

Transformers for Natural Language Processing

Concepts, Examples & Tutorials

Session Agenda



Introduction



Transfer Learning
Brief



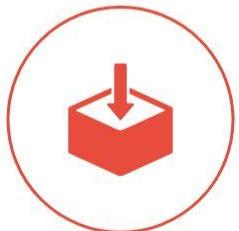
Deep Transfer
Learning for NLP



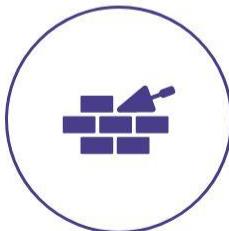
Deep Learning
Architecture
Fundamentals



Motivation for
Transformers



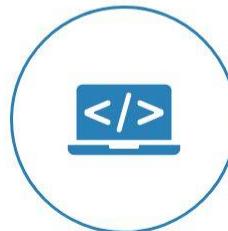
Transformer
Essentials



Transformer
Model
Architectures



Understanding
BERT & DistilBERT



Hands-on Tutorials



Introduction

About Me

Dipanjan Sarkar

Data Science Lead, Author, Google Developer Expert - ML



APPLIED
MATERIALS®

 Springboard



 Experts
Machine Learning

Slide Deck and Hands-on Examples

https://bit.ly/transformers_ds

Transfer Learning Brief



Transfer Learning

The ability to utilize knowledge learnt in prior tasks to new and novel tasks

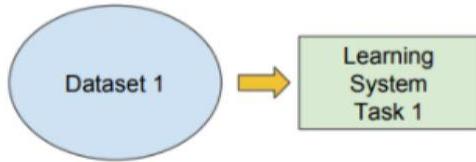


- Know how to ride a motorbike ➔ Learn how to ride a car
- Know how to play classic piano ➔ Learn how to play jazz piano
- Know math and statistics ➔ Learn machine learning

Traditional ML vs. Transfer Learning

Traditional ML

- Isolated, single task learning:
 - Knowledge is not retained or accumulated. Learning is performed w.o. considering past learned knowledge in other tasks



vs

Transfer Learning

- Learning of a new tasks relies on the previous learned tasks:
 - Learning process can be faster, more accurate and/or need less training data



Defining Transfer Learning

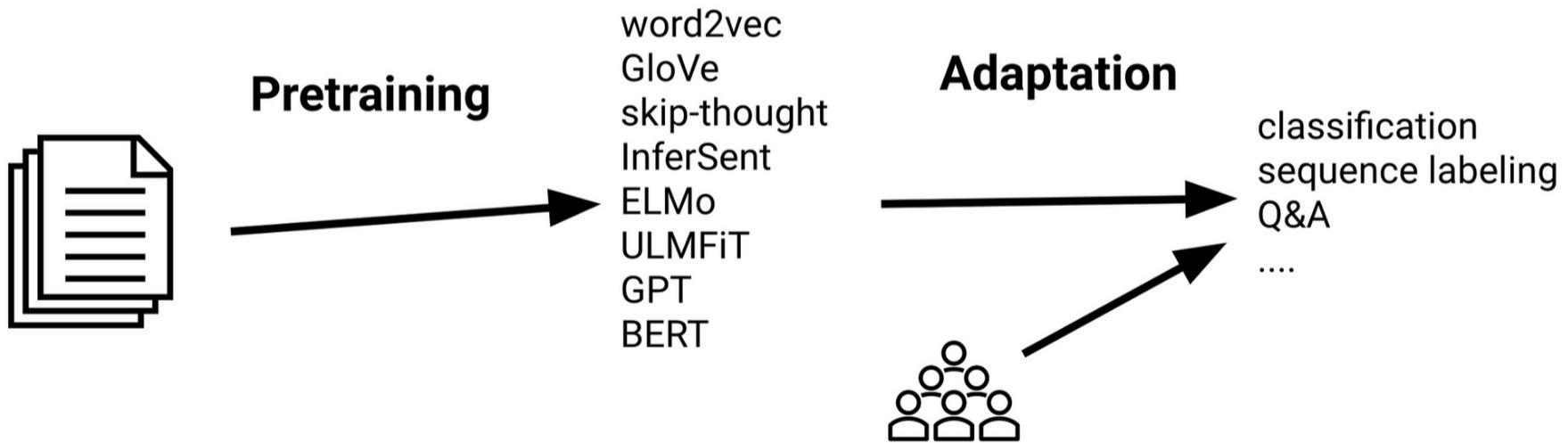
Given a source domain \mathcal{D}_S , a corresponding source task \mathcal{T}_S , as well as a target domain \mathcal{D}_T and a target task \mathcal{T}_T , the objective of transfer learning now is to enable us to learn the target conditional probability distribution $P(Y_T|X_T)$ in \mathcal{D}_T with the information gained from \mathcal{D}_S and \mathcal{T}_S where $\mathcal{D}_S \neq \mathcal{D}_T$ or $\mathcal{T}_S \neq \mathcal{T}_T$. In most cases, a limited number of labeled target examples, which is exponentially smaller than the number of labeled source examples are assumed to be available.

Deep Transfer Learning for NLP



Sequential Transfer Learning for NLP

Sequential transfer learning has led to the biggest improvements so far in the field of NLP



Sequential Transfer Learning for NLP

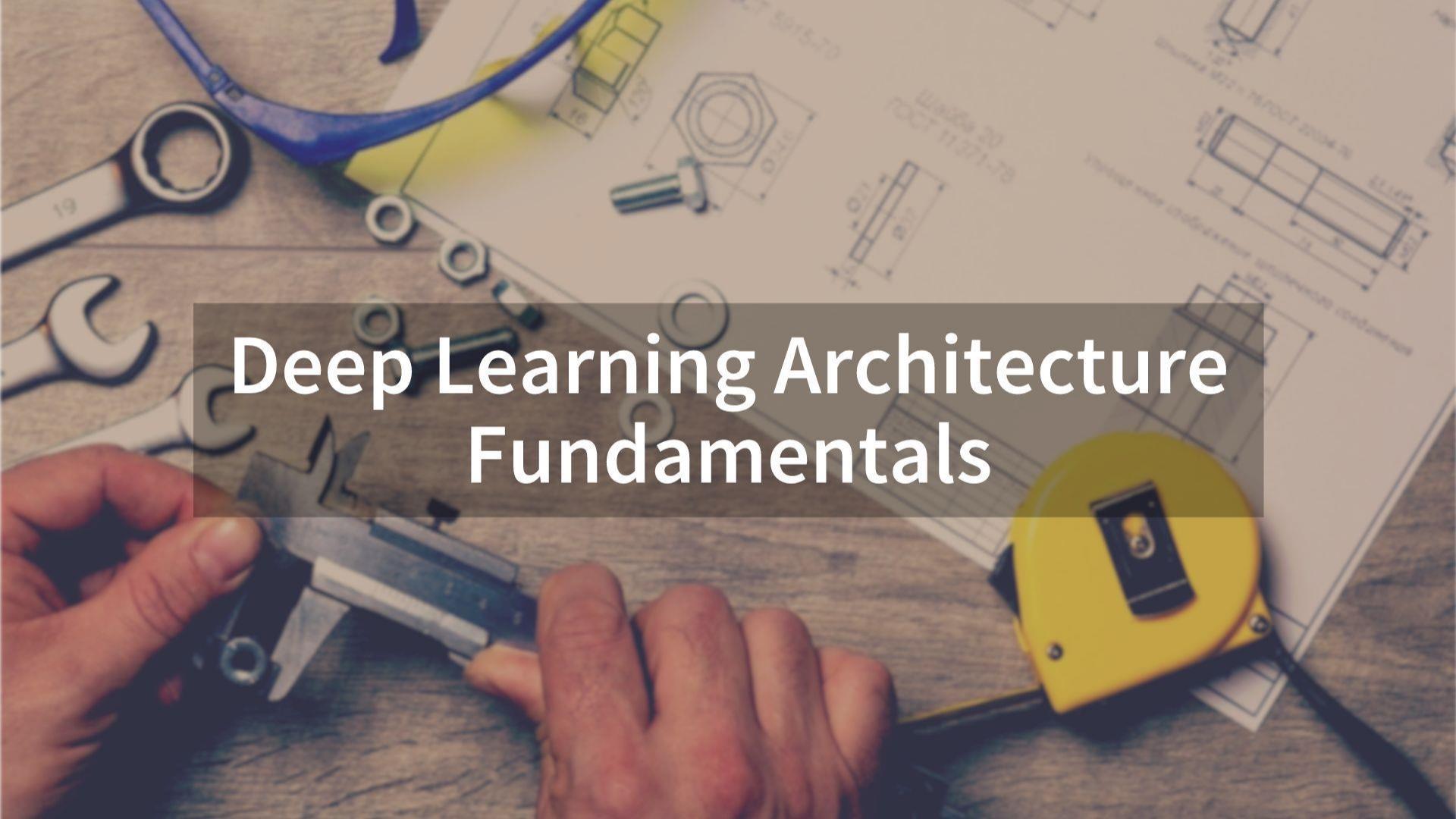
1 Pre-training

- Word and context-based representations (word2vec, fasttext, elmo)
- Simple Language Models (RNN, LSTMs, Bi-LSTMs)
- Generalized Language Modeling (ULMFiT, GPT, BERT)

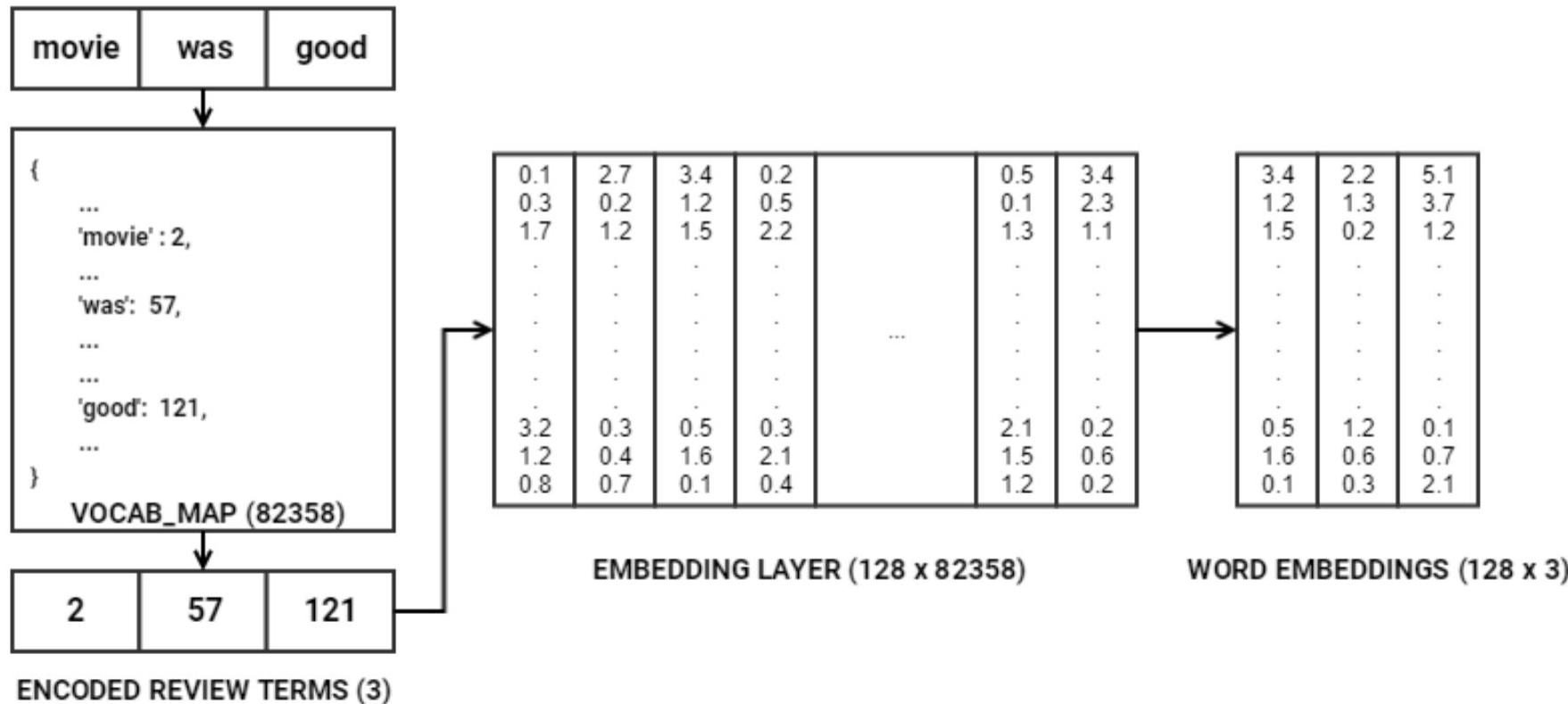
2 Adaptation

- Keep the model architecture and weights unchanged
- Change and fine-tune model architecture and weights

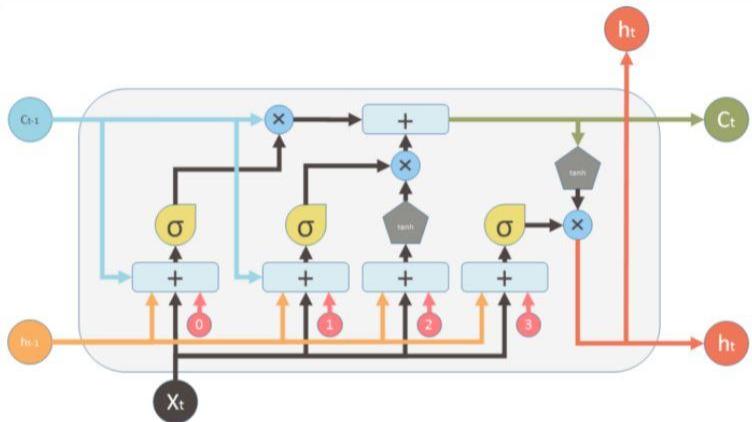
Deep Learning Architecture Fundamentals



Embedding Layer

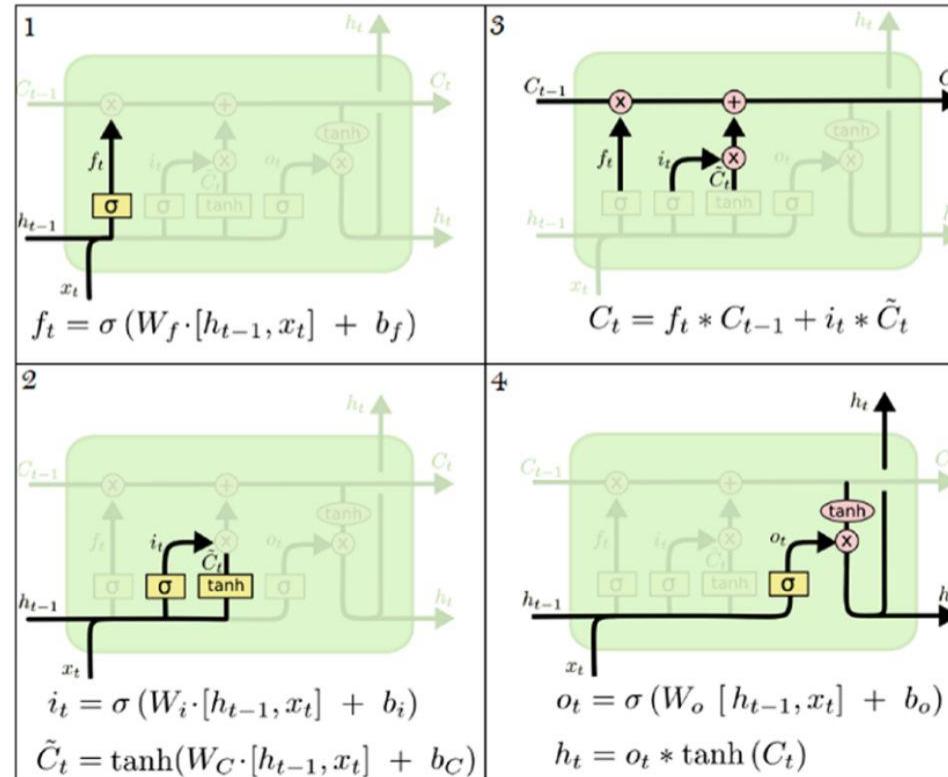


Sequential Models - LSTM

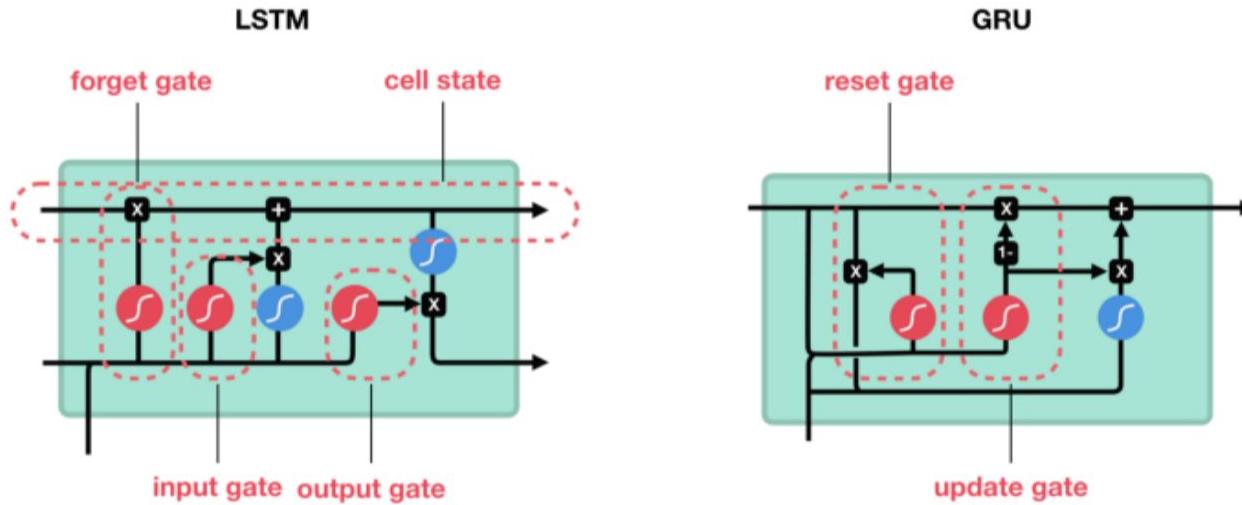


- LSTMs can remember longer sequences of data than RNNs
- The network takes three inputs
 - X_t is the input of the current time step
 - h_{t-1} is the output from the previous LSTM unit
 - C_{t-1} is the “memory” or cell state of the previous unit
- The forget gate decides what is relevant to retain from previous steps
- The input gate decides what information can be added from the current step
- The output gate decides what the next hidden state should be

LSTM Cell Operations



LSTM vs GRU



sigmoid



tanh



pointwise
multiplication



pointwise
addition

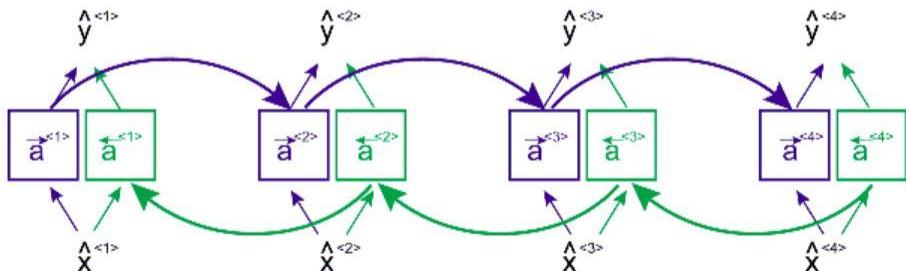


vector
concatenation

The need for Bi-directional LSTMs

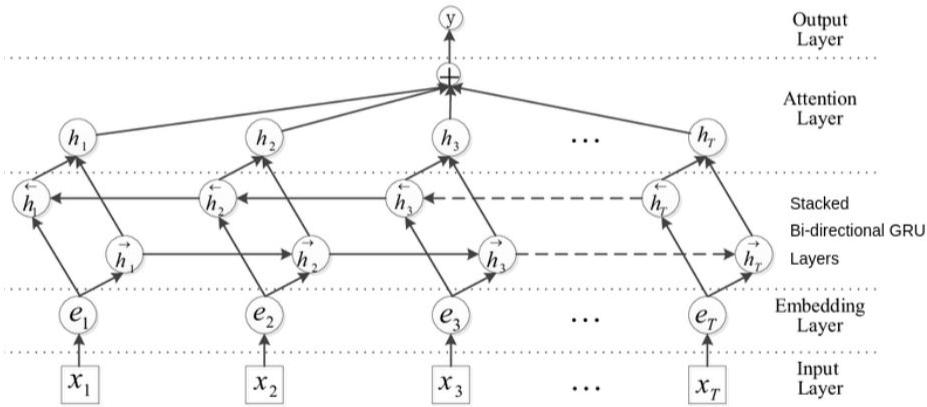
He said , "Teddy bears are on sale!"
not part of person name

He said , "Teddy Roosevelt was a great President !"
part of person name



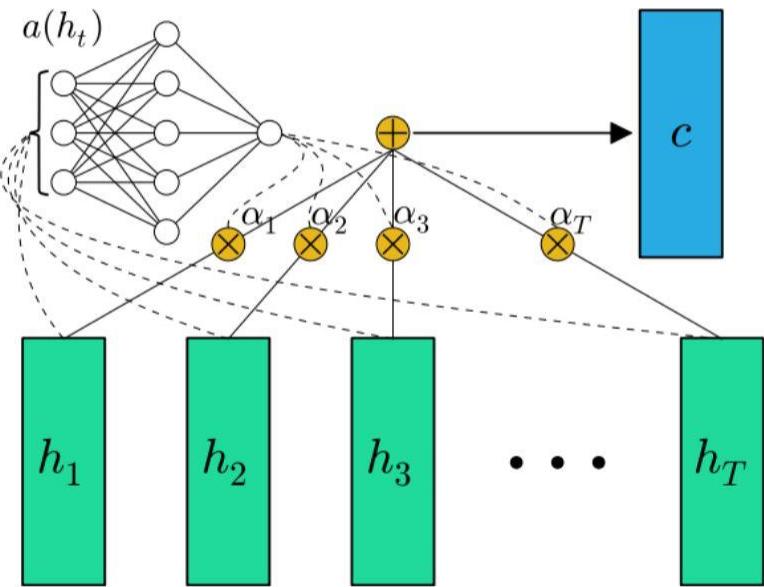
- Bi-directional LSTMs are just putting two independent LSTMs together
- The input sequence is fed in forward order for one LSTM, and in reverse order for the other
- The outputs of the two networks are usually concatenated at each time step
- Preserving information from both past and future helps understand context better

Bi-directional LSTMs + Attention



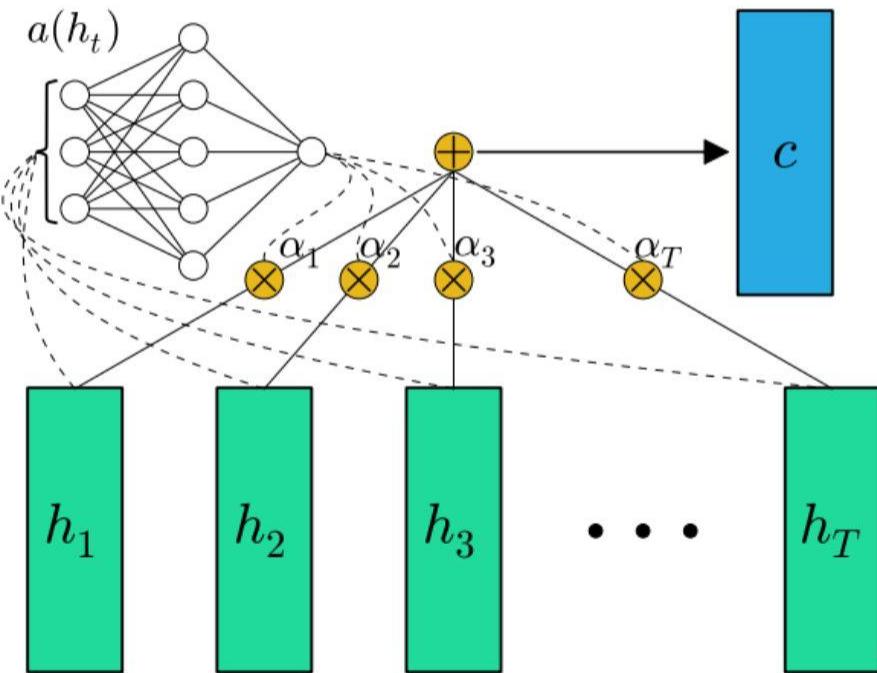
- Embedding Layer populated with pre-trained embeddings
 - FastText Embeddings
- Bi-directional GRUs or LSTMs
- Global Attention Layer
- FC Dense Layers

Attention Mechanism!



- Instead of using the output from the last GRU cell, send the entire sequence to a global attention layer
- Vectors from the hidden sequence h_t are fed into the learnable function $a(h_t)$
- Produces a probability vector α_t
- The final context vector c is a weighted average given by $\alpha_t \cdot h_t$

Attention Mechanism!



$$e_t = a(h_t) = \tanh(Wh_t + b)$$

$$\alpha_t = \text{softmax}(e_t) = \frac{\exp(e_t)}{\sum_{k=1}^T \exp(e_k)}$$

$$c = \sum_{t=1}^T \alpha_t h_t$$

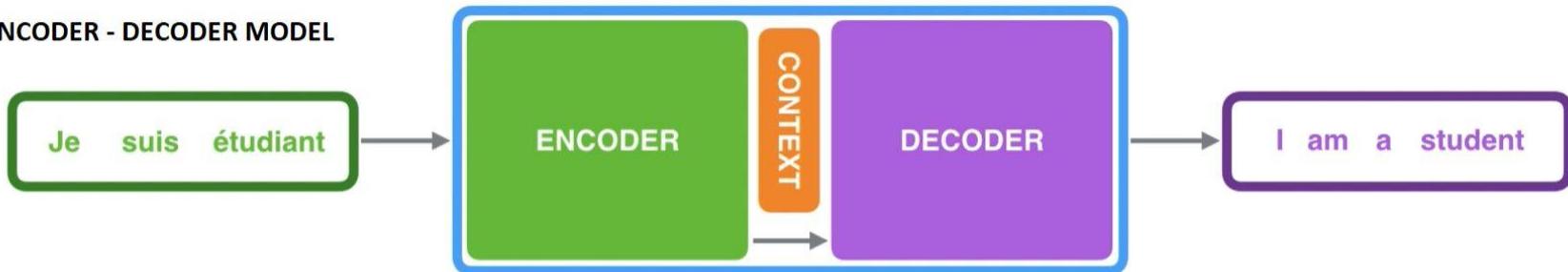
T is the total number of time steps in the input sequence



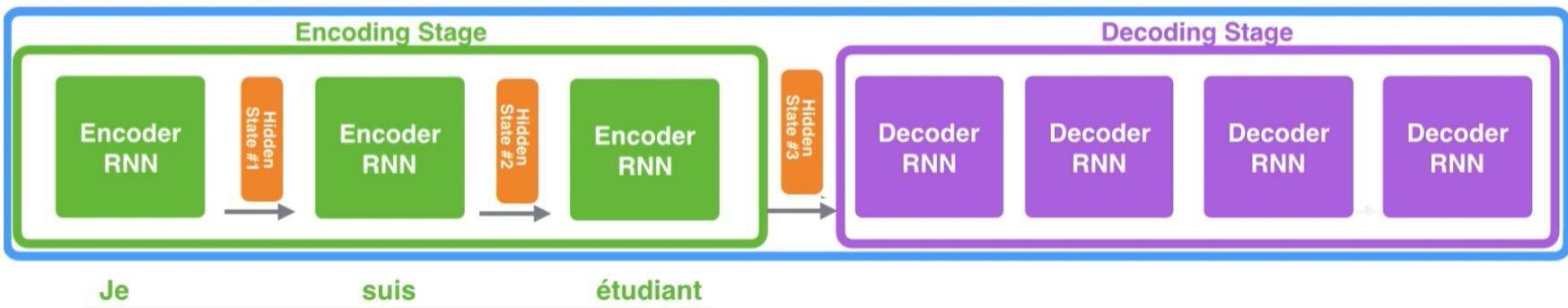
Motivation for Transformers

Encoder - Decoder Model (without Attention)

ENCODER - DECODER MODEL

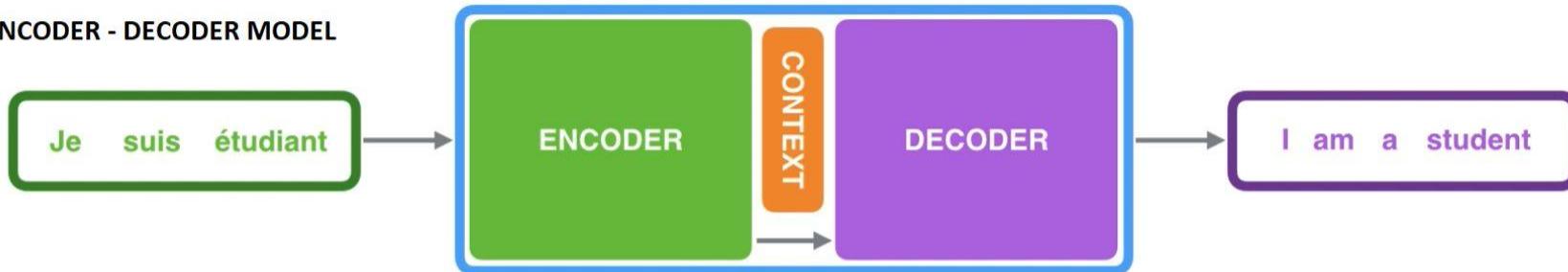


ENCODER - DECODER MODEL
(without Attention)



Encoder - Decoder Model (with Attention)

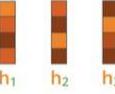
ENCODER - DECODER MODEL



ENCODER - DECODER MODEL
(with Attention)



Attention Mechanism - Illustrated

1. Prepare inputs

Encoder hidden states

Decoder hidden state at time step 4
2. Score each hidden state

13	9	9
----	---	---

scores
Attention weights for decoder time step #4
3. Softmax the scores

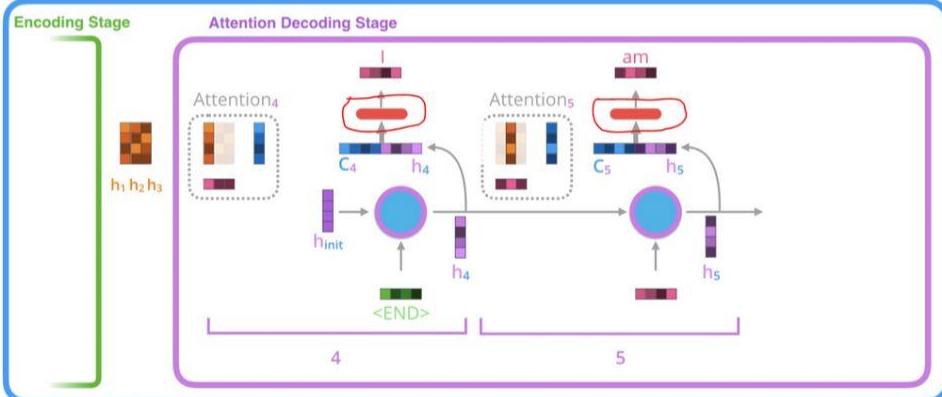
0.96	0.02	0.02
------	------	------

softmax scores
4. Multiply each vector by its softmaxed score

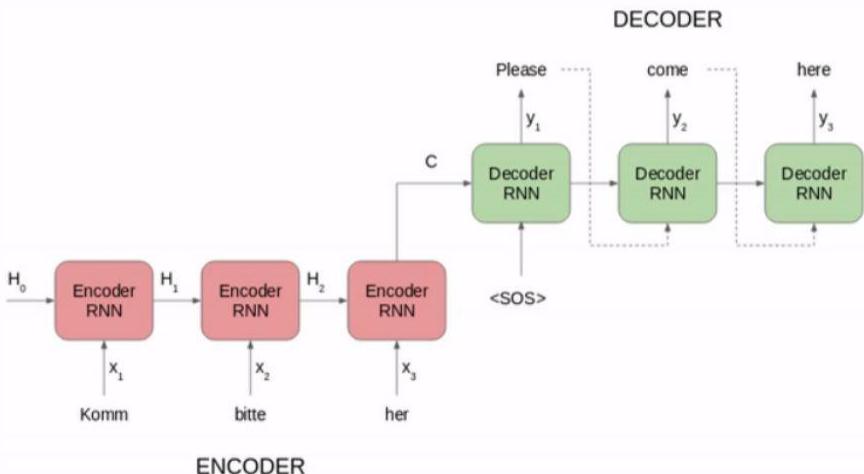
=
5. Sum up the weighted vectors

Context vector for decoder time step #4

We pass this vector through a **feedforward neural network** (one trained jointly with the model).
The **output** of the feedforward neural networks indicates the output word of this time step.

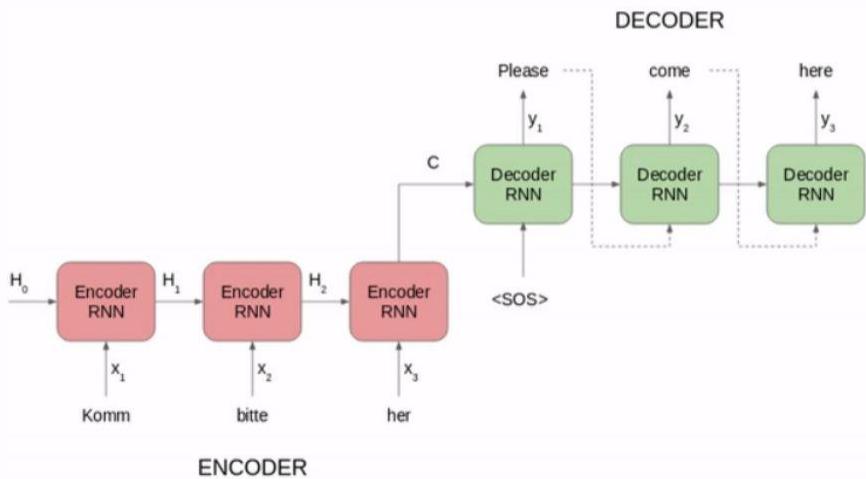


Encoder - Decoder Models - Summary



- Encoder & Decoder blocks consist of multiple RNN cells
- Hidden States are updated at each time step
- Context or Thought Vector from last step encodes information of the complete sequence
- Attention Mechanism uses a weighted sum of all hidden states instead of just the last one

Encoder - Decoder Models - Challenges



- 1 Encoding a long sequence without attention into a single context vector loses valuable information
- 2 Even with Attention, results show it is challenging to handle long sequences & dependencies
- 3 Sequential model has limitations with parallelizing training

Transformer Essentials

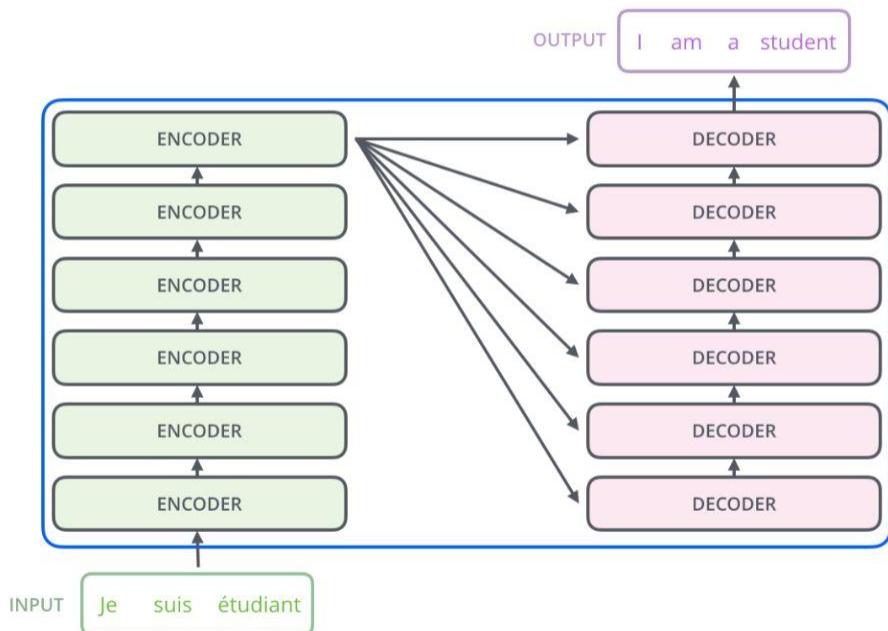


Transformer Model - Architecture



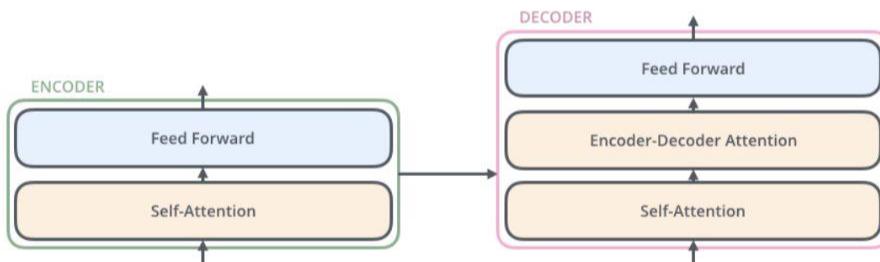
- Layered & Stacked Encoder - Decoder Model
- Relies on Multi-headed Self-Attention & Encoder-Decoder Attention
- No sequential RNN based training (completely parallelized)
- Achieved state-of-the-art performance on several NLP tasks

Transformer Model - Architecture



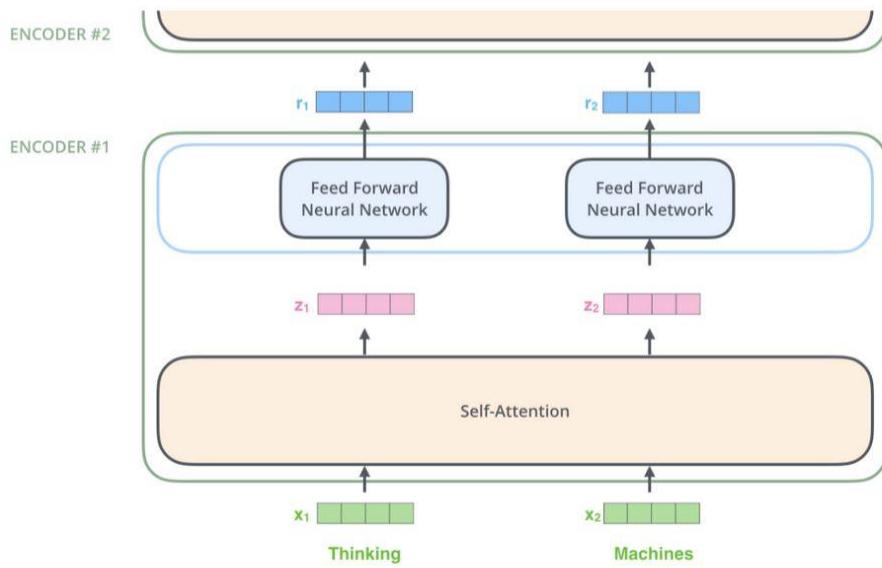
- At-heart a transformer model is a stacked encoder-decoder model
- Has 6-12 stacked encoder and decoder blocks usually based on standard architectures
- Based on the type of transformer model (only encoder \ decoder or both are used)

Transformer Model - Architecture



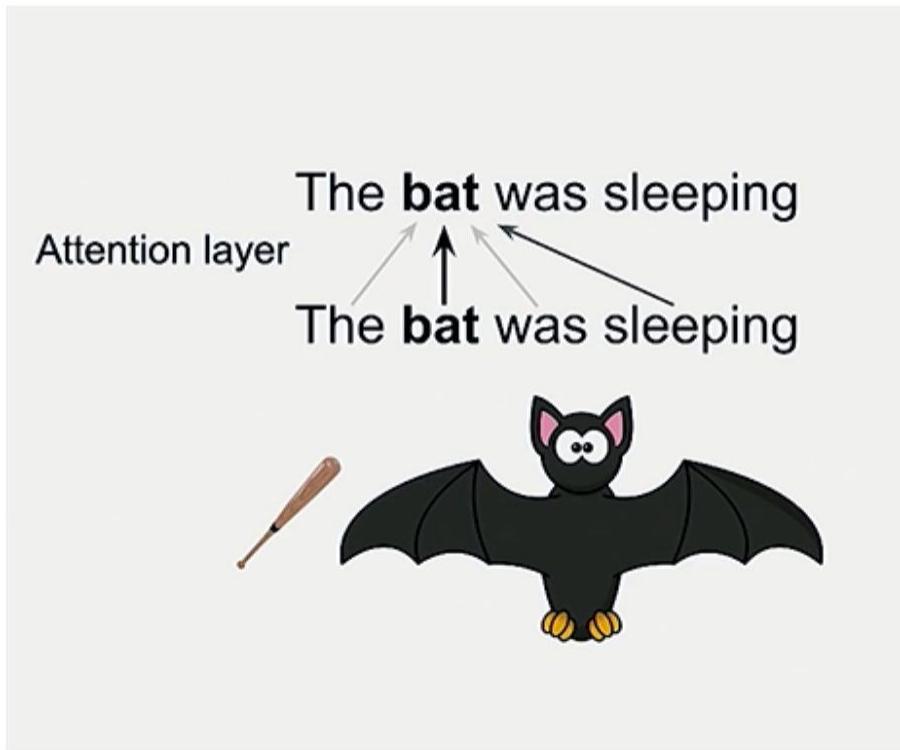
- The encoder's inputs flow through a self-attention layer
 - Each word attends (looks at) every other word in this process
- The outputs of the self-attention layer are fed to a feed-forward neural network
- The decoder has similar layers
- Between them is an attention layer that helps the decoder focus on relevant parts of the input sentence
 - This is the encoder-decoder attention layer

Transformer Model - Encoder Flow



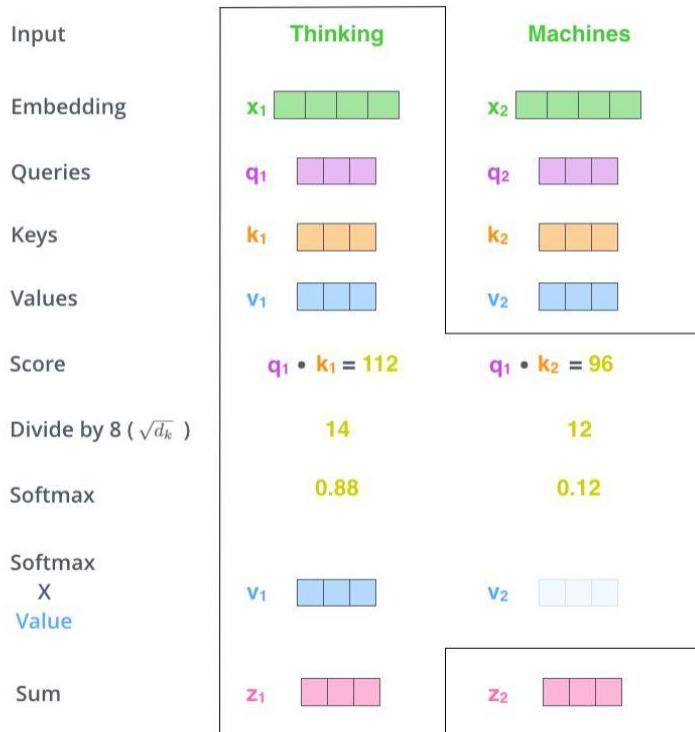
- **Self-attention layers (multi-head) at each block**
 - Each word attends (looks at) every other word in this process
- **The outputs of each self-attention layer are fed to a feed-forward neural network (FFNN)**
- **Output from each FFNN are sent upwards to the next encoder layer (with the same workflow as above)**

Self-Attention: A Simplistic Overview



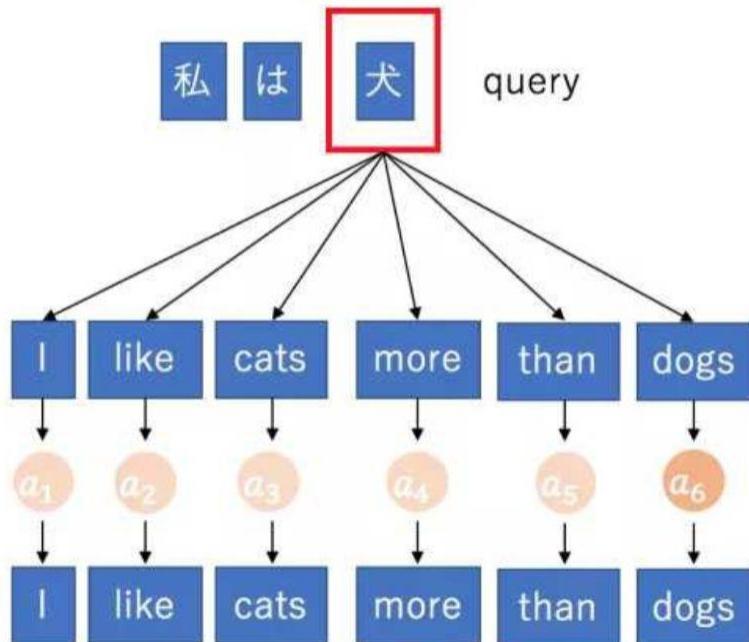
- What does “bat” in this sentence refer to? Is it referring to a baseball bat \ cricket bat or the animal?
- When the model is processing the word “bat”, self-attention allows it to associate “bat” with “sleeping” strongly
 - This gives it the indication that it could be an animal
- Self attention allows the model to look at other positions in the input sequence for clues that can help lead to a better encoding for each word

Self-Attention: A Detailed Overview



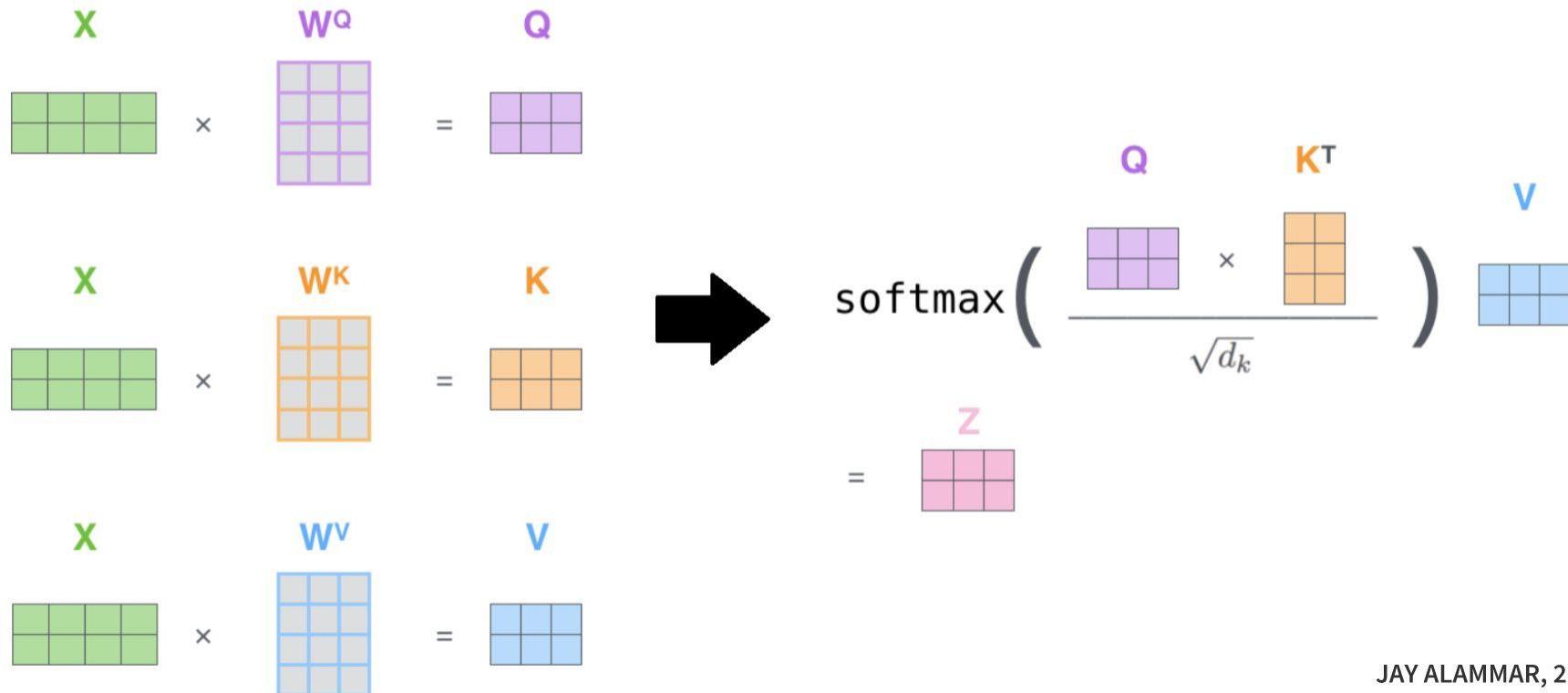
- Embeddings are generated for each word
- Query, Key and Value tensors are trained over epochs
 - Query comes from target sequence
 - Key & Value comes from source sequence
 - Initially K, V can be same
- Importance \ Alignment of each key is computed w.r.t the query to get scores
 - Scores are divided by 8 to stabilize gradients
- Softmax normalizes the scores and computes dot product with values (source sequence) to create the sum vector
- Result vector is sent to next FFNN

Self-Attention: A Detailed Overview

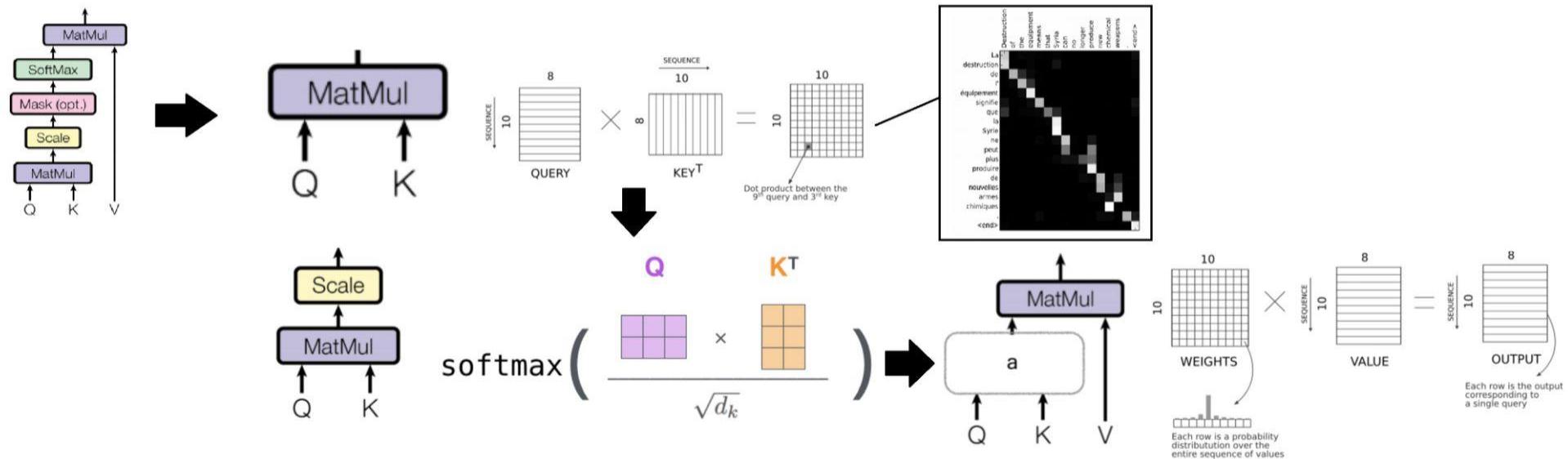


- Embeddings are generated for each word
- Query, Key and Value tensors are trained over epochs
 - Query comes from target sequence
 - Key & Value comes from source sequence
 - Initially K, V can be same
- Query and Key representations (embeddings) are multiplied (dot-product) to give dynamic alignment (attention) weights
 - Can be visualized as a heatmap
- These attention context weights are multiplied (dot-product) with the Value representation (embeddings) to give outputs for input Query
- This is then scaled using softmax

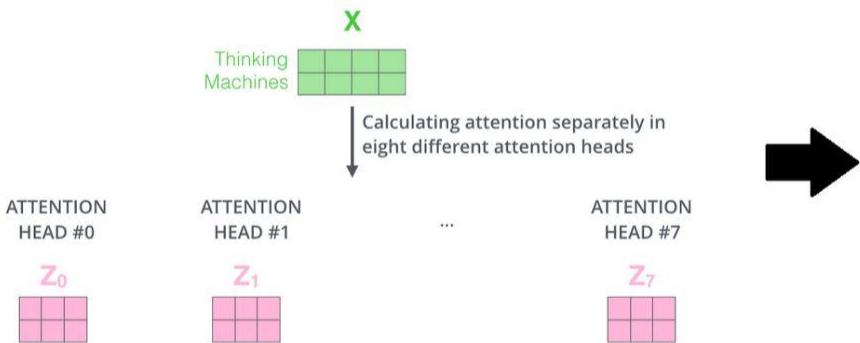
Self-Attention: A Detailed Overview - Tensor Ops



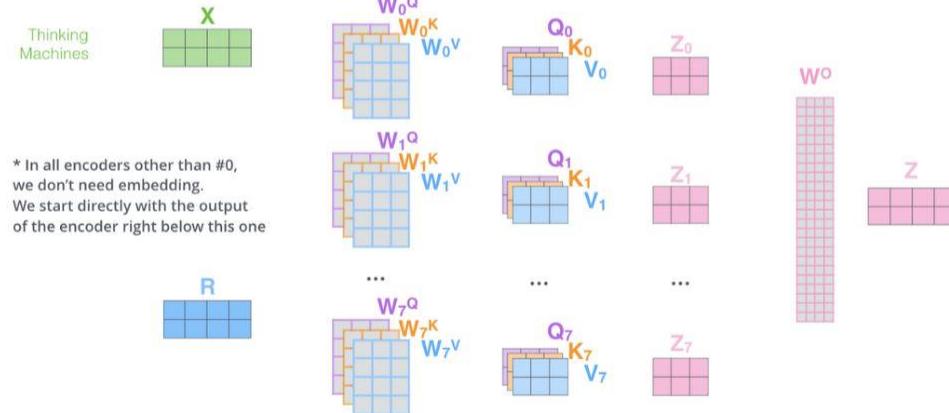
Self-Attention: A Detailed Overview - Tensor Ops



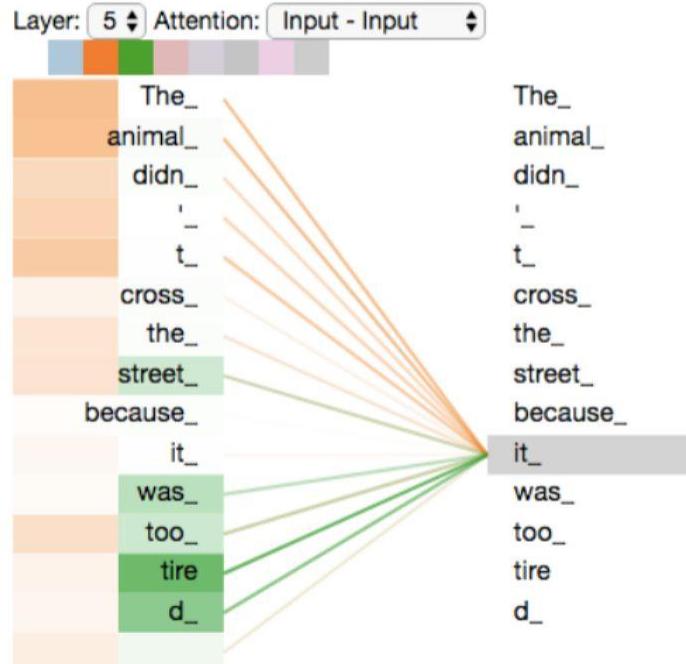
Multi-Headed Attention



- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

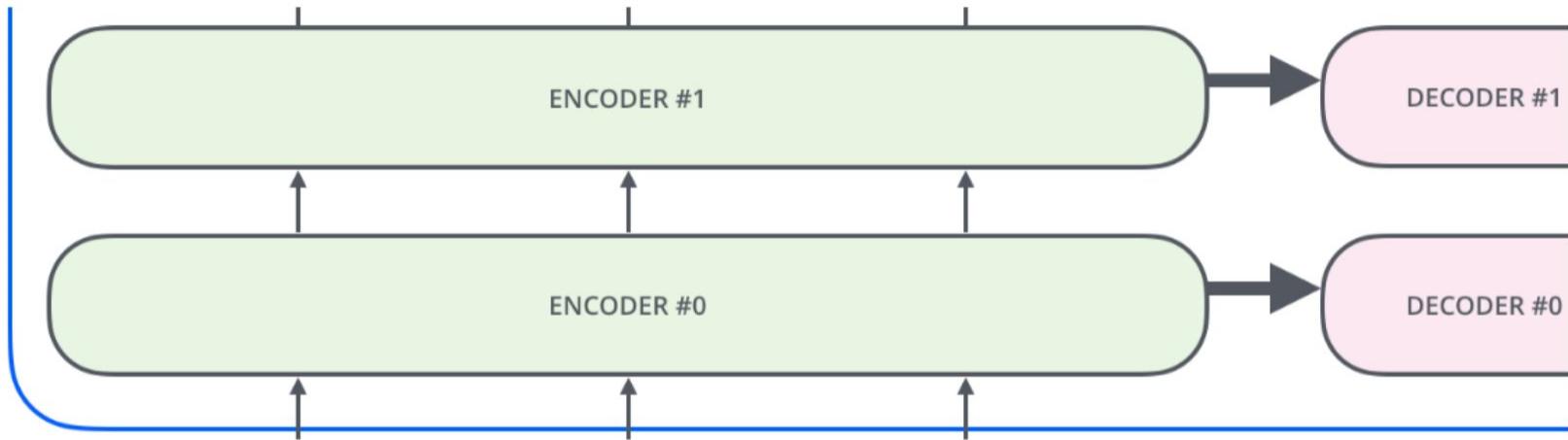


Multi-Headed Attention - Encode different aspects

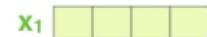


As we encode the word "it", one attention head is focusing most on "the animal", while another is focusing on "tired" -- in a sense, the model's representation of the word "it" bakes in some of the representation of both "animal" and "tired".

Encoding Sequential Representations in Transformers



EMBEDDING
WITH TIME
SIGNAL



POSITIONAL
ENCODING



EMBEDDINGS



=

=

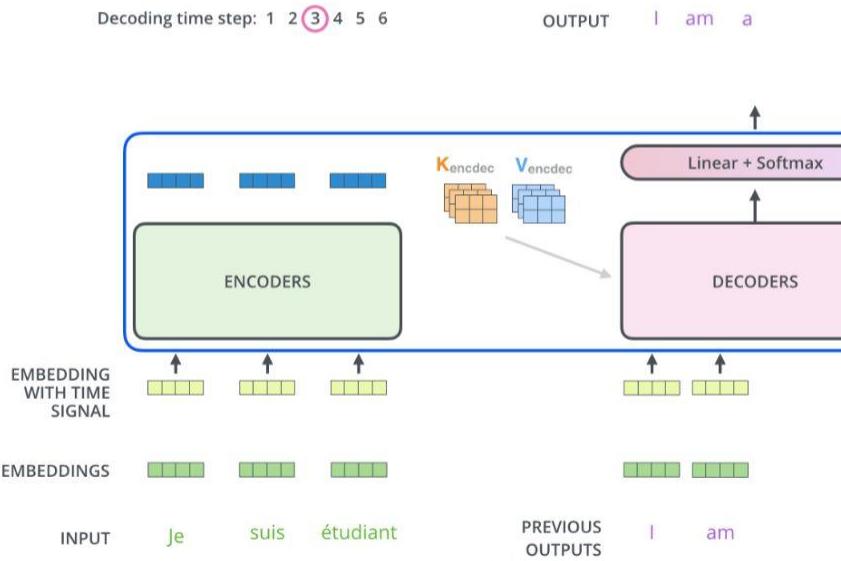
=

+

+

+

Transformer Model - Decoder Flow



- The output of the top encoder is transformed into a set of attention vectors K and V
 - Helps the decoder focus on appropriate places in the input sequence
 - Works just like multiheaded self-attention
 - Creates its Queries matrix from the layer below it
 - Takes the Keys and Values matrix from the output of the encoder stack
- The output of each step is fed to the bottom decoder in the next time step
 - The decoders bubble up their decoding results just like the encoders did
 - We embed and add positional encoding as we did before
 - In the decoder, the self-attention layer is only allowed to attend to earlier positions in the output sequence
 - This is done by masking future positions

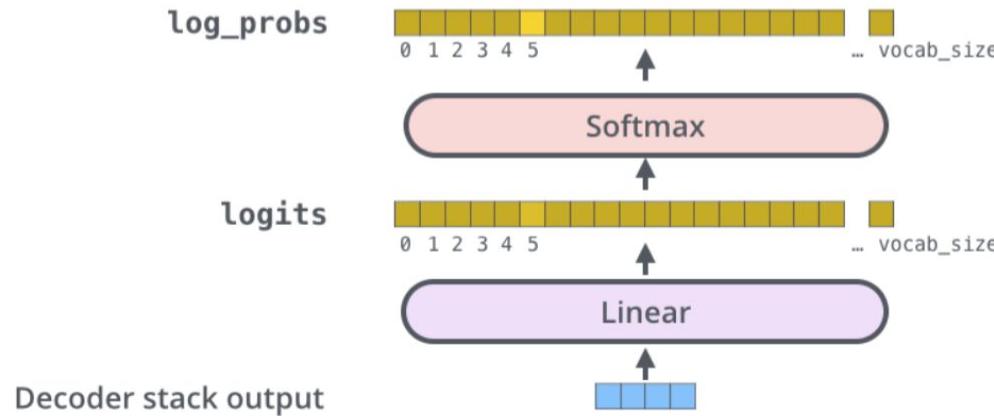
Final Layers in the Transformer Model

Which word in our vocabulary
is associated with this index?

am

Get the index of the cell
with the highest value
(argmax)

5



Transformer Model Architectures

Transformer Model Types

① Autoregressive Models

- Correspond to the decoder of the original transformer model
- Mask is used on top of the full sentence so that the attention heads can only see previous words
- Pretrained on the classic language modeling task: guess the next token
- Example - GPT family

② Autoencoding Models

- They correspond to the encoder of the original transformer model
- They get access to the full inputs without any mask
- Pretrained by corrupting the input tokens in some way and trying to reconstruct the original sentence (Masked Language Modeling)
- Example - BERT family

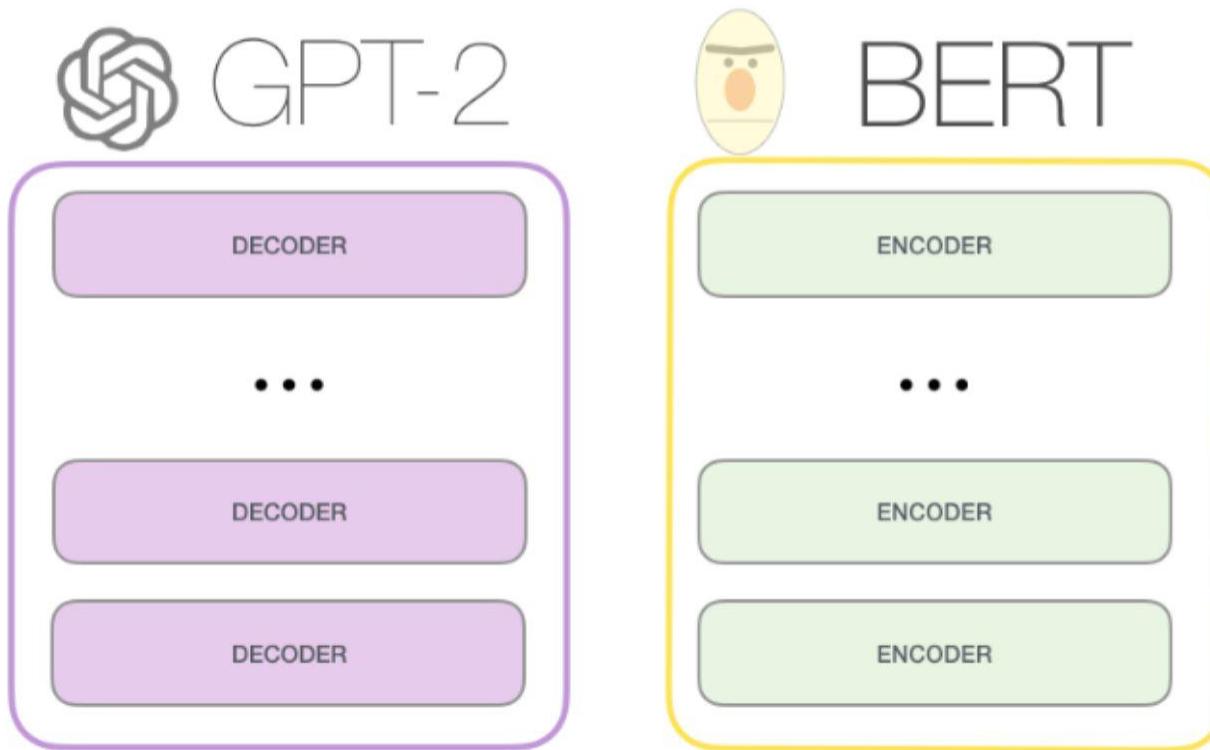
③ Sequence to Sequence Models

- Uses both the encoder and the decoder of the original transformer
- Most popular applications are translation, summarization and question answering
- Example - BART, T5

④ Multi-modal Models

- Mixes data of different modalities e.g text & images

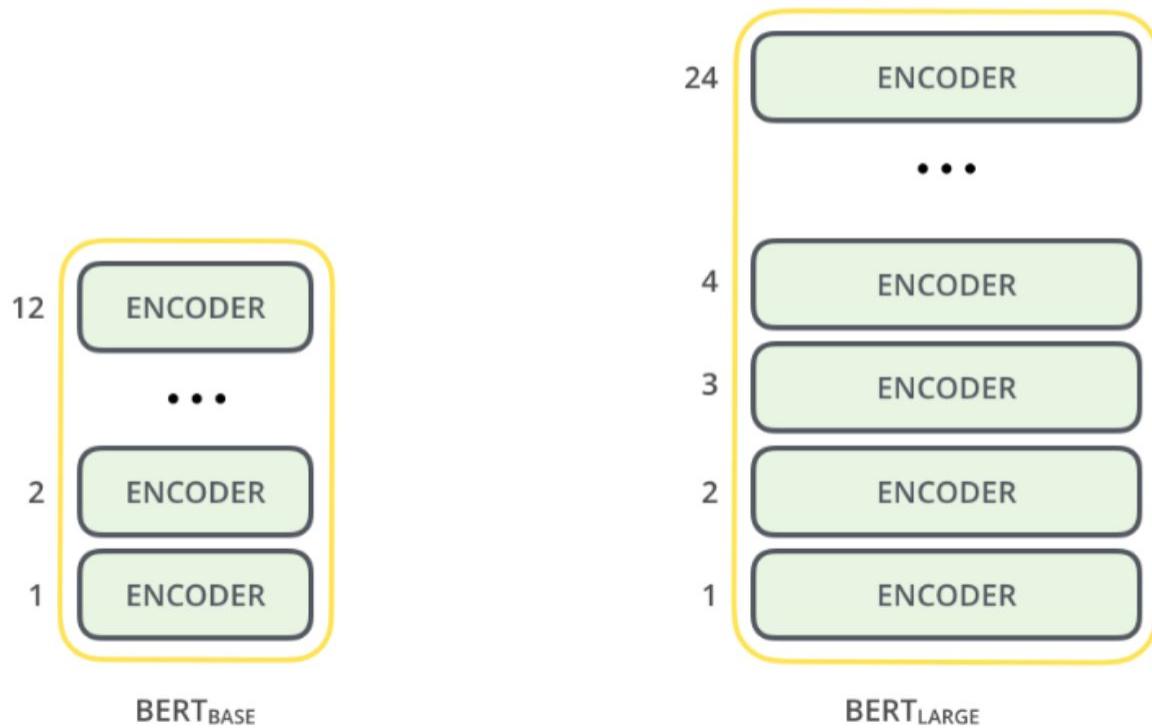
Popular Transformer Models



Understanding BERT & DistilBERT



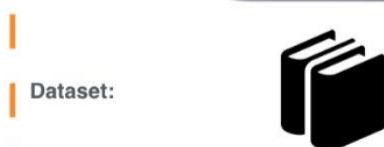
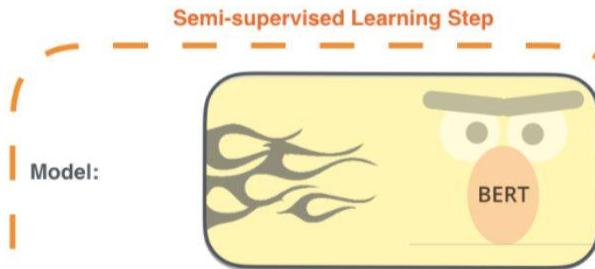
BERT - Trained Transformer Encoder Stack



BERT Training Workflow

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

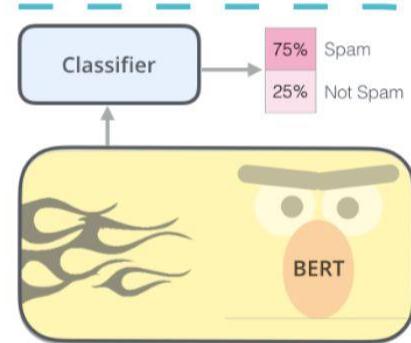
The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



Objective:
Predict the masked word
(language modeling)

2 - **Supervised** training on a specific task with a labeled dataset.

Supervised Learning Step



Dataset:

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

BERT Training Workflow - Masked Language Modeling

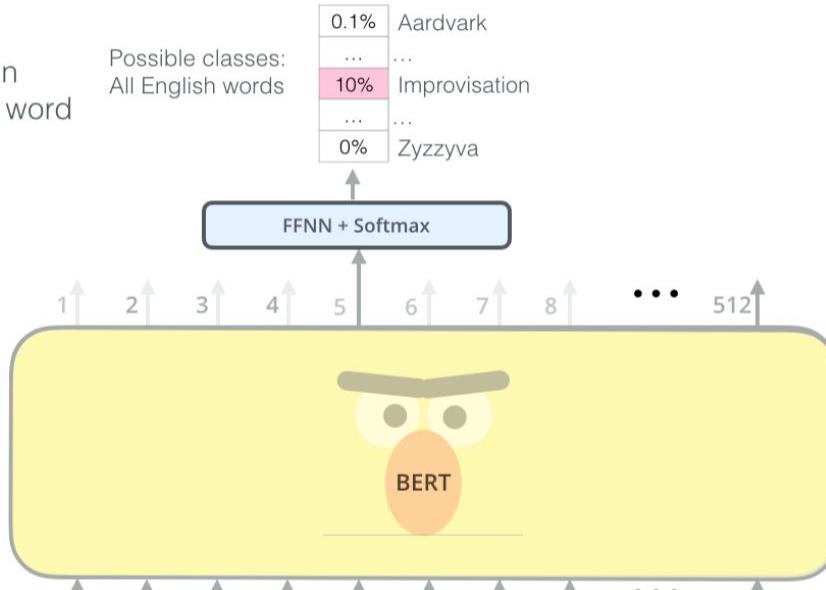
Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva

FFNN + Softmax

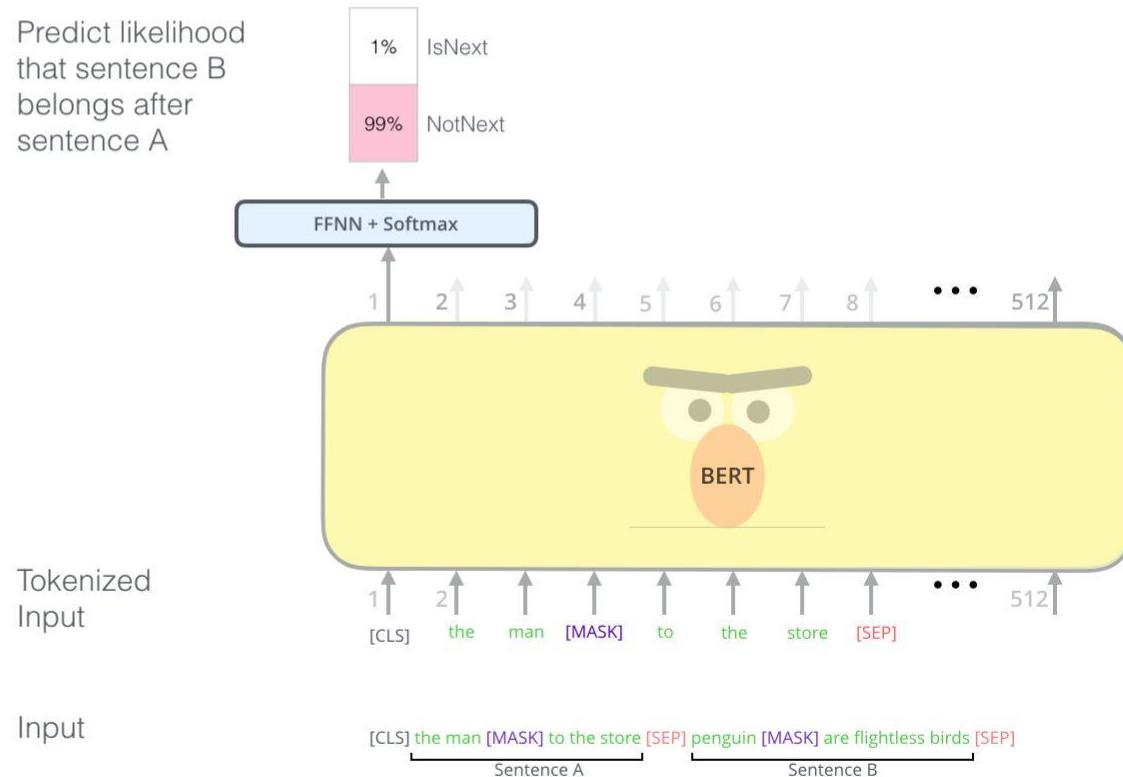
Randomly mask 15% of tokens



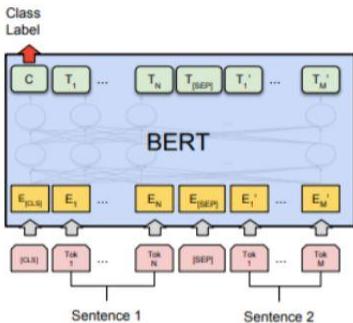
Input

[CLS] Let's stick to improvisation in this skit

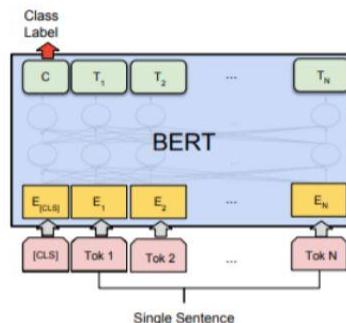
BERT Training Workflow - Next Sentence Prediction



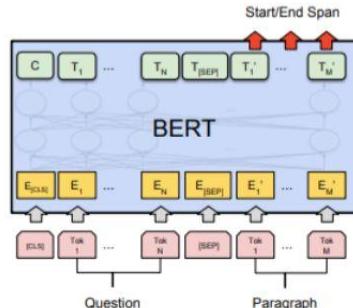
Fine-tuning BERT - Multi-Task NLP



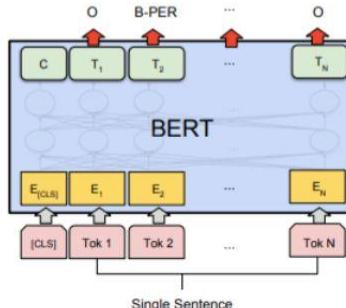
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(b) Single Sentence Classification Tasks:
SST-2, CoLA

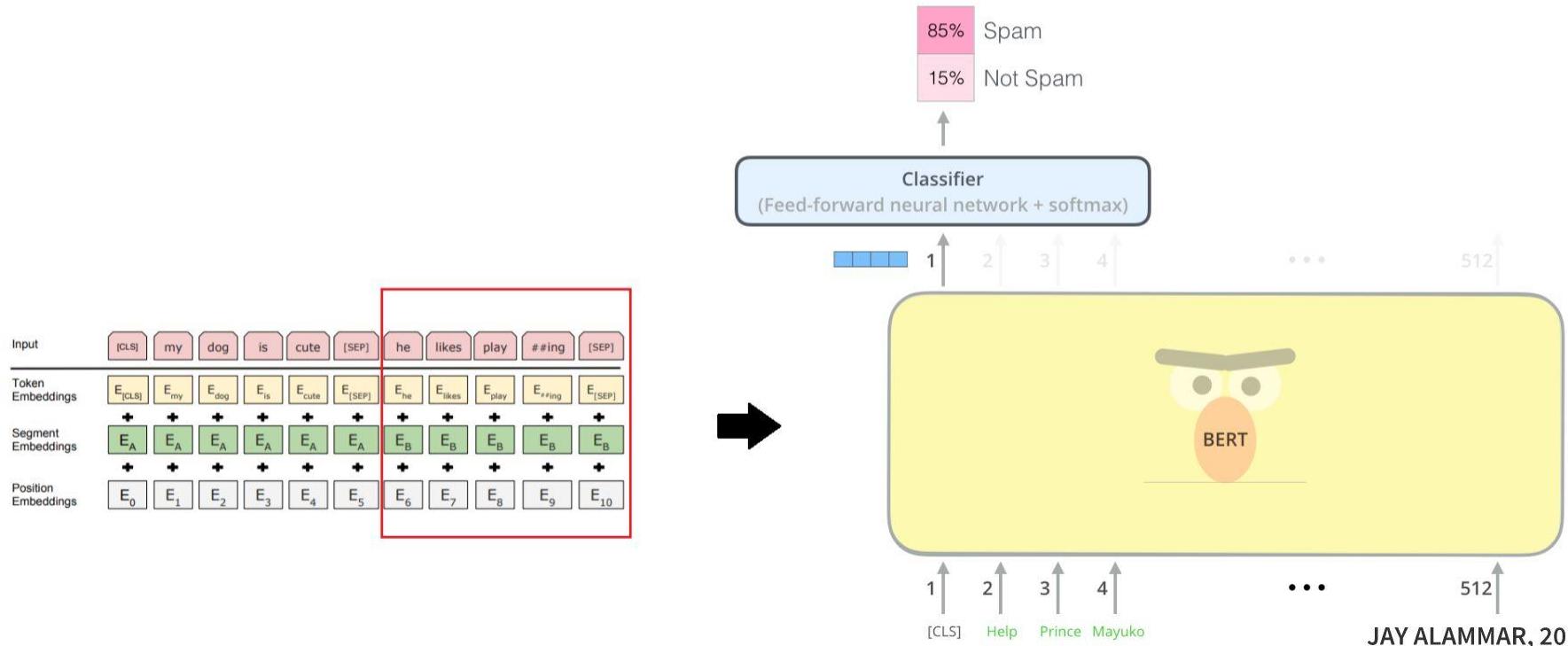


(c) Question Answering Tasks:
SQuAD v1.1

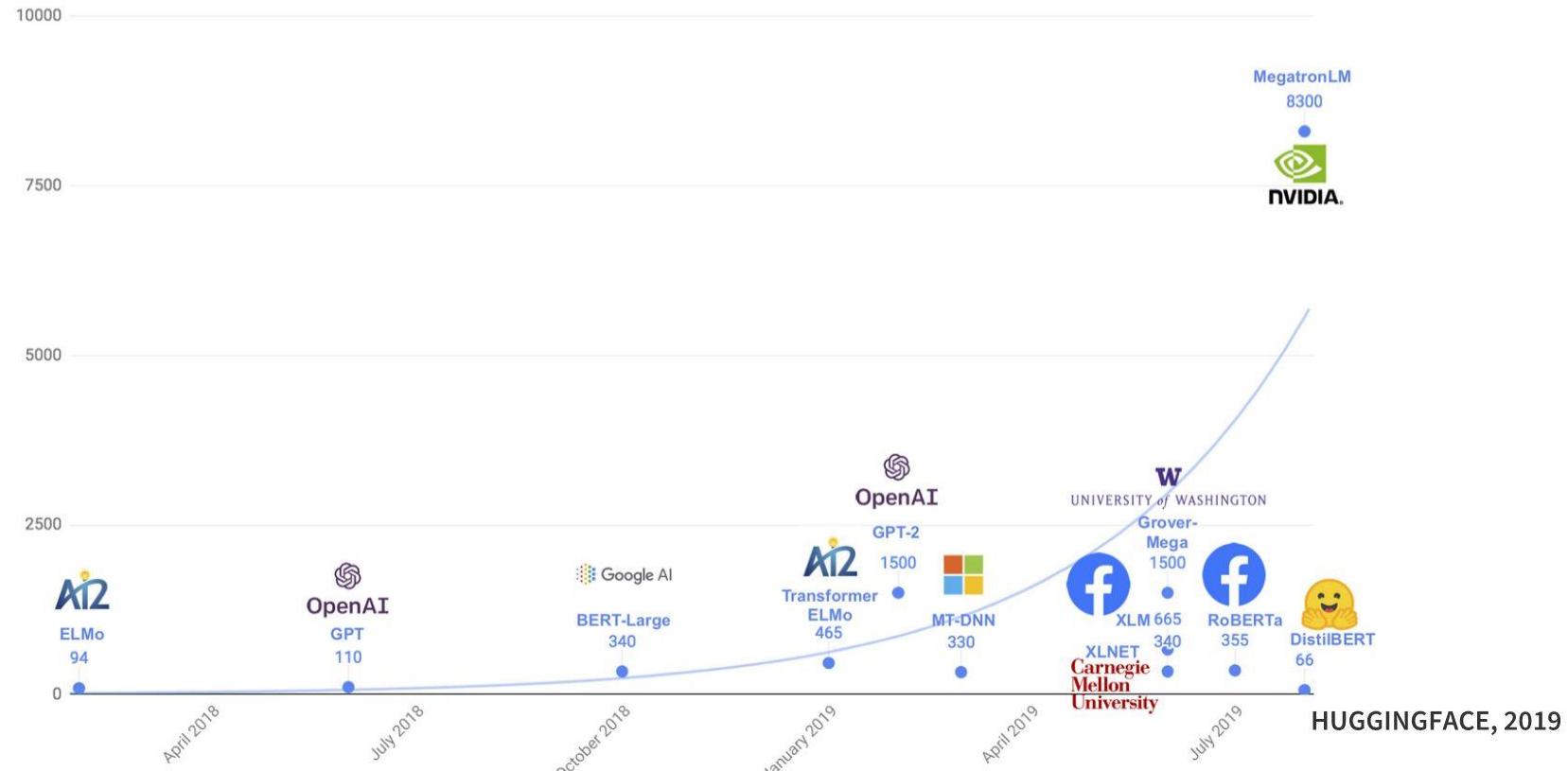


(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

Fine-tuning BERT - Text Classification



Challenges with Massive Size of Pre-trained Models



DistilBERT - Compressing BERT

- Compression technique known as knowledge distillation is used
- Teacher-Student Training is used - training student to mimic output distribution of teacher model (BERT)

Trained with cross-entropy on soft targets (probabilities from teacher model)

- HuggingFace uses KL-divergence loss to train DistilBERT
- Has almost halved the number of parameters from BERT and retains 95% of performance



Hands-on Tutorials

Hands-on Tutorials

1

Text Classification - Less Data Availability Problem

- Logistic Regression
- BERT
- DistilBERT

2

Multi-task NLP with Transformer Pipelines

- Sentiment Analysis
- NER
- Question & Answering
- Text Generation - Mask Filling
- Summarization
- Translation
- Generation
- Feature Extraction

References

- **Research Papers**

- https://github.com/dipanjanS/transformers_nlp_essentials/tree/master/papers

- **Visuals & Content**

- <https://jalammar.github.io/illustrated-transformer/>
- <https://medium.com/@b.terryjack/deep-learning-the-transformer-9ae5e9c5a190>
- [https://www.analyticsvidhya.com
\(understanding-transformers-nlp-state-of-the-art-models\)](https://www.analyticsvidhya.com/understanding-transformers-nlp-state-of-the-art-models)
- NLP Using Transformer Architectures - A. Géron
- Medium \ Towards Data Science
- <https://ruder.io/state-of-transfer-learning-in-nlp/>

Stay in Touch!



LinkedIn

<https://www.linkedin.com/in/dipanzan>



GitHub

<https://github.com/dipanjanS>



Medium

<https://medium.com/@dipanzan.sarkar>

A photograph of a sunset over a calm body of water. The sky is a gradient of orange and yellow. In the center, the sun is low on the horizon, casting a bright reflection on the water. On the left side of the frame, a small silhouette of a person is visible in a single-person kayak or canoe, facing away from the viewer towards the setting sun. The overall atmosphere is peaceful and serene.

Thank You