

# **DNABERT**

**2021-2 Advanced Bioinformatics**

**SNU Graduate School of Public Health, Genome & Health Data Lab**

**Yeojin Jeong  
2021 / 12 / 01**

# Contents

**1. Background**

**2. Attention, Transformer and BERT**

**3. DNABERT**

**4. Google Colab**

# Background

# Deep Learning in Genomics

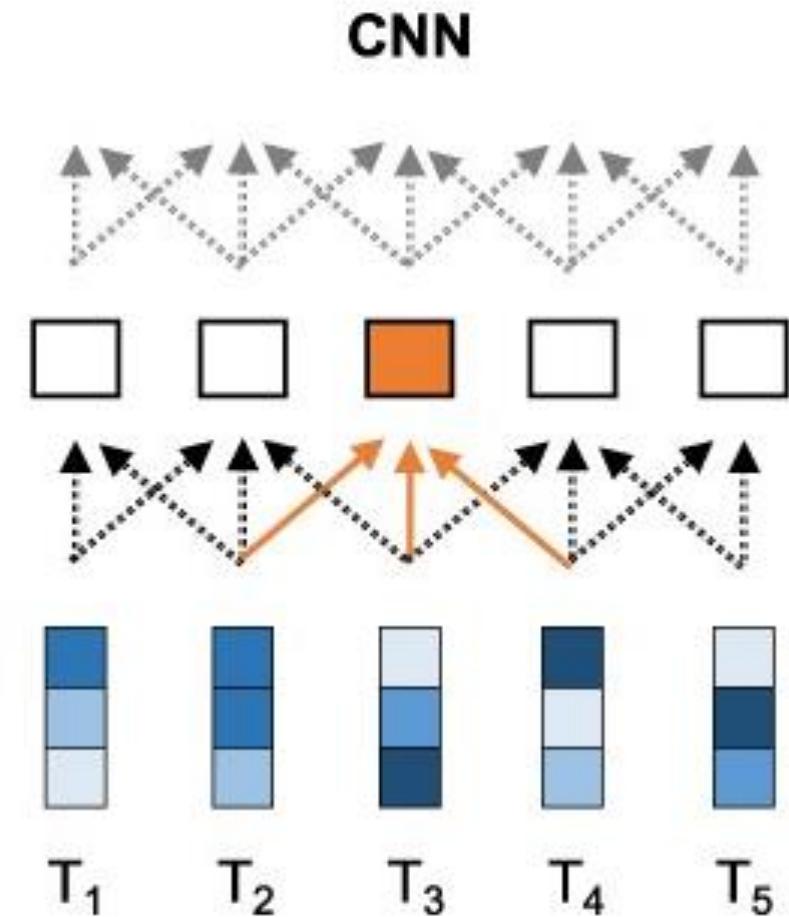
- Traditional Methods

Network type	a Fully connected	b Convolutional	c Recurrent	d Graph convolutional
Output				
Parameters				
Input				
Invariance	-	Translation	Time	Node index permutation
Input example	Predefined features such as number of k-mer matches, total conservation	<ul style="list-style-type: none"> <li>• DNA sequence</li> <li>• Amino acid sequence</li> <li>• Image</li> </ul>	<ul style="list-style-type: none"> <li>• DNA sequence</li> <li>• Amino acid sequence</li> <li>• Time series measurements</li> </ul>	<ul style="list-style-type: none"> <li>• Protein–protein interaction network</li> <li>• Citation network</li> <li>• Protein structure</li> </ul>
<p>The diagram illustrates four types of neural networks:</p> <ul style="list-style-type: none"> <li><b>Fully connected (a):</b> Shows a fully connected layer where every node in one layer is connected to every node in the next layer. Colored arrows indicate parameterized information flow between nodes.</li> <li><b>Convolutional (b):</b> Shows a convolutional layer processing a sequence of inputs. It includes a main layer of nodes and a stride layer below it, with skip connections (dashed blue arrows) bypassing intermediate layers.</li> <li><b>Recurrent (c):</b> Shows a recurrent layer where each node receives input from the previous node in the sequence, indicated by blue arrows pointing left.</li> <li><b>Graph convolutional (d):</b> Shows a graph convolutional layer processing a graph structure. It consists of multiple parallel layers of nodes, with colored arrows indicating information flow between nodes in adjacent layers across the graph structure.</li> </ul>				
<p>Legend:</p> <ul style="list-style-type: none"> <li>→ Parameterized information flow</li> <li>— Link</li> <li>○ Node in the neural network (scalar or tensor)</li> </ul>				

# Deep Learning in Genomics

- **Traditional Methods & Limitations**

- CNN (convolutional Neural Networks)
  - Takes the *local information* in developing each representation
  - Unable to capture the semantic dependency within long-range contexts
  - Its capability to extract local features is limited by the filter size

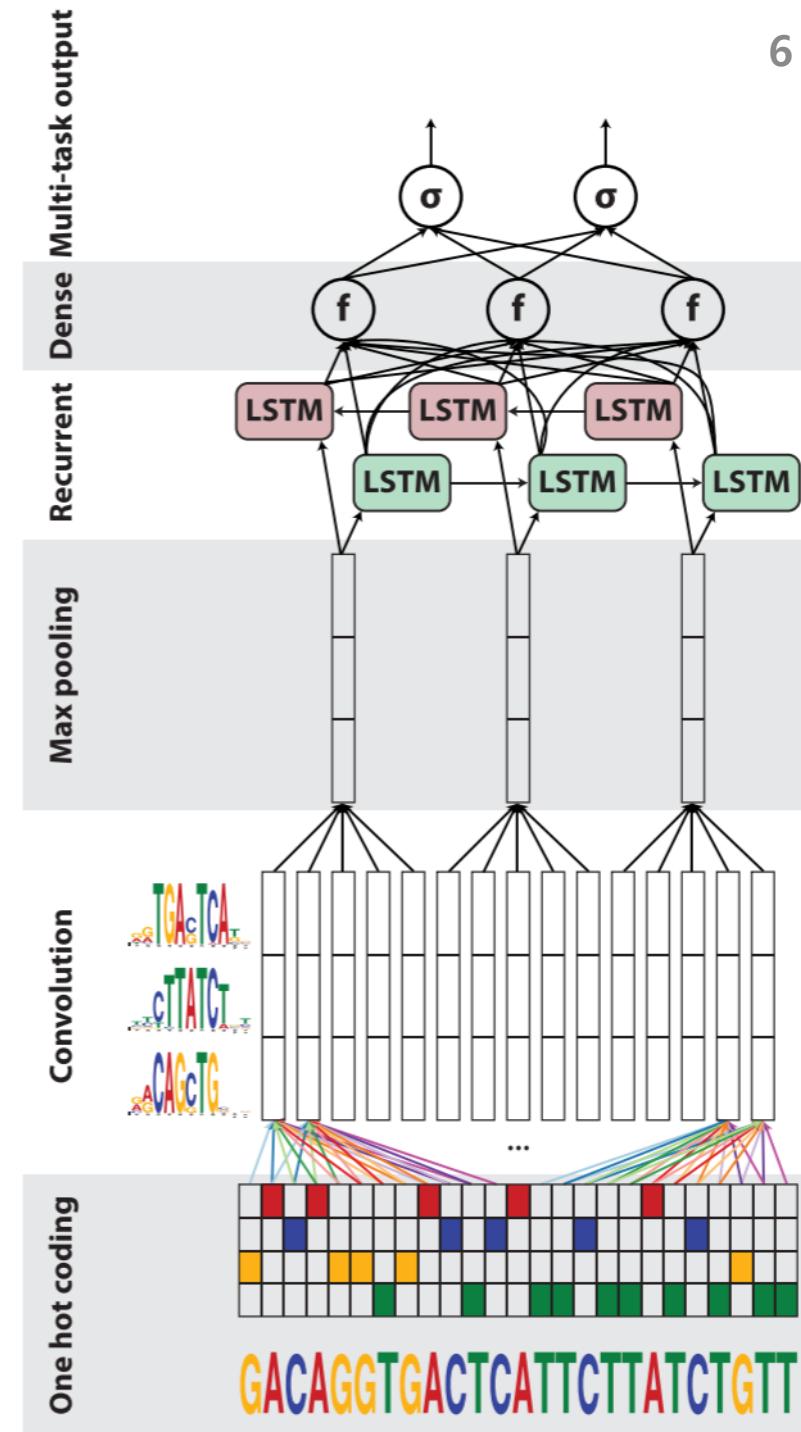


# Deep Learning in Genomics

- **Traditional Methods & Limitations**

- RNN (Recurrent Neural Networks)
  - Model long-range dependencies in sequences
  - Carry over information through long sequences via memory
  - Propagates information through all hidden states
  - Suffer from vanishing gradient and low-efficiency problem  
when it sequentially process all past states and compresses contextual information into a bottleneck w/ long input sequences

Quang D, Xie X. DanQ: a hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences. Nucleic Acids Res. 2016 Jun 20;44(11):e107.



# Representation of Word

- **Sparse Representation**

- One-hot vector
- Disadvantage
  - High dimension (curse of dimensionality)
  - Cannot encode the semantic of word

**The cat sat on the mat**

The: [0 1 0 0 0 0]

cat: [0 0 1 0 0 0]

sat: [0 0 0 1 0 0]

on: [0 0 0 0 1 0]

the: [0 0 0 0 0 1 0]

mat: [0 0 0 0 0 0 1]

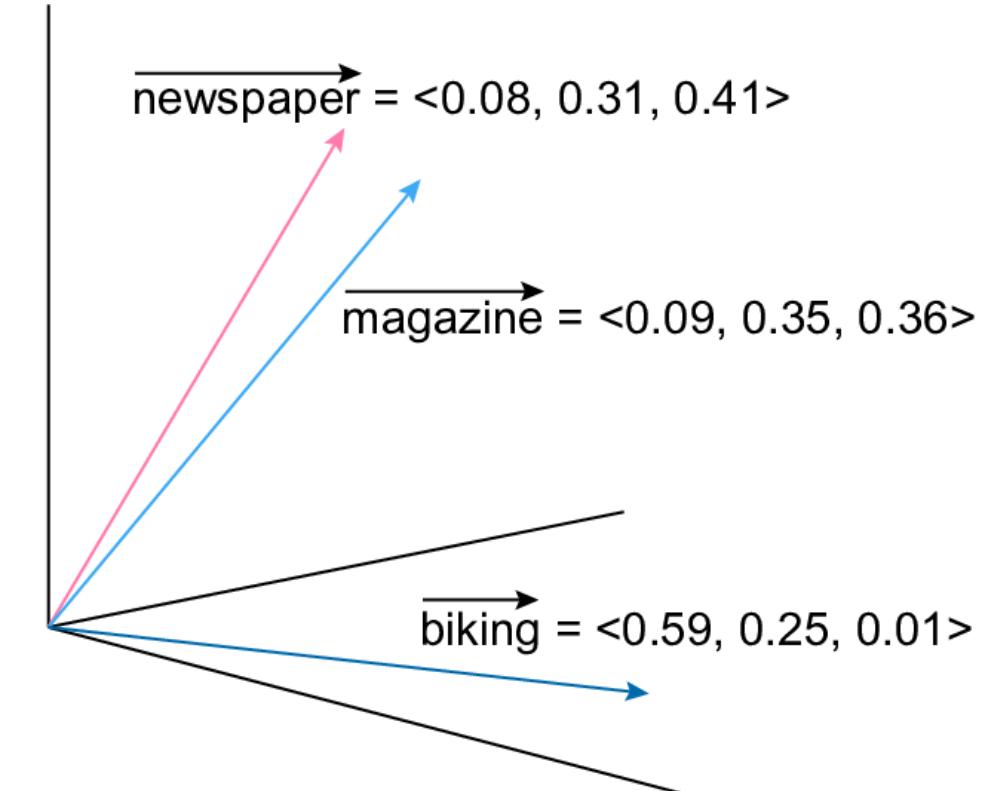
# Representation of Word

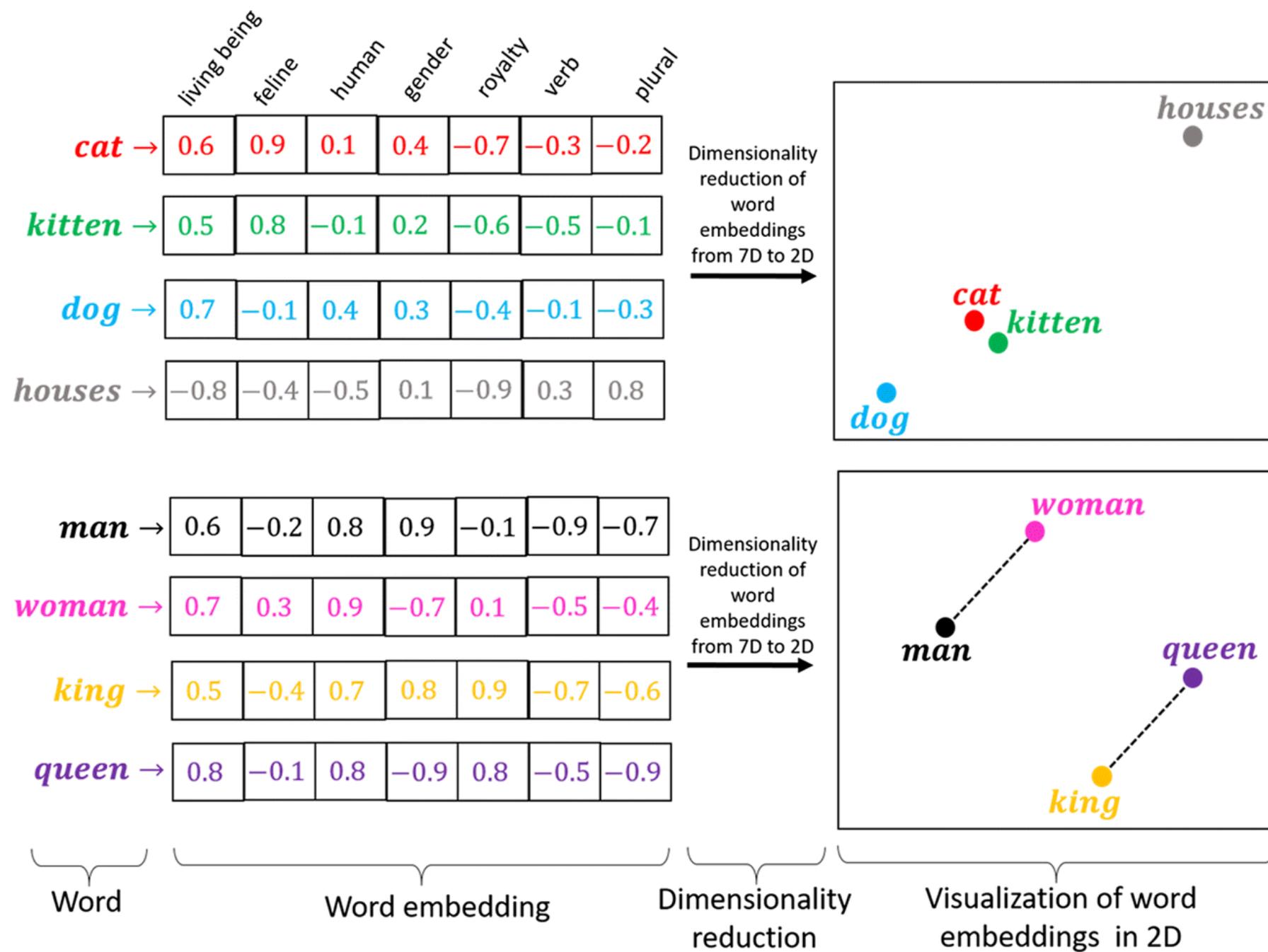
- Sparse Representation
  - One-hot vector
    - ex) DNA sequence

C	G	A	T	A	A	C	C	G	A	T	A	T
A	0	0	1	0	1	1	0	0	1	0	1	0
C	1	0	0	0	0	0	1	0	0	0	0	0
G	0	1	0	0	0	0	0	1	0	0	0	0
T	0	0	0	1	0	0	0	0	0	1	0	1

# Representation of Word

- **Dense Representation**
  - **Word embedding (representation, vector)**
    - Dense vector for each word
    - Distributed representation of word
    - Represent a point in some sort of “word” space
    - Sufficient to encode all semantics of the word
    - Each dimension would encode some meaning
  - GloVe, FastText, Word2Vec ...

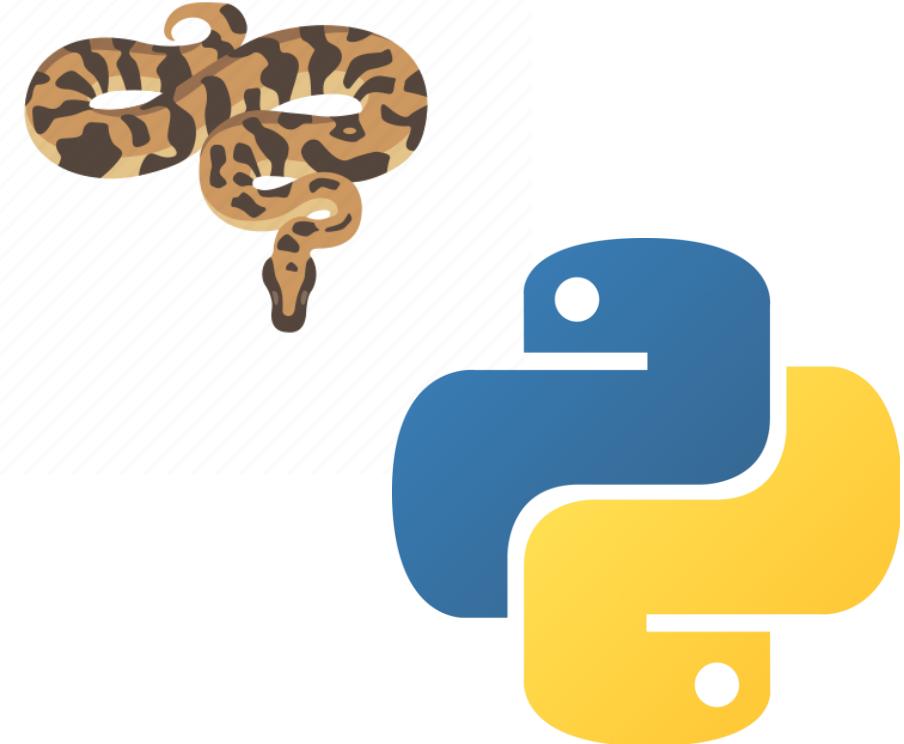




# Representation of Word

- **Dense Representation**

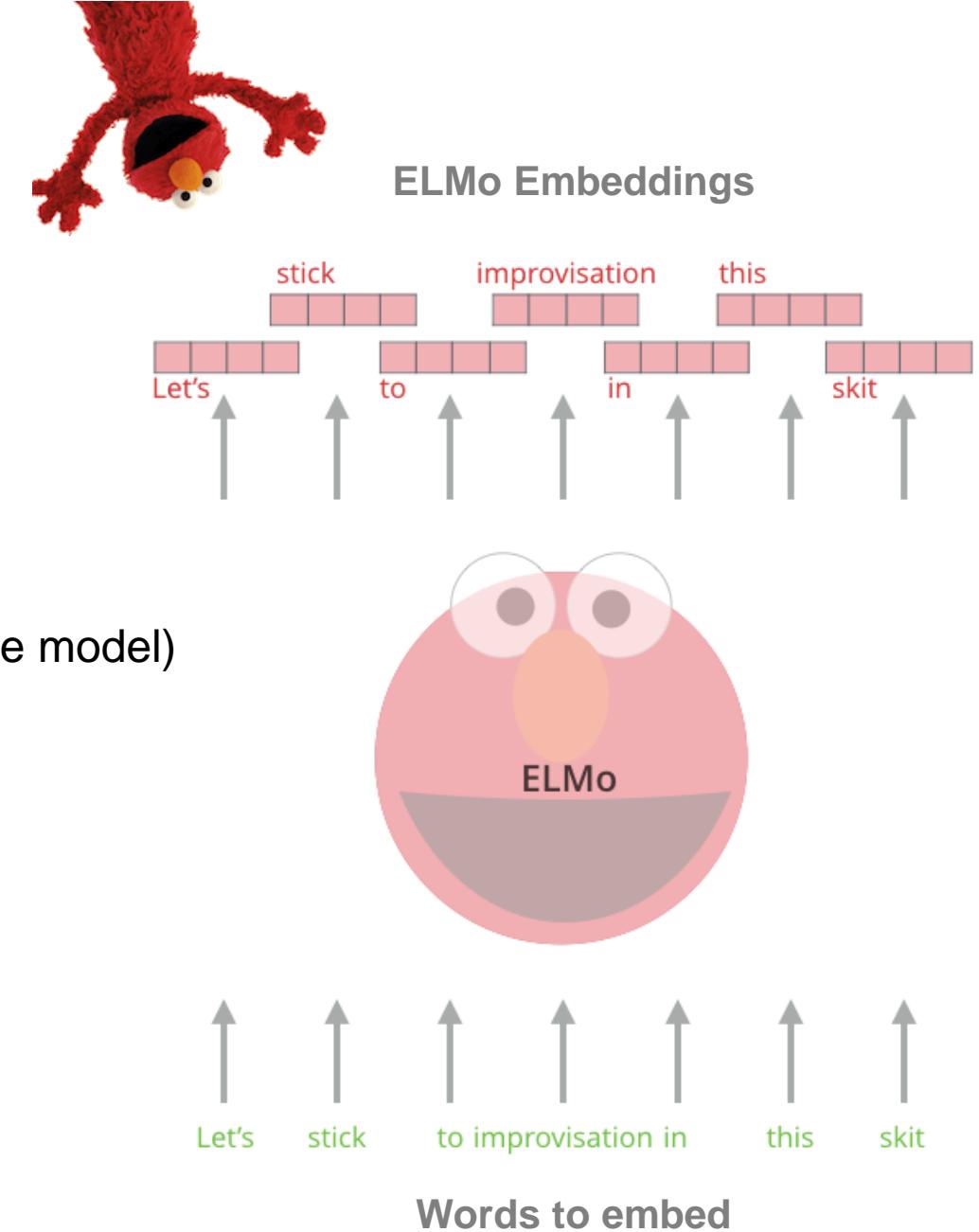
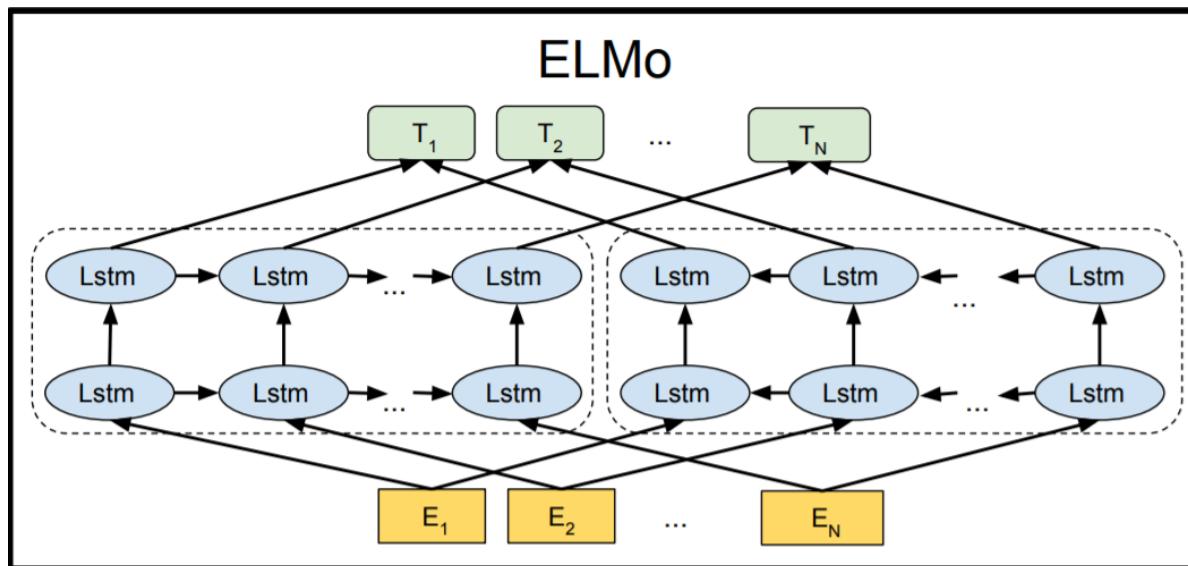
- Limitation of **static (context-free)** word embedding
  - We use pre-defined word embedding
    - Cannot distinguish homonym / polysemy
  - ex) “He got bit by Python.  
“Python is my favorite programming language.”
    - Word embedding of “Python” would be same  
in both sentences
  - It ignores the context and always give the same embedding irrespective of context



# Representation of Word

- Dense Representation**

- Dynamic (Context-based) Word Embedding**
    - ELMo (embeddings from language model)
      - From bidirectional language model  
(forward language model + backward language model)



# Attention, Transformer, and BERT

# No More RNN, Only Attention

## Attention Is All You Need

**Ashish Vaswani\***  
Google Brain  
[avaswani@google.com](mailto:avaswani@google.com)

**Noam Shazeer\***  
Google Brain  
[noam@google.com](mailto:noam@google.com)

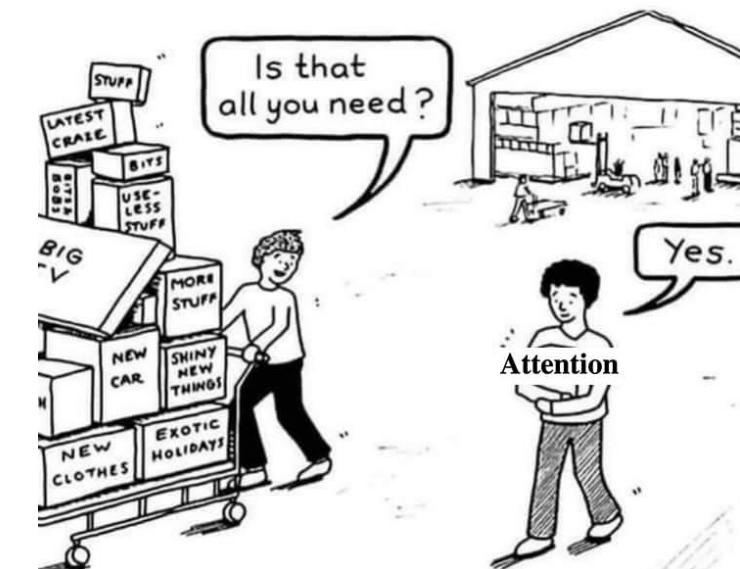
**Niki Parmar\***  
Google Research  
[nikip@google.com](mailto:nikip@google.com)

**Jakob Uszkoreit\***  
Google Research  
[usz@google.com](mailto:usz@google.com)

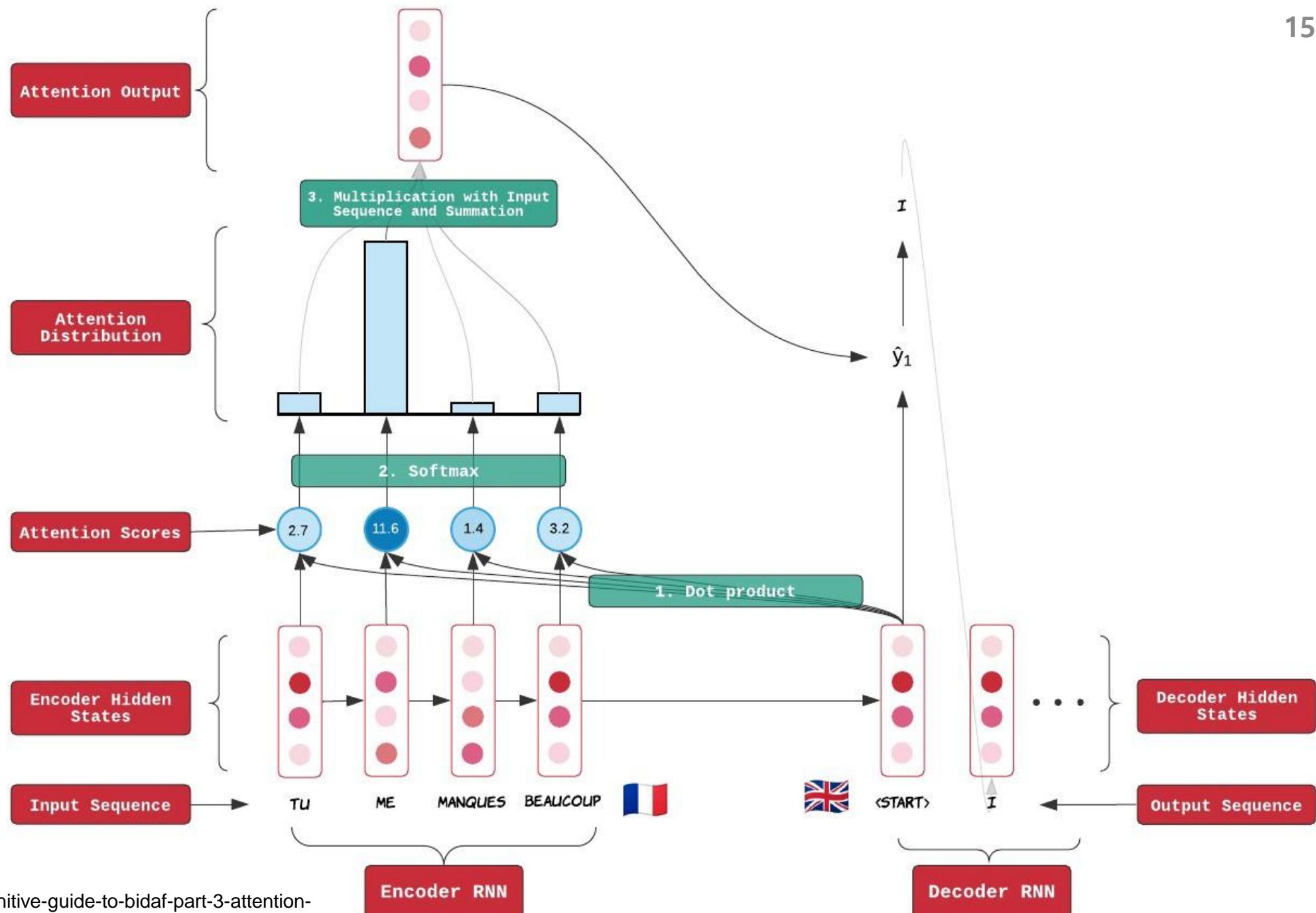
**Llion Jones\***  
Google Research  
[llion@google.com](mailto:llion@google.com)

**Aidan N. Gomez\* †**  
University of Toronto  
[aidan@cs.toronto.edu](mailto:aidan@cs.toronto.edu)

**Lukasz Kaiser\***  
Google Brain  
[lukaszkaiser@google.com](mailto:lukaszkaiser@google.com)



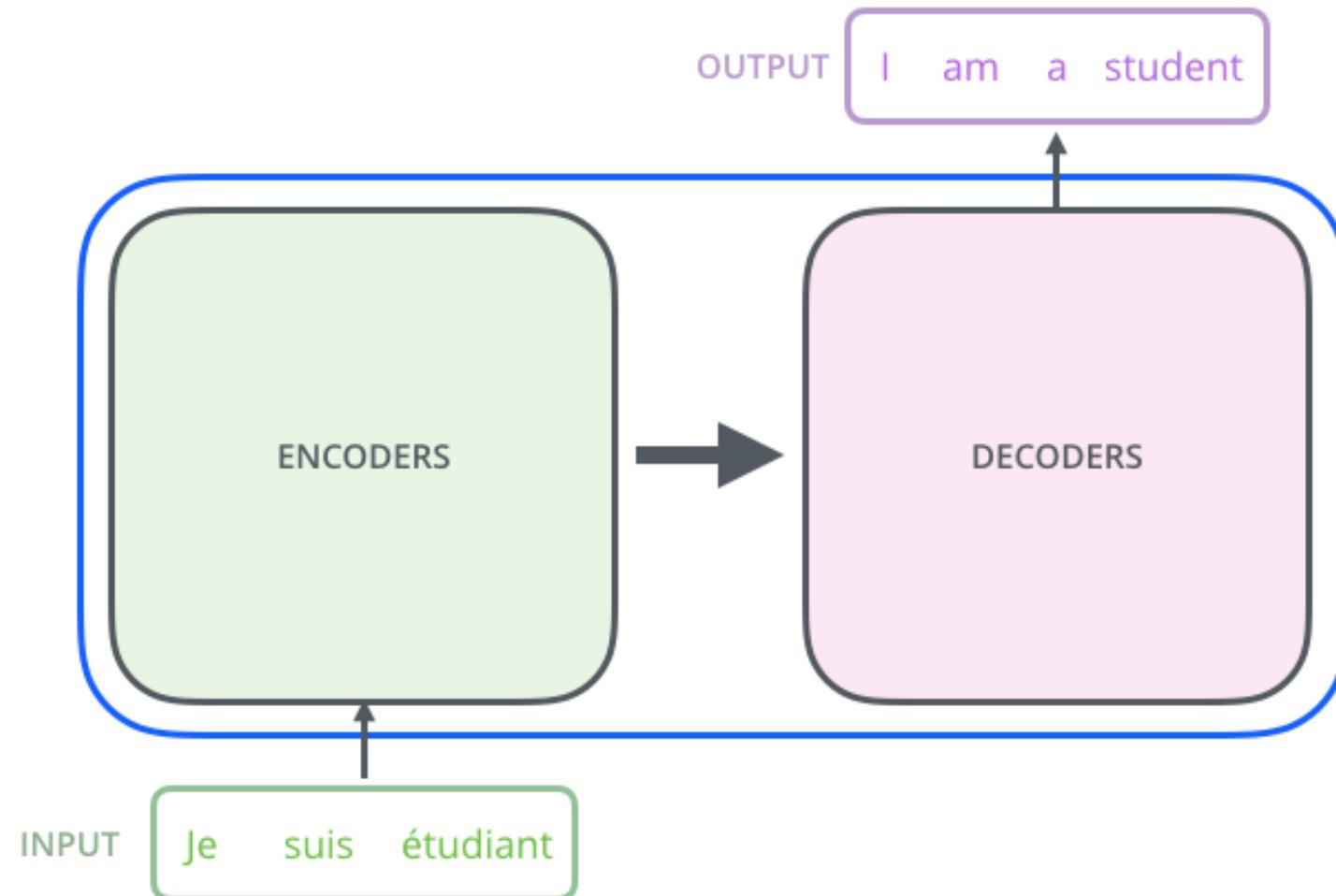
# Attention



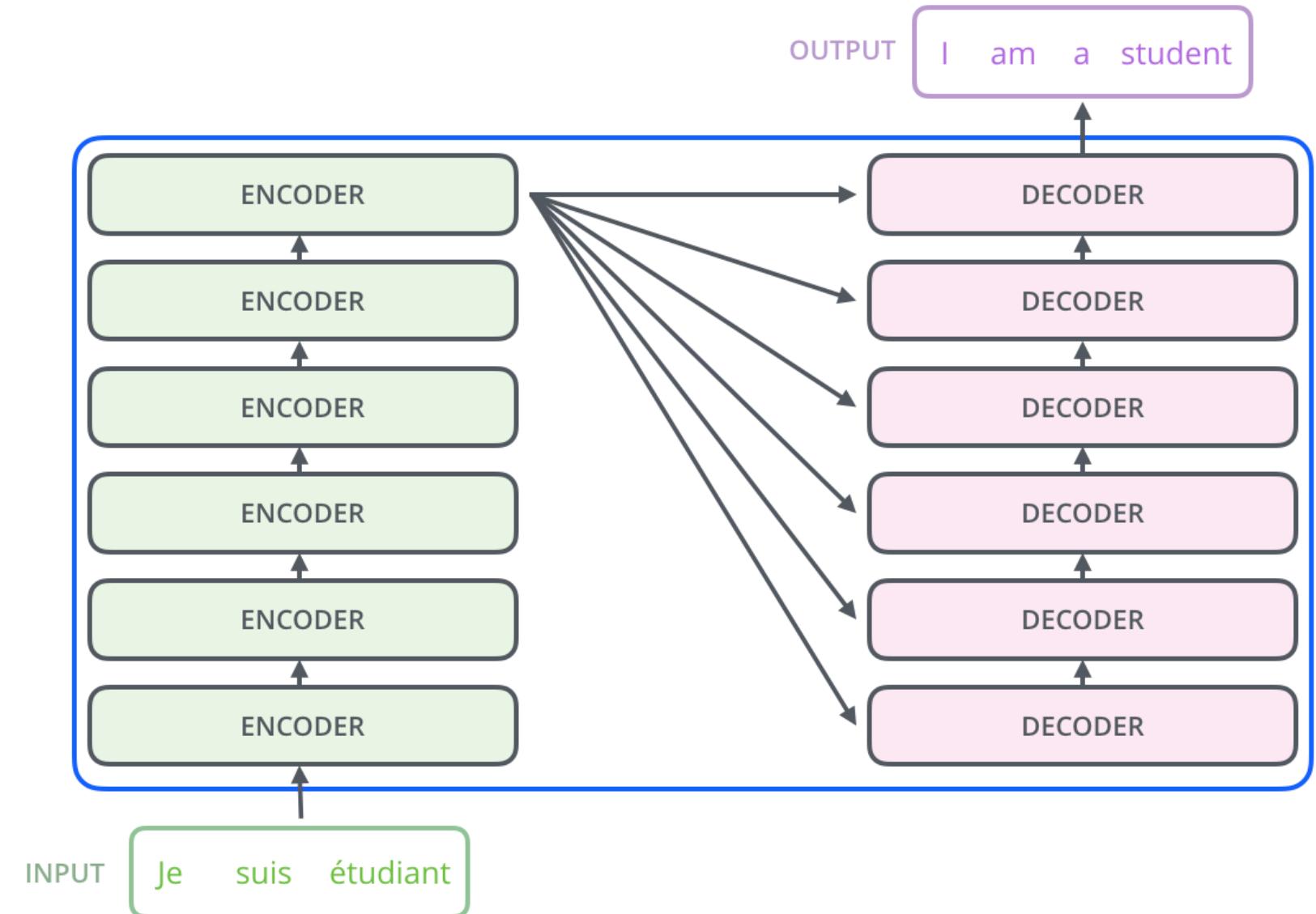
# Transformer



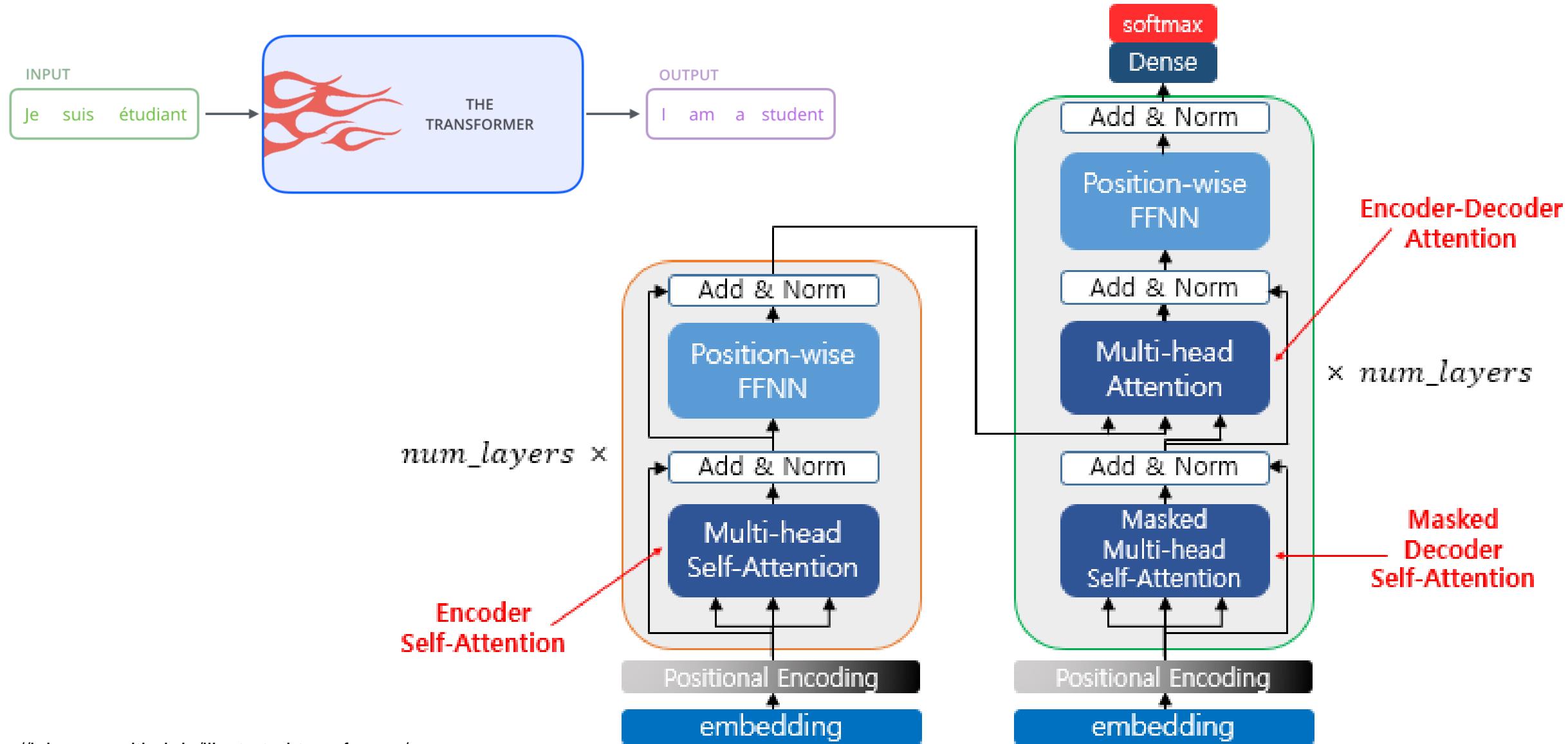
# Transformer



# Transformer

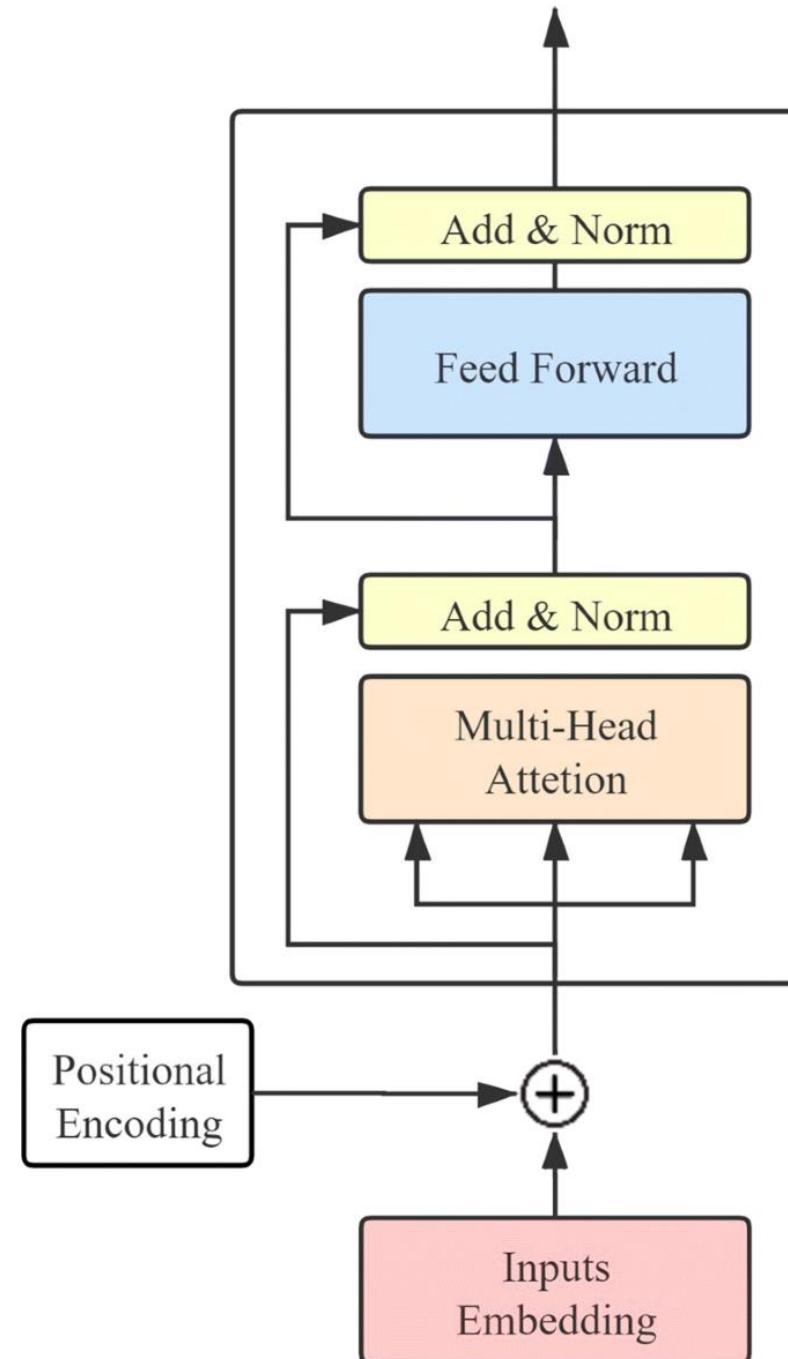


# Transformer



# Encoder of Transformer

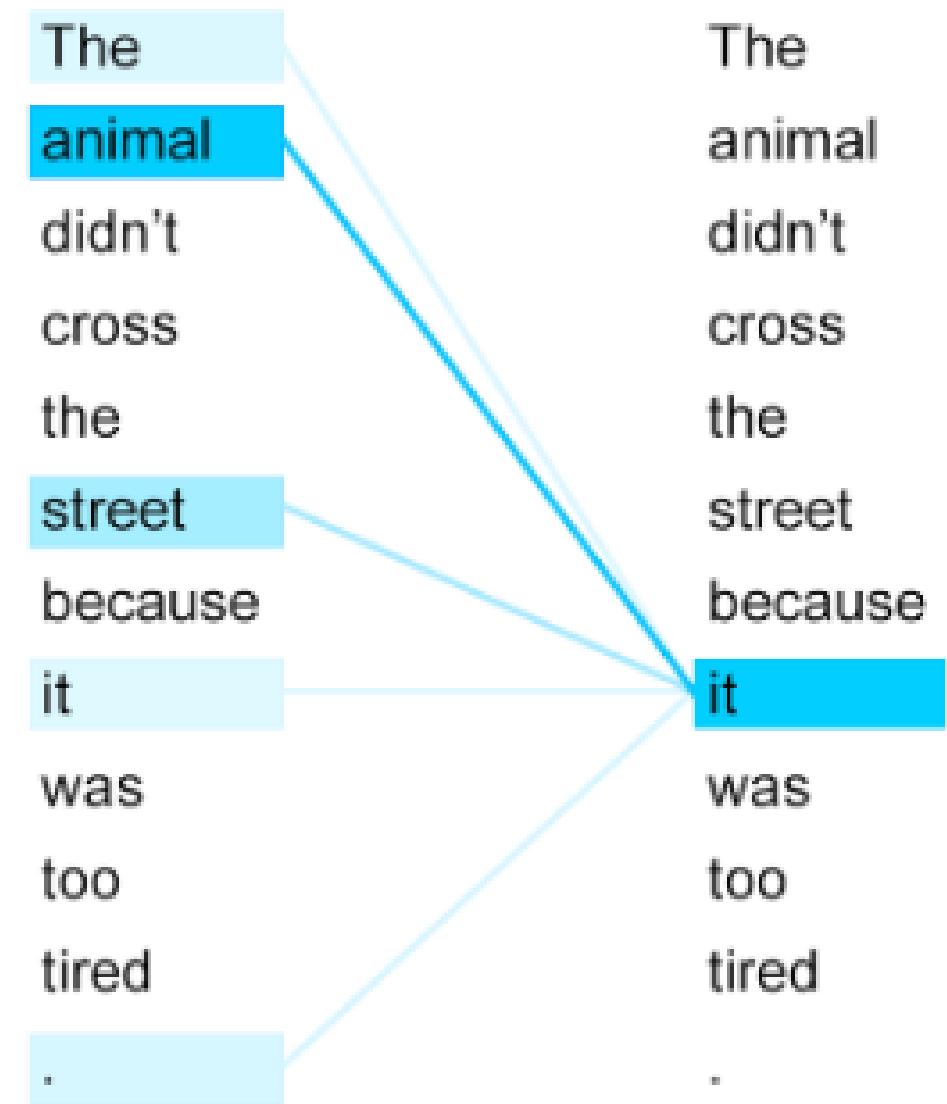
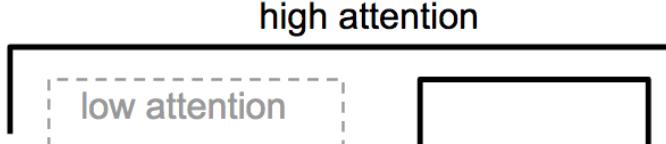
- 1) Multi-head attention
- 2) Feedforward network



# Self-Attention

- Example

She is eating a green apple.



# Self-Attention

- Input Matrix (embedding matrix / input embedding)

- ex) “I am good”
  - 1<sup>st</sup> row : embedding of “I”
  - 2<sup>nd</sup> row : embedding of “am”
  - 3<sup>rd</sup> row : embedding of “good”

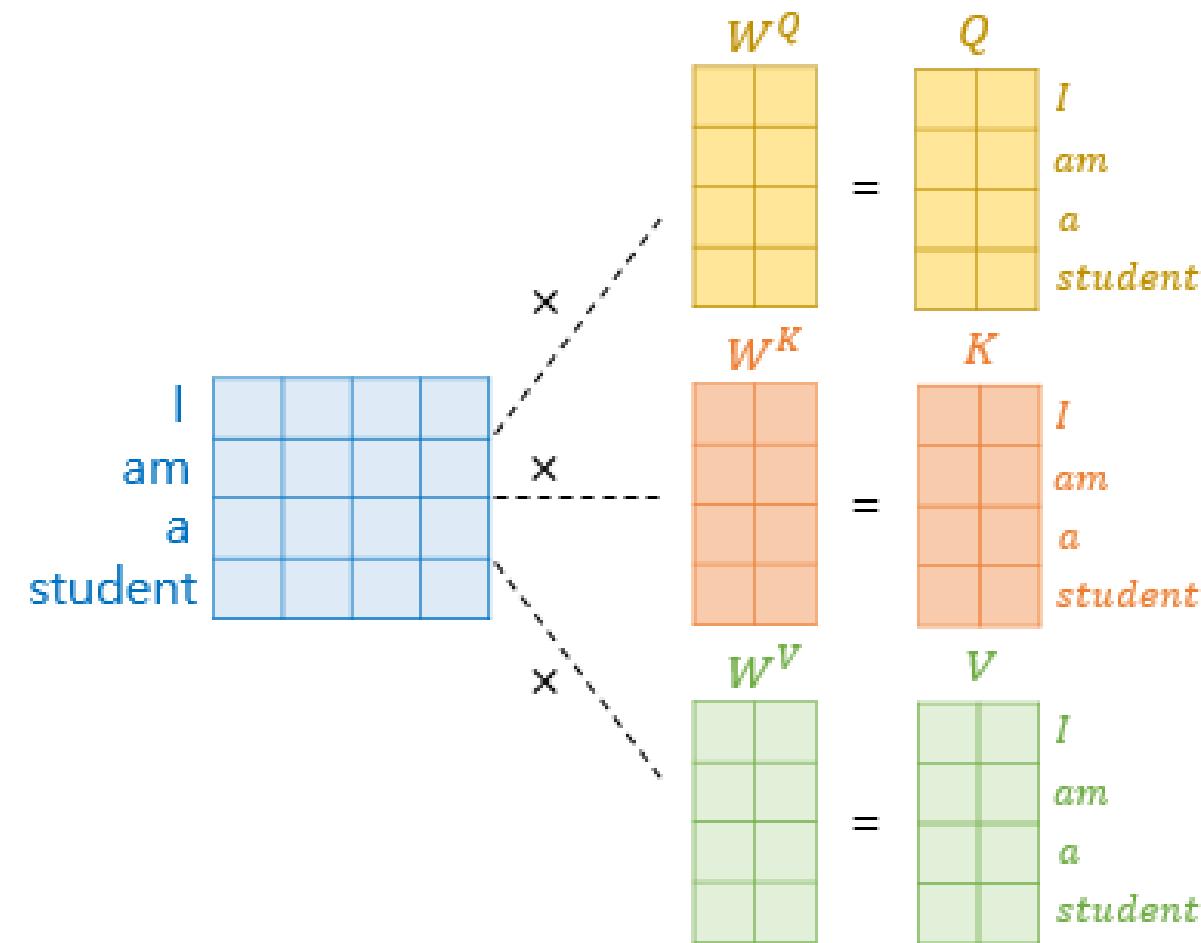
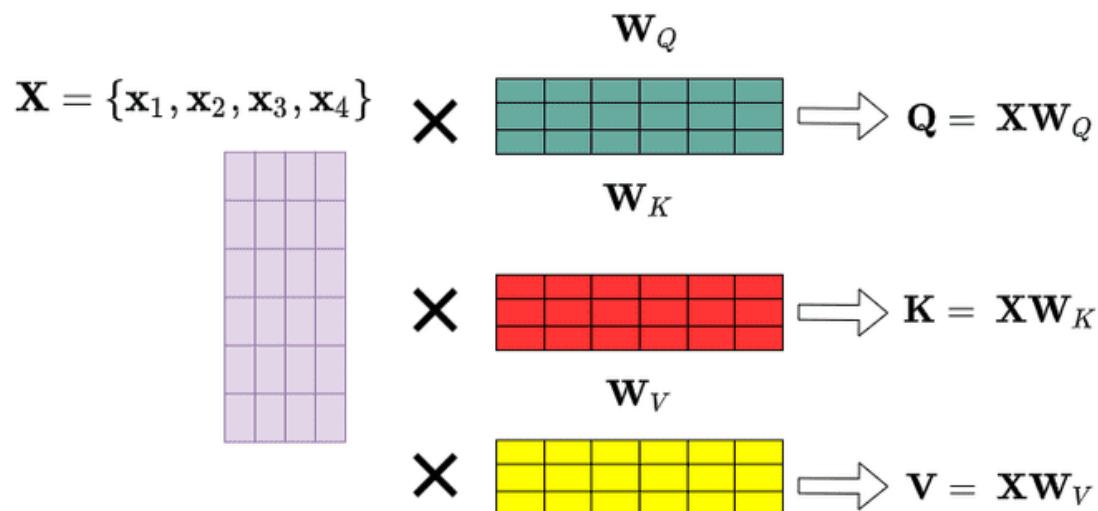
I	1.76	2.22	...	6.66	$x_1$
am	7.77	0.631	...	5.35	$x_2$
good	11.44	10.10	...	3.33	$x_3$

$\times$   
input matrix  
(embedding matrix)

3x512

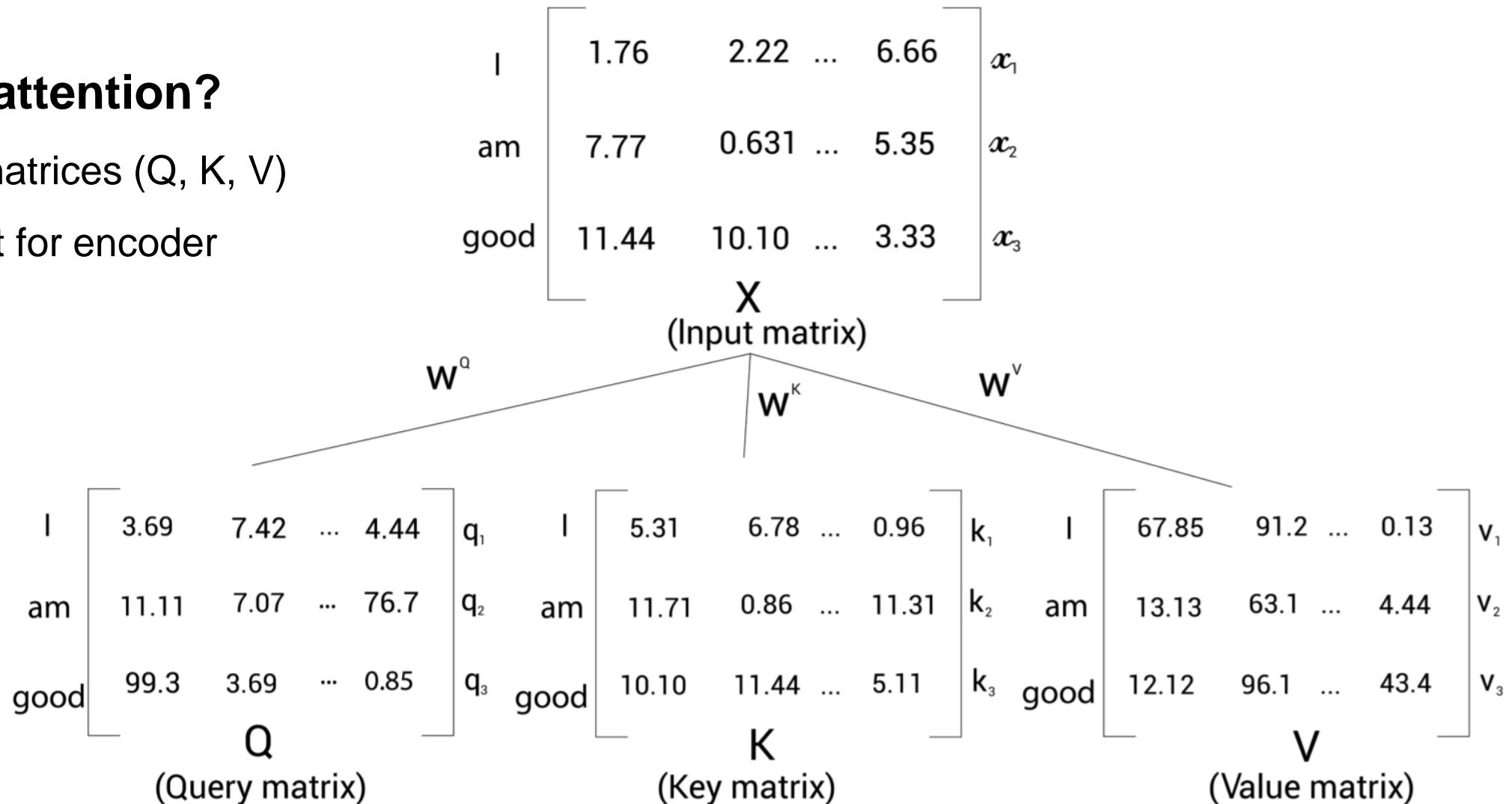
# Self-Attention

- Query (Q), Key (K), Value (V) matrix
  - $W^Q, W^K, W^V \rightarrow$  Multiplied to input matrix  $X$ 
    - 1<sup>st</sup> row : Query, Key, Value of first word
    - 2<sup>nd</sup> row : Query, Key, Value of second word
  - $Q = XW^Q, K = XW^K, V = XW^V$

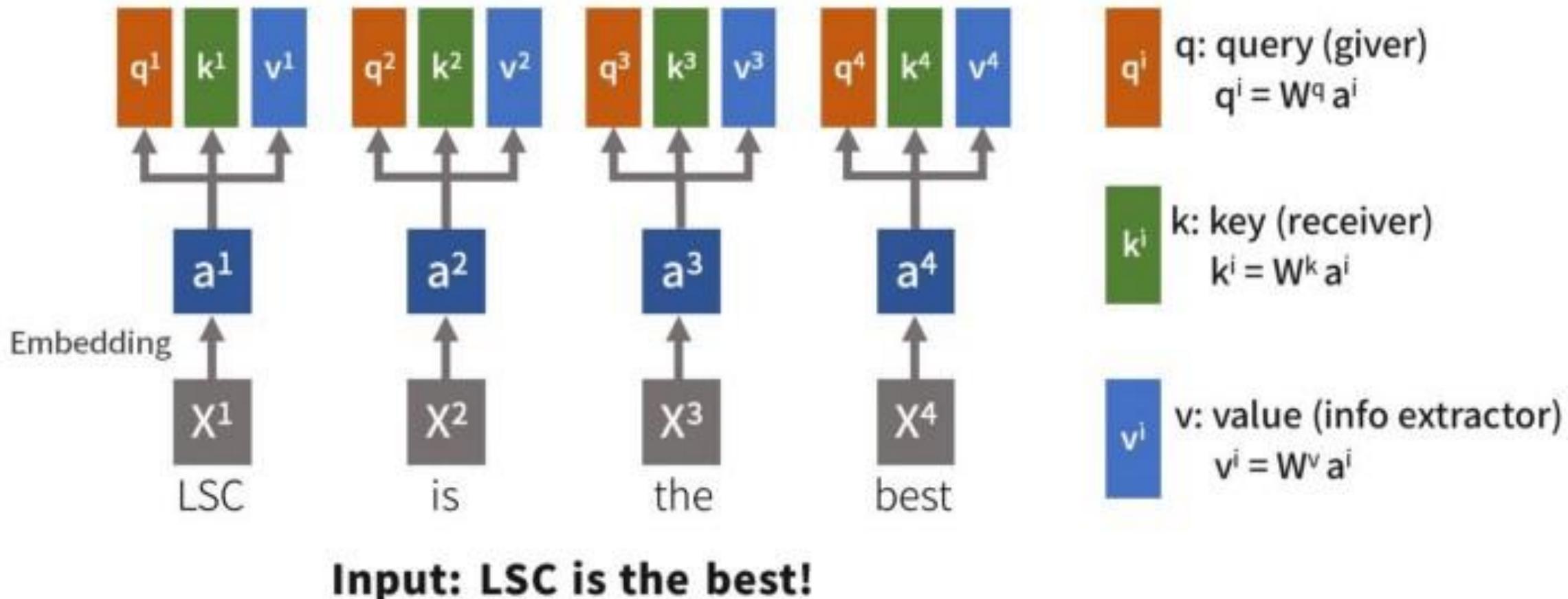


# Self-Attention

- Why **self-attention**?
  - All the matrices ( $Q$ ,  $K$ ,  $V$ ) are input for encoder



# Self-Attention



# Self-Attention

- Dot product Query & Key

- $QK^T$

$$\begin{array}{c}
 \begin{array}{c} | \\ I \\ am \\ good \end{array} \quad \begin{array}{c} | \\ q_1 \\ q_2 \\ q_3 \end{array} \\
 \begin{array}{ccccc} 3.69 & 7.42 & \dots & 4.44 & | \\ 11.11 & 7.07 & \dots & 76.7 & am \\ 99.3 & 3.69 & \dots & 0.85 & good \end{array} \quad \cdot \quad \begin{array}{c} | \\ k_1 \\ k_2 \\ k_3 \end{array} \\
 Q \qquad \qquad \qquad K^T
 \end{array} = \begin{array}{c}
 \begin{array}{c} | \\ I \\ am \\ good \end{array} \quad \begin{array}{c} | \\ q_1.k_1 \\ q_2.k_1 \\ q_3.k_1 \end{array} \quad \begin{array}{c} | \\ q_1.k_2 \\ q_2.k_2 \\ q_3.k_2 \end{array} \quad \begin{array}{c} | \\ q_1.k_3 \\ q_2.k_3 \\ q_3.k_3 \end{array} \\
 \begin{array}{ccccc} 5.31 & 11.71 & 10.10 & | \\ 6.78 & 0.86 & 11.44 & am \\ 0.96 & 11.31 & 5.11 & good \end{array}
 \end{array}$$

$$QK^T = \begin{array}{c}
 \begin{array}{c} | \\ I \\ am \\ good \end{array} \quad \begin{array}{c} | \\ 110 \\ 70 \\ 90 \end{array} \quad \begin{array}{c} | \\ 90 \\ 99 \\ 70 \end{array} \quad \begin{array}{c} | \\ 80 \\ 70 \\ 100 \end{array} \\
 \begin{array}{ccccc} & & & & | \\ & & & & am \\ & & & & good \end{array}
 \end{array}$$

# Self-Attention

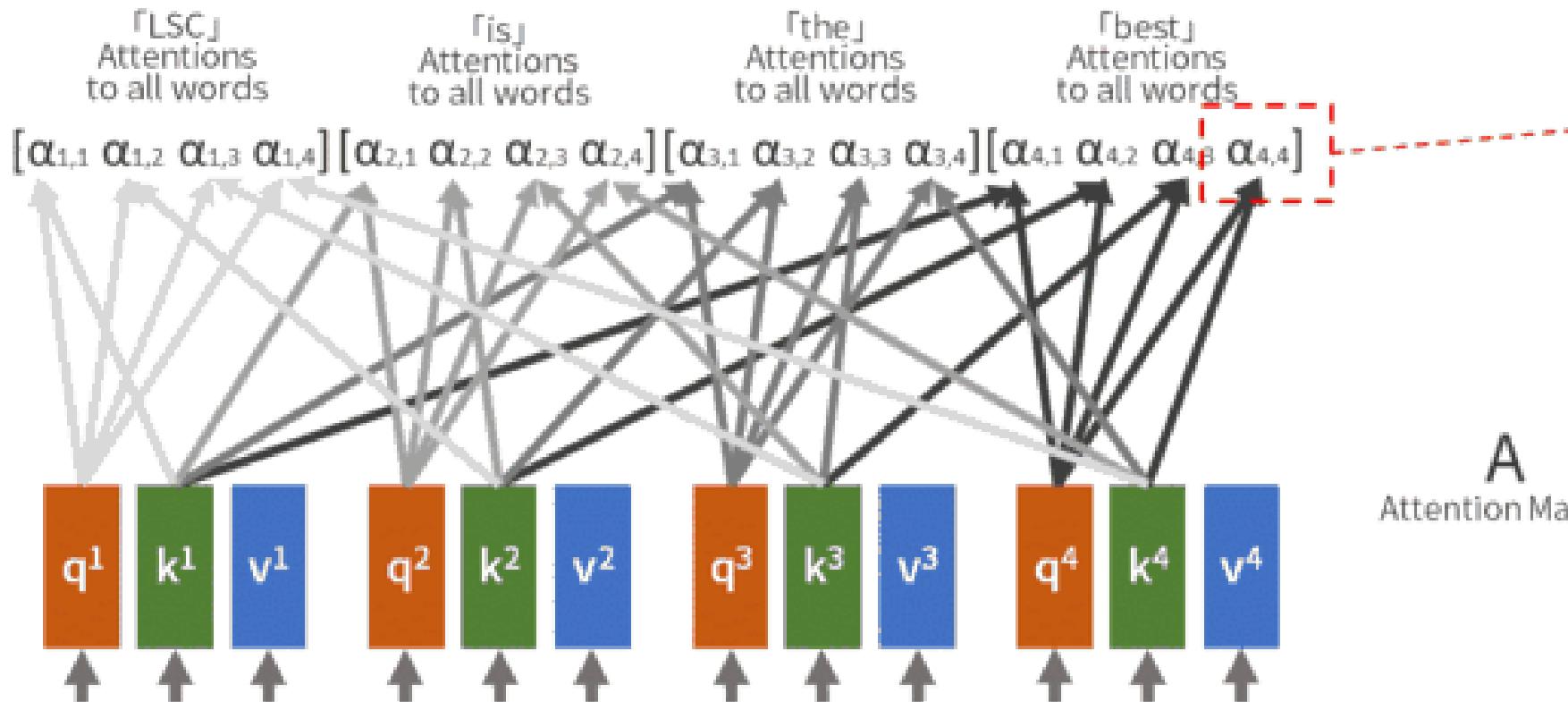
- Dot product Query & Key
  - $QK^T$ 
    - 1sr row : Similarity between  $q_1(I)$  &  $k_1(I)$ ,  $k_2(am)$ ,  $k_3(good)$

$$QK^T = \begin{array}{c|ccc} & I & am & good \\ \hline I & \boxed{110} & 90 & 80 \\ & q_1.k_1 & q_1.k_2 & q_1.k_3 \\ \hline am & 70 & 99 & 70 \\ good & 90 & 70 & 100 \end{array}$$

# Self-Attention

- Dot product Query & Key

- $QK^T$



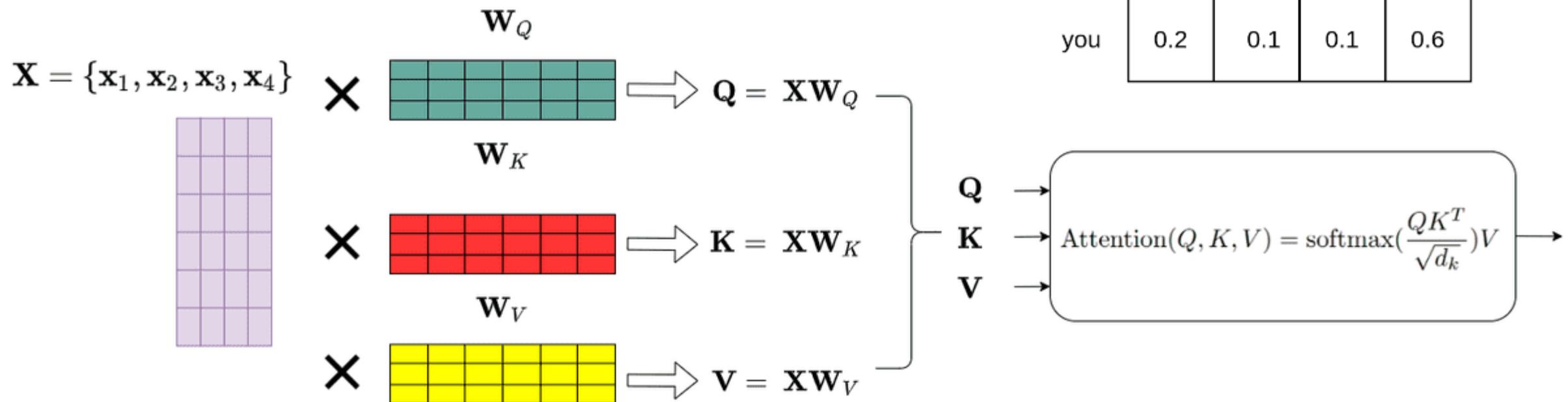
$$\alpha_{i,j} = \frac{q^i \cdot k^j}{\sqrt{d}}$$

d: dimension of q, k

$$A = \begin{bmatrix} \alpha_{1,1} & \alpha_{1,2} & \alpha_{1,3} & \alpha_{1,4} \\ \alpha_{2,1} & \alpha_{2,2} & \alpha_{2,3} & \alpha_{2,4} \\ \alpha_{3,1} & \alpha_{3,2} & \alpha_{3,3} & \alpha_{3,4} \\ \alpha_{4,1} & \alpha_{4,2} & \alpha_{4,3} & \alpha_{4,4} \end{bmatrix}$$

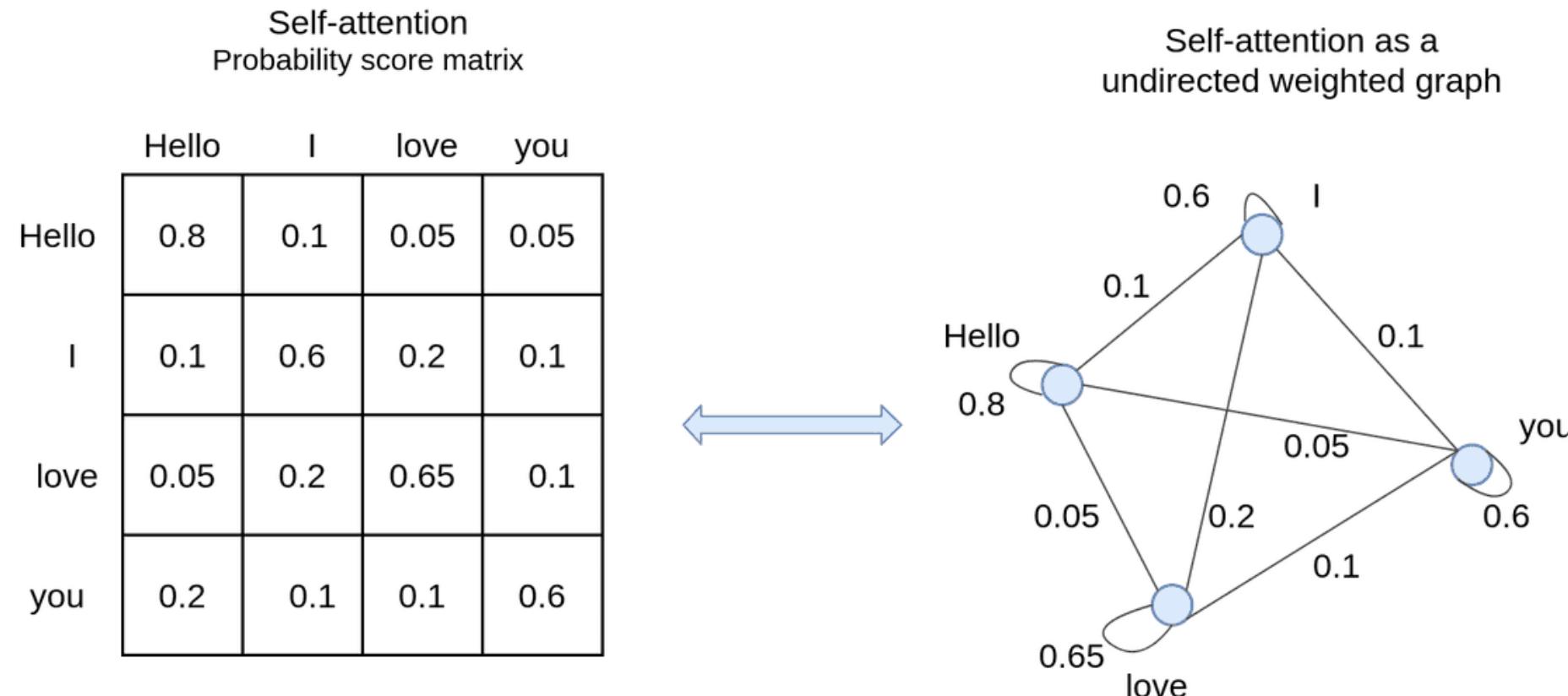
# Self-Attention

- Score matrix =  $\text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right)$



# Self-Attention

- Score matrix =  $softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)$ 
  - Self-attention as an undirected weighted graph**



# Self-Attention

- **Attention Matrix**

- = Score Matrix  $\times$  V (value)

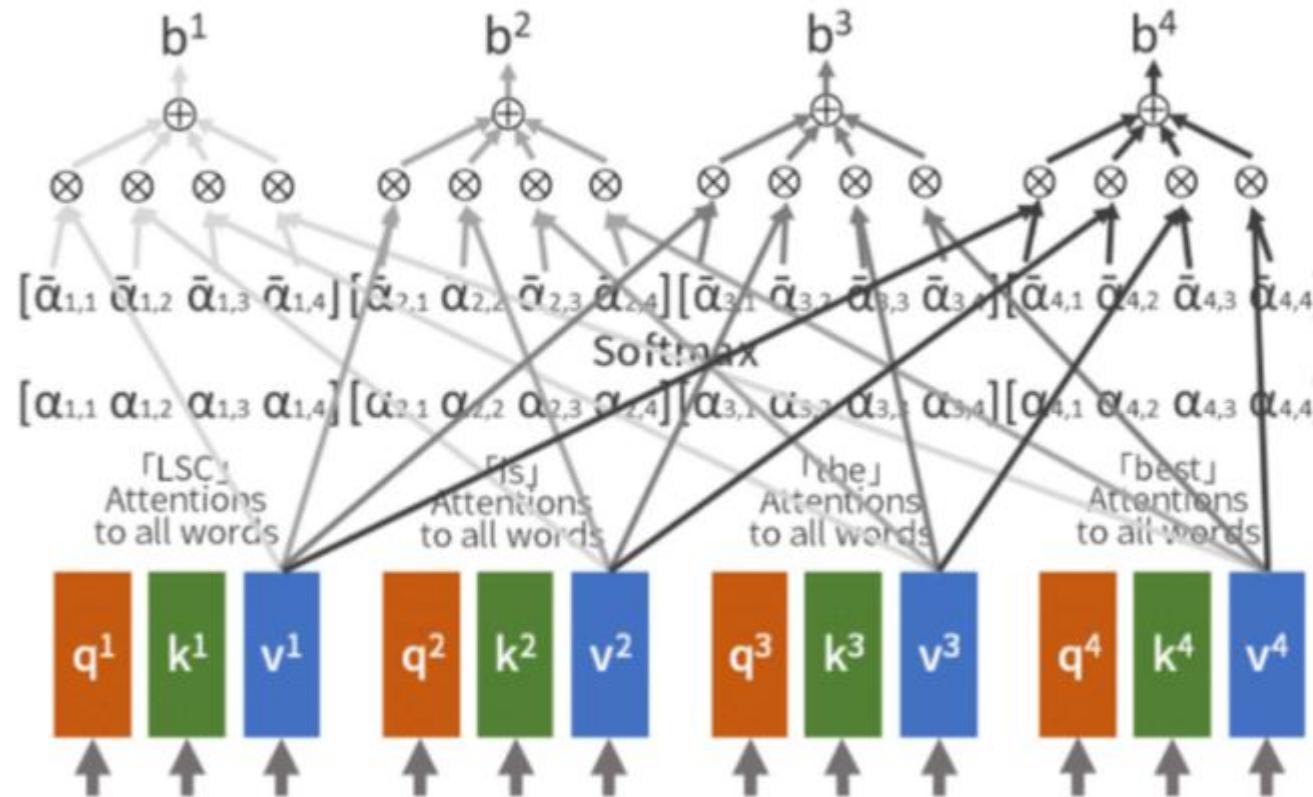
$$\text{softmax} \left( \frac{Q \times K^T}{\sqrt{d_k}} \right) \times V = \text{Attention Value Matrix } \alpha$$

- $Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V$

# Self-Attention

- **Attention Matrix  $Z$**

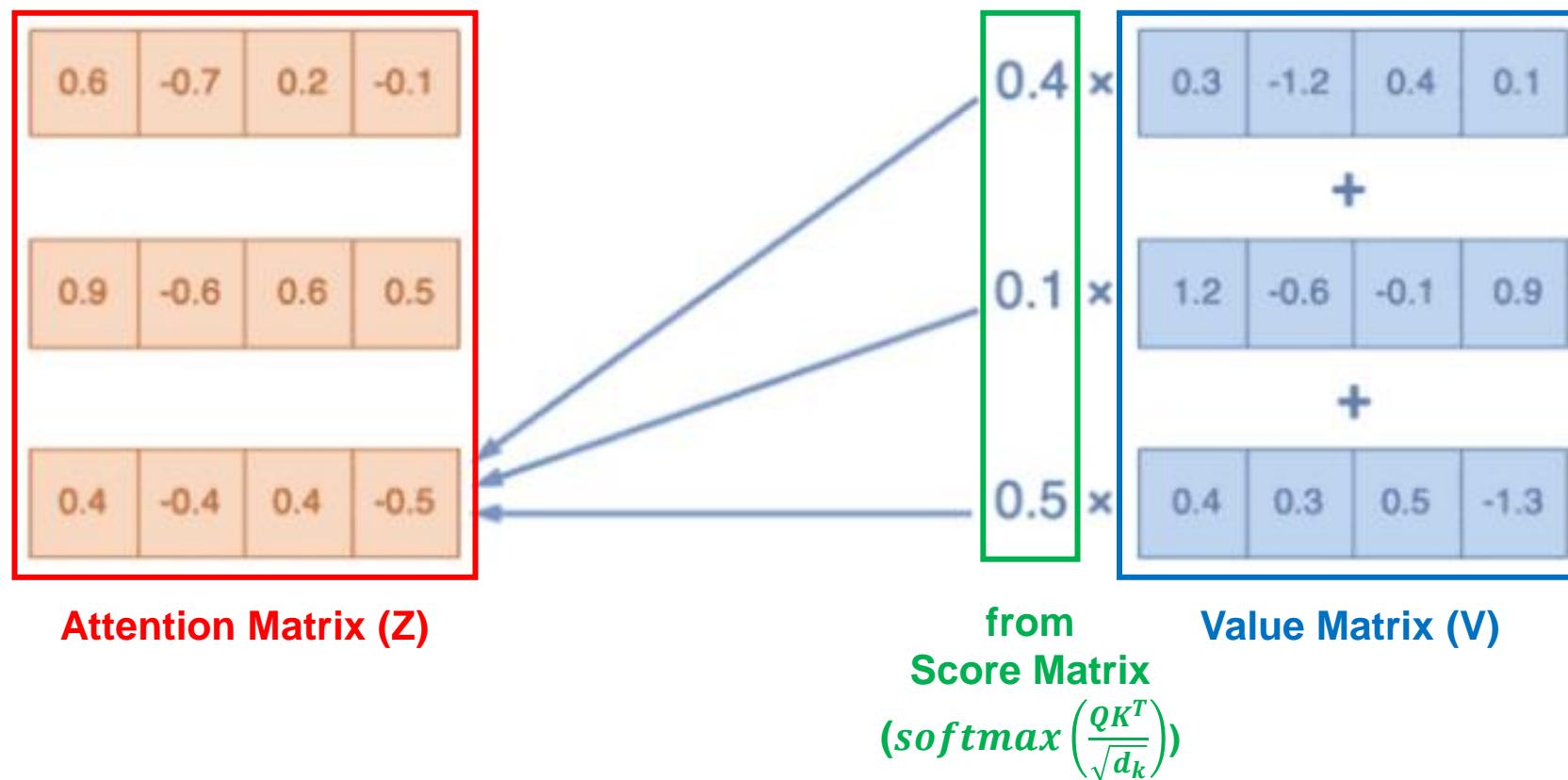
- = Score Matrix  $\times$  V (value)



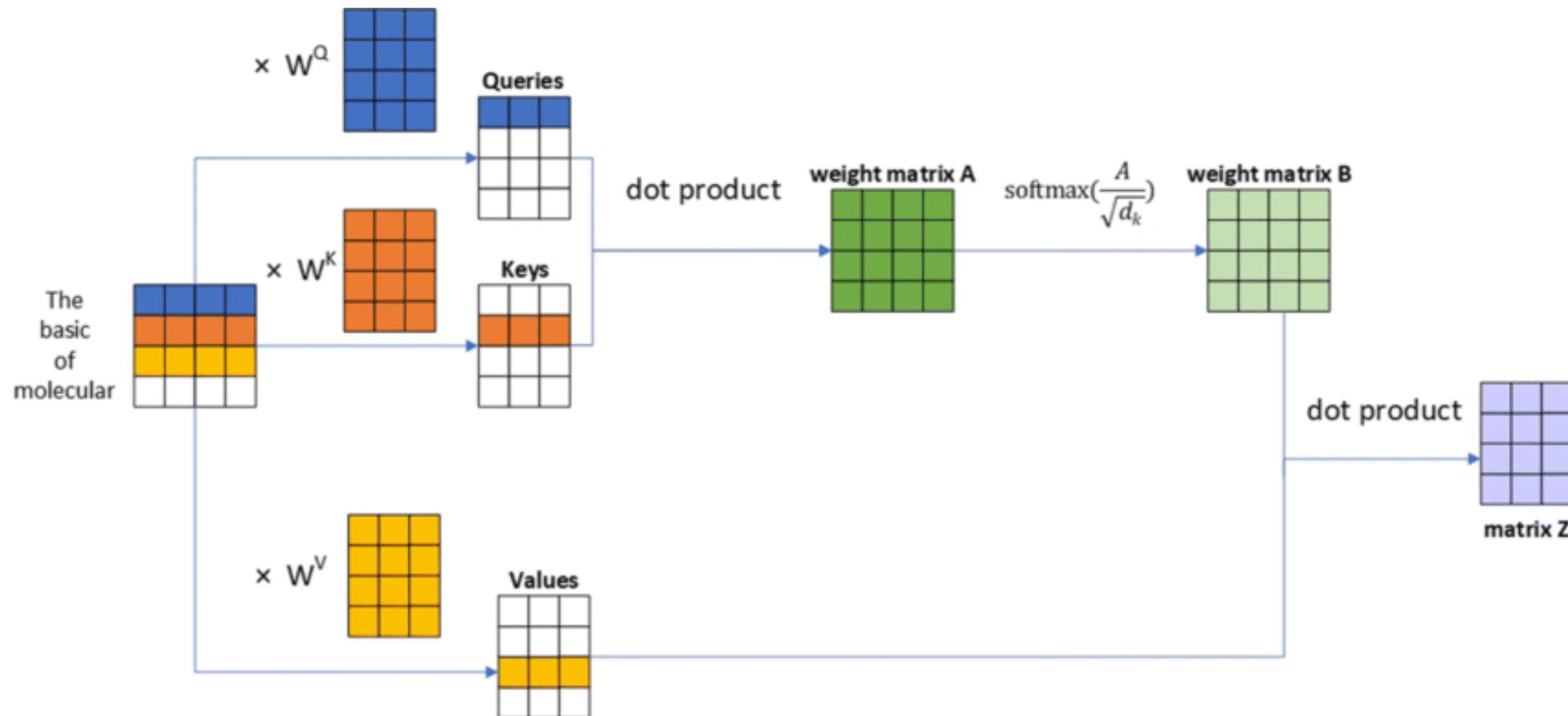
$$b^i = \sum_j \bar{\alpha}_{i,j} v^j$$

# Self-Attention

- Attention is just a fancy name for weighted average
  - **Each row of attention matrix (Z) = Weighted average of rows of Value matrix (V)**

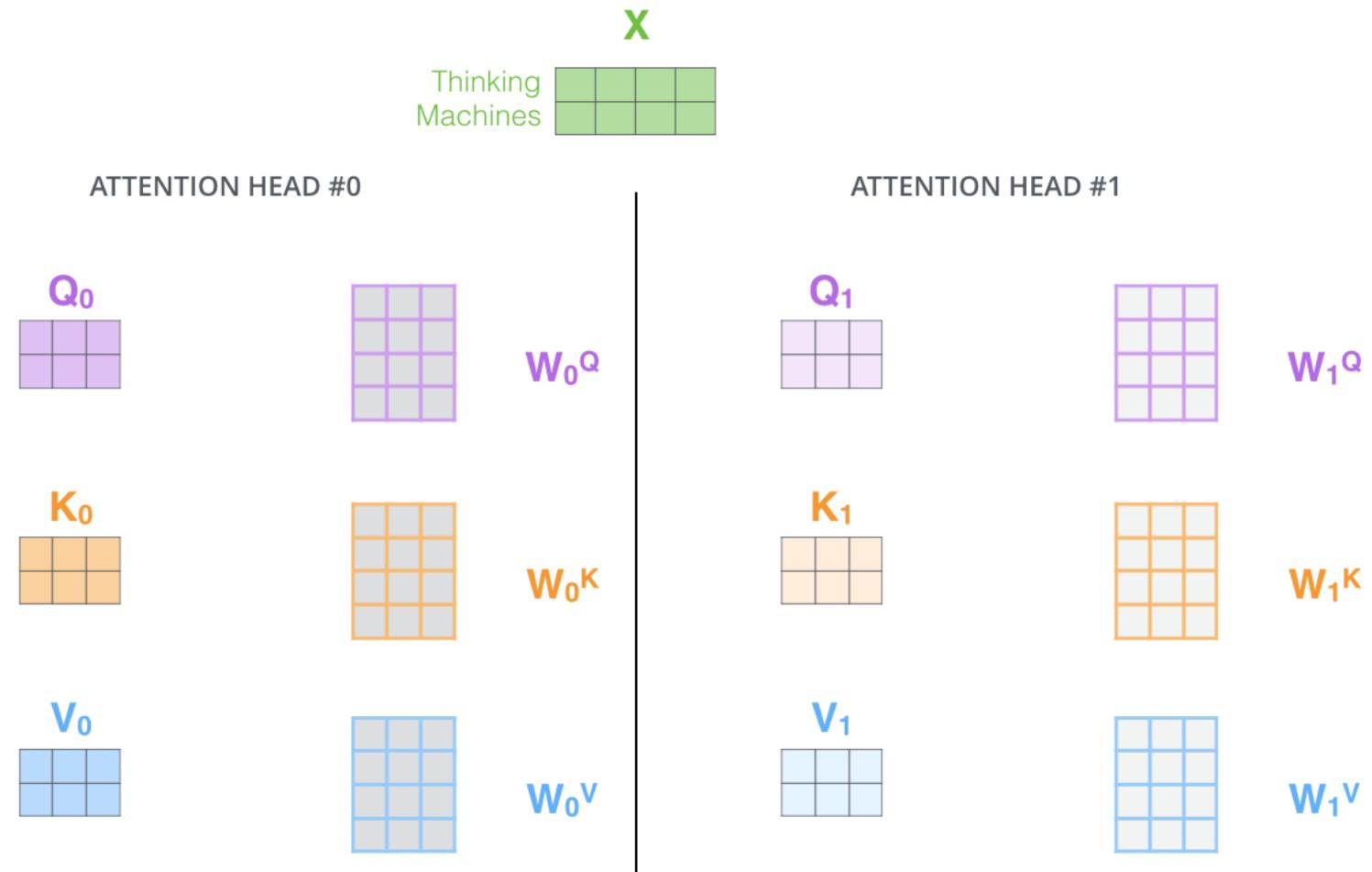


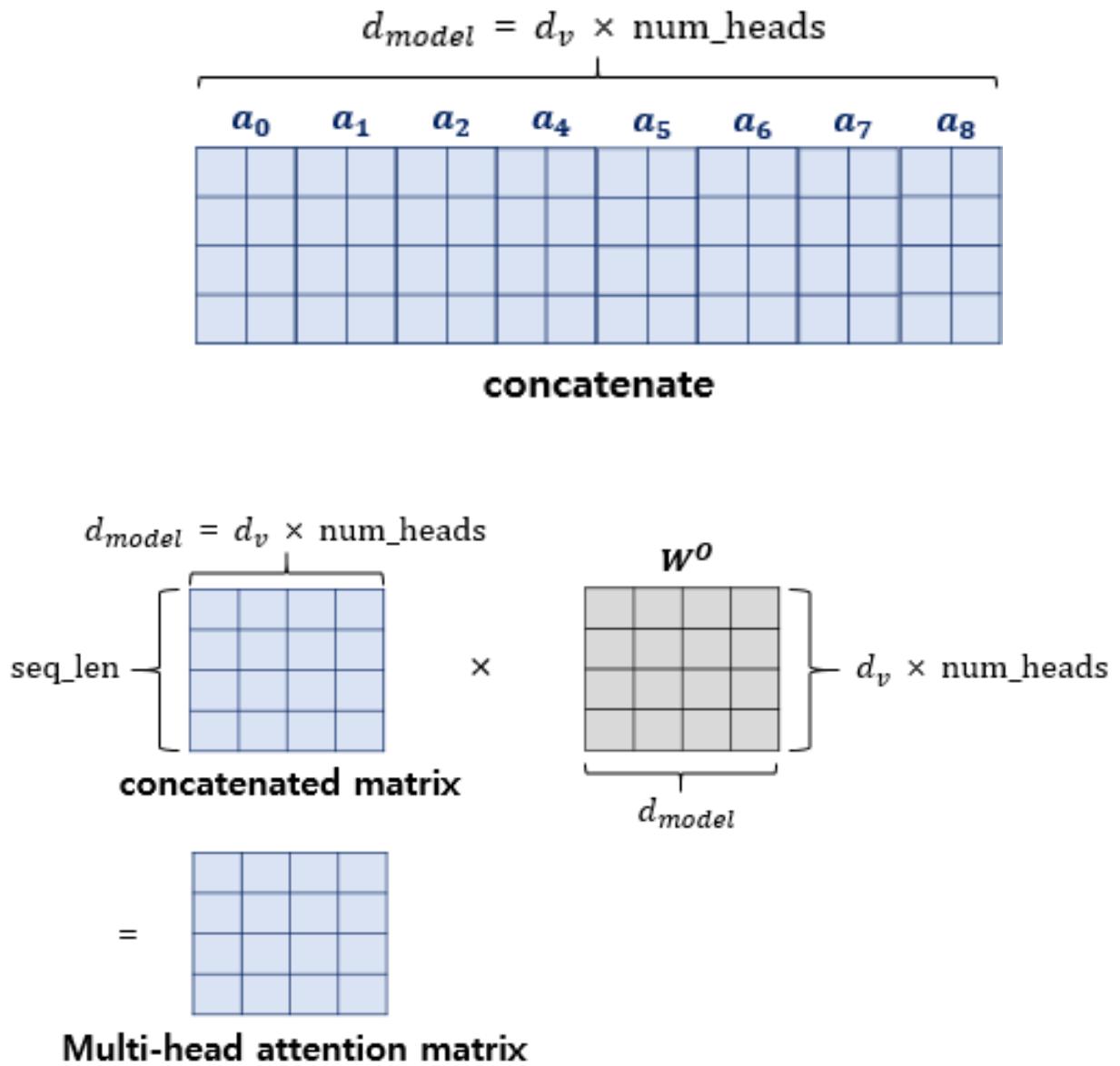
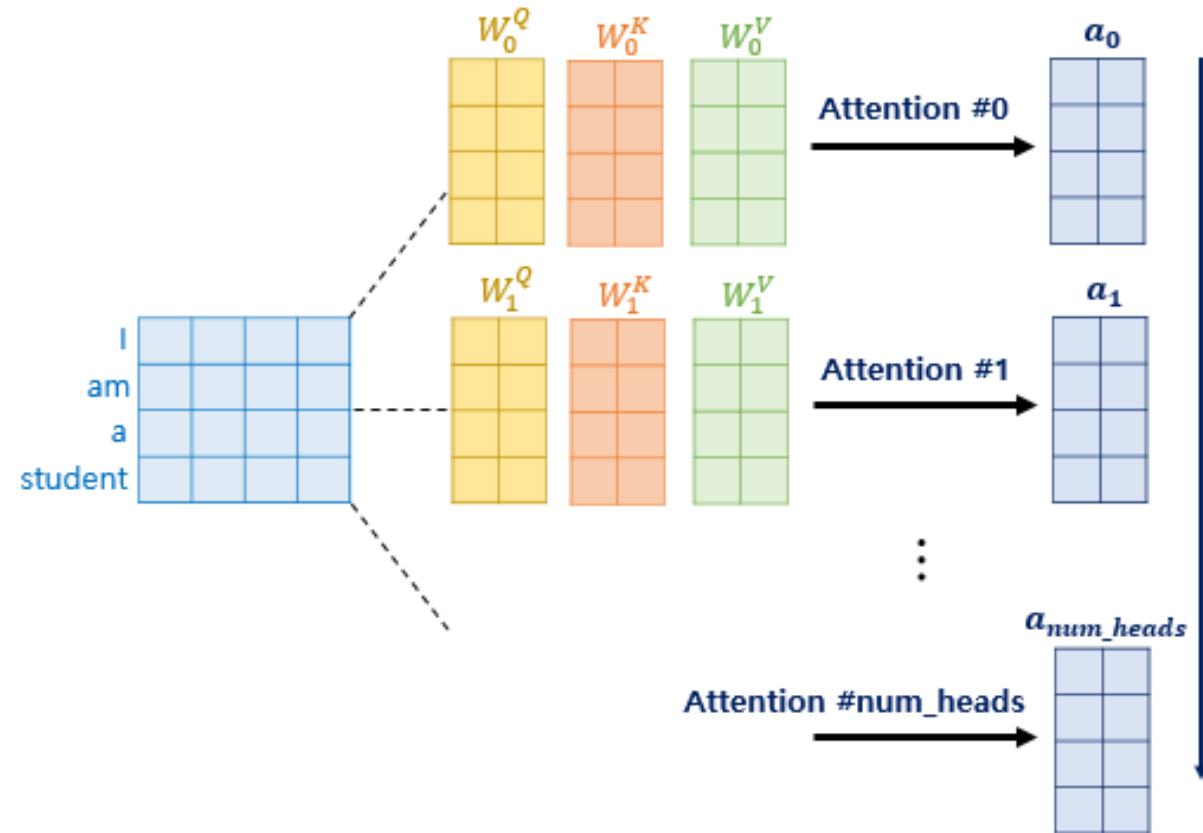
# Self-Attention



# Multi-head Self-Attention

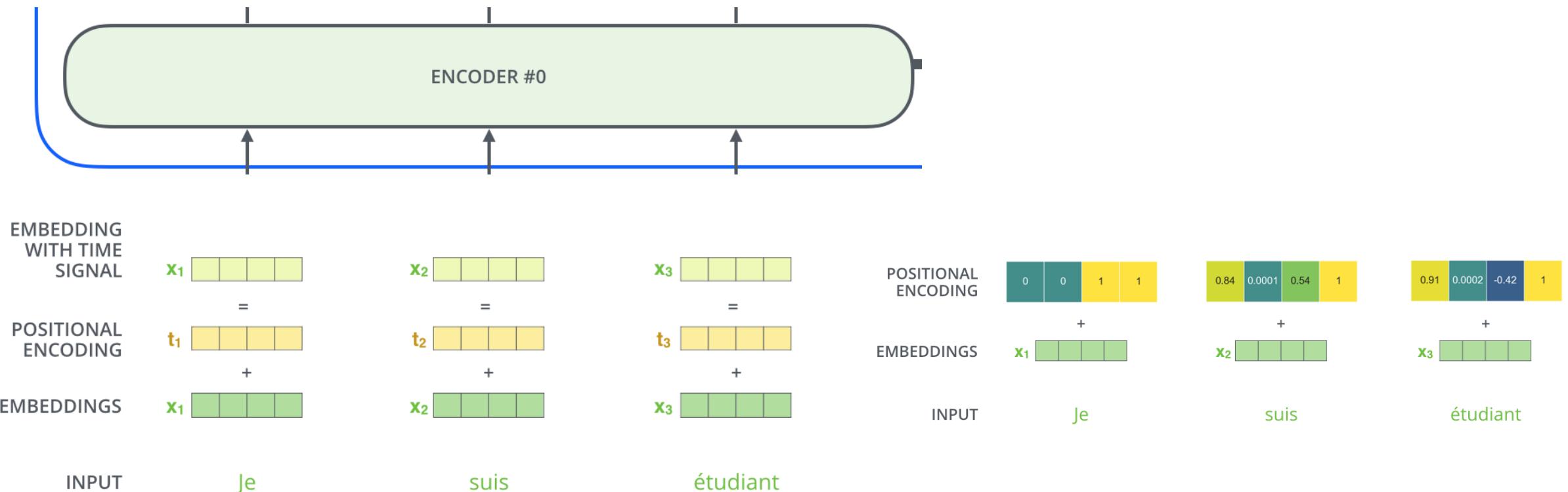
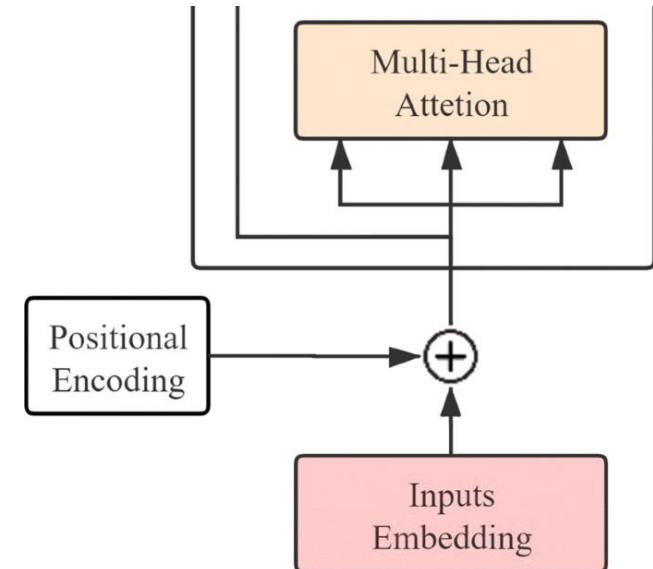
- Use multiple attention heads
  - Improve accuracy





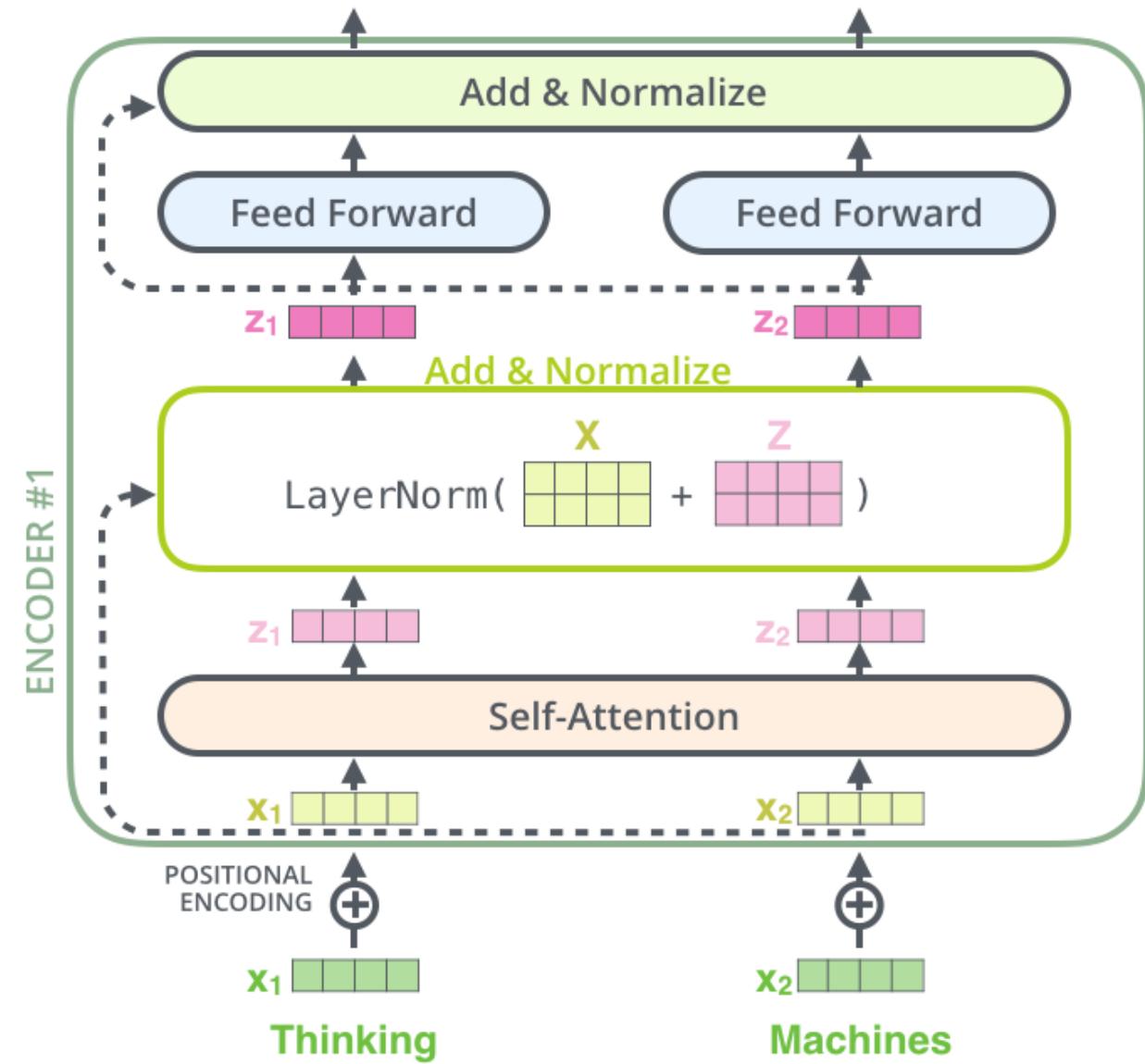
# Positional Encoding

- Positional encoding
  - Way to account for the order of the words in the input sequence
  - Determine the position of each word



# Feedforward Network

- Typical fully connected layer
  - 2 Dense layer
    - Parameters are the same over the different positions of the sentence and different over the encoder blocks
  - ReLU (activation function)



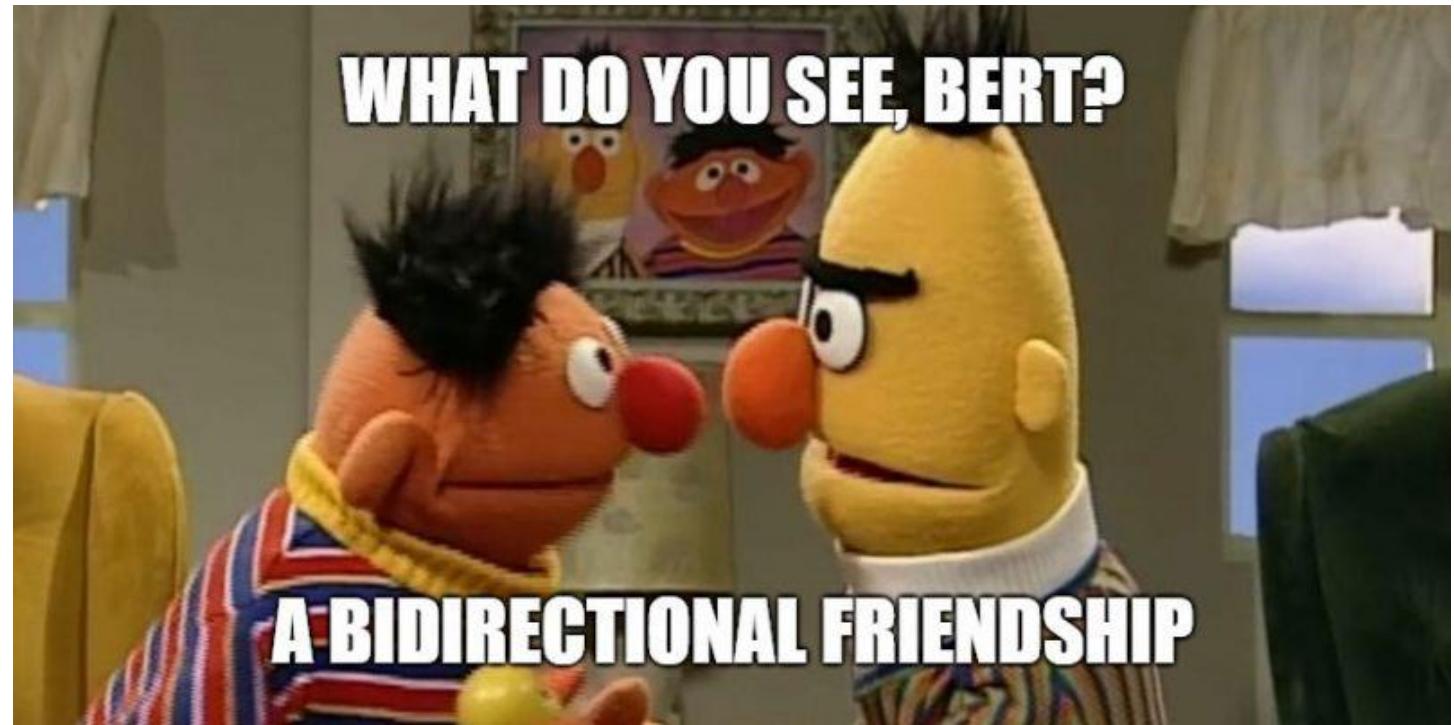
# BERT

- **Bidirectional Encoder Representation from Transformer**
  - *BERT is a method of pre-training language representations, meaning that we train a general-purpose “language understanding” model on a large text corpus (BooksCorpus and Wikipedia), and then use that model for downstream NLP tasks (fine-tuning) that we care about*



# BERT

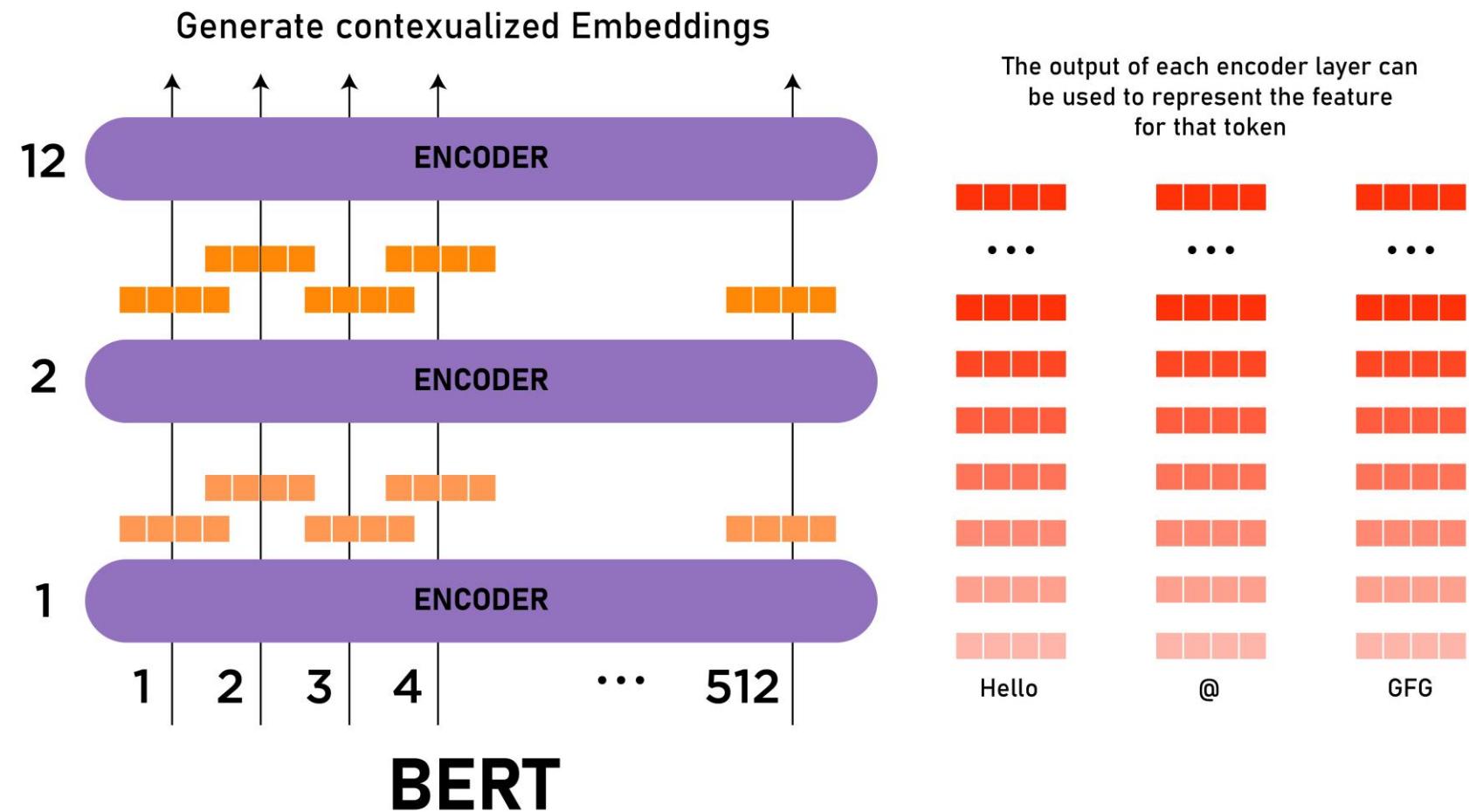
- Bidirectional Encoder Representation from Transformer
  - Transformer (like GPT)
  - Bi-directional (like ELMo)
  - Use only Encoder



# BERT

- **Architecture**

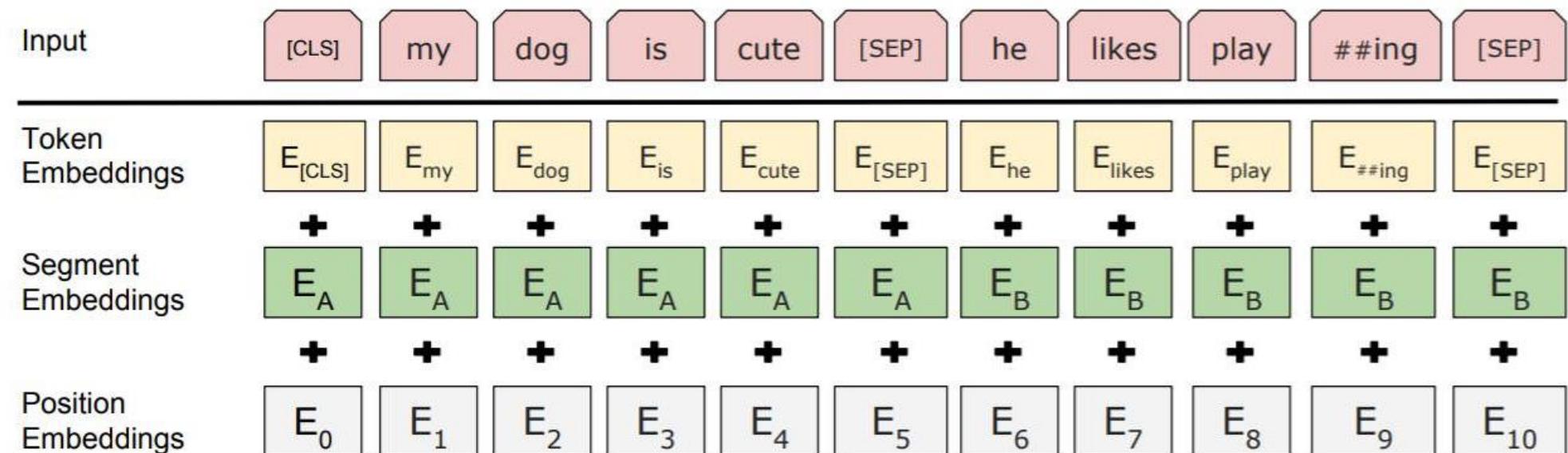
- 12 encoder layer
- 12 attention head  
(for each encoder)
- Hidden unit size 768



# BERT

- **Input Data Representation**

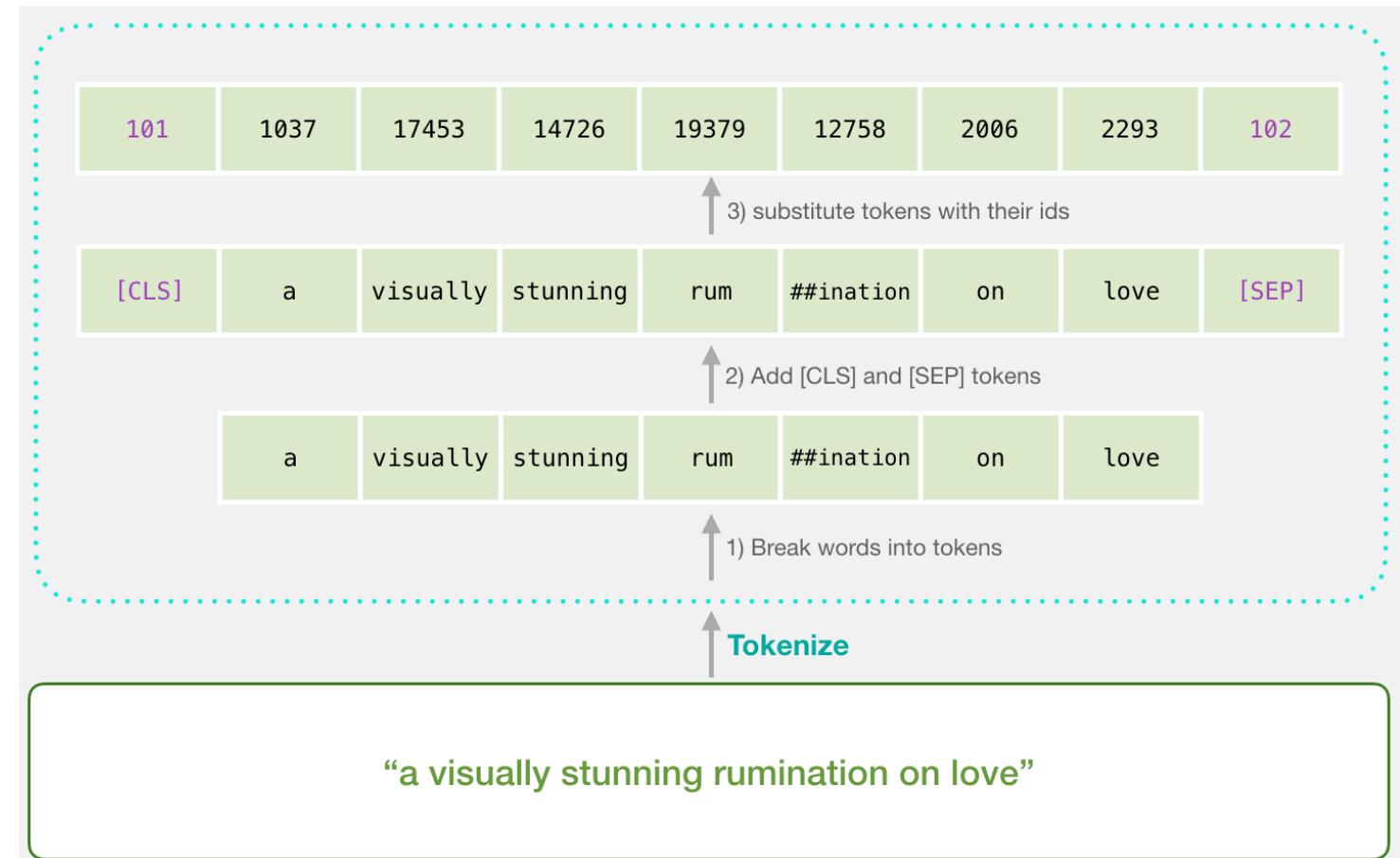
- 1) Token Embedding
- 2) Segment Embedding
- 3) Position Embedding



# BERT

- **Input Data Representation**

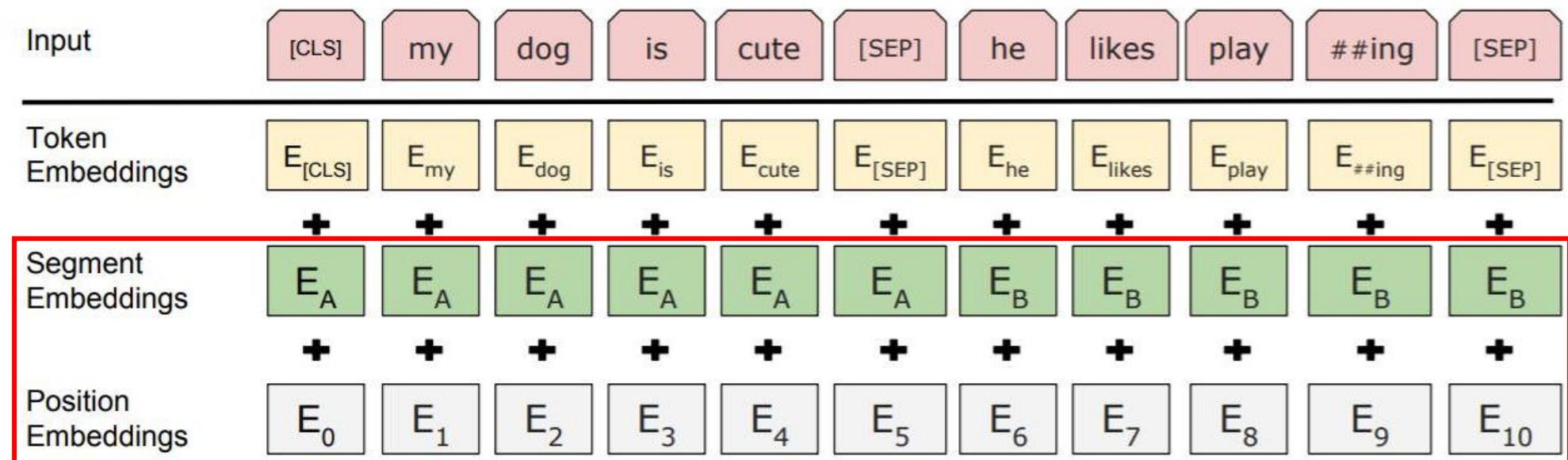
- 1) Token Embedding
  - By WordPiece embedding
  - **[CLS] token**
    - At the start of sentence
    - Represents the aggregate representation of sentence
    - Used for classification task
  - **[SEP] token**
    - At the end of sentence



# BERT

- **Input Data Representation**

- 2) Segment Embedding → Used to distinguish different sentence
- 3) Position Embedding → Used to take account the position of token in the sentence



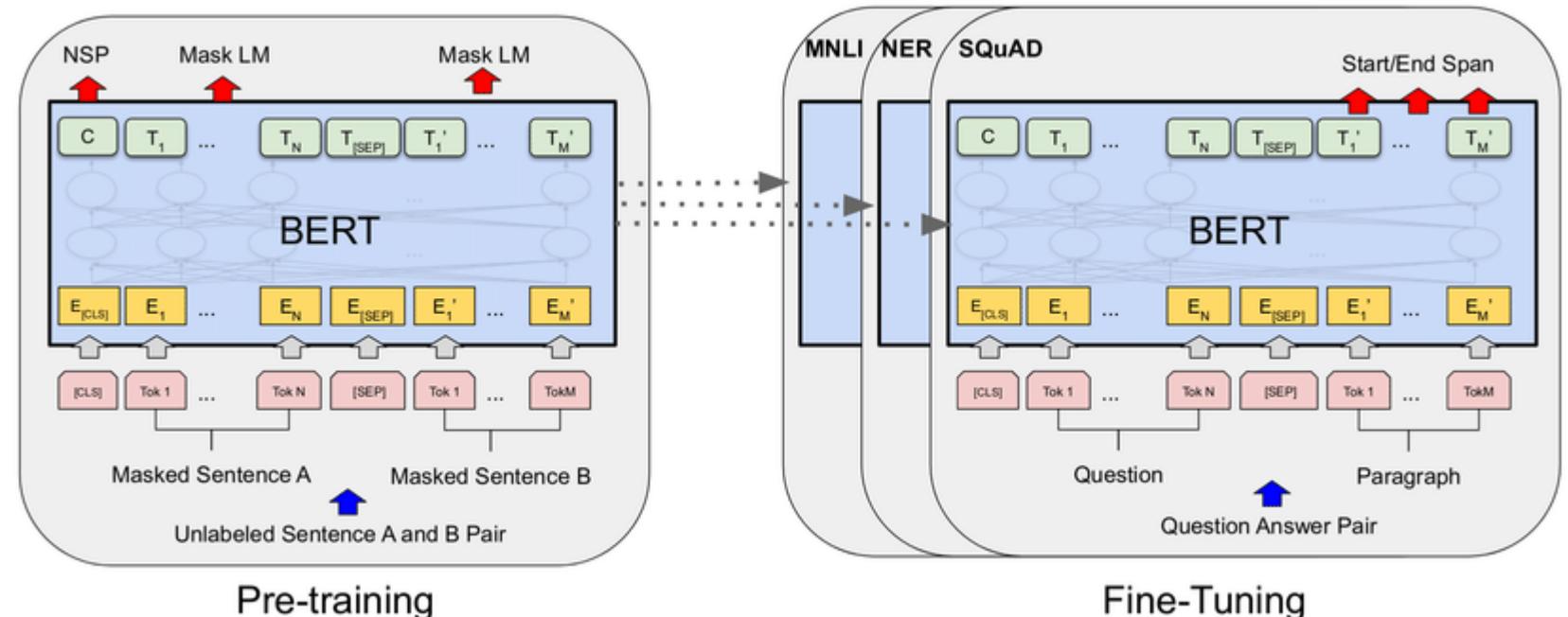
# BERT

- Pre-training & Fine-tuning



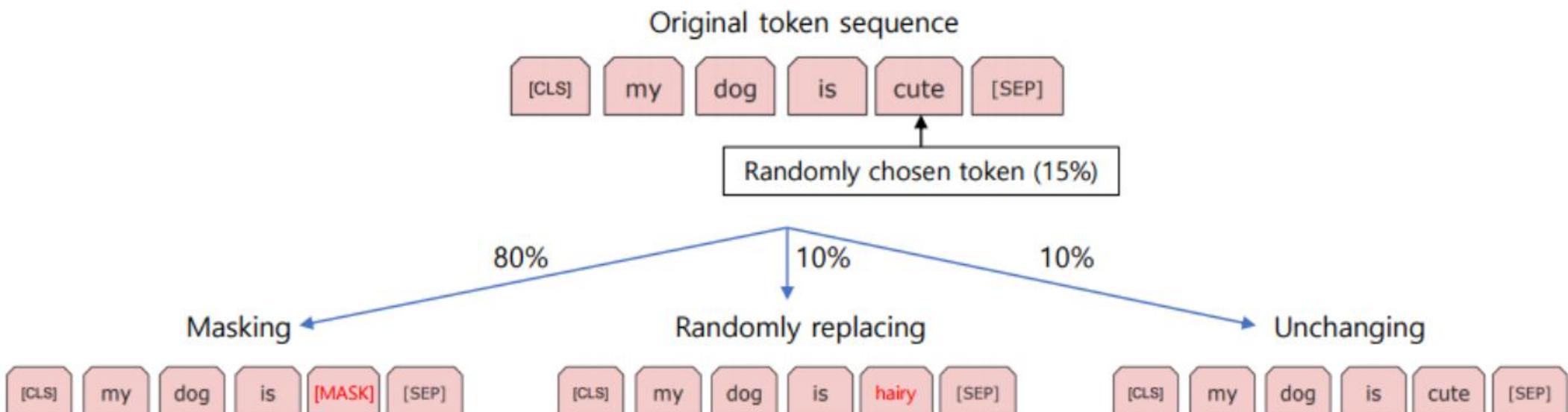
# BERT

- Pre-training
  - 1) MLM (Masked Language Modelling)
  - 2) NSP (Next Sentence Prediction)



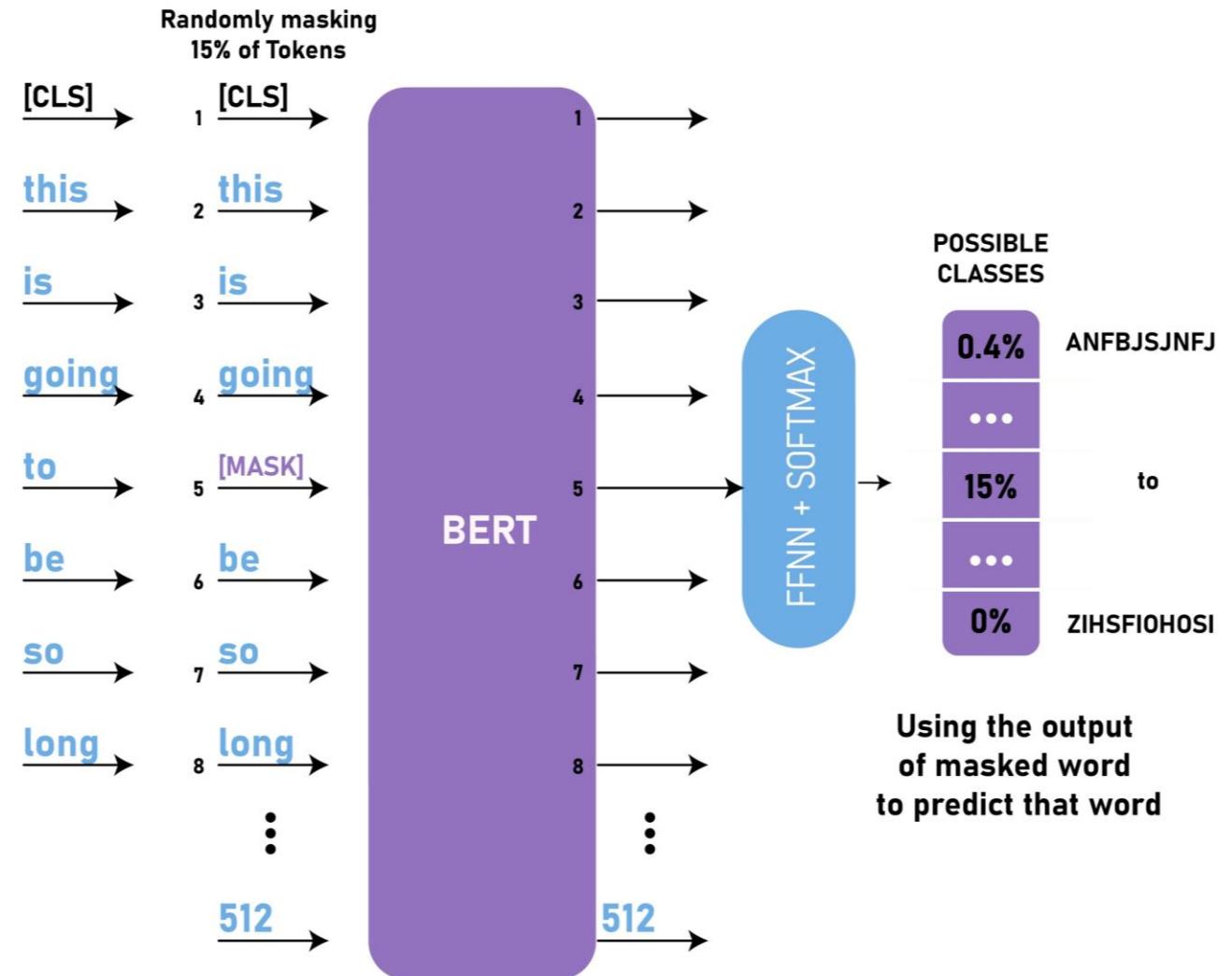
# BERT

- Pre-training
  - MLM (Masked Language Modelling)
    - Randomly mask 15% of words and train the network to **predict the masked words**
    - 80-10-10% rule (80% masking, 10% randomly replacing, 10% unchanging)



# BERT

- Pre-training
  - MLM (Masked Language Modelling)
    - Reads the sentence in both directions to predict the masked words



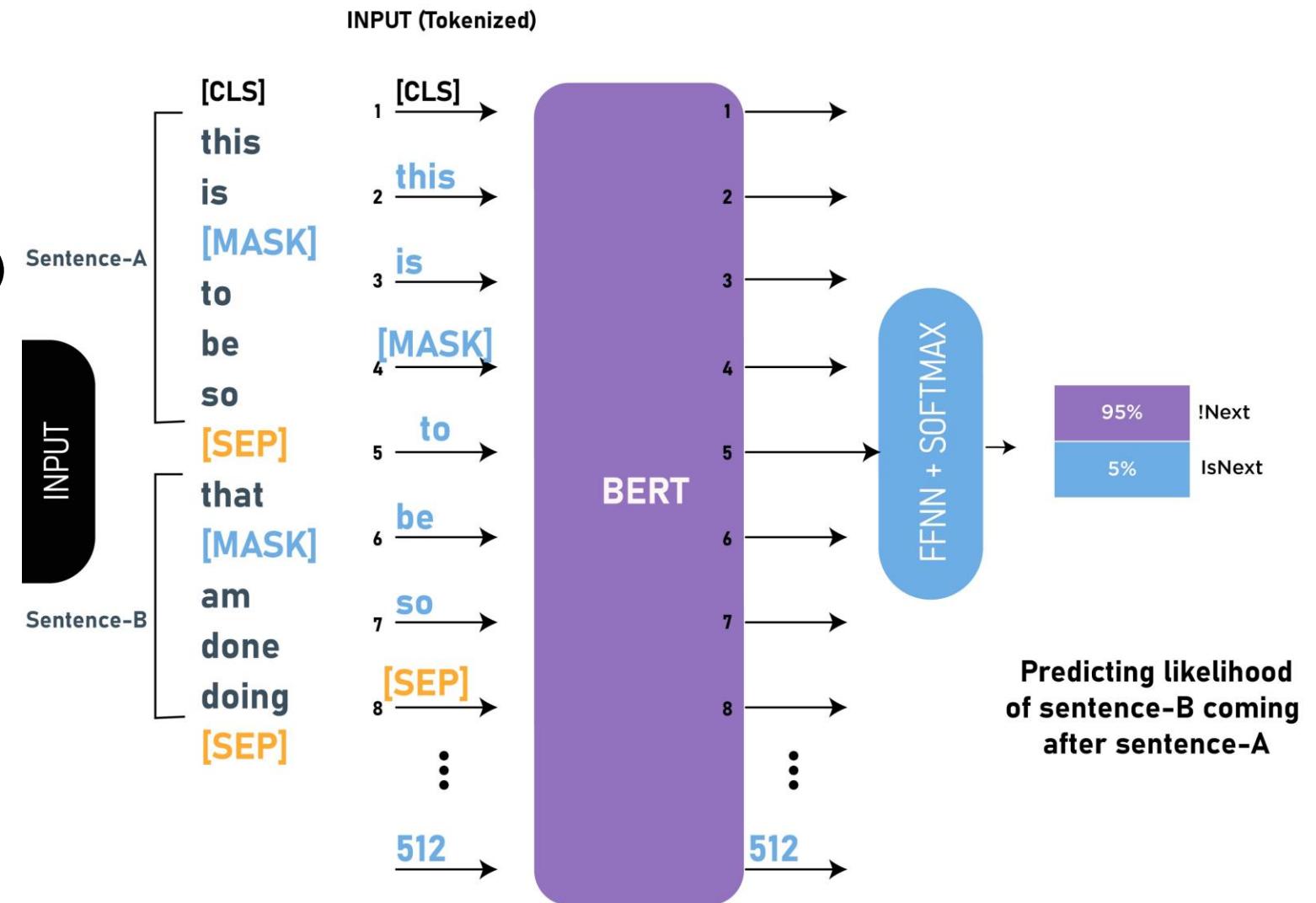
# BERT

- Pre-training
  - NSP (Next Sentence Prediction)
    - Binary classification



# BERT

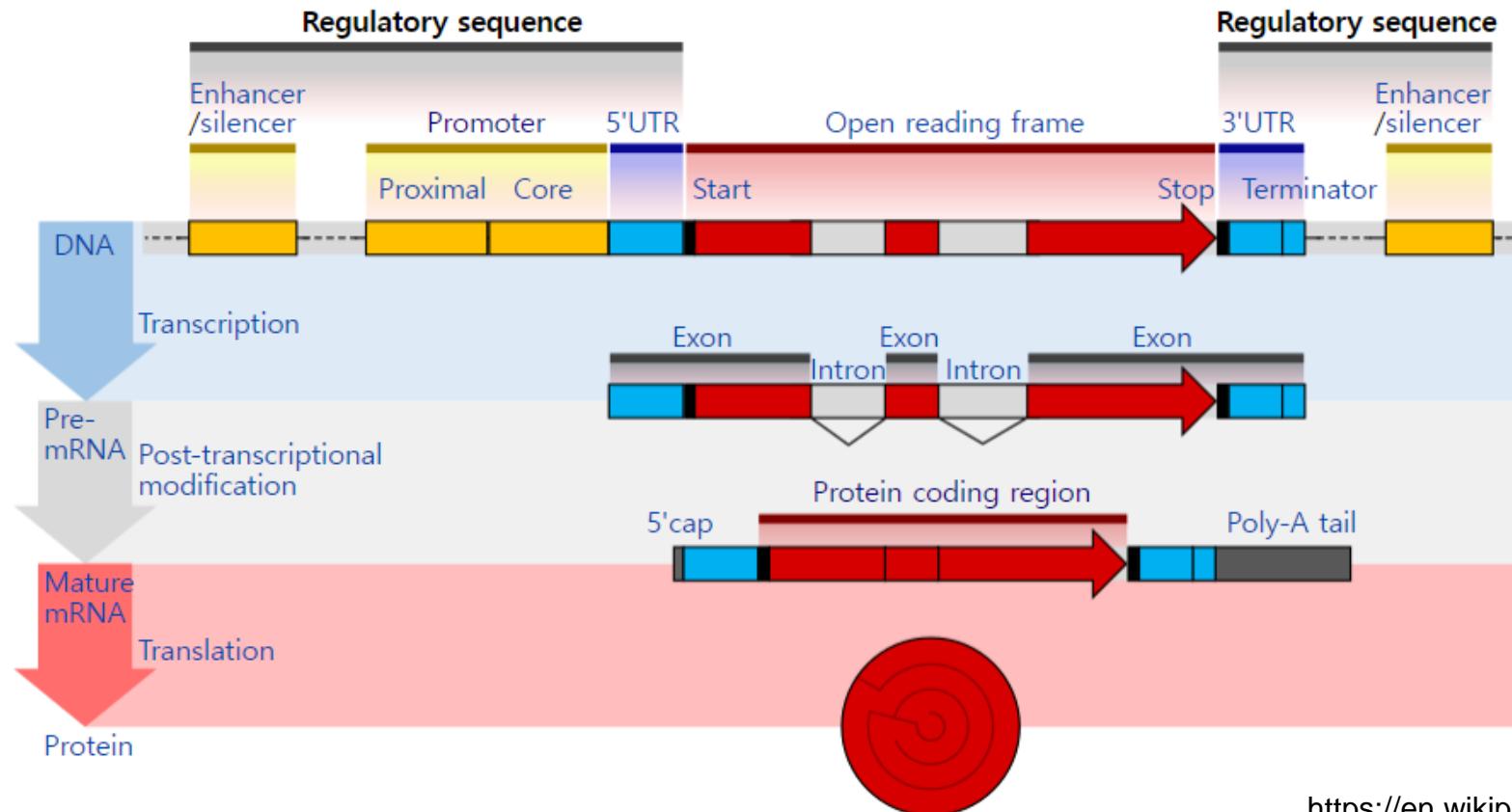
- Pre-training
  - NSP  
**(Next Sentence Prediction)**
  - Binary classification



**DNABERT**

# Introduction

- Deciphering the language of DNA for hidden instructions
  - Regulatory sequence
    - Capable of increasing or decreasing the expression of specific genes



# Introduction

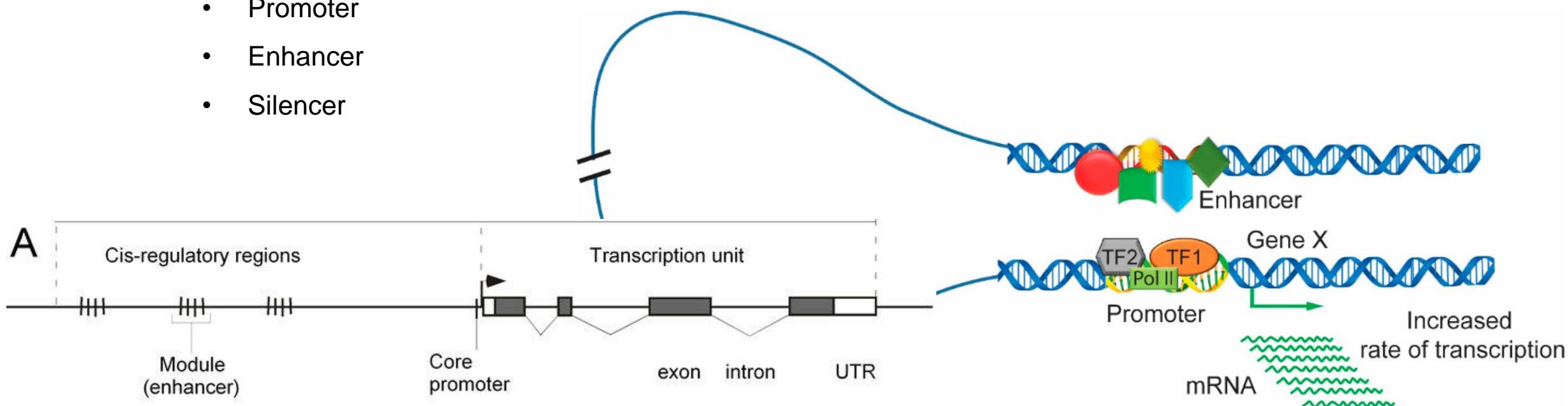
- Deciphering the language of DNA for hidden instructions

- Regulatory code

- CRE (*cis*-regulatory elements)

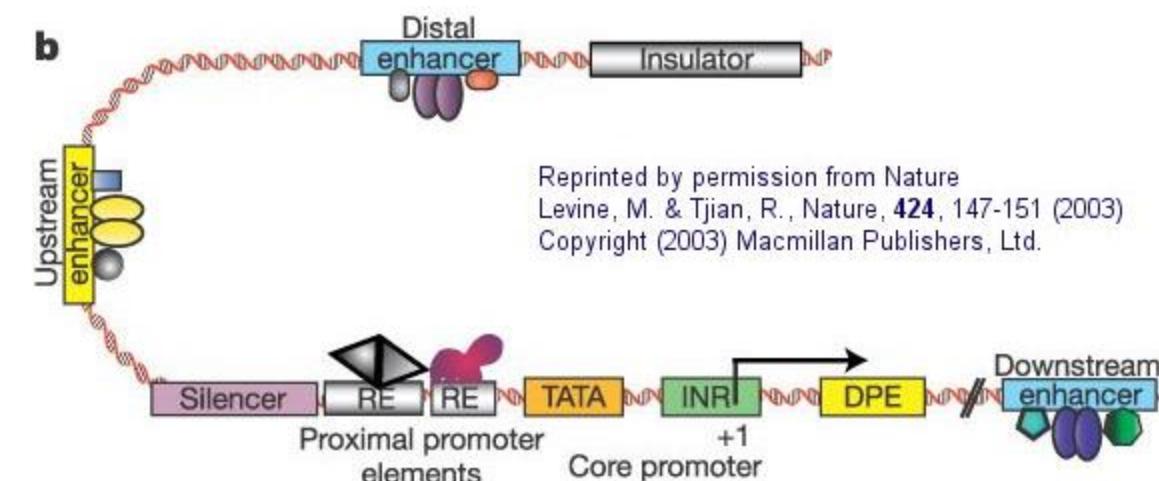
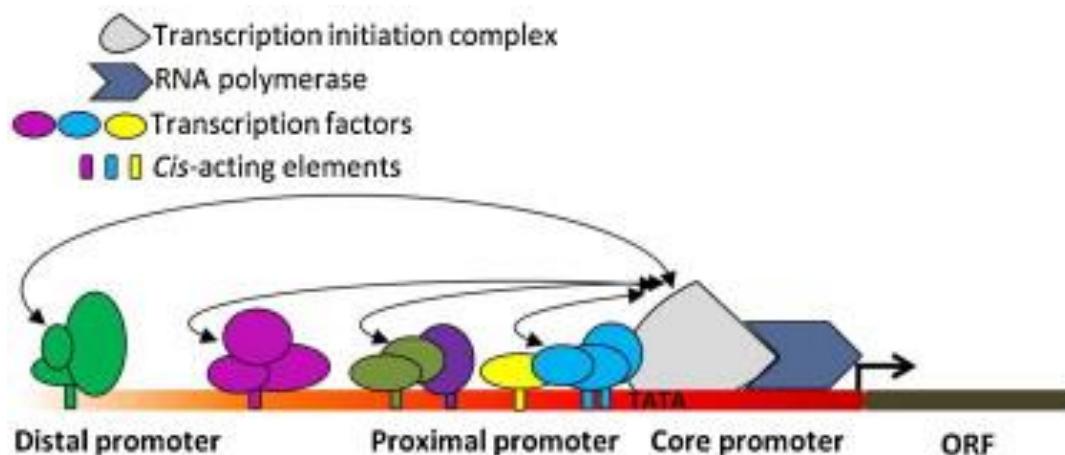
: regions of non-coding DNA which regulate the transcription of **neighboring genes**

- Promoter
- Enhancer
- Silencer



# Introduction

- Deciphering the language of DNA for hidden instructions
  - Regulatory code
    - CRE (*cis*-regulatory elements)
    - **Promoter** : Core / Proximal / Distal
      - Core promoter : Include TATA-box
      - Proximal & Distal promoter : Include enhancer / silencer / insulator / etc.

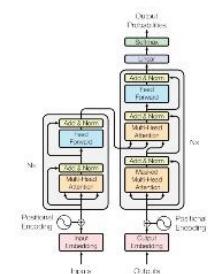


# Introduction

- **Deciphering the language of DNA for hidden instructions**
  - Regulatory code
    - Same CREs often have **distinct functions and activities in different biologic contexts**
      - e.g., different cell-types and organisms
      - Cooperation of widely spread multiple CREs → context-dependent use of alternative promoters
    - Suggest Polysemy and distant semantic relationship within sequence codes
    - How the semantics (i.e., functions) of CREs vary across different contexts

# Introduction

- Deciphering the language of DNA for hidden instructions
  - Ideal computational method
    - 1) Globally take all the **contextual information** into account to distinguish polysemous CREs
    - 2) Develop **generic understanding transferable to various tasks**
    - 3) **Generalize well** when labeled data is limited
  - ⇒ Transformer
    - Develop ***global contextual embedding*** via self-attention



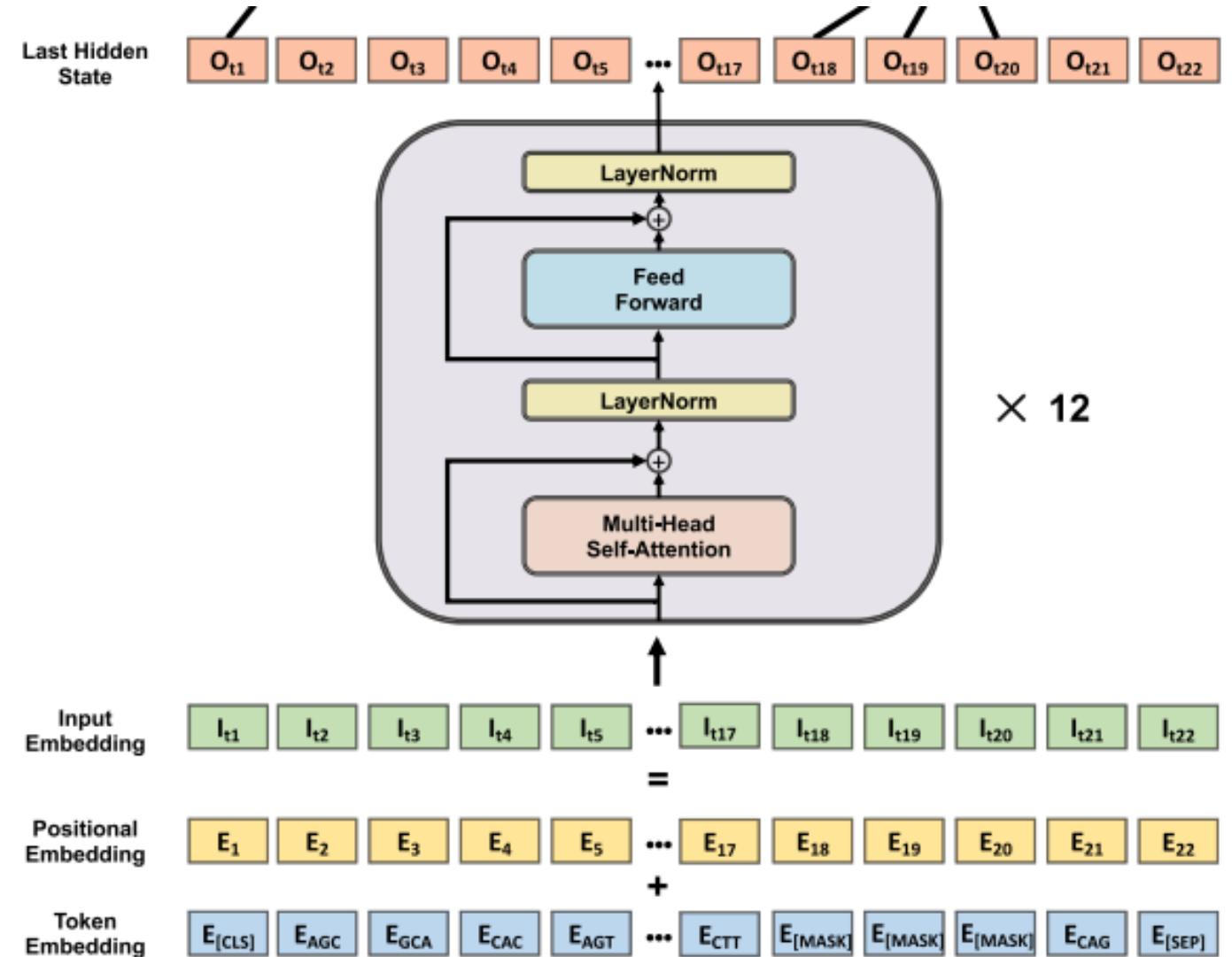
# Introduction

- **DNABERT**
  - Applies transformer (attention-based architecture)
  - It resolves
    - 1) Develop **general and transferable understanding of DNA** from the purely **unlabeled** human genome, and utilize them to generically solve various **sequence related tasks** in a “**one-model-does-it-all**” fashion
    - 2) Globally capture **contextual information** from the entire input sequence with attention mechanism
    - 3) Achieve great performance in **data-scarce scenarios**
    - 4) **Uncover important subregions and potential relationships**  
between different *cis*-elements of DNA sequence without any human guidance
    - 5) Successfully work in a **cross-organism manner**

# Materials & Method

- **DNABERT – Architecture**

- Same architecture with BERT-base
  - 12 transformer layers
  - 12 attention heads in each layer
  - 768 hidden units



# Materials & Method

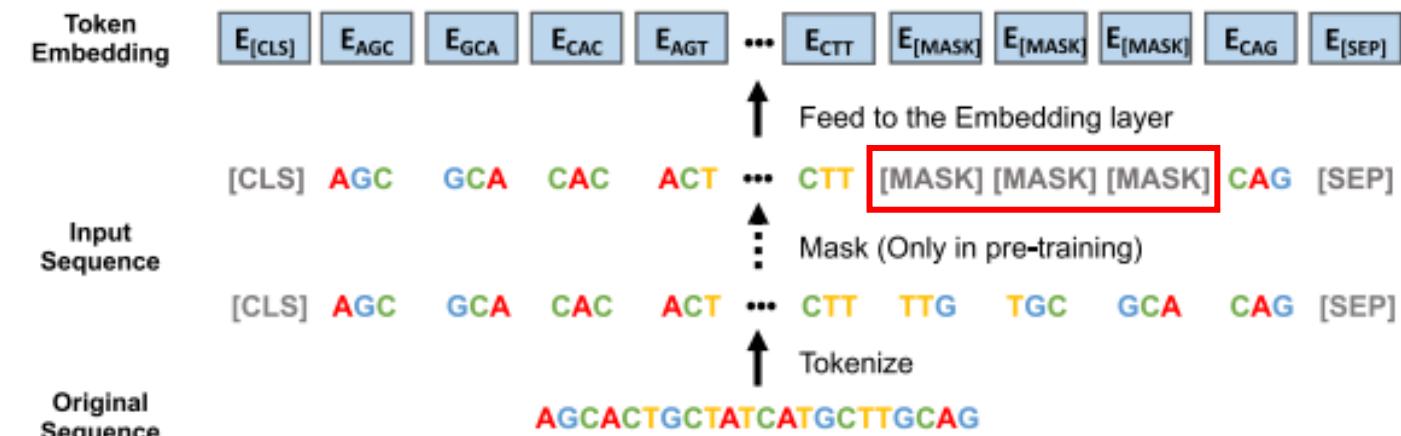
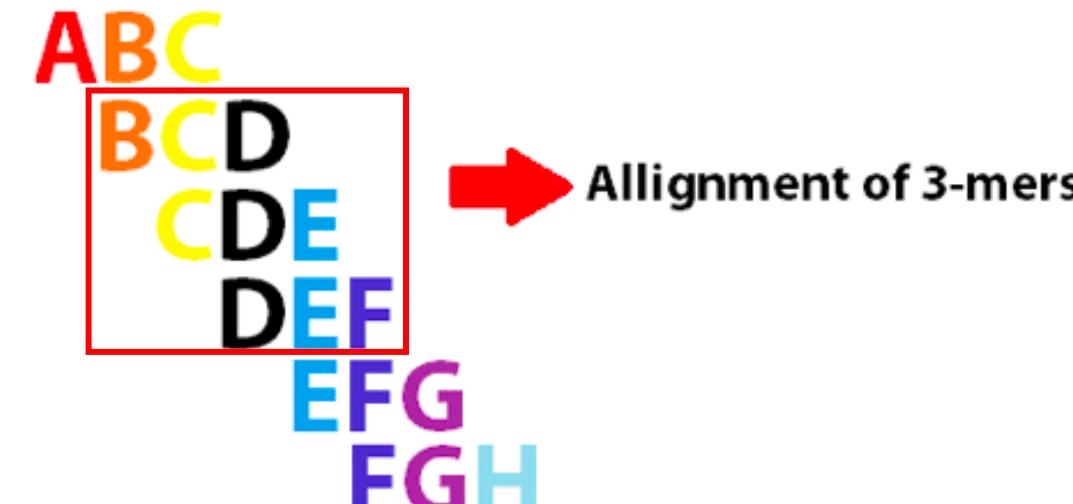
- **DNABERT – Architecture**

- Input

- **k-mer tokens ( $k = 3, 4, 5, 6$ )**

- **Masked contiguous k-length spans of k-mers**

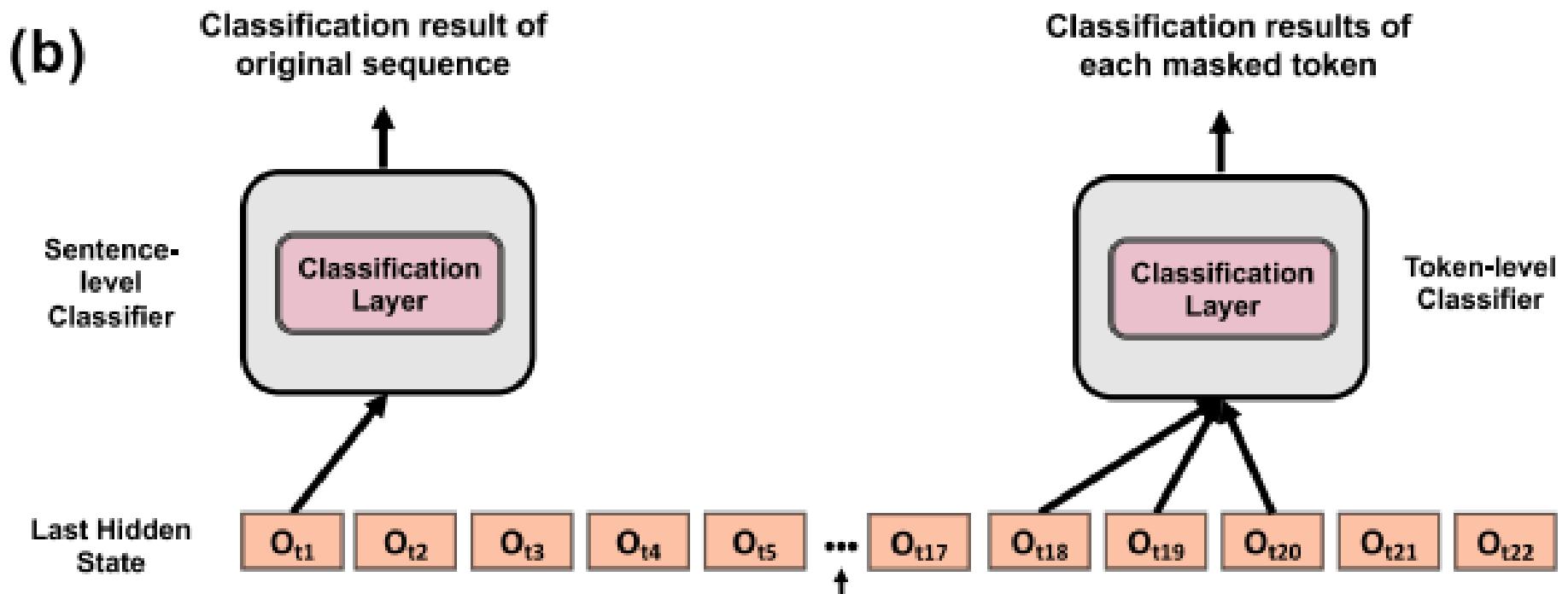
- Masked k-mer can be trivially made up based on its previous and next k-mers
- Prevent the model from learning deep semantic relations among a DNA sequence
- ex) 3-mer sequence



# Materials & Method

- **DNABERT – Architecture**

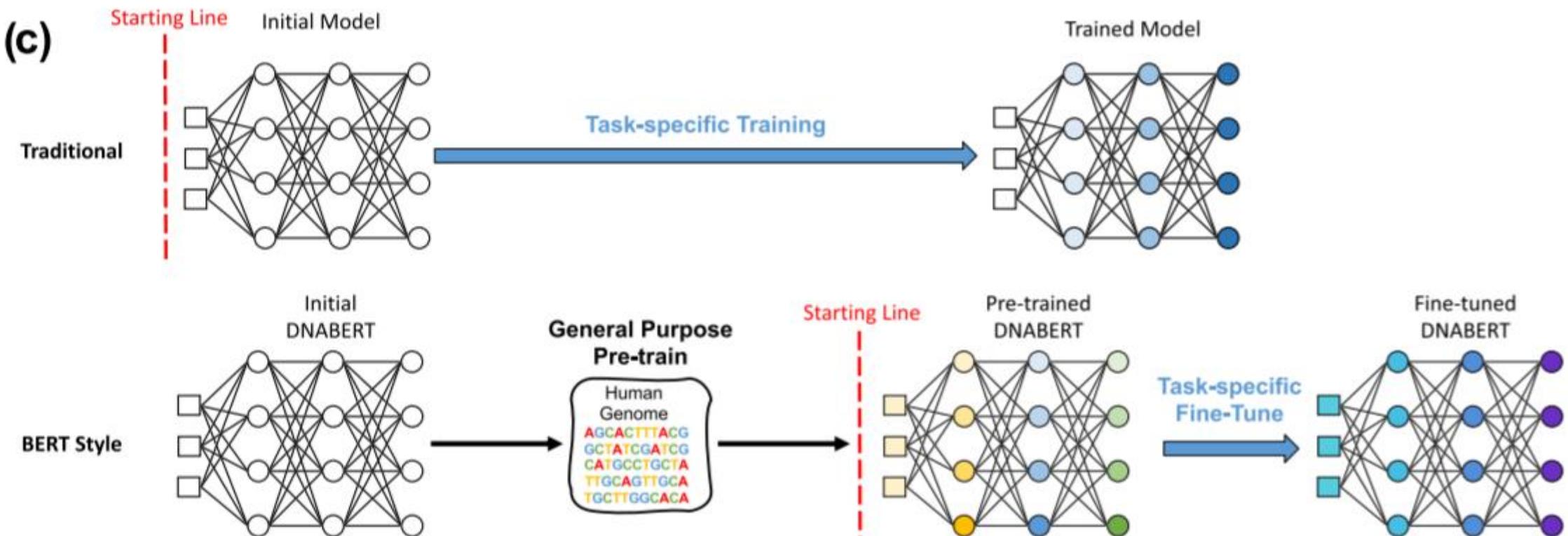
- Output
  - [CLS] token representation → Sentence-level classification
  - Individual masked token representations → token-level classification



# Materials & Method

- **DNABERT – Pre-training**

- First develops **general-purpose understandings** from massive amount of **unlabeled data**, then resolves various applications with **task-specific data** with minimal architectural modification



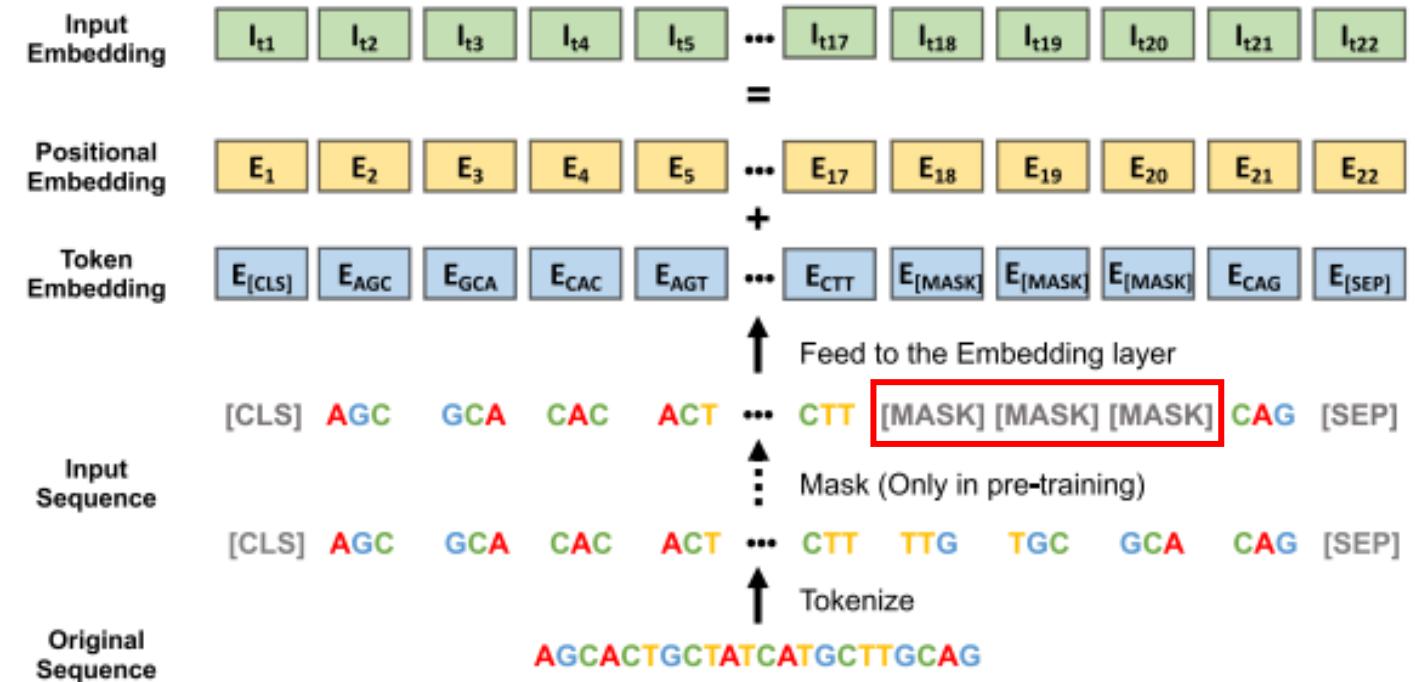
# Materials & Method

- **DNABERT – Pre-training**
  - Difference from BERT
    - Skipped pre-training with NSP task
    - Adjusted the sequence length and forced the model to predict contiguous k tokens
  - During pre-training
    - DNABERT learns basic syntax and semantics of DNA via self-supervision
    - Based on 10~510 length sequences extracted from human genome via truncation and sampling

# Materials & Method

- **DNABERT – Pre-training**

- Take a sequence
  - Maximum length 512
    - Direct non-overlap sampling
    - Random sampling
- Masking
  - Mask a contiguous k-length spans of certain k-mers
  - Total 15% of input sequence



# Materials & Method

- **DNABERT – Pre-training**

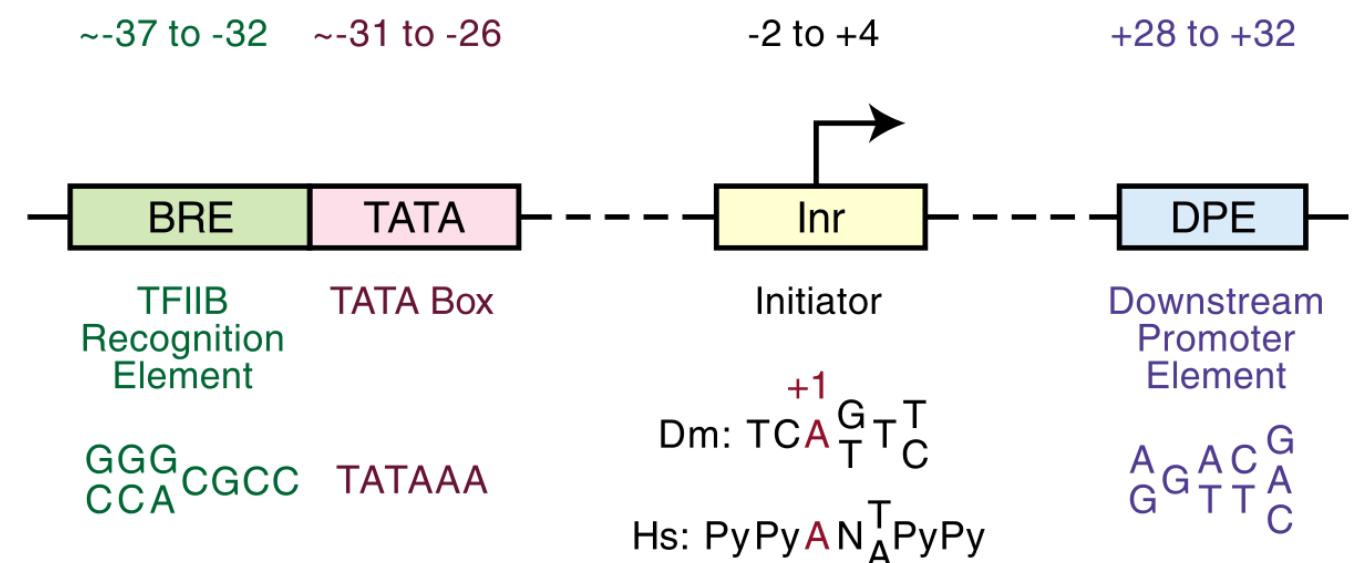
- Hyper-parameters
  - Epoch : 120k steps
    - Learning rate linearly increased (warm-up) from  $0\sim 4e^{-4}$  in the first 10k steps and then linearly decreased to 0 after 200km steps
    - Stopped training procedure after 120k steps (loss curve showed plateau)
  - Batch size 2000
  - Objective function : Cross-entropy loss calculated over all the masked k-mers
  - Optimizer : AdamW ( $\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1e^{-6}$ , weight decay = 0.01)

# Materials & Method

- **DNABERT – Fine-tuning**
  - Fine-tuning by task-specific training data
    - 1) Promoter prediction
    - 2) Transcription factor binding sites prediction
    - 3) Splice site prediction

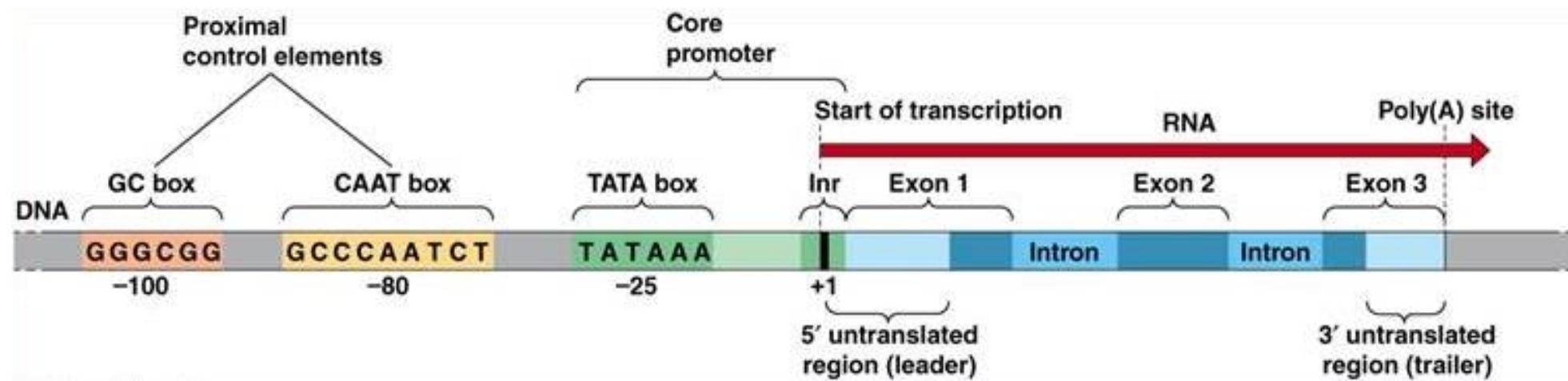
# Results

- **1. DNABERT-Prom effectively predicts proximal and core promoter regions**
  - Data
    - EPDnew (Eukaryotic Promoter Database) : Annotated non-redundant / Eukaryotic Pol II promoters
    - 3,065 human TATA & 26,533 non-TATA promoter-containing sequences
      - -5,000 ~ +5,000 bp (+1 : position of TSS)
  - Model



# Results

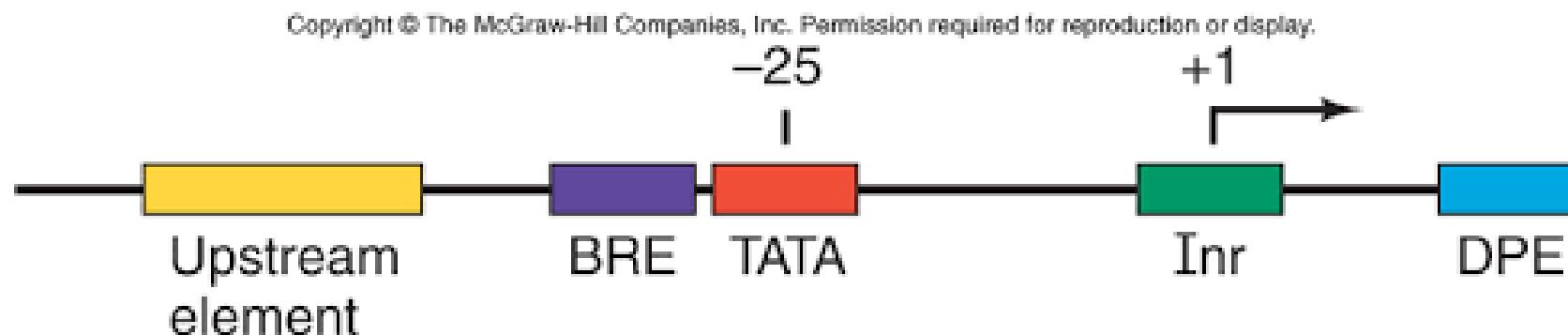
- 1. DNABERT-Prom effectively predicts proximal and core promoter regions
  - DNABERT-Prom-300
    - Binary classification
    - Positive class (promoters)
      - 300bp-long promoter sequences extracted from -249~+50bp around TSS position



© 2012 Pearson Education, Inc.

# Results

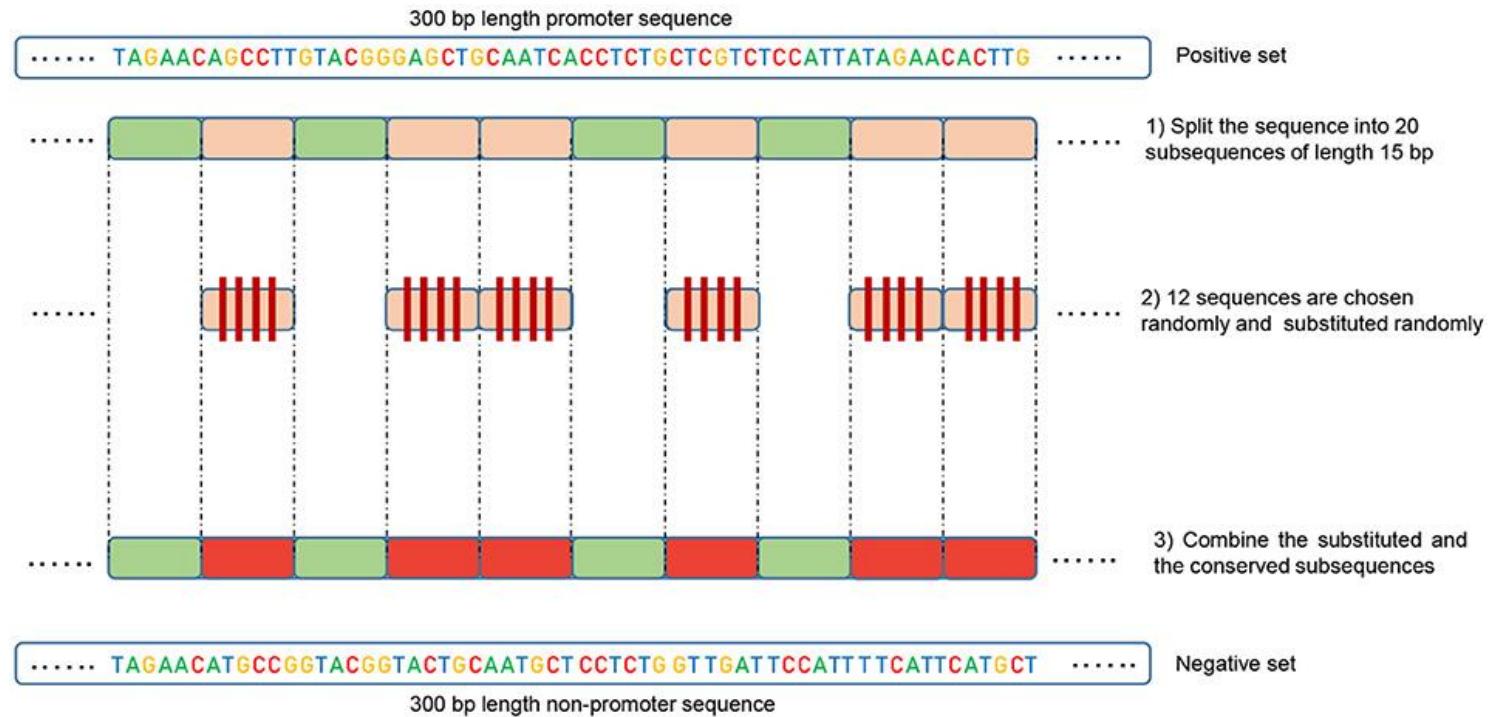
- **1. DNABERT-Prom effectively predicts proximal and core promoter regions**
    - DNABERT-Prom-300
      - Binary classification
      - Negative class (non-promoters)
        - **TATA-negative** : randomly picked 3,065 of 300bp genomic regions not within the -249~+50bp range  
*but contains the TATA motif*
- ensure that the TATA-negative is similar to TATA promoters as possible



# Results

- 1. DNABERT-Prom effectively predicts proximal and core promoter regions

- DNABERT-Prom-300
  - Binary classification
  - Negative class (non-promoters)
    - Non-TATA negative**
      - : random substitution approach  
(from DeePromoter paper)

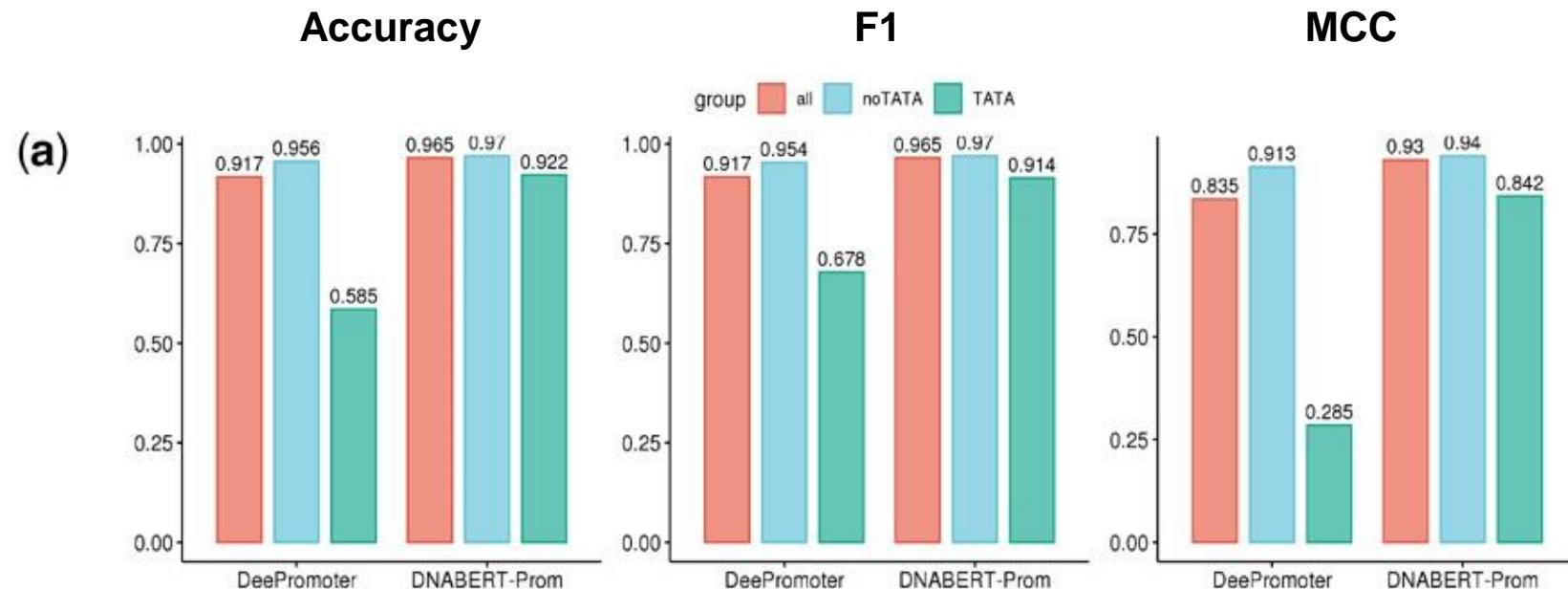


# Results

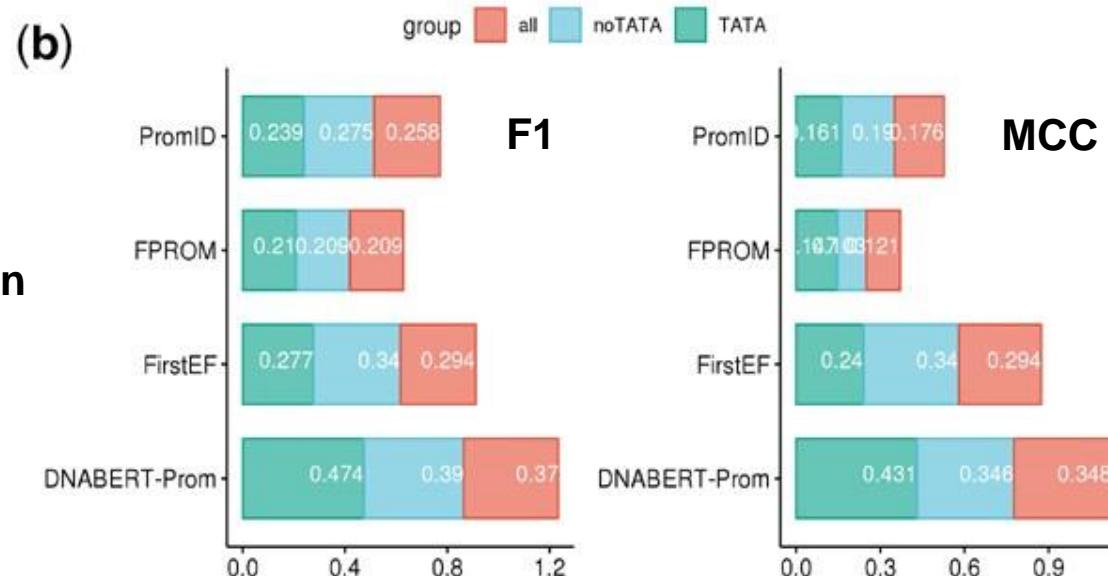
- **1. DNABERT-Prom effectively predicts proximal and core promoter regions**
  - DNABERT-Prom-scan
    - Mimics real-world situations
      - Scan very long genomic regions w/ sliding window and obtain 1001bp-long sequences for promoter region
    - Scanned all 10,000bp-long sequences from EPDnew with a step size of 100
    - Evaluation
      - If the predicted sequence has  $\geq 50\%$  (i.e., 500bp) overlap with -500~+500 bp regions of TSS  
    ⇒ True Positive
      - Otherwise ⇒ False Positive
      - Failure to make prediction in the area of -500~+500 bp ⇒ False Negative

# Results

## DNABERT-Prom-300



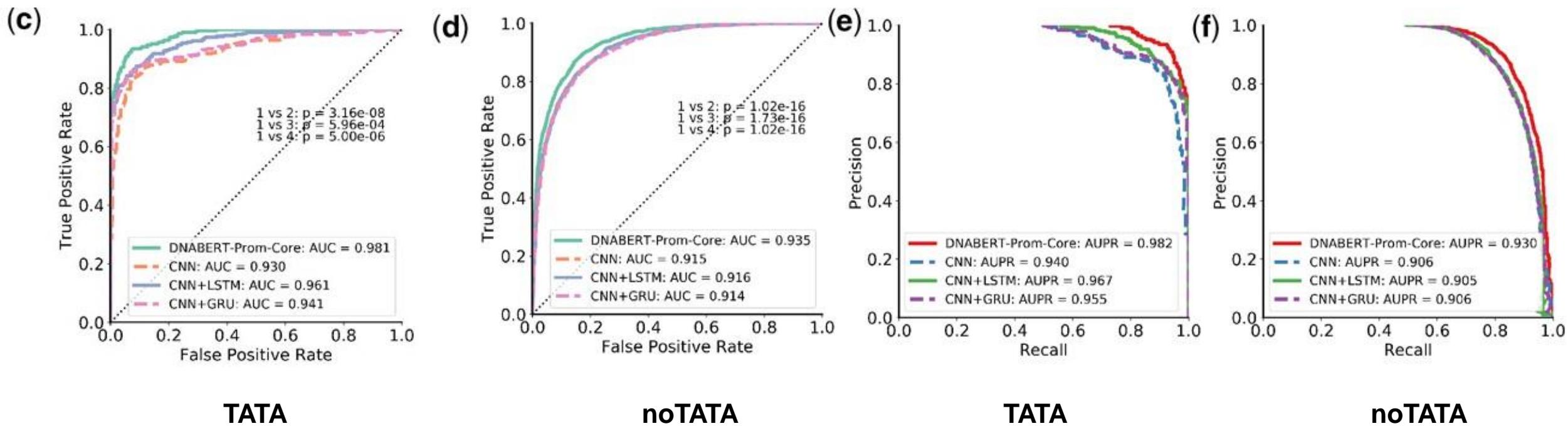
## DNABERT-Prom-scan



→ DNABERT significantly outperforms other models in identifying promoter regions

# Results

- 1. DNABERT-Prom effectively predicts proximal and core promoter regions

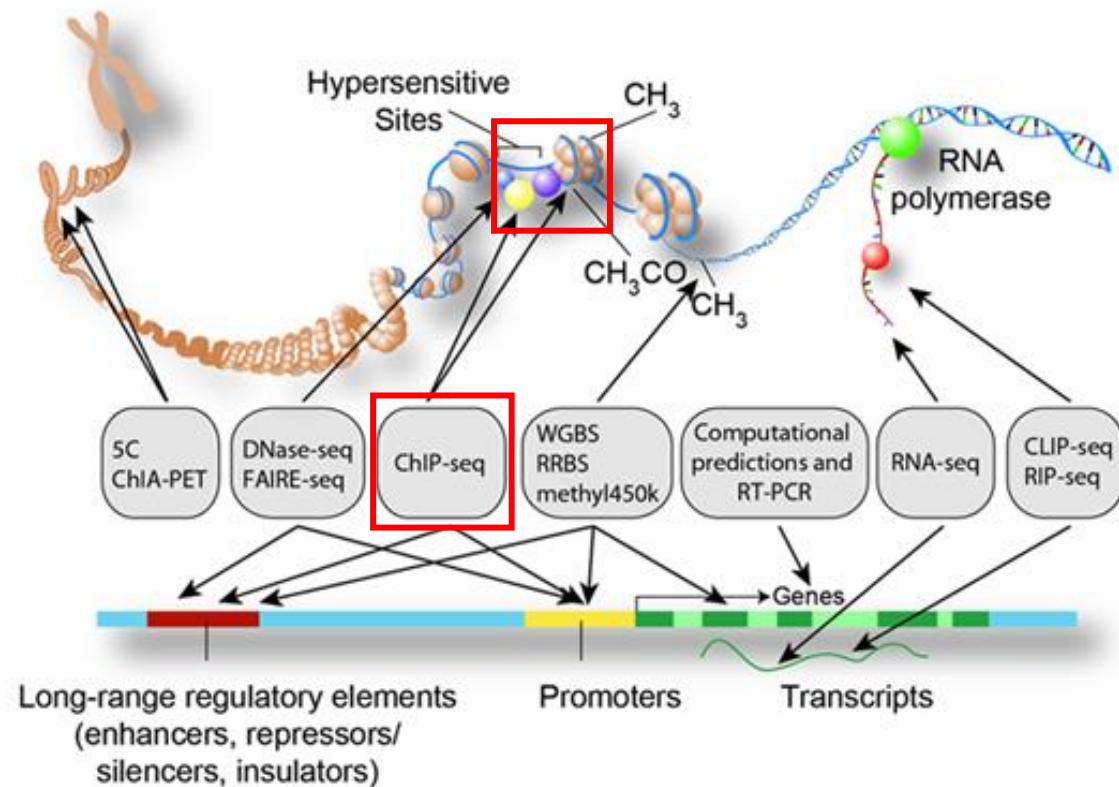


→ DNABERT significantly outperforms other models in identifying promoter regions

# Results

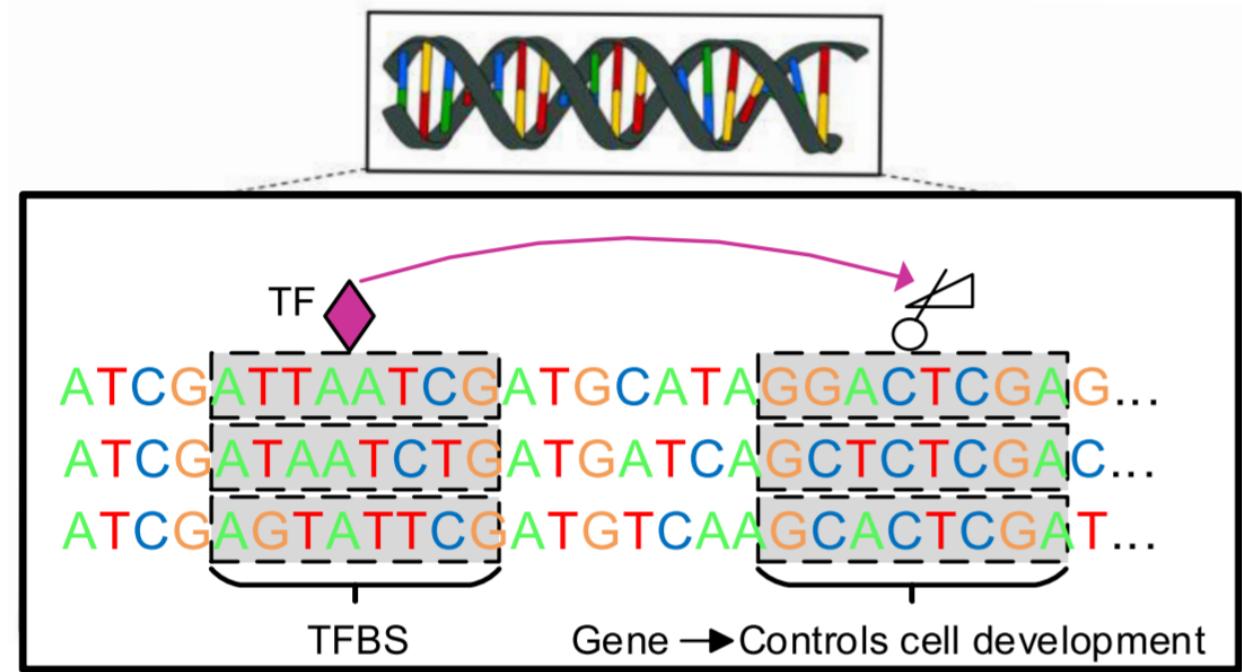
- **2. DNABERT-TF accurately identifies transcription factor binding sites**

- Fine-tuned DNABERT to predict TFBS in ChIP-seq enriched regions
  - Using 690 TF ChIP-seq uniform peak profiles from ENCODE database
  - Positive set
    - extracted the centering 101-bp region around each ChIP-seq peak
  - Negative set
    - Similar approach with DESSO
    - Pick actual 101-bp long sequences not overlapping with any peaks with same GC content



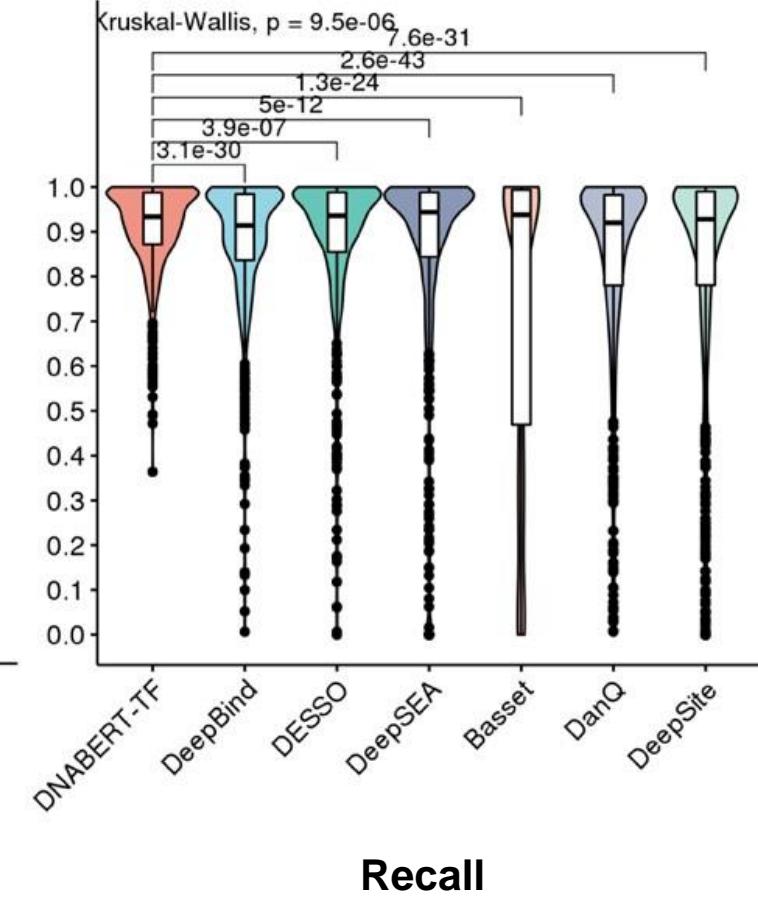
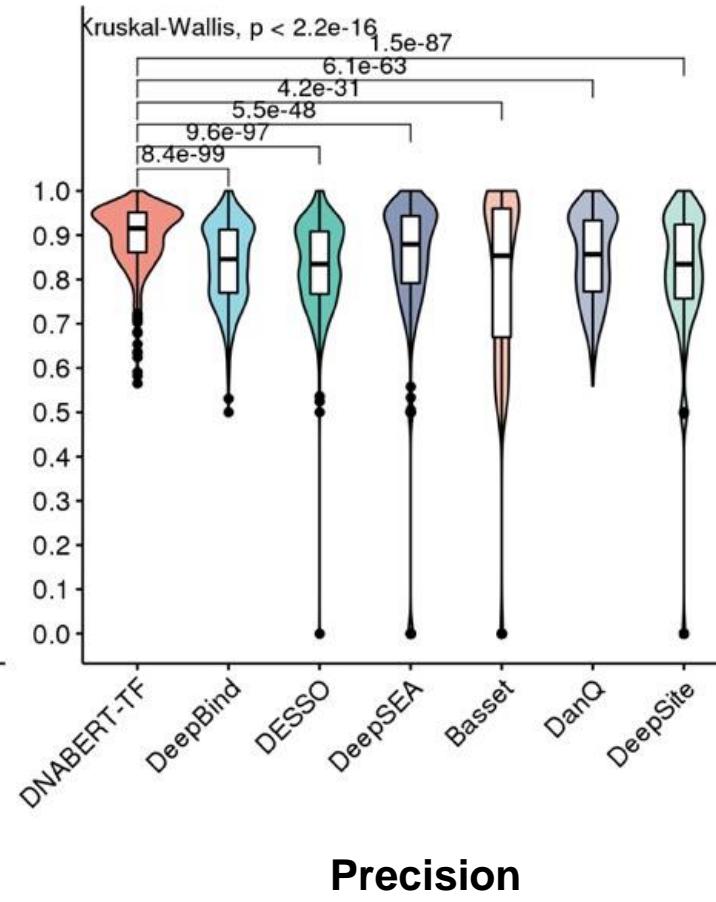
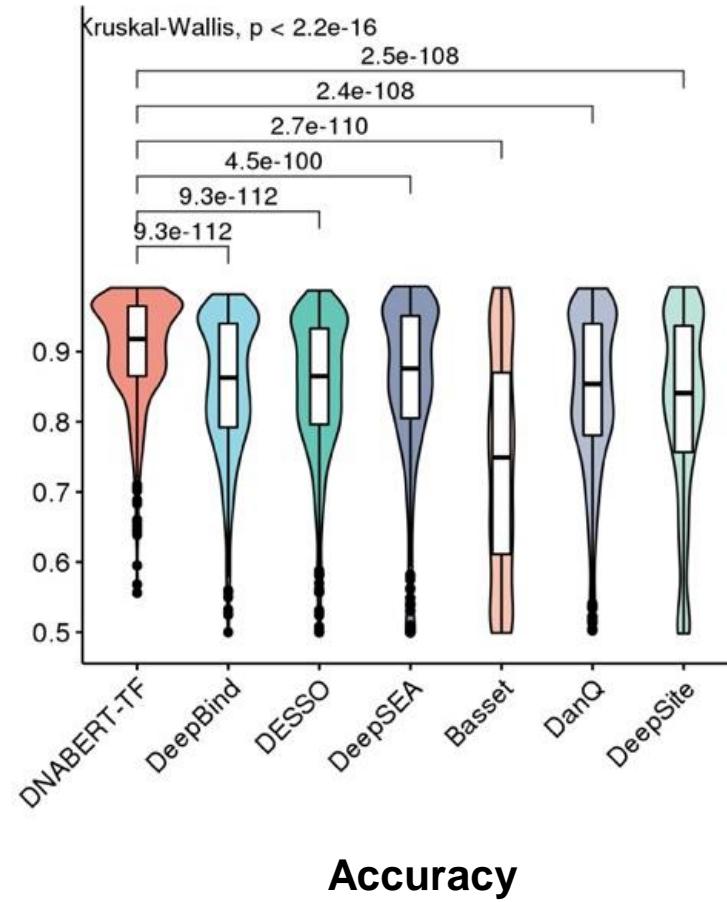
# Results

- **2. DNABERT-TF accurately identifies transcription factor binding sites**
  - Fine-tuned DNABERT to predict TFBS in ChIP-seq enriched regions
    - Binary classification using the top 500 even-numbered peaks as test set
    - Compared with well-known and previous published TFBS prediction tools
      - DeepBind, DeepSEA, DeepSite, DESSO

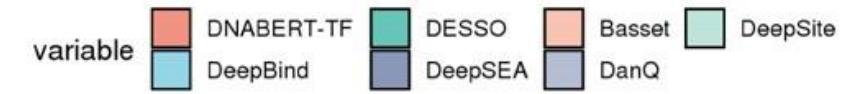


# Results

- 2. DNABERT-TF accurately identifies TFBS



→ DNABERT accurately identifies TFBSs

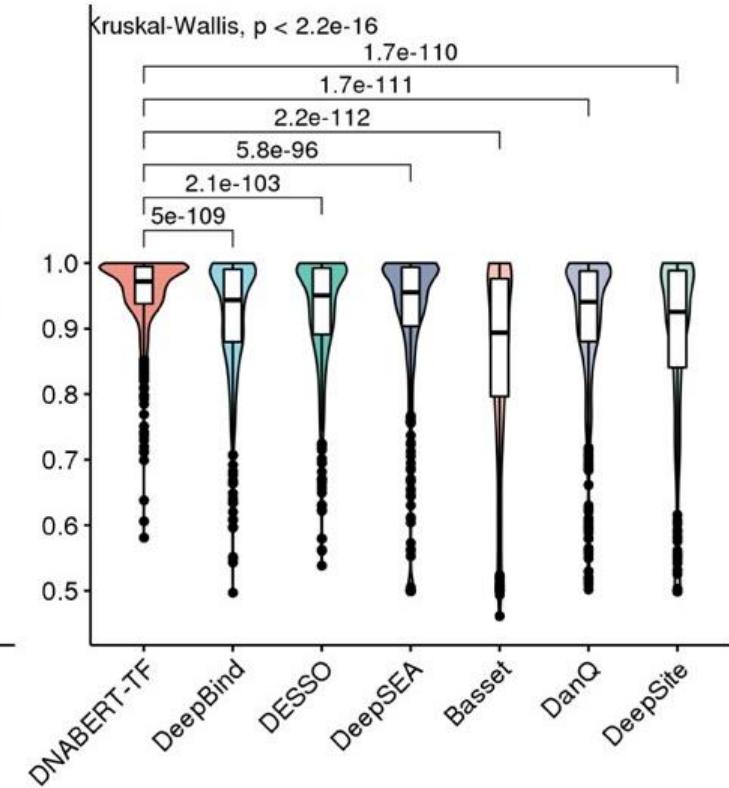
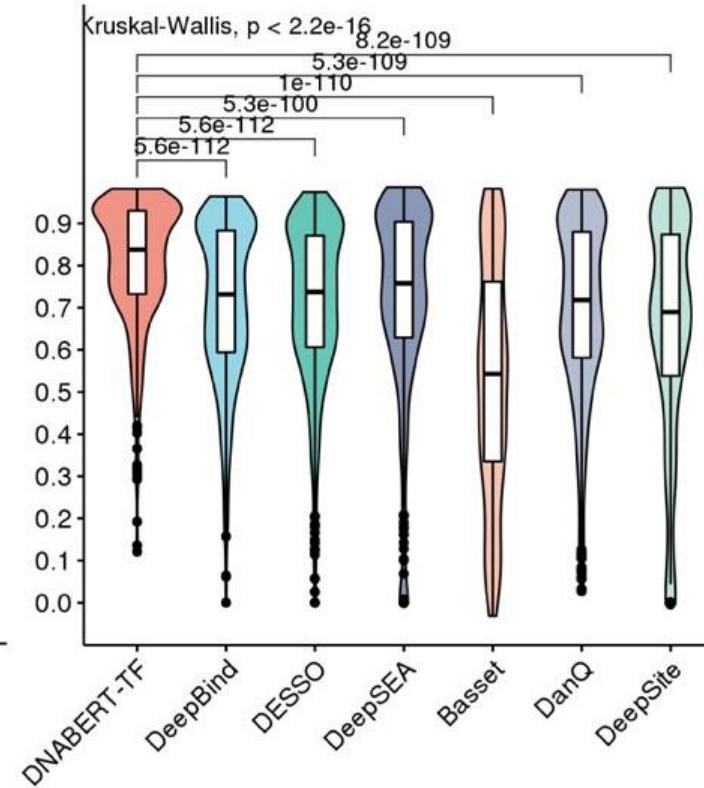
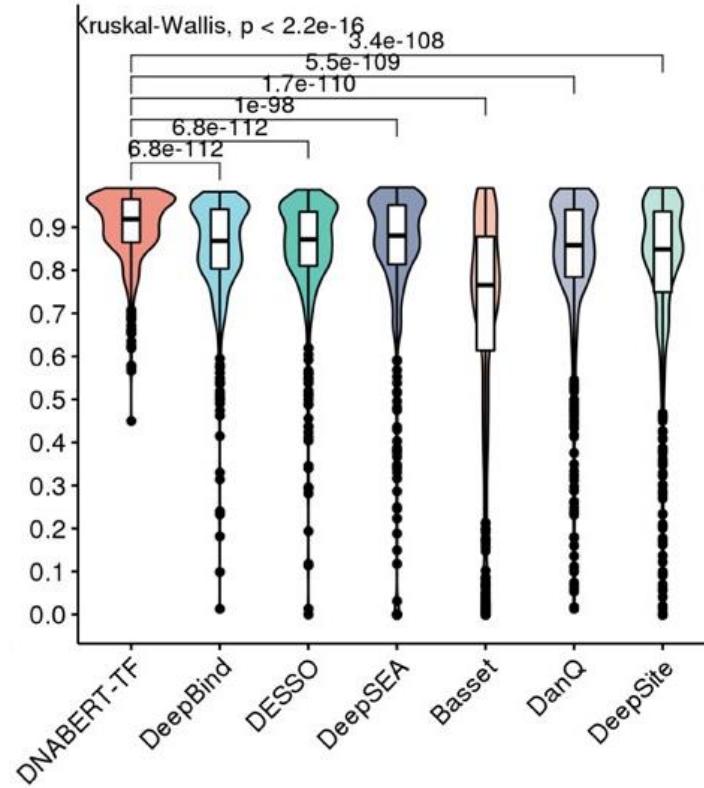


# Results

- 2. DNABERT-TF accurately identifies TFBS

variable

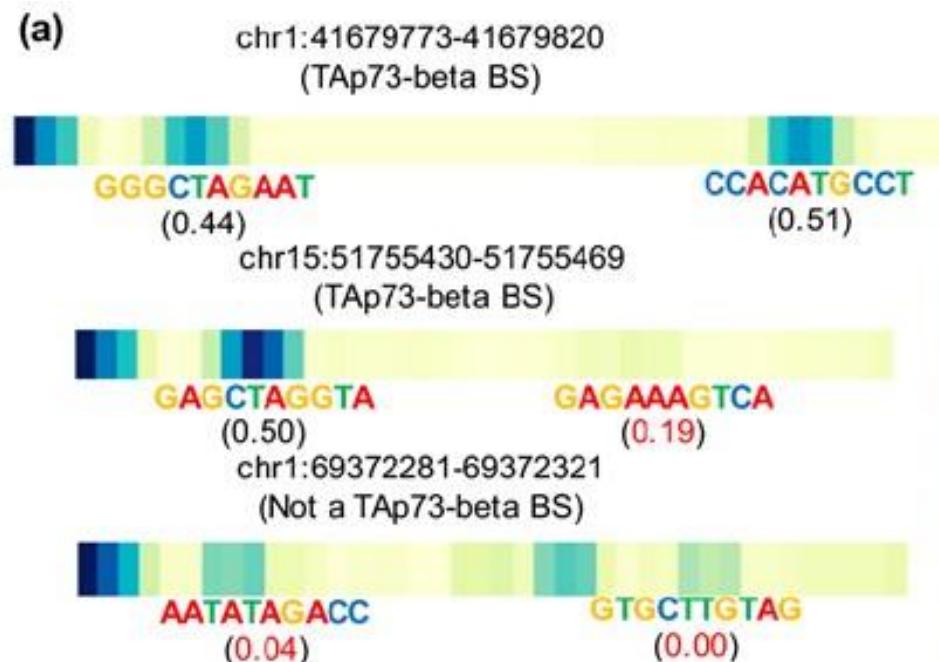
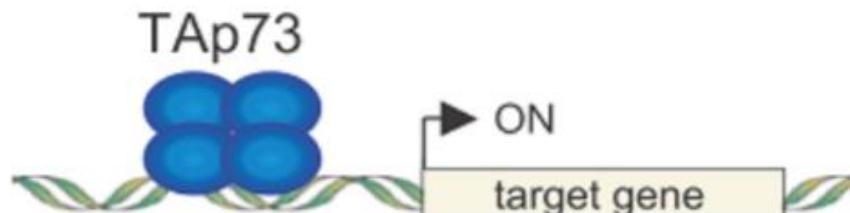
DNABERT-TF	DeepBind	DESSO	Basset	DeepSite
DeepSEA	DanQ			



→ DNABERT accurately identifies TFBSs

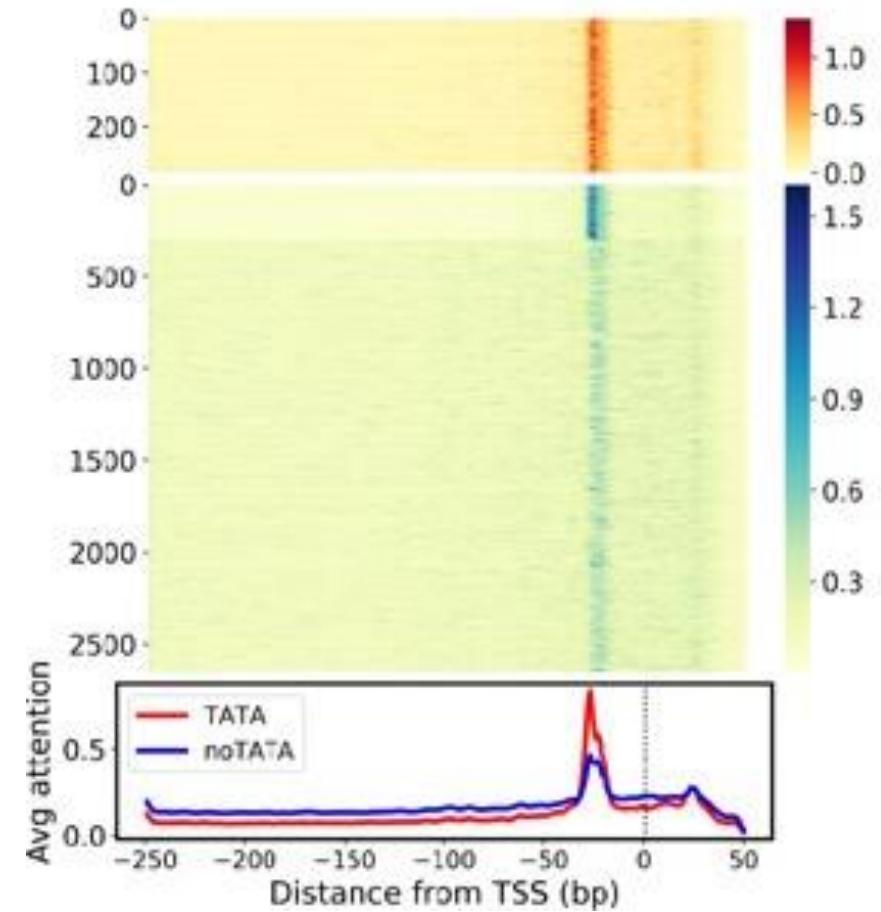
# Results

- **3. Visualization of important regions, contexts and sequence motifs**
  - DNABERT-viz module
    - Since it uses attention mechanism, we can directly interpret the scores
    - ex) Attention maps
      - 2 ChIP-seq-validated TAp73-beta binding sites (top, middle)
      - Non-binding site (bottom)



# Results

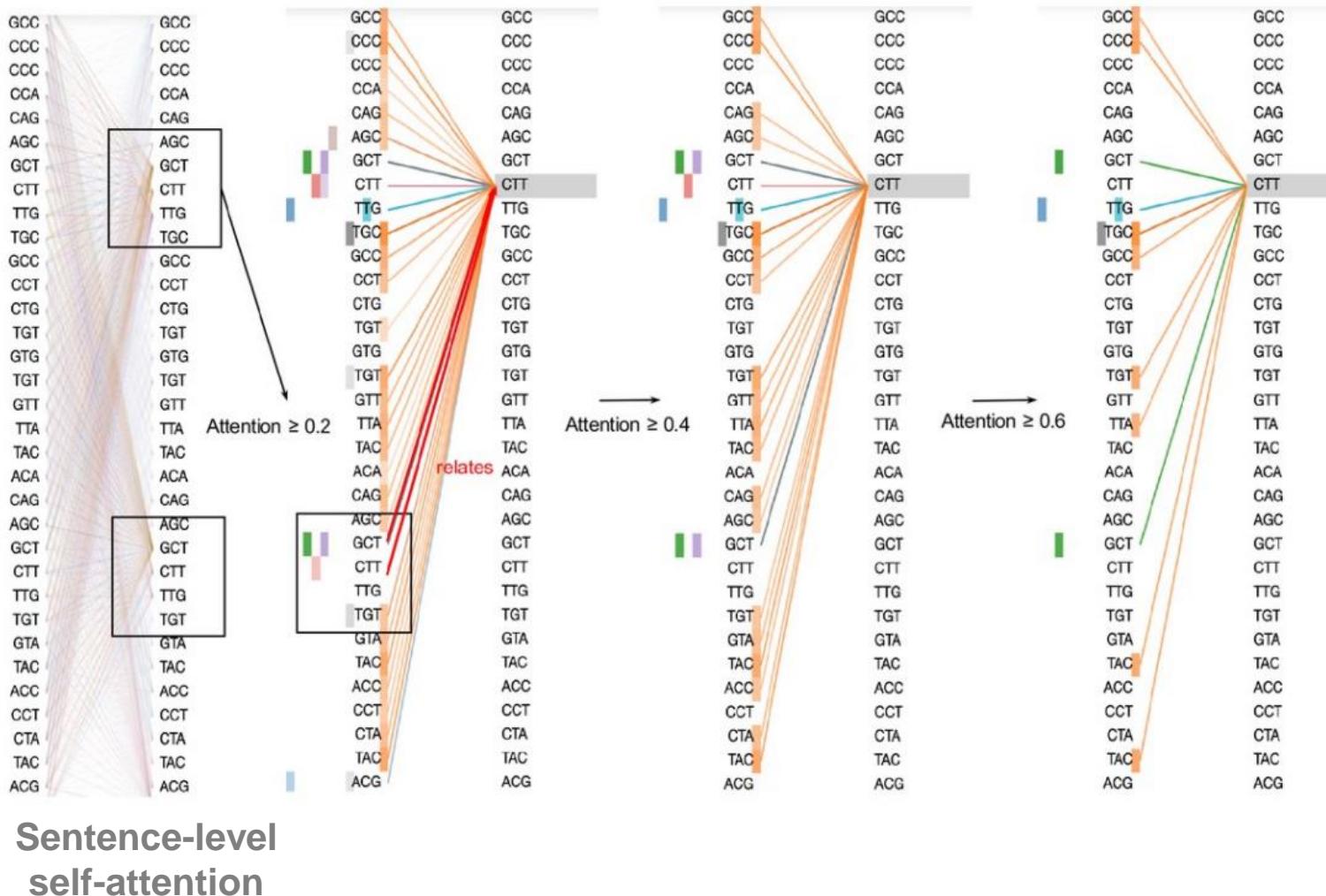
- **3. Visualization of important regions, contexts and sequence motifs**
  - DNABERT-viz module
    - ex) Attention landscape  
(in Prom-300 test set)
      - TATA promoters (top)  
→ DNABERT consistently put high attention upon -20 to -30 bp region upstream of TSS where TATA box is located
      - noTATA promoters (bottom)  
→ Majority of non-TATA promoters show a more scattered attention pattern



# Results

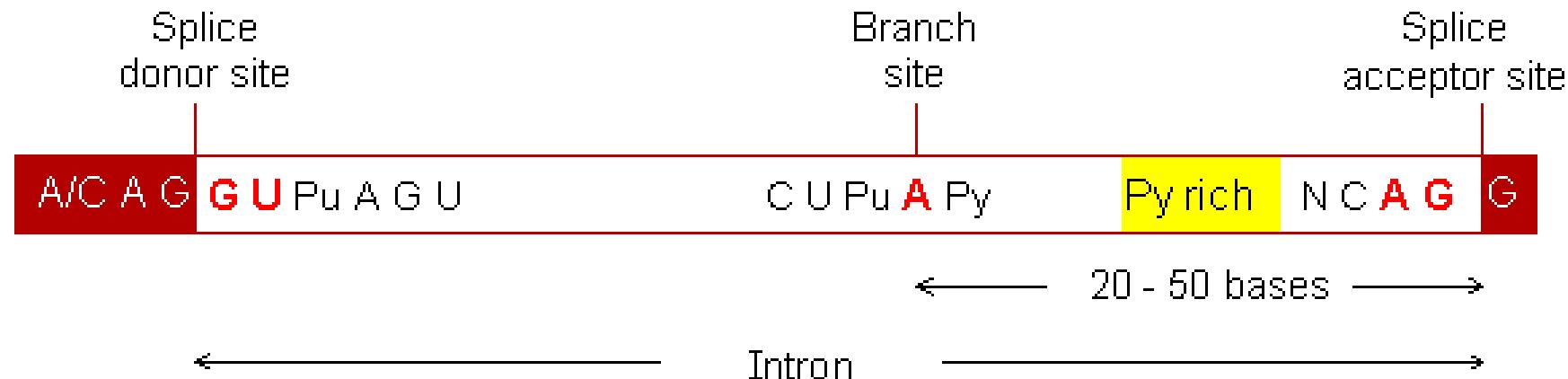
- **3. Visualization of important regions, contexts and sequence motifs**

- DNABERT-viz module
  - Direct visualization of **contextual relationship** **within any input sequence**
  - ex) Attention-head plots of p53 binding site



# Results

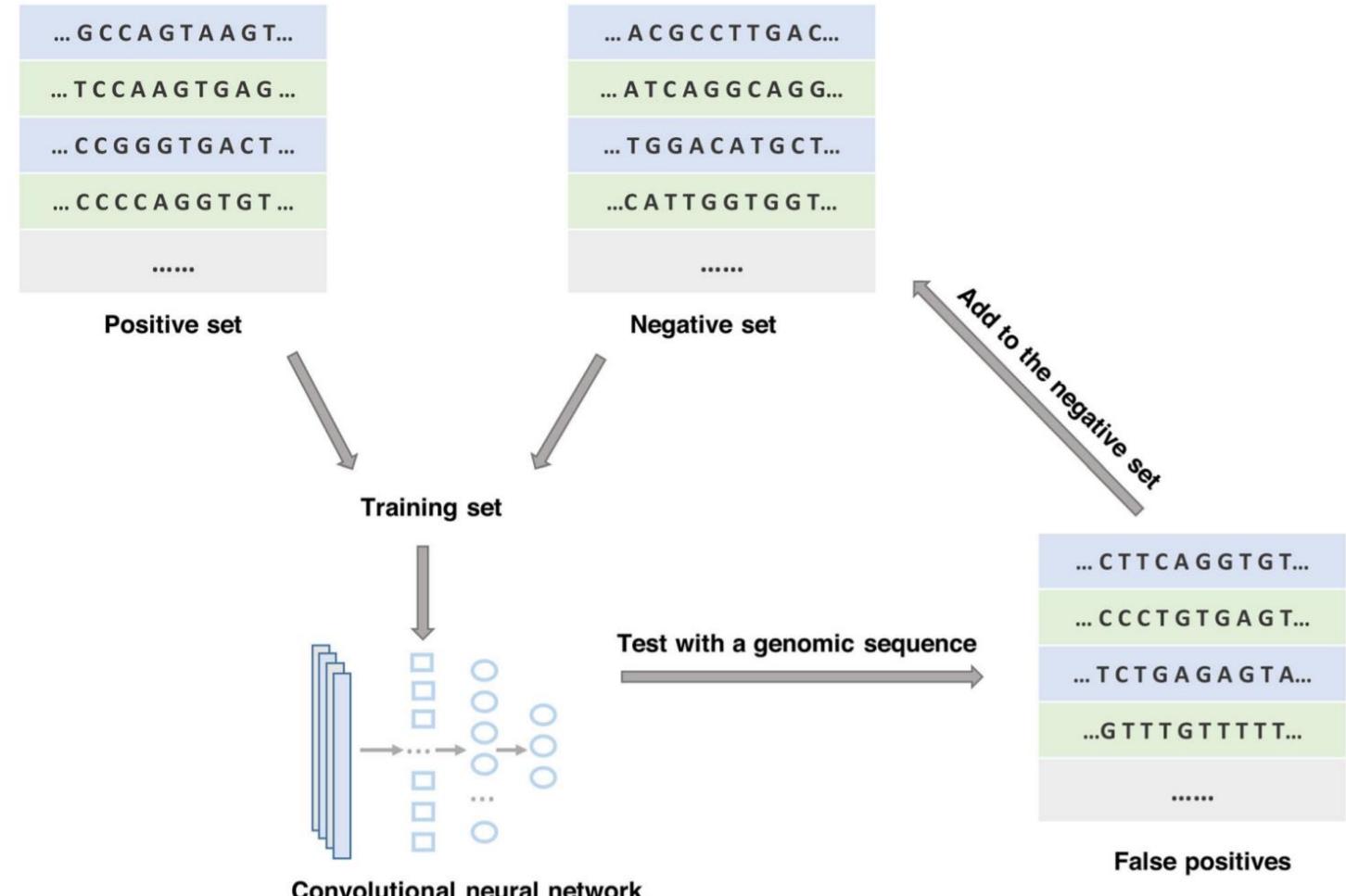
- **4. DNABERT-Splice accurately recognized canonical and non-canonical splice sites**
  - Fine-tuned DNABERT with GRCh38 FASTA file from Ensemble release 99
    - Positive sequences : Extracted from 400-bp long sequences around the donor and acceptor sites of randomly selected exons



# Results

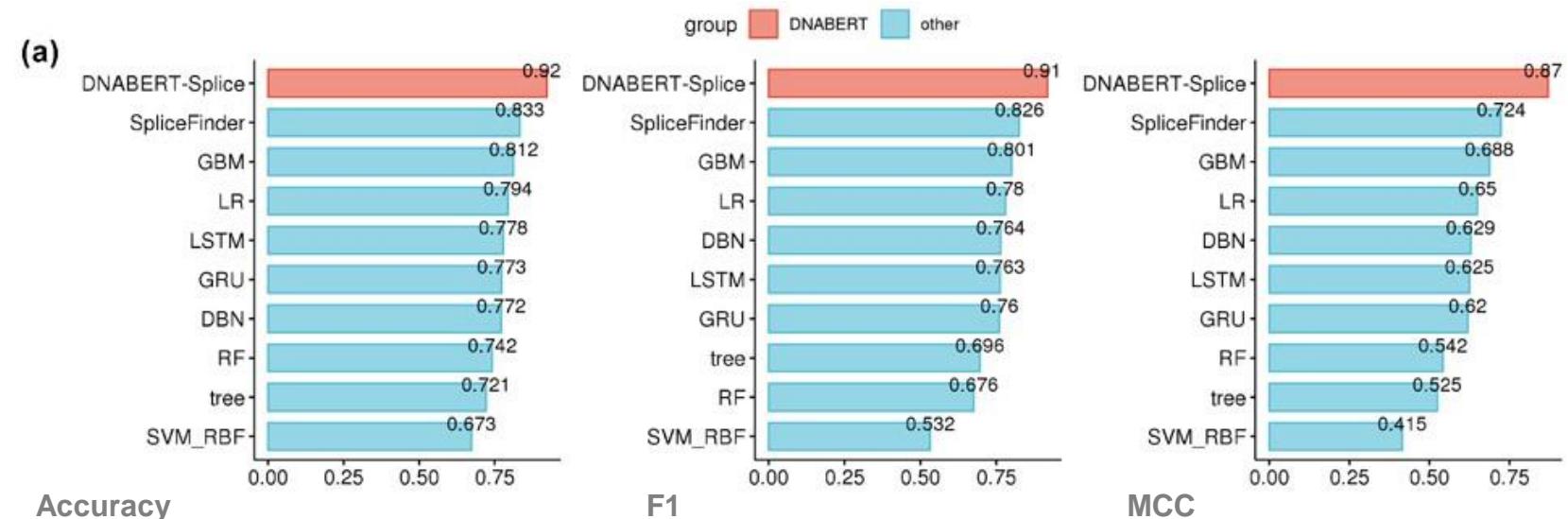
- 4. DNABERT-Splice accurately recognized canonical and non-canonical splice sites**

- To detect non-canonical splice sites without GT or AG dimers
- : Iterative data augmentation
  - Repeatedly built new models to predict on sliding-window-based scans



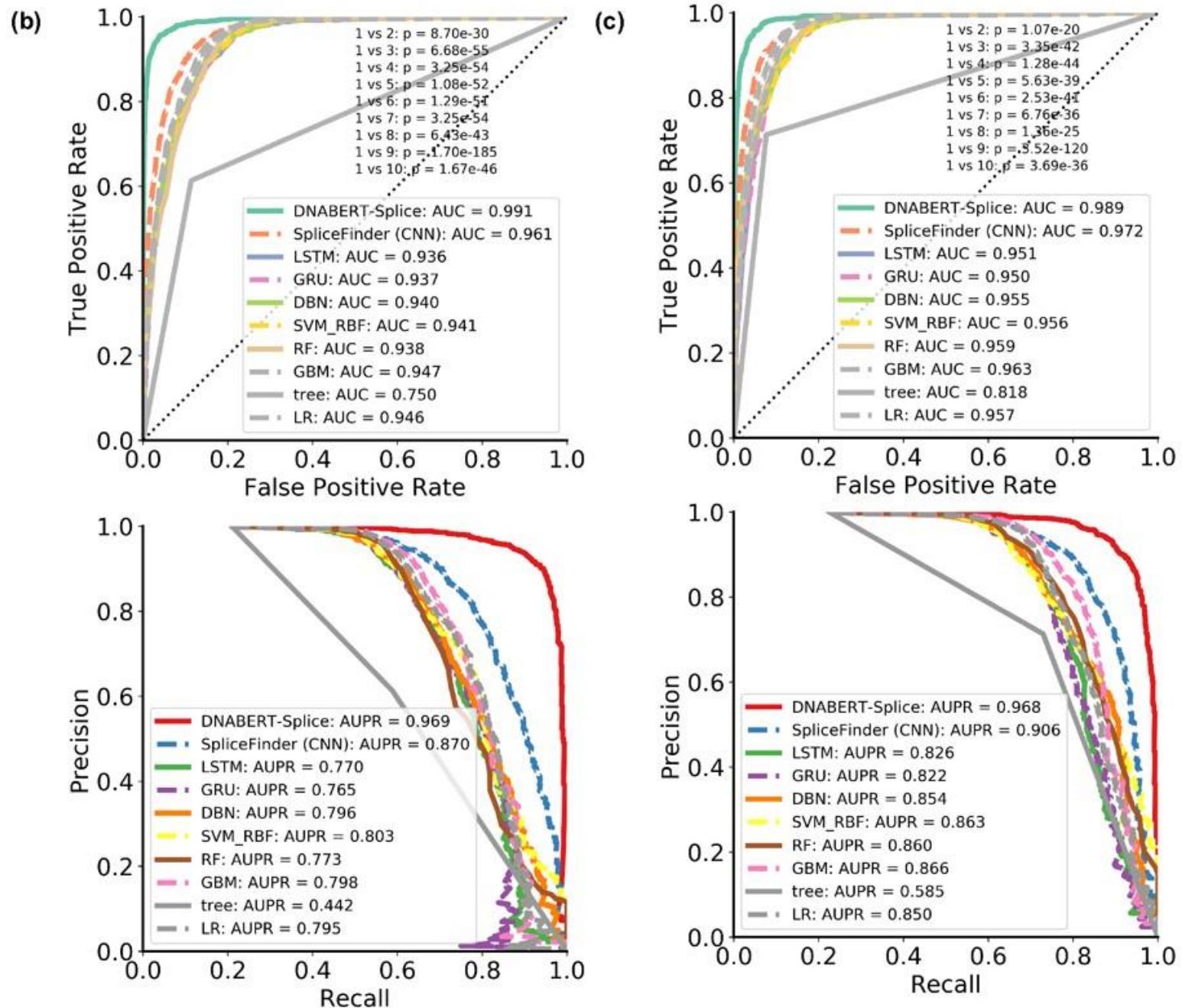
# Results

- **4. DNABERT-Splice accurately recognized canonical and non-canonical splice sites**
  - Fine-tuned DNABERT with GRCh38 FASTA file from Ensemble release 99
    - DNABERT-splice showed globally high attention upon intronic regions, highlighting the presence and functional importance of various intronic splicing enhancers (ISEs) and silencers (ISSs) acting as CREs for splicing



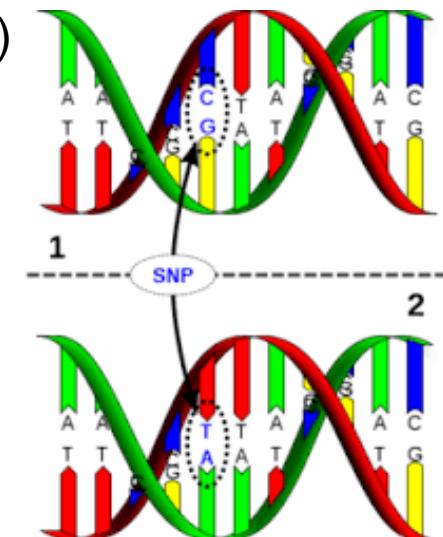
# Results

- DNABERT significantly outperforms other models in finding splice sites
  - (b) : Splice donor dataset
  - (c) : Splice acceptor dataset



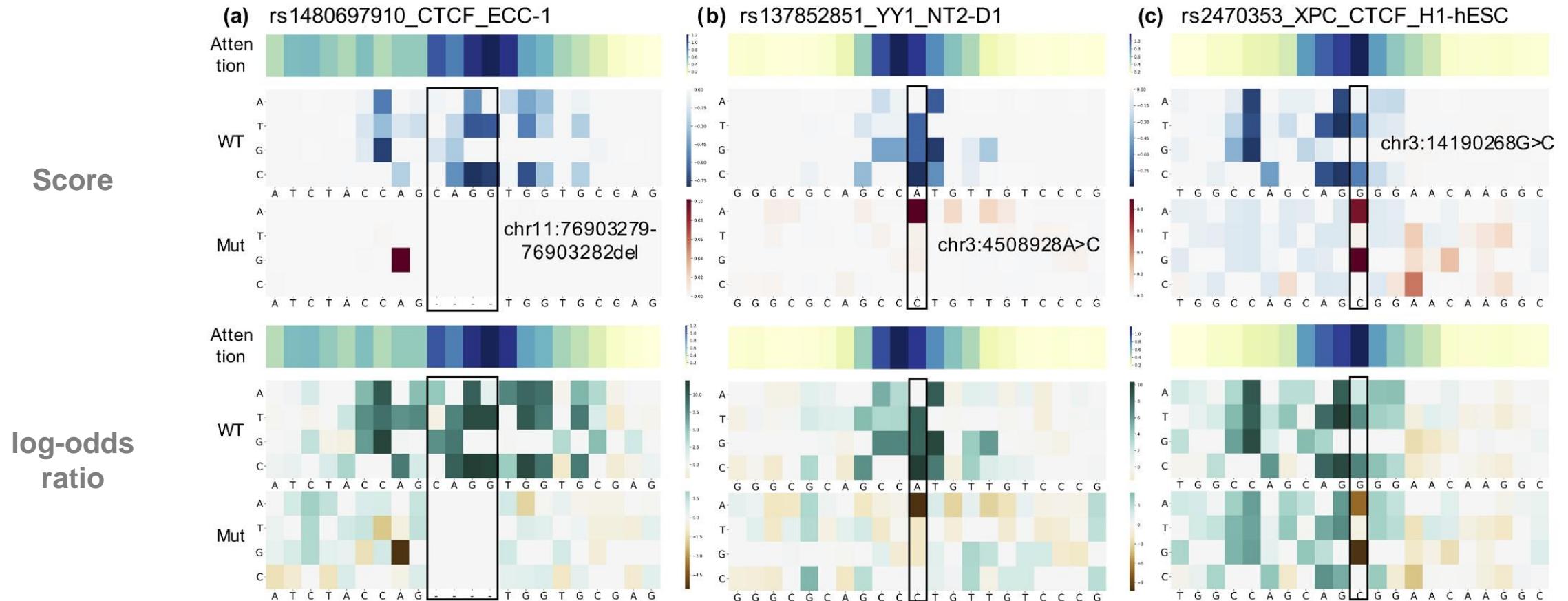
# Results

- **5. Identifying functional genetic variants with DNABERT**
  - Fine-tuned DNABERT with dbSNP release 153 : 700 million functional variants
    - Selected only those variants located inside DNABERT-predicted high-attention regions and repeated the predictions, using sequence with altered alleles ( $p(s) \rightarrow p(s')$ )
      - Score change =  $\Delta p = (p(s') - p(s)) \max(p(s'), p(s))$
      - Log odds ratio =  $\log_2 OR = \log \frac{p(s)}{1-p(s)} - \log \frac{p(s')}{1-p(s')}$
    - Candidate variants resulting in significant changes in prediction probability  
→ queried in ClinVar, GRASP, and NHGRI-EBI GWAS catalog
    - In Prom-300 dataset
      - 24.7% and 34.1% of dbSNP common variants we identified using TATA and non-TATA promoters are present in at least one of the three databases



# Results

- 5. Identifying functional genetic variants with DNABERT



→ DNABERT identifies functional genetic variants  
DNABERT consistently shows highest attention at/around the variants of interest

- (a) : Rare, pathogenic 4bp deletion  
Cause Usher syndrome
- (b) : Rare SNV  
Multiple sulfates deficiency
- (c) : Common risk variant of  
pancreatic cancer

# Discussion

- DNABERT achieved superior performance across various downstream DNA sequence prediction tasks by largely surpassing existing tools
- DNABERT tackles the problem of sequence specificity prediction with a “top-down” approach by first developing general understanding of DNA language via self-supervised pre-training and then applying it to specific tasks, in contrast to the traditional “bottom-up” approach using task-specific data
- Effectively learn from DNA context with great flexibility adapting to multiple situations, and enhanced performance with limited data

# Google Colab

# Thank You!