

---

---

# PROPOR 2020 - Évora, Portugal

— Modelos de Linguagem —

Joaquim Santos, Renata Vieira

---

---

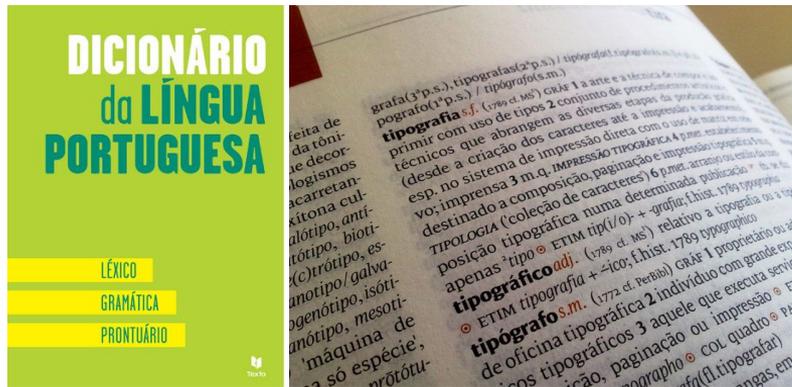
# Roteiro

1. Introdução
2. Fundamentos matemáticos
3. Histórico
4. Modelos WE iniciais
5. Modelos dinâmicos, sensíveis ao contexto
6. Modelos para a língua portuguesa
7. Modelos específicos de domínio
8. Geração de modelos
9. Aplicações
10. Avaliação de modelos
11. Limitações

# Introdução

# LÍNGUA E SIGNIFICADO

## O que as palavras significam?



# Exemplo

**Mangueira**

Árvore da família das anacardiáceas.

Tubo flexível, de tecido, borracha ou plástico, destinado a conduzir líquidos ou gases; manga, borracha.

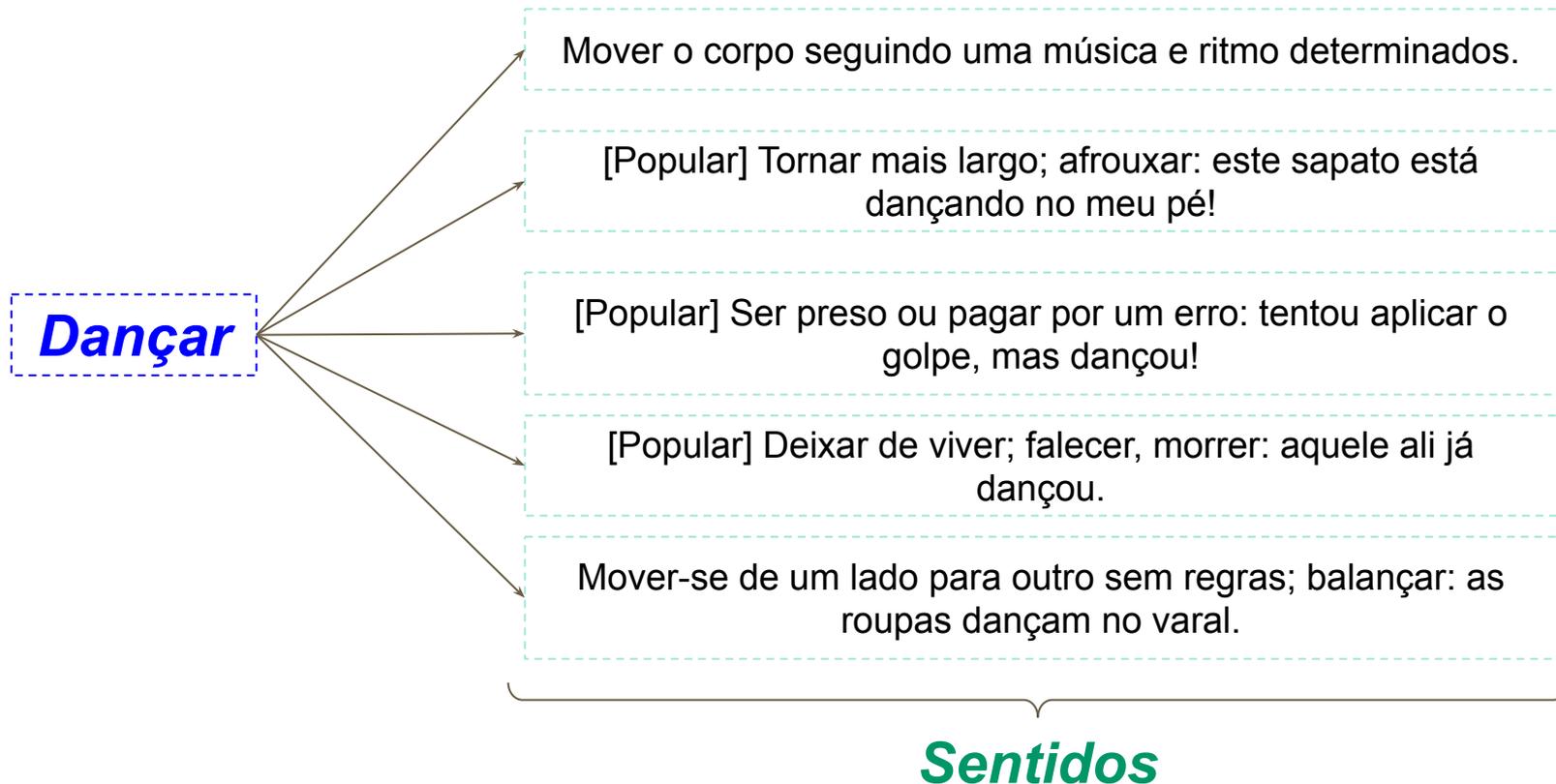
Tradicional escola de samba brasileira da cidade do Rio de Janeiro

Bairro da cidade de Salvador-BA, situa-se na cidade baixa entre massaranduba e bonfim.

Curral grande. Local de trabalhos com gado manso e bravo, localizado próximo à casa principal e feito de pedra, pau a pique, varas etc.

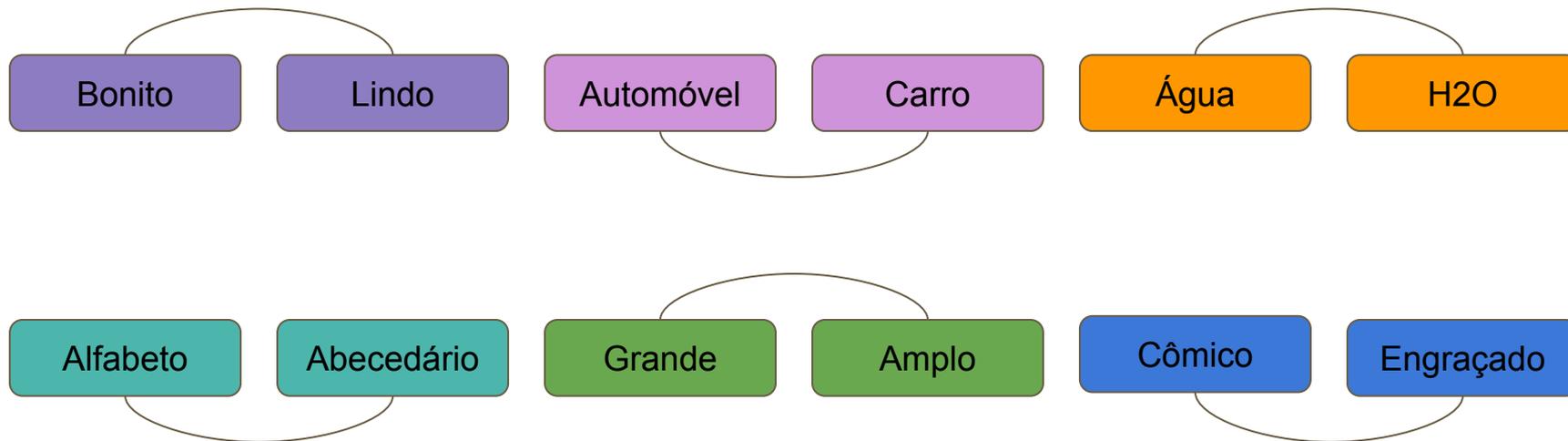
**Sentidos**

# Exemplo



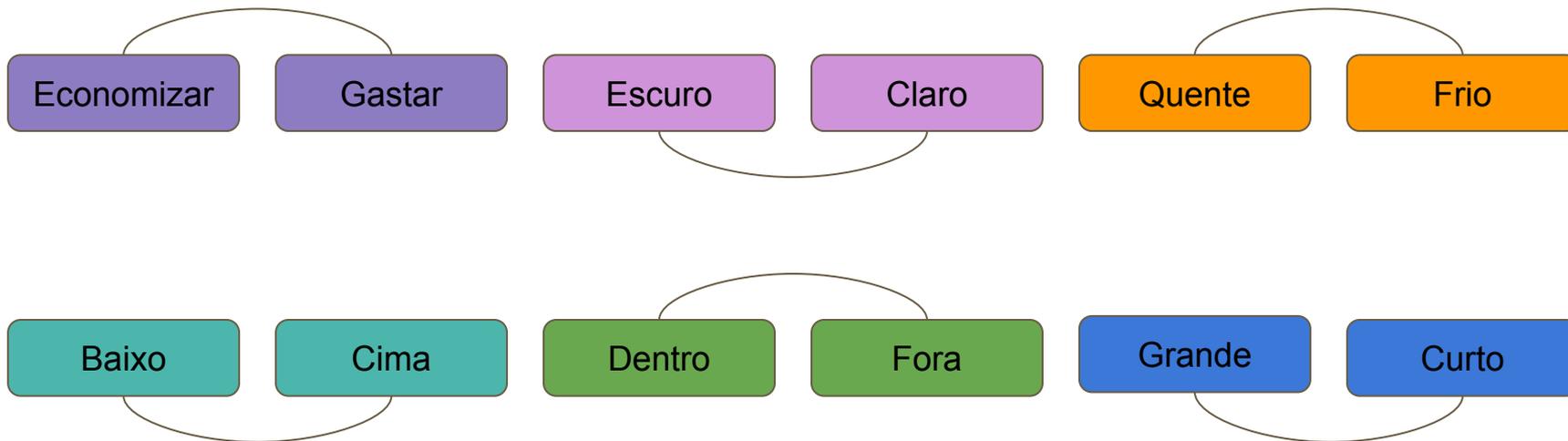
Existem relações entre os *Sentidos* das palavras.

## Sinonímia



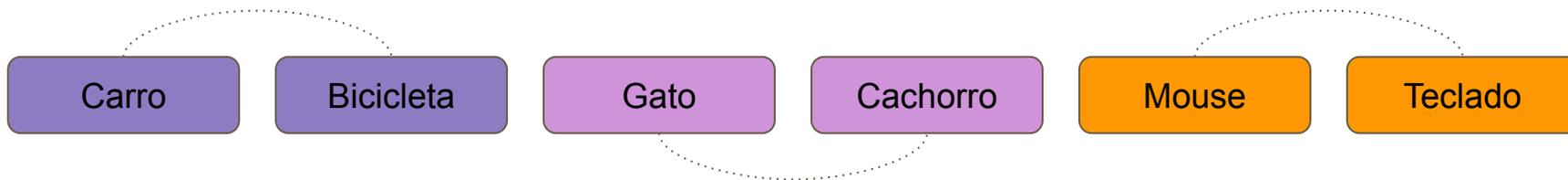
Existem relações entre os *Sentidos* das palavras.

## Antonímia



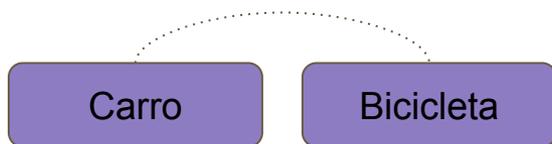
Existem relações entre os *Sentidos* das palavras.

**Similaridade** Não são sinônimos, mas compartilham algum elemento de significado.

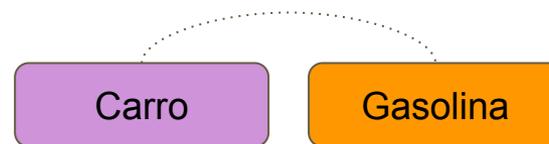


Existem relações entre os *Sentidos* das palavras.

## Palavras Relacionadas



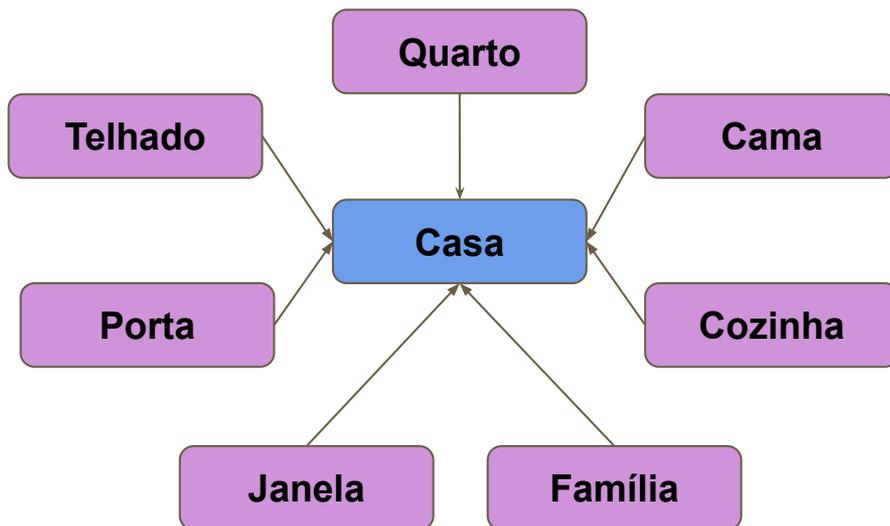
**Similares**



**Relacionadas, mas não  
similares**

Existem relações entre os *Sentidos* das palavras.

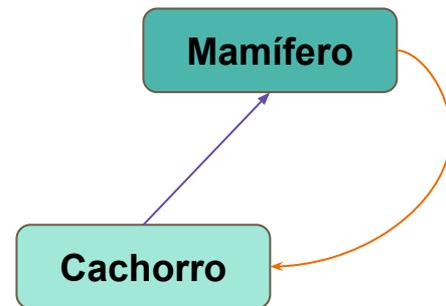
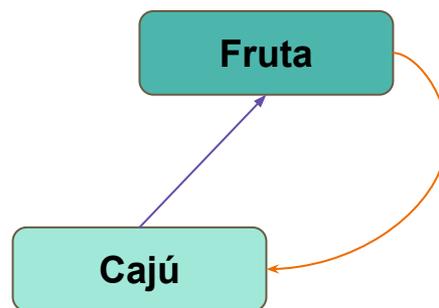
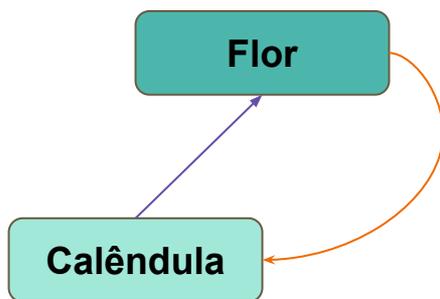
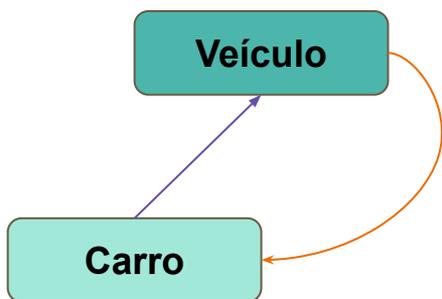
## Campos Semânticos



Existem relações entre os *Sentidos* das palavras.

**Hipônimo**

**Hiperônimo**



Mas como definir um **conceito**?

Os **conceitos** das palavras têm associações complexas.

**Sinonímia**

**Antonímia**

**Similaridade**

**Campos  
Semânticos**

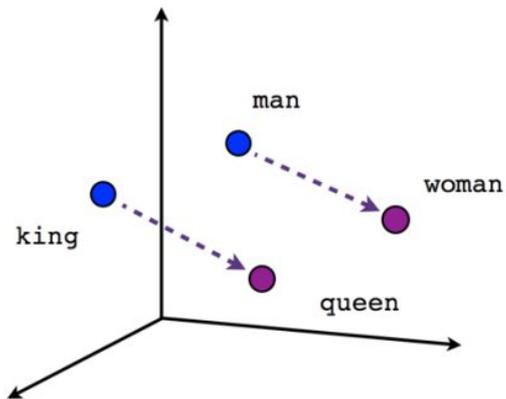
**Hiperônimo/  
Hipônimo**

Os **conceitos** dependem do **contexto**.

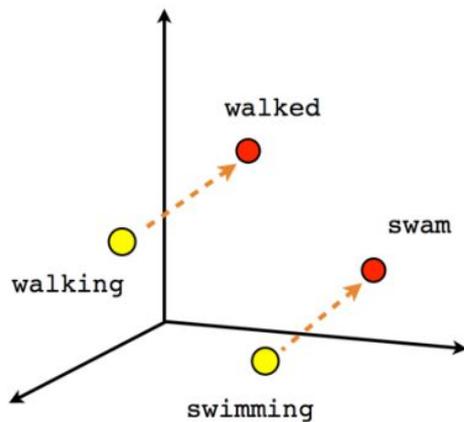
- Seria possível representar uma linguagem natural por um modelo genérico (computacional) que capturasse essas características através da análise de muitos contextos?

- Seria possível representar uma linguagem natural por um modelo genérico (computacional) que capturasse essas características através da análise de muitos contextos?
- SIM, podemos representar uma linguagem natural por um modelo de natureza matemática!

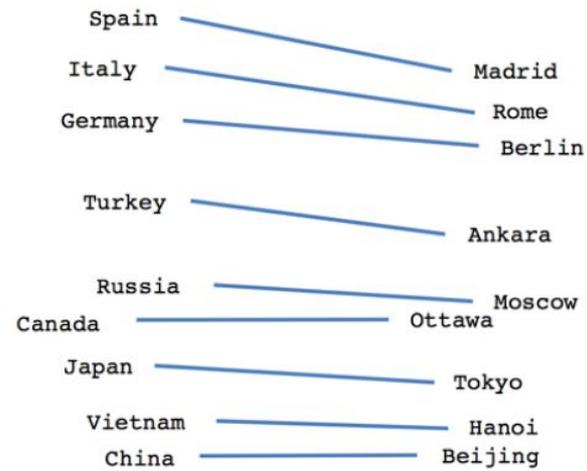




Male-Female



Verb tense



Country-Capital

# Fundamentos matemáticos

Van Gogh, o ???

O gato sentou no ???

O astrônomo olha para o ???

Van Gogh, o ???

O gato sentou no ???

O astrônomo olha para o ???

**PROBABILIDADE!**

Van Gogh, o ???

O gato sentou no ???

O astrônomo olha para o ???

$P(?|o, astrônomo, olha, para, o)$

$P(w_t|w_1, \dots, w_n), n \in \mathbb{N}$

$$\begin{aligned} P(\text{céu} | \text{o}, \text{astrônomo}, \text{olha}, \text{para}, \text{o}) &= \\ &= \frac{\#(\text{o astrônomo olha para o céu})}{\#(\text{o astrônomo olha para o})} \end{aligned}$$

$$\begin{aligned} P(X_1 \dots X_n) &= P(X_1)P(X_2|X_1)P(X_3|X_1^2) \dots P(X_n|X_1^{n-1}) \\ &= \prod_{k=1}^n P(X_k|X_1^{k-1}) \end{aligned}$$

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2|w_1)P(w_3|w_1^2) \dots P(w_n|w_1^{n-1}) \\ &= \prod_{k=1}^n P(w_k|w_1^{k-1}) \end{aligned}$$

$$\begin{aligned} P(X_1 \dots X_n) &= P(X_1)P(X_2|X_1)P(X_3|X_1^2) \dots P(X_n|X_1^{n-1}) \\ &= \prod_{k=1}^n P(X_k|X_1^{k-1}) \end{aligned}$$

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2|w_1)P(w_3|w_1^2) \dots P(w_n|w_1^{n-1}) \\ &= \prod_{k=1}^n P(w_k|w_1^{k-1}) \end{aligned}$$

Mas e se  $w_1^n$  nunca acontece?

$$P(\text{céu} | o, \text{astrônomo}, \text{olha}, \text{para}, o) \approx P(\text{céu} | o)$$

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-1})$$

$$P(\text{céu}|o, \text{astrônomo}, \text{olha}, \text{para}, o) \approx P(\text{céu}|o)$$

$$P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-1})$$

Uma sequência pode ser então estimada por:

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2|w_1)P(w_3|w_1^2) \dots P(w_n|w_1^{n-1}) \\ &= \prod_{k=1}^n P(w_k|w_1^{k-1}) \\ P(w_1^n) &\approx \prod_{k=1}^n P(w_k|w_{k-1}) \end{aligned}$$

$$P(\text{céu}|o, \text{astrônomo}, \text{olha}, \text{para}, o) \approx P(\text{céu}|o)$$

$$P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-1})$$

Uma sequência pode ser então estimada por:

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2|w_1)P(w_3|w_1^2) \dots P(w_n|w_1^{n-1}) \\ &= \prod_{k=1}^n P(w_k|w_1^{k-1}) \\ P(w_1^n) &\approx \prod_{k=1}^n P(w_k|w_{k-1}) \end{aligned}$$

Olhando para o **Histórico**, onde estamos?

# Histórico

2.500 sentenças  
14 milhões de tokens

### Feedforward Neural Network



**Slava M. Katz**

**Stanley and Joshua**

**Xu and Rudnick**

**Bengio et al.**

“Estimation of Probabilities from Sparse Data for the Language Model Component of a Speech Recognizer”

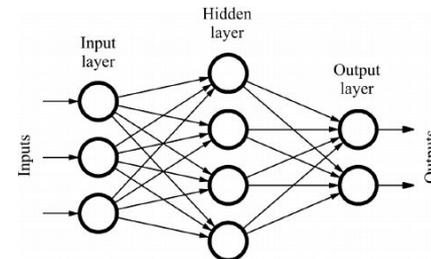
“An Empirical Study of Smoothing Techniques for Language Modeling”

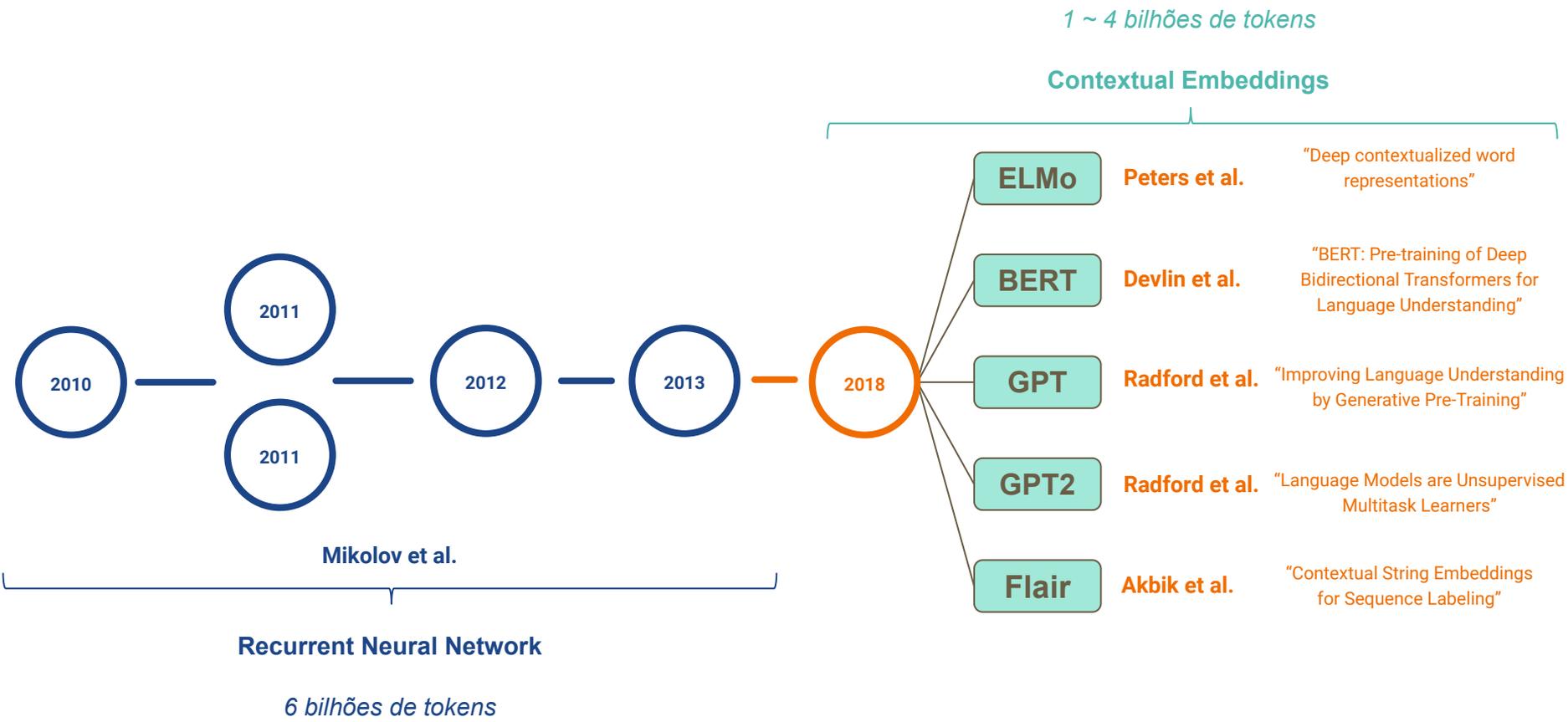
“Can Artificial Neural Networks Learn Language Models?”

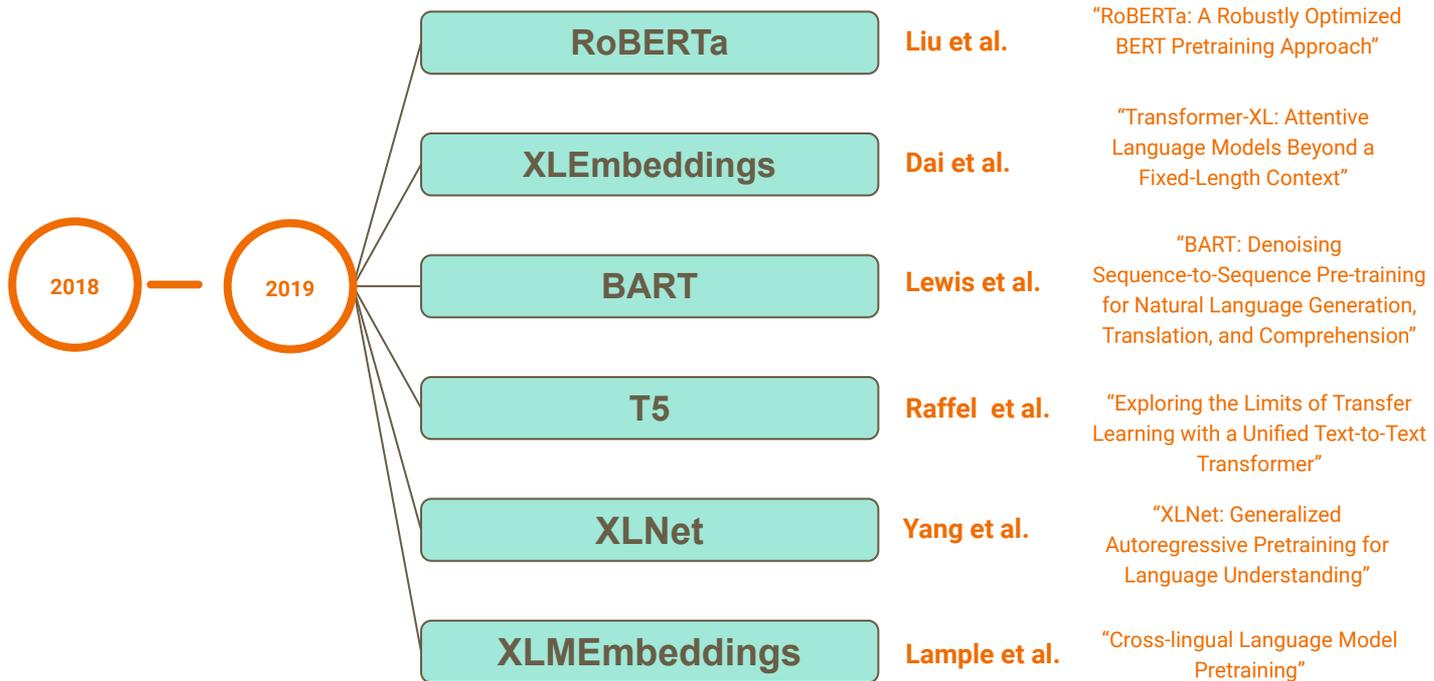
“A Neural Probabilistic Language Model”

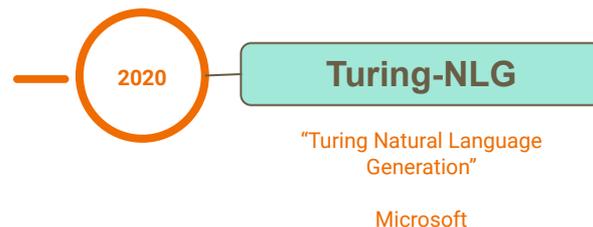
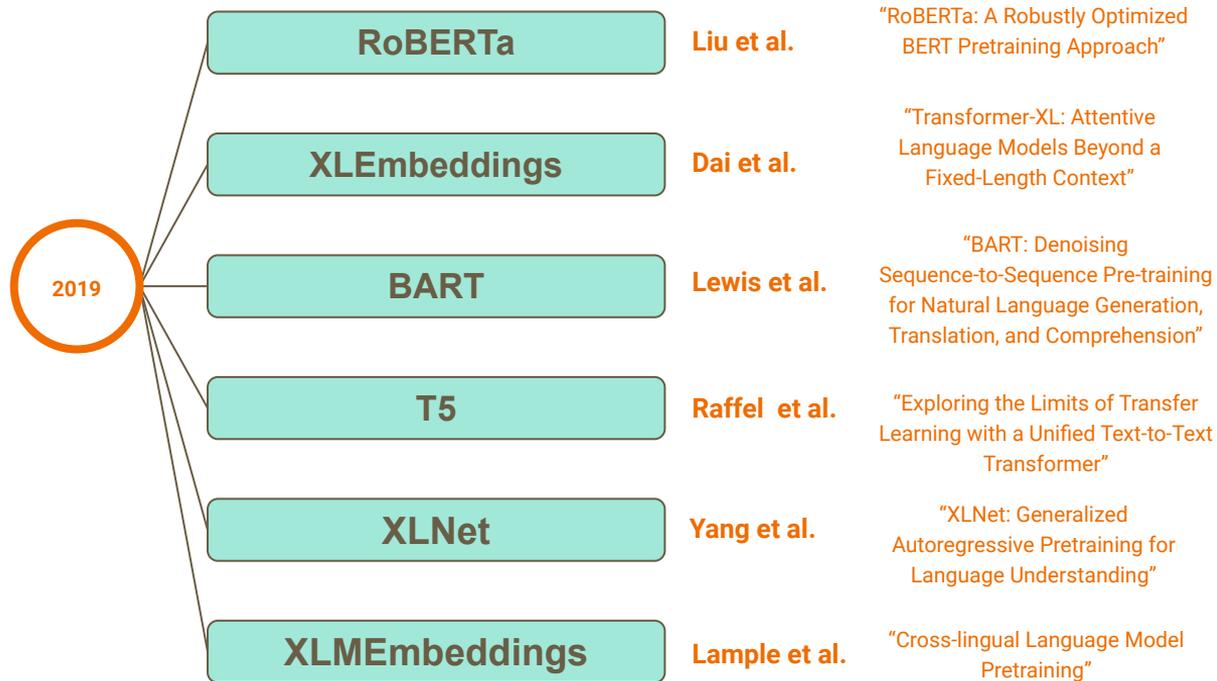
### N-gram Language Models

*2,3-gram  
750.000 mil ~ 250 milhões  
tokens*



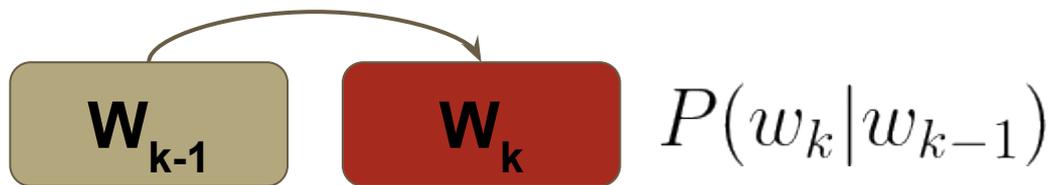




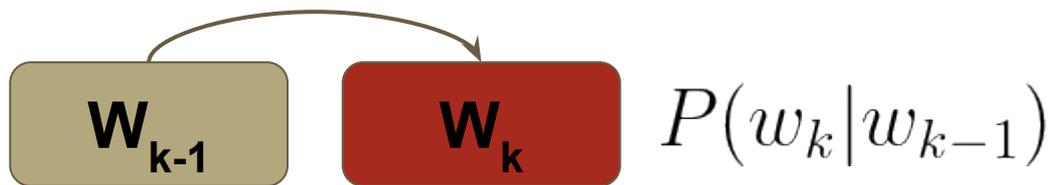


# Modelos WE iniciais

Nós vimos que uma sequência pode ser então estimada por:  $P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-1})$

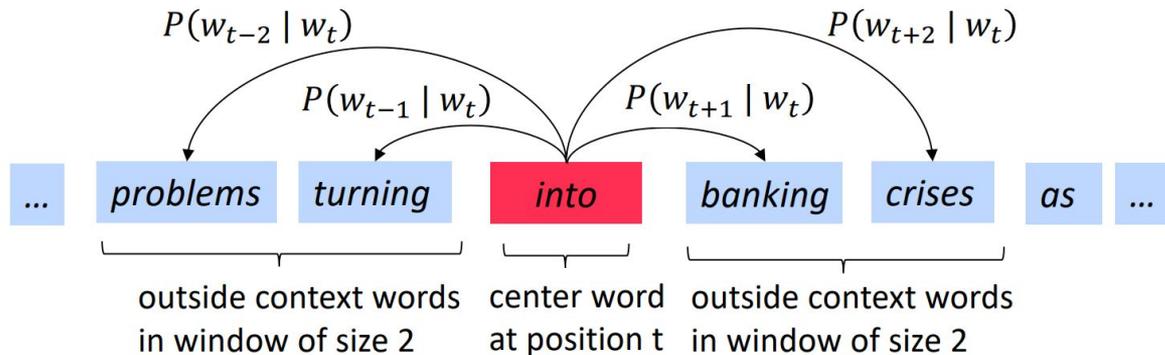


Nós vimos que uma sequência pode ser então estimada por:  $P(w_1^n) \approx \prod_{k=1}^n P(w_k|w_{k-1})$



Nós vimos que uma sequência pode ser então estimada por:  $P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-1})$

$$P(w_t | w_1, \dots, w_n), n \in \mathbb{N} \quad \mathbf{e} \quad P(w_1, \dots, w_n | w_t), n \in \mathbb{N}$$



De outro modo: Para cada posição  $t$ , prevemos palavras de contexto.

$$L(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m} P(w_{t+j} | w_t; \theta)$$

Temos que maximizar essa probabilidade.

Negative Log Likelihood (NLL)  $J(\theta) = -\frac{1}{T} \log L(\theta)$

$$L(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m} P(w_{t+j} | w_t; \theta)$$

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m} \log P(w_{t+j} | w_t; \theta)$$

$\therefore$

$$\frac{\partial}{\partial v_c} J(\theta)$$

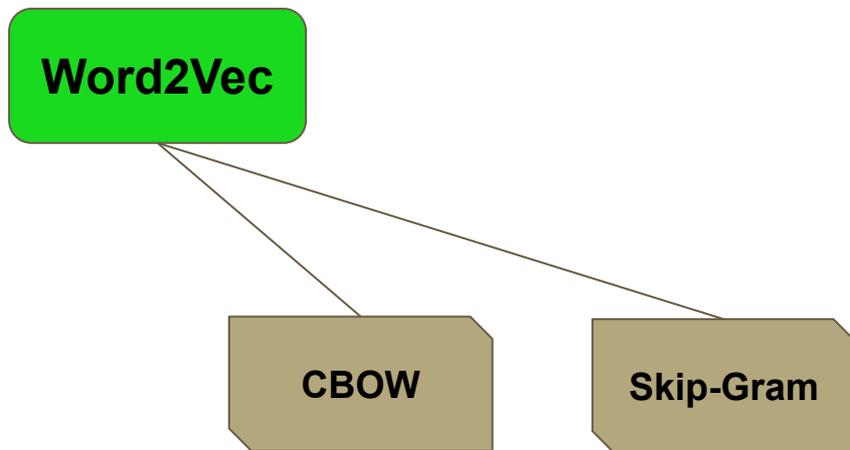
**Minimizar a função objetivo**  
 $\Leftrightarrow$   
**Maximizar a precisão preditiva**



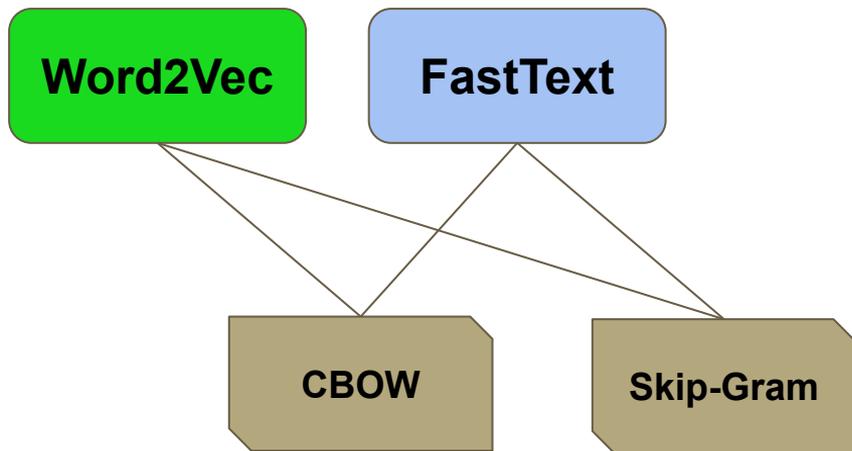
- Existem várias arquiteturas de Rede Neural que geram Word Embeddings:

**Word2Vec**

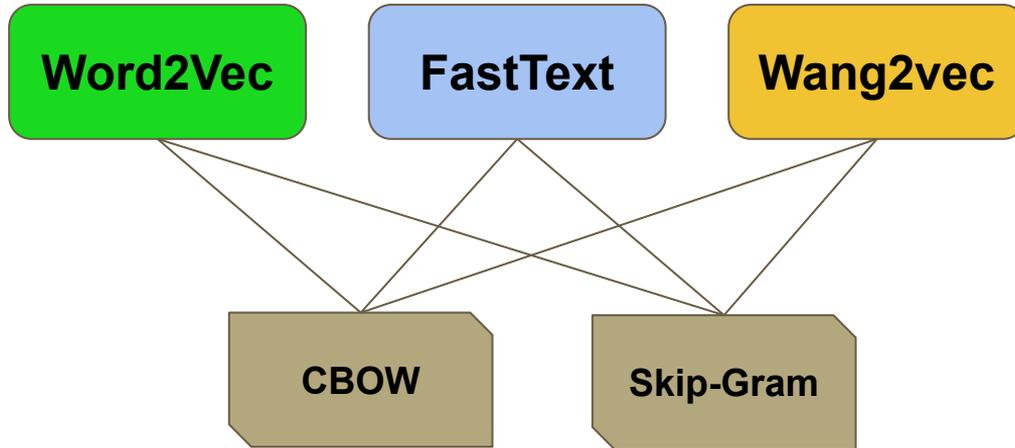
- Existem várias arquiteturas de Rede Neural que geram Embeddings:



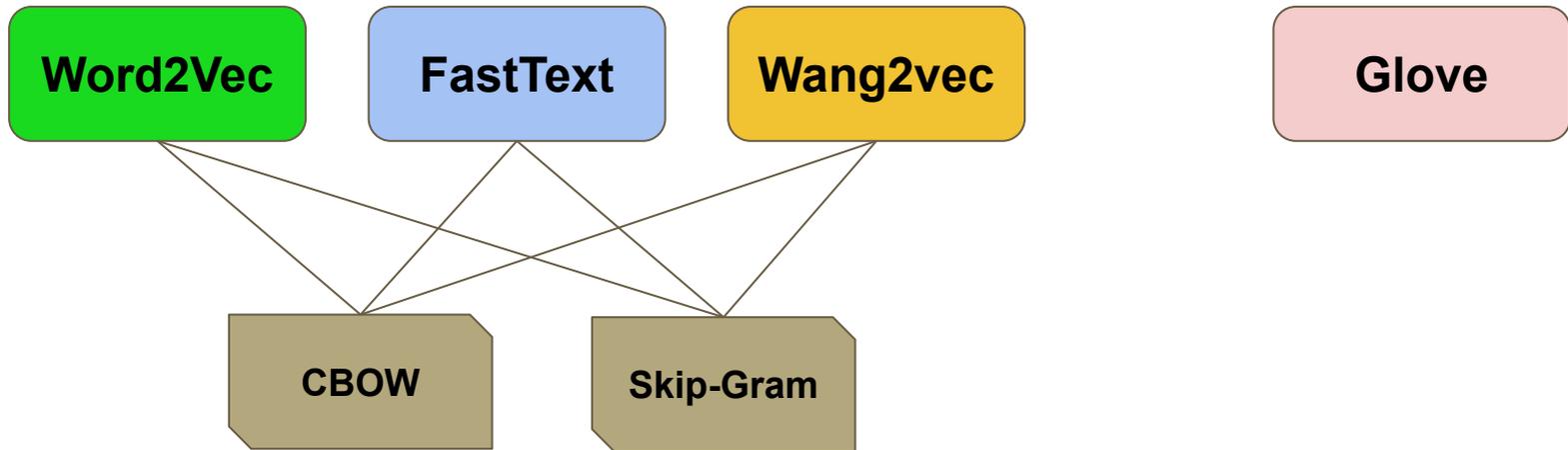
- Existem várias arquiteturas de Rede Neural que geram Embeddings:



- Existem várias arquiteturas de Rede Neural que geram Embeddings:

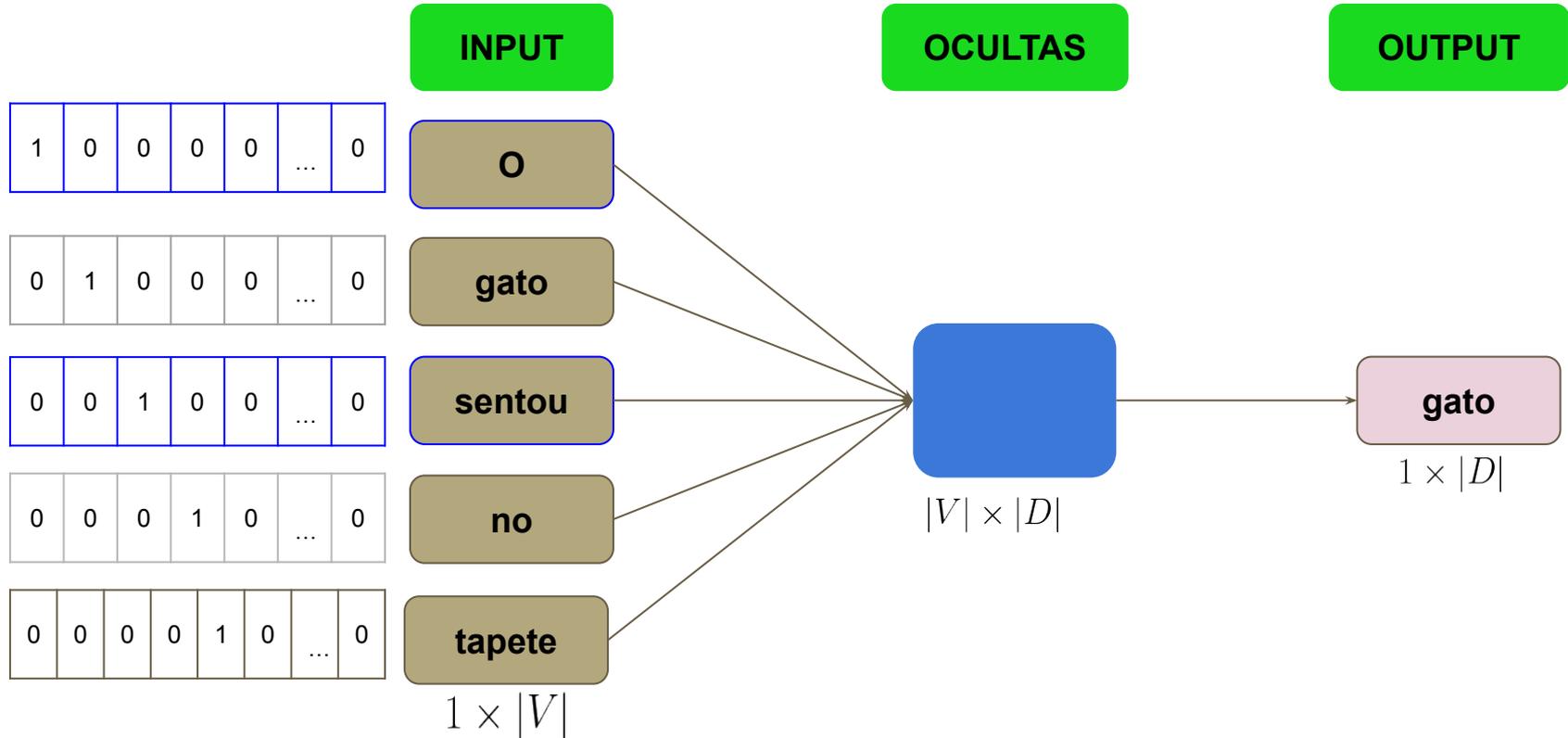


- Existem várias arquiteturas de Rede Neural que geram Embeddings:



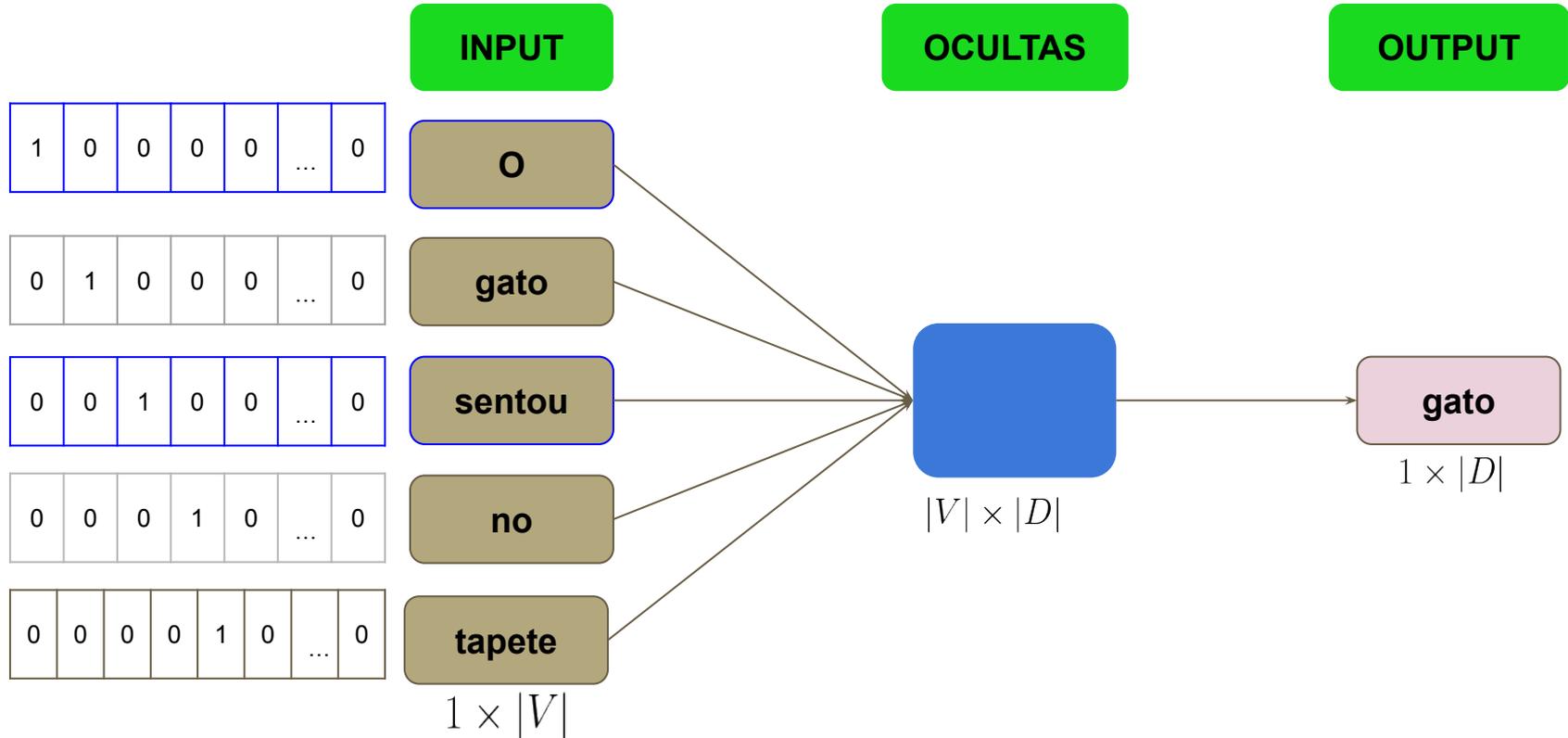
CBOW - Continuous Bag-Of-Words

O gato sentou no tapete



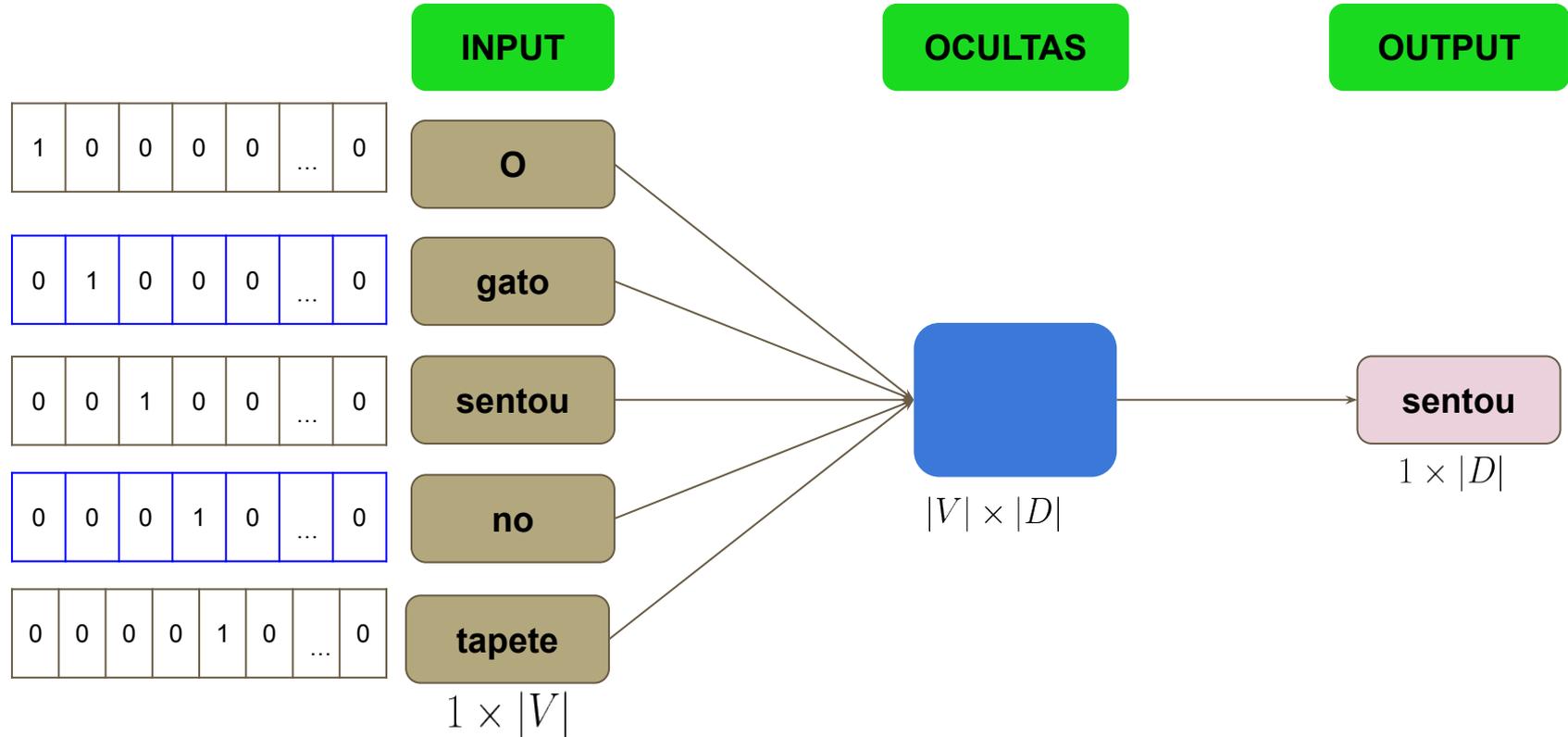
CBOW - Continuous Bag-Of-Words

O gato sentou no tapete



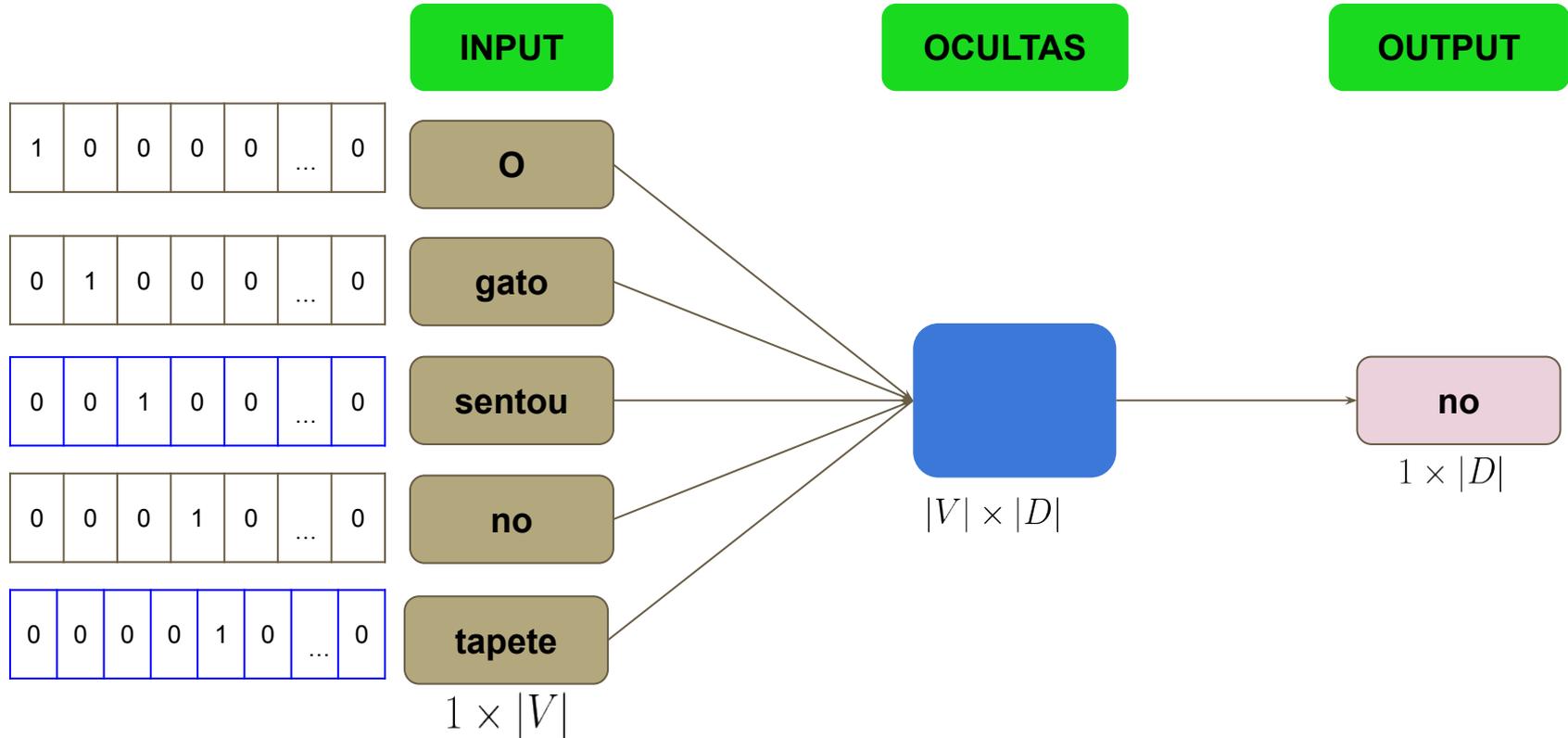
CBOW - Continuous Bag-Of-Words

O gato sentou no tapete



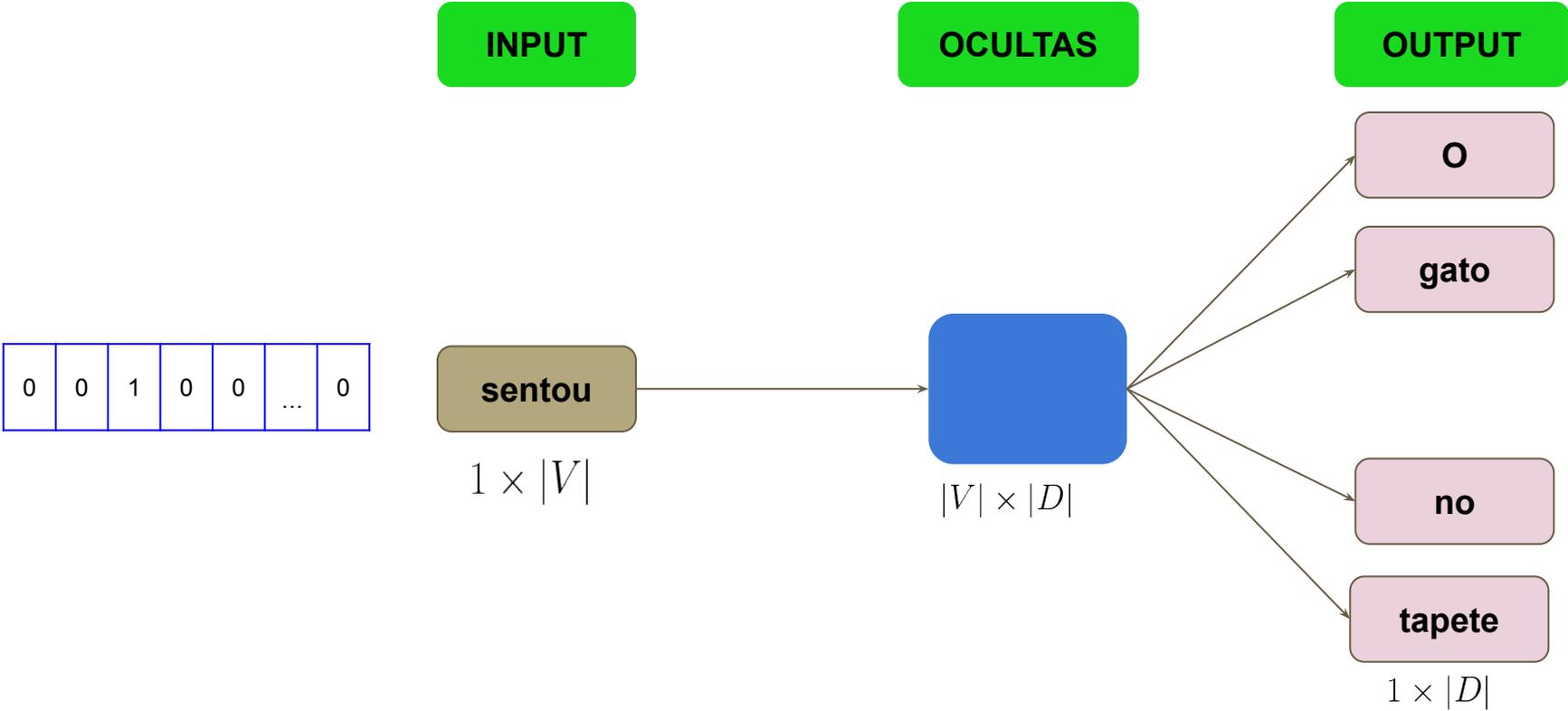
CBOW - Continuous Bag-Of-Words

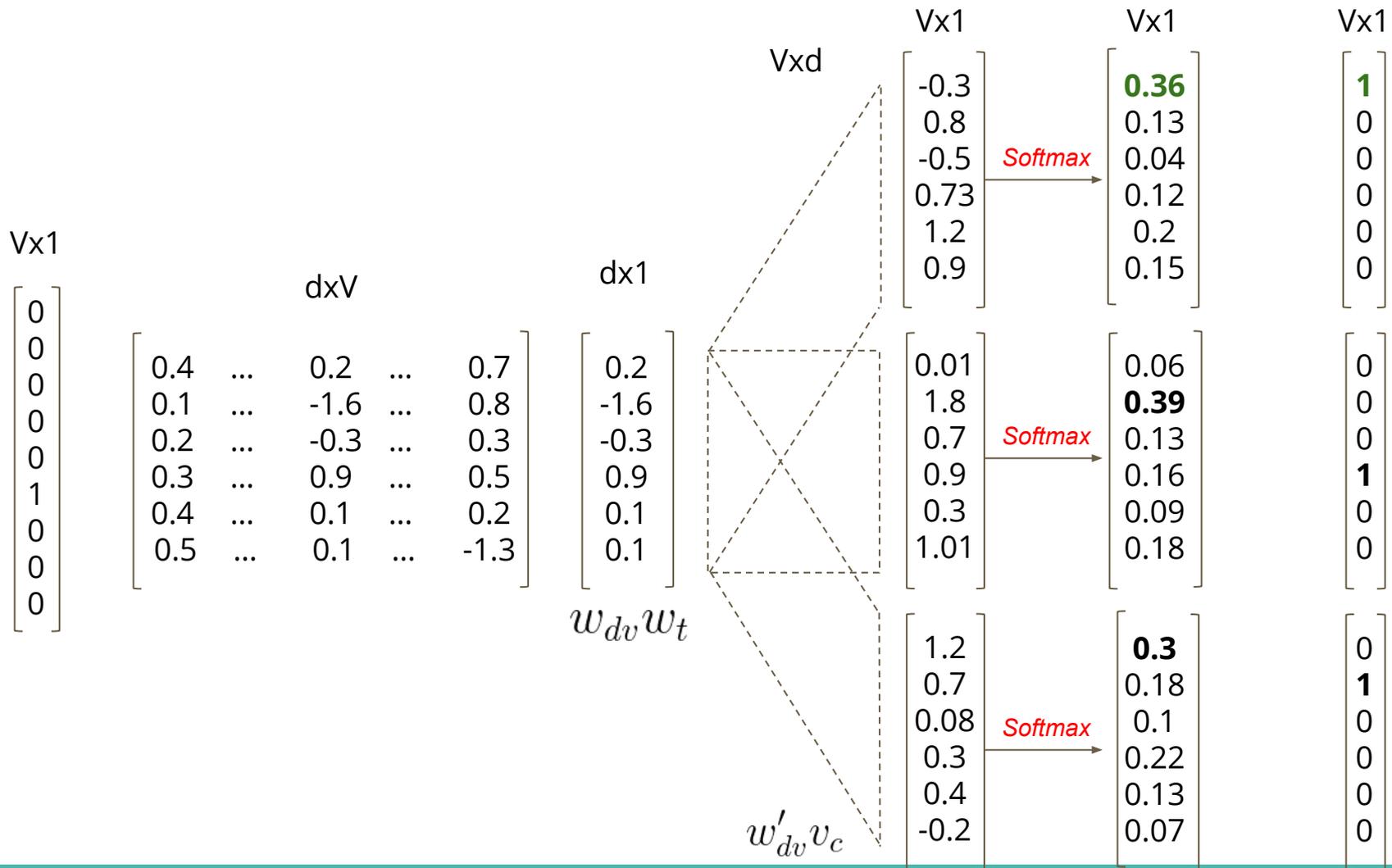
O gato sentou no tapete



Skip-Gram

O gato **sentou** no tapete





# Processo de construção do modelo

- Construir o vocabulário do corpus
- Gerar o input para a arquitetura (context, target)
- Construir a arquitetura do modelo
- Treinar o Modelo
- Obter os embeddings

## CBOW e Skip-Gram



**avó** -0.35911 0.64988 -0.41623 0.57159 -0.38013 0.38985 -0.73346 -0.30124  
**espelho** -0.16847 1.087 -0.34621 0.2506 -0.38312 -0.49646 -0.32659  
**carro** -0.47797 1.1822 -0.30762 0.51227 -0.74199 0.29385 -0.19629 -0.21463  
**arte** 0.18662 0.8713 -0.90179 0.26126 -0.19692 -0.12022 -0.34776 -0.067401

## Word2Vec

- Desenvolvido por pesquisadores da Google;
- Possibilidade de dois tipos de treinamento:
  - Skip-Gram
  - CBOW

## FastText

- Desenvolvido por pesquisadores do Facebook;
- Consegue aproximar um valor para palavras não contidas no vocabulário;
- Cada palavra é representado por uma coleção de n-gramas:

renoir = <re, ren, eno, noi, oir, ir>

- FastText tem sido usado para diversas tarefas, entre elas classificação de texto e REN [Bojanowski, 2016];

## Wang2Vec

- Desenvolvido por pesquisadores do INESC-ID (Portugal) e CMU (EUA);
- Modificação no Word2Vec (Skip-Gram e CBOW) para tarefas envolvendo sintaxe;
- Impulsionaram resultados para Part-Of-Speech Tagging e Dependency Parsing [Ling, 2015];

## **Glove**

- Desenvolvido por pesquisadores de Stanford;
- Utiliza uma Matriz de co-ocorrência entre palavras em um contexto para aprender seu significado;
- As probabilidades de predição são calculadas com base na matriz de co-ocorrência [Pennington, 2014];

# Modelos dinâmicos, sensíveis ao contexto

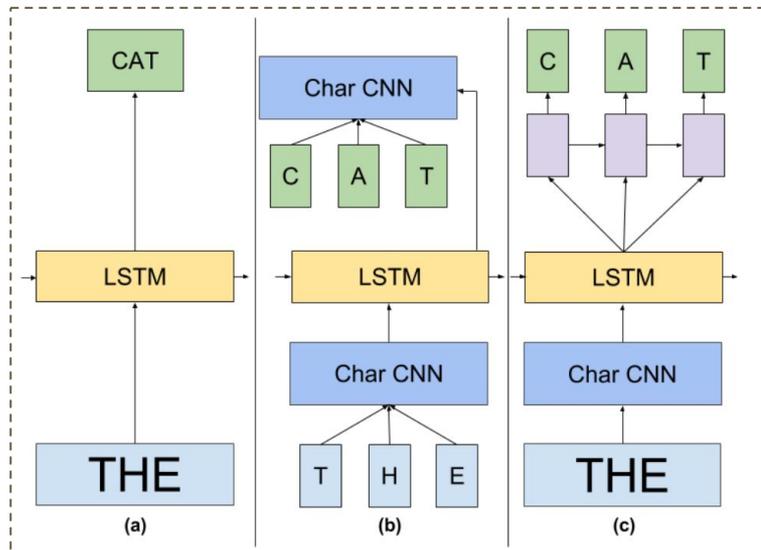
**Mangueira** `tensor([-0.0365, -0.1508, 0.0065, ..., 0.1447, 0.1526, -0.1346])`

A --> `tensor([ 0.0234, -0.0244, -0.0039, ..., 0.1887, 0.0441, -0.0241])`  
*mangueira* --> `tensor([ 0.0610, -0.1249, -0.0066, ..., 0.1447, 0.1526, -0.1346])`  
 está --> `tensor([ 0.0017, -0.0064, 0.0015, ..., -0.1571, -0.2378, -0.0496])`  
 cheia --> `tensor([ 0.1113, -0.0066, 0.0196, ..., -0.1280, 0.1558, -0.2580])`  
 de --> `tensor([ 0.2375, 0.0074, -0.0143, ..., 0.0438, 0.0729, -0.1007])`  
 frutos! --> `tensor([-0.0756, -0.0022, -0.0008, ..., 0.0000, 0.0000, 0.0000])`

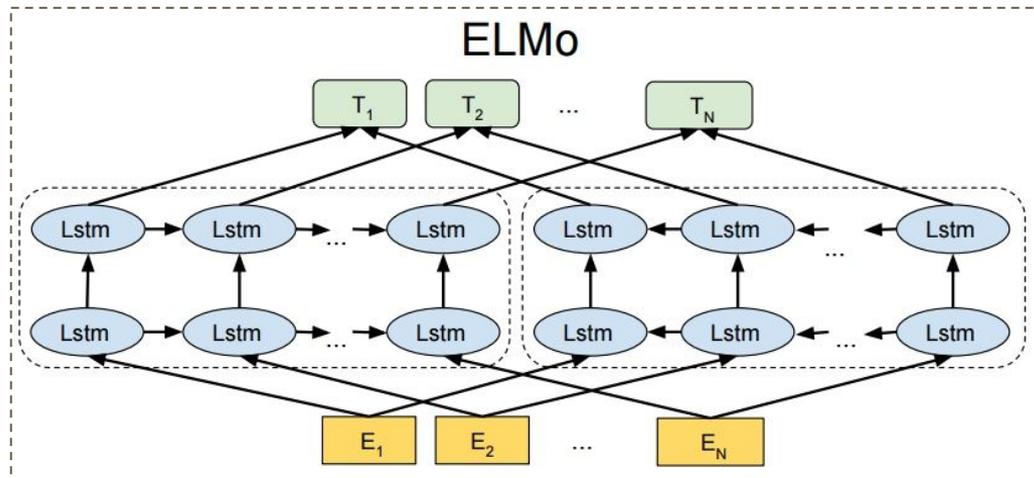
Quanto --> `tensor([ 0.0316, 0.0026, 0.0140, ..., -0.3645, 0.1062, -0.1547])`  
 custa --> `tensor([ 2.4523e-01, -1.5203e-05, 9.3146e-03, ..., 9.8593e-02, 1.5959e-01, 2.2675e-01])`  
 a --> `tensor([-0.1728, -0.0290, 0.0043, ..., 0.1887, 0.0441, -0.0241])`  
*mangueira* --> `tensor([-0.0845, -0.1070, -0.0008, ..., 0.1447, 0.1526, -0.1346])`  
 do --> `tensor([ 0.0919, 0.0257, -0.0442, ..., 0.1748, -0.0800, -0.1993])`  
 radiador? --> `tensor([ 2.0074e-01, -1.4746e-02, 1.7419e-04, ..., 0.0000e+00, 0.0000e+00, 0.0000e+00])`

Hoje --> `tensor([ 0.0342, -0.0047, 0.0275, ..., 0.0495, -0.0992, -0.0210])`  
 será --> `tensor([-0.0833, -0.0167, 0.0071, ..., -0.2242, -0.0804, -0.1034])`  
 o --> `tensor([-0.0836, -0.0149, 0.0045, ..., 0.1768, 0.1693, -0.1223])`  
 desfile --> `tensor([-0.0160, -0.0048, 0.0005, ..., -0.0340, 0.2525, -0.3314])`  
 da --> `tensor([-0.0571, 0.0065, -0.0605, ..., 0.1263, -0.1481, -0.1220])`  
*Mangueira* --> `tensor([ 0.0593, -0.0889, 0.0104, ..., 0.1447, 0.1526, -0.1346])`

# ELMo - Embeddings from Language Models



Jozefowicz et al. 2016



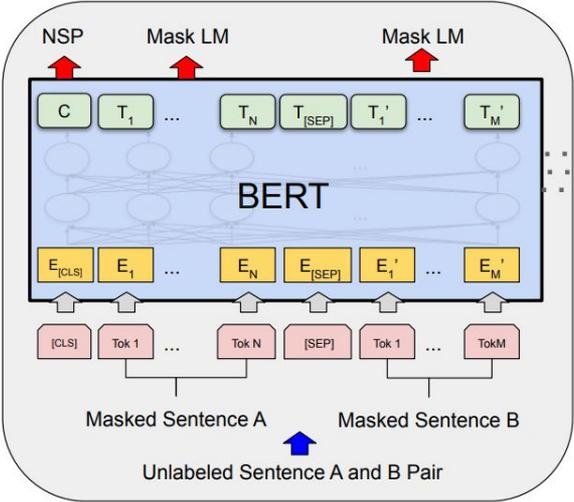
2 BiLSTM Layers com 4096 unidades  
1B Tokens para Treinamento

Peters et al. 2018

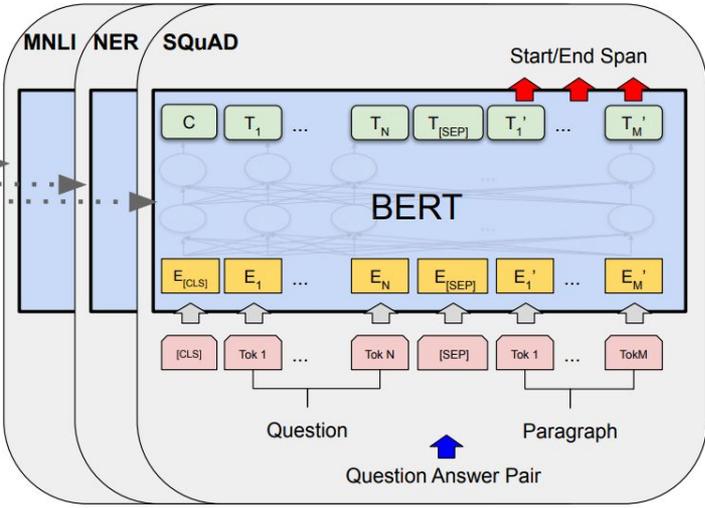
## BERT - **B**idirectional **E**ncoder **R**epresentations from **T**ransformers [Devlin, 2018]

- Recente modelo de representação de linguagem desenvolvido pelo Google;
- No artigo original de apresentação do BERT são apresentadas 11 tarefas de PLN onde o BERT foi aplicado; entre elas:
  - **Multi-Genre Natural Language Inference:** dado um par de sentenças deve-se decidir se a segunda sentença é uma continuação, contradição ou neutra em relação a primeira sentença;
  - **Semantic Textual Similarity Benchmark:** classifica a similaridade entre duas sentenças;
  - **Named Entity Recognition:** reconhecer em um texto determinadas menções e classificá-las em categorias;

BERT



Pre-training



Fine-Tuning

BERT

- Treinar um modelo BERT envolve duas tarefas:

**Masked LM**

**Next Sentence Prediction**

**Masked LM:** o objetivo é "mascarar" 15% dos tokens de uma sentença para predição. Algumas regras:

- 80% das vezes: a palavra é substituída pelo símbolo [MASK], por exemplo:

meu cachorro é cabeludo. → meu [MASK] é cabeludo.

- 10% das vezes: a palavra é substituída por uma palavra aleatória, por exemplo:

meu cachorro é cabeludo. → meu cachorro é maçã.

- 10% das vezes: a palavra é mantida, por exemplo:

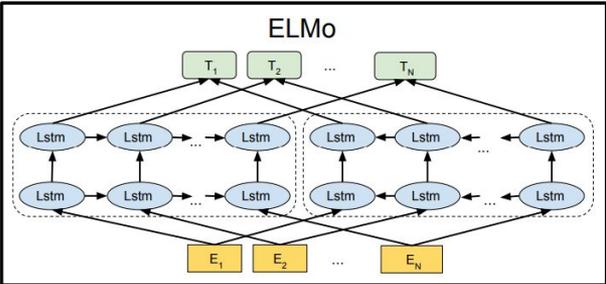
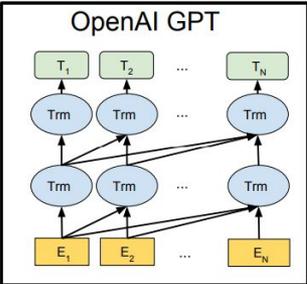
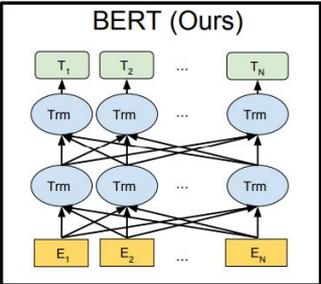
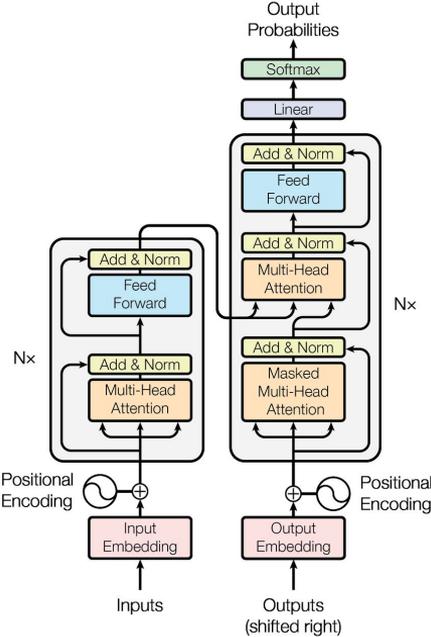
meu cachorro é cabeludo. → meu cachorro é cabeludo.

**Next Sentence Prediction:** dado um par de sentenças  $(s_1, s_2)$  o modelo deve prever se a sentença  $s_2$  é a subsequente a  $s_1$ . Algumas regras:

- 50% dos pares de sentenças são tal que  $s_2$  é de fato a sentença subsequente;
- 50% dos outros pares são tal que  $s_2$  é uma sentença aleatória do corpus;

BERT

A Rede Neural responsável por realizar as tarefas **Masked LM** e **Next Sentence Prediction** é composta por várias camadas do neurônio **Transformer** [Vaswani, 2017].



## BERT

- Treinar um modelo de linguagem BERT requer muito recurso computacional;
- Por exemplo, para o Inglês foram usadas 16 Cloud TPUs com um corpora de 3,3Bi tokens;
- Existe um modelo Multi-lingue que contempla o **Português**;

Tipo	Idioma	Camadas	Camadas Ocultas	Self-Attention heads
BERT <sub>BASE</sub>	Inglês	12	768	12
BERT <sub>LARGE</sub>	Inglês	24	1024	16
BERT <sub>BASE</sub>	Multi-lingue	12	768	12

**Flair Embedding** - é um recente modelo de embedding que permite a modelagem da linguagem com base na distribuição das sequências de caracteres e palavras [Akbik, 2018];

- A geração de um Flair Embedding **não só depende** do contexto das **palavras** vizinhas, mas também do nível de **caractere** das palavras vizinhas;
- O modelo é gerado a partir de duas redes **LSTM** que capturam e incorporam as informações aprendidas durante o treinamento;

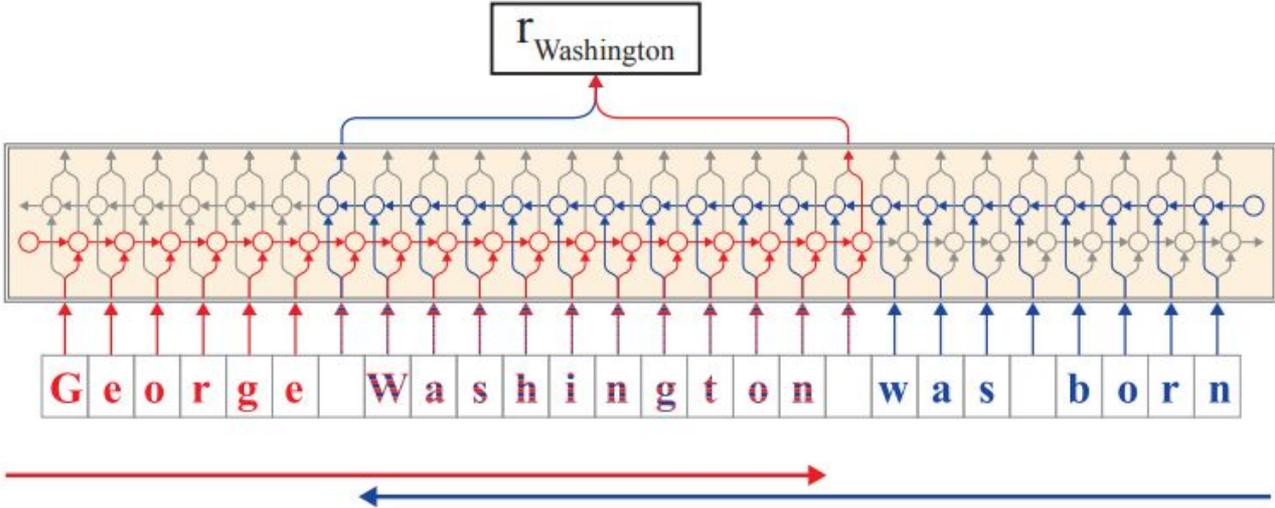
**LSTM - Long Short-Term Memory** - é um tipo de Rede Neural Recorrente, com uma estrutura computacional mais complexa, que tem tido sucesso na resolução de tarefas sequenciais [Tai, 2015];

- Uma LSTM permite manter, alterar ou descartar informações anteriores para relacionar com uma informação atual;

Há uma variação das redes LSTM, que são as **Bidirectional LSTM (Bi-LSTM)**;

- As redes Bi-LSTM consistem de duas LSTM que funcionam em paralelo;

Flair



- Dado  $X_{0:T}$  uma sequência de caracteres, que produz uma linguagem natural;

$$X_{0:T} := (x_0, x_1, \dots, x_{t-1}, x_t)$$

- Será possível prever  $\{x_t\} \in X_{0:T} | X_{0:t-1} \subset X_{0:T}$ ?

- Dado  $X_{0:T}$  uma sequência de caracteres, que produz uma linguagem natural;

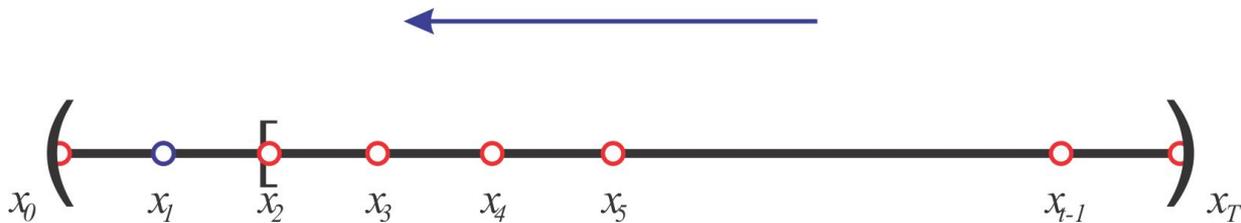
$$X_{0:T} := (x_0, x_1, \dots, x_{t-1}, x_t)$$

- Será possível prever  $\{x_t\} \in X_{0:T} | X_{0:t-1} \subset X_{0:T}$ ?
- **Sim**, podemos aprender como os caracteres dessa linguagem estão distribuídos, abstraindo um modelo de linguagem [Rosenfeld, 2000];
- Então quando dado  $X_{0:x_{t-1}}$  queremos prever  $\{x_t\} \in X_{0:T}$ :

$$P(x_t | x_0, x_1, \dots, x_{t-1}) = P(x_{0:T}) \approx \prod_{t=0}^T P(x_t | h_t; \theta)$$

- Então quando dado  $X_{0:x_{t-1}}$  queremos prever  $\{x_t\} \in X_{0:T}$ :

$$P(x_t | x_0, x_1, \dots, x_{t-1}) = P(x_{0:T}) \approx \prod_{t=0}^T P(x_t | h_t; \theta)$$



- Como Flair Embedding é treinado por duas LSTM temos dois modelos:

$$p^f(x_t|X_{t+1:T}) \approx \prod_{t=0}^T p^f(x_t|h_t^f; \theta)$$

$$h_t^f = f_h^f(x_{t-1}, h_{t-1}^f, c_{t-1}^f; \theta)$$

$$c_t^f = f_c^f(x_{t-1}, h_{t-1}^f, c_{t-1}^f; \theta)$$

**Forward**

$$p^b(x_t|X_{t+1:T}) \approx \prod_{t=0}^T p^b(x_t|h_t^b; \theta)$$

$$h_t^b = f_h^b(x_{t-1}, h_{t-1}^b, c_{t-1}^b; \theta)$$

$$c_t^b = f_c^b(x_{t-1}, h_{t-1}^b, c_{t-1}^b; \theta)$$

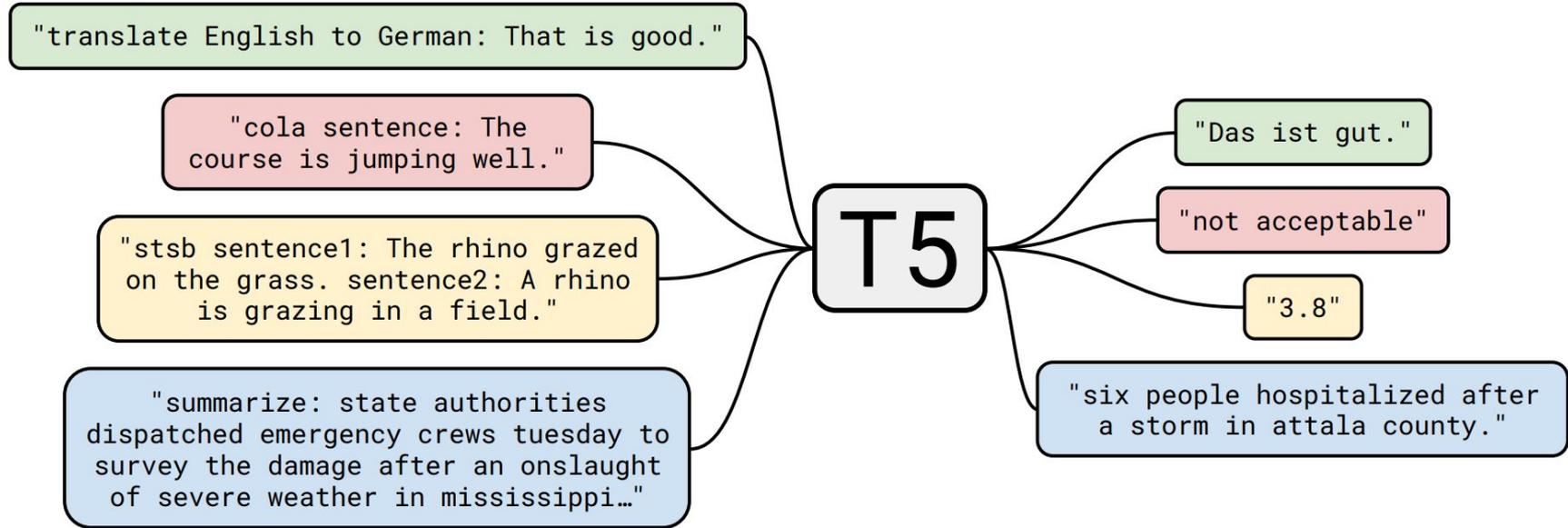
**Backward**

- Os processos de **Forward** e **Backward** produzem duas saídas para cada palavra:  $h_t^f$  e  $h_t^b$  ;
- Essas saídas são empilhadas em uma matriz, formando o embedding final:

$$w^{CharLM} := \begin{bmatrix} h_{t_{i+1}-1}^f \\ h_{t_i-1}^b \end{bmatrix}$$

- Treinamento Supervisionado e não-supervisionado (unsupervised);
  - NLC e CV (ImageNet)
- Um dos pontos mais fortes do aprendizado não-supervisionado é não precisar de dado anotado;
- Algoritmos não-supervisionados para geração de Modelos de Linguagem;
- Um grande corpus: C4 (Colossal Clean Crawled Corpus)

T5



## Text-to-Text Transfer Transformer

T5

- Classificação de Texto (GLUE e SuperGLUE)
- Sumarização (CNN/Daily Mail)
- QA (SQuAD)
- MT (WMT)

- Classificação de Texto (GLUE e SuperGLUE)
- Sumarização (CNN/Daily Mail)
- QA (SQuAD)
- MT (WMT)

Um problema é se o modelo gera “hambúrguer” em vez dos labels aceitos.

- Classificação de Texto (GLUE e SuperGLUE)
- Sumarização (CNN/Daily Mail)
- QA (SQuAD)
- MT (WMT)

Um problema é se o modelo gera “hamburger” em vez dos labels aceitos.

Mas nunca observaram esse comportamento

- Pre-Training
  - Steps = 524,288
  - Max\_length = 512
  - Batch\_size = 128 sequencias
  - Tokens =  $128 * 512 = 65.536$ /batch
  - Steps \* batch\_tokens  $\approx$  34Bi tokens
- Fine-tuning
  - Steps = 262,144
  - Max\_length = 512
  - Batch\_size = 128 sequencias
  - Tokens =  $128 * 512 = 65.536$ /batch
  - Learning\_rate = 0.001

- Mas e....

**RoBERTa**

**XLEmbeddings**

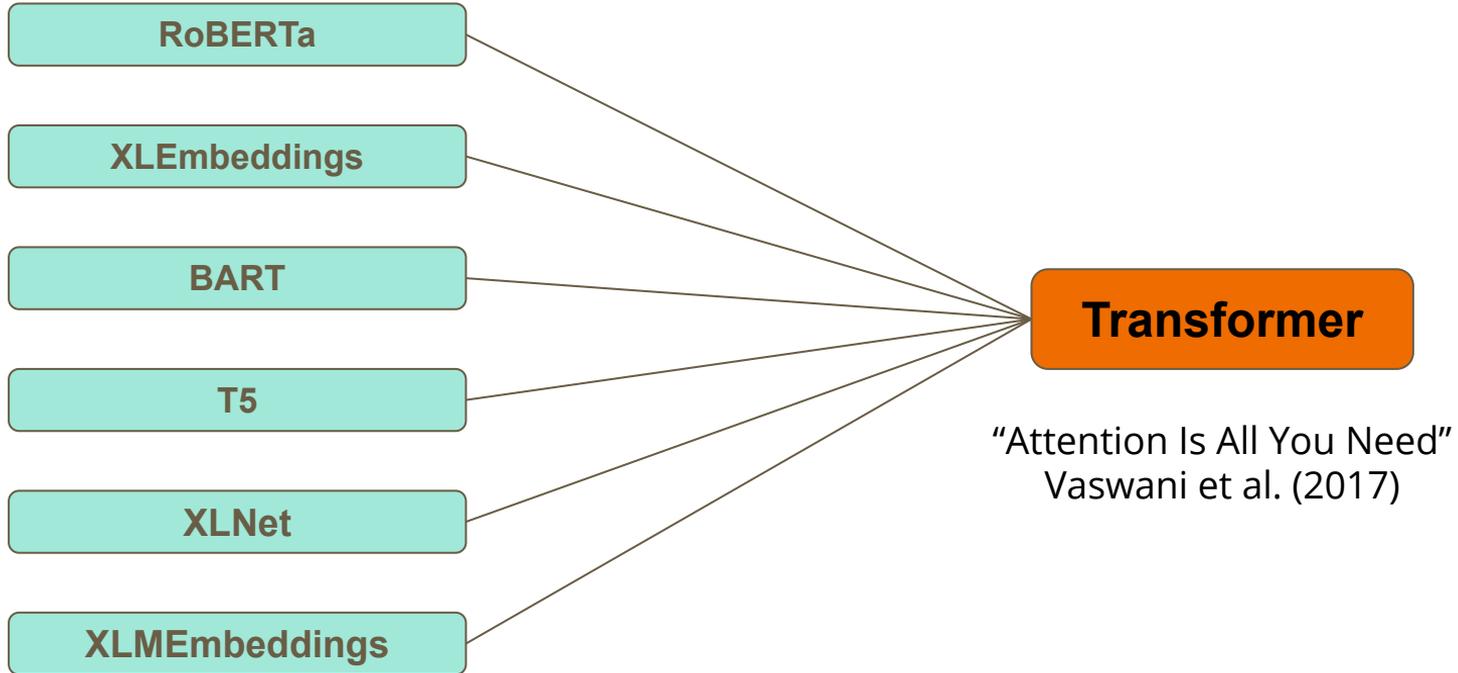
**BART**

**T5**

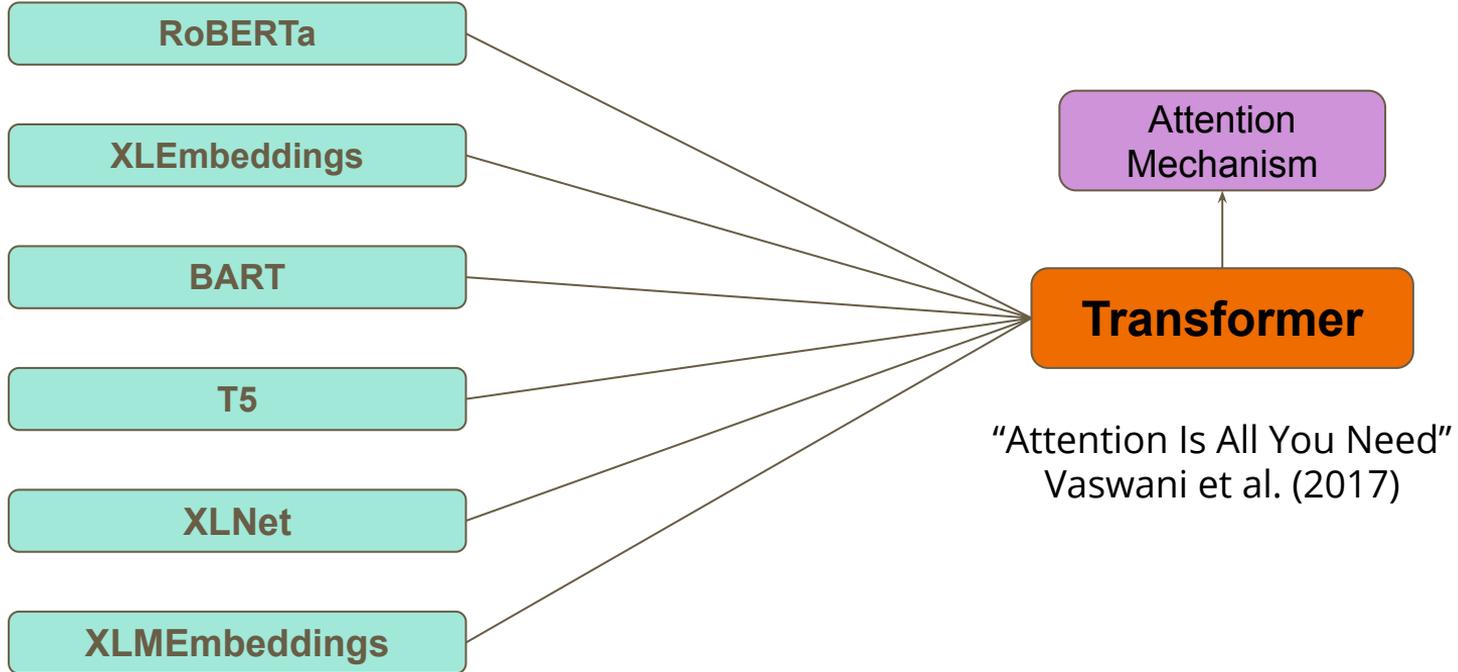
**XLNet**

**XLME embeddings**

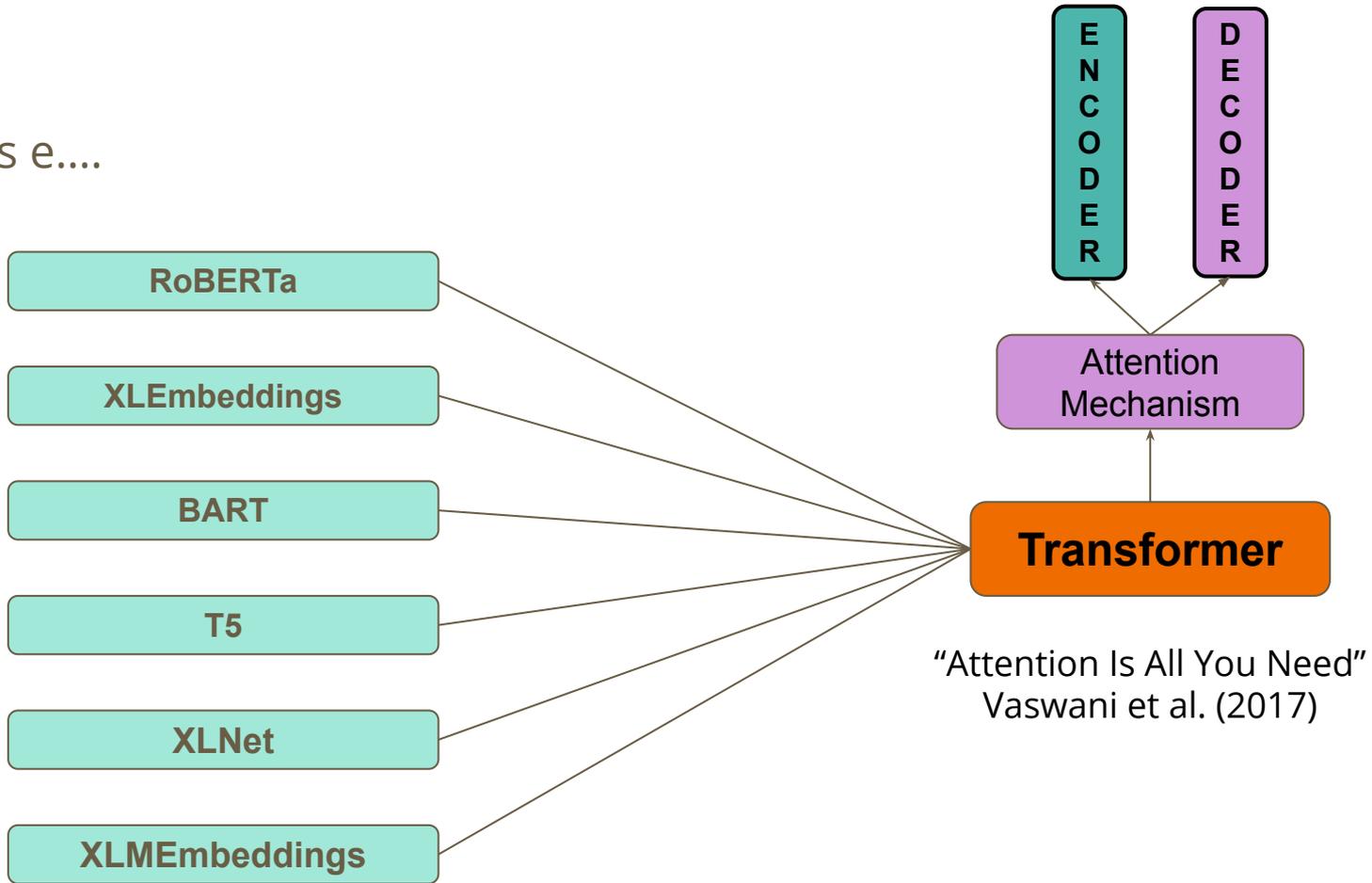
- Mas e....



- Mas e....



- Mas e....

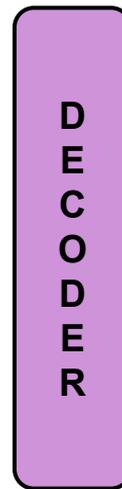
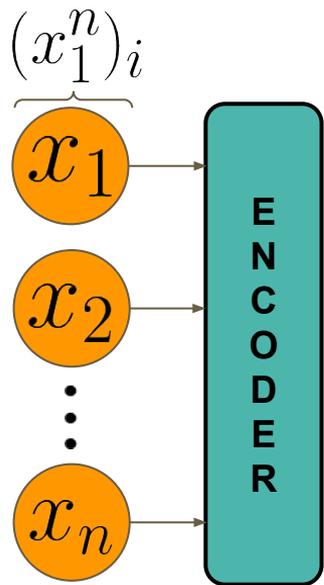


Attention  
Mechanism

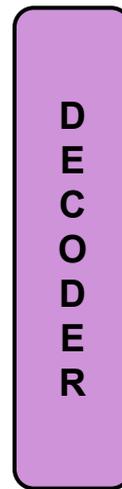
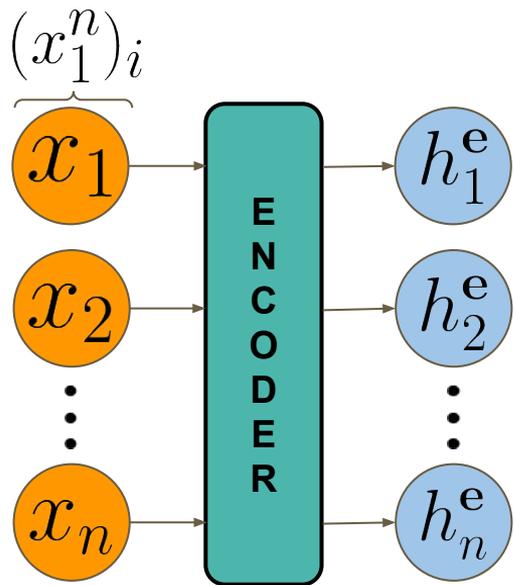
**E  
N  
C  
O  
D  
E  
R**

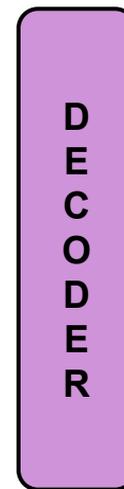
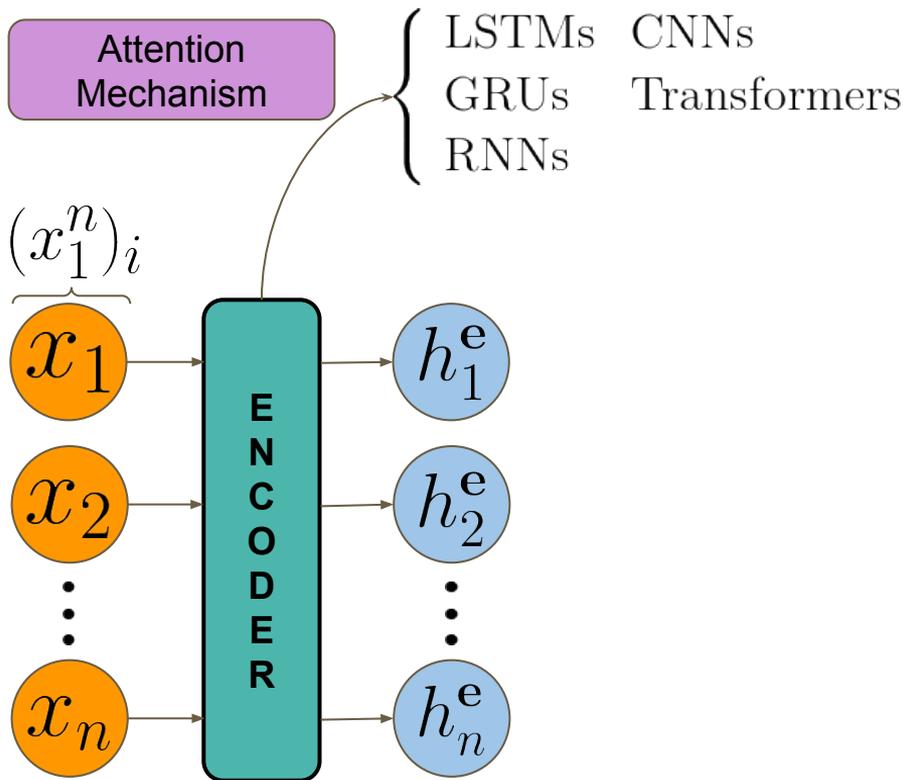
**D  
E  
C  
O  
D  
E  
R**

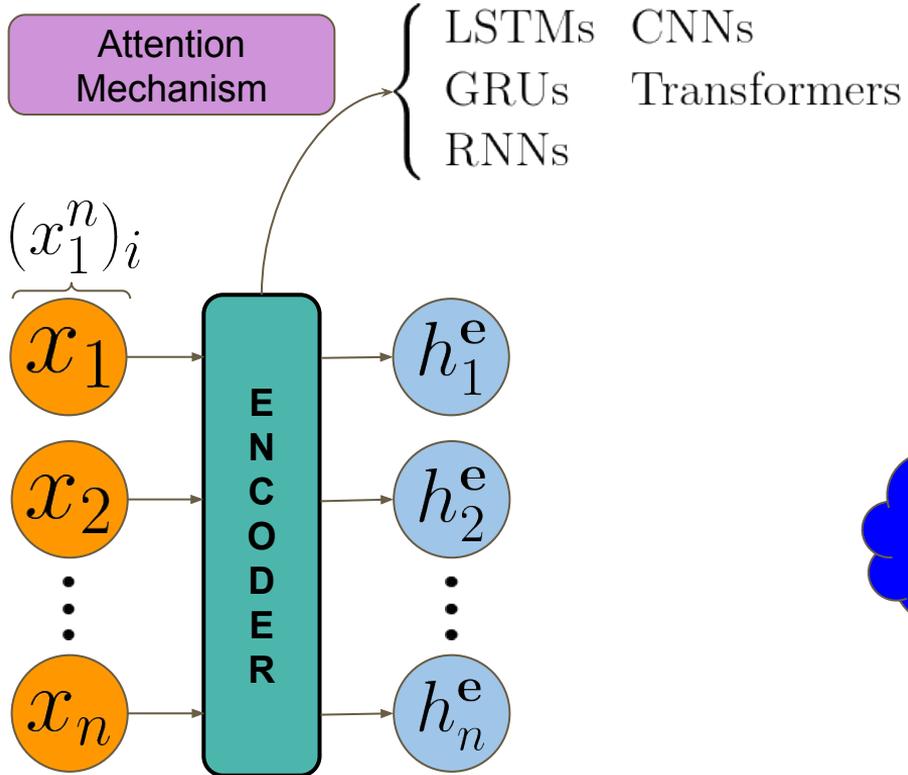
# Attention Mechanism

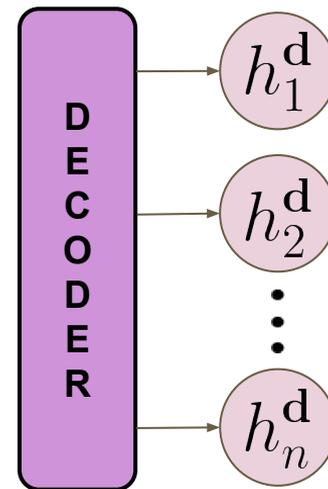
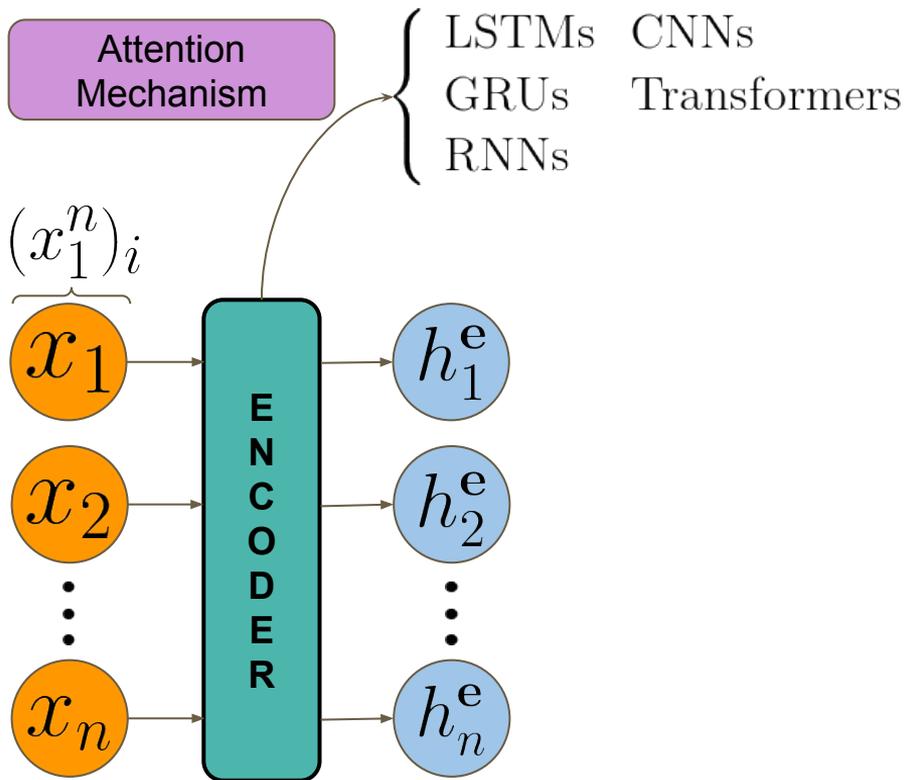


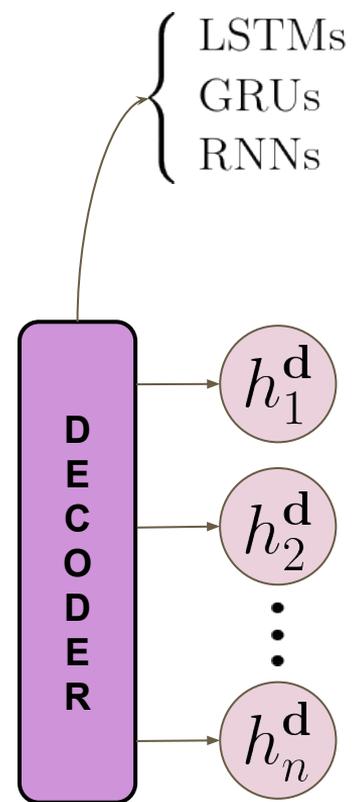
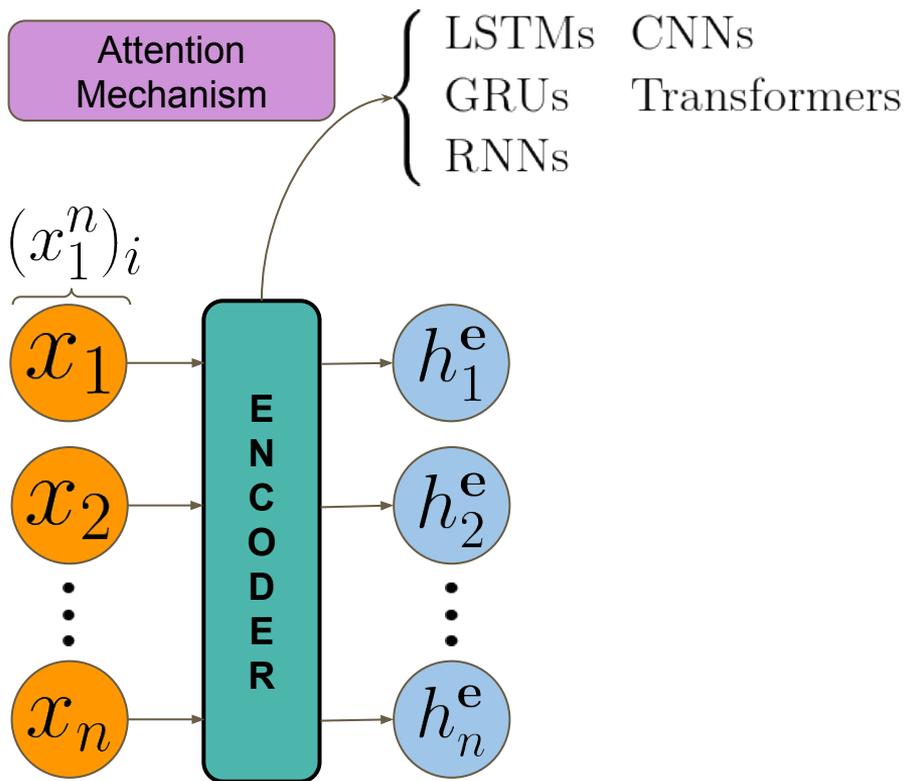
Attention  
Mechanism

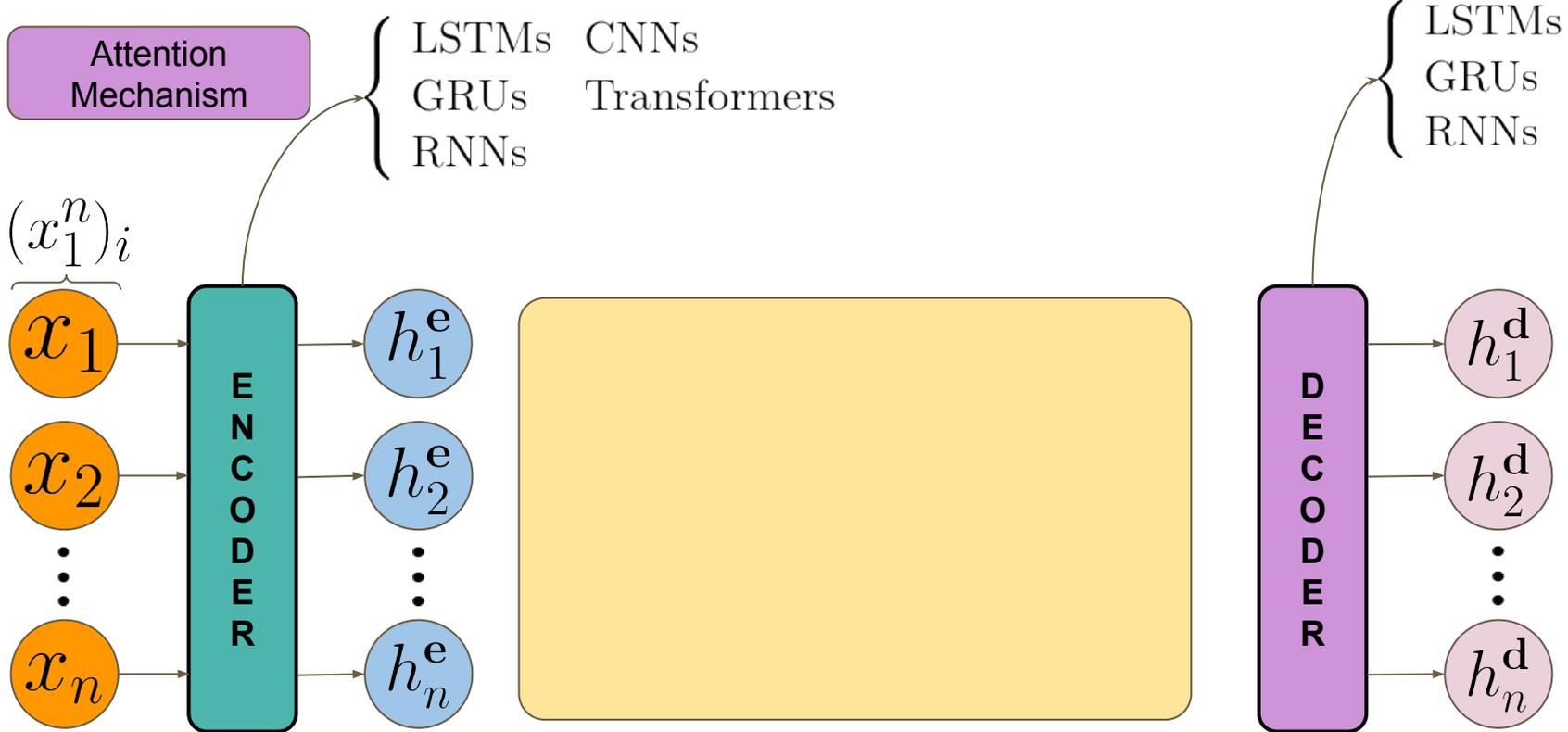


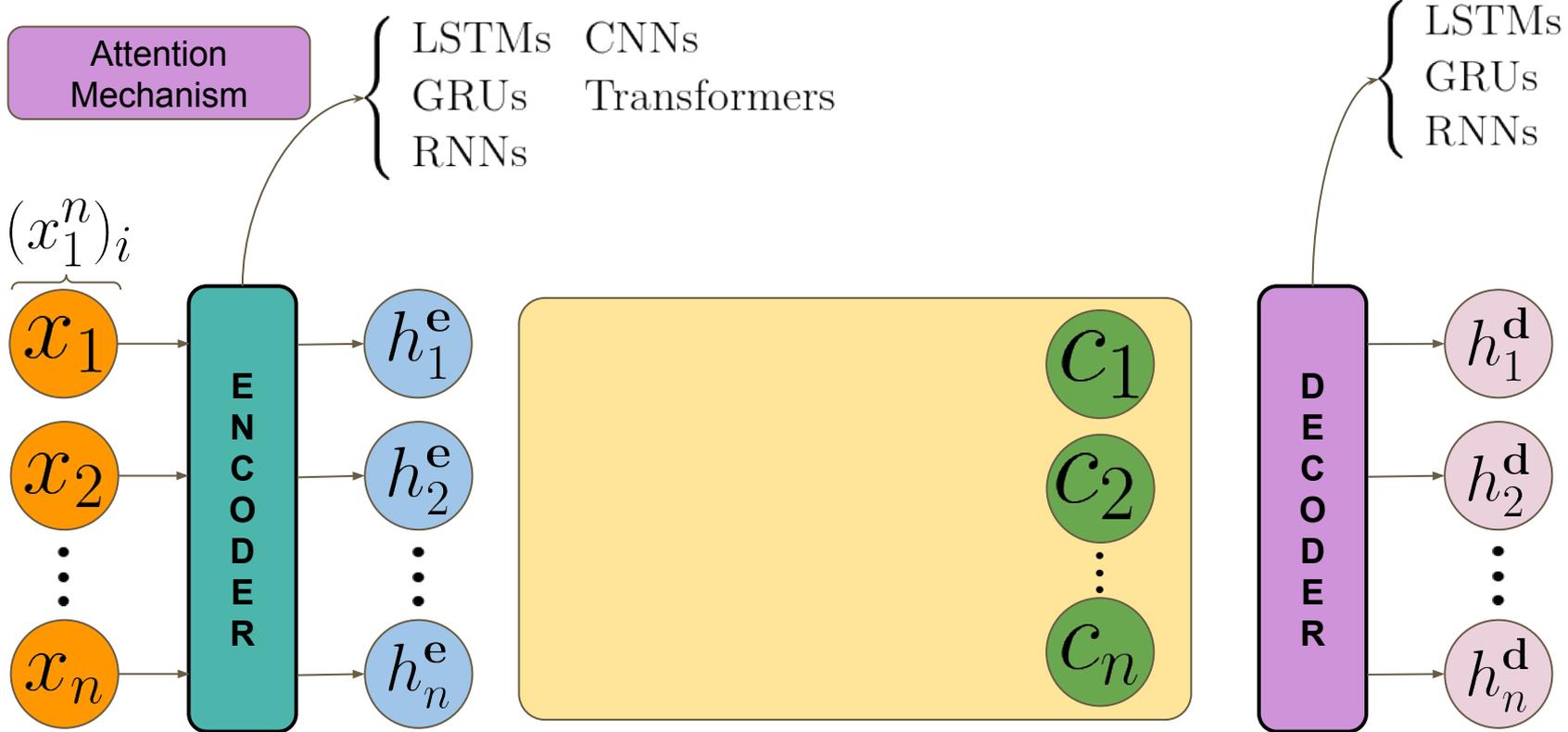


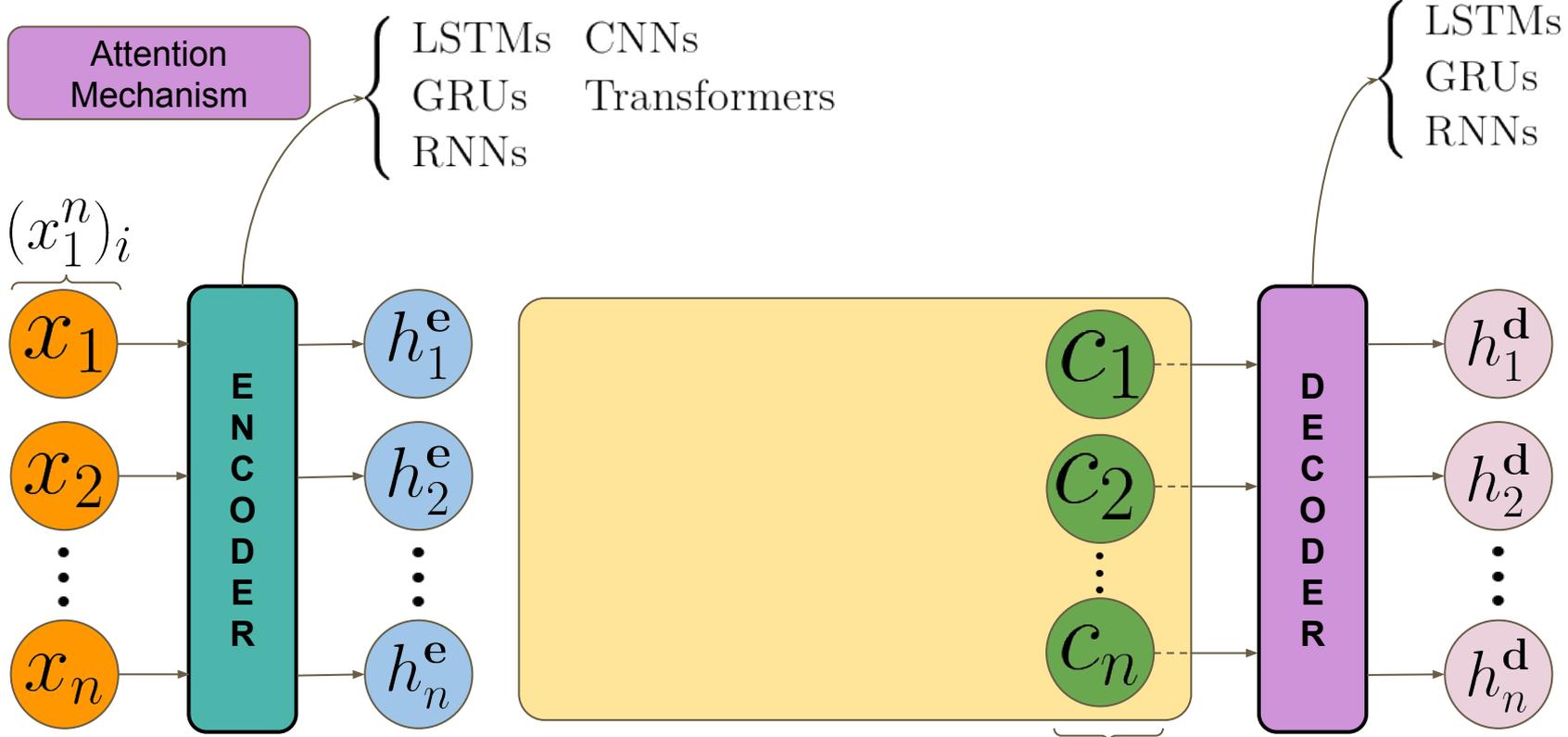




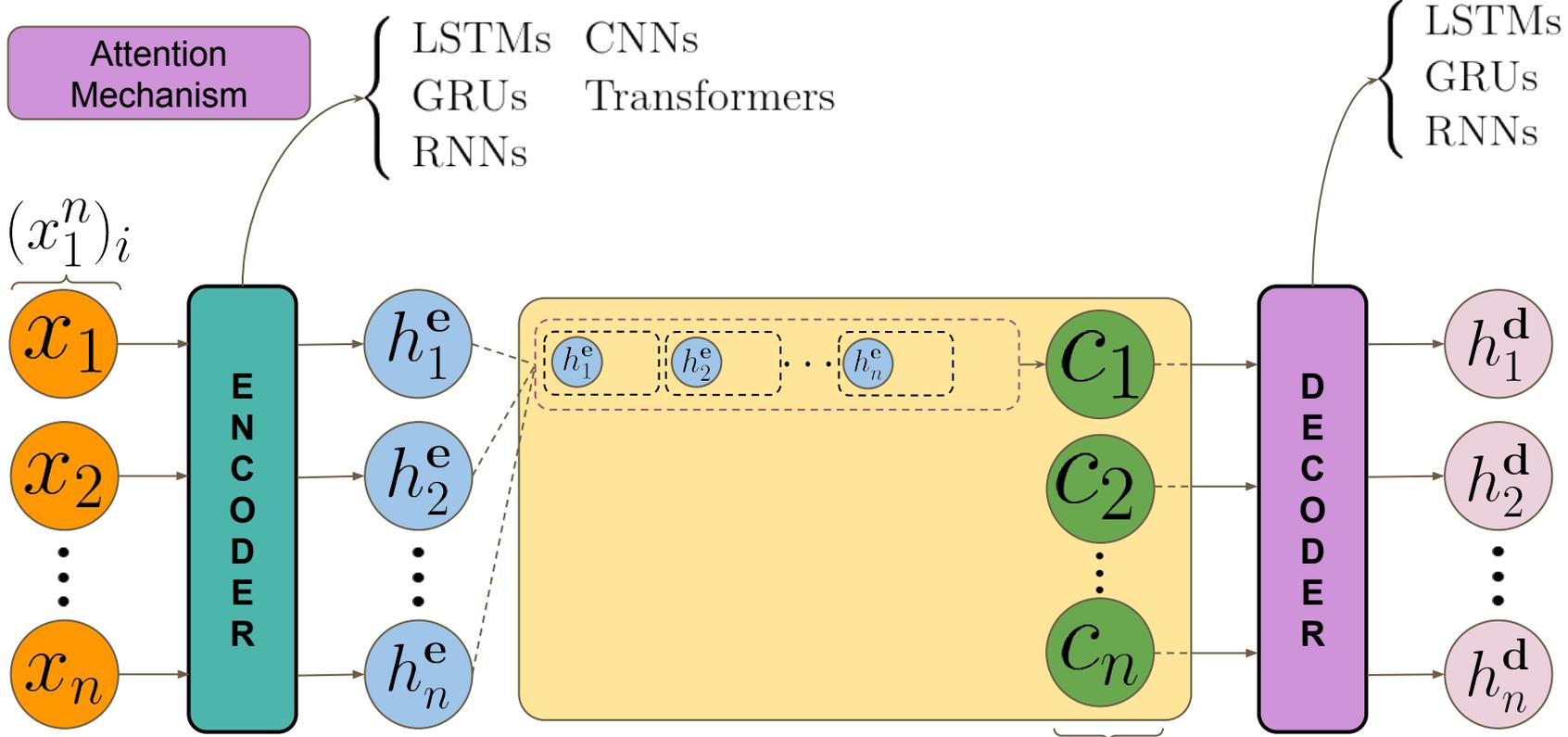




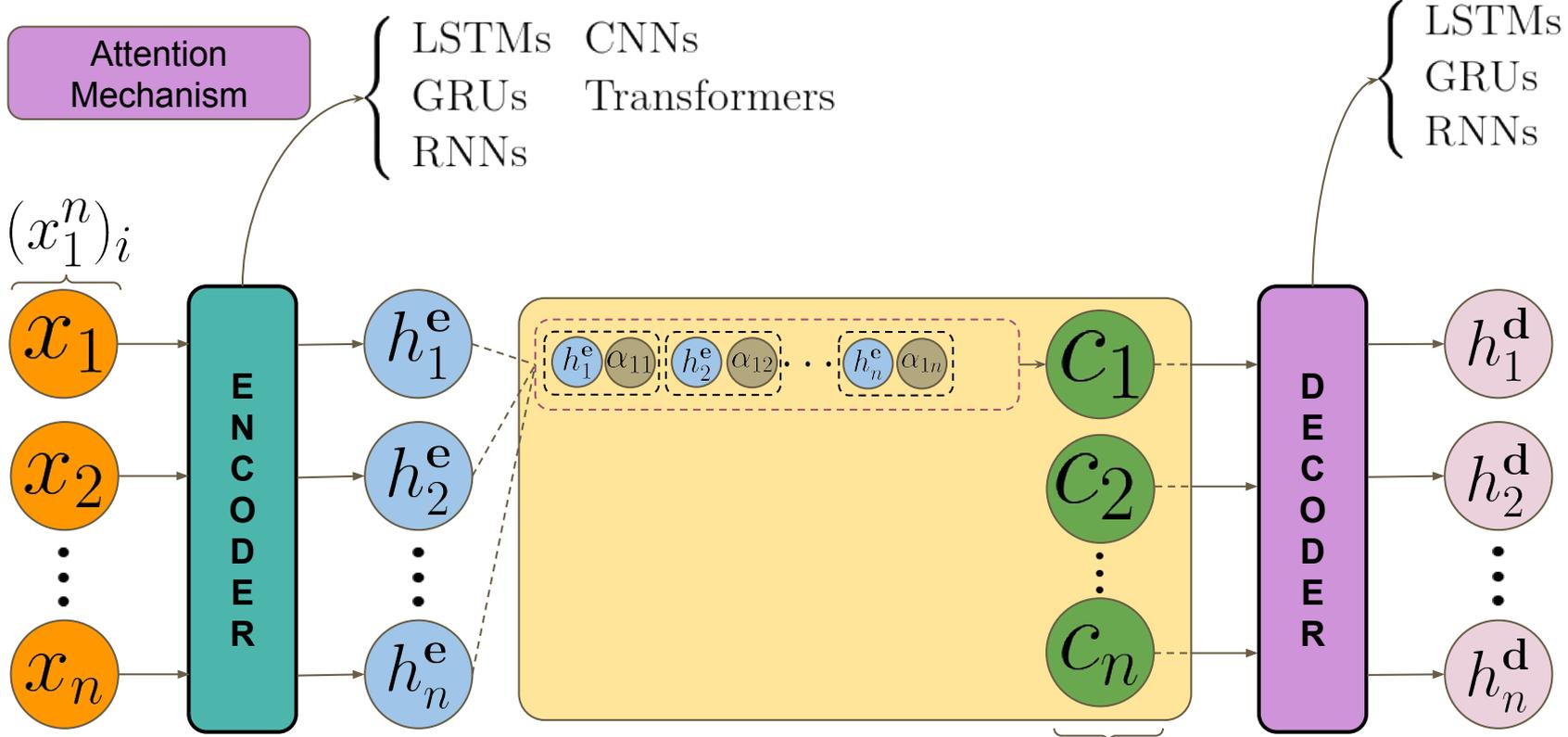




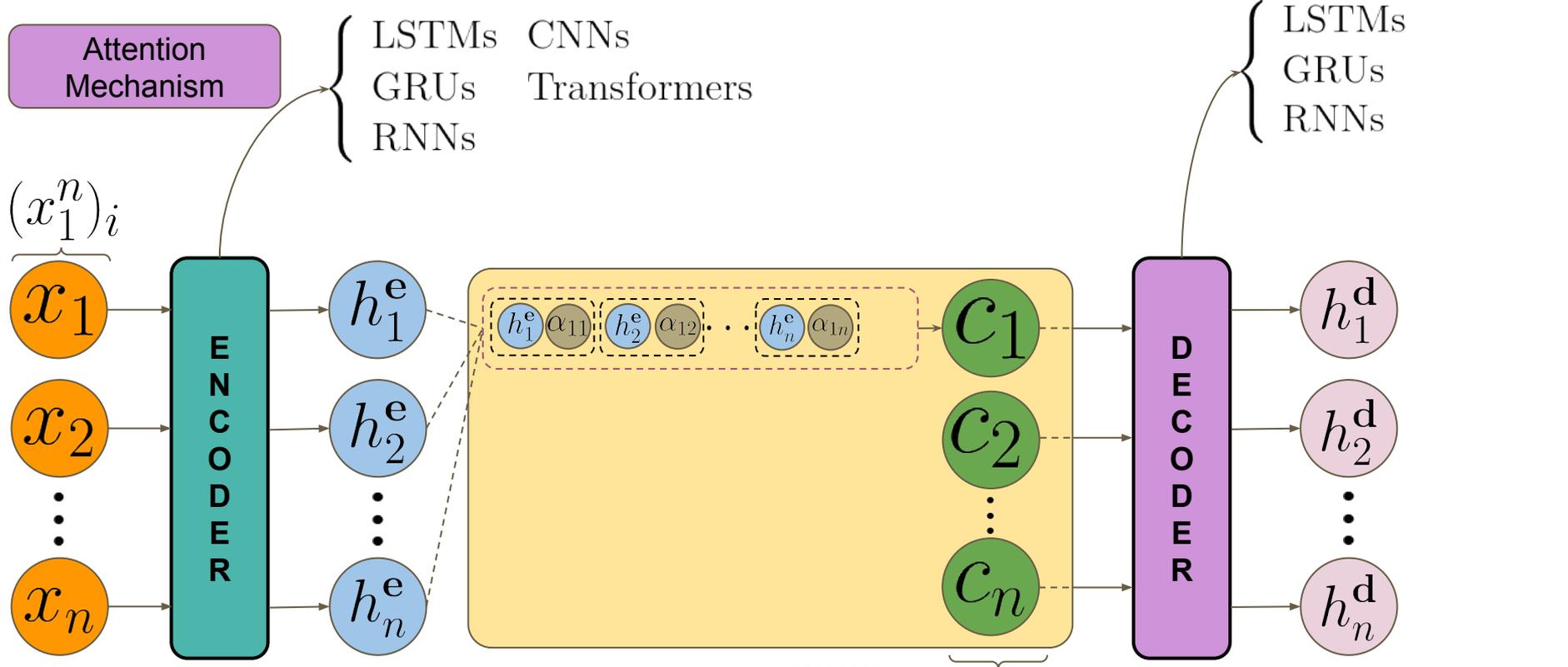
$$c_i = \sum_j \alpha_{ij} h_j^e$$



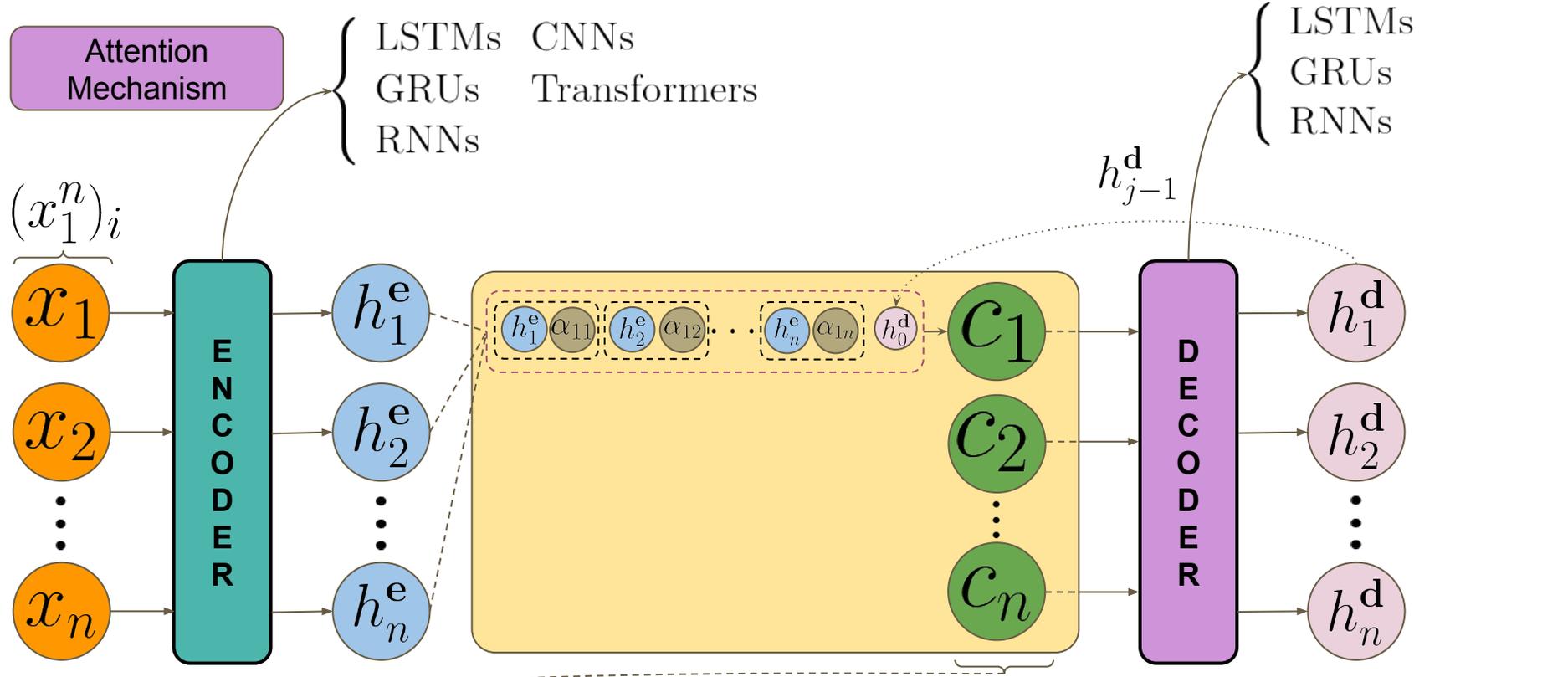
$$c_i = \sum_j \alpha_{ij} h_j^e$$



$$c_i = \sum_j \alpha_{ij} h_j^e$$



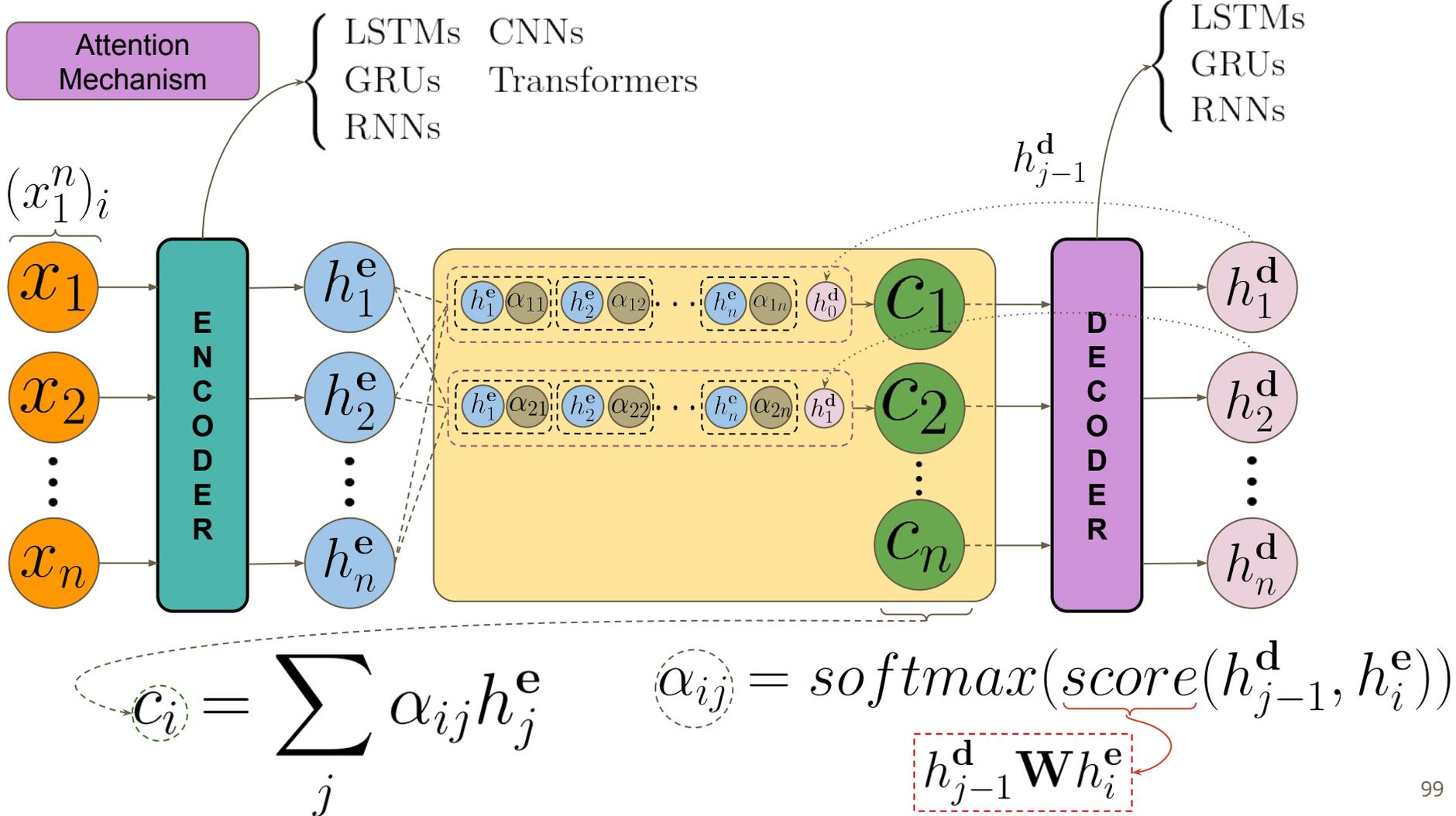
$$c_i = \sum_j \alpha_{ij} h_j^e \quad \alpha_{ij} = \text{softmax}(\text{score}(h_{j-1}^d, h_i^e))$$

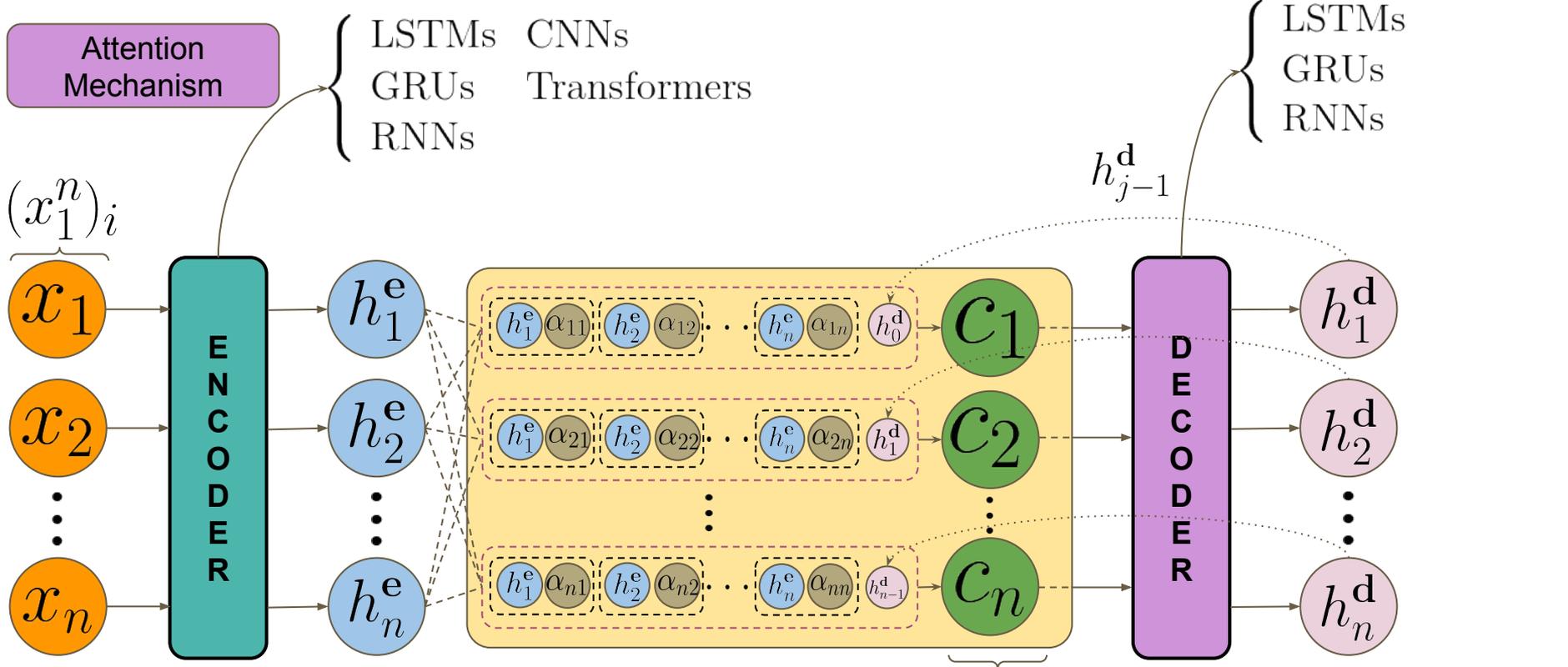


$$c_i = \sum_j \alpha_{ij} h_j^e$$

$$\alpha_{ij} = \text{softmax}(\text{score}(h_{j-1}^d, h_i^e))$$

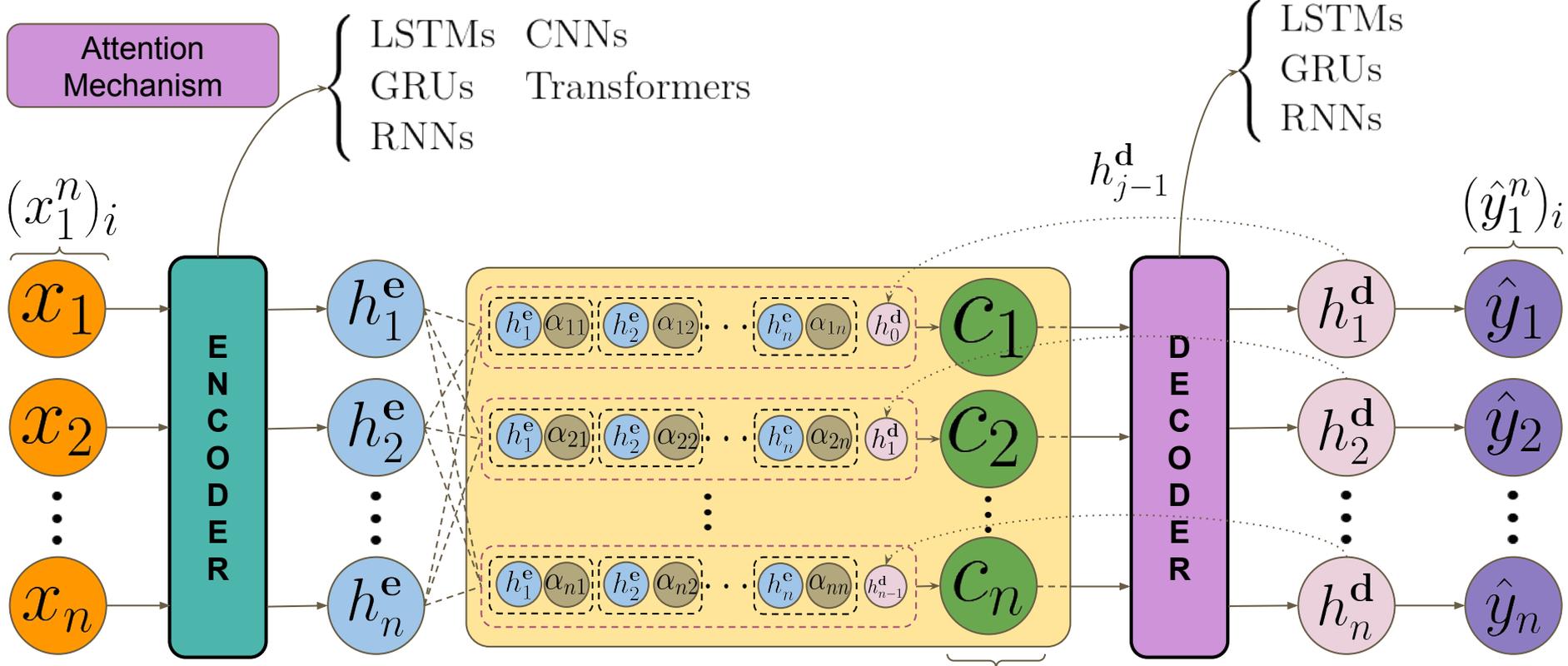
$$h_{j-1}^d \mathbf{W} h_i^e$$





$$c_i = \sum_j \alpha_{ij} h_j^e \quad \alpha_{ij} = \text{softmax}(\text{score}(h_{j-1}^d, h_i^e))$$

$h_{j-1}^d \mathbf{W} h_i^e$



$$c_i = \sum_j \alpha_{ij} h_j^e \quad \alpha_{ij} = \text{softmax}(\text{score}(h_{j-1}^d, h_i^e))$$

$$h_{j-1}^d \mathbf{W} h_i^e$$

Attention  
Mechanism

$$\text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{1}{\sqrt{d_k}} \mathbf{Q}\mathbf{K}^T\right) \mathbf{V}$$

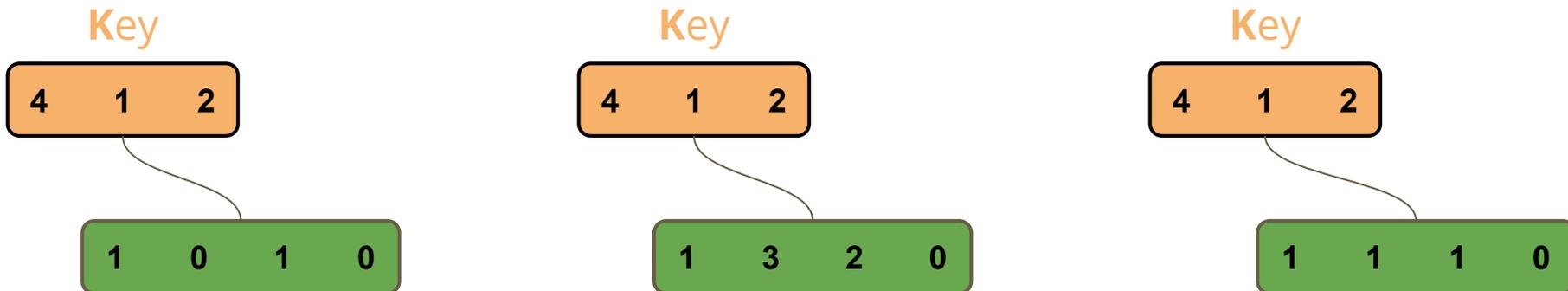
1 0 1 0

1 3 2 0

1 1 1 0

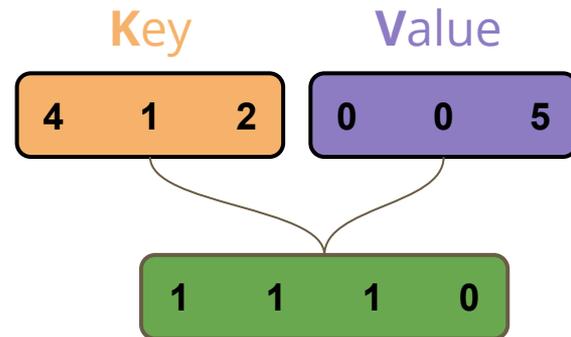
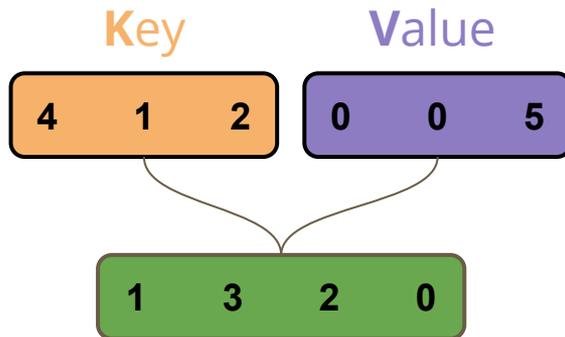
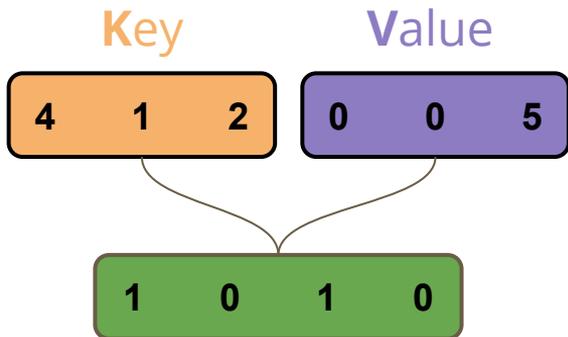
## Attention Mechanism

$$\text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{1}{\sqrt{d_k}} \mathbf{Q}\mathbf{K}^T\right) \mathbf{V}$$



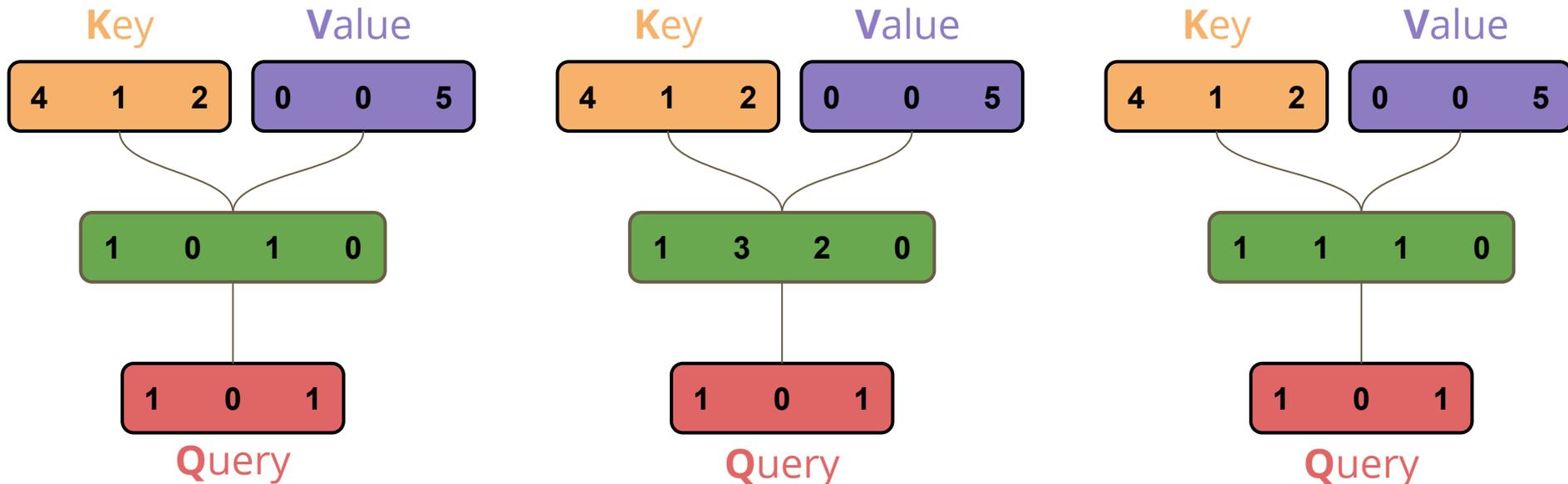
# Attention Mechanism

$$\text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{1}{\sqrt{d_k}} \mathbf{Q}\mathbf{K}^T\right) \mathbf{V}$$



# Attention Mechanism

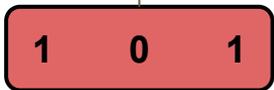
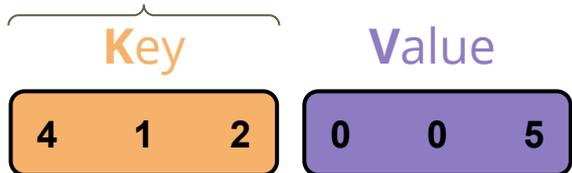
$$\text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{1}{\sqrt{d_k}} \mathbf{Q}\mathbf{K}^T\right) \mathbf{V}$$



Attention Mechanism

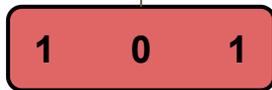
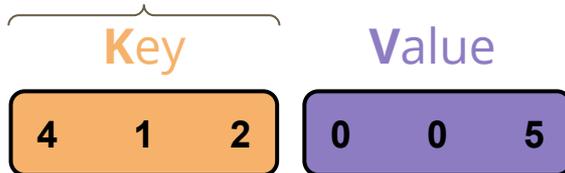
$$\text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{1}{\sqrt{d_k}} \mathbf{Q}\mathbf{K}^T\right) \mathbf{V}$$

$\dim(\mathbf{K}) = d_k$



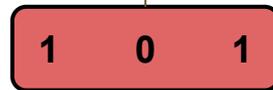
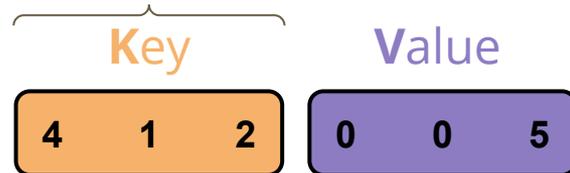
Query  
 $\dim(\mathbf{Q}) = d_k$

$\dim(\mathbf{K}) = d_k$



Query  
 $\dim(\mathbf{Q}) = d_k$

$\dim(\mathbf{K}) = d_k$

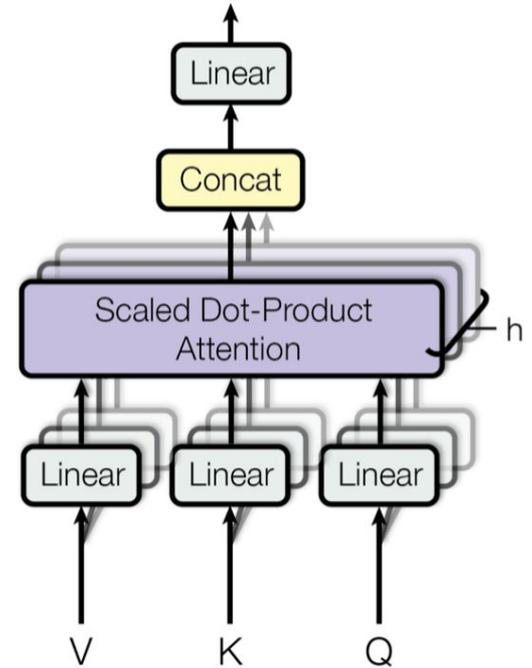


Query  
 $\dim(\mathbf{Q}) = d_k$

$$\text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{1}{\sqrt{d_k}} \mathbf{Q}\mathbf{K}^T\right) \mathbf{V}$$

$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) \mathbf{W}^O$$

$$\text{head}_i = \text{Att}(\mathbf{Q}\mathbf{W}_i^Q, \mathbf{K}\mathbf{W}_i^K, \mathbf{V}\mathbf{W}_i^V)$$

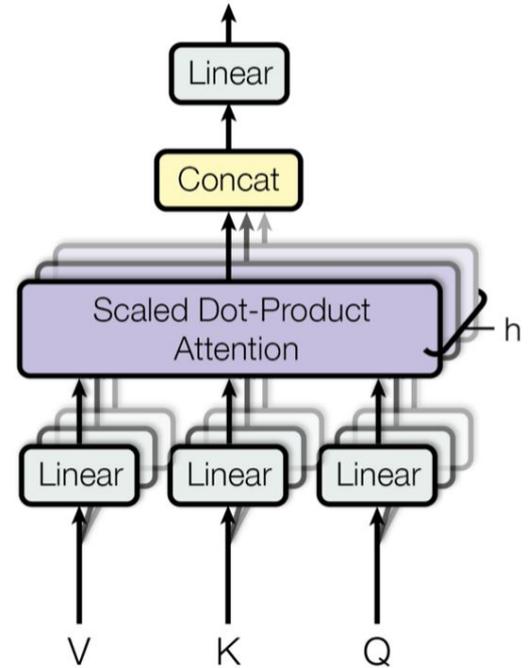


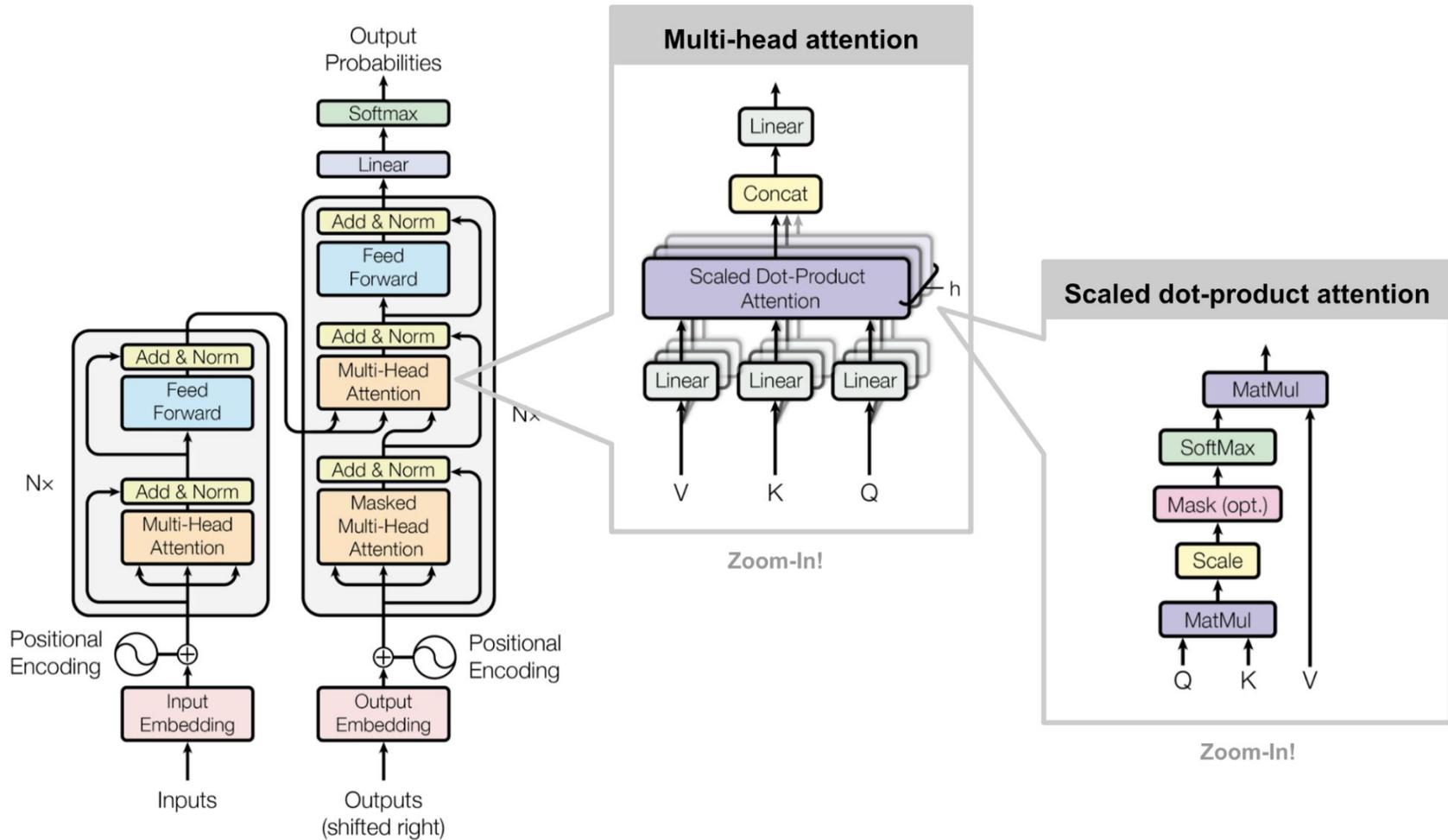
$$\text{Att}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{1}{\sqrt{d_k}} \mathbf{Q}\mathbf{K}^T\right) \mathbf{V}$$

$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) \mathbf{W}^O$$

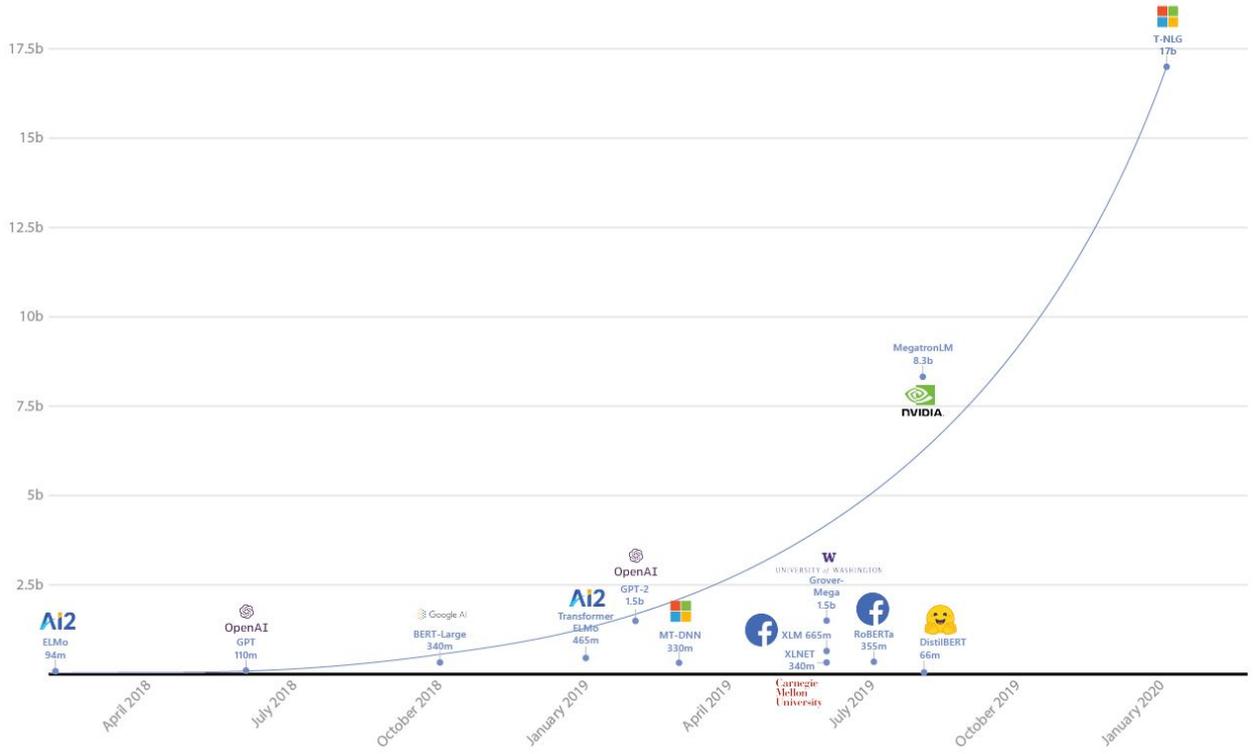
$$\text{head}_i = \text{Att}(\mathbf{Q}\mathbf{W}_i^Q, \mathbf{K}\mathbf{W}_i^K, \mathbf{V}\mathbf{W}_i^V)$$

**Matrizes de parâmetros a serem aprendidas**





# Curiosidade



- Dicas de Leitura

- Jurafsky
- <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html#summary>
- Jay Alammar - <http://jalammar.github.io/illustrated-transformer/>
- “Attention is all you need!” - <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>
- <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html#summary>

# Modelos para a língua portuguesa

## Word Embeddings para o Português

### LX - Center Portugal

1,7 *Bilhões de Tokens*

Word2Vec

LX Center - Disponível em:

<http://lxcenter.di.fc.ul.pt/datasets/en/index.html>

### NILC Embeddings

1,3 *Bilhões de Tokens*

Word2Vec

FastText

Wang2vec

Glove

NILC Embeddings - Disponível em:

<http://nilc.icmc.usp.br/embeddings>

### PLN-PUCRS

4,9 *Bilhões de Tokens*

Word2Vec

FastText

PLN-PUCRS - Disponível em:

<https://github.com/jneto04/ner-pt>

## Flair Embeddings para o Português

### Charles University FlairEL

*0,9 Bilhões de  
Tokens*

<i>Forward</i>	2.78
<i>Backward</i>	2.81

**Eric Lief - Disponível em:**  
<https://github.com/zalandoresearch/flair>

### PUCRS FlairBBP

*4,9 Bilhões de  
Tokens*

<i>Forward</i>	2.76
<i>Backward</i>	2.80

**FlairBBP - Disponível em:**  
<https://github.com/jneto04/ner-pt>

## WE Português - Tweeter

**U. Porto, Portugal**

**Tweeter**

*5 milhões de tweets*

**Bengio, 2003**

**UPorto - Disponível em:**

<https://github.com/saleiro/embedpt>

BERT Multilingue  
(inclui o Português)

**BERT - PyTorch**

*Dump Wikipédia  
Top 100*

**PyTorch - Disponível em:**

<https://github.com/huggingface/pytorch-pretrained-BERT>

**BERT - TensorFlow**

*Dump Wikipédia  
Top 100*

**TensorFlow - Disponível em:**

<https://github.com/google-research/bert>

# Modelos específicos de domínio

## Word Embeddings de domínio

**PLN-PUCRS**

**Saúde**

*603 milhões de  
tokens*

**Word2Vec**

**FastText**

**PLN-PUCRS - Disponível em:**

<https://github.com/nlp-pucrs/cci-regression>

## Word Embeddings de domínio

**Petrobrás**

**Petróleo e Gás**

*10,5 milhões de  
Tokens*

**Word2Vec**

Disponível em:

<https://github.com/diogosmg/wordEmbeddingsOG/tree/master/geracaoEmbeddings/modelos>

**Norway - Sirius  
Univ. de Oslo**

**Petróleo e Gás**

*8,2 milhões de  
Sentenças*

**Word2Vec**

**FastText**

Disponível em:

<http://vectors.nlpl.eu/repository/11/75.zip>

# Geração de modelos



## TensorFlow

<https://www.tensorflow.org>

## PyTorch

<https://pytorch.org/>

## Gensim

<https://radimrehurek.com/gensim/>

# Aplicações

# Aplicações no grupo de PLN PUCRS

Classificação de textos (prontuários eletrônicos)

Recuperação de informação em Geologia

Reconhecimento de Entidades nomeadas:

- Extração de informação em investigação policial

- Extração de informação em inteligência de negócios

- Extração de informação em geologia

- Extração de informação em saúde

# Avaliação de modelos

# Tipos de avaliação

Intrínseca

Extrínseca

# Tipos de avaliação

Modelos WE

Diferentes idiomas

Diferentes coberturas:

Cobertura geral

Cobertura em domínios específicos

# Avaliação

Modelos de cobertura geral da língua

Testes de analogia semântica:

- capitais e países comuns
- todas as capitais e países
- países e moedas
- cidades e estados
- gênero em relações familiares (brother, sister) - (irmão, irmã)

# Avaliação intrínseca modelo geral

Cobertura geral

Testes de analogía sintática:

- adjetivo ao advérbio (rápido, rapidamente)
- oposto (possível, impossível)
- comparativo (grande, maior)
- superlativo (easy, easiest) - (fácil, o mais fácil)
- participípio presente (pensar, pensante) - gerúndio (pensar, pensando)

# Avaliação em domínio específico

Avaliações de WE em domínios específicos, ex. área biomédica

Bruni 2014:

- Dois pares são apresentados a especialistas que devem julgar qual par é mais semelhante entre si.
- Cada par é avaliado com outros 50 pares escolhidos randomicamente, que então recebem um score dado pelo número de vezes que o par foi julgado como o mais similar.
- O modelo WE deve reproduzir a intuição de similaridade dos especialistas.

# Avaliação intrínseca em domínio específico

Modelo prévio de domínio Geológico

Farhad (inglês)

Petrobrás (português)

# Avaliação em domínio específico

Par 1			Par 2		
Cientista	Hidrocarboneto		Montmorilonita	Argilomineral	X
Jurássico	Mesozóico	X	Geofísico	Pacote	

# Avaliação

Geração da lista de pares de termos do domínio

Termos:

- Lista gerada por especialistas

- Termos do tesouro

# Lista do especialista

Geofísica	Paleontologia	Ambientes Depositionais	Petrologia Sedimentar	Petrologia Ígnea	Petrologia Metamórfica	Mineralogia
sísmica	plâncton	fluvial	diagênese	biotita	xisto	lâminas
log	diatomáceas	deltaico	sal	quartzo	ardósia	petrografia
sismologia	algas vermelhas	carbonático	carbonato	ortoclásio	gnaisse	minerais
gravimetria	algas verdes	depósito	óleo	lâmina petrográfica	clorita	argila
seção	dinoflagelados	marinho	dolomita	textura	biotita	biotita
perfil	foraminíferos	subambientes	água	mineralogia	silimanita	piroxênio
magnetometria	cianobactérias	plataformas	soterramento	classificação petrográfica	granada	ortopiroxênio
refração	estampa	glacial	porosidade	augita	cianita	quartzo
reflexão	paleocorrente	aluvial	esmectita	hornblenda	feldspato alcalino	plagioclásio
perfilagem		maré	ilita	feldspato alcalino	migmatito	albita
sismoestratigrafia		ondas	arredondado	feldspato potássico		granada
prospecção		costeiro	trama	olivina		zircão
magnetismo		litologia	fábrica	magma		epidoto
seção		sistema	depósito	câmara		arsenopirita
topografia		deposicional	folhelho	magmático		bornita

# Tesauro (skos)

```
<skos:Concept rdf:about="http://bs/#GRAPH">
```

```
  <skos:prefLabel xml:lang="en">GRAPH</skos:prefLabel>
```

```
  <skos:prefLabel xml:lang="pt-BR">Gráfico</skos:prefLabel>
```

```
  <skos:altLabel xml:lang="en">CURVE (GRAPH)</skos:altLabel>
```

```
  <skos:altLabel xml:lang="en">GRAPHICAL CURVE</skos:altLabel>
```

# Tesauro

```
<skos:broader rdf:resource="http://bs/#CHART"/>
```

```
<skos:narrower rdf:resource="http://bs/#AMPLITUDE+CURVE"/>
```

```
<skos:related rdf:resource="http://bs/#CURVING"/>
```

```
<skos:related rdf:resource="http://bs/#DIAGRAM"/>
```

# Geração dos pares

A partir de termos presentes em tesauro

Termos similares - "broader" (termos mais amplo), ou "narrower" (termo menos amplo)

Termos levemente similares - "related" (termo relacionado)

Termos menos similares foram estabelecidos como sendo aqueles sem relação

# Lista de pares

Combinação de pares

3900 entradas, indicação do mais similar

Par 1			Par 2		
Cientista	Hidrocarboneto		Montmorilonita	Argilomineral	X
Jurássico	Mesozóico	X	Geofísico	Pacote	

$$(300 \div 25) \times (25 \times ((25 + 1) \div 2)) \text{ ou } (np \div nr) \times (nr \times ((nr + 1) \div 2))$$

# Avaliação extrínseca

Aplicação dos modelos em tarefas de outra natureza

<b>Tarefa</b>	<b>WE</b>	<b>ELMo</b>	<b>BERT</b>	<b>Flair</b>
<b>REN</b>	✓ <i>Lample, 2016</i>	✓ <i>Perters, 2018</i>	✓ <i>Devlin, 2018</i>	✓ <i>Akbik, 2018</i>
<b>ER</b>	✓ <i>Lin, 2016</i>		✓ <i>Soares, 2019</i>	
<b>COREF</b>	✓ <i>Lee, 2017</i>	✓ <i>Perters, 2018</i>	✓ <i>Joshi, 2019</i>	
<b>CT</b>	✓ <i>Joulin, 2016</i>			✓ <i>Akbik, 2018</i>
<b>SS</b>	✓ <i>Hartmann, 2017</i>		✓ <i>ASSIN 2</i>	

# Limitações

- A semântica refletida no modelo incorpora todas as diferentes nuances do significado mas não as diferencia
  - (e.g. casos em que é sinônimo, casos em que é antônimo)
- Custo computacional
- Caixa preta - ideal é IA explicável (explainable AI)
- Semântica na perspectiva da palavra ou termo (discurso, pragmática)
- Limitação sobre o conhecimento da língua por parte de desenvolvedores

**Obrigada!**

# Referências

**[Akbik, 2018]** Akbik, A., Blythe, D., & Vollgraf, R. (2018). Contextual string embeddings for sequence labeling. In Proceedings of the 27th International Conference on Computational Linguistics (pp. 1638-1649).

**[Bojanowski, 2016]** Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2016). Enriching Word Vectors with Subword Information. arXiv preprint arXiv:1607.04606.

**[Devlin, 2018]** Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

**[Levy, 2014]** Levy, O.; Goldberg, Y. "Dependency-based word embeddings". In: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), 2014, pp. 302-308.

**[Lee, 2017]** Lee, K., He, L., Lewis, M., & Zettlemoyer, L. (2017). End-to-end neural coreference resolution. arXiv preprint arXiv:1707.07045.

# Referências

**[Ling, 2015]** Ling, W., Dyer, C., Black, A., and Trancoso, I. (2015). Two/too simple adaptations of word2vec for syntax problems. In Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.

**[Lin, 2016]** Lin, Y., Shen, S., Liu, Z., Luan, H., & Sun, M. (2016, August). Neural relation extraction with selective attention over instances. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (pp. 2124-2133).

**[Lample, 2016]** Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., & Dyer, C. (2016). Neural architectures for named entity recognition. arXiv preprint arXiv:1603.01360.

**[Mikolov, 2013]** Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G. S.; Dean, J. “Distributed representations of words and phrases and their compositionality”. In: Advances in neural information processing systems, 2013, pp. 3111–3119.

**[Mikolov and Zweig, 2012]** T. Mikolov and G. Zweig. Context dependent recurrent neural network language model. In IEEE SLT, pages 234–239, 2012.

# Referências

**[Mikolov et al., 2010]** T. Mikolov, M. Karafiat, L. Burget, J. Cernocky, and S. Khudanpur. Recurrent neural network based language model. In INTERSPEECH, pages 1045–1048, 2010.

**[Mikolov et al., 2011a]** T. Mikolov, S. Kombrink, L. Burget, J. Cernocky, and S. Khudanpur. Extensions of recurrent neural network language model. In Proc. of IEEE ICASSP, pages 5528–5531, 2011.

**[Mikolov et al., 2012]** T. Mikolov, I. Sutskever, A. Deoras, H. Le, and S. Kombrink. Subword Language Modeling With Neural Networks. 2012.

**[Mikolov et al., 2013]** T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. CoRR, abs/1301.3781, 2013.

**[Pennington, 2014]** Pennington, J., Socher, R., & Manning, C. (2014). Glove: Global vectors for word representation. In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP) (pp. 1532-1543).

# Referências

**[Perters, 2018]** Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (2018). Deep contextualized word representations. arXiv preprint arXiv:1802.05365.

**[Rosenfeld, 200]** Rosenfeld, R. (2000). Two decades of statistical language modeling: Where do we go from here?. Proceedings of the IEEE, 88(8), 1270-1278.

**[Soares, 2019]** Soares, Livio Baldini, et al. "Matching the Blanks: Distributional Similarity for Relation Learning." arXiv preprint arXiv:1906.03158 (2019).

**[Jozefowicz, 2015]** Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. 2016. Exploring the limits of language modeling. CoRR abs/1602.02410.

**[Joulin, 2016]** Joulin, A., Grave, E., Bojanowski, P., & Mikolov, T. (2016). Bag of tricks for efficient text classification. arXiv preprint arXiv:1607.01759.

# Referências

**[Joshi, 2019]** Joshi, M., Levy, O., Weld, D. S., & Zettlemoyer, L. (2019). BERT for Coreference Resolution: Baselines and Analysis. arXiv preprint arXiv:1908.09091.

**[Zhang, 2015]** Zhang, D.; Xu, H.; Su, Z.; Xu, Y. "Chinese comments sentiment classification based on word2vec and svm perf", Expert Systems with Applications, vol. 42-4, 2015, pp.1857-1863.