

Trilha AI / ML

Quando ter atenção é melhor que ter memória?



Lúcio Sanchez Passos
Data Science Manager, Santander



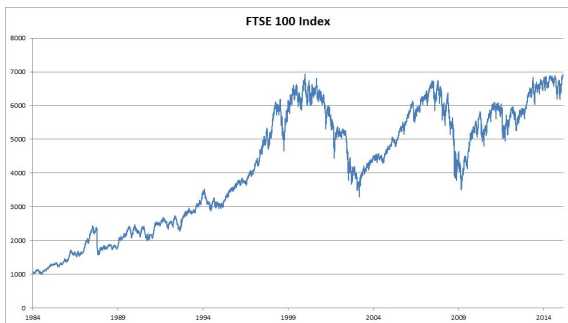
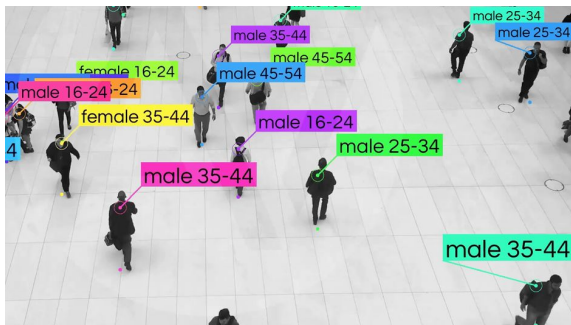
Leonardo Piedade
Solutions Architect, AWS

Agenda



- Background: Sequences
- Recurrent Process Units: RNN, LSTM and GRU
- Seq2Seq Overview
- Attention & Transformers: How, When, and Why?
- Demo

Sequences - When order matters!

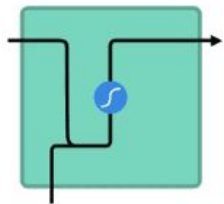


**THIS IS A TEST
IF YOU CAN
READ THIS
ALL THE WAY
DOWN TO HERE**

PLEASE
DO NOT READ

Recurrent Process Units

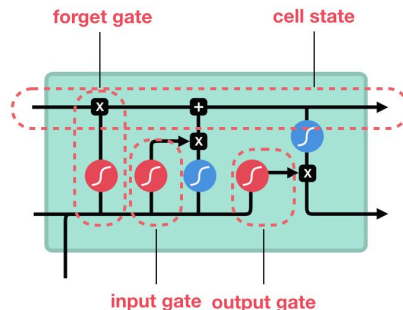
Recurrent Neural Network (RNN)



Good at Modeling Sequence Data

Short-Term Memory Problem

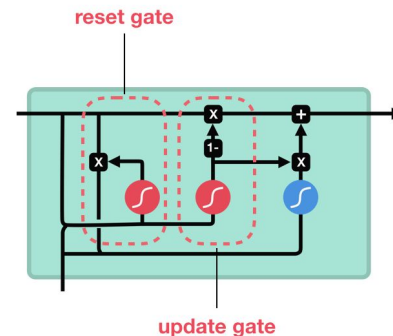
Long Short-Term Memory (LSTM)



Short-Term Memory Problem

More Complex Training Process

Gated Recurrent Units (GRU)



More Complex Training Process



sigmoid



tanh



pointwise
multiplication



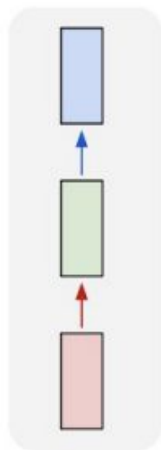
pointwise
addition



vector
concatenation

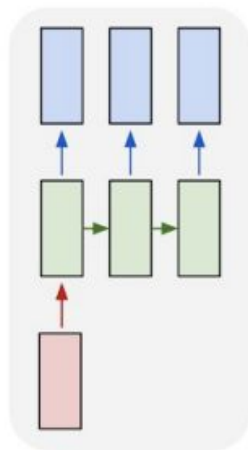
Applications of RNNs

One-to-one



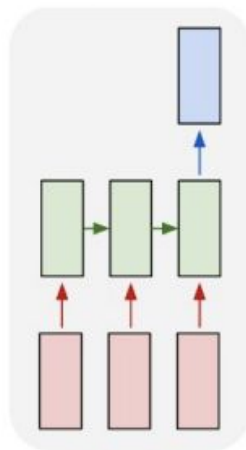
Object
Classification

One-to-many



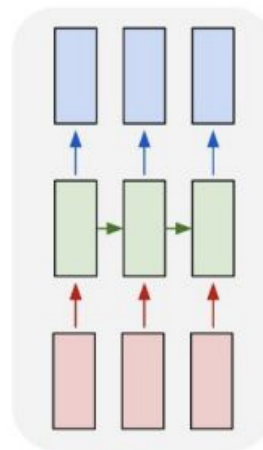
Music
generation

Many-to-one



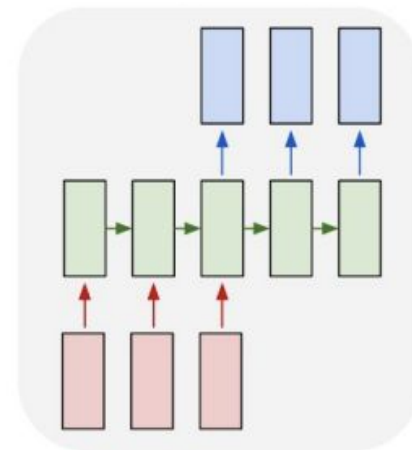
Sentiment
analysis

Many-to-many



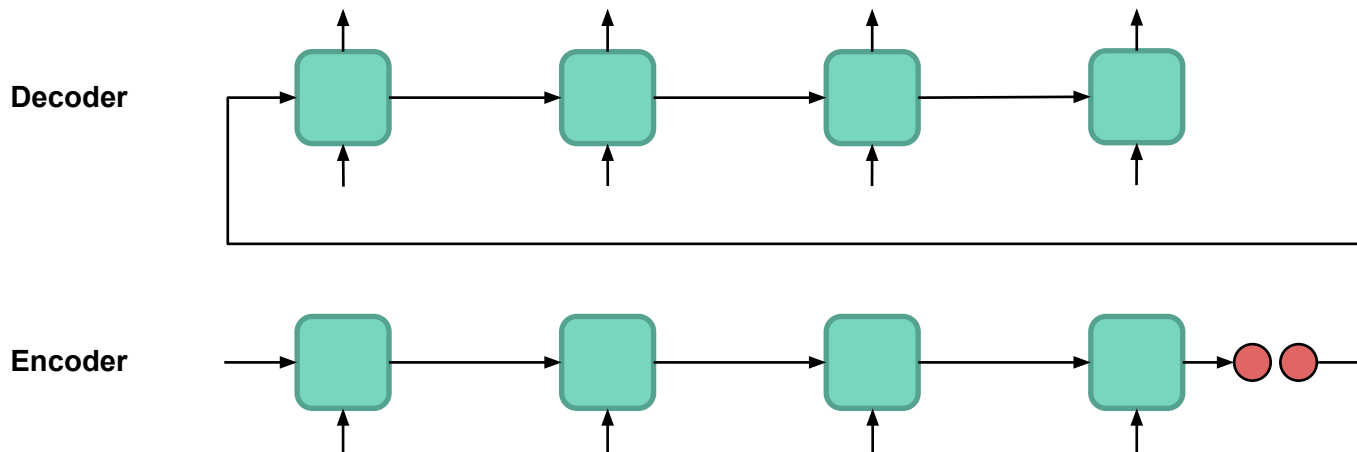
Name entity
recognition

Many-to-many

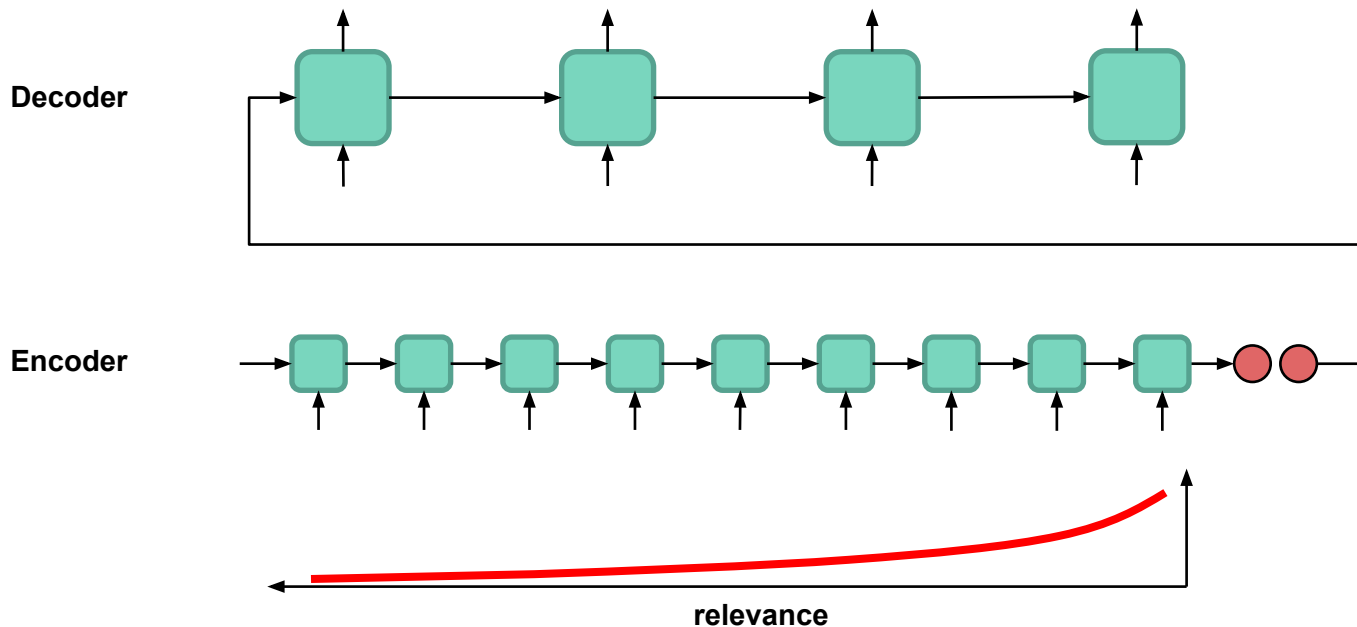


Machine
translation

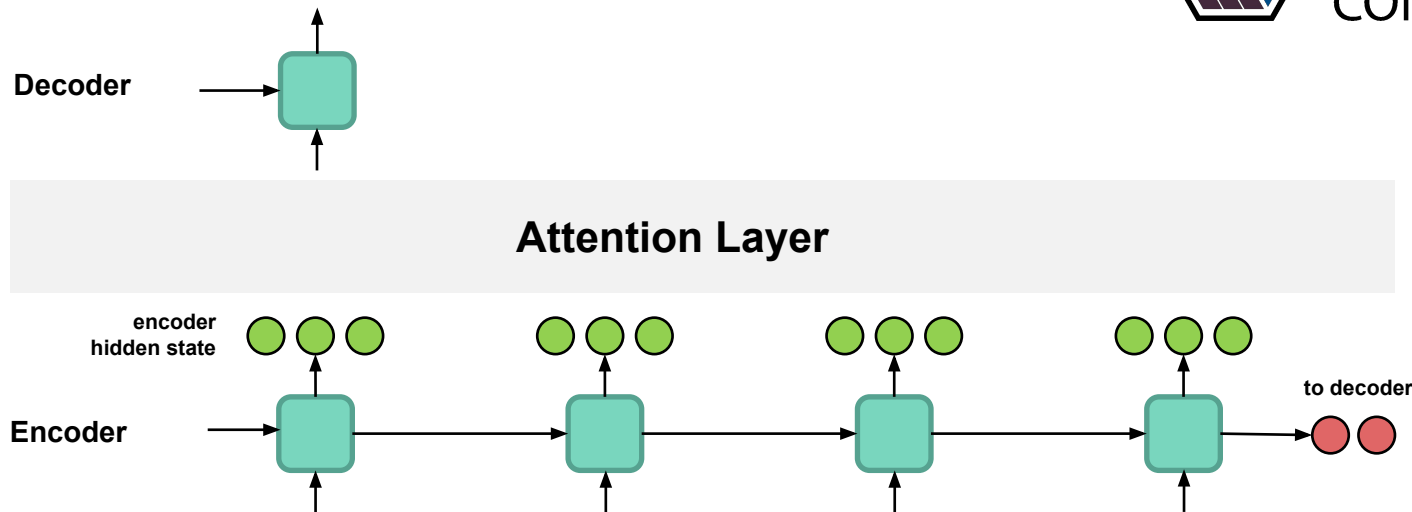
Seq2Seq (Many-to-many)



Seq2Seq – Bottleneck Problem

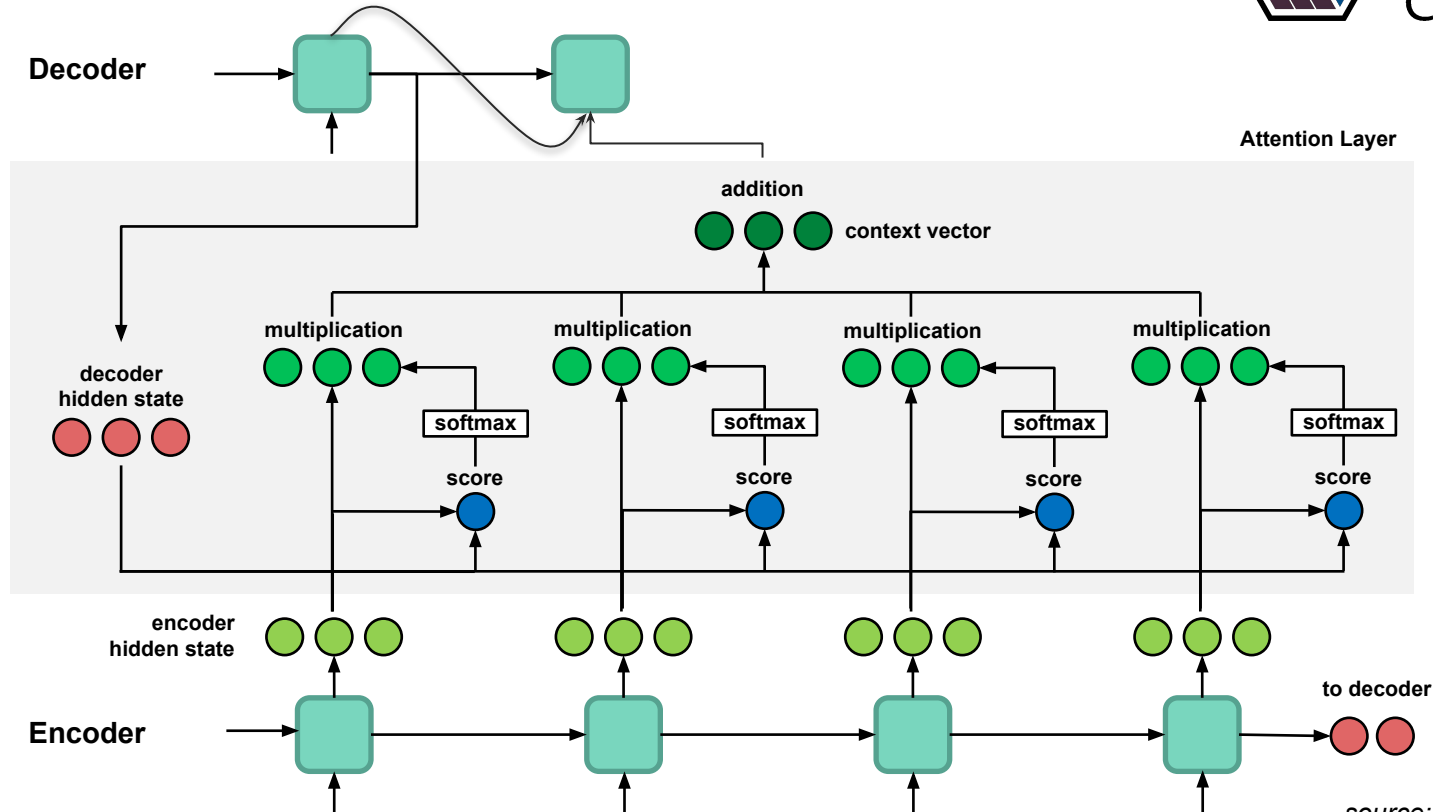


Attention – Definition



“Attention is an interface between the encoder and decoder that provides the decoder with information from every encoder hidden state”

Attention Mechanism



Attention is great...



- Attention significantly improves performance (in many applications)
- Attention solves the bottleneck problem
- Attention helps with vanishing gradient problem
- Attention provides some interpretability

Seq2Seq + Attention - Drawback



- Sequential computation of data prevents parallelism
- Even with LSTM/GRU + Attention, the gradient vanishing problem is not completely solved

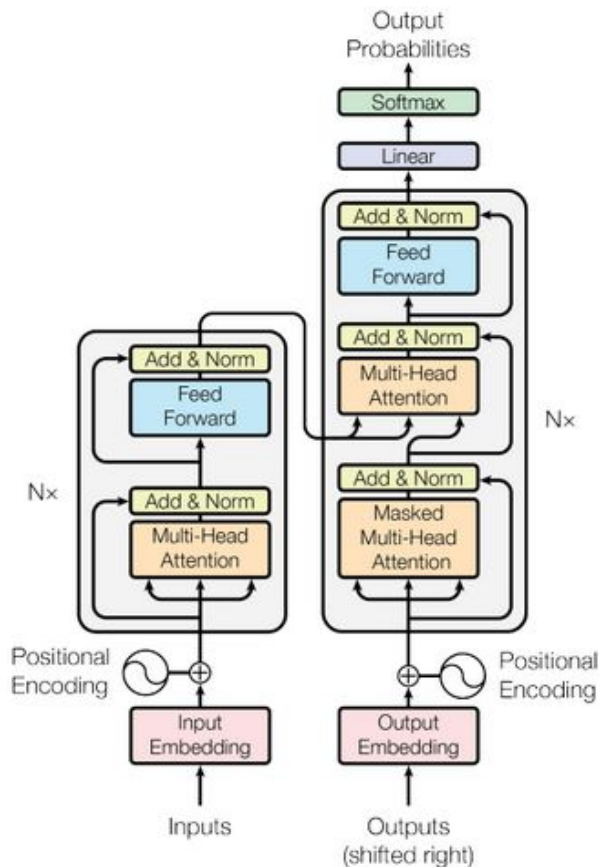
But if we have all states with Attention...why use RNN?

Transformers



THE
DEVELOPER'S
CONFERENCE

Encoder →



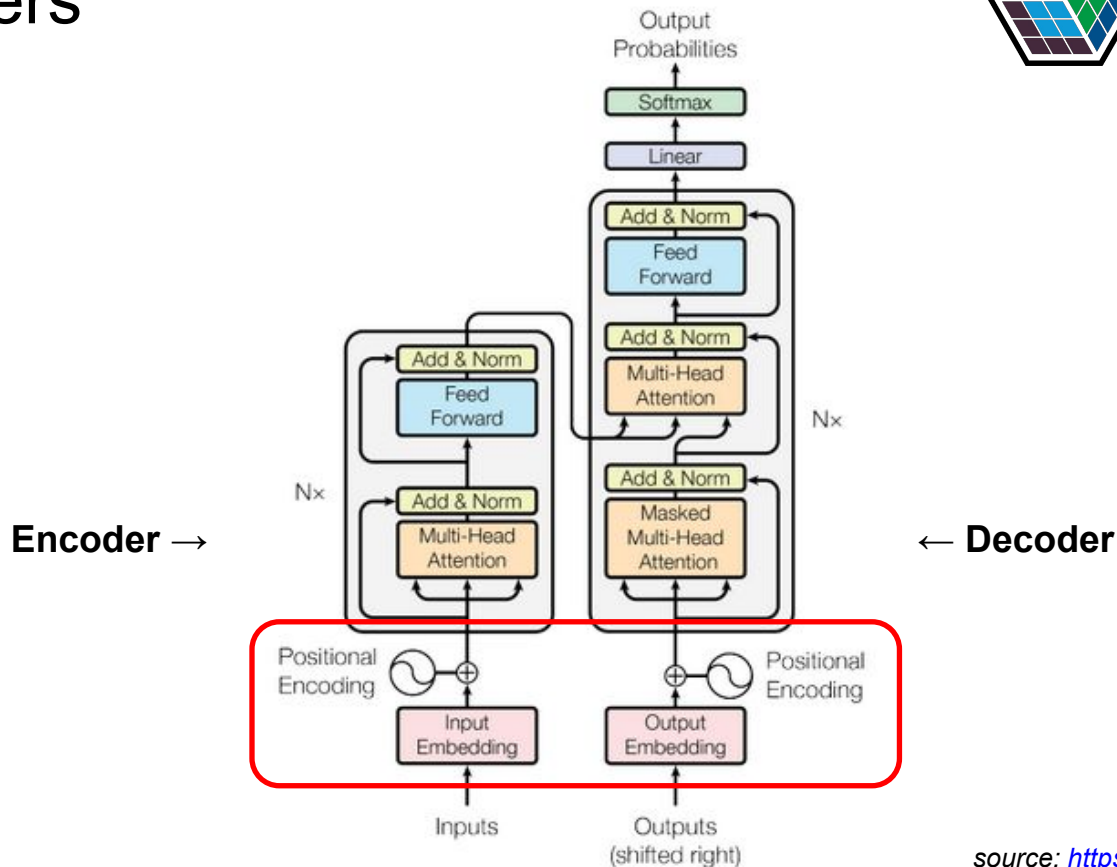
← Decoder

source: <https://arxiv.org/pdf/1706.03762.pdf>

Transformers



THE
DEVELOPER'S
CONFERENCE



source: <https://arxiv.org/pdf/1706.03762.pdf>

Positional Encoding



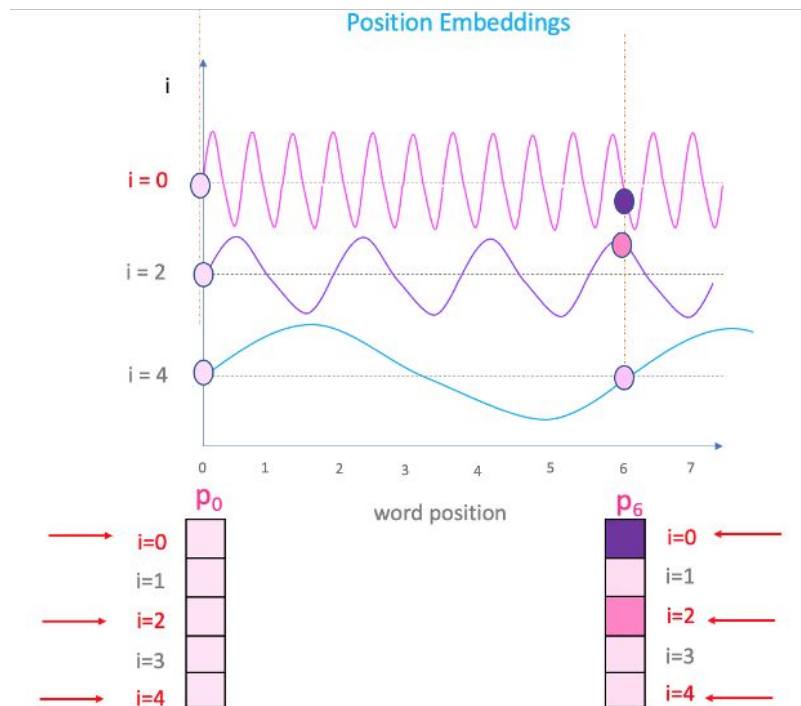
THE
DEVELOPER'S
CONFERENCE

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

pos – position in the sequence

d – size of token vector

i – position in the tokenized vector

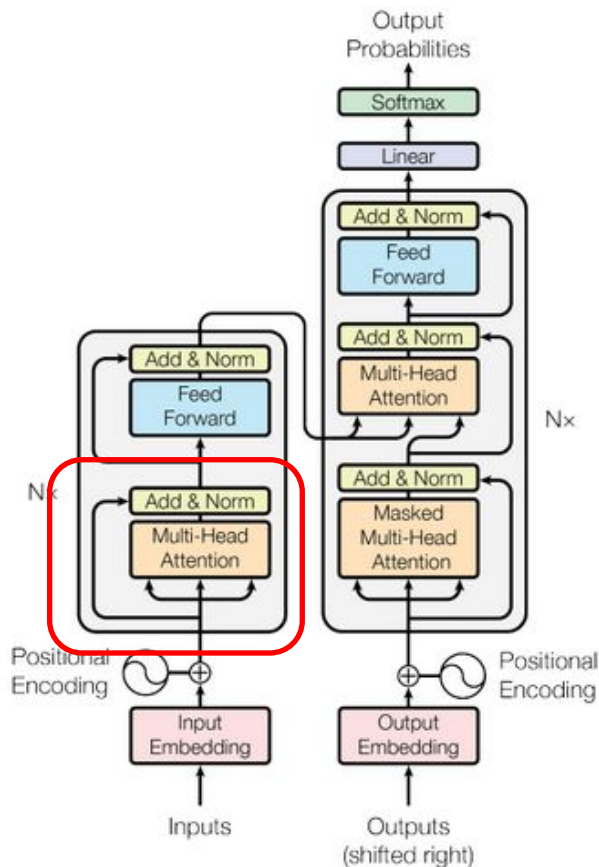


Transformers



THE
DEVELOPER'S
CONFERENCE

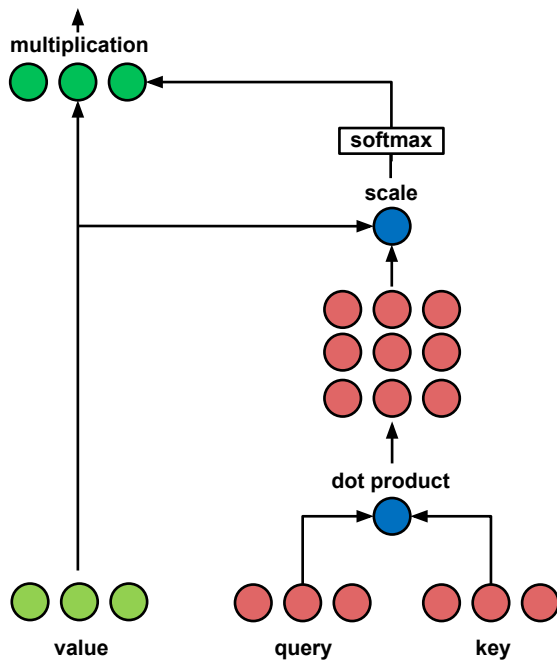
Encoder →



← Decoder

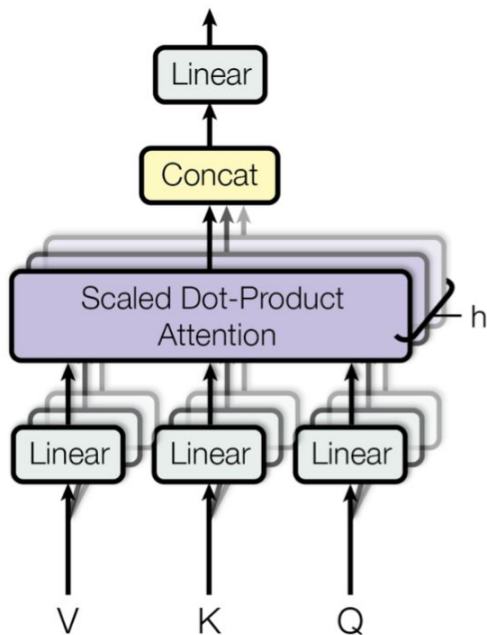
source: <https://arxiv.org/pdf/1706.03762.pdf>

Self-Attention



Self-attention measures the relevance of interaction among all inputs.

Multi-headed Attention



“Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions.”

Transformers Summary



- Easier to train (parallel training)
- No gradient vanishing and explosion
- Allows Transfer Learning

Demo...

Original Papers and Presentations...



- [Attention Is All You Need](#)
- [Long Short-Term Memory](#)
- [Attn: Illustrated Attention](#)
- [Illustrated Guide to Transformers](#)
- [Attentional Neural Network Model](#)
- [Transcoder: Facebook's Unsupervised Programming Language Translator](#)

Quando ter atenção é melhor que ter memória?



Obrigado!



linkedin.com/in/luciopassos/



linkedin.com/in/leoap/



THE DEVELOPER'S CONFERENCE