

LARGE LANGUAGE MODELS FOR GENOMICS

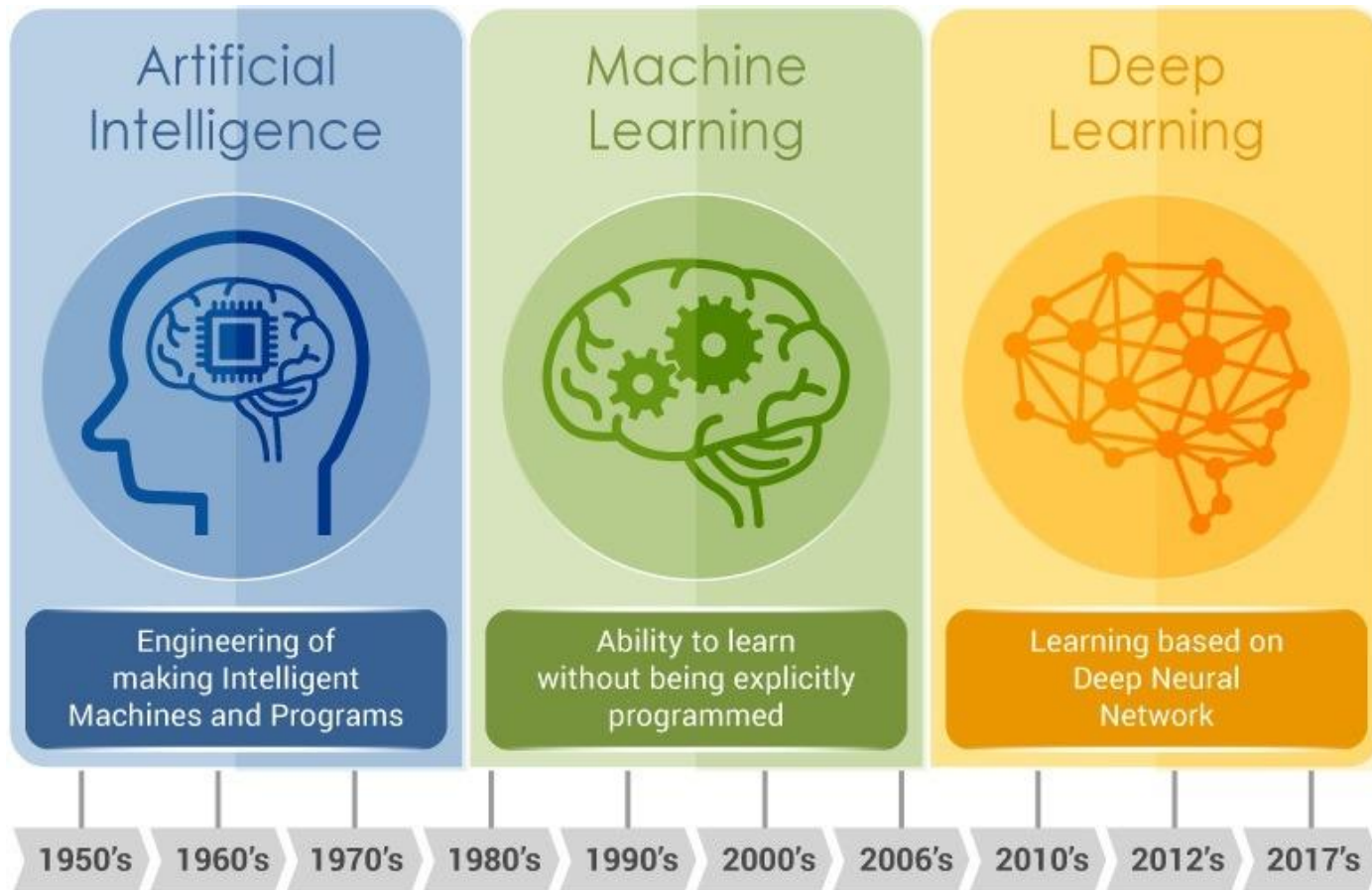
Raphaël MOURAD, Associate Prof.

MIAT INRAe

Université Paul Sabatier, Toulouse III

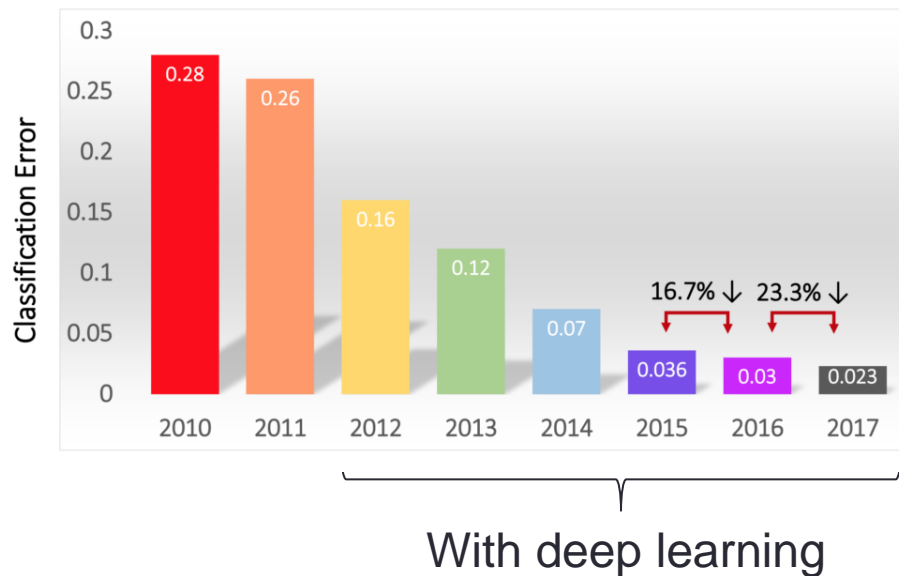
QUICK INTRO TO LARGE LANGUAGE MODELS

Deep learning as a branch of AI



Success of deep learning since 2012: Example of computer vision

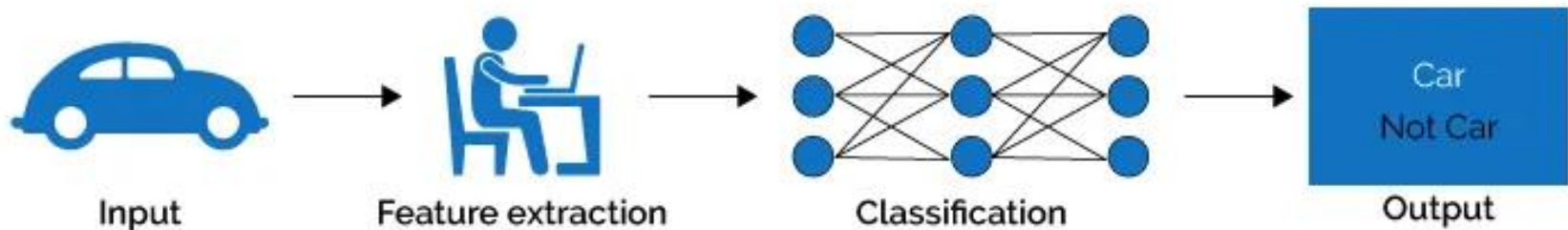
Image classification (ImageNet challenge)



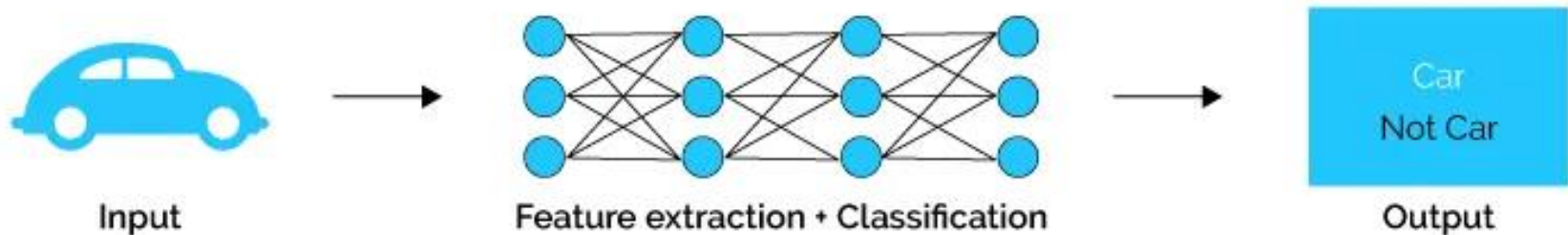
- 2012: AlexNet (convNet)
- 2013: ZFNet
- 2014:
 - VGGNet (deeper, simpler)
 - InceptionNet (faster)
- 2015: ResNet (deeper)
- 2016: Ensemble networks

Difference between machine and deep learning

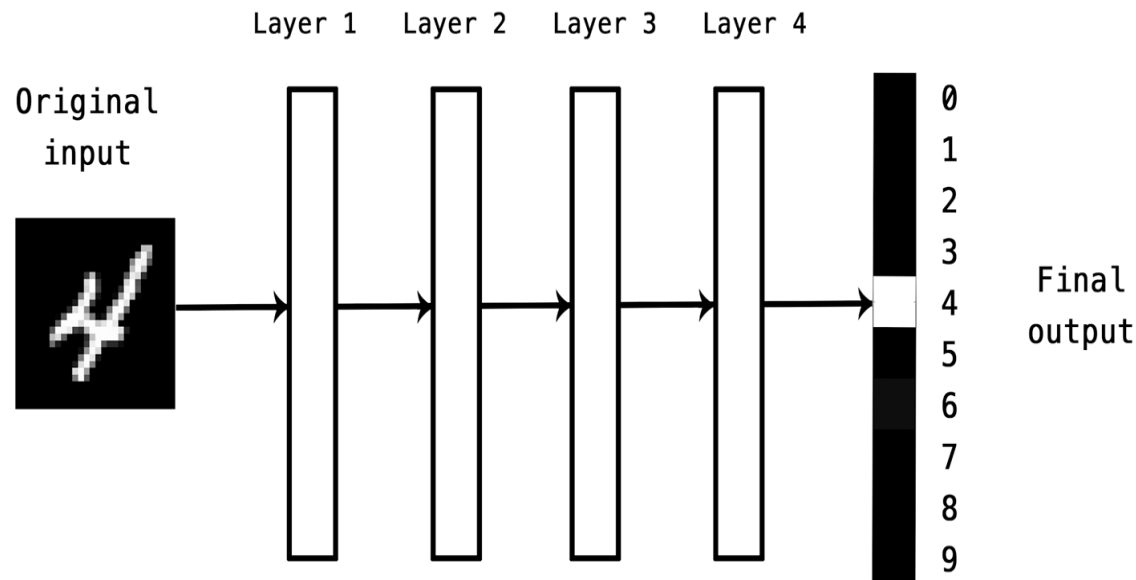
Machine Learning



Deep Learning

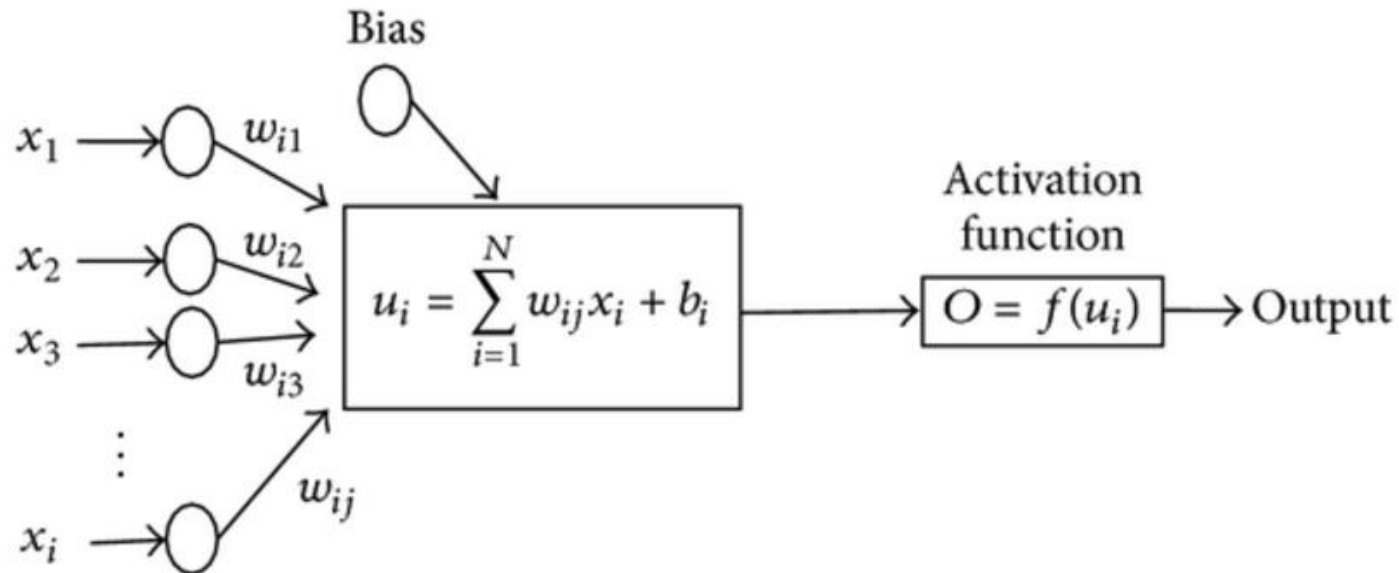


Deep learning as neural networks



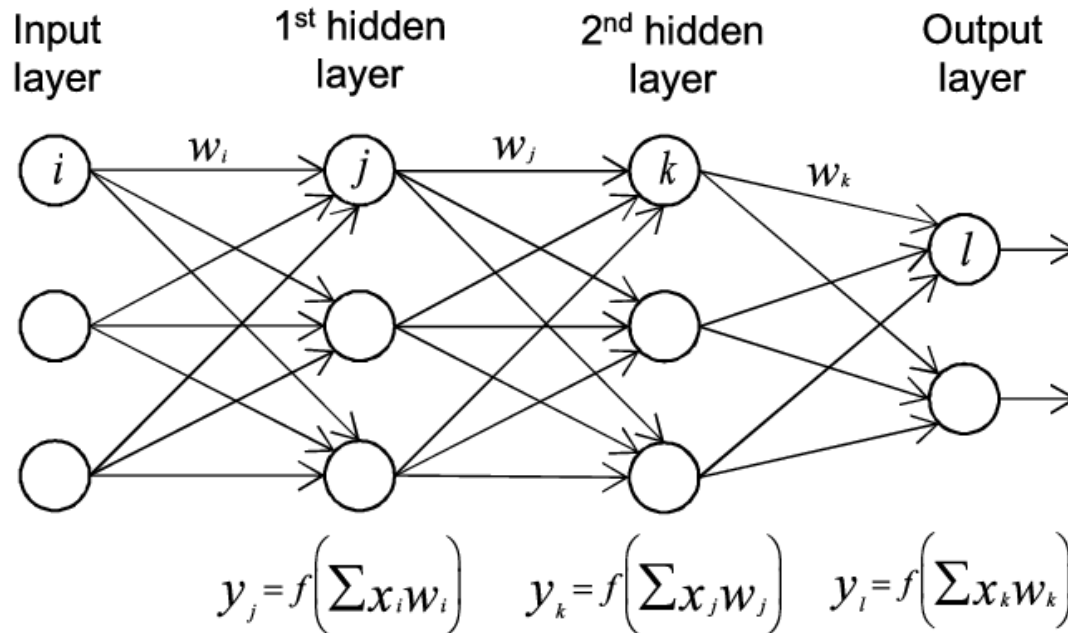
- Deep learning is based on a **deep** neural network which is the stacking of different neuronal layers to predict a final output.

Neural networks



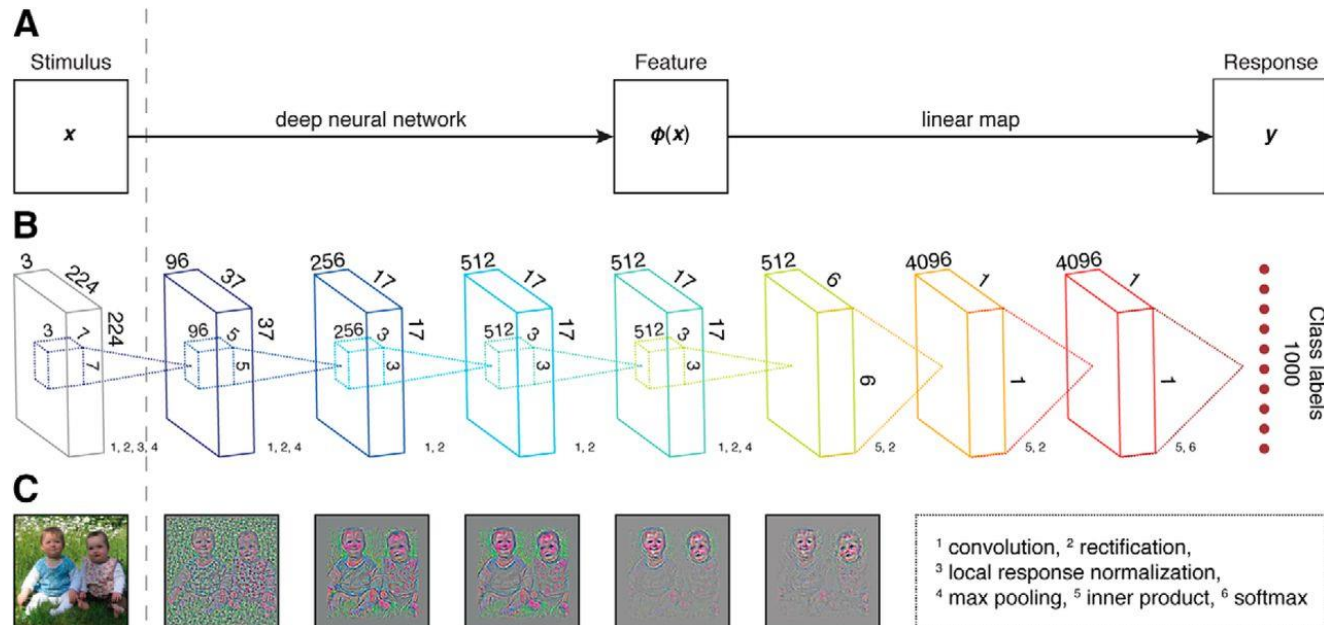
- In a neural network, multiple inputs x_i are combined through a linear combination (with weights w_i), and then an activation function is used for a non-linear transformation to obtain the output.

Deep neural networks



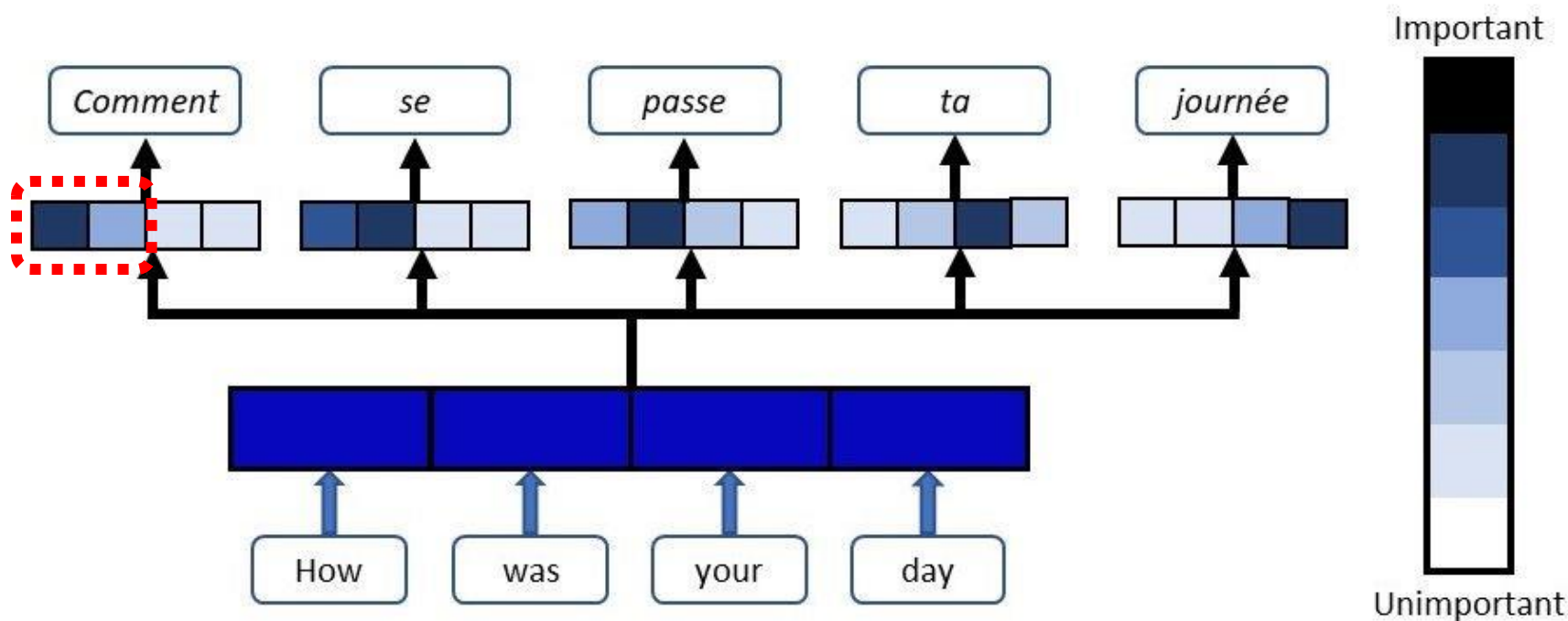
- A deep neural network (DNN) is an neural network (NN) with multiple layers between the input and output layers. Each hidden layer linearly combines the output from the previous layer and then does a non-linear transformation.

Principle of deep learning : stacking many different sorts of layers



- Building a deep neural network is like assembling a lego toy where every lego brick is a layer.

Attention



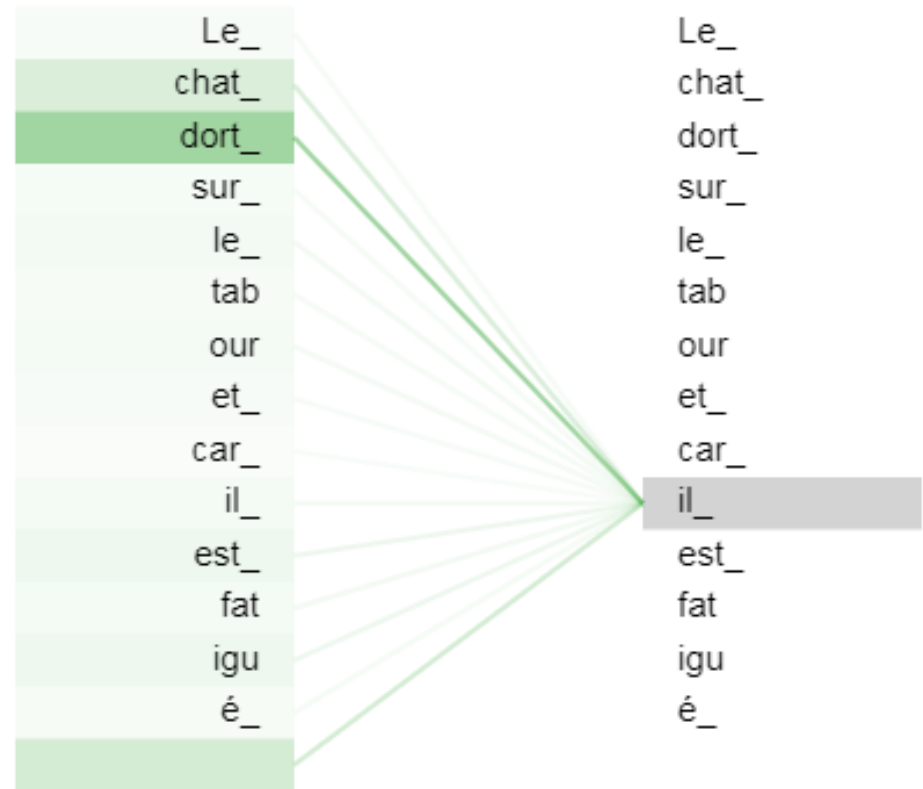
- Here, attention is used to weight different words from English to translate into French.
- For instance, to translate “how” to “comment”, you don’t only need the word “comment” (high weight) but you need other words such as the word “was” (moderate weight).

Self-attention in transformer model

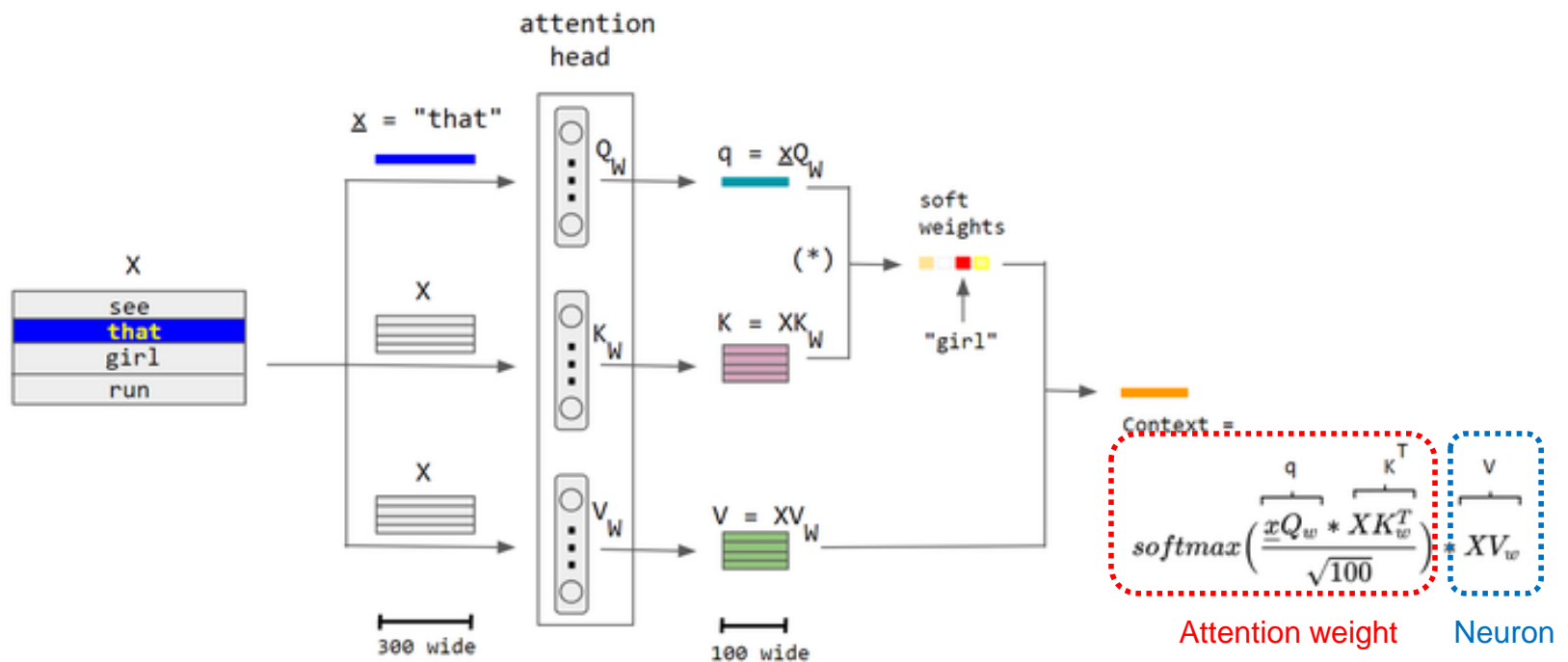
Self-attention is similar to attention, except that it is applied to the input sequence itself.

Self-attention allows to model long-range dependencies between any word in a sequence.

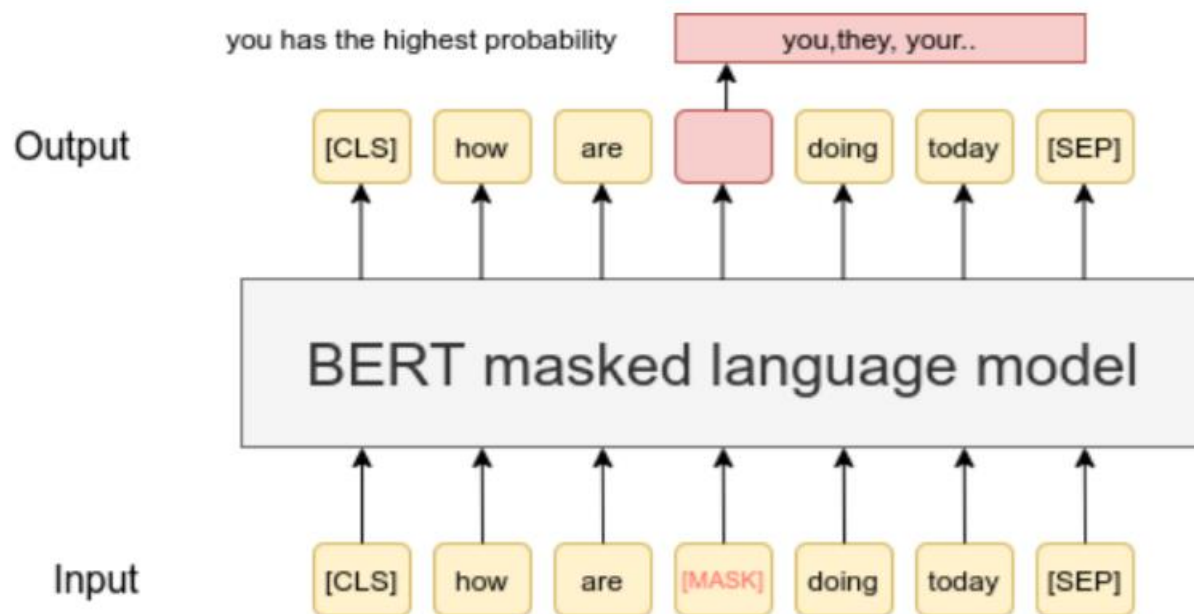
Self-attention is not directional as compared to RNN (LSTM or GRU), allowing parallel computing.



How to compute self-attention weights



Masked language model



- In a masked language model, the task is to predict some words that were masked using the context of the words.
- Together with self-attention, it was used for the BERT model (Bidirectional Encoder Representations from Transformers).

GPT-1 model (Masked language model)

- GPT-1 implements the transformer architecture (self-attention).

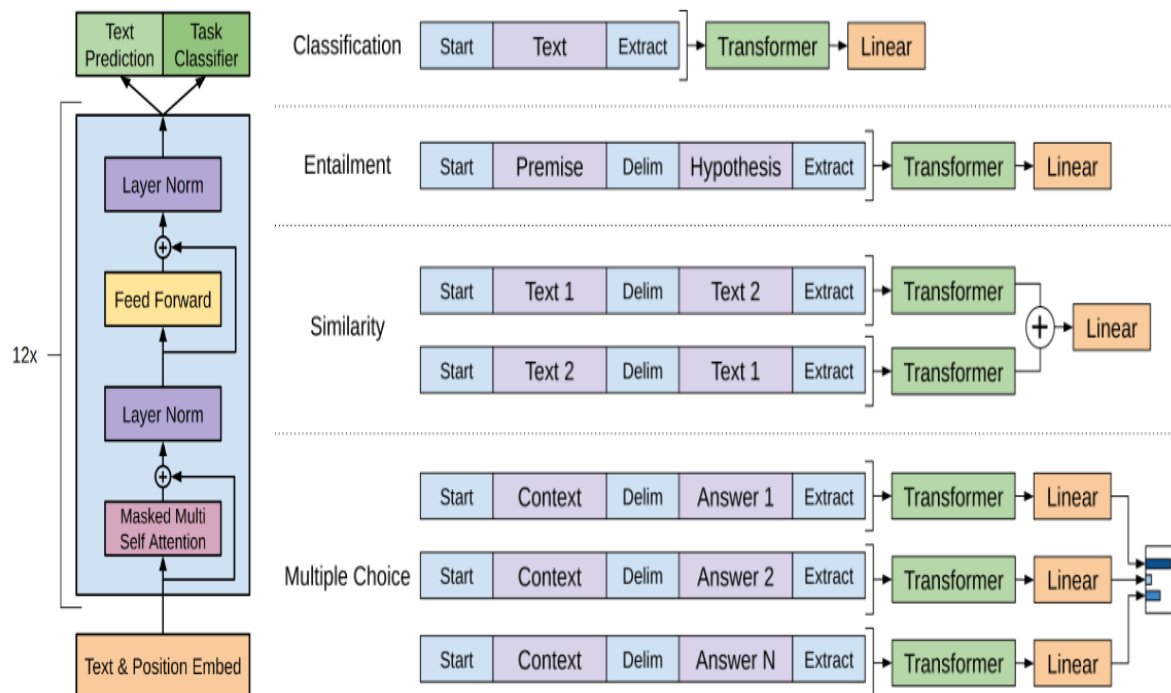


Figure 1: **(left)** Transformer architecture and training objectives used in this work. **(right)** Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

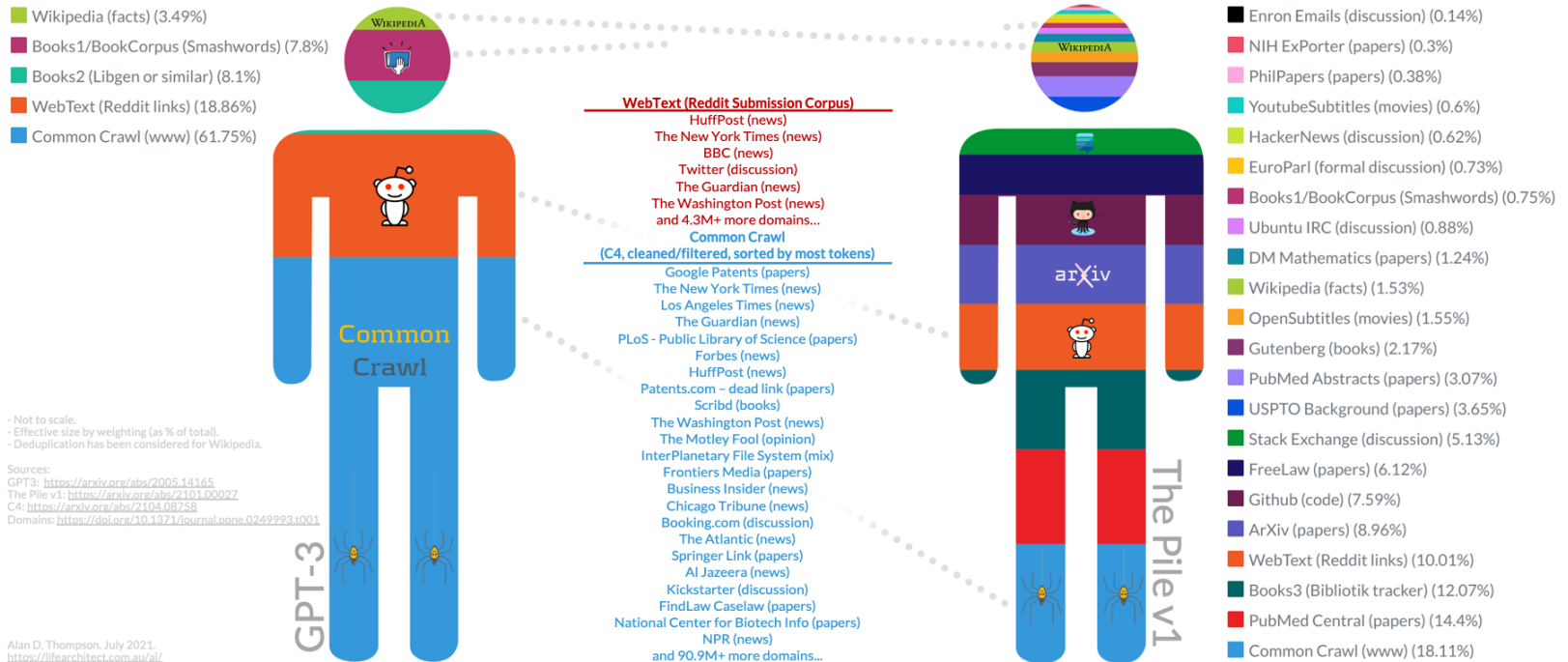
GPT-1, GPT-2 and GPT-3

	GPT-1	GPT-2	GPT-3
Parameters	117 million	1.5 billion	175 billion
Decoded Layers	12	48	96
Hidden Layer	768	1600	12288
Context Token size	512	1024	2048

- Number of parameters increasing over time.

GPT-3 training data

CONTENTS OF GPT-3 & THE PILE V1 ELEUTHER'S GPT-NEO, GPT-J, GPT-NEOX, BAAI'S WUDAO 2.0, AND MORE...



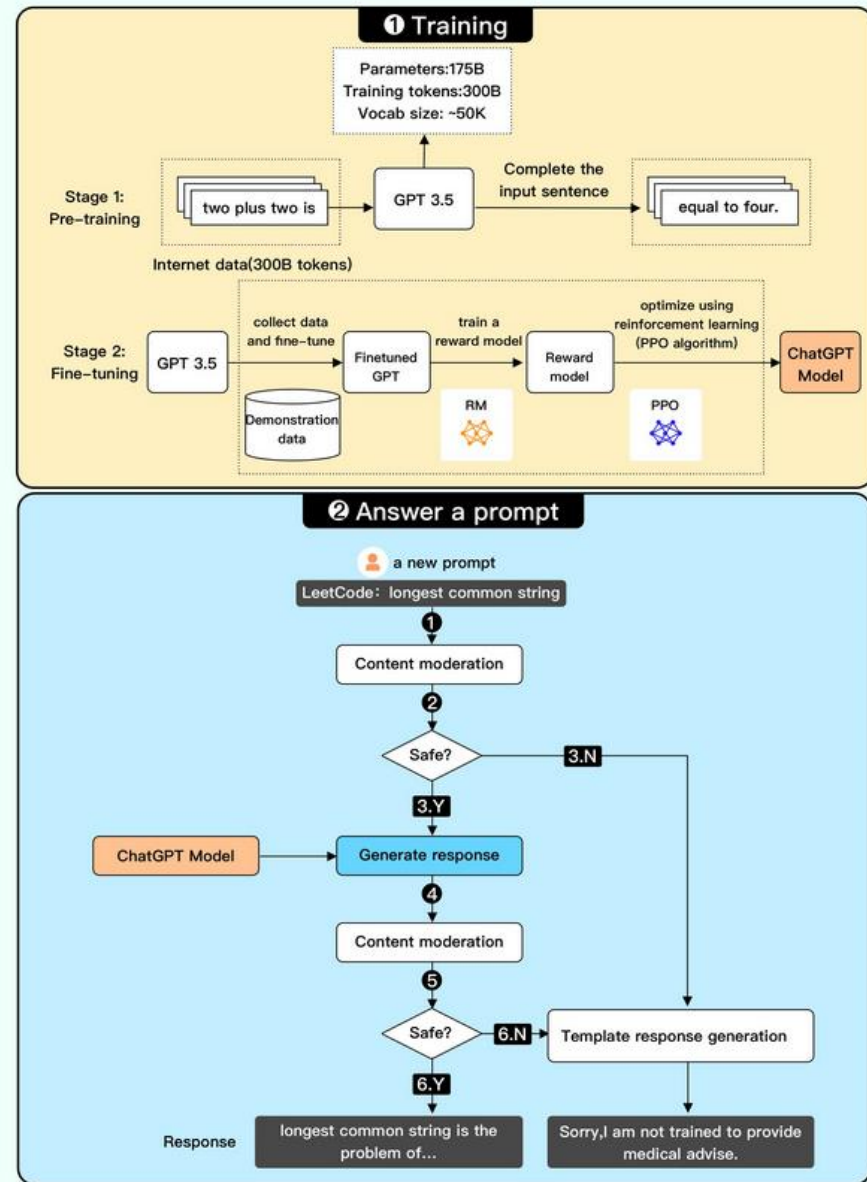
LifeArchitect.ai/models

chatGPT

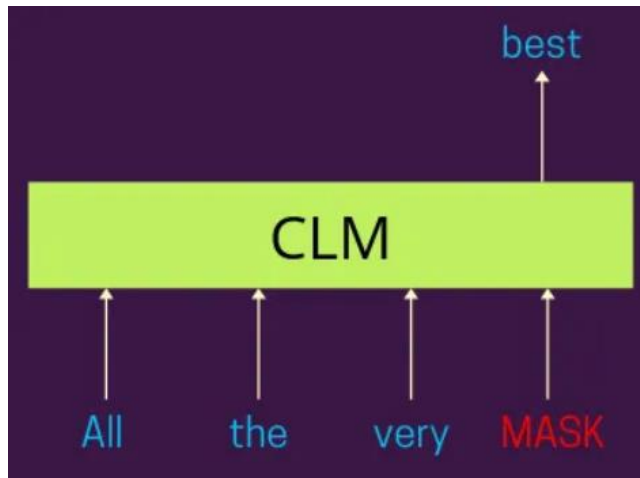
- Based on gpt3.5 (pretraining), with fine-tuning (stage 2) and reinforcement learning with human feedback (in blue).

How does ChatGPT-like System Work?

ByteByteGo.com



Causal language modeling (CLM) vs masked language modeling (MLM)



DEEP LEARNING FOR GENOMICS

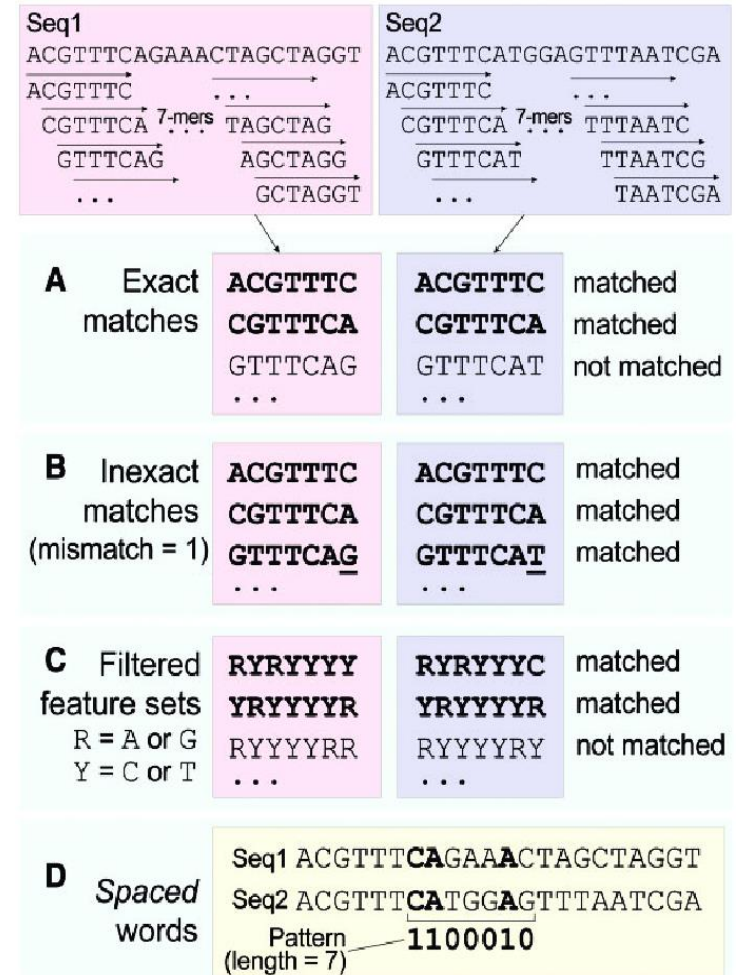
Human genome

- The human genome is composed on 22 autosomal chromosomes and two sex chromosomes. It is composed of 3,3 billion DNA letters (A, T, G, C).
- Below is an example of a DNA sequence:

T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C
C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C
T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C C T A A C C
G A G G A G A A C G C A A C T C C G C C G T T G C A A A G G C G C G C C G C G C C G G C G C A G G C G C A G A G A G G C G C
C A C A T G C T A G C G C G T C G G G G T G G A G G C G T G G C G C A G G C G C A G A G A G G C G C G C C G C G C C G G C G
A A G C C T A C G G G C G G G G T T G G G G G G C G T G T G T T G C A G G A G C A A A G T C G C A C G G C G C C G G G C
G C T T G C T C A C G G T G C T G T G C C A G G G C G C C C C C T G C T G G C G A C T A G G G C A A C T G C A G G G C T C T C T
G C A C G C C C A C C T G C T G G C A G C T G G G G A C A C T G C C G G G C C C T C T T G C T C C A A C A G T A C T G G C G G A

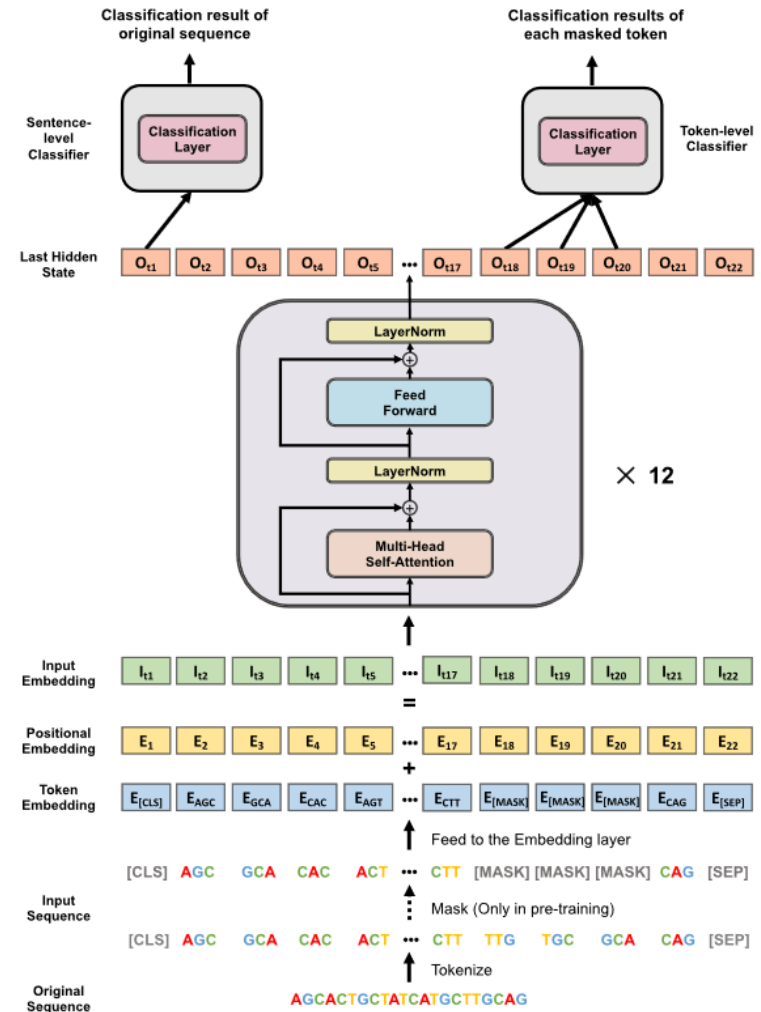
How to represent a DNA sequence for a large language model?

- K-mers (for instance: ATCTC or ATTTC, ...): very powerful approach since almost no prior information (except k) is needed to build the features.



DNABERT

- The self-attention model DNABERT is trained by masking some kmers in the DNA sequence and then by trying to predict them using the other k-mers in the DNA sequence (context).
- At the end, the model provides features that encode DNA sequences in a very efficient way for any predictive task.

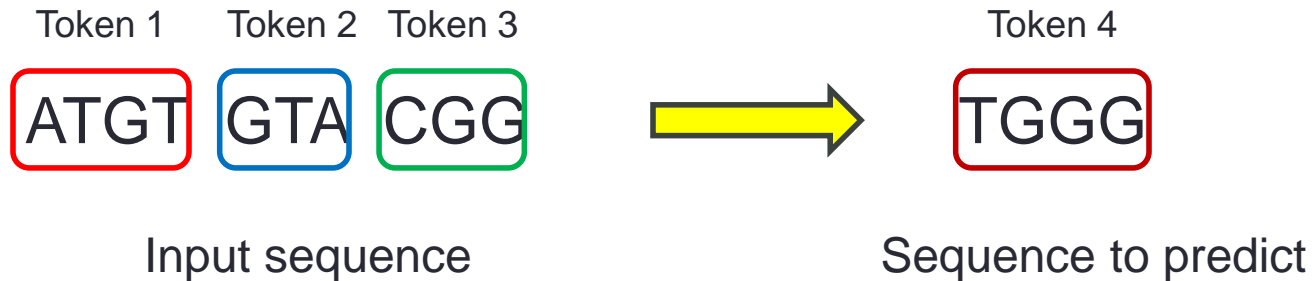


Mistral-DNA

- Mistral-DNA is an LLM based on Mixtral-8x7B-v0.1 model.
- Causal language model.
- Based on Byte pair tokenizer:
 - AGCCTTTCTCT -> AGC**Z**TT**ZZ**, where **Z**=CT
 - Allows to identify most frequent k-mers
- Sparse mixture of Experts: reduces number of parameters used during inference (prediction)

Mistral-DNA

- Prediction:

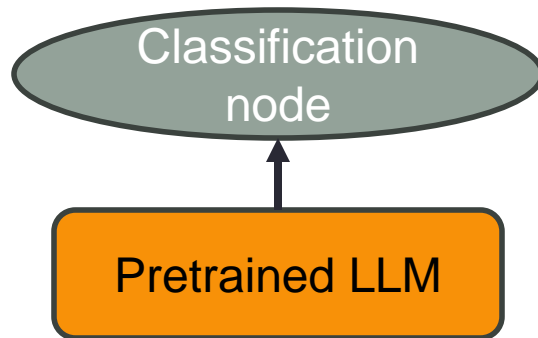


- The model is « pretrained » using this strategy. Given a DNA sequence, it tries to predict the next k-mer (=next word).

APPLICATIONS

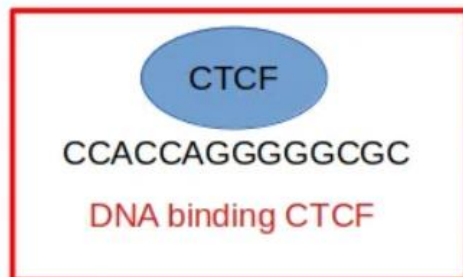
1) Model modification for classification

- 1st step: modify the model



Add classification node
at the top of the LLM

- 2nd step: train the model on labeled sequences (=finetuning).



2) Assessing the impact of a SNP

- Easy way to assess the impact of a SNP:
 - Predict the value for a sequence with the reference allele **C**:
 - ATGTAGTGGGTACCC**C**TGTGTAGAAGCCA
 - Predict the value for a sequence with the reference allele **T**:
 - ATGTAGTGGGTACCT**T**TGTGTAGAAGCCA

3) Synthetic DNA sequence generation

- Predict a synthetic DNA sequence:



- Repeat same task iteratively:

