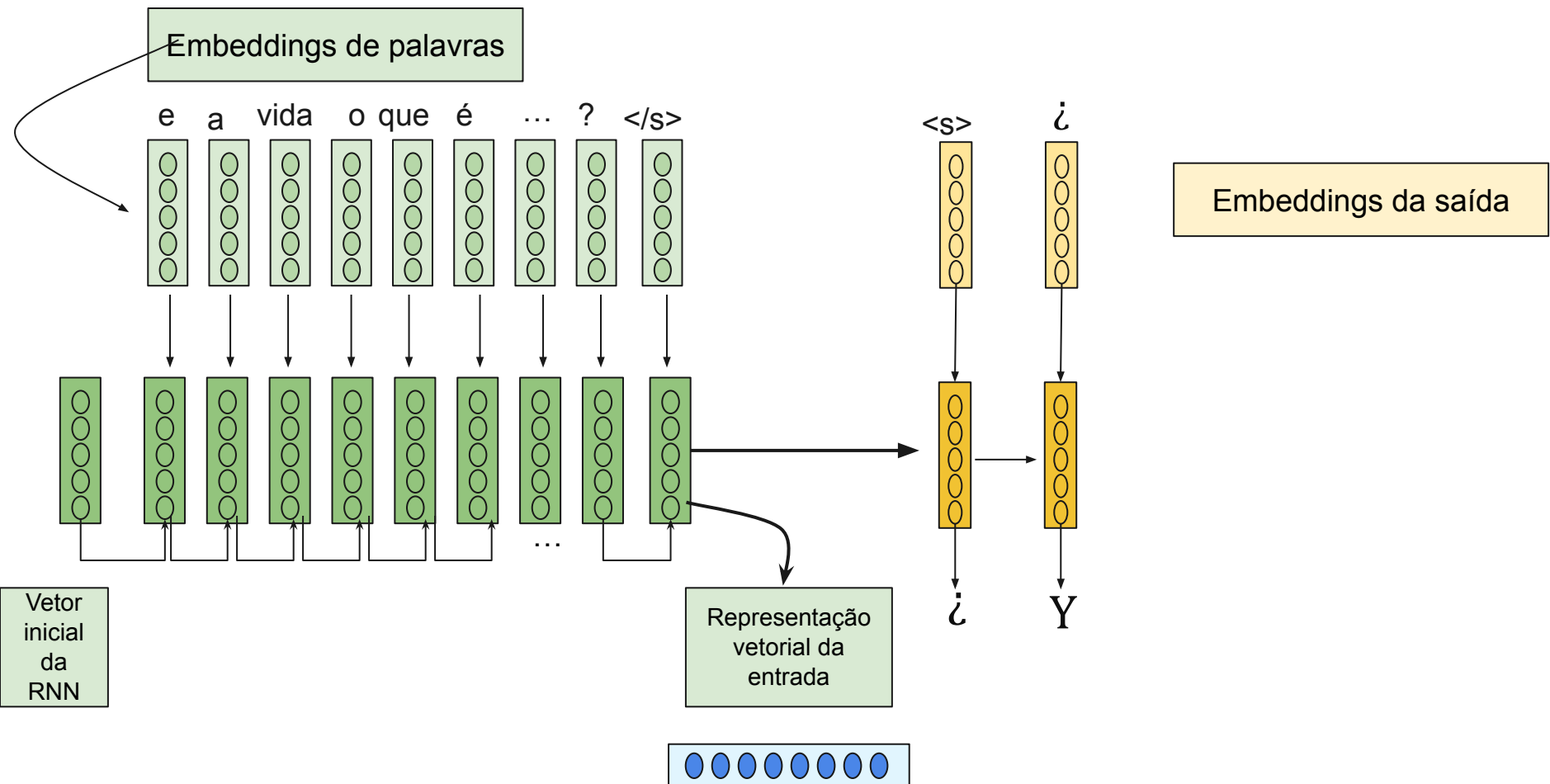
# Encoder-Decoder com RNN



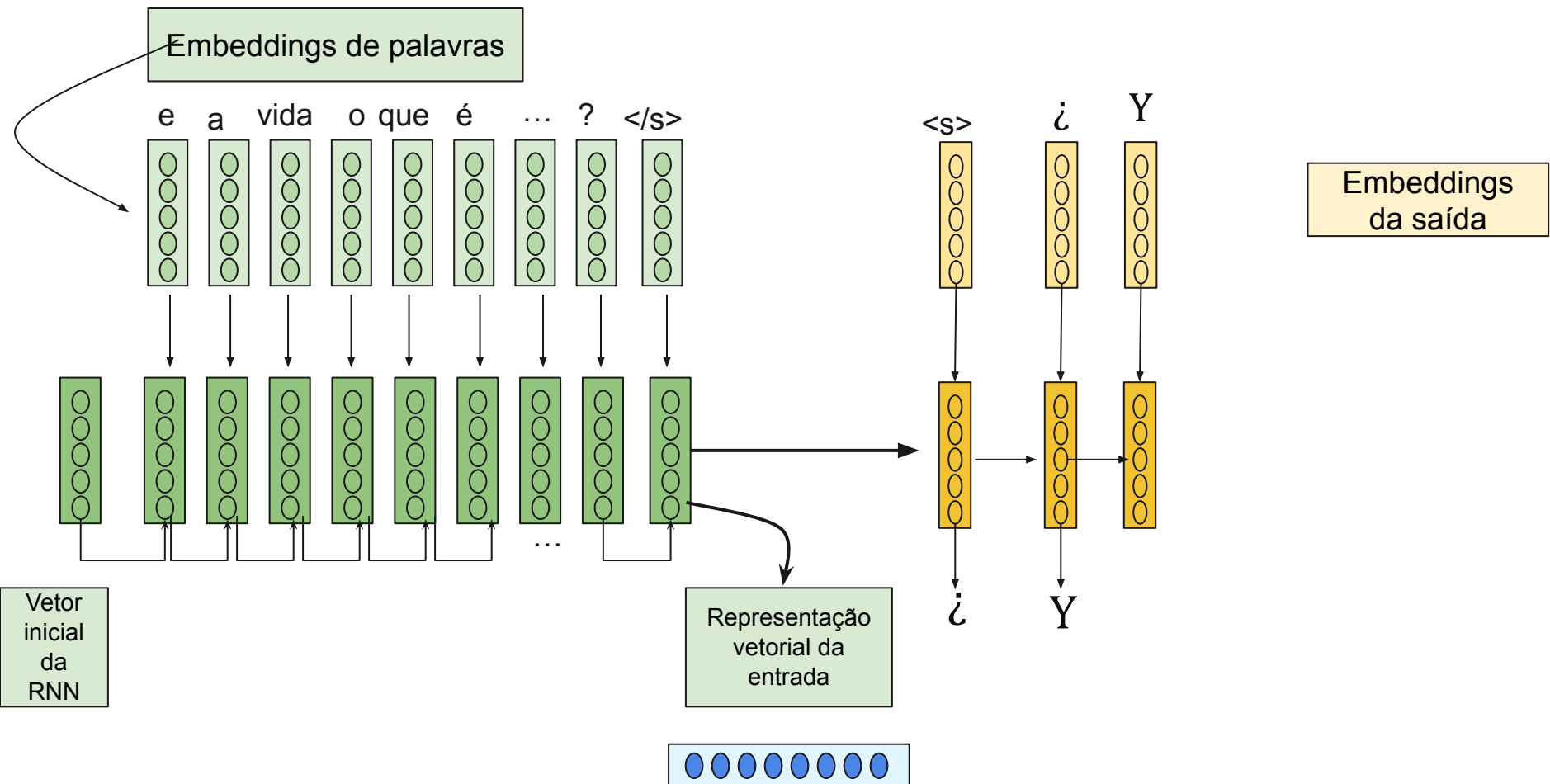
# Encoder-Decoder com RNN



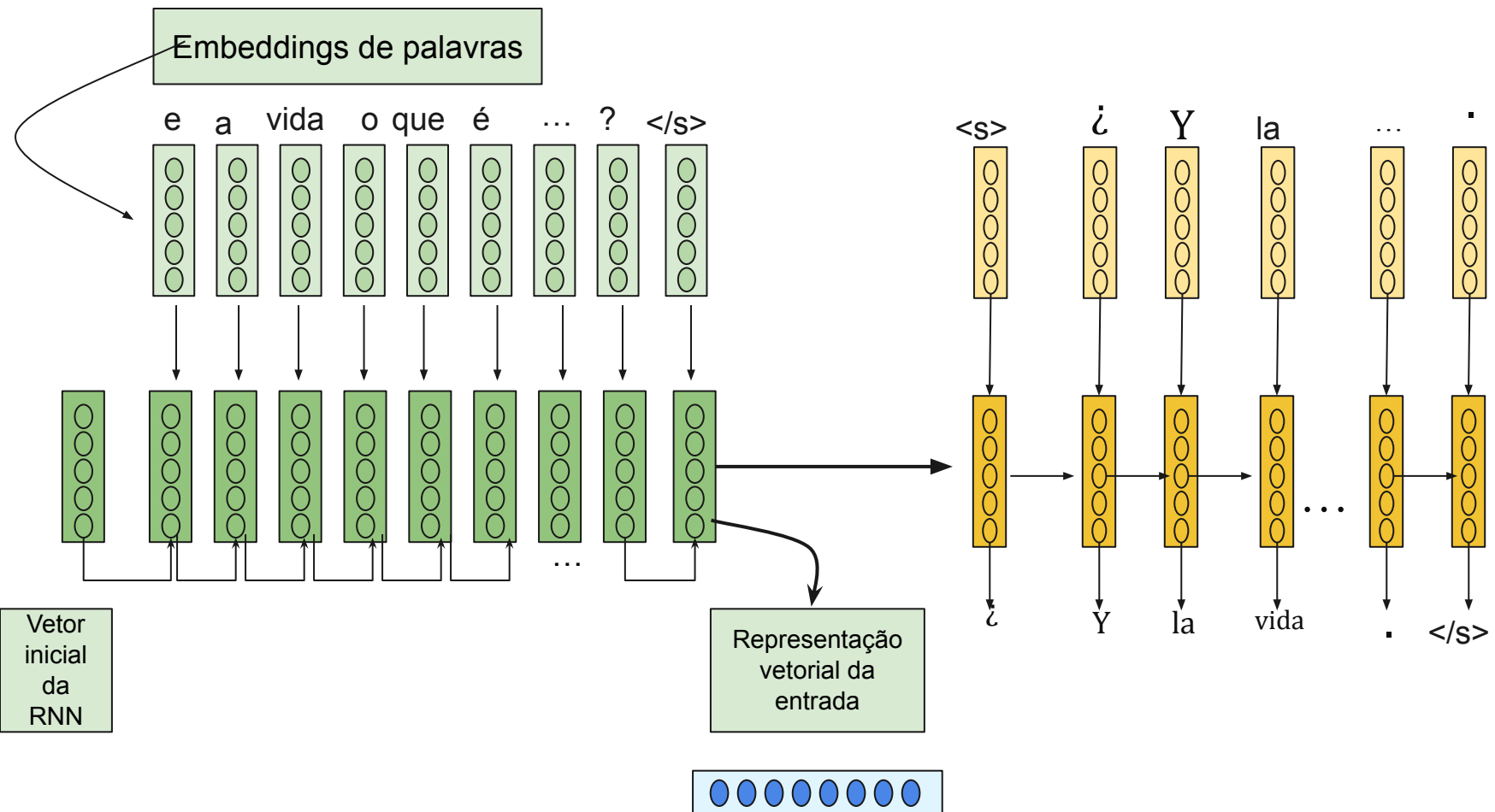
# Encoder-Decoder com RNN



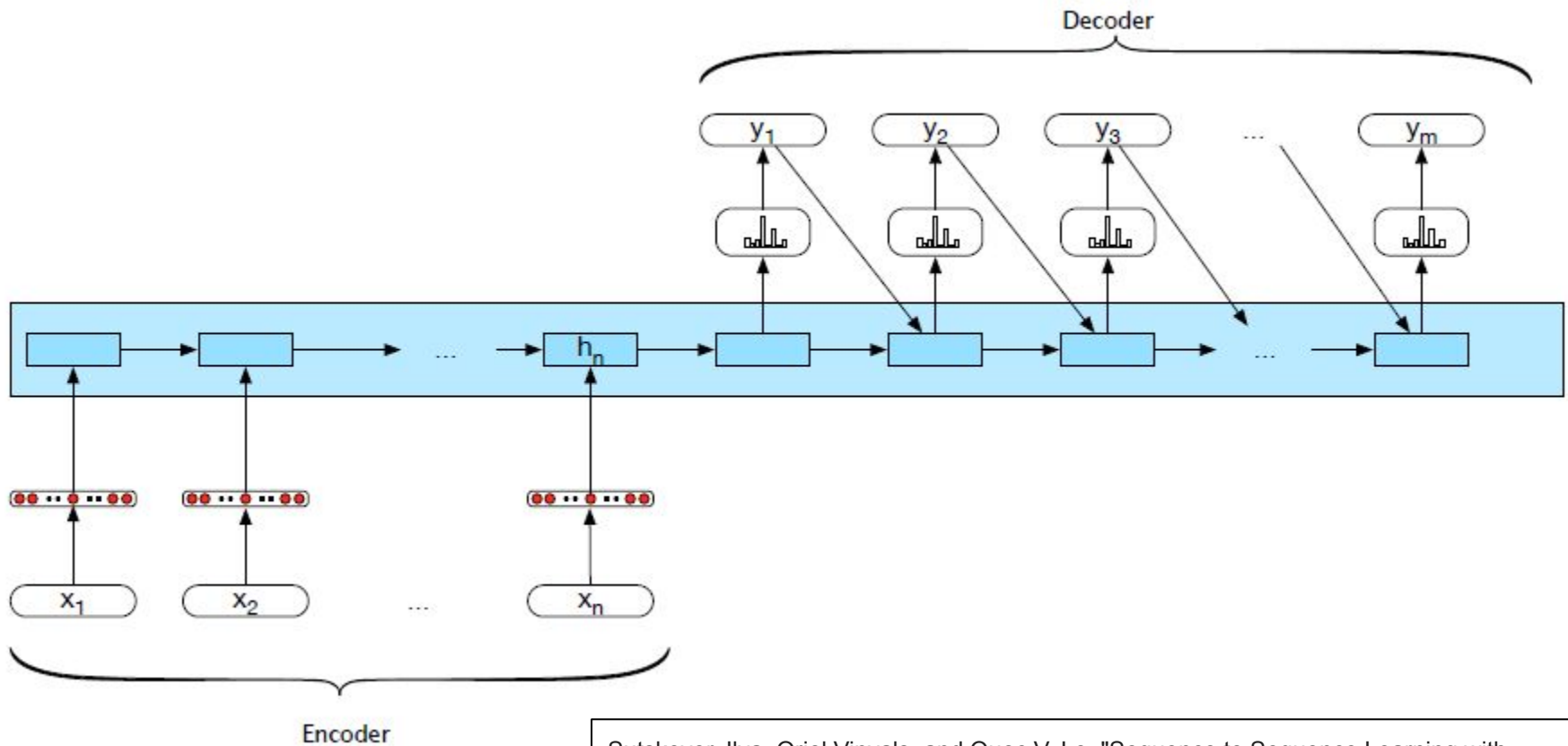
# Encoder-Decoder com RNN



# Encoder-Decoder com RNN



# Encoder-decoder básico



Speech & language Processing

Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. "Sequence to Sequence Learning with Neural Networks." *Advances in Neural Information Processing Systems* 27 (2014): 3104-3112.

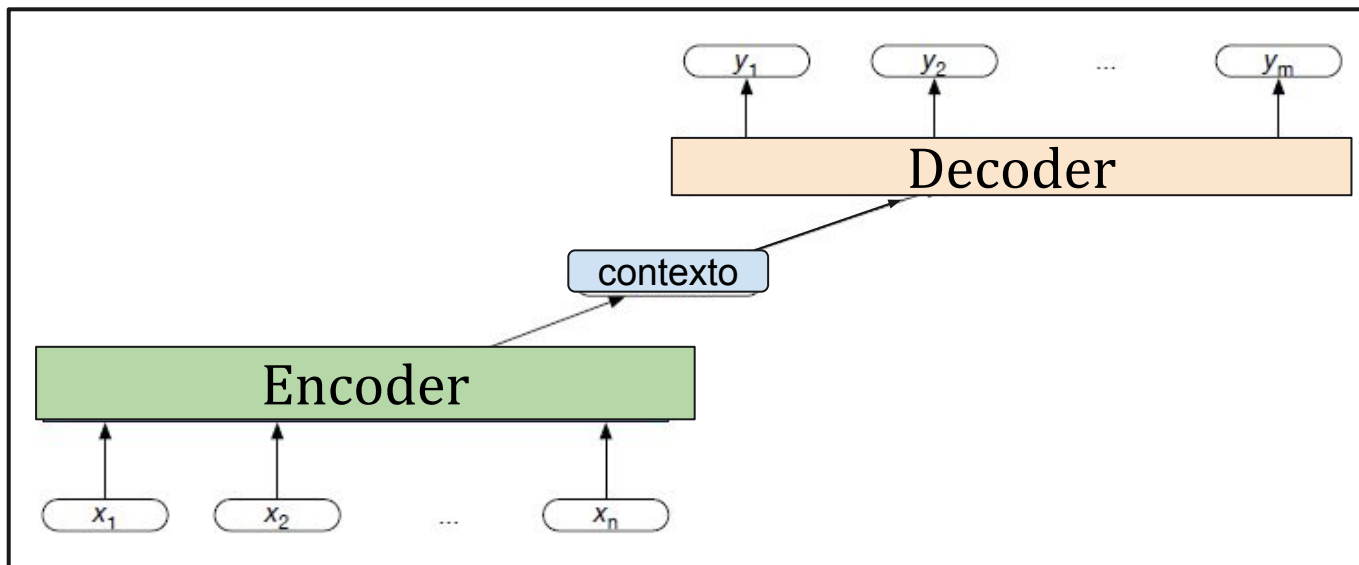
# Pergunta

Como aprender os pesos do encoder?

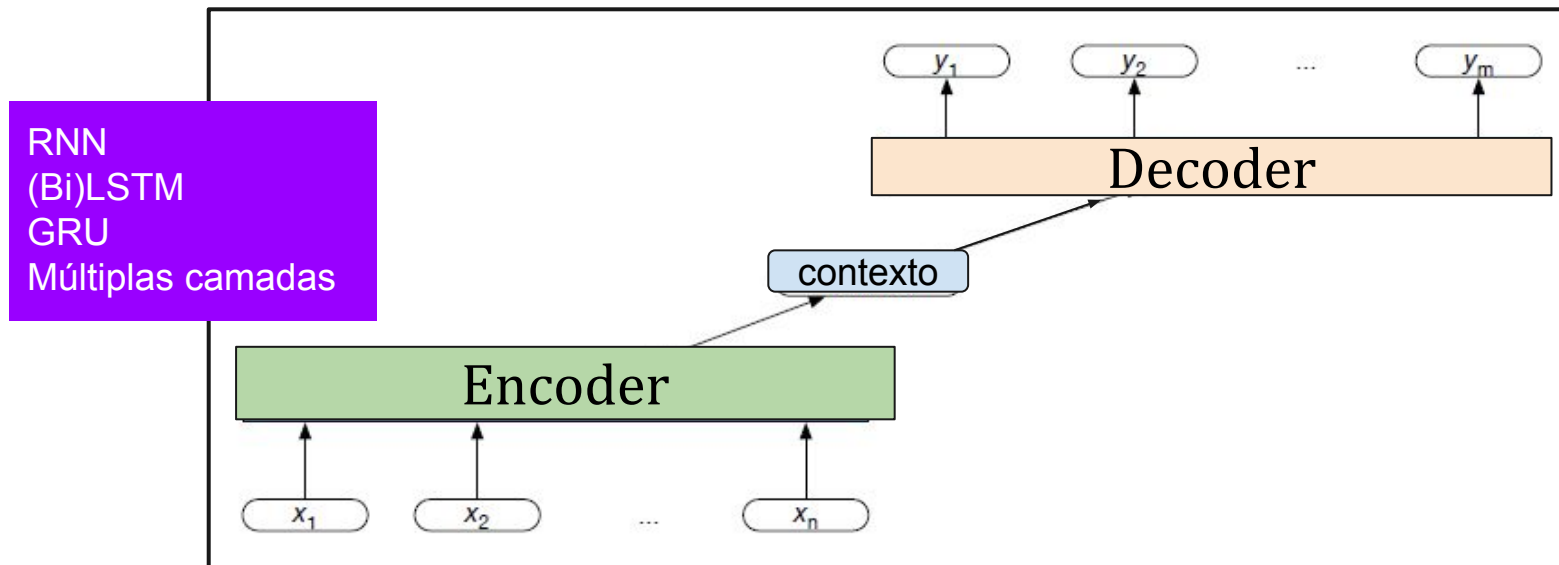
- a) Usando uma função de custo para o contexto
- b) A partir da função de custo calculada no fim do decoder
- c) Usando uma função de custo para cada componente do encoder
- d) Descobrendo numericamente o máximo da função



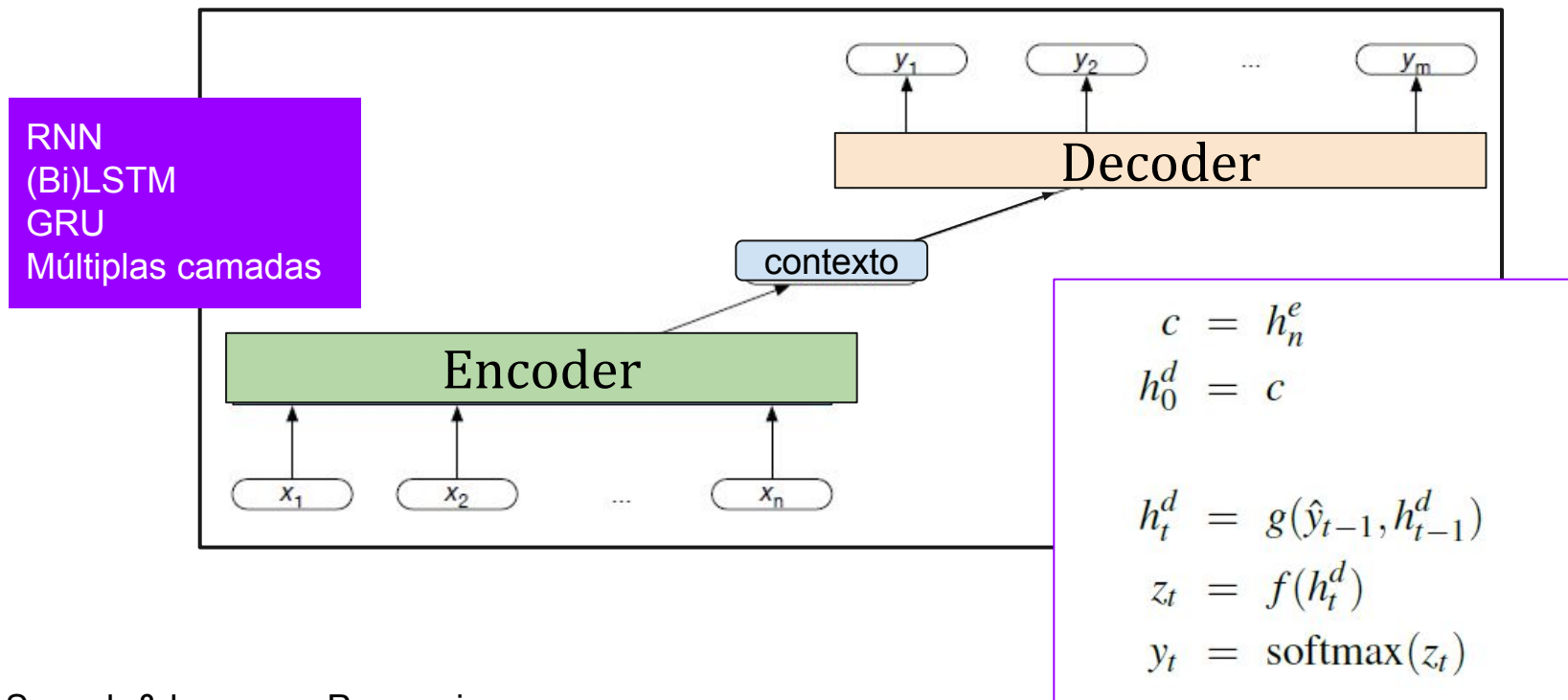
# Encoder-decoder abstração



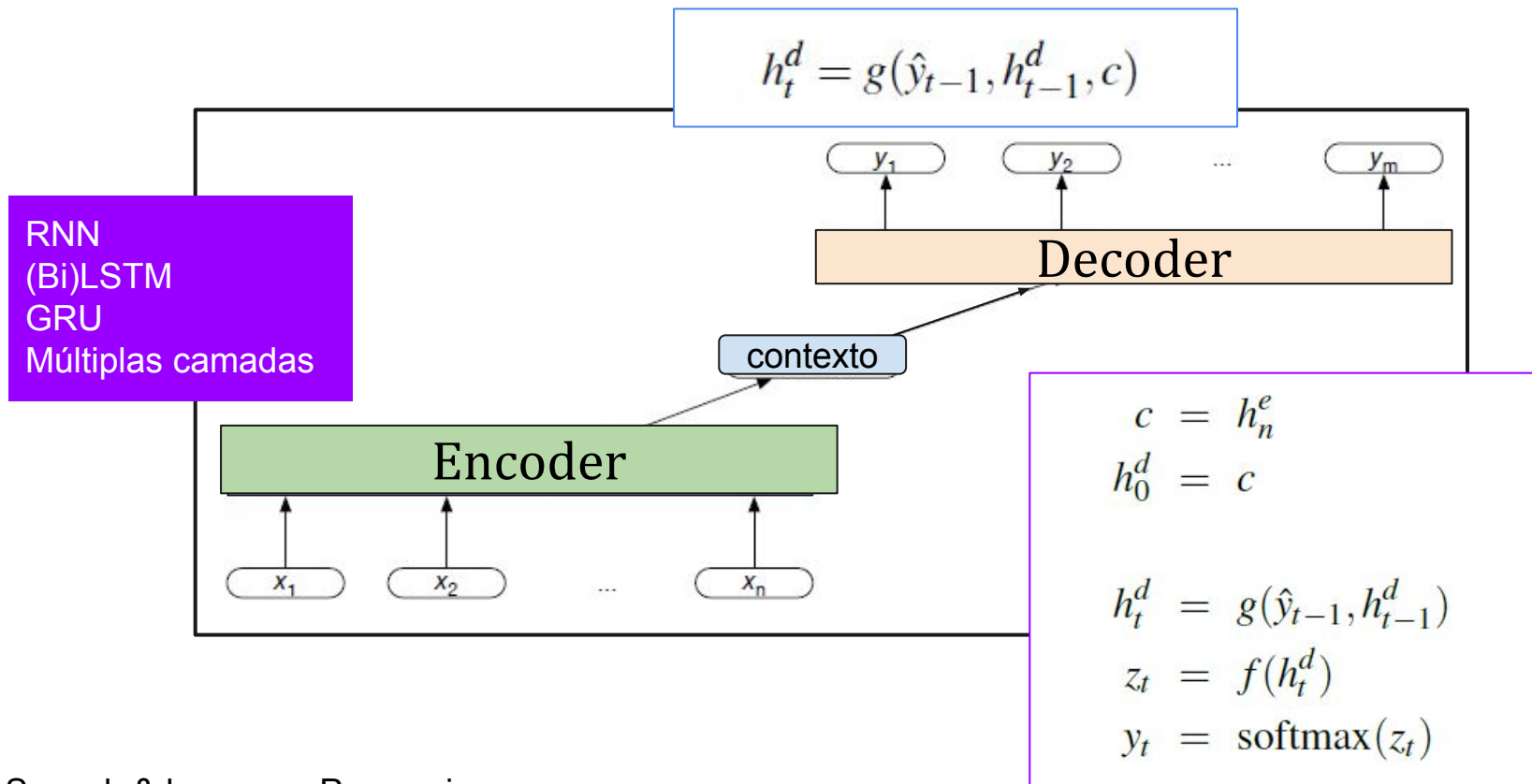
# Encoder-decoder abstração



# Encoder-decoder abstração



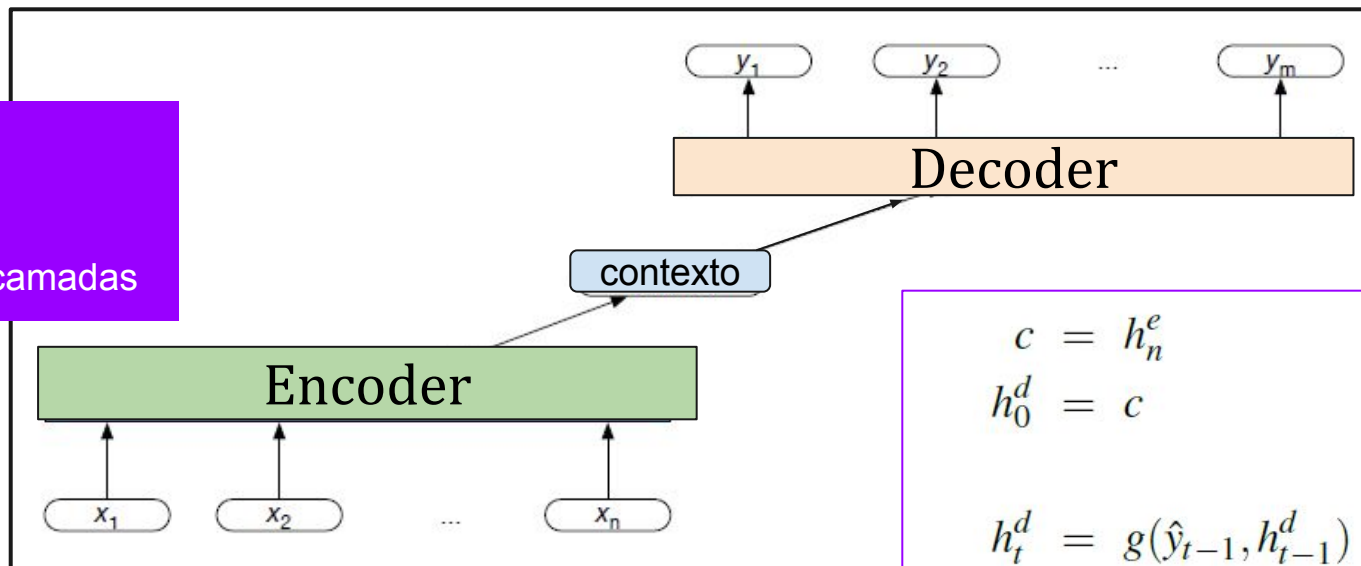
# Encoder-decoder abstração



# Encoder-decoder abstração

Decoder olha apenas o estado escondido e a última saída

RNN  
(Bi)LSTM  
GRU  
Múltiplas camadas



$$c = h_n^e$$

$$h_0^d = c$$

$$h_t^d = g(\hat{y}_{t-1}, h_{t-1}^d)$$

$$z_t = f(h_t^d)$$

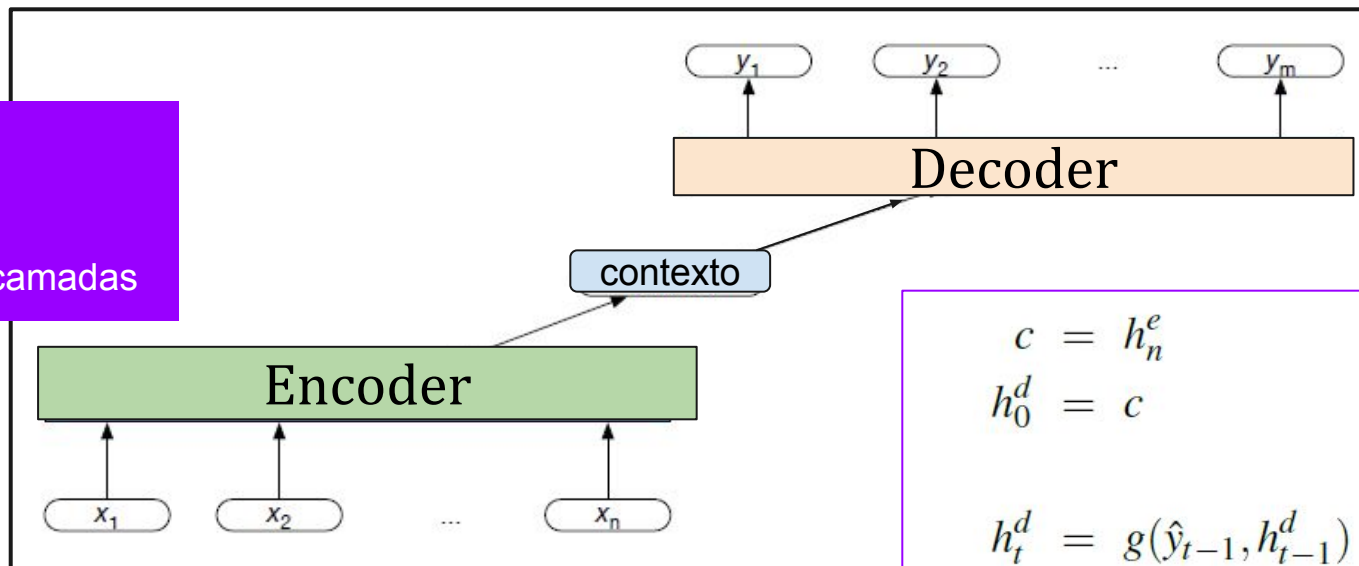
$$y_t = \text{softmax}(z_t)$$

# Encoder-decoder abstração

Decoder olha apenas o estado escondido e a última saída

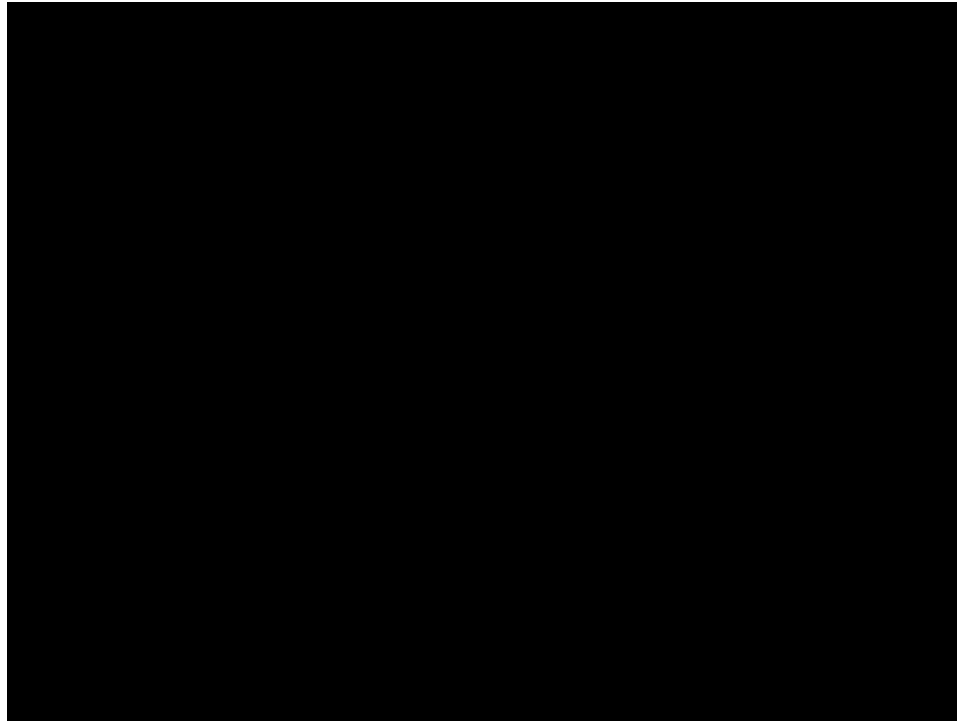
$$y_t = \text{softmax}(\hat{y}_{t-1}, z_t, c)$$

RNN  
(Bi)LSTM  
GRU  
Múltiplas camadas



$$\begin{aligned} c &= h_n^e \\ h_0^d &= c \\ h_t^d &= g(\hat{y}_{t-1}, h_{t-1}^d) \\ z_t &= f(h_t^d) \\ y_t &= \text{softmax}(z_t) \end{aligned}$$

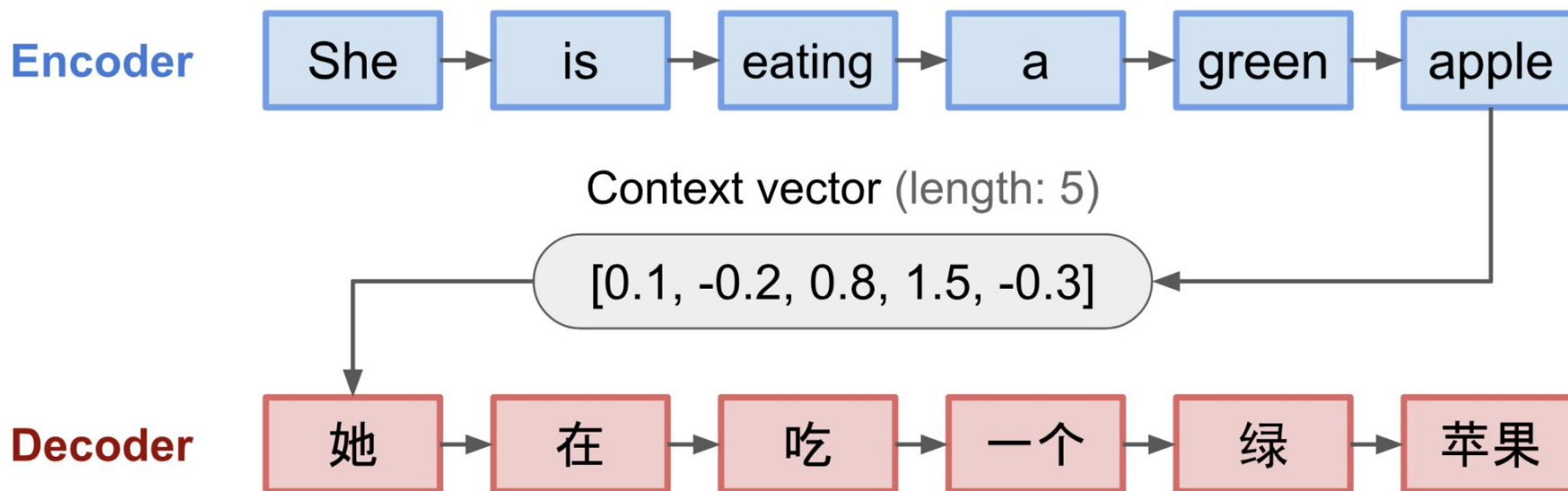
# Bottleneck



<http://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>

# Bottleneck

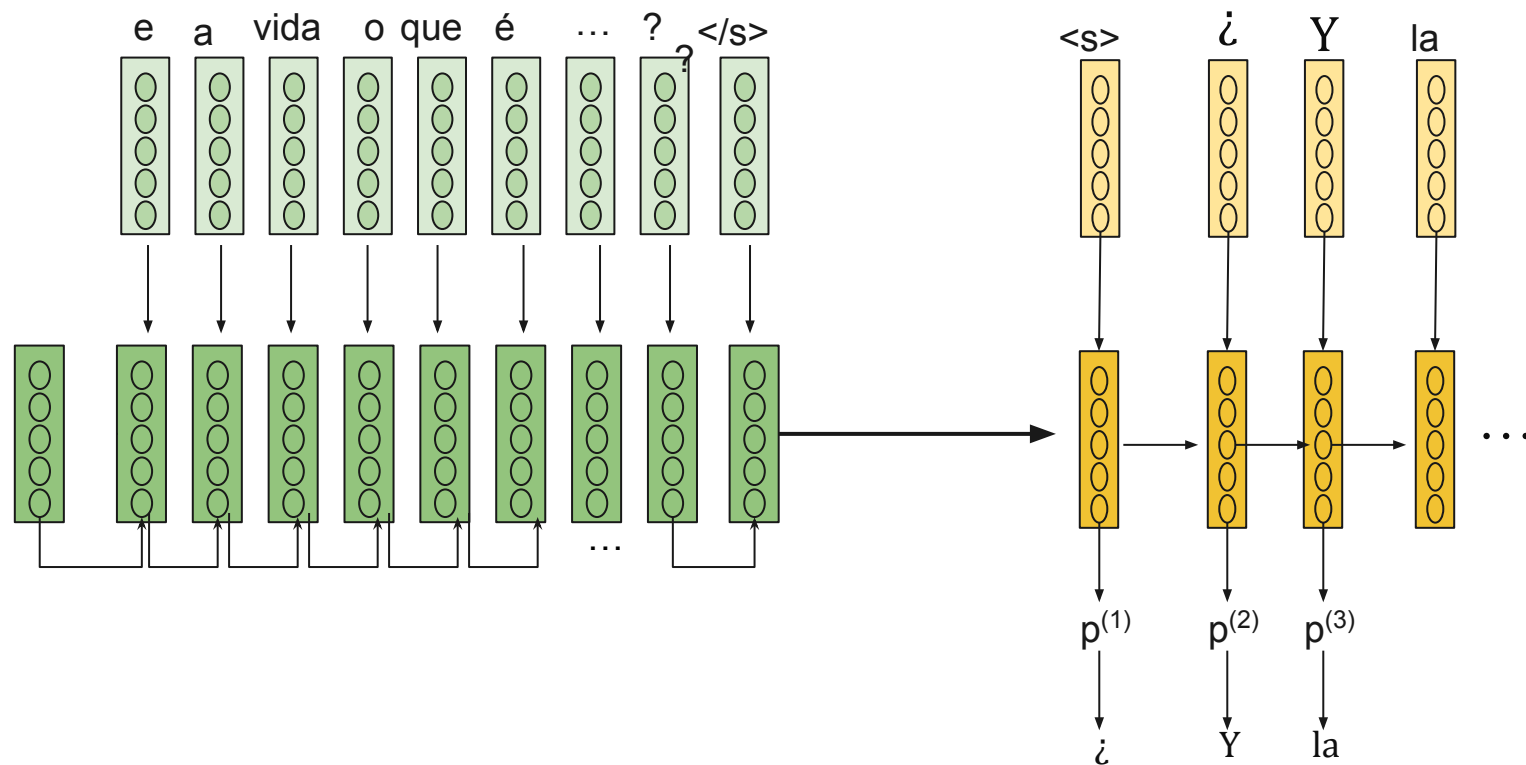
*Processa a sequência de entrada e compacta as informações em um vetor de contexto de comprimento fixo. Espera-se que esta representação seja um bom resumo do significado de toda a sequência fonte.*



*Um decodificador é inicializado com o vetor de contexto para emitir a saída transformada. Os primeiros trabalhos usavam apenas o último estado da rede do codificador como estado inicial do decodificador : como lembrar de tudo aqui??.*

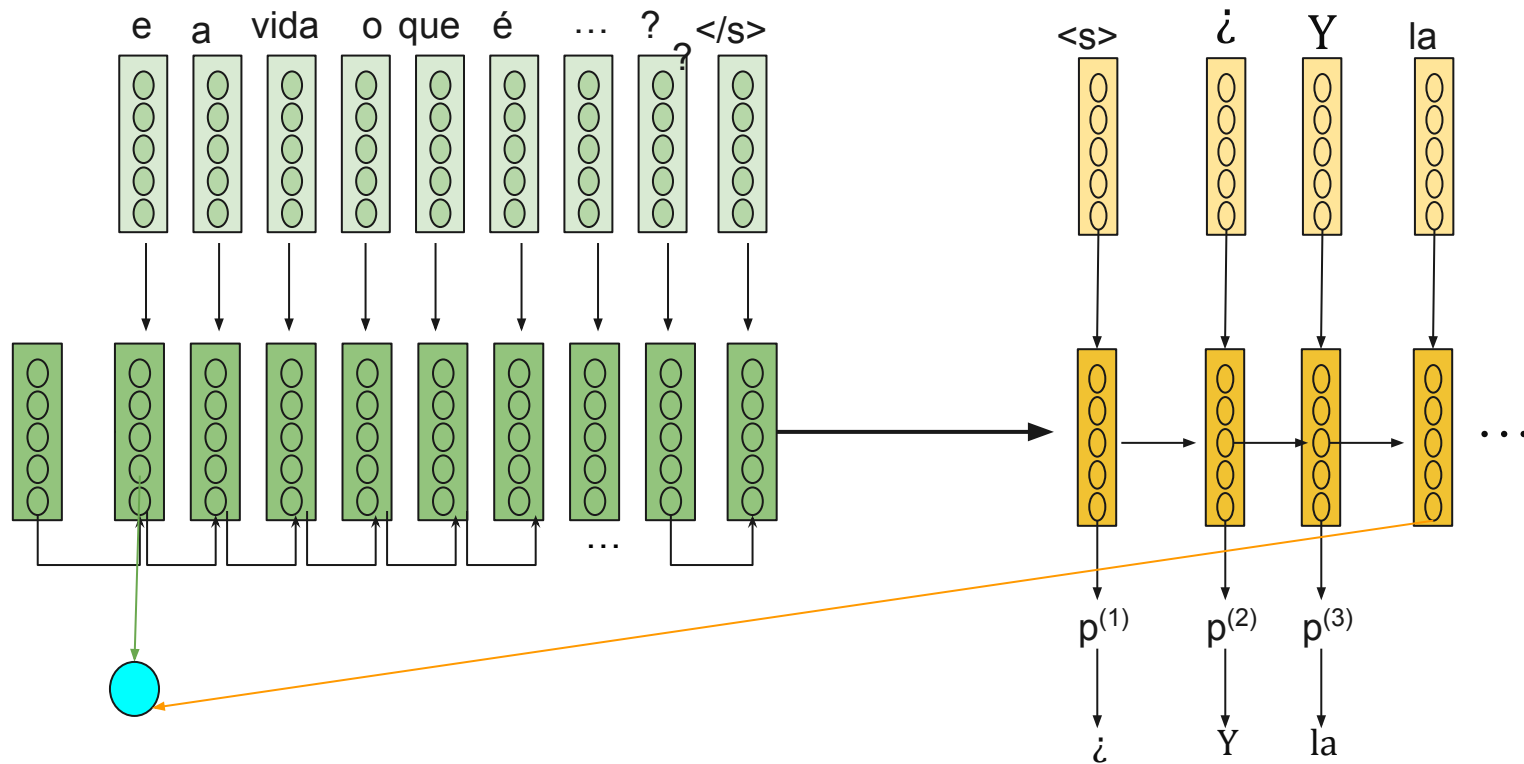


# ATENÇÃO



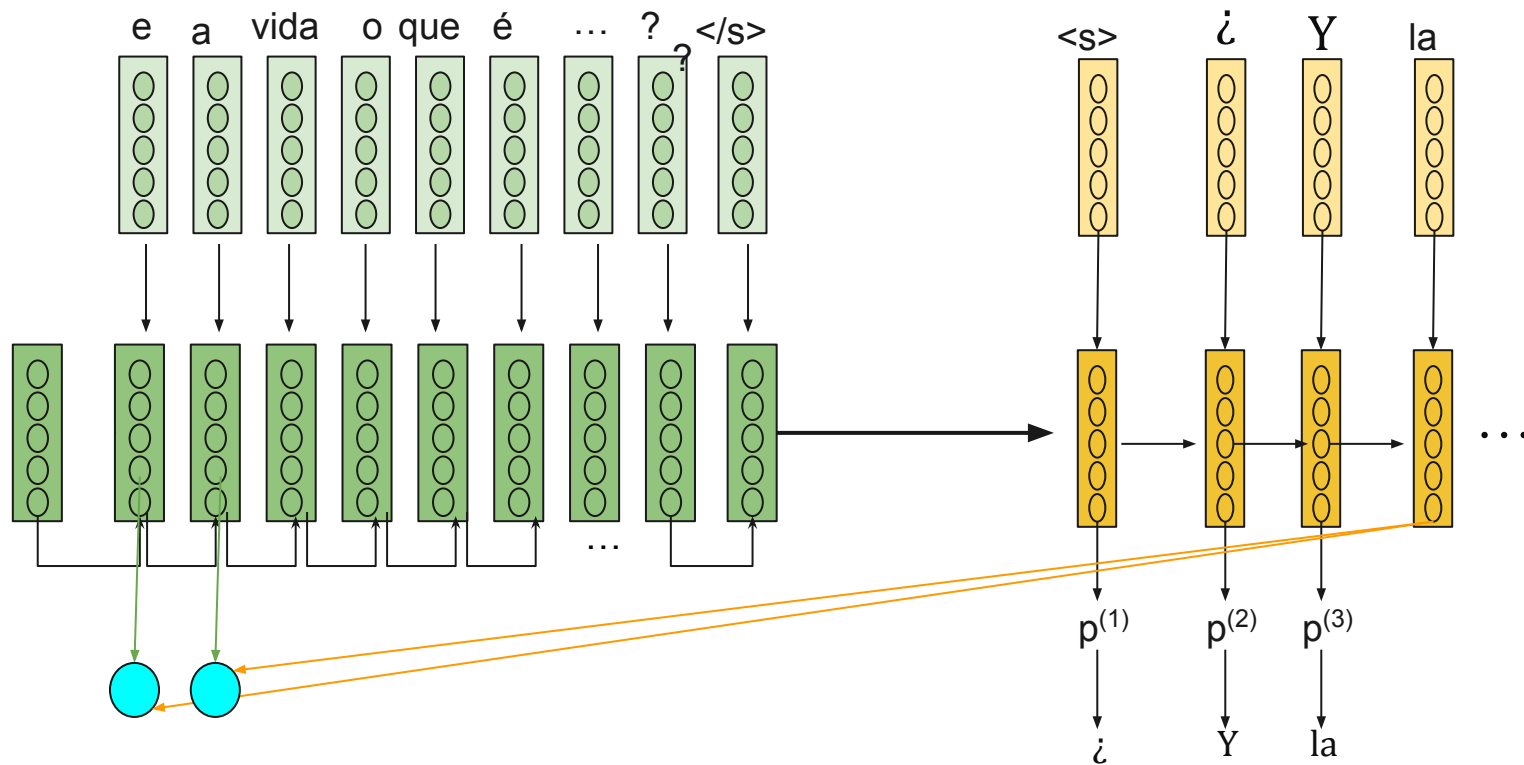
FAÇA O MODELO FOCAR EM DIFERENTES PARTES DA ENTRADA A CADA PASSO

# ATENÇÃO



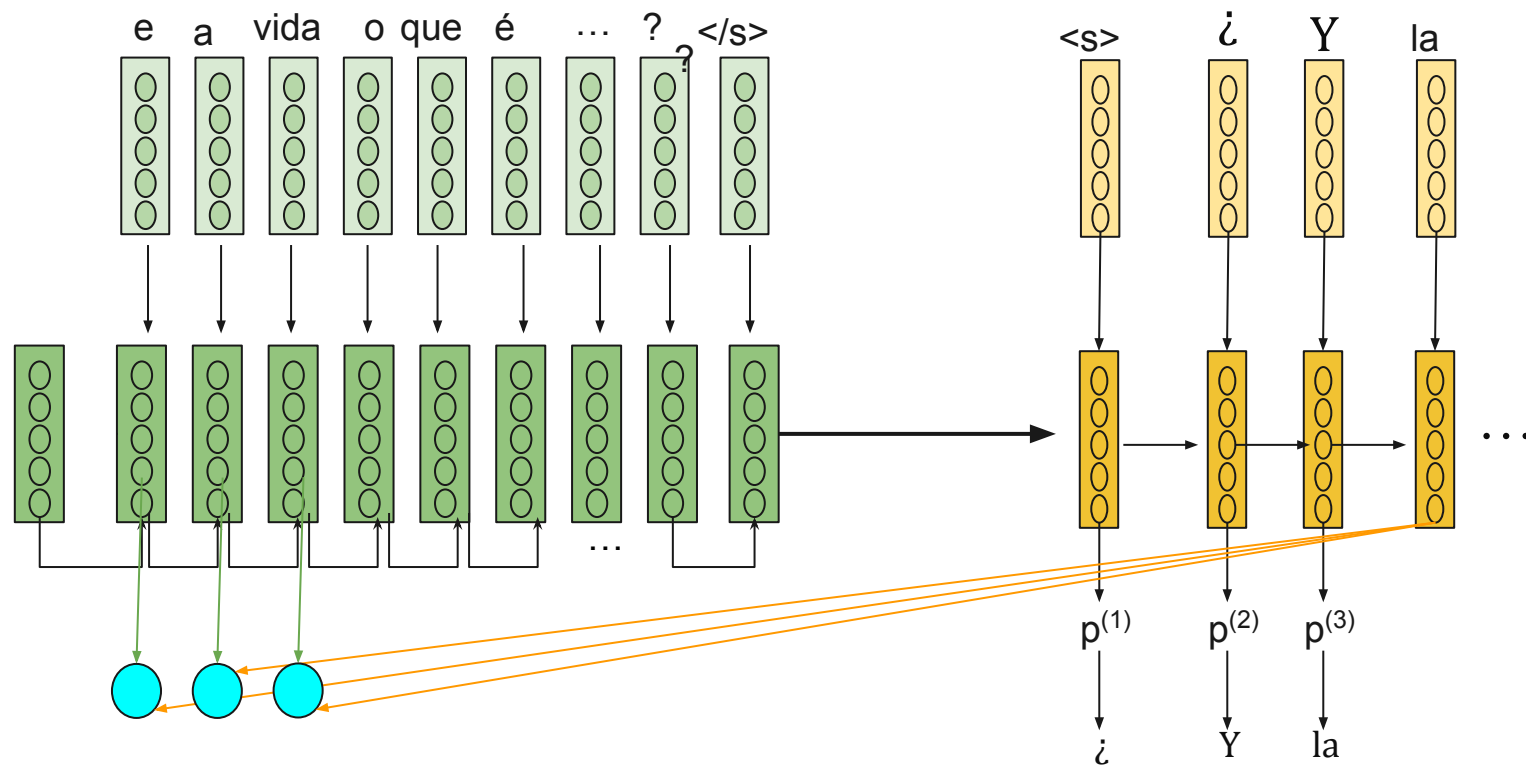
FAÇA O MODELO FOCAR EM DIFERENTES PARTES DA ENTRADA A CADA PASSO

# ATENÇÃO



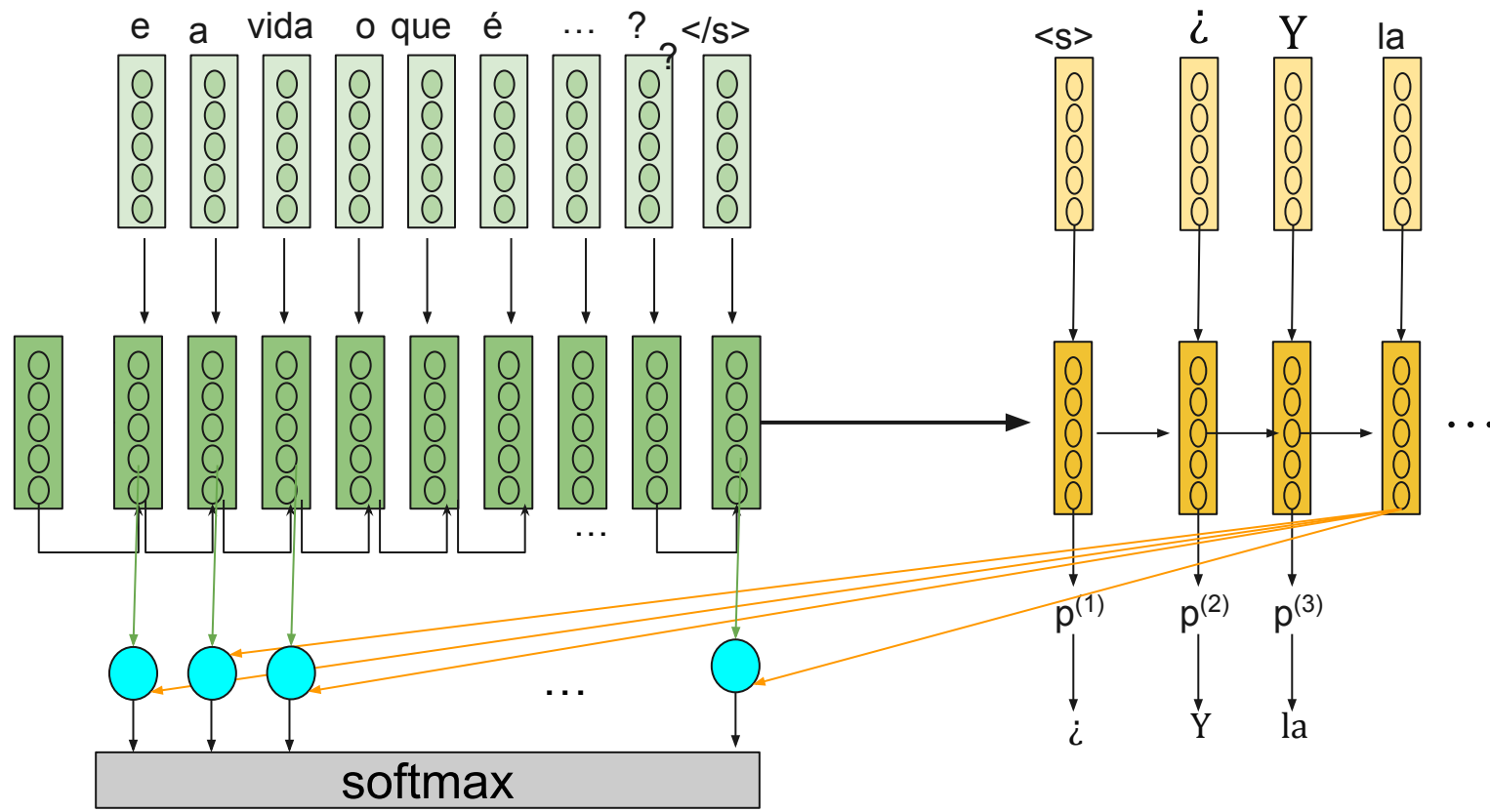
FAÇA O MODELO FOCAR EM DIFERENTES PARTES DA ENTRADA A CADA PASSO

# ATENÇÃO

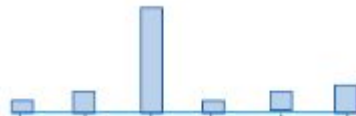


FAÇA O MODELO FOCAR EM DIFERENTES PARTES DA ENTRADA A CADA PASSO

# ATENÇÃO

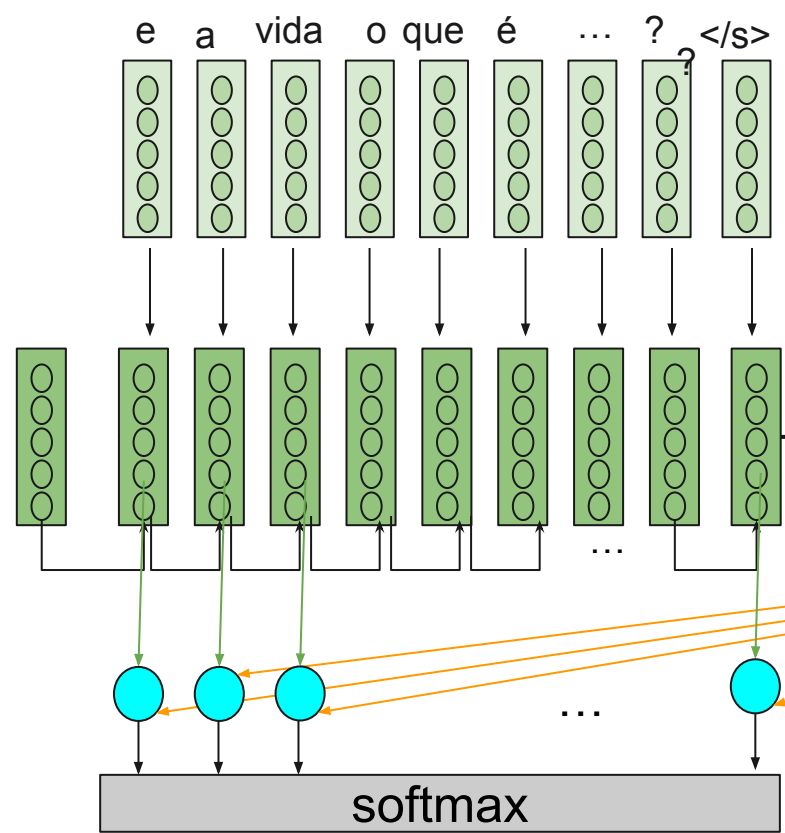


Pesos de  
atenção



Saída da atenção: soma ponderada dos estados do codificador com pesos de atenção.

# ATENÇÃO



Pesos de atenção

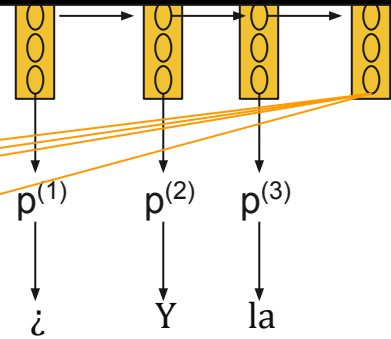


O quão relevante é o token  $k$  na entrada para o passo  $t$  na saída?

$\text{score}(h_t, s_k)$

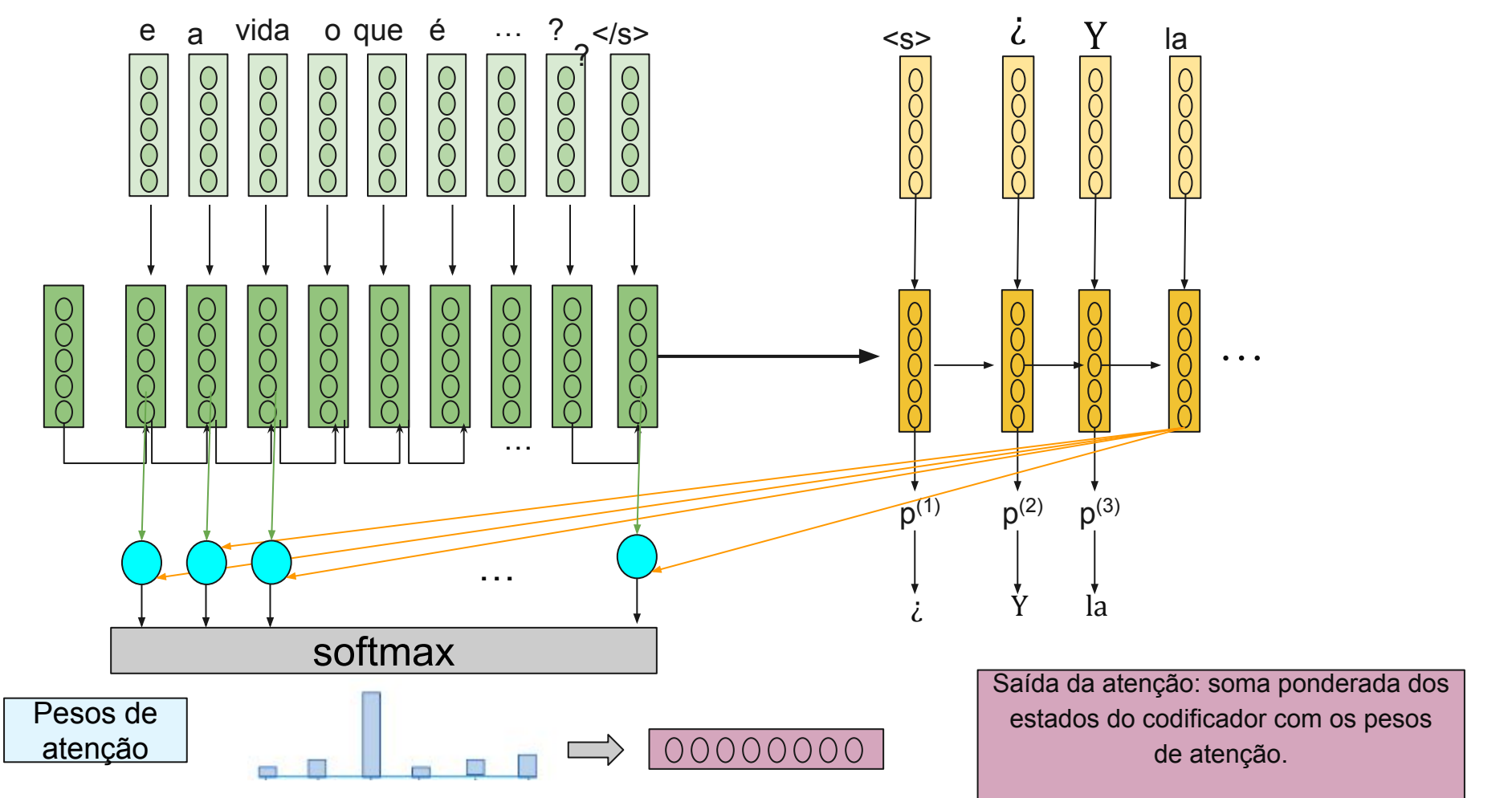
$\text{Fn attn}$

Encoder para token  $k$       Decoder no passo  $t$



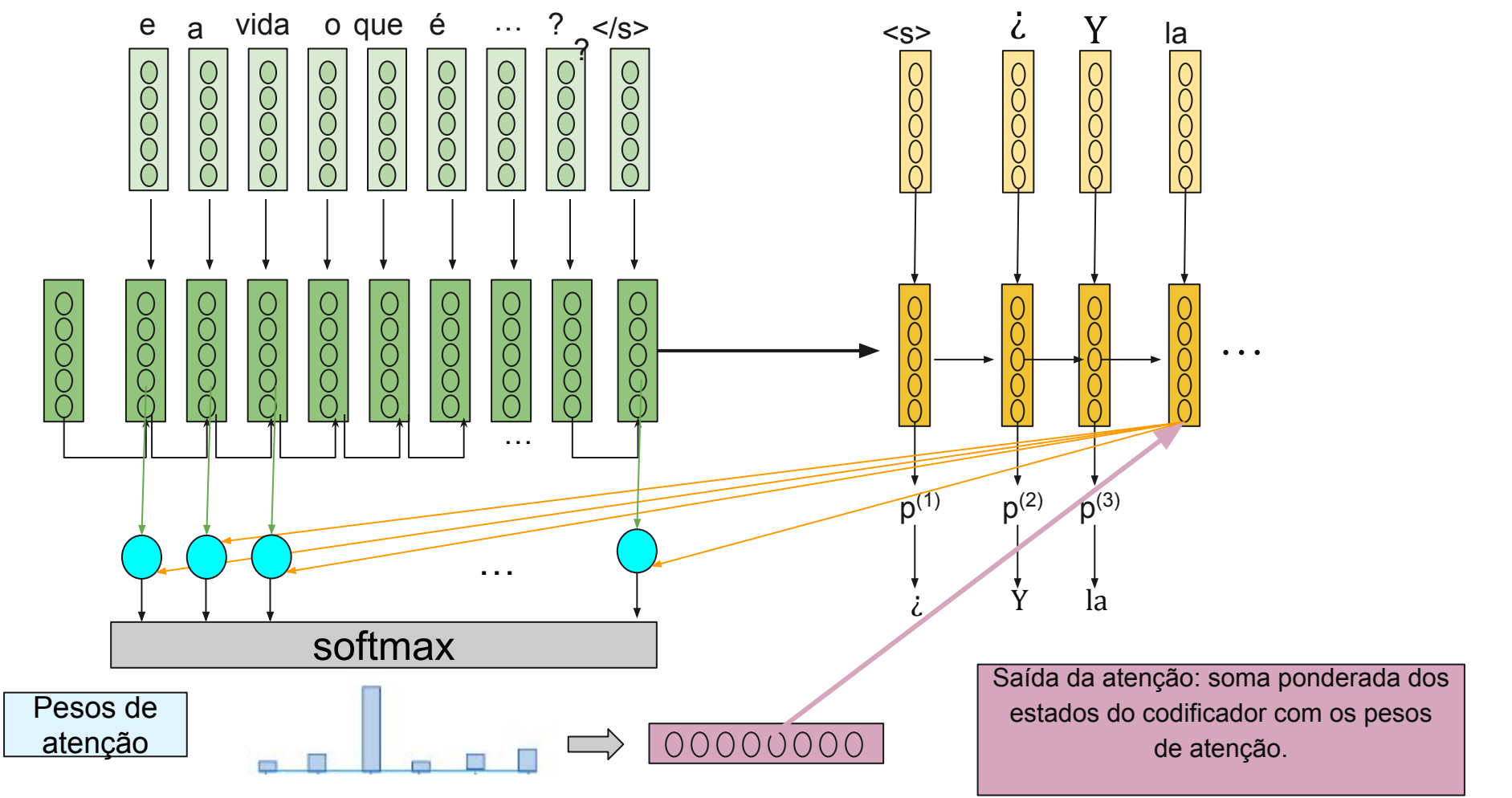
Saída da atenção: soma ponderada dos estados do codificador com pesos de atenção.

# ATENÇÃO



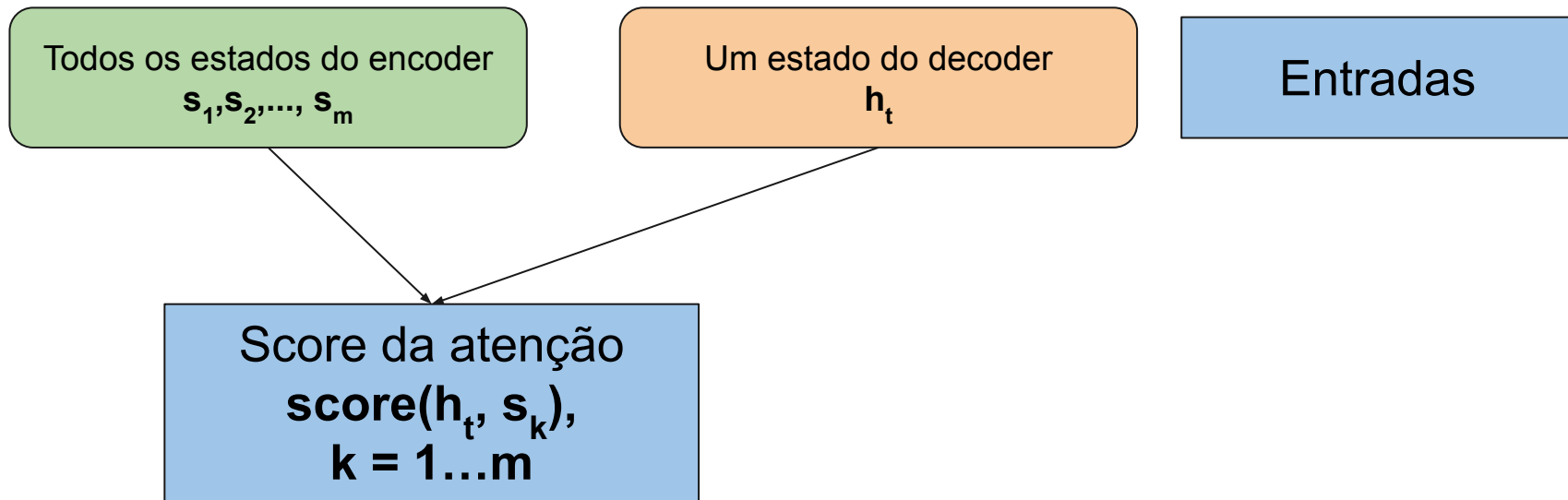
Saída da atenção: soma ponderada dos estados do codificador com pesos de atenção.

# ATENÇÃO

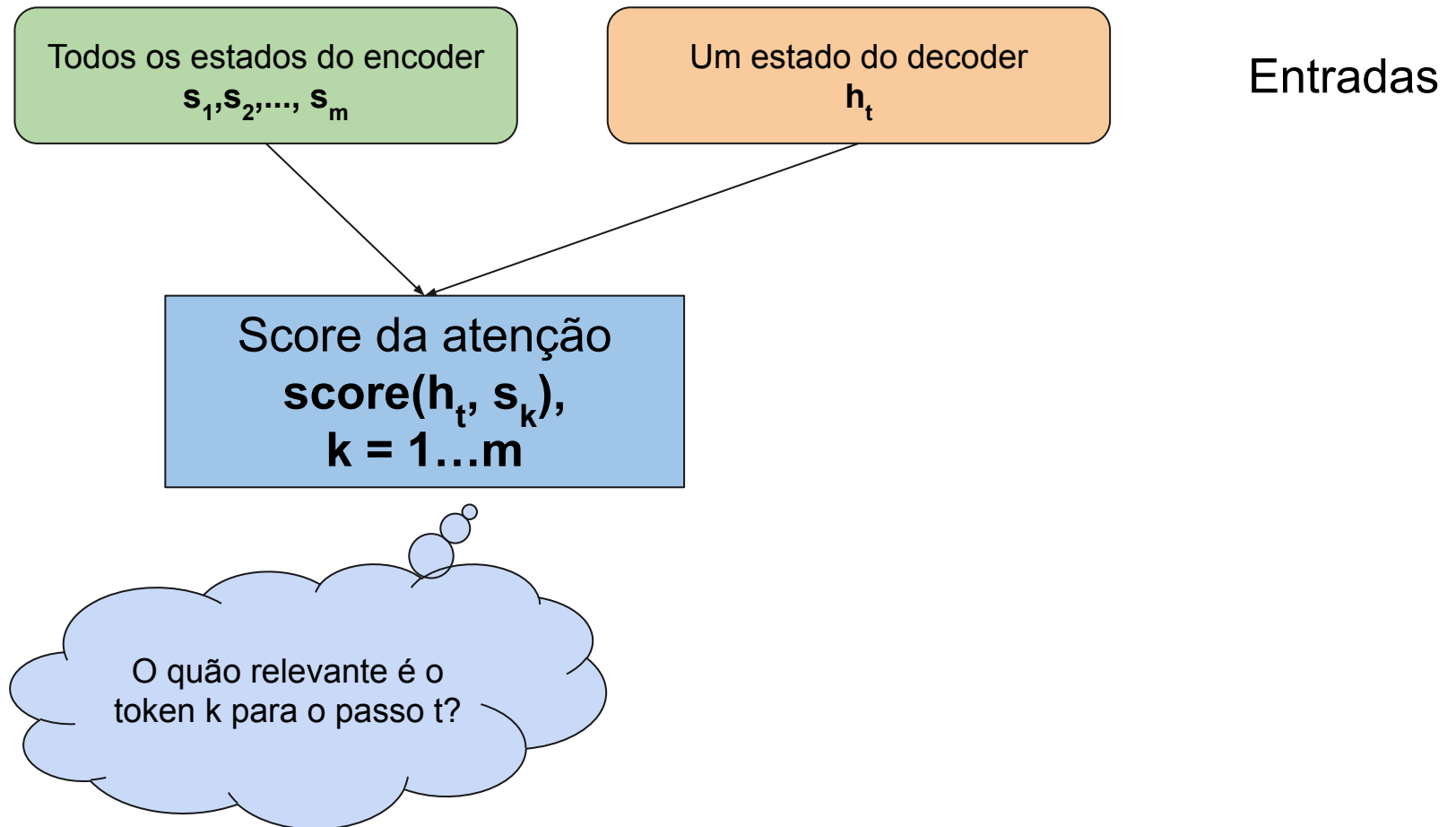




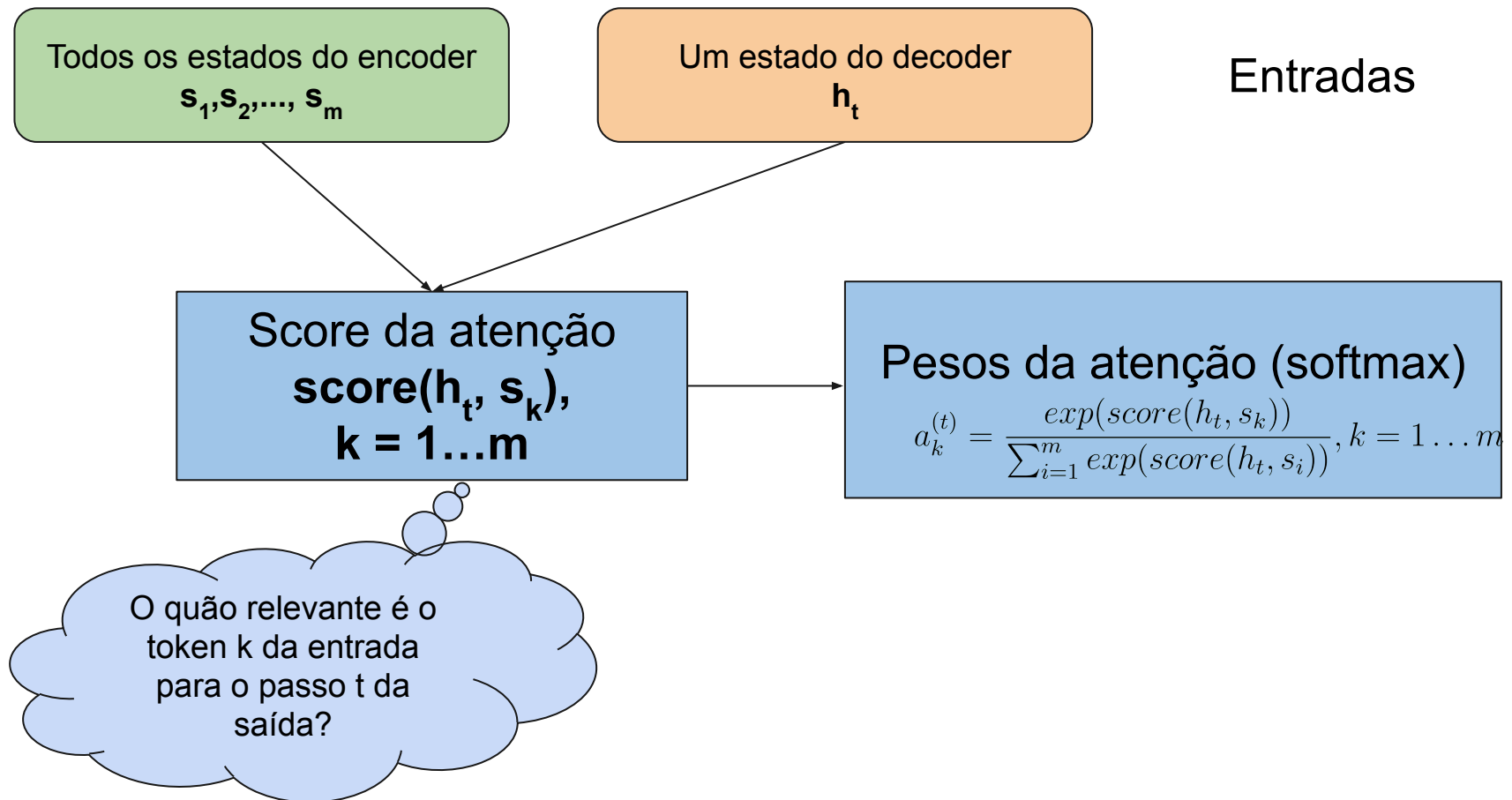
# Pipeline



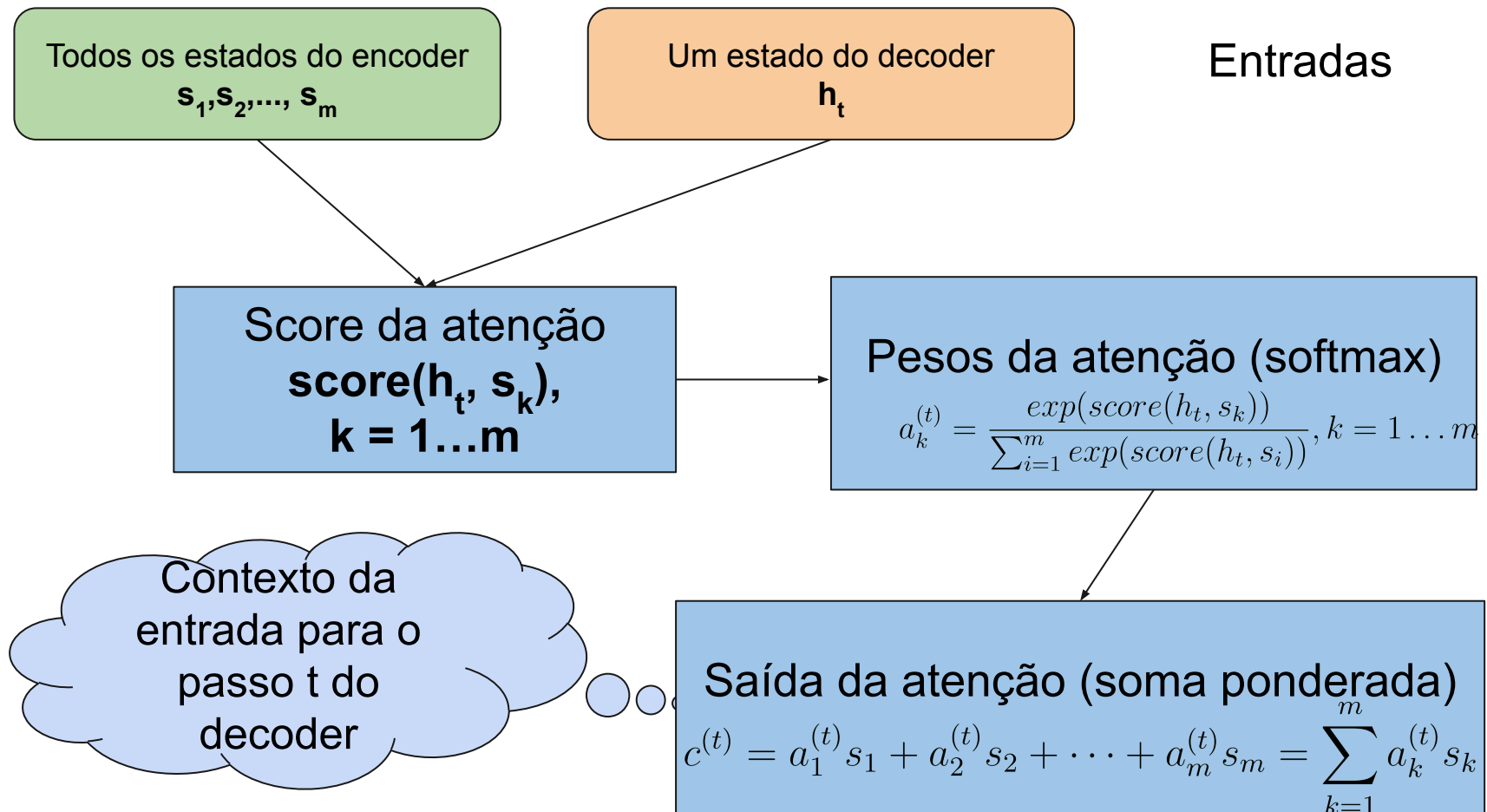
# Pipeline



# Pipeline

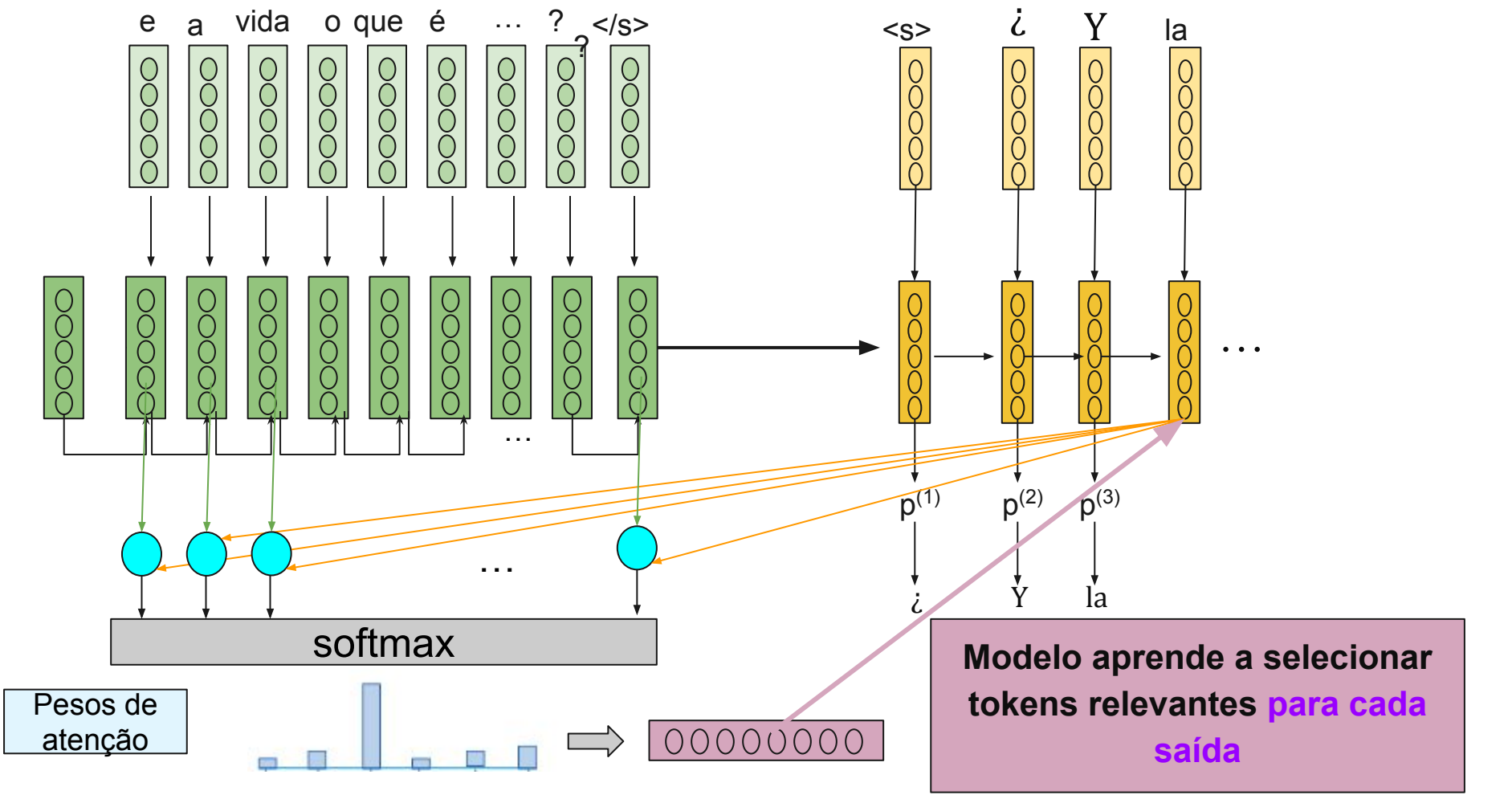


# Pipeline



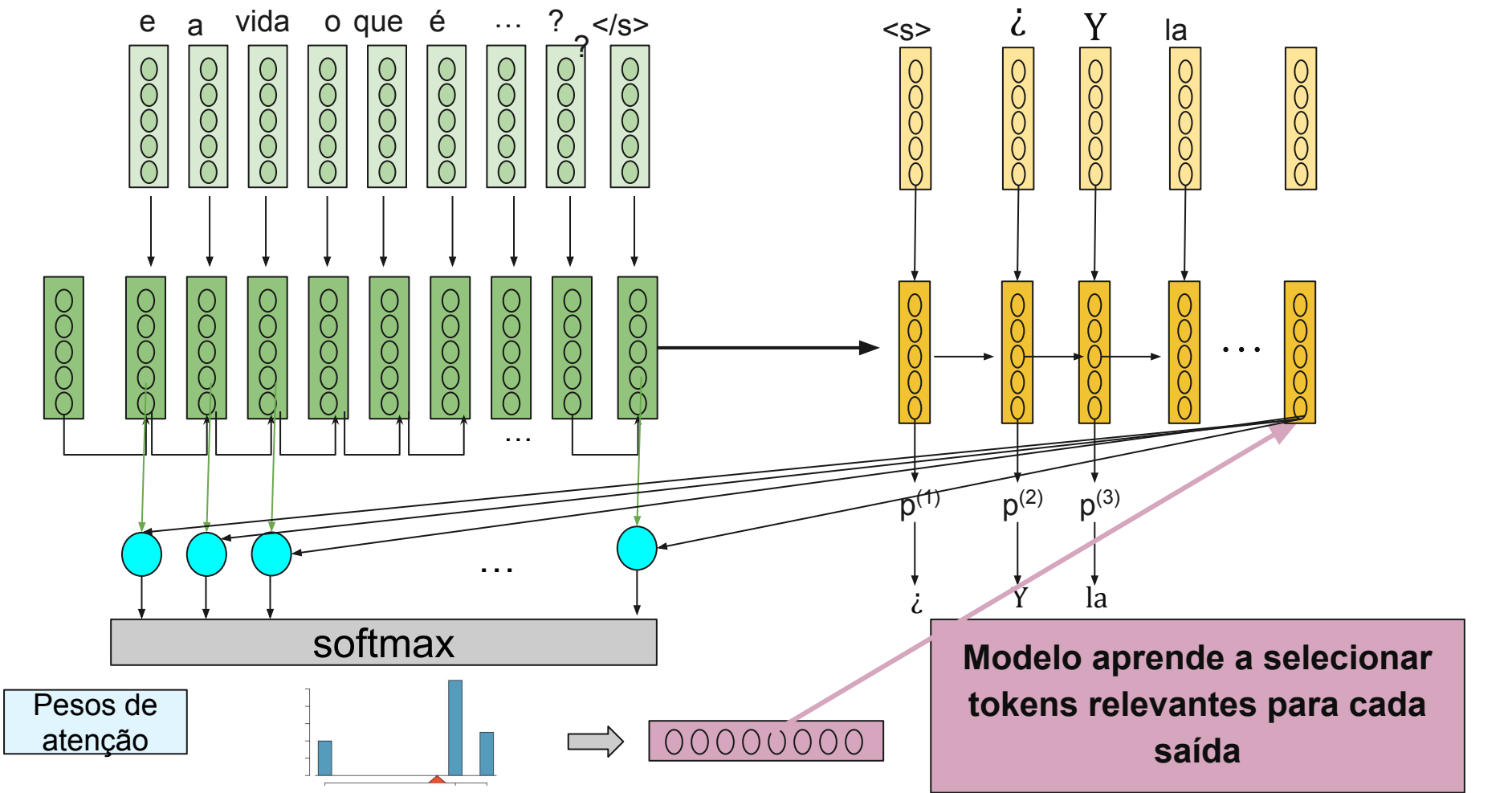
Saída da atenção: soma ponderada dos estados do codificador com pesos de atenção.

# ATENÇÃO

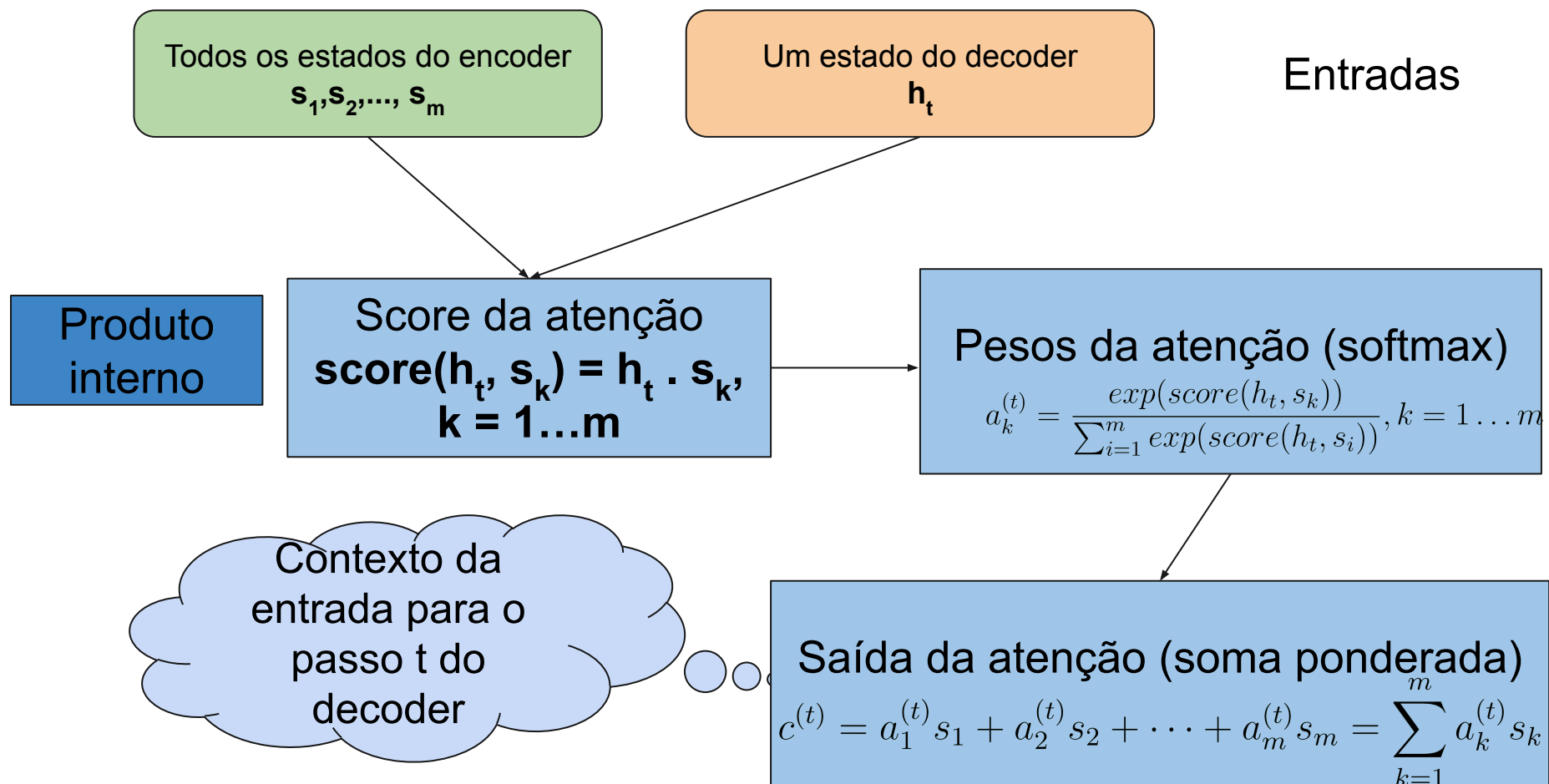


Saída da atenção: soma ponderada dos estados do codificador com pesos de atenção.

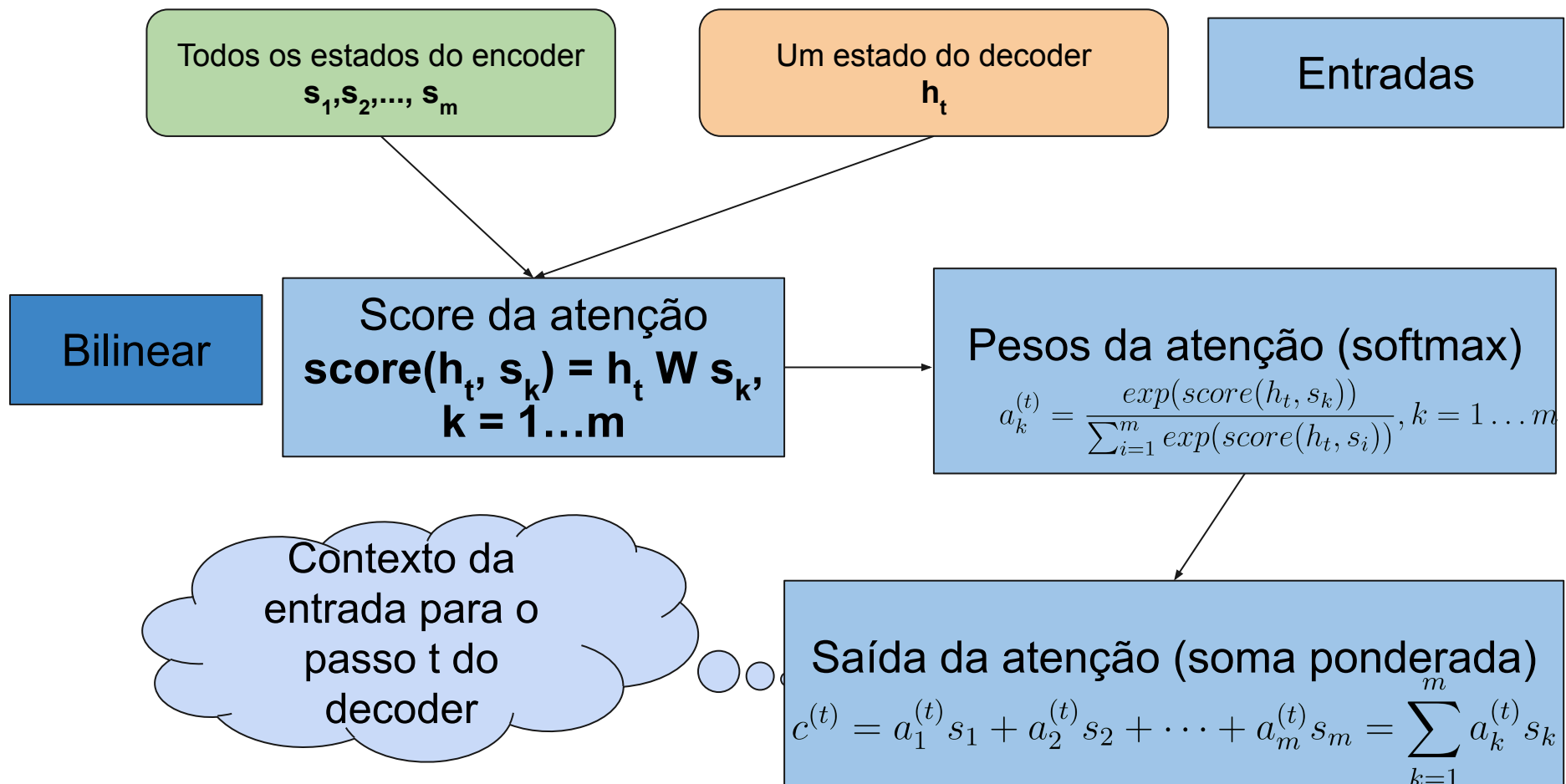
# ATENÇÃO



# Pipeline

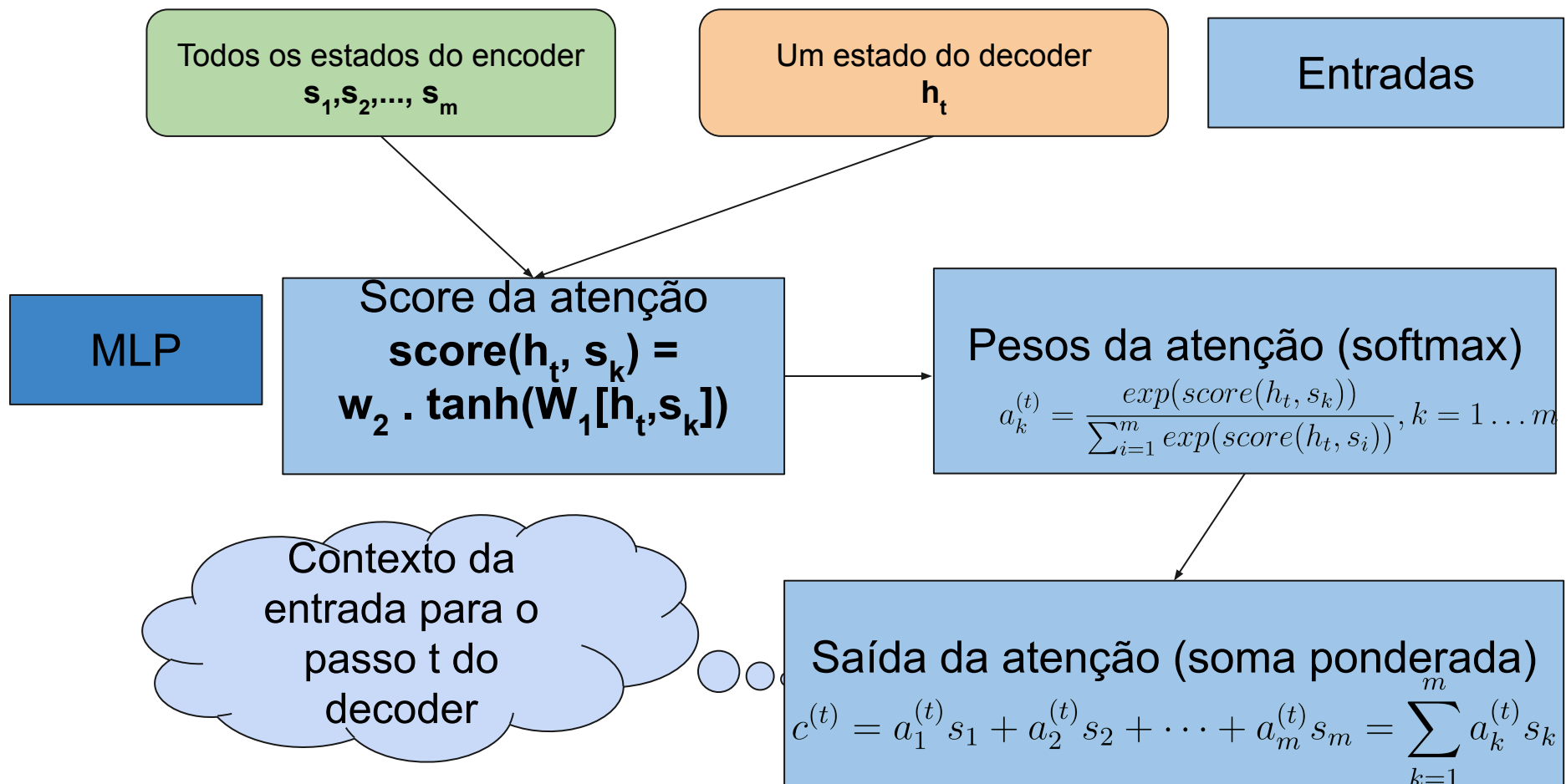


# Pipeline

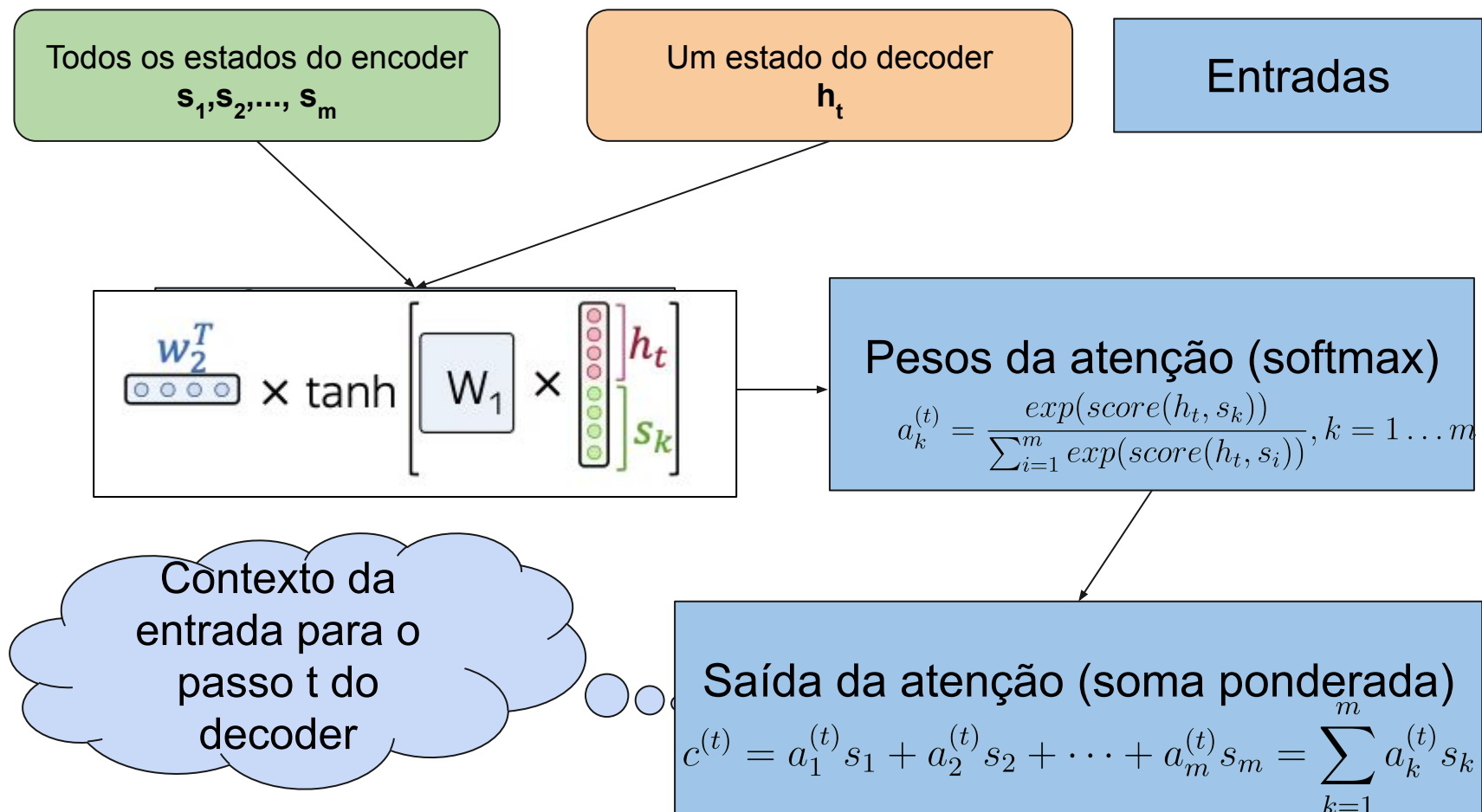




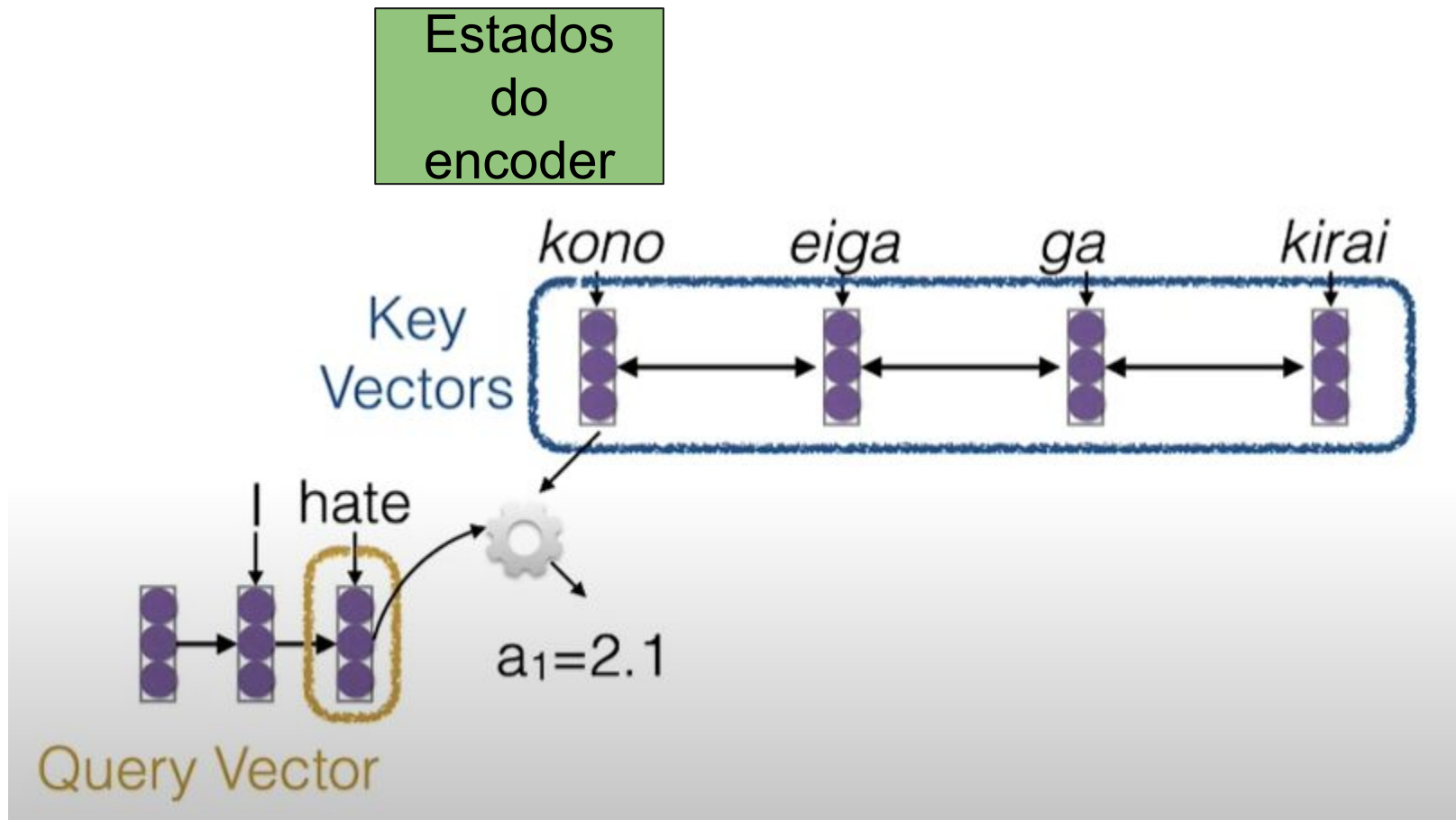
# Pipeline



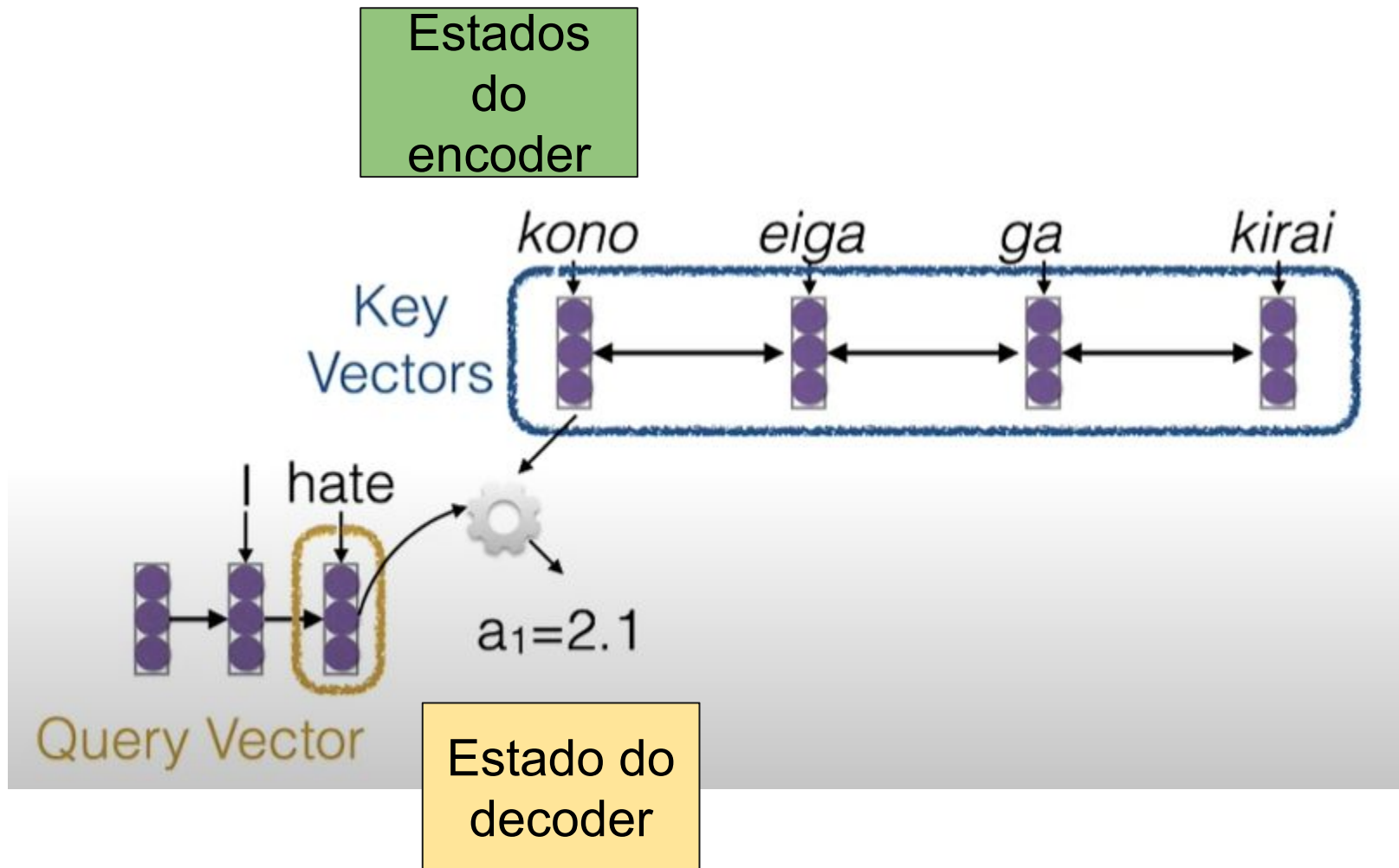
# Pipeline



# Attention



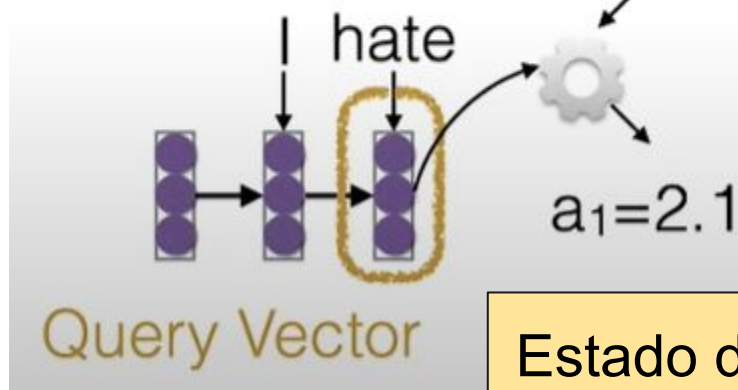
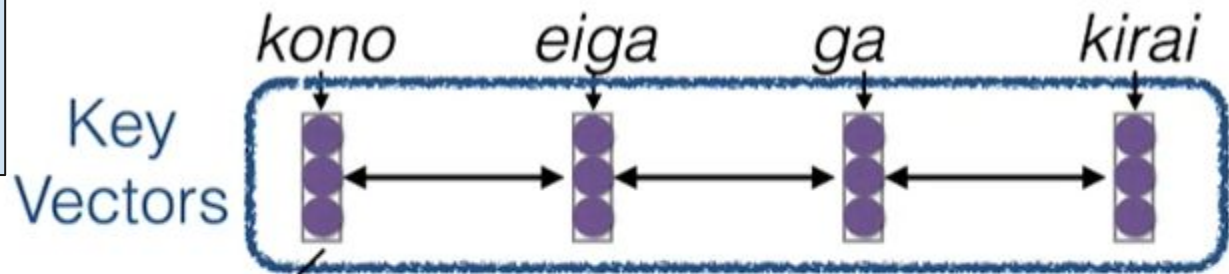
# Attention



# Attention

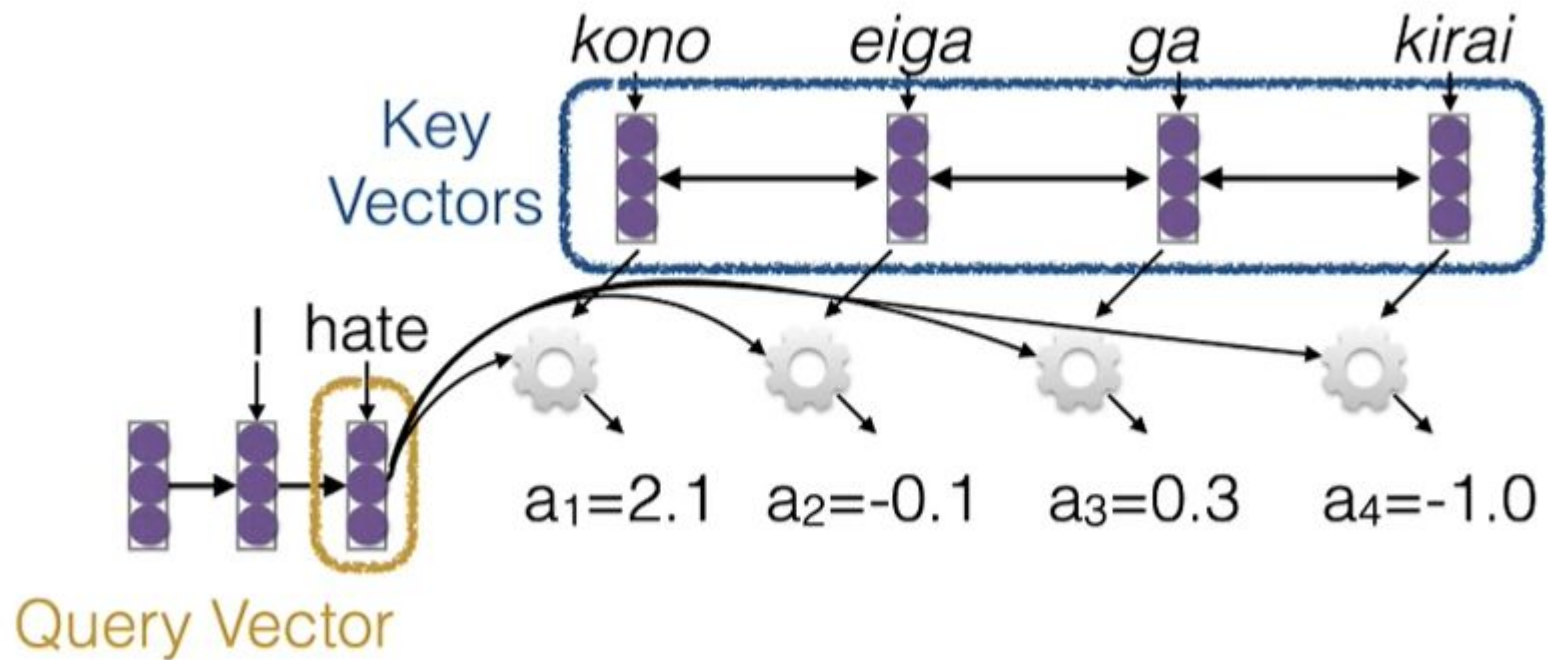
ponte entre as queries  
(o que procuramos) e  
o valor (o que  
queremos): elas vão  
apontar em como a  
atenção deve pesar  
os valores, dada a  
consulta

Estados  
do  
encoder

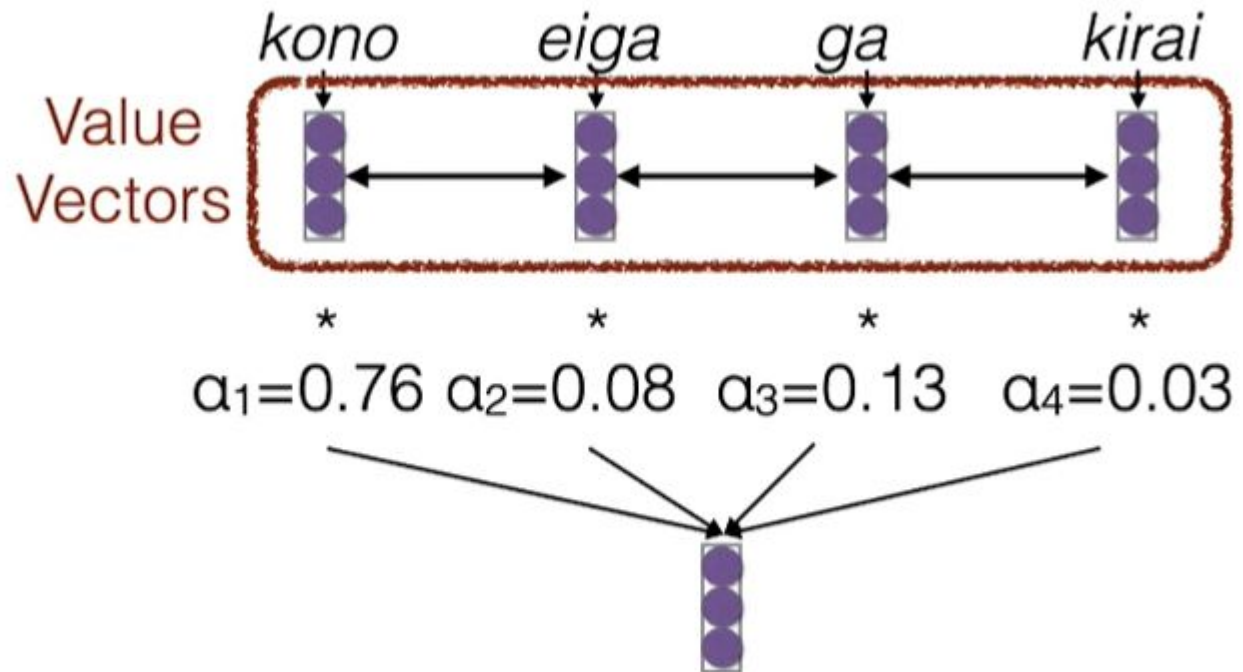


Estado do  
decoder

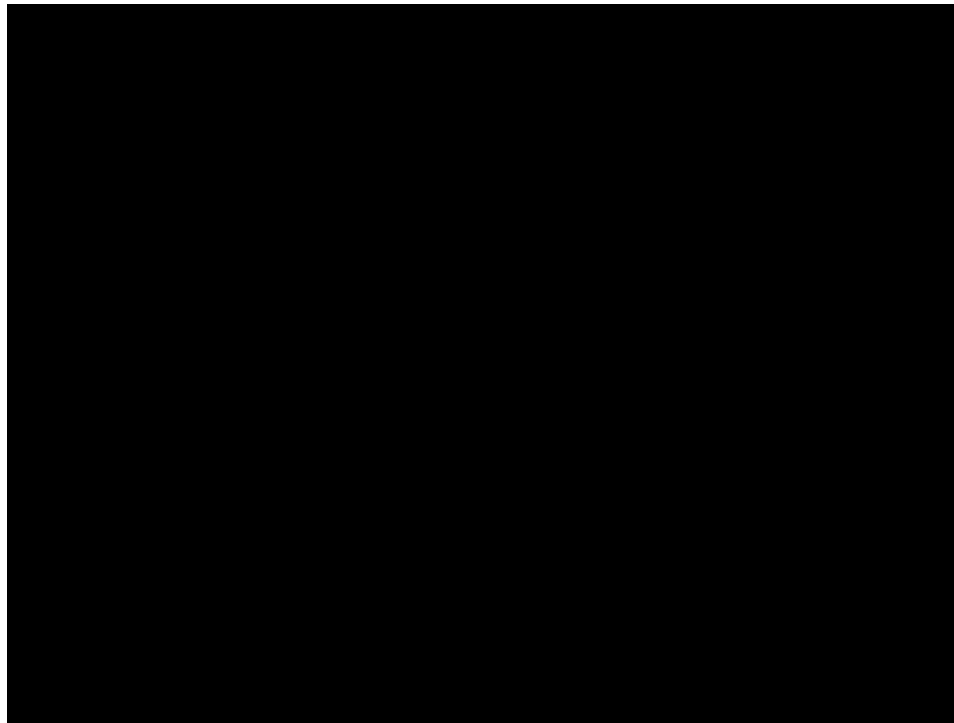
# Attention



# Attention



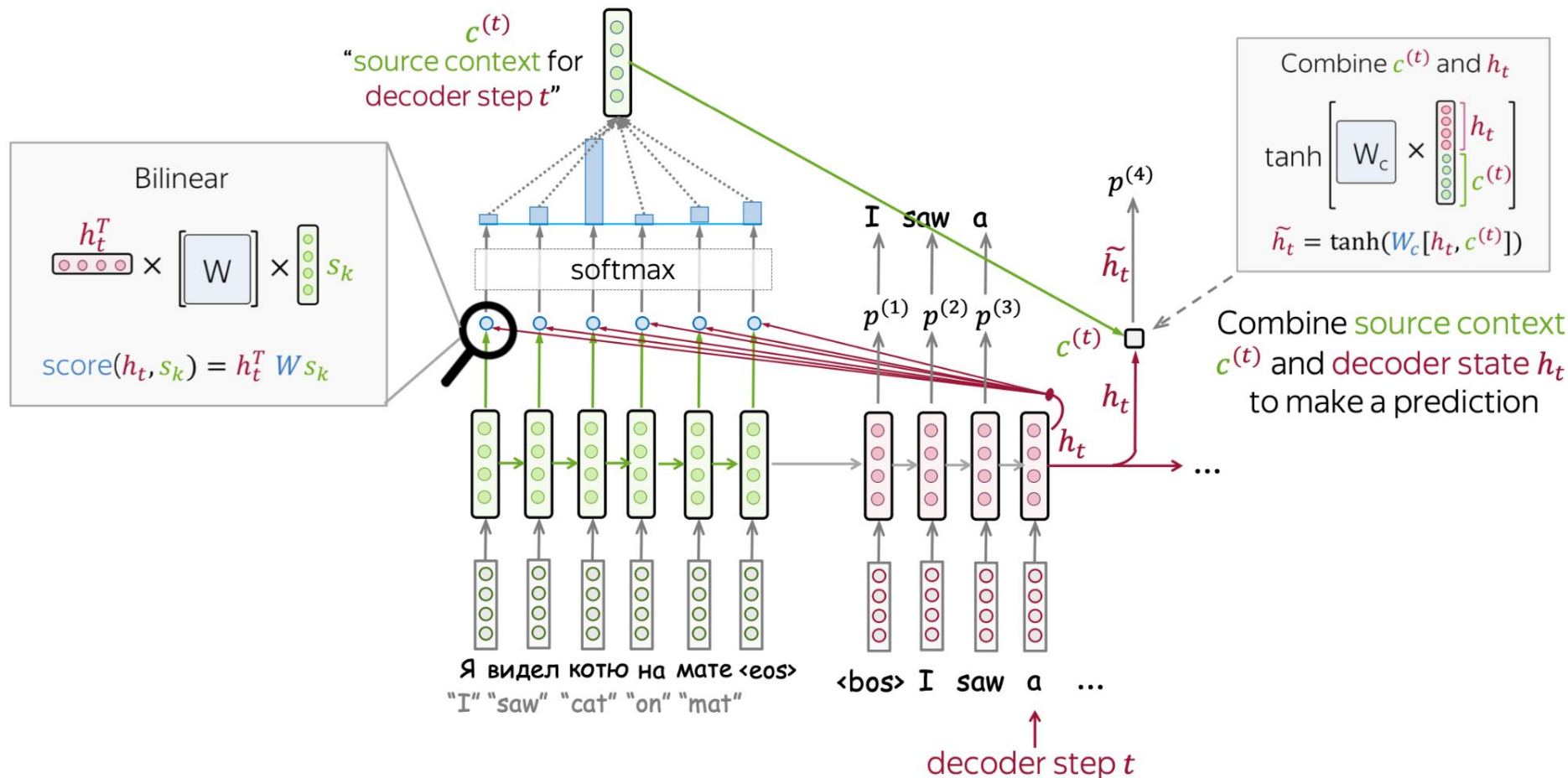
# Atenção



<http://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>



# Modelo de Luong



# Alinhamento

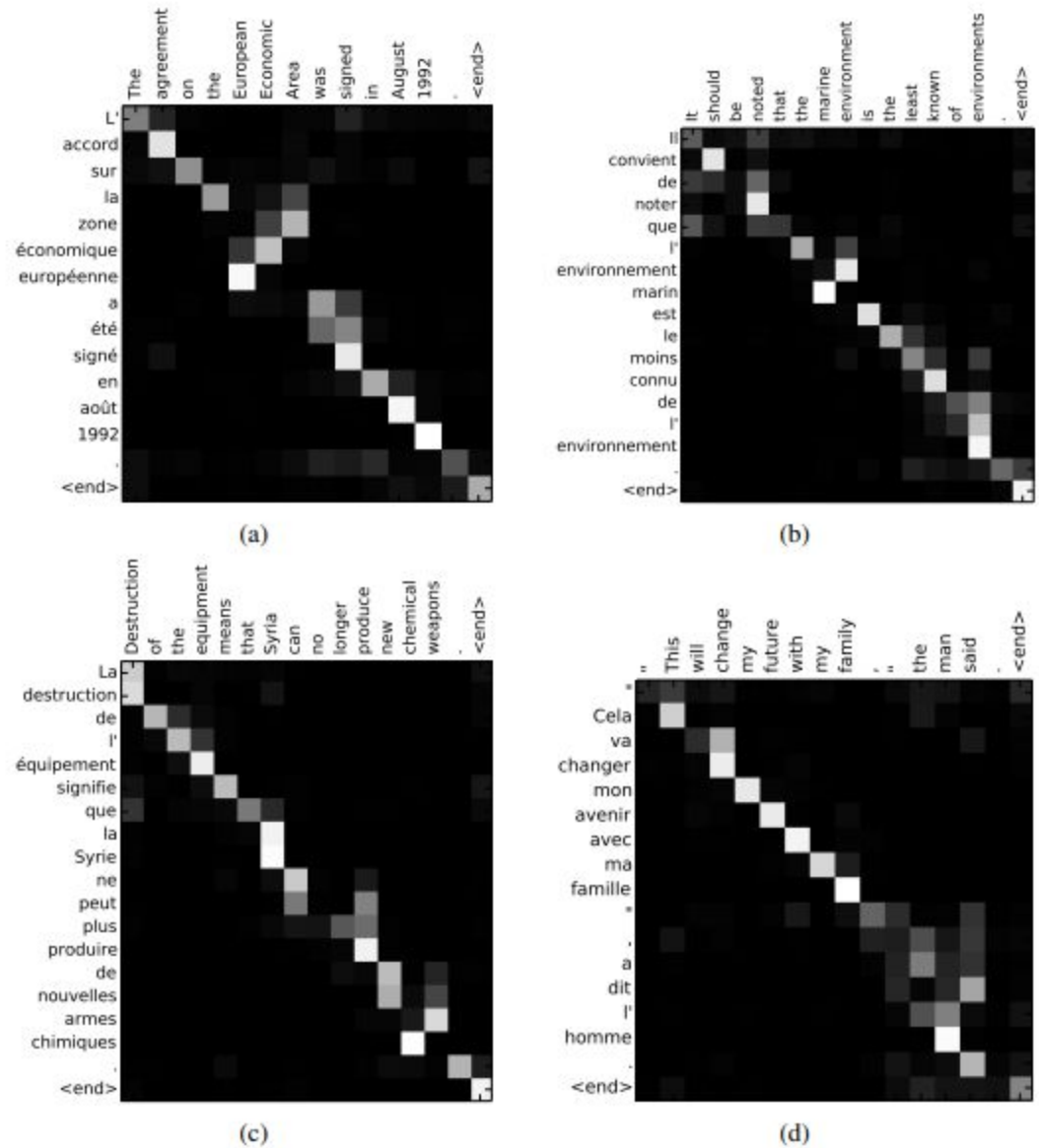
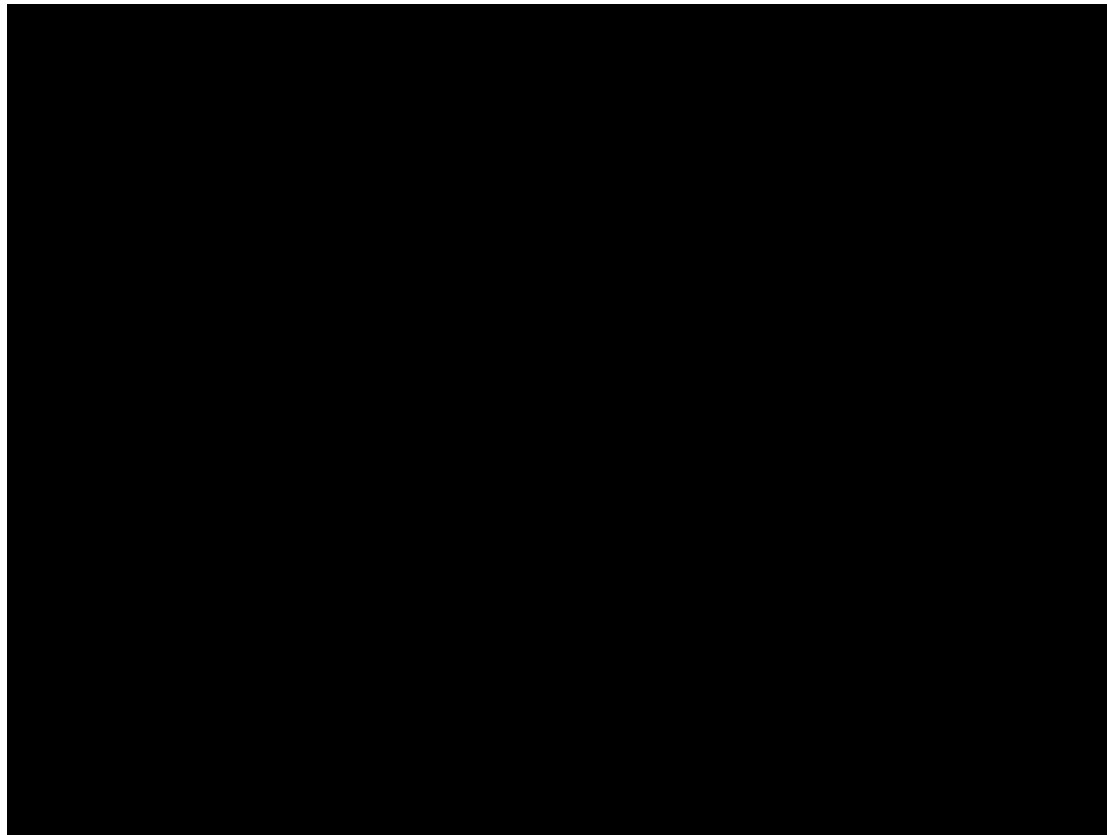
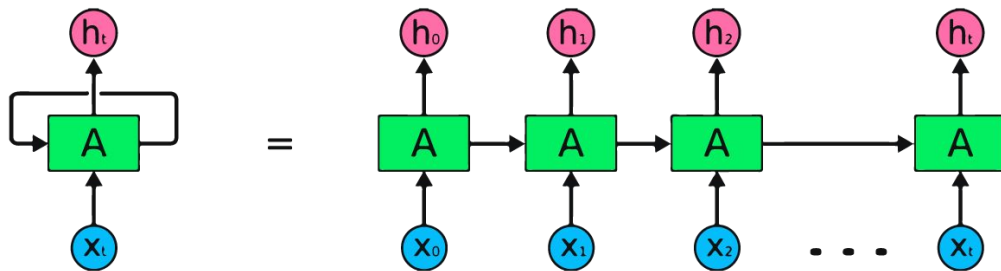


Figure 3: Four sample alignments found by RNNsearch-50. The x-axis and y-axis of each plot correspond to the words in the source sentence (English) and the generated translation (French), respectively. Each pixel shows the weight  $\alpha_{ij}$  of the annotation of the  $j$ -th source word for the  $i$ -th target word (see Eq. (6)), in grayscale (0: black, 1: white). (a) an arbitrary sentence. (b–d) three randomly selected samples among the sentences without any unknown words and of length between 10 and 20 words from the test set.

# Atenção



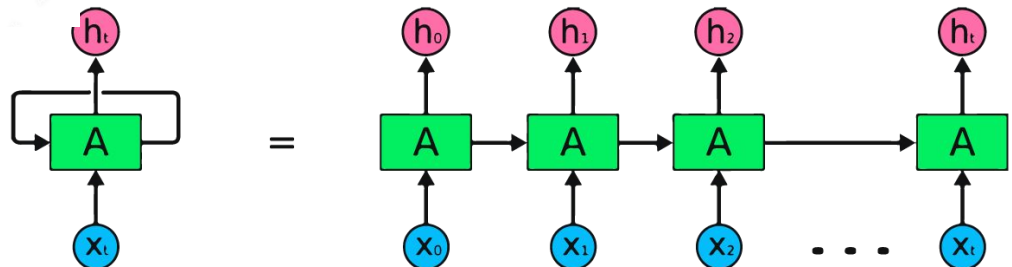
# All you need is attention



encoder



decoder



# All you need is attention



encoder



decoder



# All you need is attention



Todos os tokens  
Olham uns para os outros  
Atualizam representações

$\mathbf{N} \times \mathbf{N}$

encoder



decoder

self-attention



# All you need is attention



Todos os tokens  
**Olham uns para os outros**  
Atualizam representações

$N \times N$

encoder



decoder

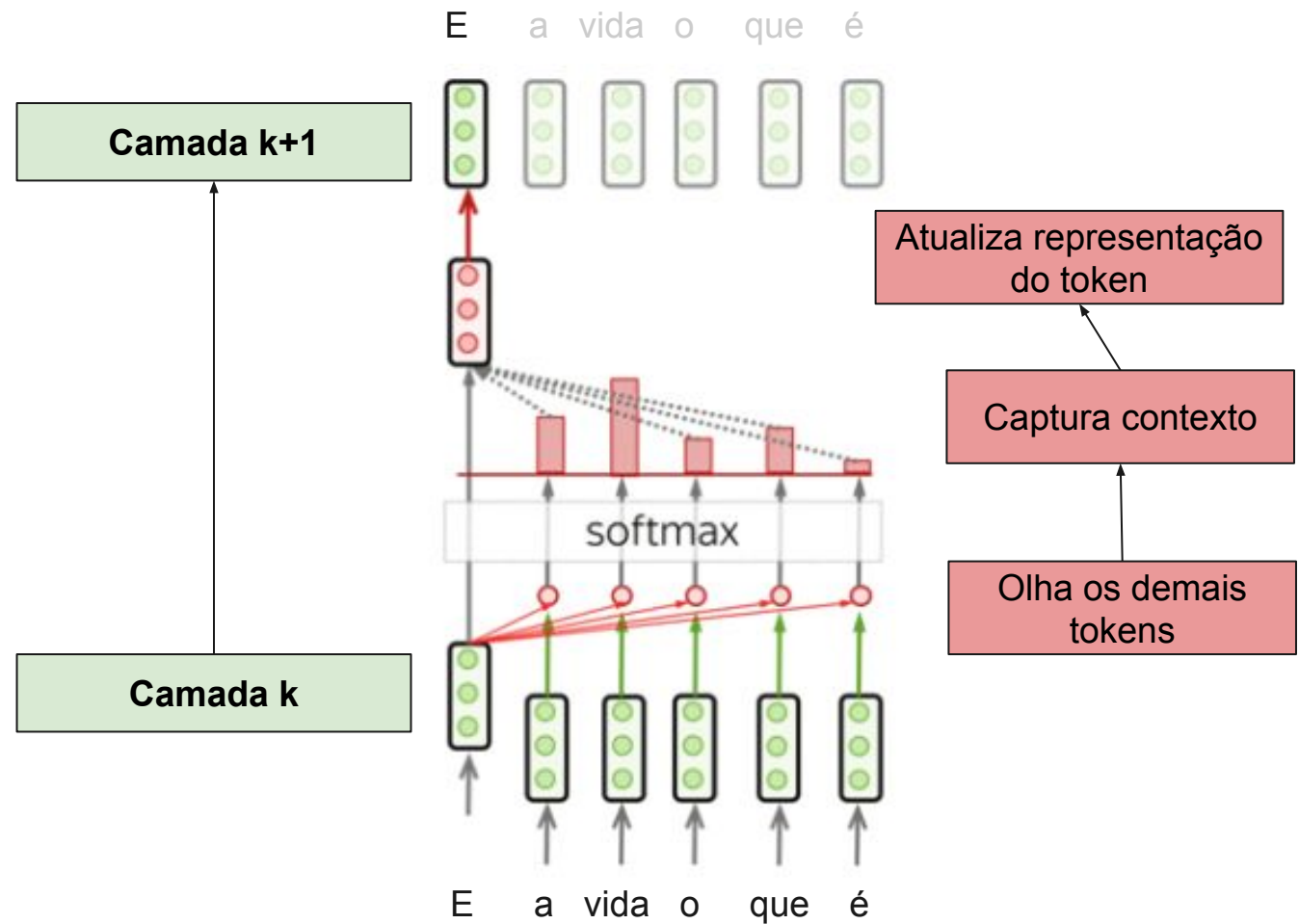
Token na saída corrente  
Olha para os tokens anteriores  
Olha para as representações da entrada  
Atualiza representações



$N \times N$

[Google animation](#)

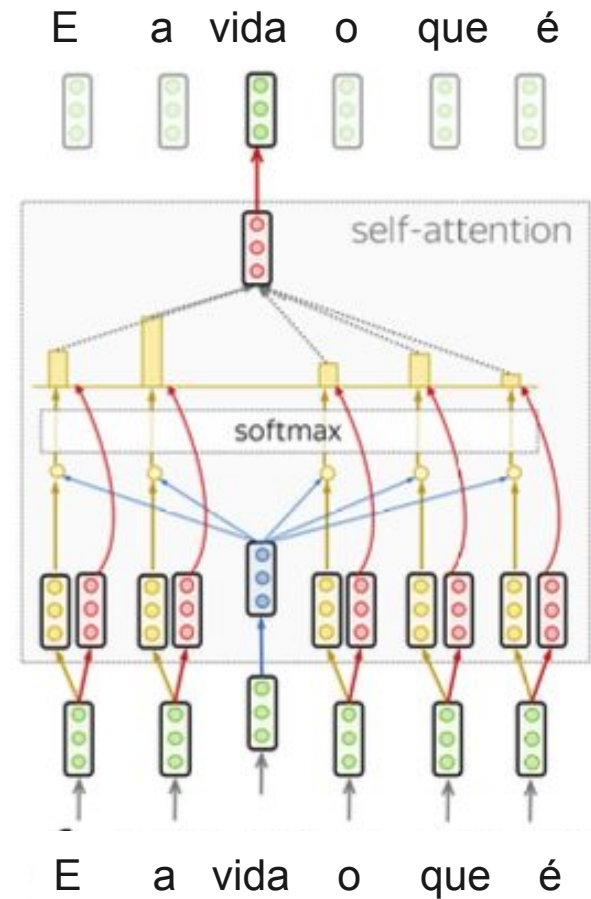
# Self-attention





# Self-attention, Q, K, V

Três "representações" para cada token

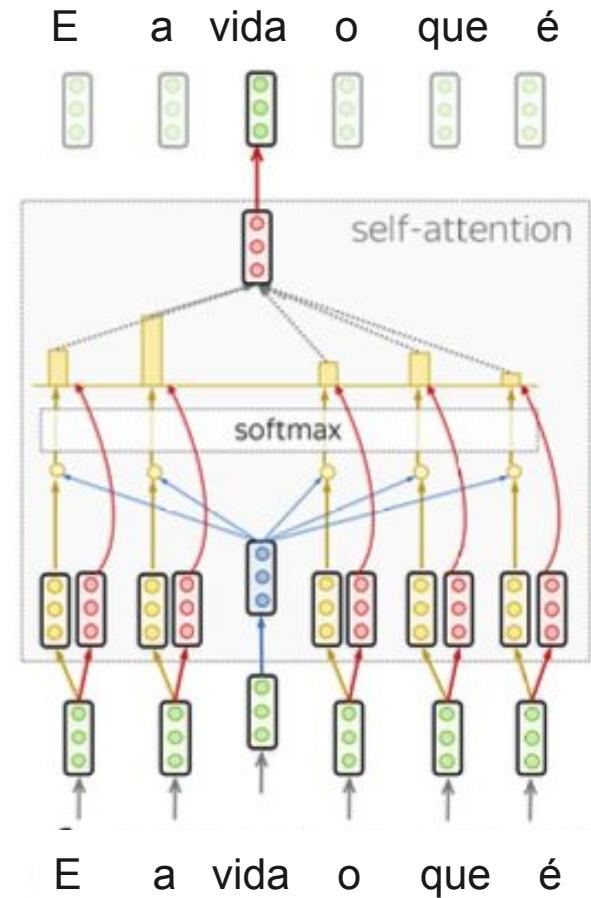


# Self-attention, Q, K, V

Três "representações" para cada token

$$[W_Q] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix}$$

Query: "Você tem  
essa informação?"



# Self-attention, Q, K, V

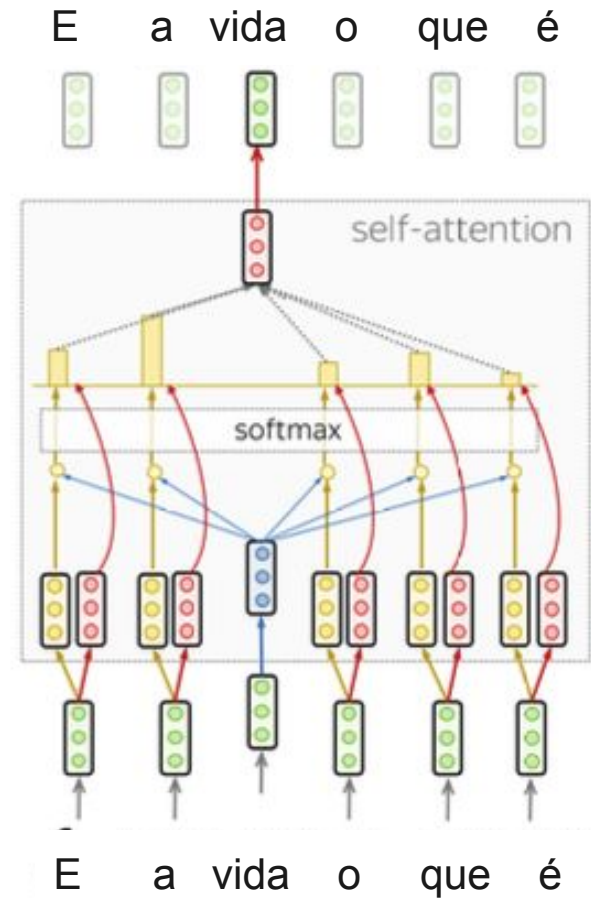
Três "representações" para cada token

$$[W_Q] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix}$$

Query: "Você tem  
essa informação?"

$$[W_K] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{yellow} \\ \text{yellow} \\ \text{yellow} \end{bmatrix}$$

Key: "Opa, eu tenho a  
informação que você quer.  
Me dê relevância (mais  
peso)"



# Self-attention, Q, K, V

Três "representações" para cada token

$$[W_Q] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix}$$

Query: "Você tem  
essa informação?"

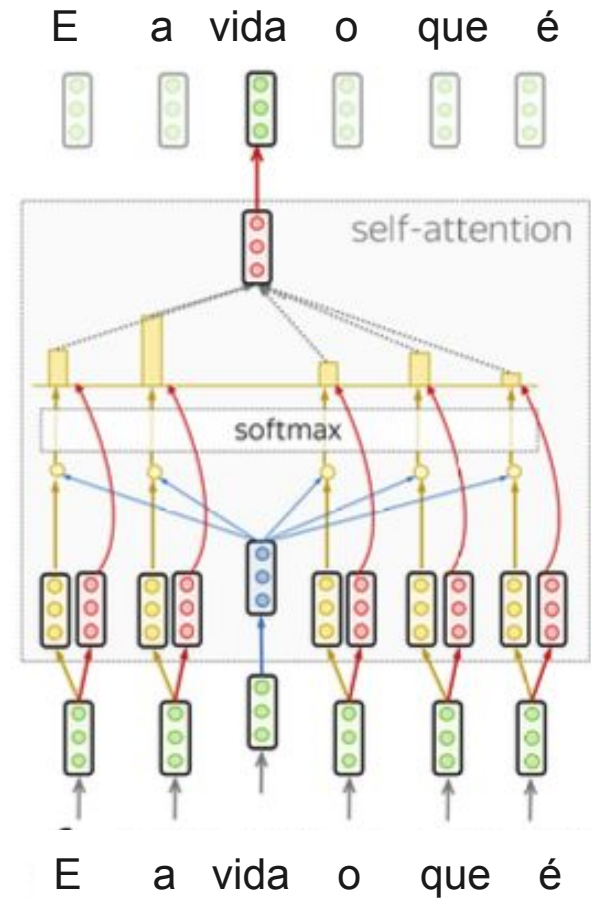
$$[W_K] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{yellow} \\ \text{yellow} \\ \text{yellow} \end{bmatrix}$$

Key: "Opa, eu tenho a  
informação que você quer.  
Me dê relevância (mais  
peso)"

$$[W_V] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{red} \\ \text{red} \\ \text{red} \end{bmatrix}$$

Value: "Aqui está a  
informação."

Soma ponderada

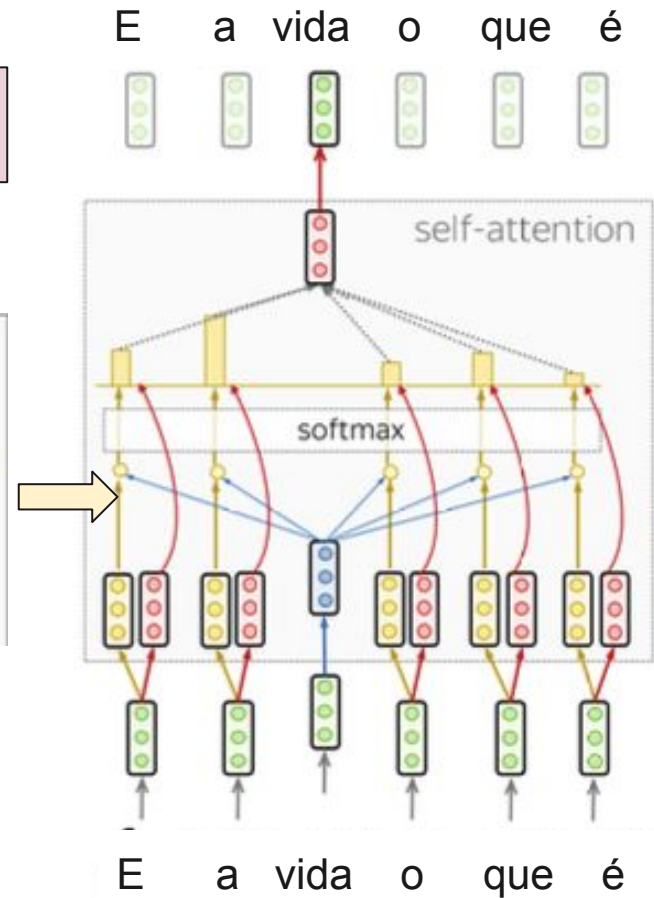


# Self-attention, Q, K, V

Três "representações" para cada token

Produto interno  
"normalizado"

$$\frac{1}{\sqrt{d_k}} \cdot q \times k^T$$
$$\text{score}(q, k) = \frac{qk^T}{\sqrt{d_k}}$$



# Self-attention, Q, K, V

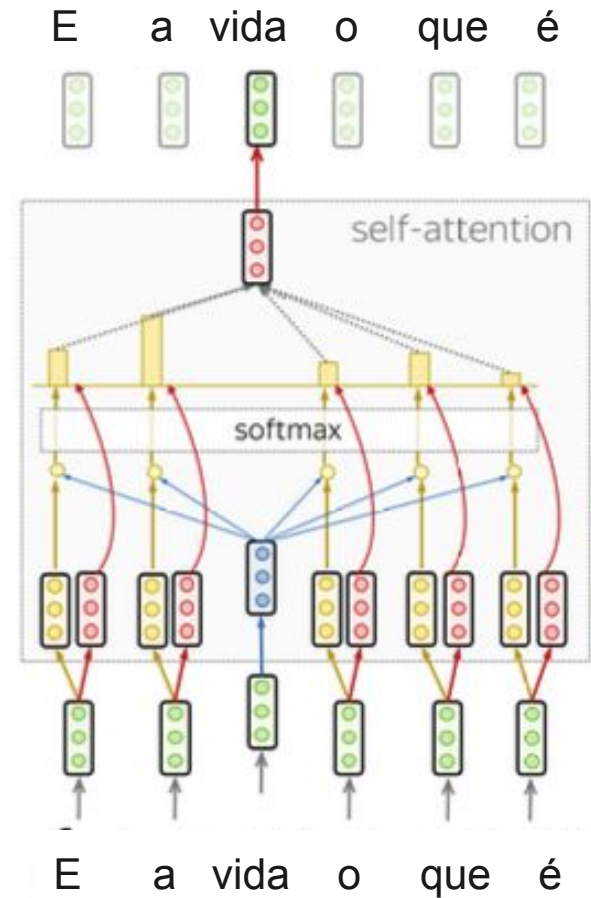
Três "representações" para cada token

$$\text{Attention}(q, k, v) = \text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)v$$

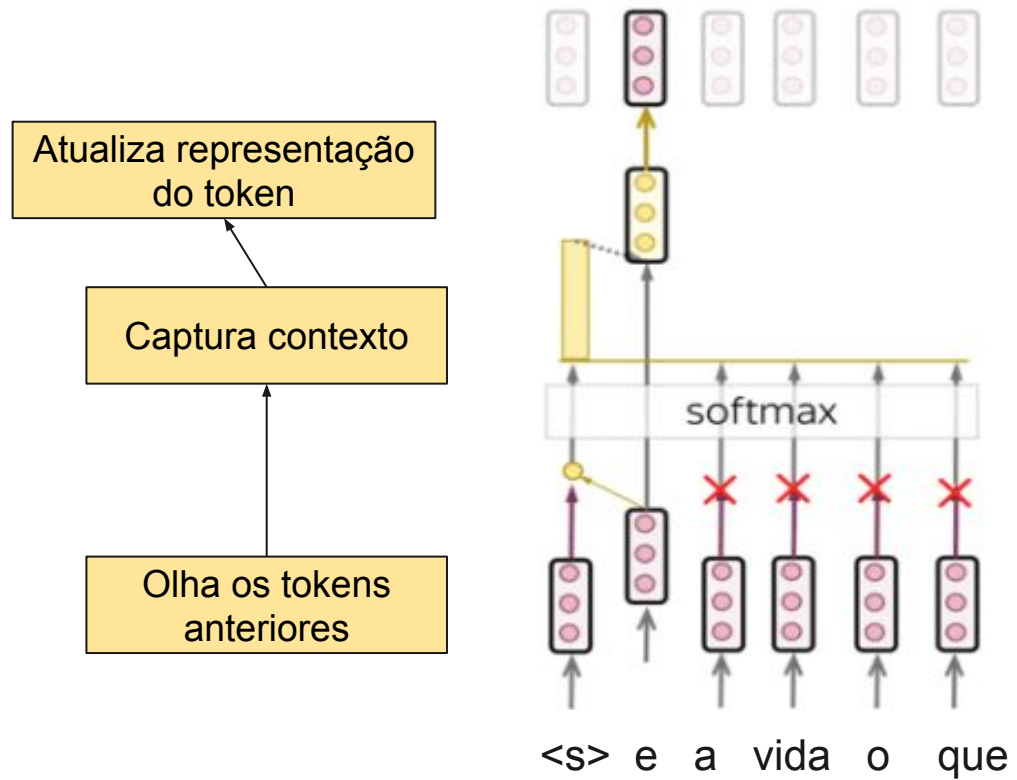
Pesos da atenção

Dimensão de K, V

de para

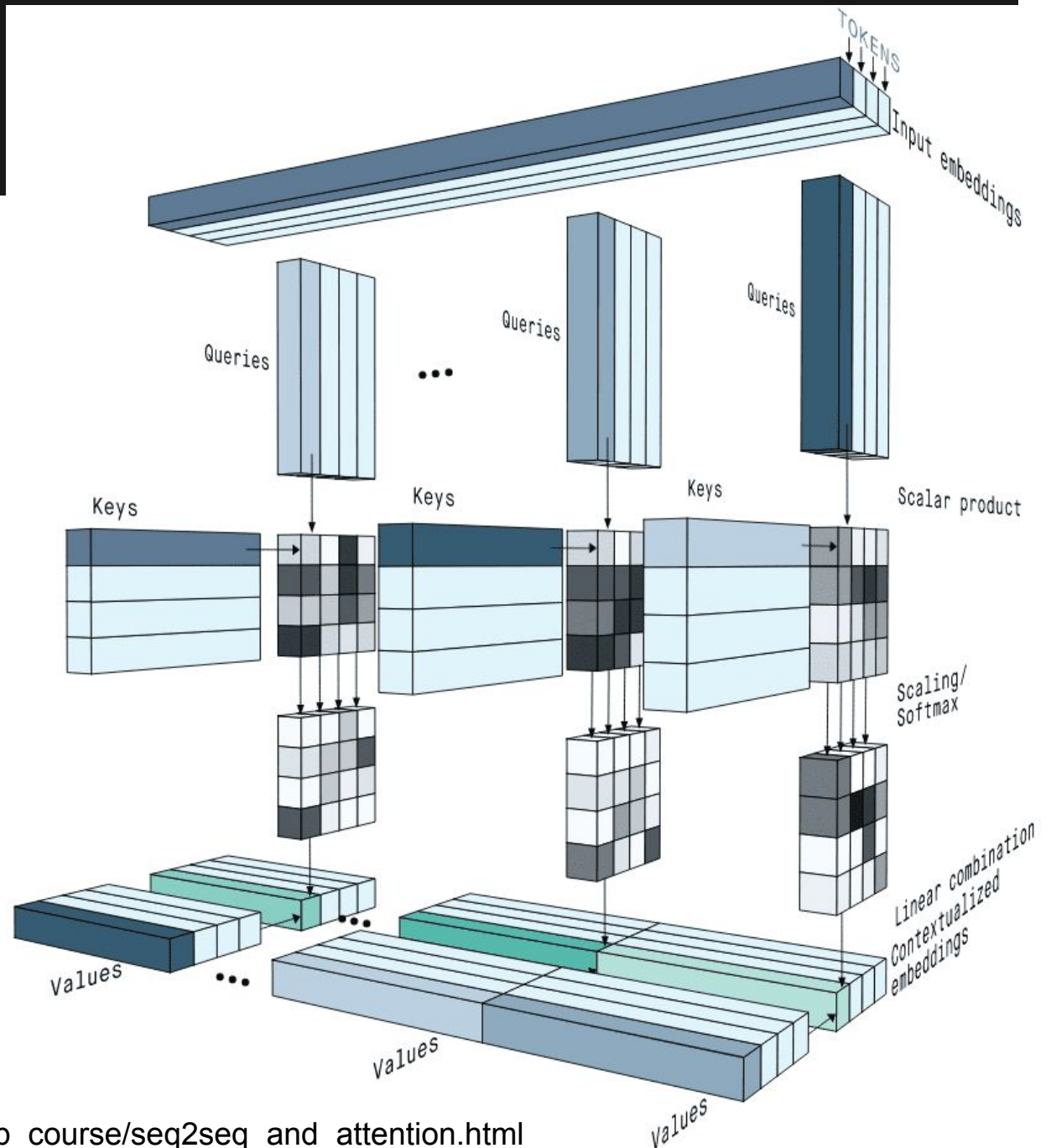
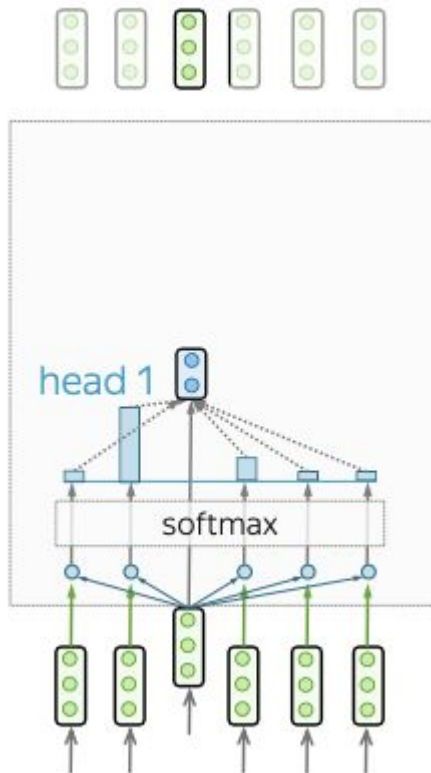


# Atenção mascarada do decoder



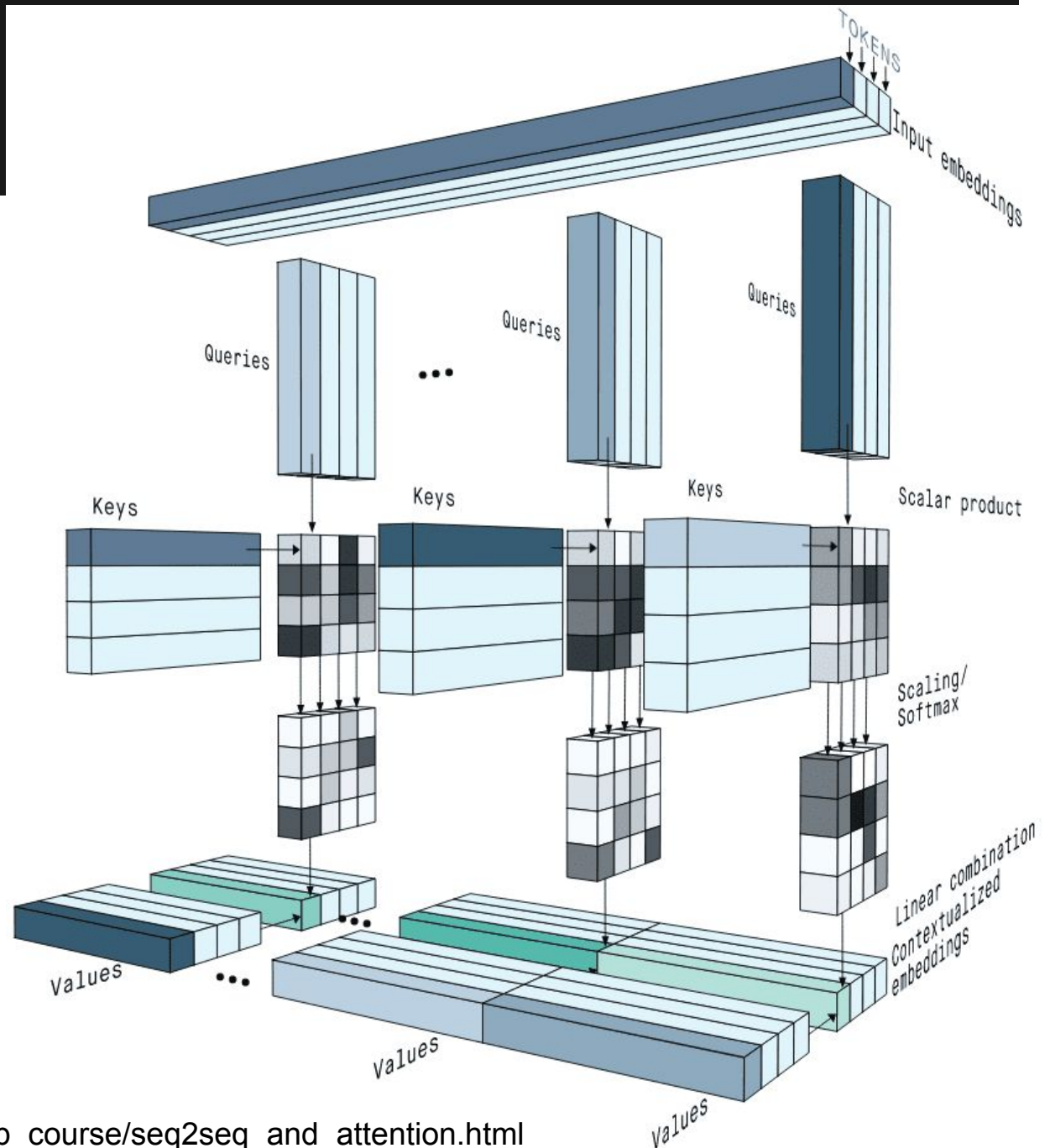
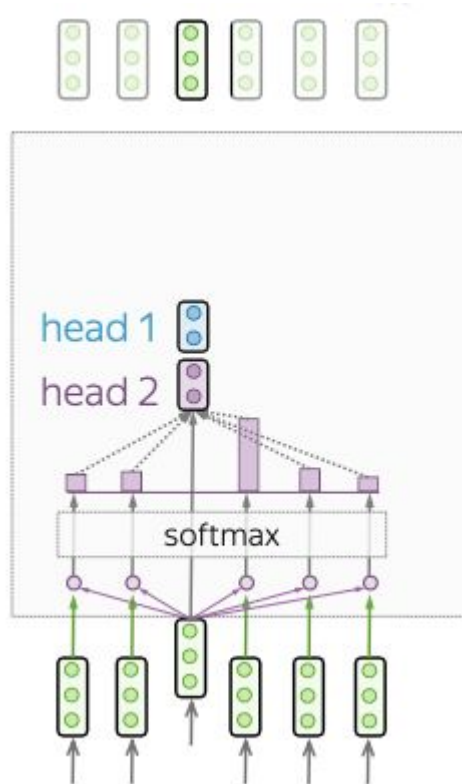


# Multi-Head attention

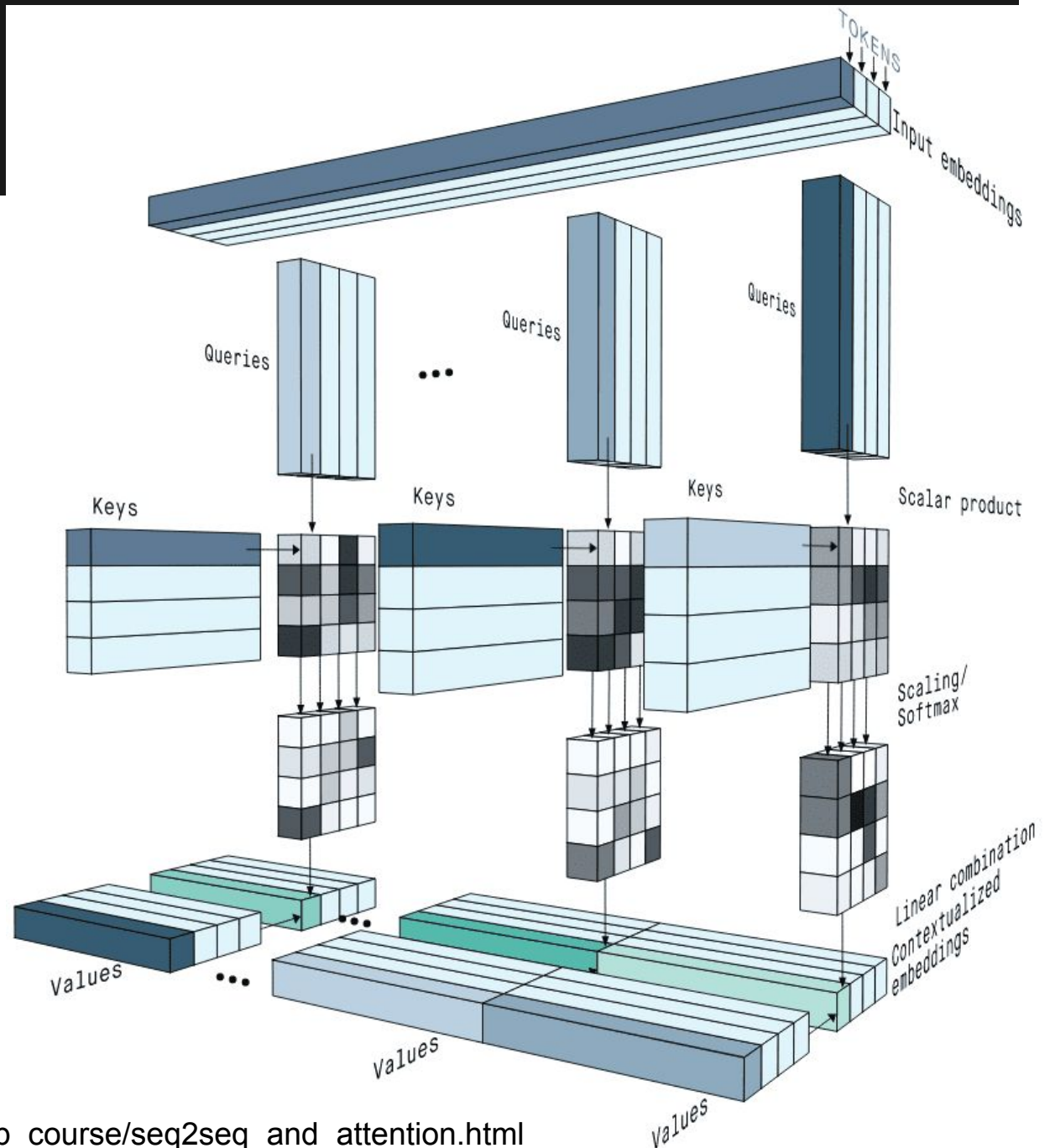
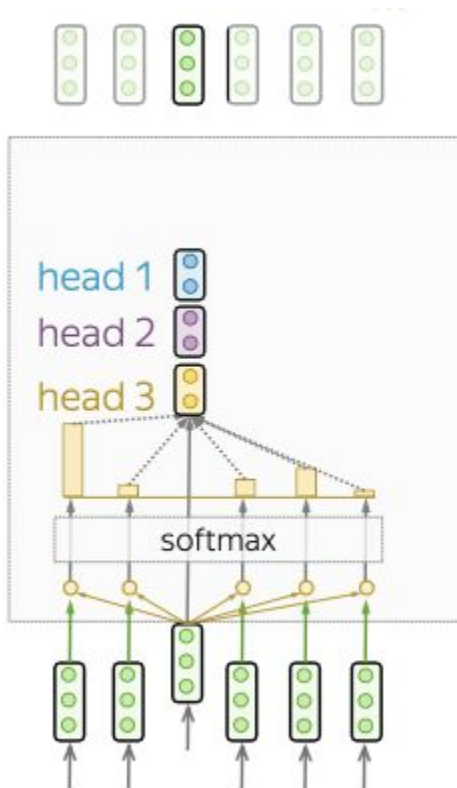




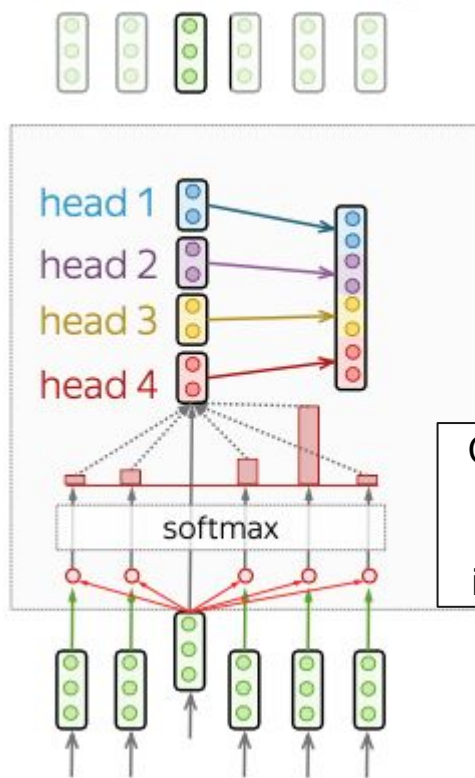
# Multi-Head attention



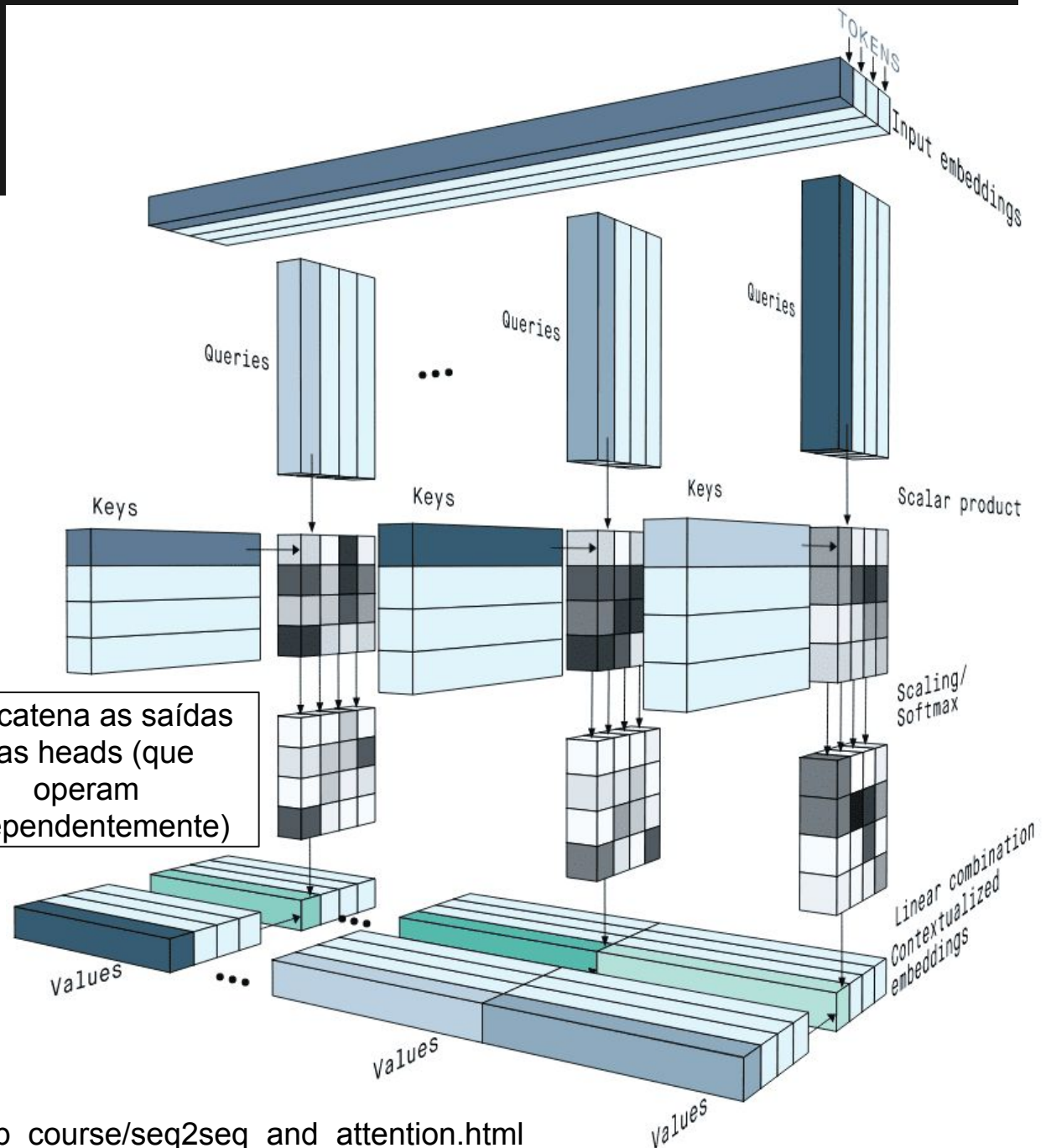
# Multi-Head attention



# Multi-Head attention

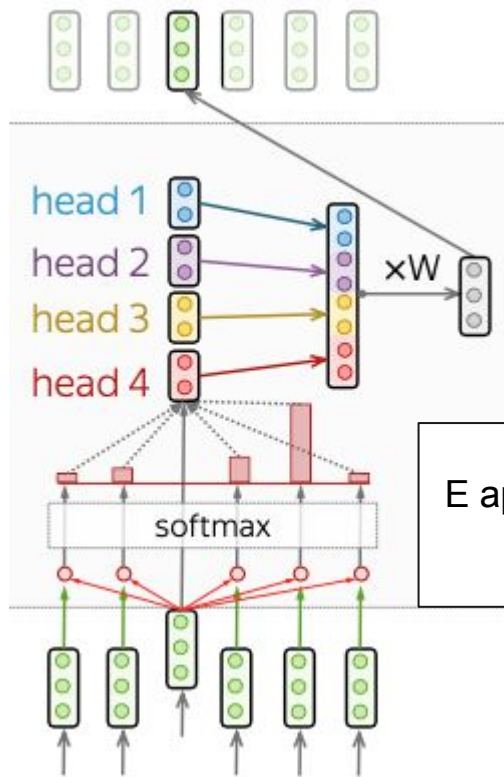


Concatena as saídas das heads (que operam independentemente)

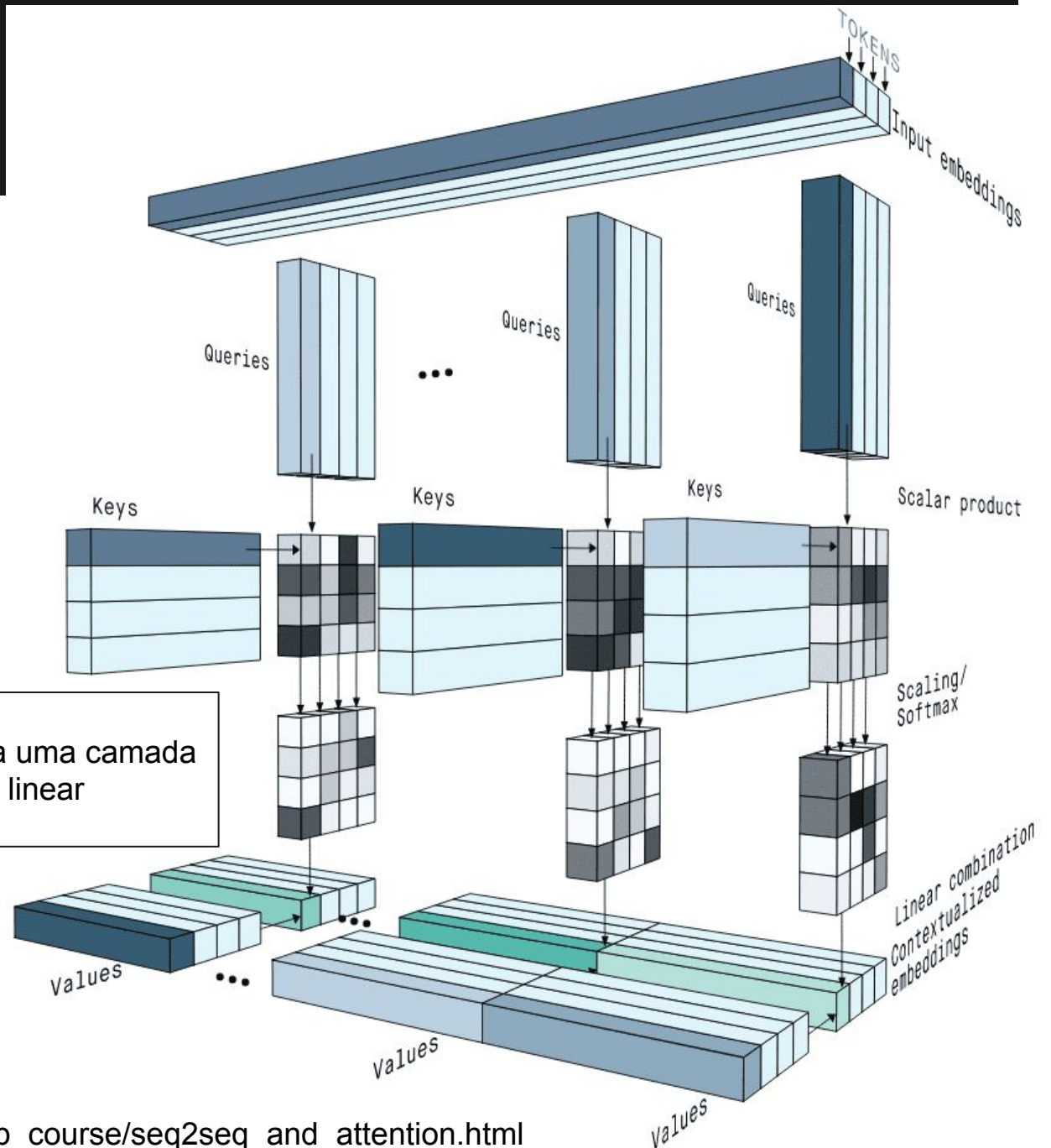




# Multi-Head attention



E aplica uma camada linear



# All you need is attention



encoder



decoder

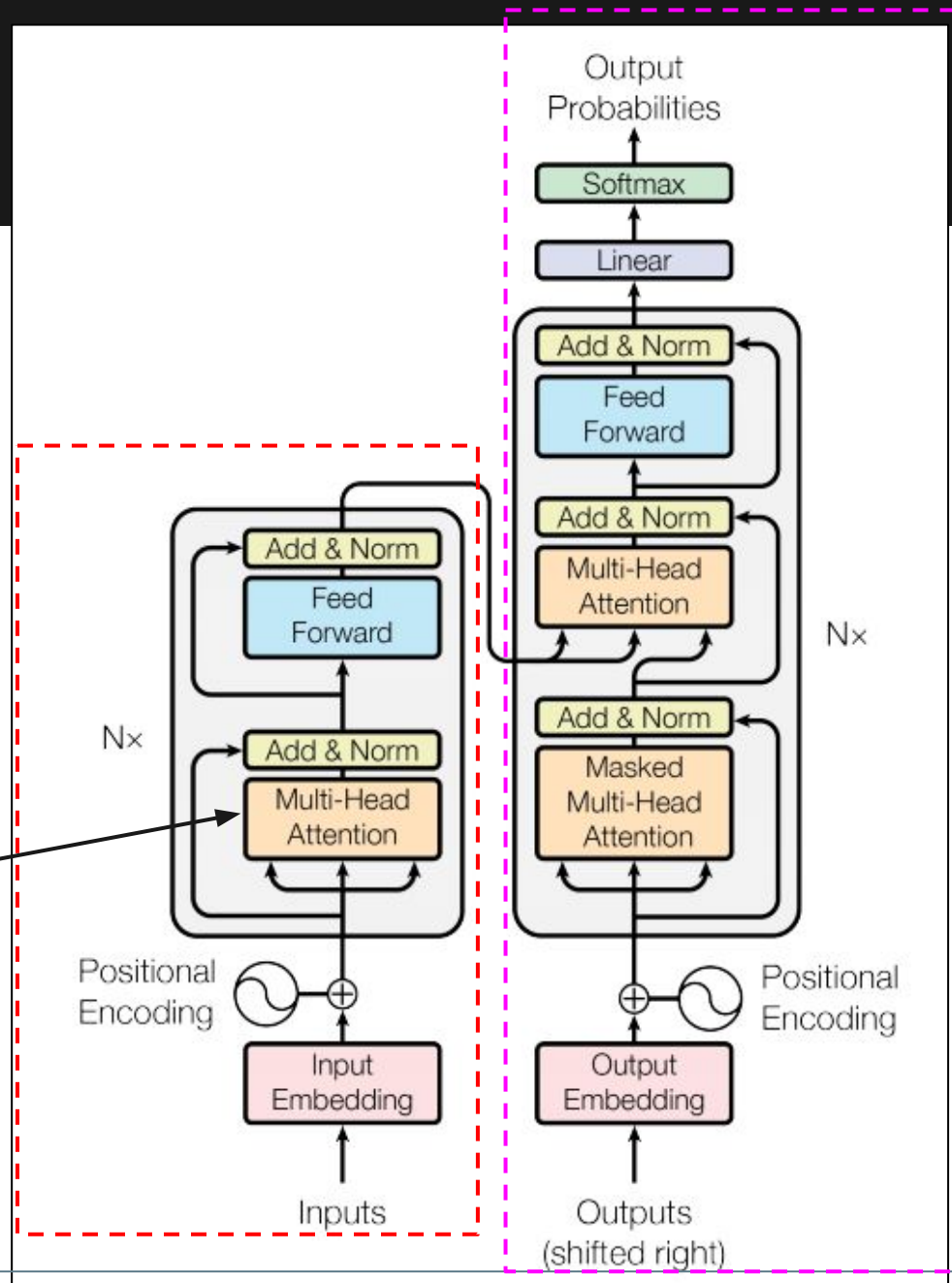


# Só vantagens?

- a) Ver todos os tokens "de uma vez".
- b) Paralelizar por tokens.
- c) Ordem dos tokens.
- d) Tamanho da sequência de entrada.

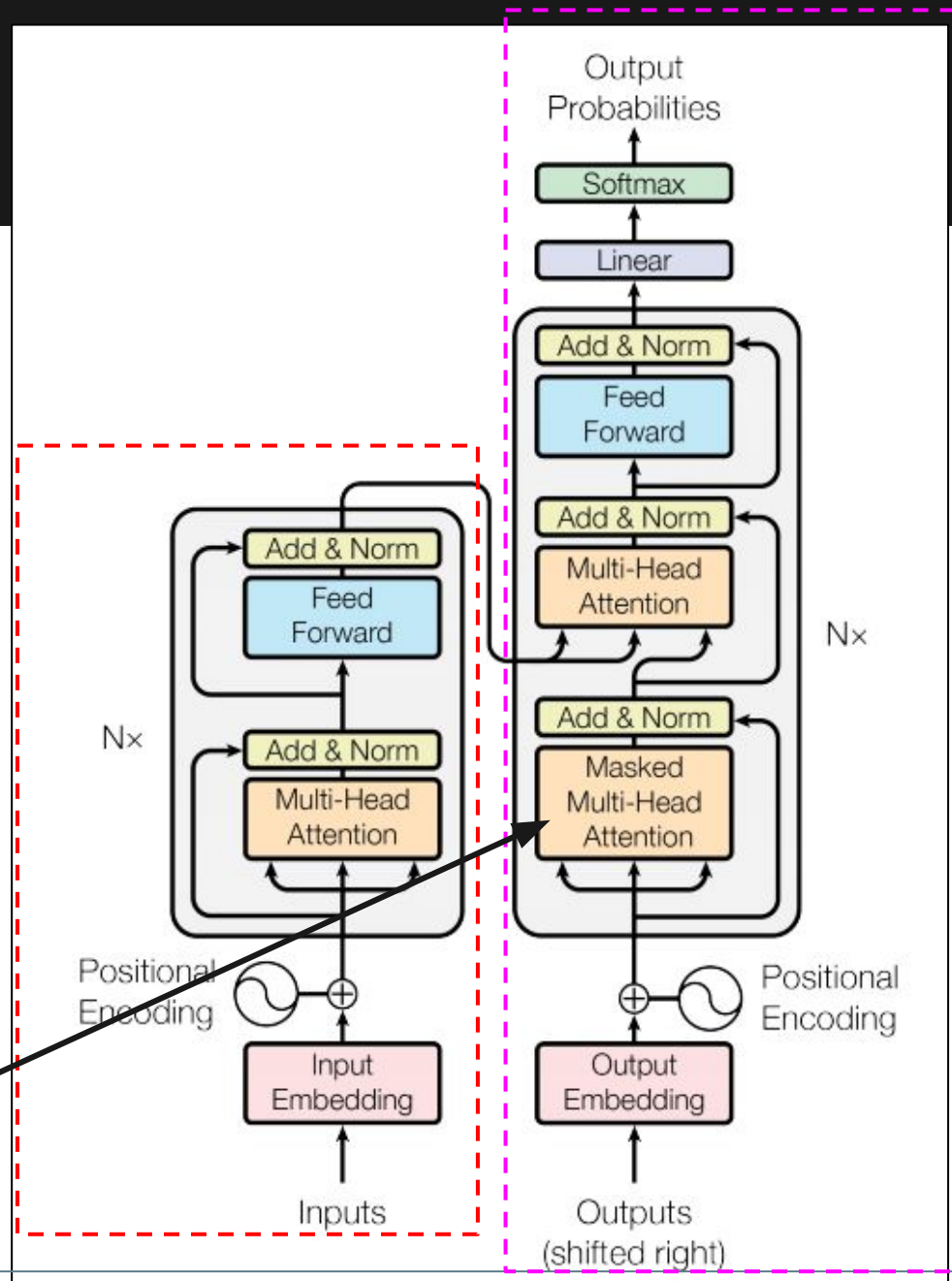
# Transformers

Self attention



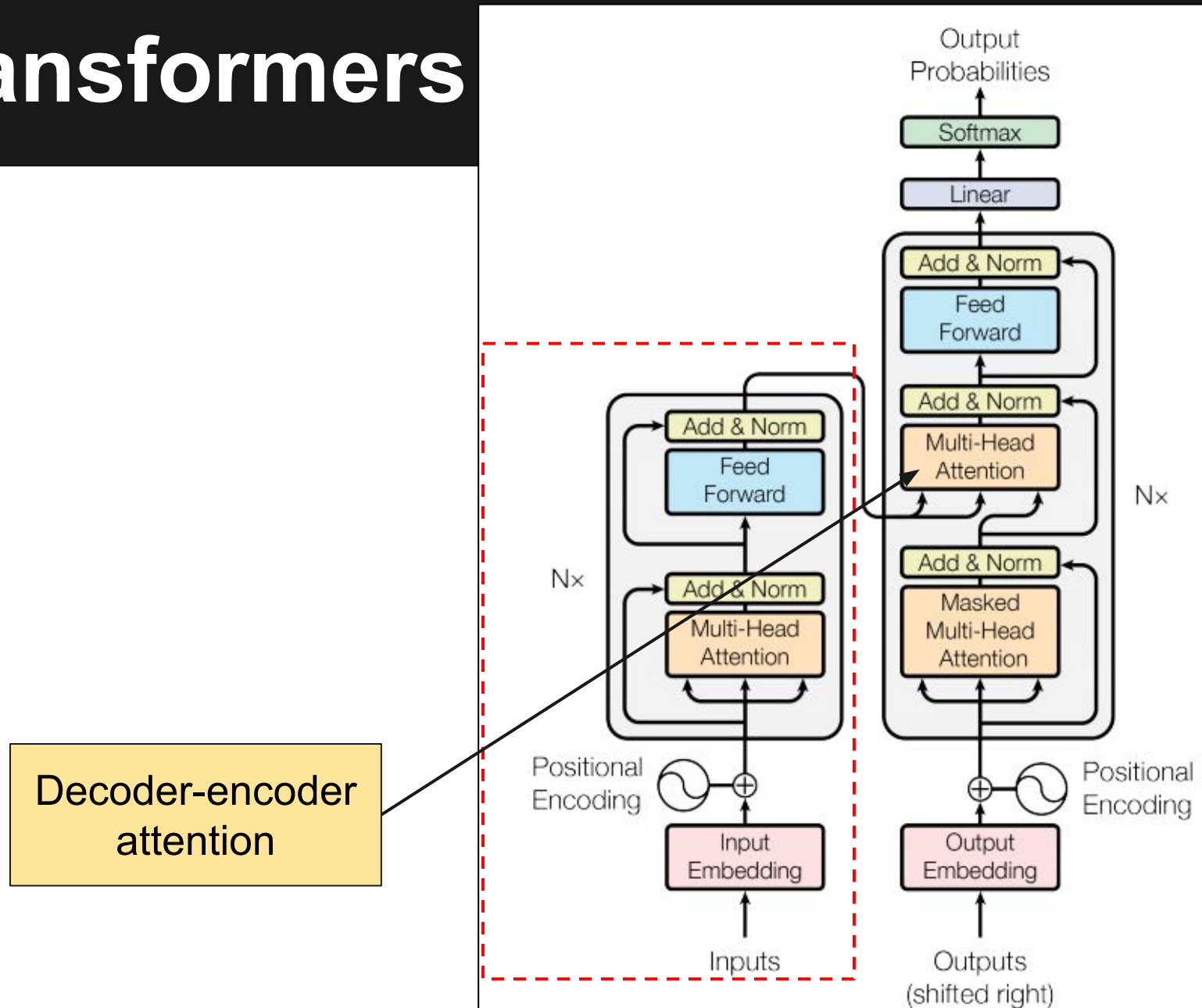
# Transformers

Masked  
attention

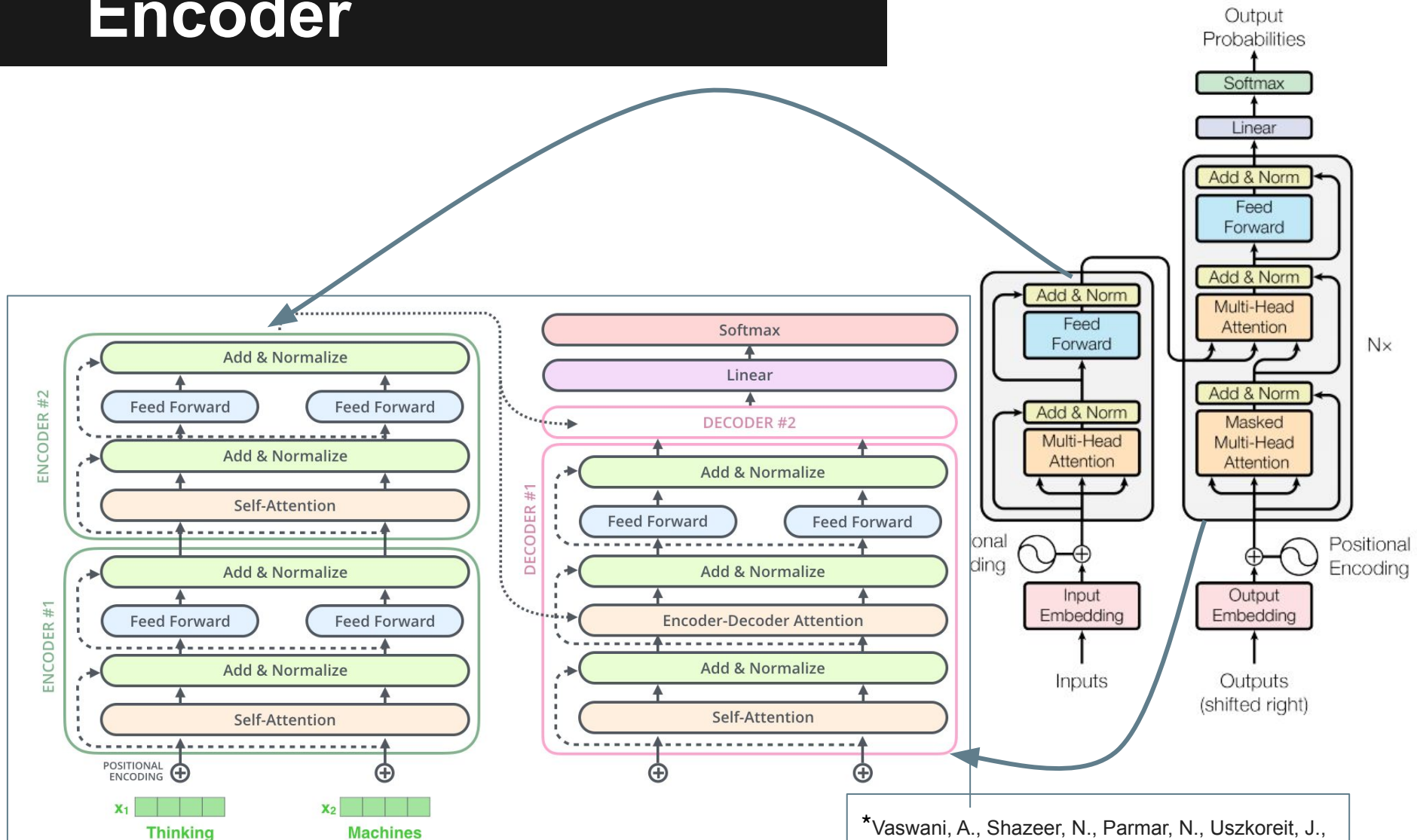




# Transformers



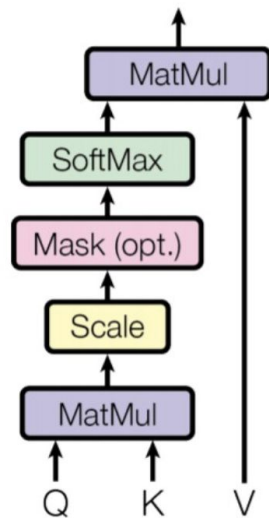
# Encoder



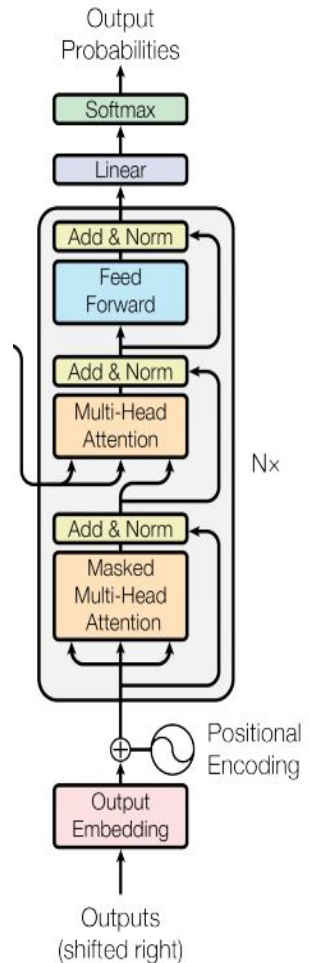
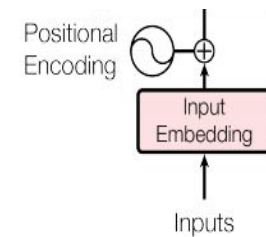
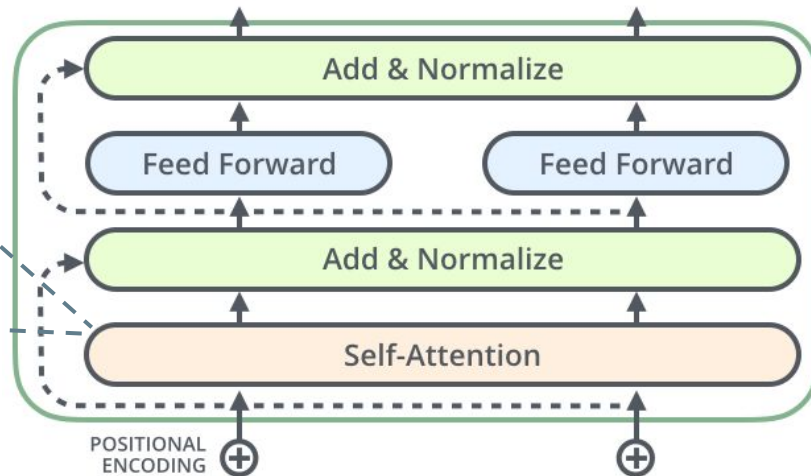
\* <http://jalamar.github.io/illustrated-transformer/>

\*Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). **Attention is all you need.** In *Advances in neural information processing systems* (pp. 5998-6008).

# Self-attention

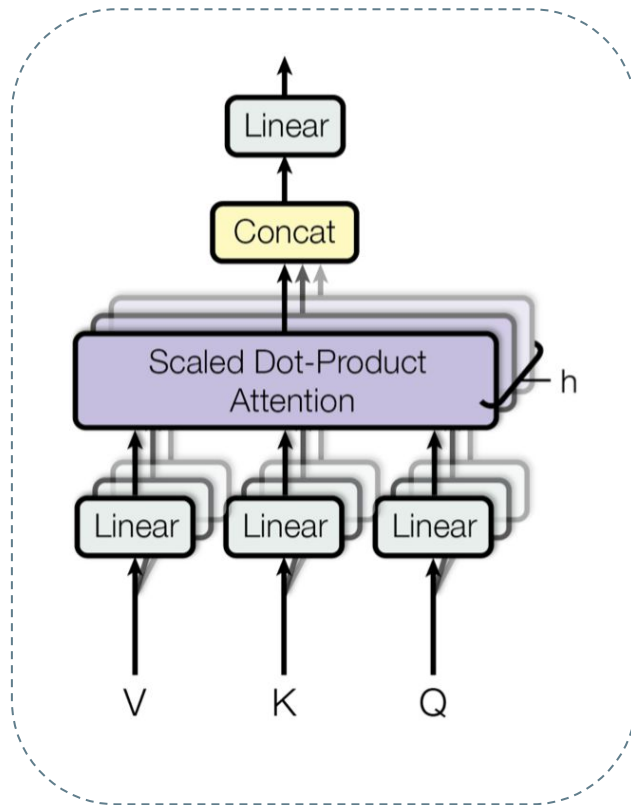


Scaled Dot-Product  
attention

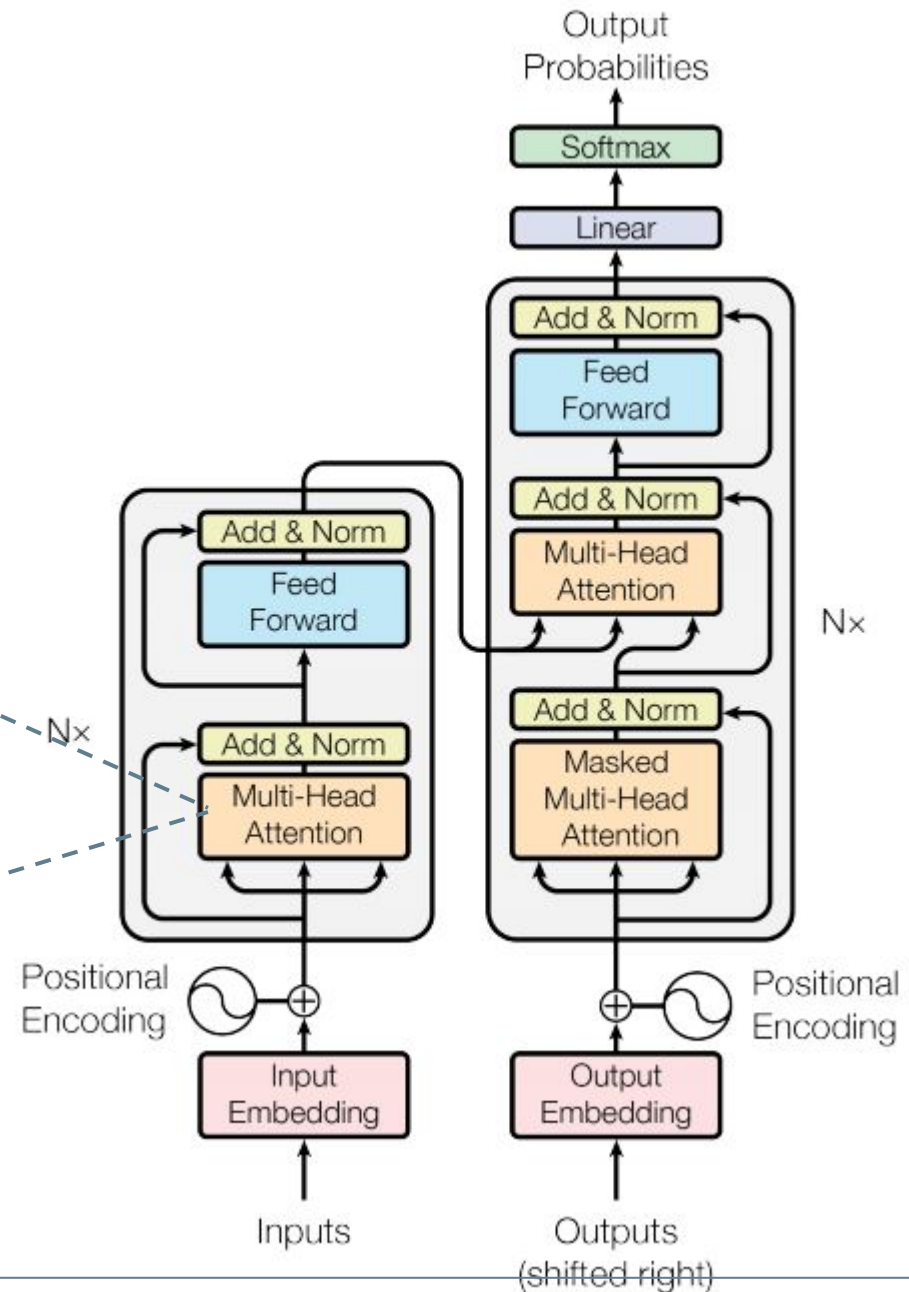


\*Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).

# Multi-head attention

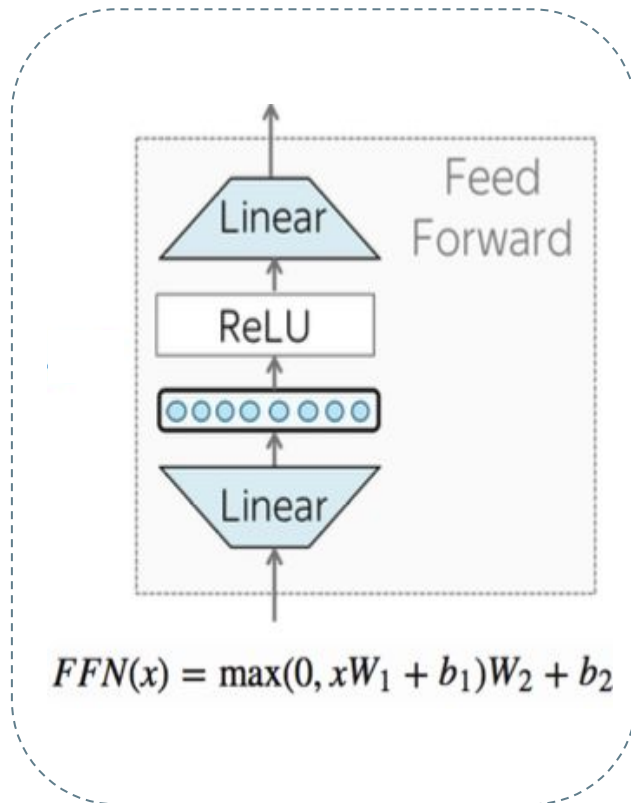


Multi-head attention

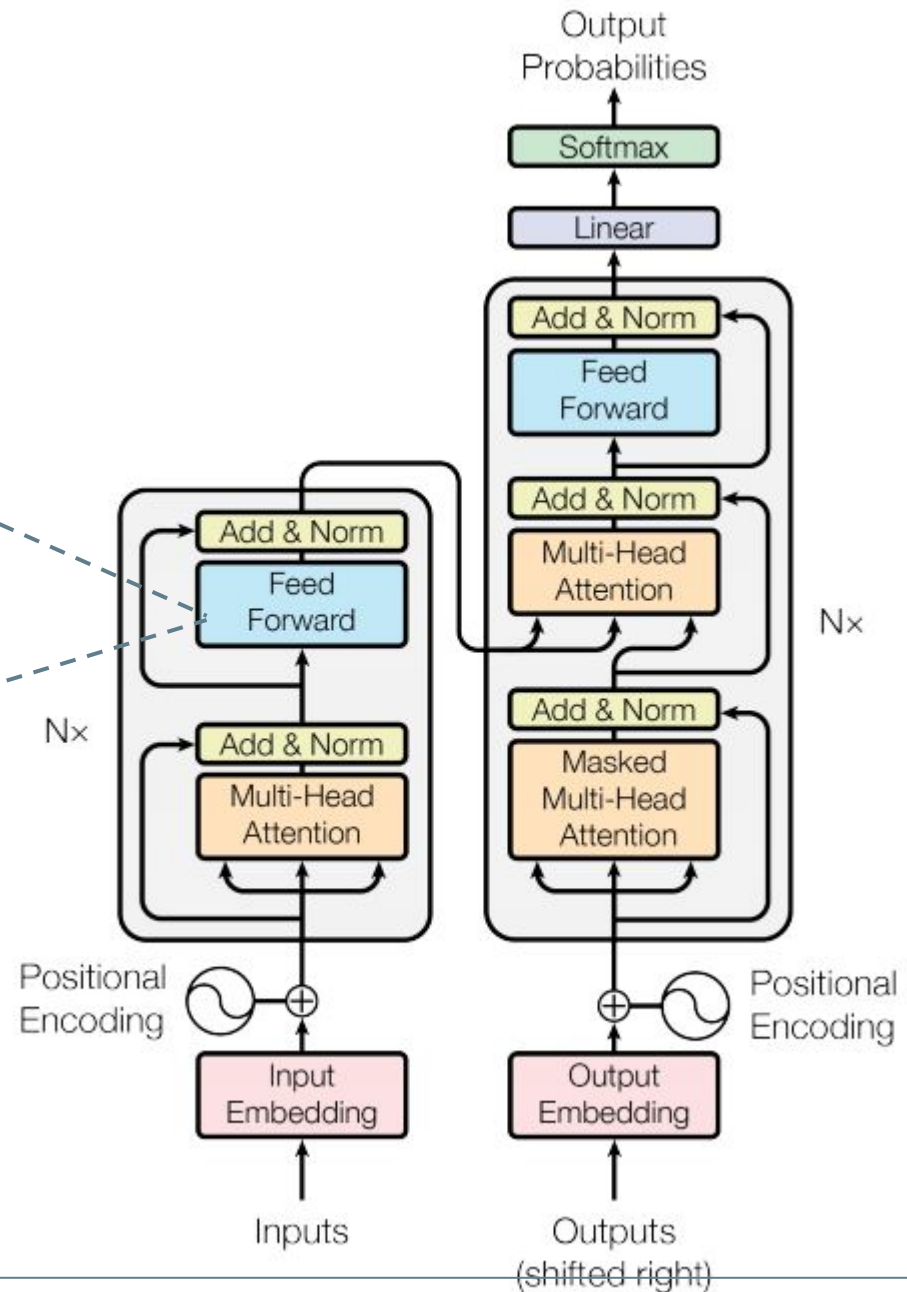


\*Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).

# Rede Neural

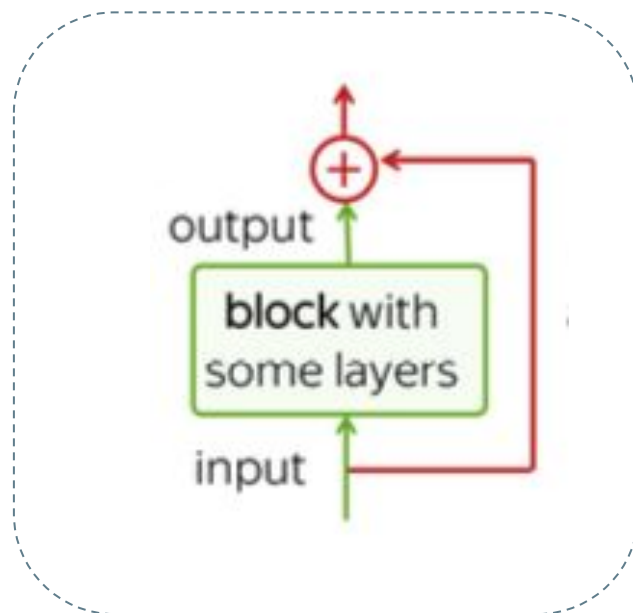


Camada de rede neural

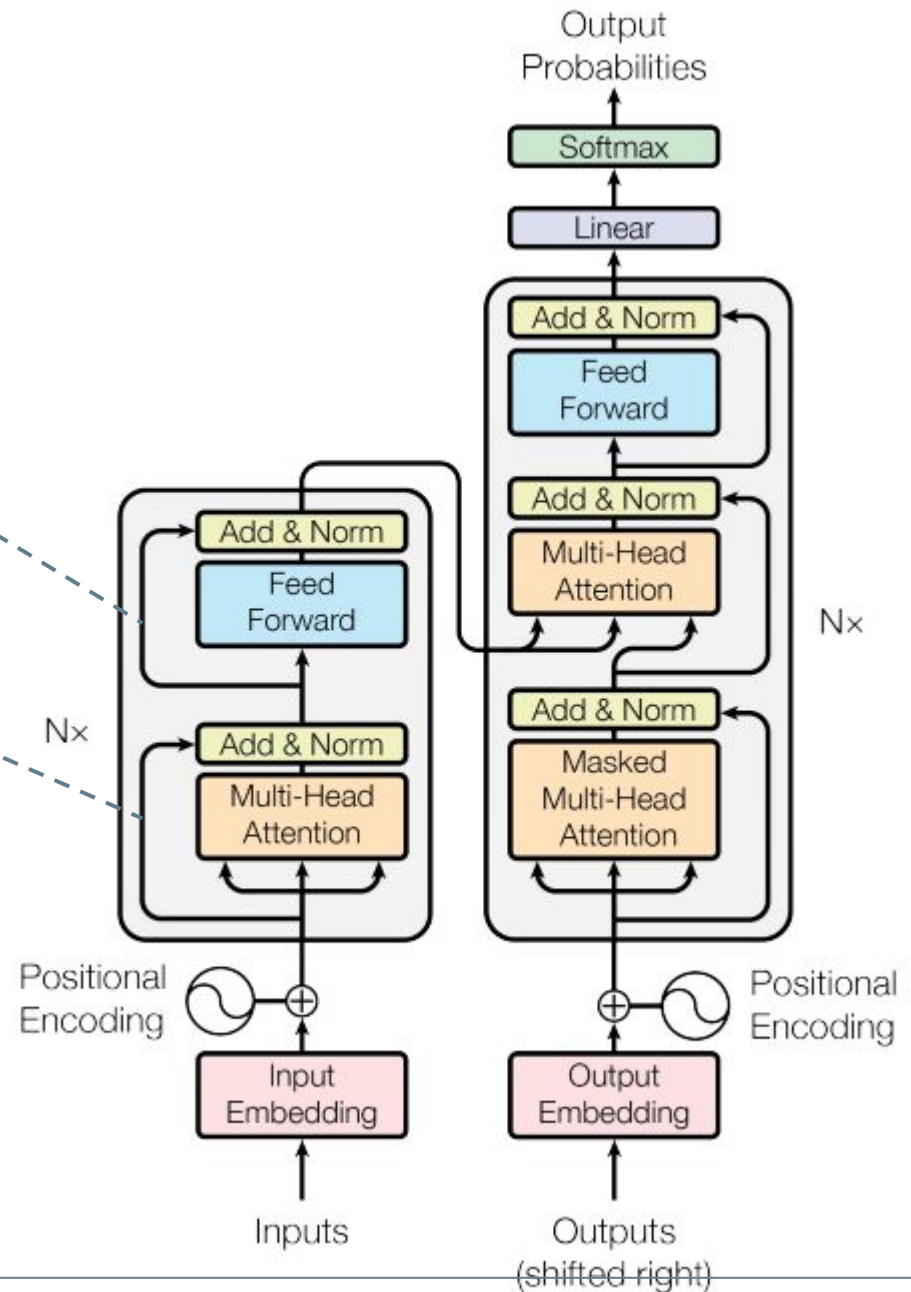


\*Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).

# Residual



Conexões residuais



\*Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).



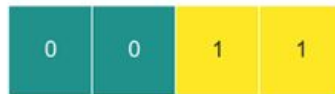
# Positional encoding

Posição normalizada

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

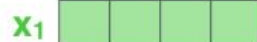
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

POSITIONAL  
ENCODING



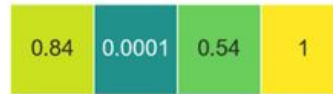
+

EMBEDDINGS

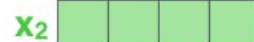


INPUT

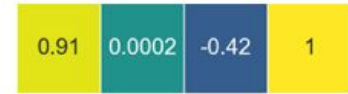
Je



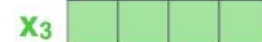
+



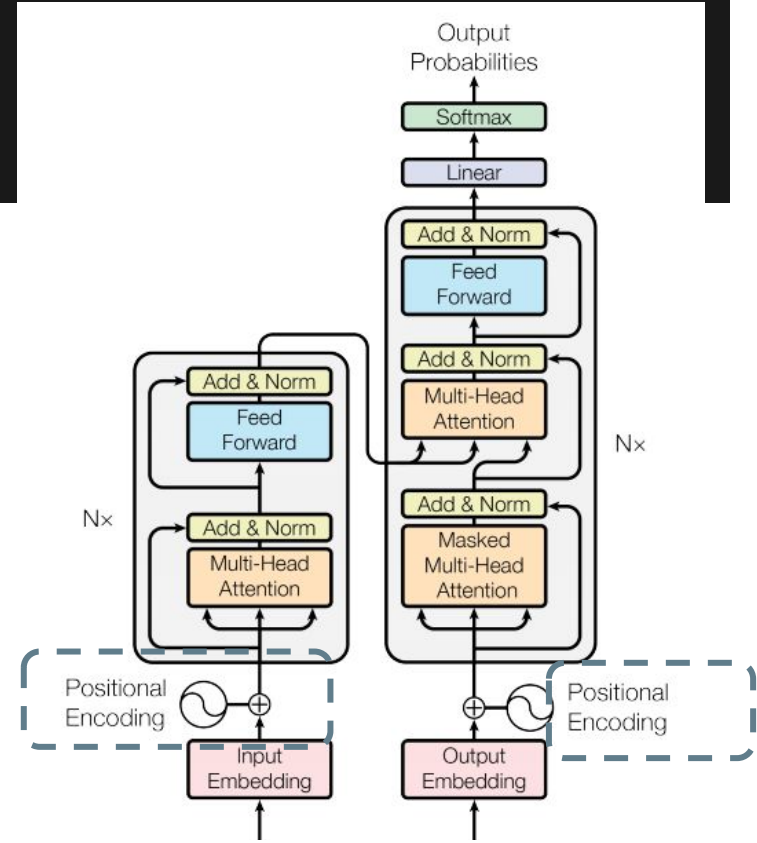
suis



+



étudiant

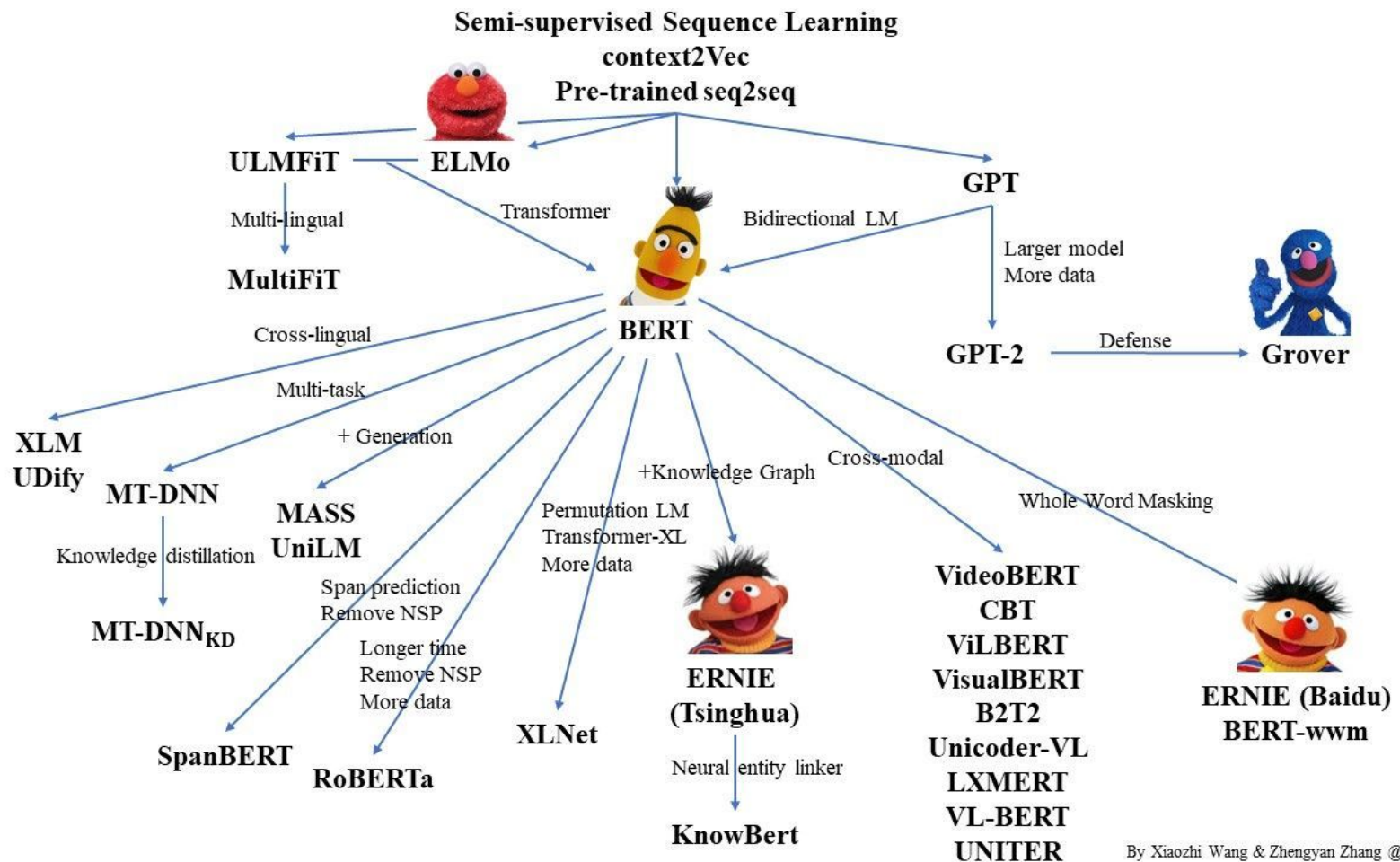


# Attention is all that you need?

- Inferência lenta
- Hardware pesado
- Não necessariamente melhores que RNN
- Não convergem com poucos dados

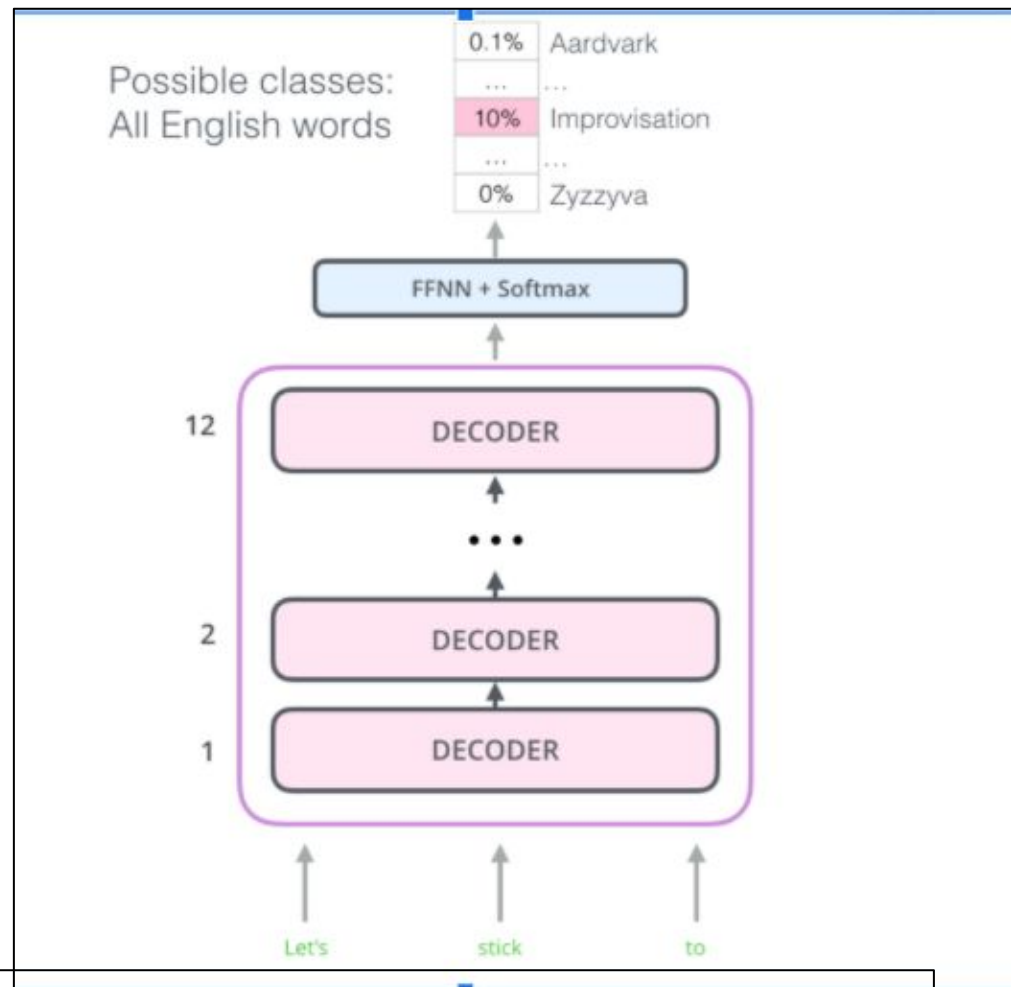
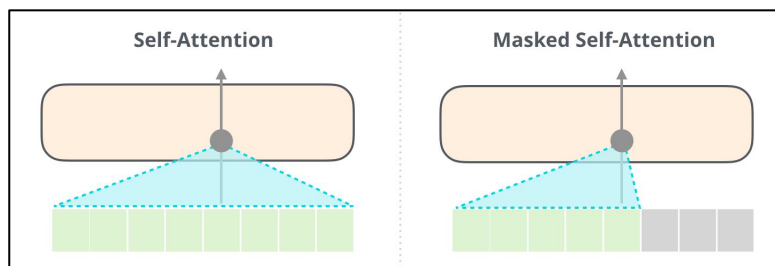


# BERT e seus amigos



# GPT

- Transformer-decoder
- Unidirectional

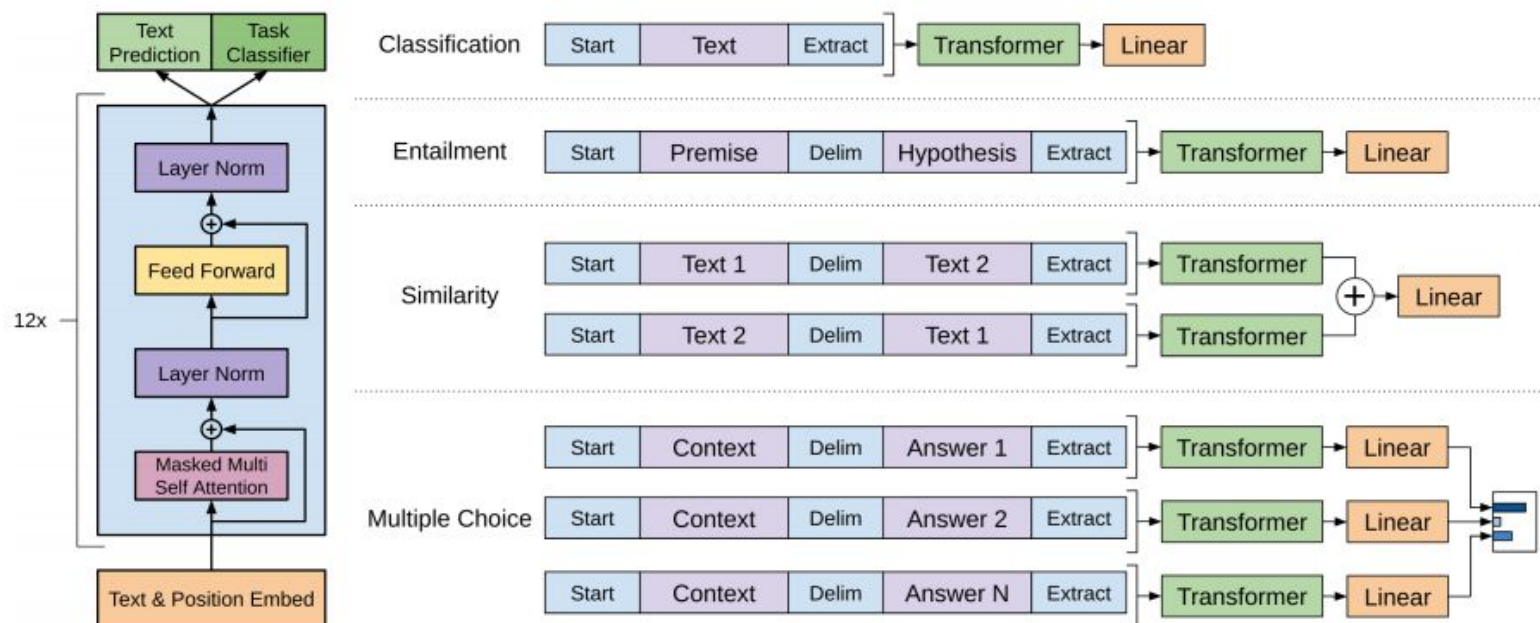


Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018) Improving Language Understanding by Generative Pre-Training.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019) Language Models are Unsupervised Multitask Learners.

# GPT

## Auto-regressive language model

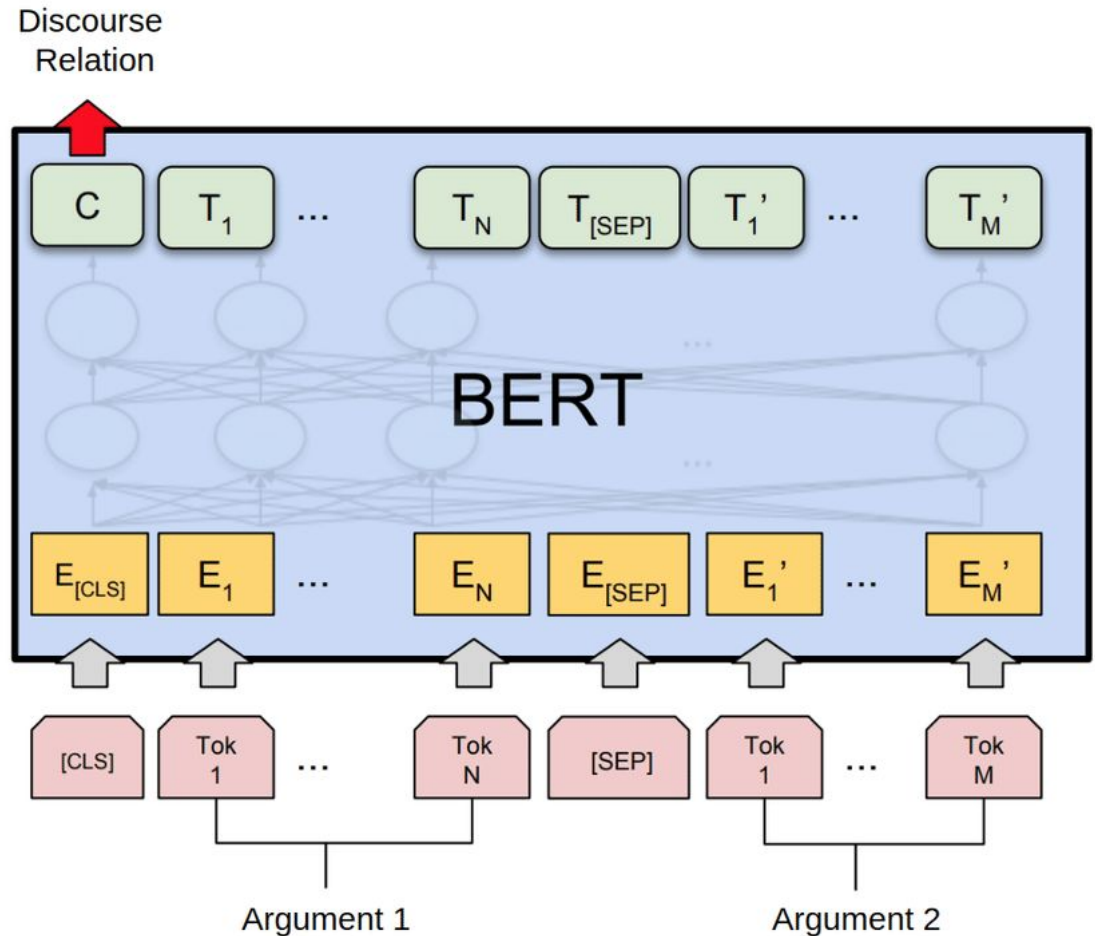
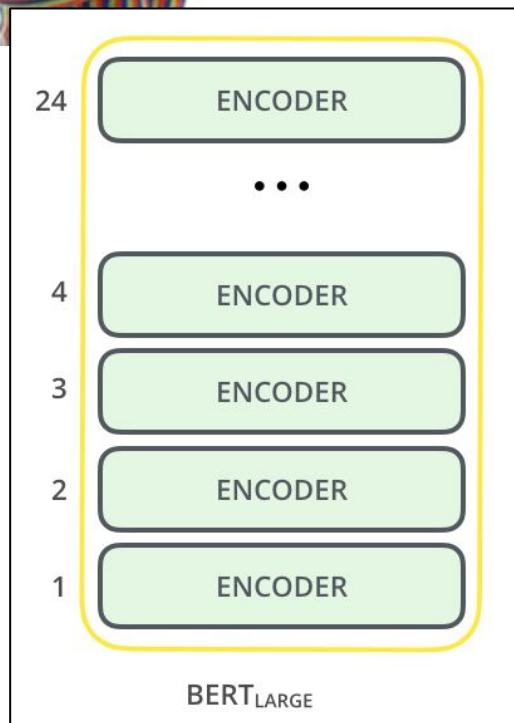


\* P. J. Liu, M. Saleh, E. Pot, B. Goodrich, R. Sepassi, L. Kaiser, and N. Shazeer. Generating wikipedia by summarizing long sequences. ICLR, 2018

\* Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.



# BERT



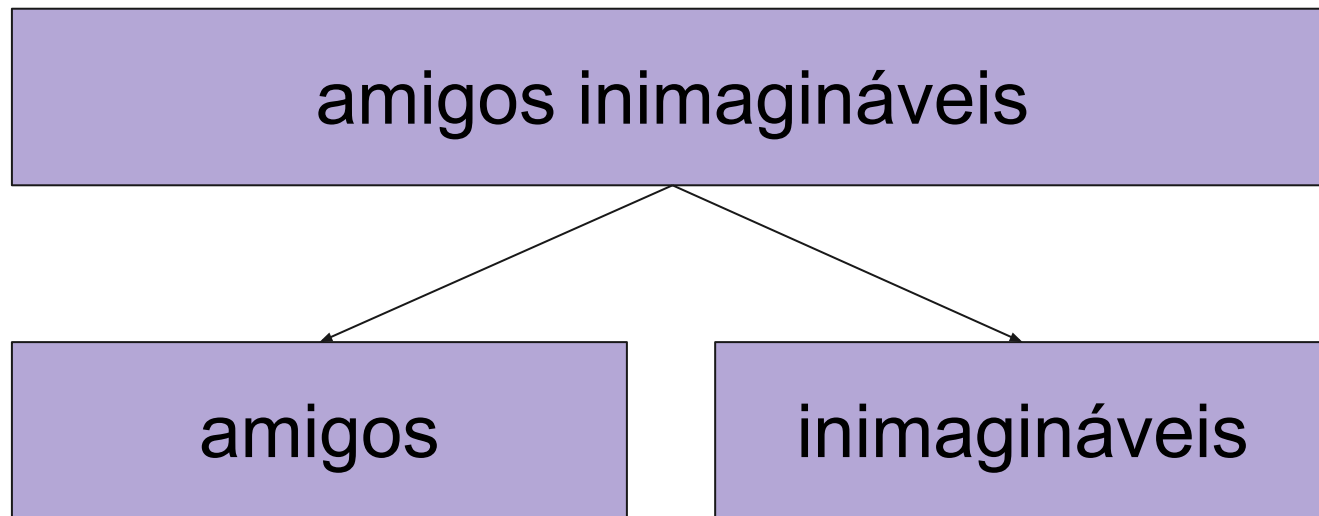
# Lembram das palavras desconhecidas?

- a) Assume que todas elas são um símbolo UNK
- b) Nunca teremos palavras desconhecidas, pois todos os modelos conseguem lidar com o vocabulário inteiro
- c) Ignora e segue em frente

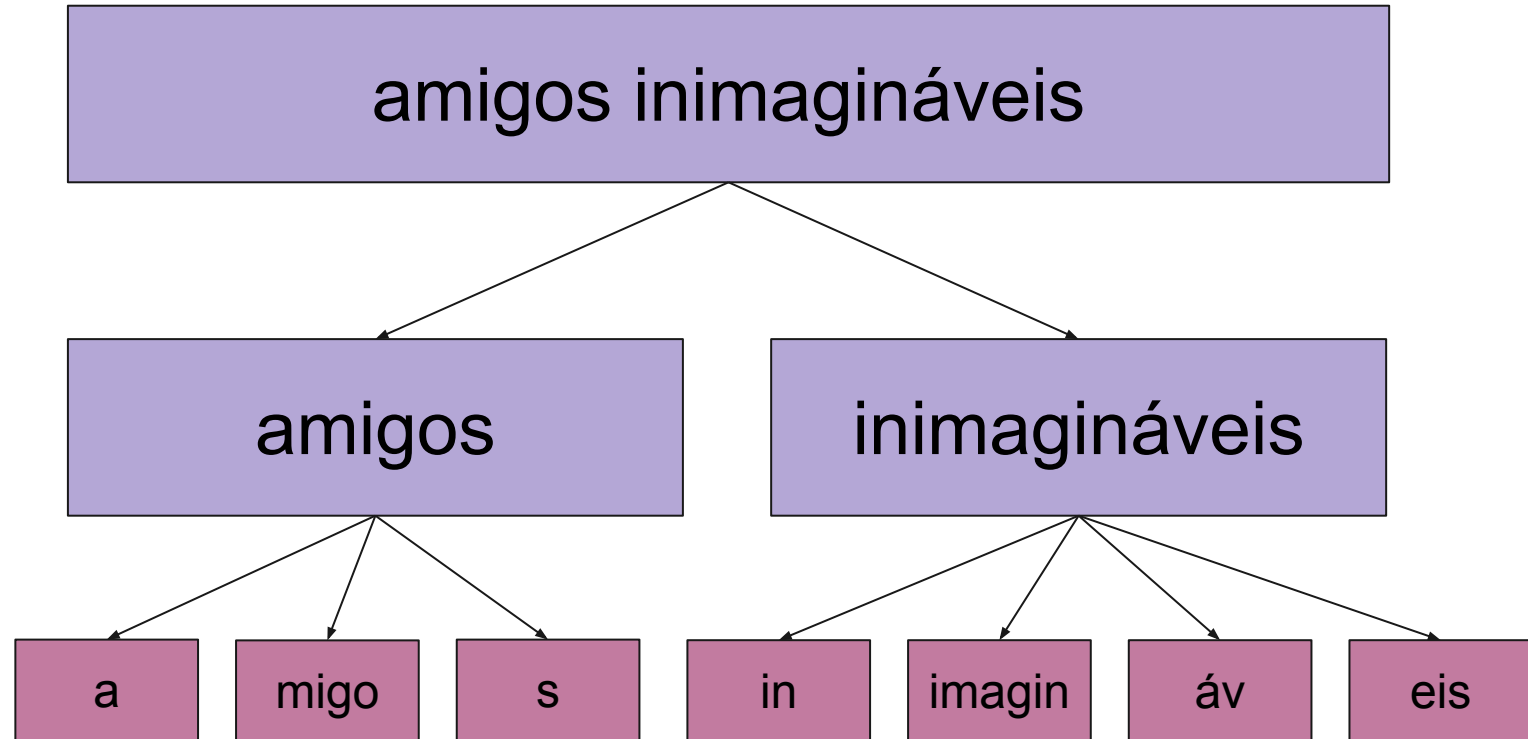
# O que modelos de linguagem neurais fazem para diminuir UNKs?

amigos inimagináveis

# O que modelos de linguagem neurais fazem para diminuir UNKs?

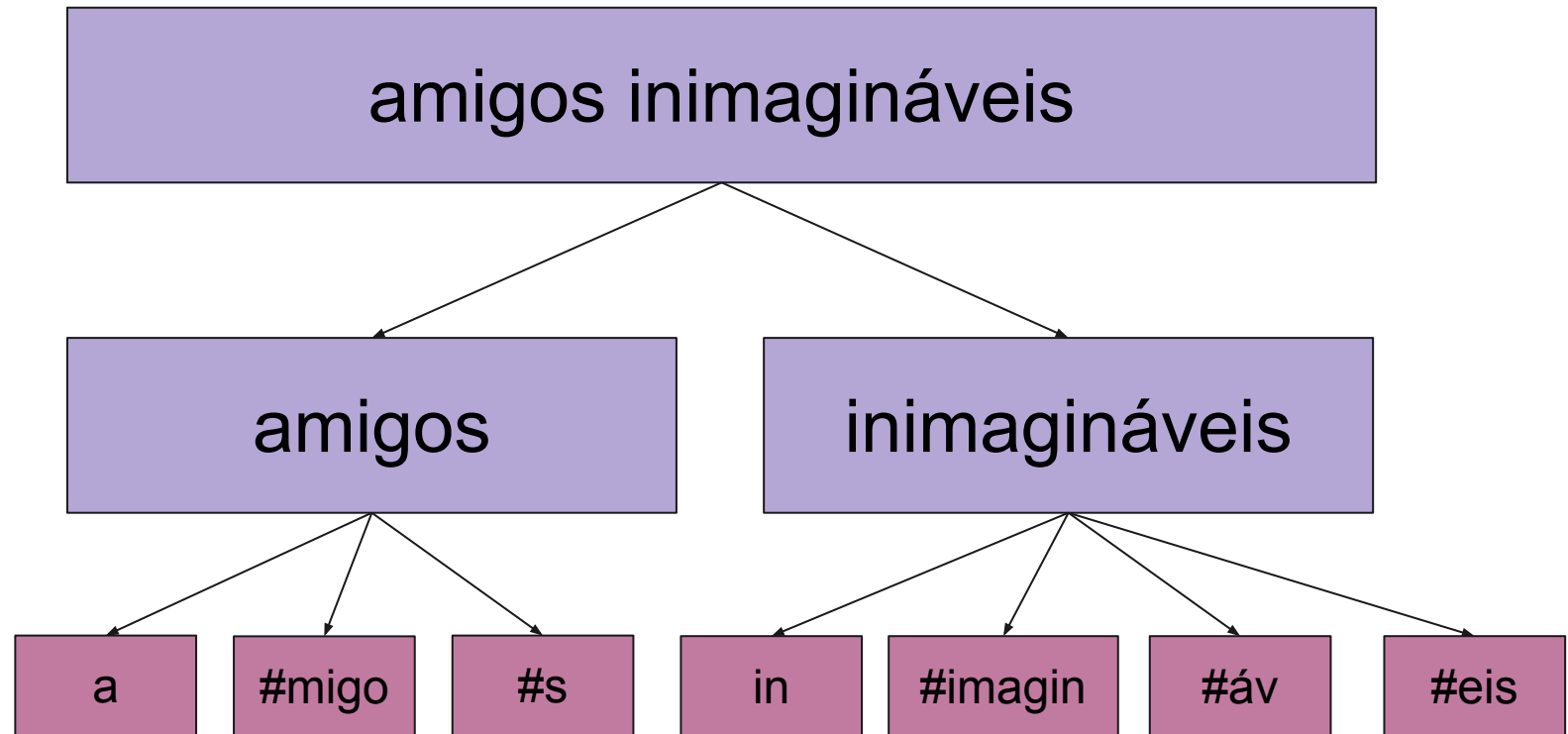


# O que modelos de linguagem neurais fazem para diminuir UNKs?

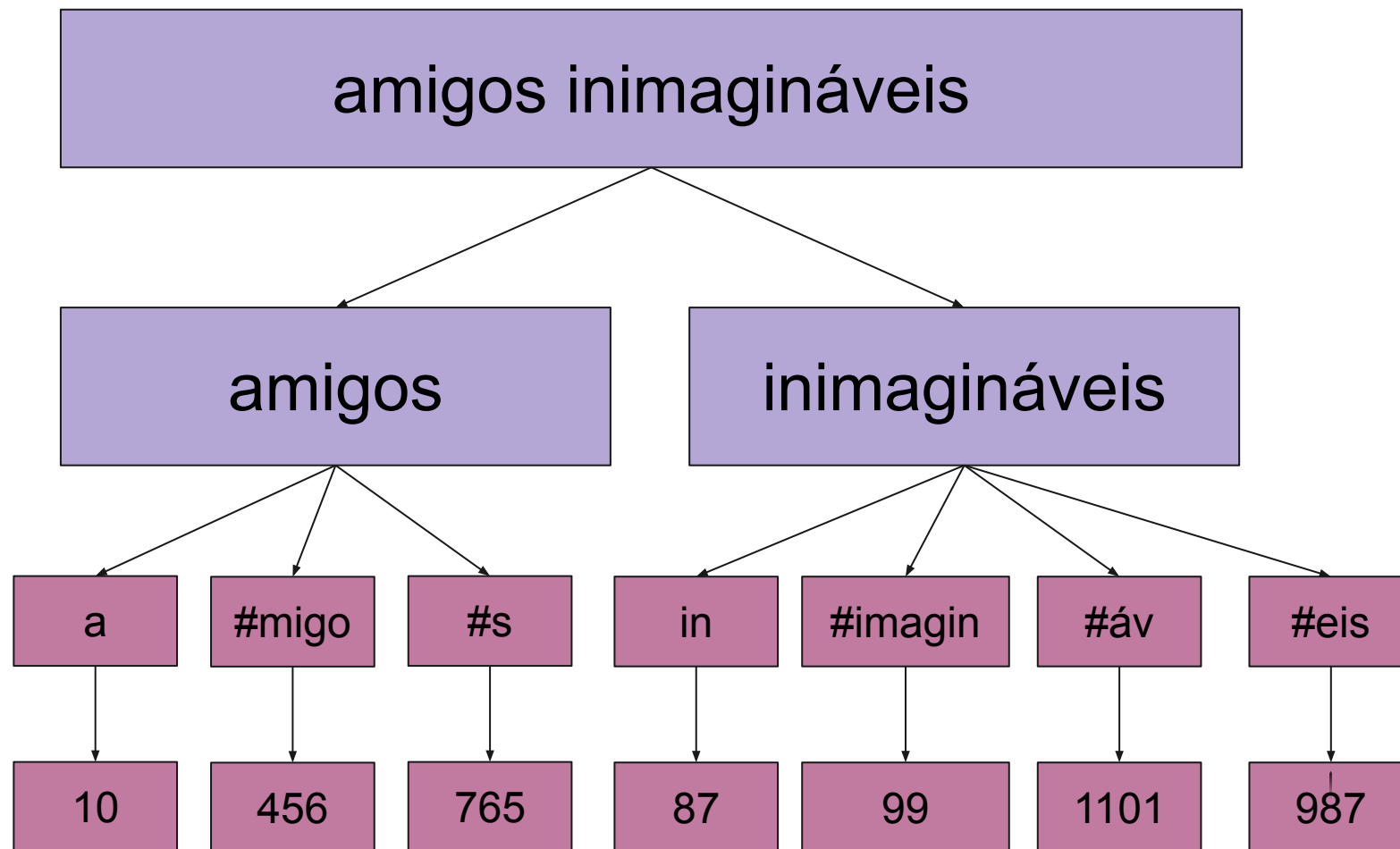




# O que modelos de linguagem neurais fazem para diminuir UNKs?



# O que modelos de linguagem neurais fazem para diminuir UNKs?



# Tokenizador

## Tiktokenizer

microsoft/phi-2

amigos inimagináveis

Token count  
9

amigos inimagináveis

321, 328, 418, 287, 320, 23183, 6557, 303, 271

☐ Show whitespace

# Tokenizador

Spaces | xu-song/tokenizer-arena | like 56 | Running

App | Files | Communi

## Tokenization Arena

Playground | Compression Leaderboard | Character Statistics

### Input Text

Examples ▾

E a vida o que é, diga lá meu irmão

### Tokenization

Tokenizer 1

Meta/llama3

Organization

Meta

Vocab Size

128256

Overlap Tokens

40094

VS

Tokenizer 2

openai/gpt-4o

Organization

OpenAI

Vocab Size

200019

Overlap Tokens

40094

Tokens: 13

E a vida o que é , dig a lá meu irm ão

TokenID	Token	Text	UTF8 Bytes
36	E	E	b' E '
264	Êa	a	b' \xc4\xa0a '
25994	Êvida	vida	b' \xc4\xa0vida '
297	Êo	o	b' \xc4\xa0o '
1744	Êque	que	b' \xc4\xa0que '
4046	ÊÃ	é	b' \xc4\xa0xc3\x83\xc2\xa9 '

Tokens: 11

E a vida o que é , diga lá meu irmão

TokenID	Token	Text	UTF8 Bytes
36	E	E	b' E '
261	a	a	b' a '
8417	vida	vida	b' vida '
293	o	o	b' o '
661	que	que	b' que '
1212	é	é	b' \xc3\xa9 '

# Notebook embeddings