

TEXT SUMMARIZE & BERTVIZ

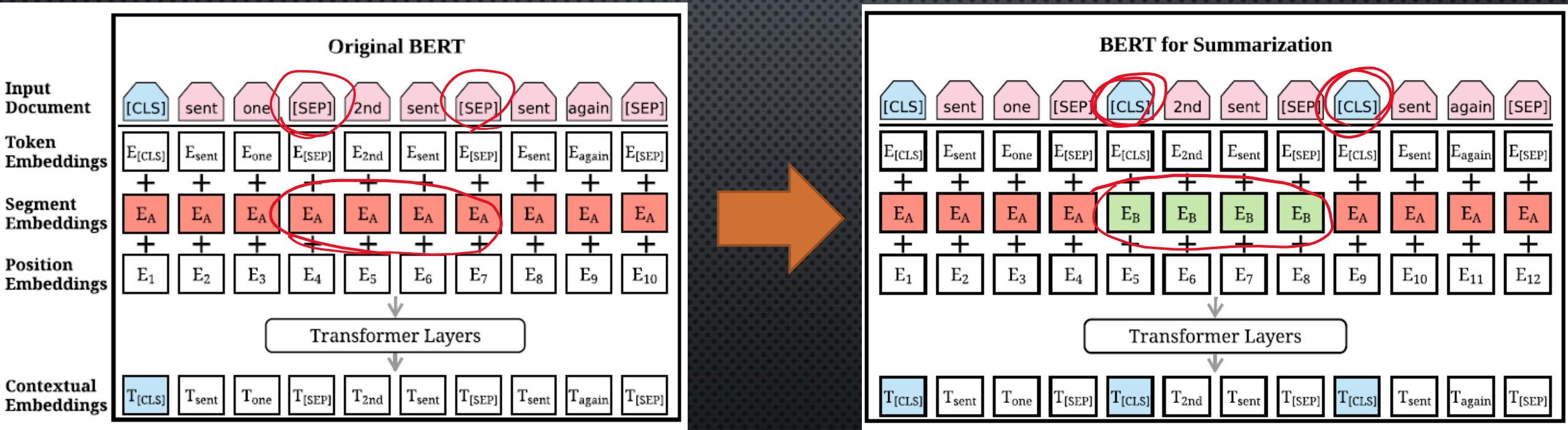
# OUTLINE

- 使用TRANSFORMER擷取文章摘要
- 使用BERTVIZ視覺化模型的ATTENTIONS



# 從 Bert Finetune來處理摘要擷取(text summarization)

CLS: sentence embedding>用來後續的預測任務  
E<sub>A</sub>, E<sub>B</sub>: segmentation embeddings>根據奇數偶數句來區分



# BERT-SUM-TEXT

[1908.08345] TEXT SUMMARIZATION WITH PRETRAINED ENCODERS (ARXIV.ORG)

- 什麼是extractive summarization?
  - 對每個句子做Binary classification
    - 是重點/不是重點

簡單來說，既然前面已經對句子做embedding，  
我們丟進去分類器的單位就不是token為單位，而是sentence

優點> 取出的摘要都是完整的句子

缺點> 沒有辦法重組句子(精煉)



# Language Model : Unigram

$$\circ p(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

$$= p(w_1) \cdot p(w_2) \cdot p(w_3) \dots p(w_n)$$

$$\circ p(\text{今天, 是, 春节, 我们, 都, 休息}) \quad \checkmark$$

$$= p(\underline{\text{今天}}) \cdot p(\underline{\text{是}}) \cdot p(\underline{\text{春节}}) \cdot p(\underline{\text{我们}}) \cdot p(\underline{\text{都}}) \cdot p(\underline{\text{休息}})$$

$$\circ p(\text{今天, 春节, 是, 都, 我们, 休息}) \quad -$$

$$= p(\underline{\text{今天}}) \cdot p(\underline{\text{春节}}) \cdot p(\underline{\text{是}}) \cdot p(\underline{\text{都}}) \cdot p(\underline{\text{我们}}) \cdot p(\underline{\text{休息}})$$

Bigram: 詞語出現的機率和前一個詞有關  
(考慮兩個詞的順序)

Unigram: 假設詞語之間是互相獨立 (不考慮順序)

## Language Model : Bigram $\leftarrow$ 1st order

markov assumption

$$\circ p(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

$$= p(w_1) \cdot p(w_2|w_1) \cdot p(w_3|w_2) \dots p(w_n|w_{n-1}) = p(w_1) \cdot \prod_{i=2}^n p(w_i|w_{i-1})$$

$$\circ p(\text{今天, 是, 春节, 我们, 都, 休息})$$

$$= p(\underline{\text{今天}}) \cdot p(\underline{\text{是}}|\underline{\text{今天}}) \cdot p(\underline{\text{春节}}|\underline{\text{是}}) \cdot p(\underline{\text{我们}}|\underline{\text{春节}}) \cdot p(\underline{\text{都}}|\underline{\text{我们}}) \cdot p(\underline{\text{休息}}|\underline{\text{都}})$$

$$\circ p(\text{今天, 春节, 是, 都, 我们, 休息})$$

$$= p(\underline{\text{今天}}) \cdot p(\underline{\text{春节}}|\underline{\text{今天}}) \cdot p(\underline{\text{是}}|\underline{\text{春节}}) \cdot p(\underline{\text{都}}|\underline{\text{是}}) \cdot p(\underline{\text{我们}}|\underline{\text{都}}) \cdot p(\underline{\text{休息}}|\underline{\text{我们}})$$

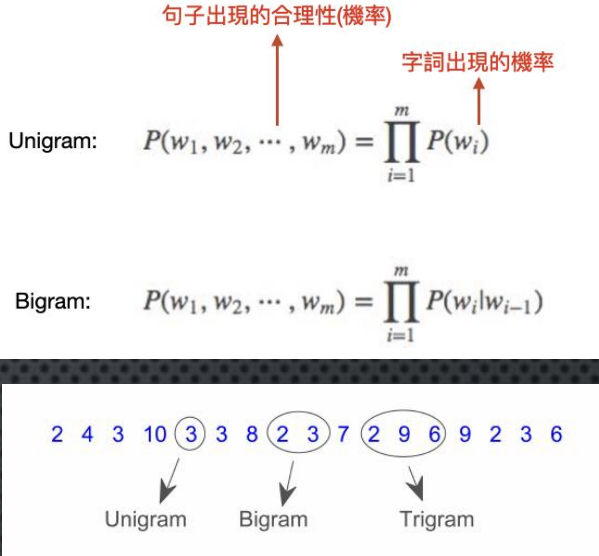
# R1: Rouge 1

# R2: Rouge 2

# RL: RougeL

- + R1 : police killed the gunman.
- + R2 : the gunman was shot down by police.
- C1 : police ended the gunman.
- C2 : the gunman murdered police.

$$ROUGE - N = \frac{\sum_{S \in References} \sum_{gram_n \in S} Count_{match}(gram_n)}{\sum_{S \in References} \sum_{gram_n \in S} Count(gram_n)}$$
$$R_{lcs} = \frac{LCS(X,Y)}{m}$$
$$P_{lcs} = \frac{LCS(X,Y)}{n}$$
$$ROUGE - L = F_{lcs} = \frac{(1+\beta^2)R_{lcs}P_{lcs}}{R_{lcs}+\beta^2P_{lcs}}$$



Scoring	相對於誰	Calc	
R1	C1>R1,R2	(3+3)/(4+7)=6/11	C1和R1,R2個重複了3個詞 R1,R2分別有4,7個詞
R1	C2>R1,R2	(3+3)/(4+7)=6/11	
R2	C1>R1,R2	(1+1)/(3+6)=2/9	重疊的bigram只有1個 (the gunman) R1,R2分別 有3,6個bigram
R2	C2>R1,R2	(1+1)/(3+6)=2/9	
RL	C1>R1	Avg(3/4, 3/4)=3/8	LCS=police the gunman(3) R1,C1長度=4



前3句  
 $R_1, R_2 \rightarrow$  overlap  
 higher is better  
 $R_1$ : unigram  
 $R_2$ : bigram

Model	R1	R2	RL
ORACLE	52.59	31.24	48.87
LEAD-3	40.42	17.62	36.67
Extractive			
SUMMARUNNER (Nallapati et al., 2017)	39.60	16.20	35.30
REFRESH (Narayan et al., 2018b)	40.00	18.20	36.60
LATENT (Zhang et al., 2018)	41.05	18.77	37.54
NEUSUM (Zhou et al., 2018)	41.59	19.01	37.98
SUMO (Liu et al., 2019)	41.00	18.40	37.20
TransformerEXT	40.90	18.02	37.17
Abstractive			
PTGEN (See et al., 2017)	36.44	15.66	33.42
PTGEN+COV (See et al., 2017)	39.53	17.28	36.38
DRM (Paulus et al., 2018)	39.87	15.82	36.90
BOTTOMUP (Gehrmann et al., 2018)	41.22	18.68	38.34
DCA (Celikyilmaz et al., 2018)	41.69	19.47	37.92
TransformerABS	40.21	17.76	37.09
BERT-based			
BERTSUMEXT	43.25	20.24	39.63
BERTSUMEXT w/o interval embeddings	43.20	20.22	39.59
BERTSUMEXT (large)	43.85	20.34	39.90
BERTSUMABS	41.72	19.39	38.76
BERTSUMEXTABS	42.13	19.60	39.18

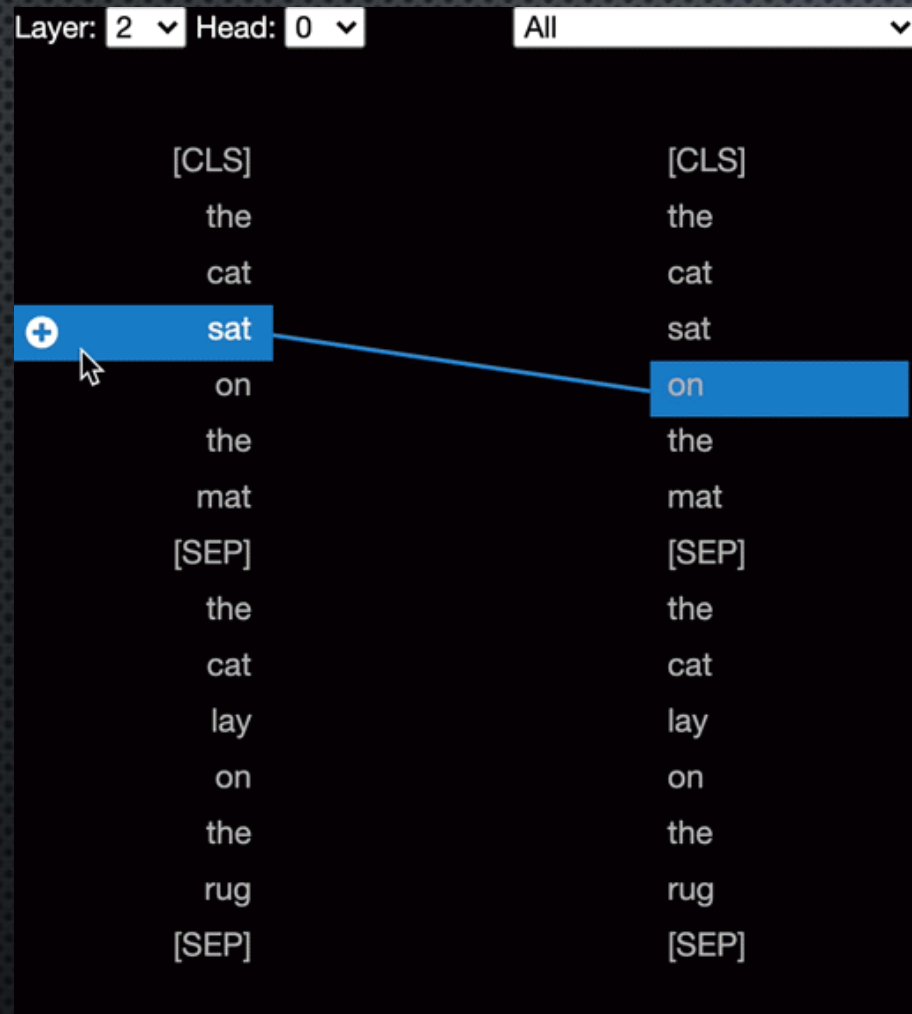
DEMO

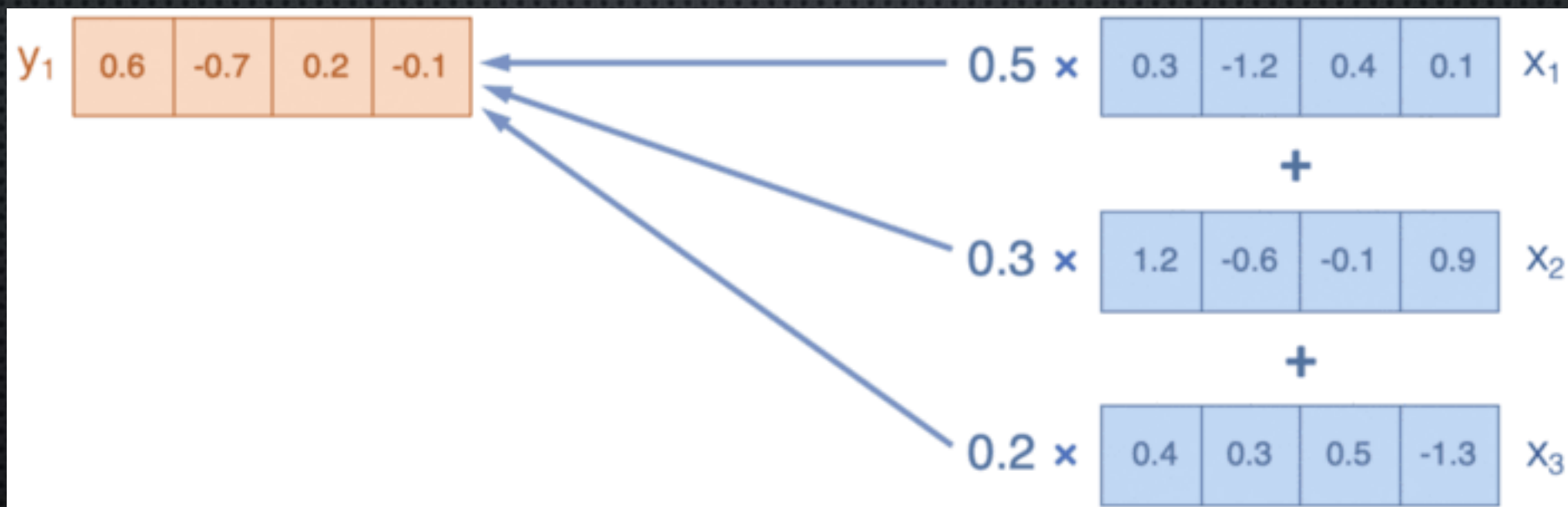
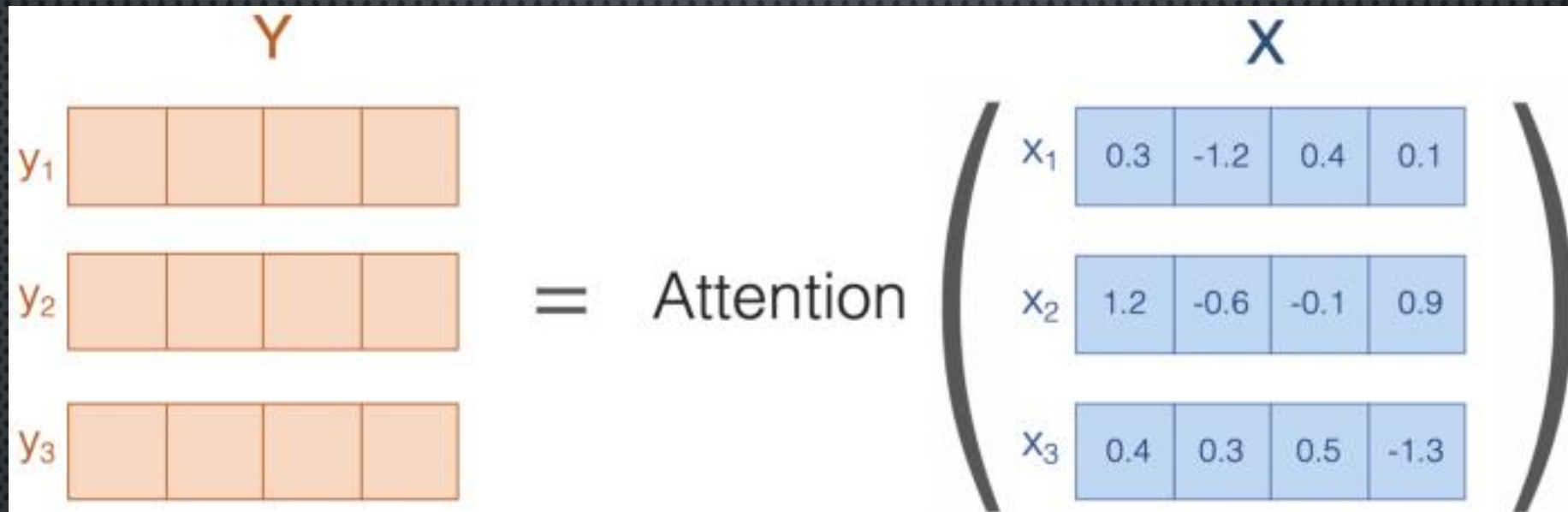
使用TRANSFORMER來判斷新聞摘要



# BERTVIZ

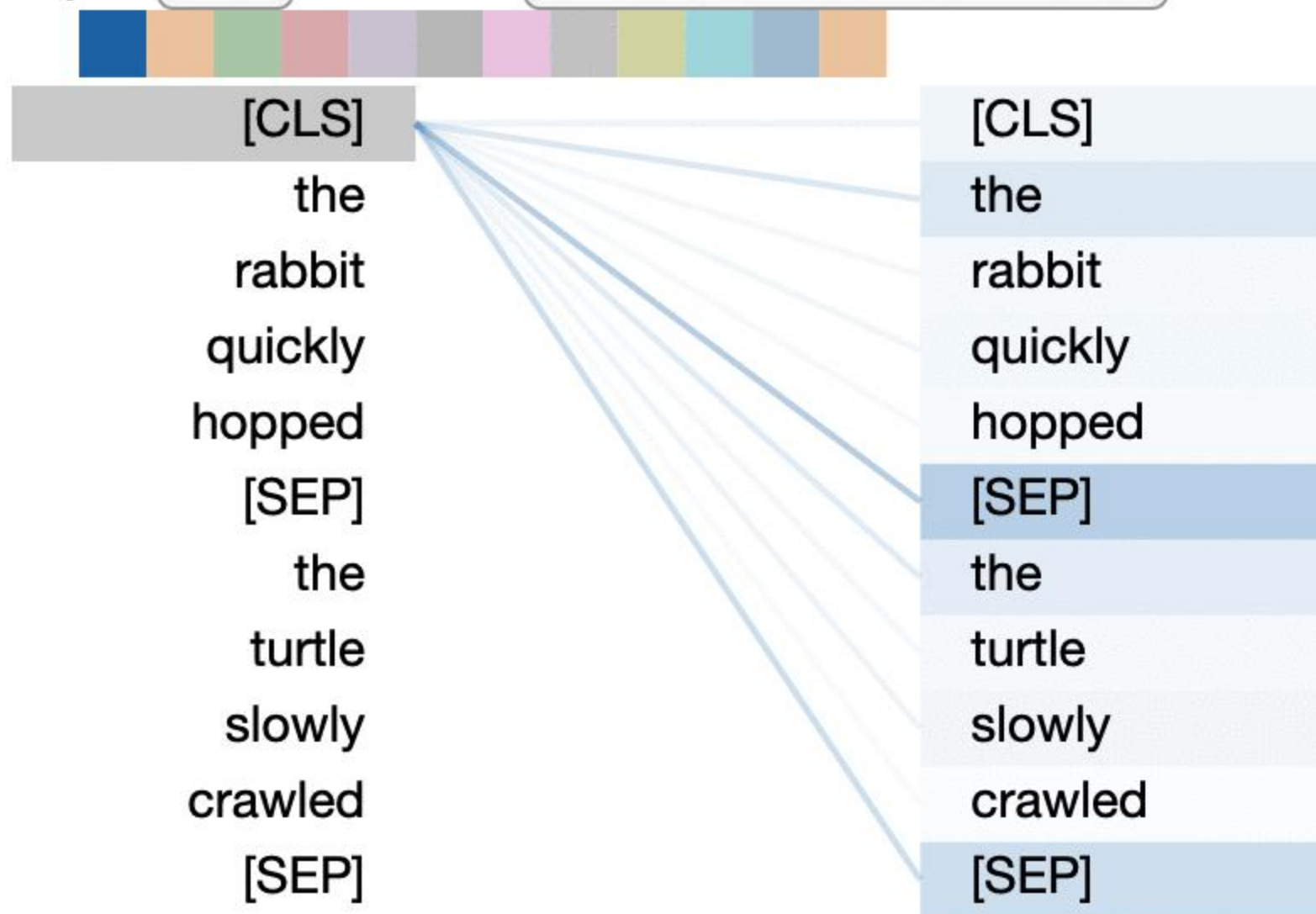
[JESSEVIG/BERTVIZ: BERTVIZ: VISUALIZE ATTENTION IN NLP MODELS \(BERT, GPT2, BART, ETC.\) \(GITHUB.COM\)](https://github.com/jessevig/bertviz)







Layer: 0 ▾ Attention: All ▾



	query				key					score	softmax
dog	0.3	-0.2	0.4	•	0.5	-0.9	0.2	The	=	0.4	0.4
dog	0.3	-0.2	0.4	•	1.1	-0.3	0.5	dog	=	0.6	0.5
dog	0.3	-0.2	0.4	•	-1.0	0.3	-0.7	ran	=	-0.6	0.1

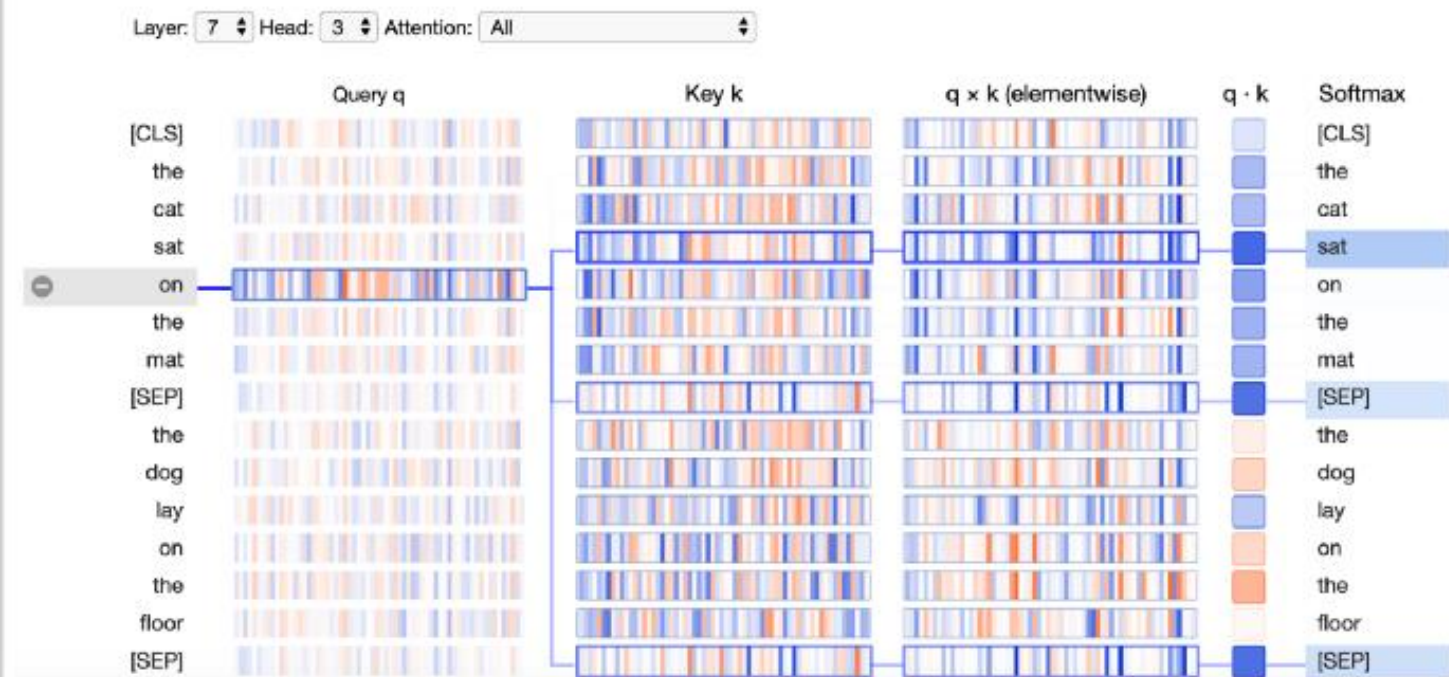
**Query  $q$ :** the query vector  $q$  encodes the word on the left that is paying attention, i.e. the one that is “querying” the other words. In the example above, the query vector for “on” (the selected word) is highlighted.

**Key  $k$ :** the key vector  $k$  encodes the word on the right to which attention is being paid. The key vector and the query vector together determine a compatibility score between the two words.

**$q \times k$  (elementwise):** the elementwise product between the query vector of the selected word and each of the key vectors. This is a precursor to the dot product (the sum of the elementwise product) and is included for visualization purposes because it shows how individual elements in the query and key vectors contribute to the dot product.

**$q \cdot k$ :** the scaled dot product (see above) of the selected query vector and each of the key vectors. This is the unnormalized attention score.

**Softmax:** the softmax of the scaled dot product. This normalizes the attention scores to be positive and sum to one.





DEMO

使用BERTVIZ視覺化ATTENTION