

# Arquitetura Transformer e seus derivados: BERT e GPT

Redes Neurais e Aprendizado Profundo

Moacir A. Ponti

[www.icmc.usp.br/~moacir](http://www.icmc.usp.br/~moacir) — [moacir@icmc.usp.br](mailto:moacir@icmc.usp.br)

São Carlos-SP/Brasil

# Agenda

Transformer Network

BERT: pré-treinamento de encoders de transformers bidirecionais

GPT: generative pretrained transformer

# Agenda

Transformer Network

BERT: pré-treinamento de encoders de transformers bidirecionais

GPT: generative pretrained transformer

# RNNs vs Transformer Networks

## RNNs

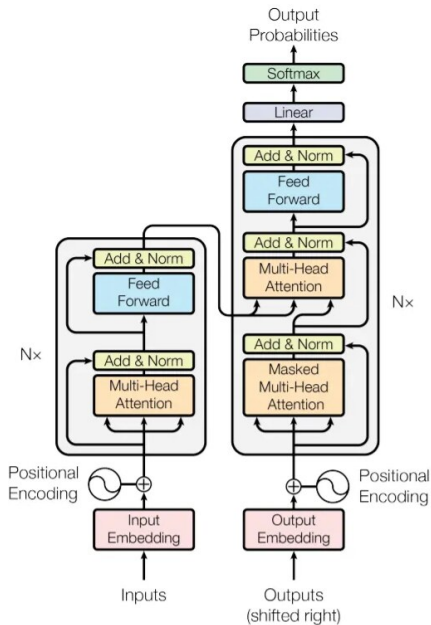
- ▶ podem não funcionar com dependências longas
- ▶ recorrência dificulta computação paralela (pontos da sequência não podem ser processados em paralelo)
- ▶ podem sofrer com explosão ou desaparecimento de gradiente

## Transformer Networks

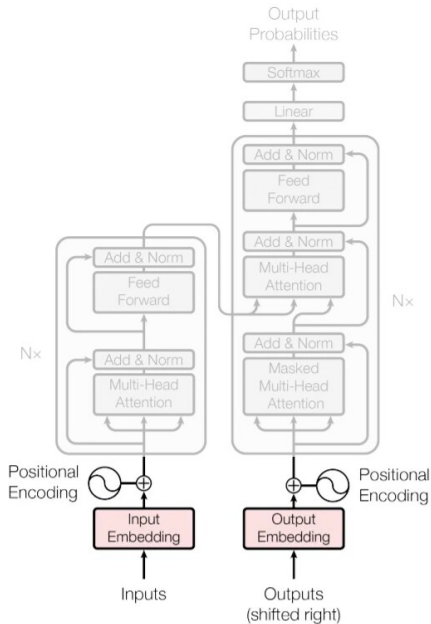
- ▶ não possui recorrência, apenas atenção
- ▶ facilita capturar dependências longas
- ▶ construído para ser sequence-to-sequence

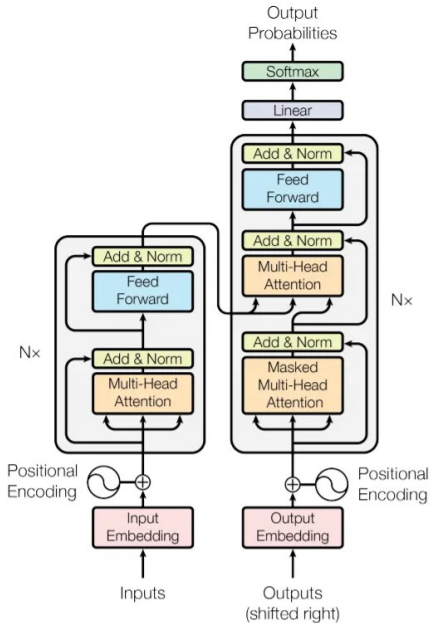
Artigo: "Attention is all you need" Vaswani, NeurIPS 2017.

# Arquitectura Transformer

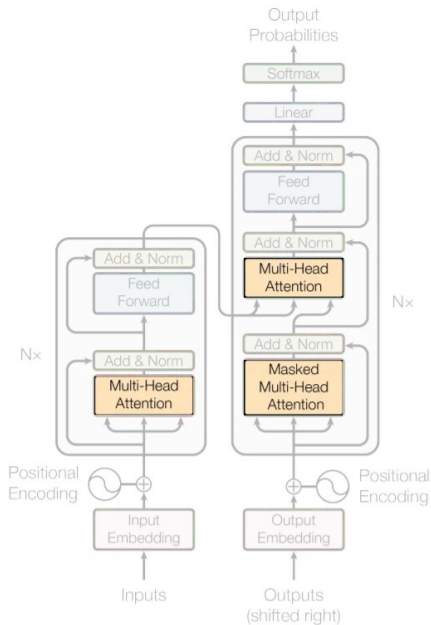


# Positional Encoding: adiciona informação de posicionamento



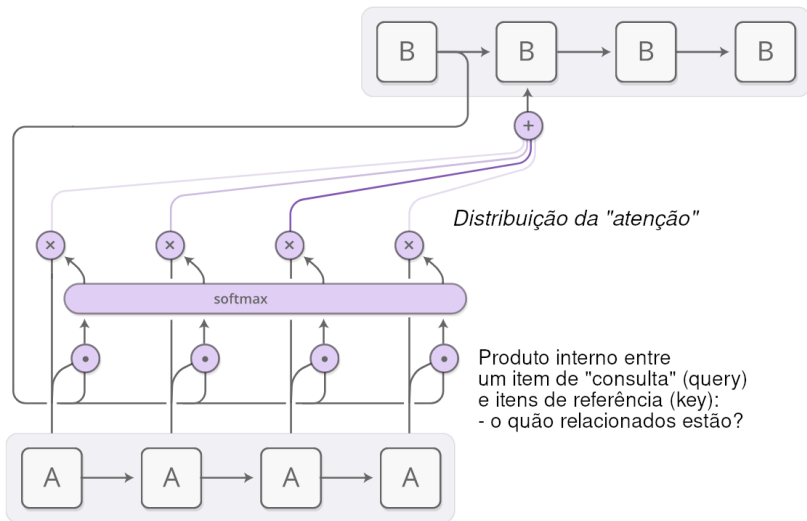


# Arquitetura Transformer: multi-atenção





# Distribuição da atenção com relação ao estado atual



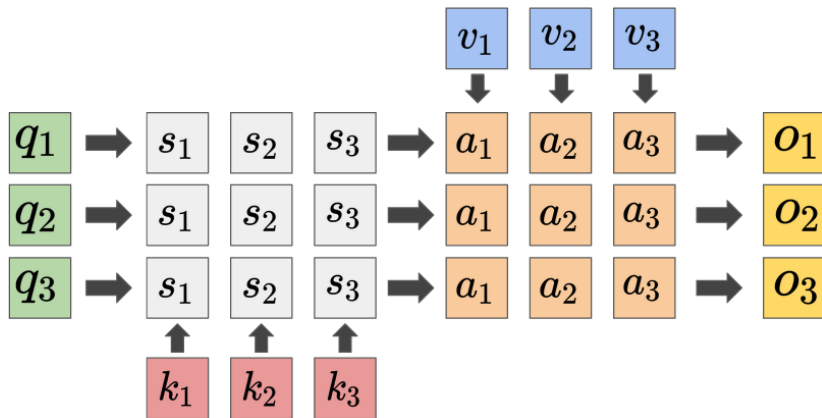
# Mecanismo de atenção em Transformer Networks

- ▶ Recuperar um valor  $v_i$  para uma consulta/query  $q$  baseada numa chave/key  $k_i$

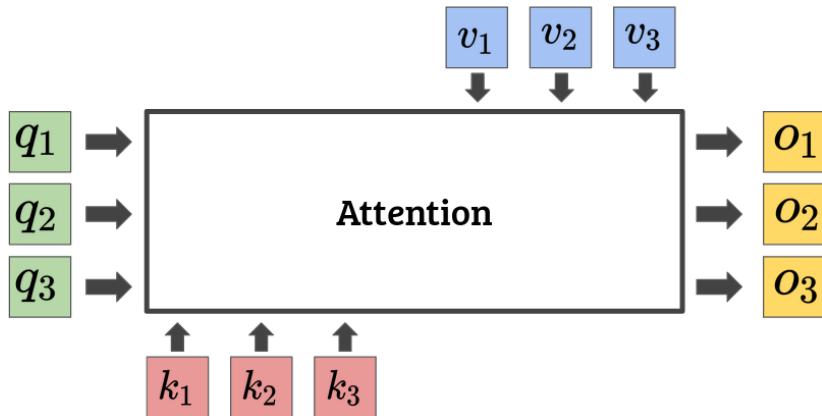
$$Attention(q, k, v) = \sum_i similarity(q, k_i) \times v_i$$

- ▶ A **similaridade** entre uma consulta e todas as chaves, ponderadas pelos valores
- ▶ Somar ao longo de todas as chaves/valores, produz uma **distribuição** de pesos relacionando consulta e todos os valores

# Atenção

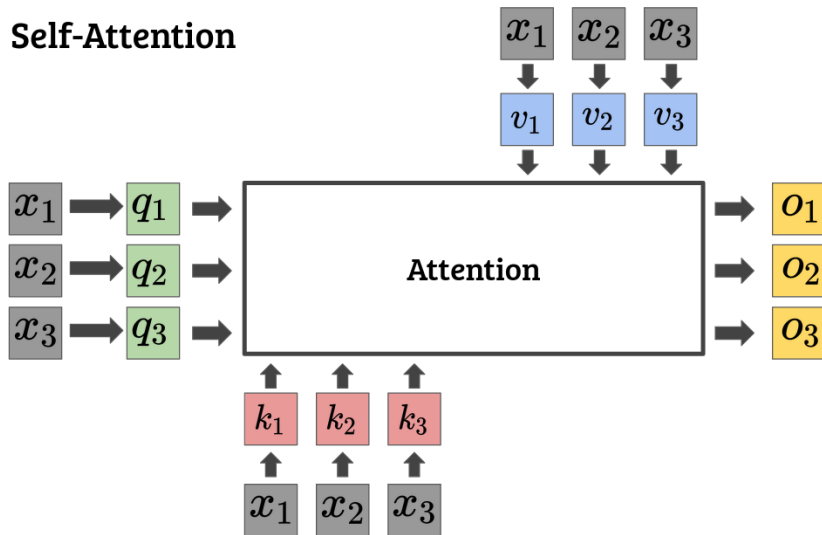


# Atenção



# Atenção (auto-atenção)

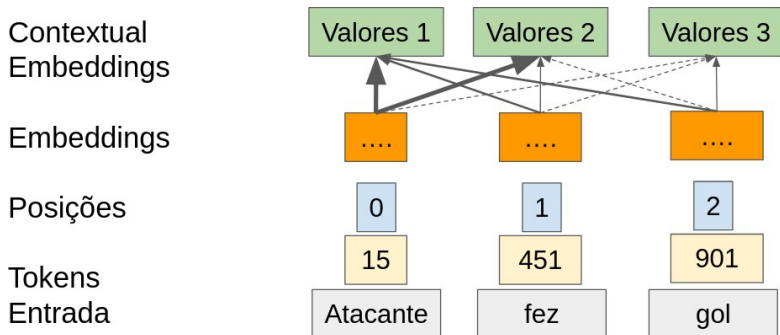
## Self-Attention



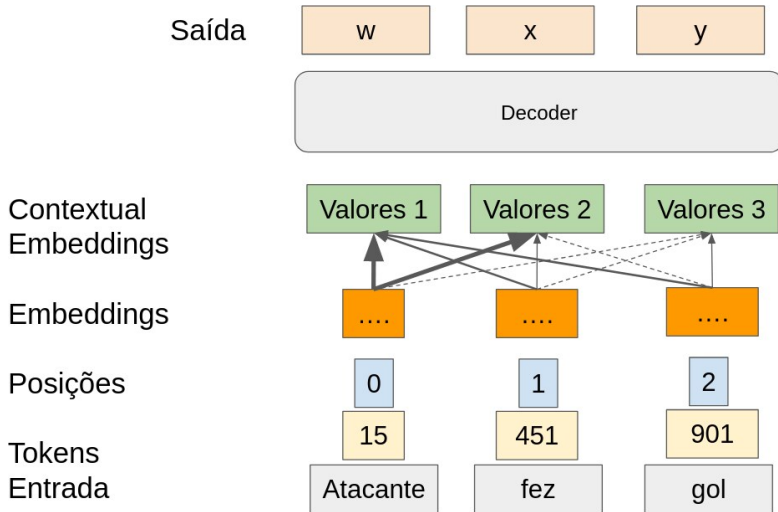
## Para texto: representações para cada palavra

1. Word embeddings:  $w$
2. Word embeddings + Positional encoding:  $w + p = e$
3. Embeddings contextuais com base na saída da atenção

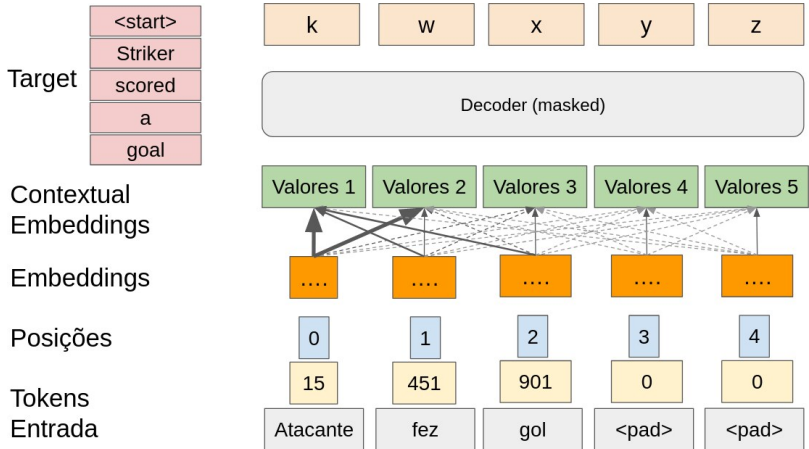
# Encoder

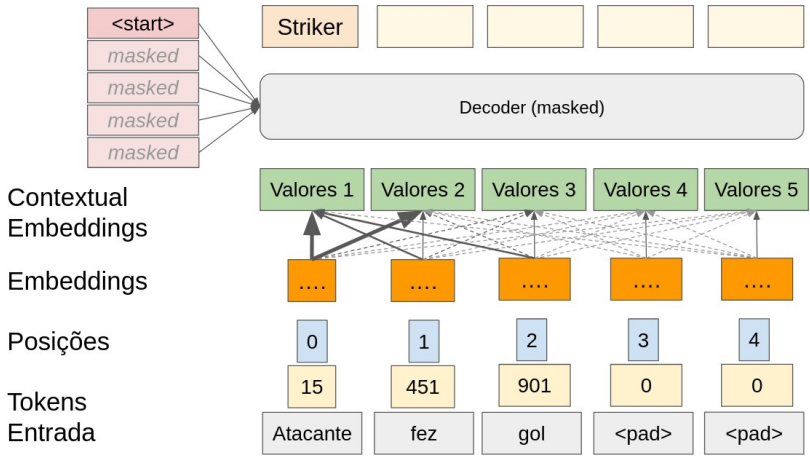


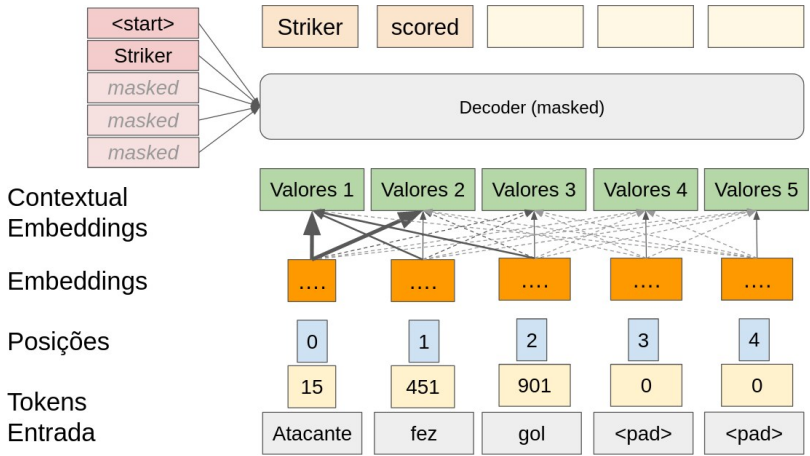
# Encoder-Decoder

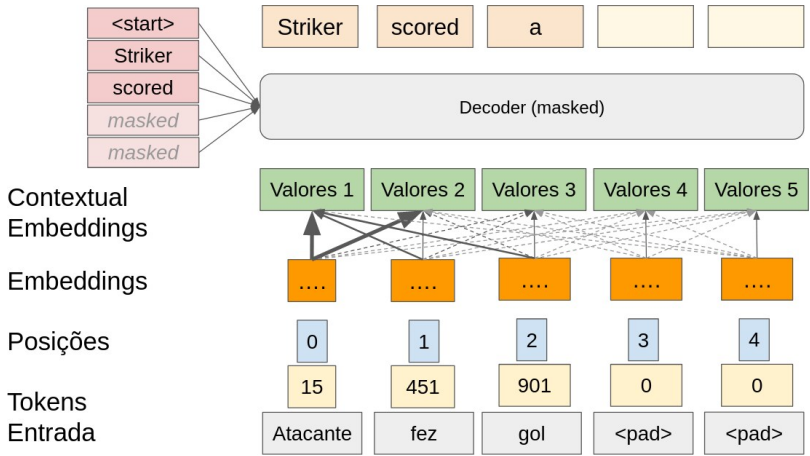


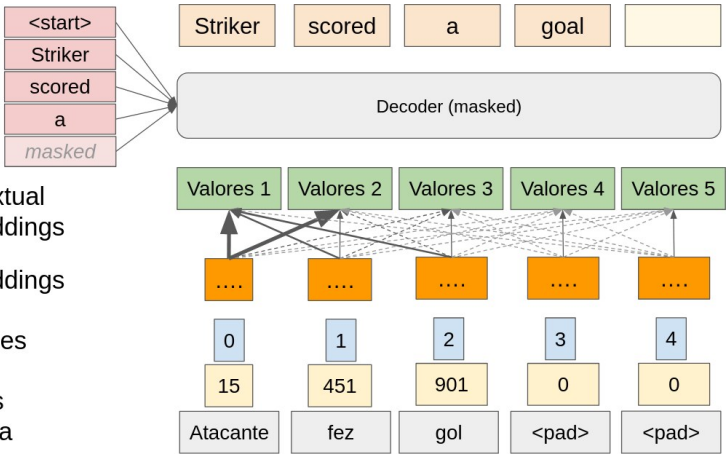


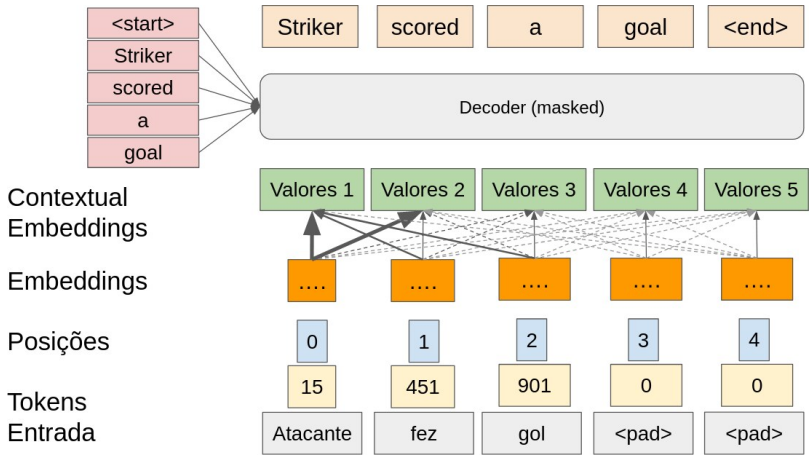




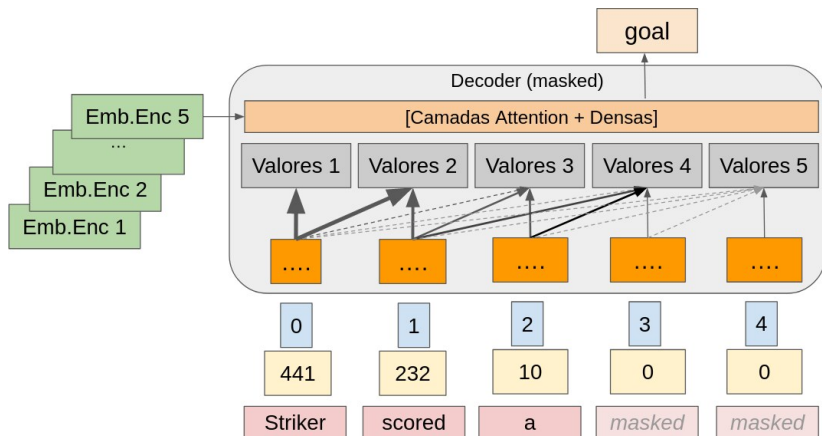








# Visão do Decoder



# Attention e Masked Attention

## Attention

$$\text{attention}(Q, K, V) = \text{softmax} \left( \frac{Q^T K}{\sqrt{d_k}} \right) V$$

$d_k$  dimensões da chave

## Masked Attention

- ▶ o decoder deve anular atenção com relação a palavras de entrada futuras, senão não há aprendizado

$$\text{maskedattention}(Q, K, V) = \text{softmax} \left( \frac{Q^T K + M}{\sqrt{d_k}} \right) V$$

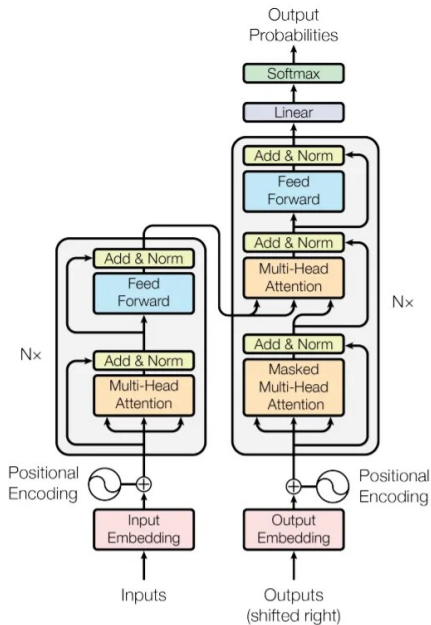
$M$  é uma matriz com  $-\infty$  nas posições de palavras futuras



## Representações para cada palavra

1. Word embeddings:  $w$
2. com positional encoding:  $w + p = e$
3. Query:  $Q_h = W_Q \cdot e_h$
4. Key:  $K_h = W_K \cdot e_h$
5. Value:  $V_h = W_V \cdot e_h$
6. Attention heads:  $Z_h = \text{softmax} \left( \frac{Q_h^T K_h + M}{\sqrt{d_k}} \right) V_h$
7. Attention output:  $Z = W_O \text{concat}((Z_1, Z_2, \dots, Z_H))$

# Arquitectura Transformer



# Agenda

Transformer Network

BERT: pré-treinamento de encoders de transformers bidirecionais

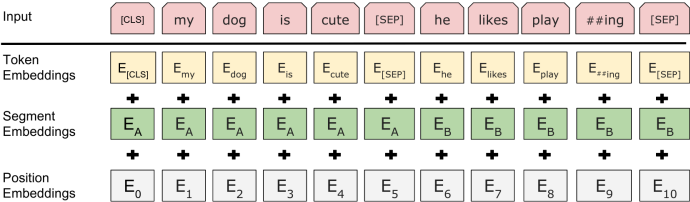
GPT: generative pretrained transformer

Método para **pré-treinar** encoders to tipo Transformer

— BERT Base e BERT Large, modelos com blocos Transformer:

- ▶ Base:
  - ▶ 12 camadas (blocos Transformer)
  - ▶ Embedding com 768 dimensões
  - ▶ 110 milhões de parâmetros
- ▶ Large:
  - ▶ 24 camadas (blocos Transformer)
  - ▶ Embedding com 1024 dimensões
  - ▶ 336 milhões de parâmetros

# BERT Treinamento e Embeddings

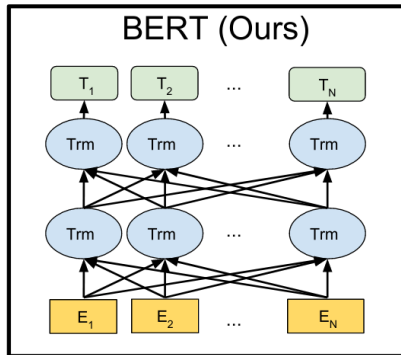


# BERT: pré-treinamento inspirado em ELMo

- ▶ GloVe e similares
  - ▶ vetor fixo por palavra independente do contexto
  - ▶ ex. manga, bateria, pilha.
- ▶ Olha para toda a sentença antes de atribuir vetor
  - ▶ aprende (sem labels) a prever a próxima palavra (e a anterior)

vamos deixar o menino jogar bola  
→ vamos deixar o menino jogar [mask]  
[mask] deixar o menino jogar bola ←  
**vamos deixar o → [mask] ←jogar bola**

# BERT Bidirectional Transformer



# BERT: segunda tarefa de pré-treinamento

**Input** = [CLS] the man went to [MASK] store [SEP]  
he bought a gallon [MASK] milk [SEP]

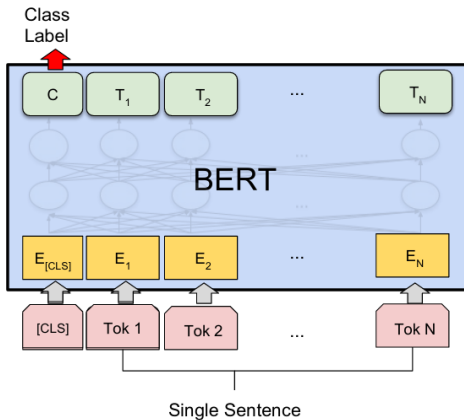
**Label** = IsNext

**Input** = [CLS] the man [MASK] to the store [SEP]  
penguin [MASK] are flight ##less birds [SEP]

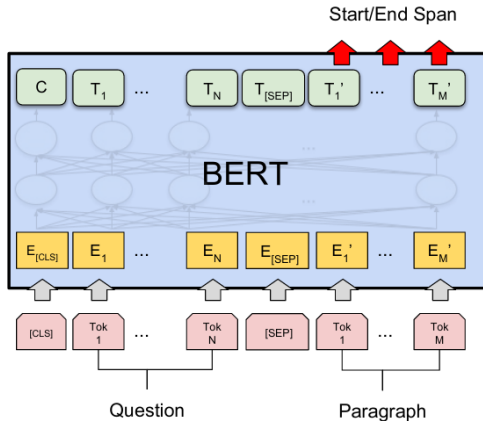
**Label** = NotNext



# BERT: para classificação



# BERT: para encontrar resposta em texto



# Agenda

Transformer Network

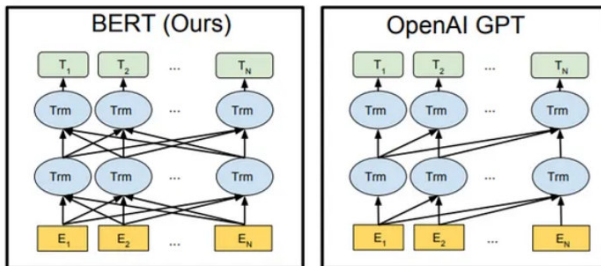
BERT: pré-treinamento de encoders de transformers bidirecionais

GPT: generative pretrained transformer

# GPT e Transformer

Método para **pré-treinar** decoders to tipo Transformer

Modelo auto-regressivo: usa a saída atual como entrada



O atacante fez um

$$p(\text{"gol"} | \text{"O atacante fez um"}) = 0.203$$

$$p(\text{"jogo"} | \text{"O atacante fez um"}) = 0.196$$

$$p(\text{"golaço"} | \text{"O atacante fez um"}) = 0.141$$

$$p(\text{"pedido"} | \text{"O atacante fez um"}) = 0.105$$

$$p(\text{"comentário"} | \text{"O atacante fez um"}) = 0.083$$

$$p(\text{"chute"} | \text{"O atacante fez um"}) = 0.070$$

$$p(\text{"bolo"} | \text{"O atacante fez um"}) = 0.009$$

O atacante fez um gol **na partida de ontem**

# Temperatura

- O atacante fez um gol na partida de ontem
- O atacante fez um gol invalidado pelo VAR
- O atacante fez um gol de fora da área
- O atacante fez um jogo ruim no segundo tempo
- O atacante fez um pedido à torcida

## GPT 3 e superiores

1. Coletar grandes quantidades de dados
2. Pré-treinar com bases massivas para gerar tokens (auto-supervisionado)
3. Ajustar com tarefas específicas e controladas (ex. perguntas e respostas): fine-tuning
4. Realizar ajuste com aprendizado por reforço usando humanos (RLHF - Reinforcement Learning Human Feedback)



1. Word / sentence embeddings: espaços com coerência e contexto entre textos
2. Transformer: modelo sequence-to-sequence
3. BERT: modelo que olha todo o texto para gerar a predição, possível de ajustar para diferentes tarefas alvo
4. GPT: modelo autoregressivo para geração de texto

# References

- ▶ Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin. Attention is all you need, NeurIPS 2017
- ▶ Dzmitry Bahdanau, KyungHyun Cho and Yoshua Bengio. Neural Machine Translation By Jointly Learning To Align And Translate. ICLR 2015.
- ▶ A. Karpathy. Understanding LSTM Networks.  
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- ▶ C. Olah. Understanding LSTM Networks  
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>