

BERT - Experiments

```
[40] text = "[CLS] 等到潮水 [MASK] 了，就知道誰沒穿褲子。"
      tokens = tokenizer.tokenize(text)
      ids = tokenizer.convert_tokens_to_ids(tokens)

      print(text)
      print(tokens[:10], '...')
      print(ids[:10], '...')

⇒ [CLS] 等到潮水 [MASK] 了，就知道誰沒穿褲子。
[['[CLS]', '等', '到', '潮', '水', '[MASK]', '了', ' ', ' ', '就', '知'] ...  
[101, 5023, 1168, 4060, 3717, 103, 749, 8024, 2218, 4761] ...
```

Masked Language Model (MLM)

```
[8] # 除了 tokens 以外還需要辨別句子的 segment ids
tokens_tensor = torch.tensor([ids]) # (1, seq_len)
segments_tensors = torch.zeros_like(tokens_tensor) # (1, seq_len)
maskedLM_model = BertForMaskedLM.from_pretrained(PRETRAINED_MODEL_NAME)
clear_output()
```

```
[9] # 使用 masked LM 估計 [MASK] 位置所代表的實際 token
maskedLM_model.eval()
with torch.no_grad():
    outputs = maskedLM_model(tokens_tensor, segments_tensors)
    predictions = outputs[0]
    # (1, seq_len, num_hidden_units)
del maskedLM_model
```

```
[11] # 將 [MASK] 位置的機率分佈取 top k 最有可能的 tokens 出來
masked_index = 5
k = 3
probs, indices = torch.topk(torch.softmax(predictions[0, masked_index], -1), k)
predicted_tokens = tokenizer.convert_ids_to_tokens(indices.tolist())
```

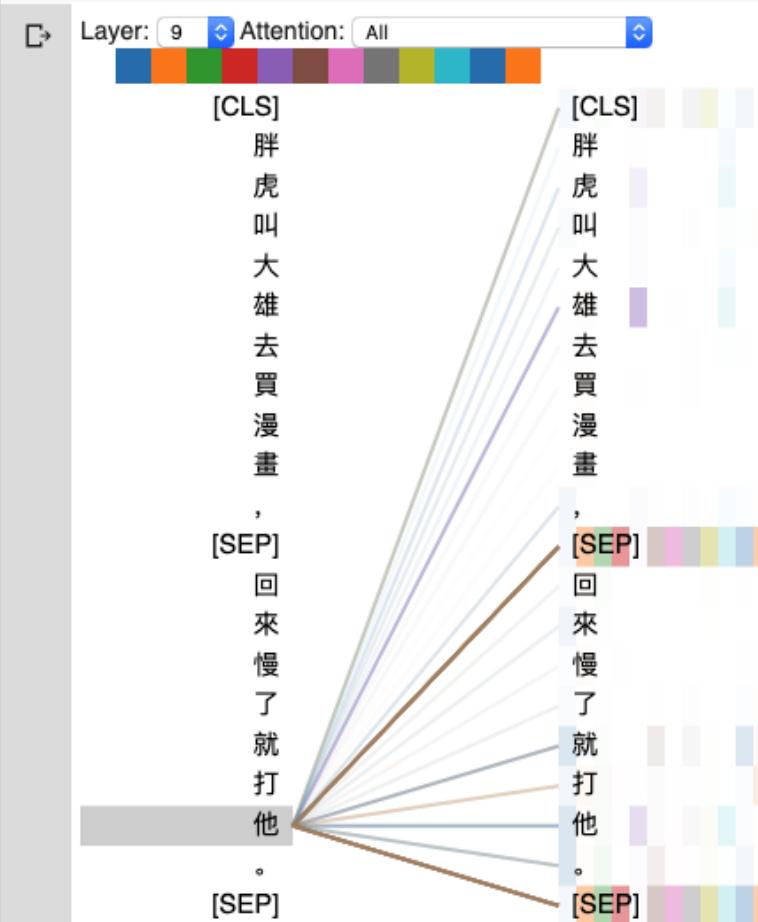
```
▶ # 顯示 top k 可能的字。取 top 1 當作預測值
print("輸入 tokens : ", tokens[:10], '...')
print('-' * 50)
for i, (t, p) in enumerate(zip(predicted_tokens, probs), 1):
    tokens[masked_index] = t
    print("Top {} ({:2}%): {}".format(i, int(p.item() * 100), tokens[:10]), '...')
```

```
⇒ 輸入 tokens : [['[CLS]', '等', '到', '潮', '水', '過', '了', ' ', ' ', '就', '知'] ...
Top 1 (67%): [['[CLS]', '等', '到', '潮', '水', '來', '了', ' ', ' ', '就', '知'] ...
Top 2 (25%): [['[CLS]', '等', '到', '潮', '水', '濕', '了', ' ', ' ', '就', '知'] ...
Top 3 ( 2%): [['[CLS]', '等', '到', '潮', '水', '過', '了', ' ', ' ', '就', '知'] ...
```

Natural Language Inference (NLI)

Attention Visualization

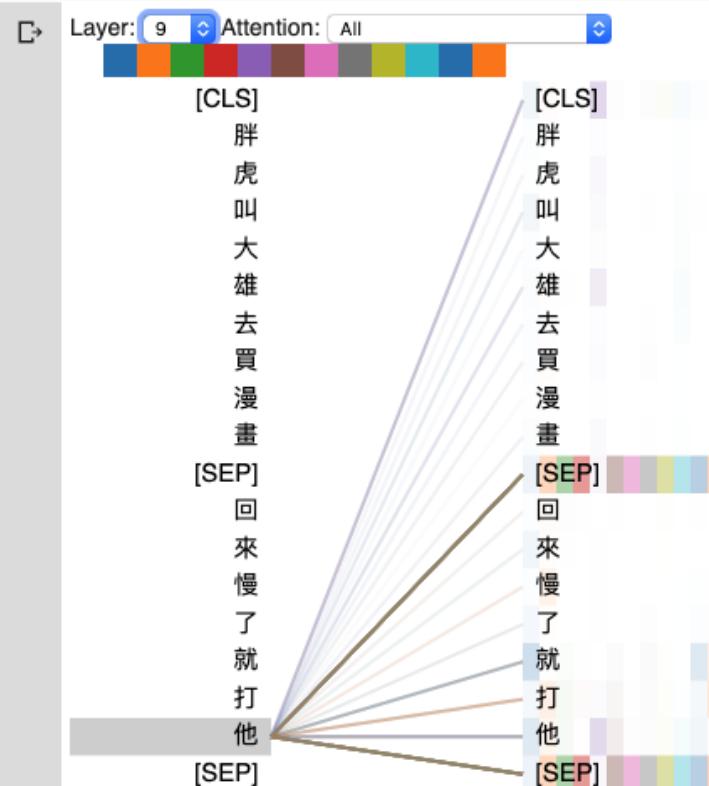
```
model_type = 'bert'  
bert_version = 'bert-base-chinese'  
  
bertviz_model = BertModel.from_pretrained(bert_version)  
bertviz_tokenizer = BertTokenizer.from_pretrained(bert_version)  
  
# 情境 1  
sentence_a = "胖虎叫大雄去買漫畫，"  
sentence_b = "回來慢了就打他。"  
call_html()  
show(bertviz_model, model_type, bertviz_tokenizer, sentence_a, sentence_b)
```



Natural Language Inference (NLI)

Attention Visualization

```
model_type = 'bert'  
bert_version = 'bert-base-chinese'  
  
bertviz_model = BertModel.from_pretrained(bert_version)  
bertviz_tokenizer = BertTokenizer.from_pretrained(bert_version)  
  
# 情境 1  
sentence_a = "胖虎叫大雄去買漫畫"  
sentence_b = "回來慢了就打他"  
call_html()  
show(bertviz_model, model_type, bertviz_tokenizer, sentence_a, sentence_b)
```



Source Code Analysis

- `input_file`
- `output_file`
- `vocab_file` : 指定字典路徑 (Google提供, Tokenization 會用到)
- `do_lower_case`
- `max_seq_length` : 每一條訓練數據 (兩句話) 相加後的最大長度限制 (128)
- `max_predictions_per_seq` : 每一條訓練數據mask的最大數量
- `random_seed` : 隨機種子 (12345)
- `dupe_factor` : 對文檔多次重複隨機產生訓練集 · 隨機的次數 (10)
- `masked_lm_prob` : 訓練數據產生mask的概率 (0.15)
- `max_predictions_per_seq*masked_lm_prob` 數量的mask
- `short_seq_prob` : 為了縮小預訓練和微調過程的差距 · 從而概率產生小於 `max_seq_length` 的訓練數據

```
37
38 flags = tf.flags
39
40 FLAGS = flags.FLAGS
41
42 flags.DEFINE_string("input_file", None,
43                      "Input raw text file (or comma-separated list of files.)")
44
45 flags.DEFINE_string(
46     "output_file", None,
47     "Output TF example file (or comma-separated list of files.)")
48
49 flags.DEFINE_string("vocab_file", None,
50                      "The vocabulary file that the BERT model was trained on.")
51
52 flags.DEFINE_bool(
53     "do_lower_case", True,
54     "Whether to lower case the input text. Should be True for uncased "
55     "models and False for cased models.")
56
57 flags.DEFINE_bool(
58     "do_whole_word_mask", False,
59     "Whether to use whole word masking rather than per-WordPiece masking.")
60
61 flags.DEFINE_integer("max_seq_length", 128, "Maximum sequence length.")
62
63 flags.DEFINE_integer("max_predictions_per_seq", 20,
64                      "Maximum number of masked LM predictions per sequence.")
65
66 flags.DEFINE_integer("random_seed", 12345, "Random seed for data generation.")
67
68 flags.DEFINE_integer(
69     "dupe_factor", 10,
70     "Number of times to duplicate the input data (with different masks.).")
71
72 flags.DEFINE_float("masked_lm_prob", 0.15, "Masked LM probability.")
73
74 flags.DEFINE_float(
75     "short_seq_prob", 0.1,
76     "Probability of creating sequences which are shorter than the6"
77     "maximum length.")
```

預訓練 → 核心模型構建 → 訓練

create_pretraining_data.py → modeling.py → run_pretraining.py

```
447
448 def main(_):
449     tf.logging.set_verbosity(tf.logging.INFO)
450
451     tokenizer = tokenization.FullTokenizer(
452         vocab_file=FLAGS.vocab_file, do_lower_case=FLAGS.do_lower_case)
453
454     input_files = []
455     for input_pattern in FLAGS.input_file.split(","):
456         input_files.extend(tf.gfile.Glob(input_pattern)) # 獲得輸入文件列表
457
458     tf.logging.info("*** Reading from input files ***")
459     for input_file in input_files:
460         tf.logging.info(" %s", input_file)
461
462     rng = random.Random(FLAGS.random_seed)           # random_seed
463     instances = create_training_instances(          # 創建訓練實例
464         input_files, tokenizer, FLAGS.max_seq_length, FLAGS.dupe_factor,
465         FLAGS.short_seq_prob, FLAGS.masked_lm_prob, FLAGS.max_predictions_per_seq,
466         rng)
467
```

def create_training_instances

創建訓練實例

Pre-training

Tokenization

Input File Format Restriction

Tokenization

```
191 def create_training_instances(input_files, tokenizer, max_seq_length,
192                               dupe_factor, short_seq_prob, masked_lm_prob,
193                               max_predictions_per_seq, rng):
194     """Create `TrainingInstance`'s from raw text."""
195     all_documents = []
196
197     # Input file format:
198     # (1) One sentence per line. These should ideally be actual sentences, not
199     # entire paragraphs or arbitrary spans of text. (Because we use the
200     # sentence boundaries for the "next sentence prediction" task).
201     # (2) Blank lines between documents. Document boundaries are needed so
202     # that the "next sentence prediction" task doesn't span between documents.
203     for input_file in input_files:
204         with tf.gfile.GFile(input_file, "r") as reader:
205             while True:
206                 line = tokenization.convert_to_unicode(reader.readline())
207                 if not line:
208                     break
209                 line = line.strip()
210
211                 # Empty lines are used as document delimiters
212                 if not line:
213                     all_documents.append([])
214                     tokens = tokenizer.tokenize(line)
215                     if tokens:
216                         all_documents[-1].append(tokens)
217
218             # Remove empty documents
219             all_documents = [x for x in all_documents if x]
220             rng.shuffle(all_documents)
221
222             vocab_words = list(tokenizer.vocab.keys())
223             instances = []
224             for _ in range(dupe_factor):
225                 for document_index in range(len(all_documents)):
226                     instances.extend(
227                         create_instances_from_document(
228                             all_documents, document_index, max_seq_length, short_seq_prob,
229                             masked_lm_prob, max_predictions_per_seq, vocab_words, rng))
230
231             rng.shuffle(instances)
232
233     return instances
```



Pre-training

Next Sentence Prediction (NSP)

[CLS], [SEP], [SEP]

short_seq_prob

current_chunk

current_chunk → token

current_chunk → token

```
235 def create_instances_from_document(
236     all_documents, document_index, max_seq_length, short_seq_prob,
237     masked_lm_prob, max_predictions_per_seq, vocab_words, rng):
238     """Creates `TrainingInstance`'s for a single document."""
239     document = all_documents[document_index]
240
241     # Account for [CLS], [SEP], [SEP]
242     max_num_tokens = max_seq_length - 3
243
244     # We *usually* want to fill up the entire sequence since we are padding
245     # to `max_seq_length` anyways, so short sequences are generally wasted
246     # computation. However, we *sometimes*
247     # (i.e., `short_seq_prob == 0.1 == 10% of the time) want to use shorter
248     # sequences to minimize the mismatch between pre-training and fine-tuning.
249     # The `target_seq_length` is just a rough target however, whereas
250     # `max_seq_length` is a hard limit.
251     target_seq_length = max_num_tokens
252     if rng.random() < short_seq_prob: # 如果生成的隨機數小於 short_seq_prob 則產生一個較短的訓練序列
253         target_seq_length = rng.randint(2, max_num_tokens)
254
255     # We DON'T just concatenate all of the tokens from a document into a long
256     # sequence and choose an arbitrary split point because this would make the
257     # next sentence prediction task too easy. Instead, we split the input into
258     # segments "A" and "B" based on the actual "sentences" provided by the user
259     # input.
260     instances = []
261     current_chunk = [] # 產生訓練集的候選集
262     current_length = 0
263     i = 0
264     while i < len(document):
265         current_chunk.append(segment) # 候選集將文本中的段落添加到 list
266         current_length += len(segment) # 目前長度也隨著 segment 增加
267     # 當存入候選集已經達到文檔長度或已經大於目標序列長度，會從候選集中隨機取出一個文檔作為句子1的截止文檔，並且把目前候選集都放到 a 的 token
268     if i == len(document) - 1 or current_length >= target_seq_length:
269         if current_chunk:
270             # 'a_end' is how many segments from 'current_chunk' go into the 'A'
271             # (first) sentence.
272             a_end = 1
273             if len(current_chunk) >= 2:
274                 a_end = rng.randint(1, len(current_chunk) - 1) # 從current_chunk中隨機選出一個文檔做為句子1的截止文檔
275
276             tokens_a = []
277             for j in range(a_end):
278                 tokens_a.extend(current_chunk[j]) # 將截止文檔之前的文檔都加到 tokens_a
279
280             tokens_b = []
```

ELMo, BERT, GPT



- Word Embedding

One-Hot-Encoding

Bag-of-Words (BoW)

- 一詞多義

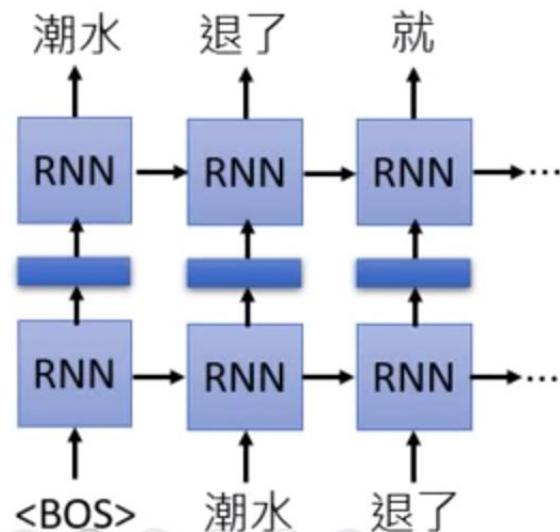
1 type → 1 embedding 1 type → 固定數量的 embedding

每個 token → 不同的embedding

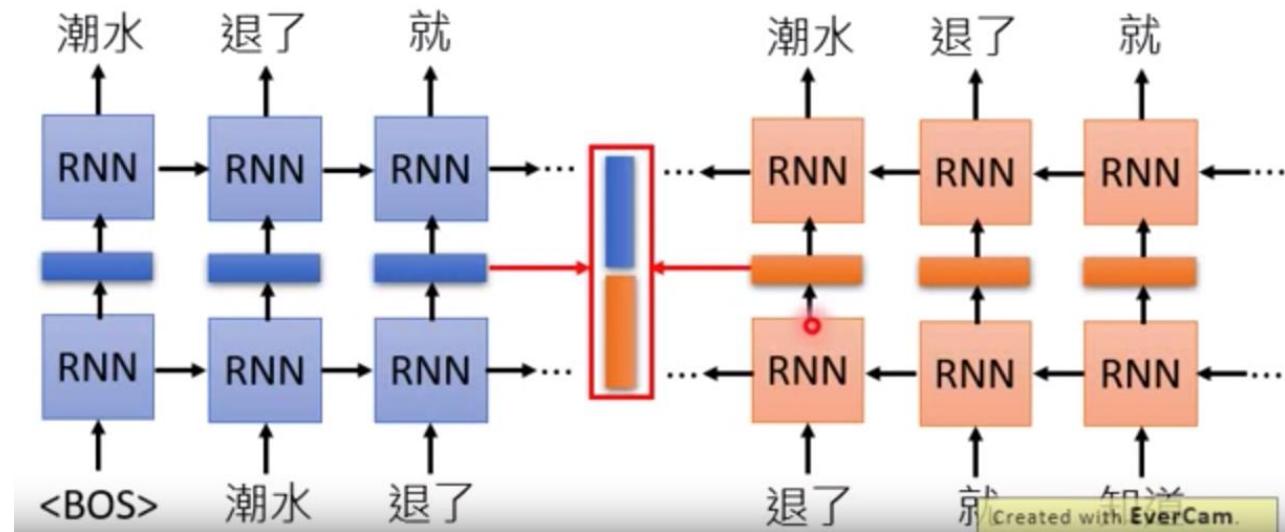
- Contextualized Word Embedding

- Embedding from Language Model (ELMo)

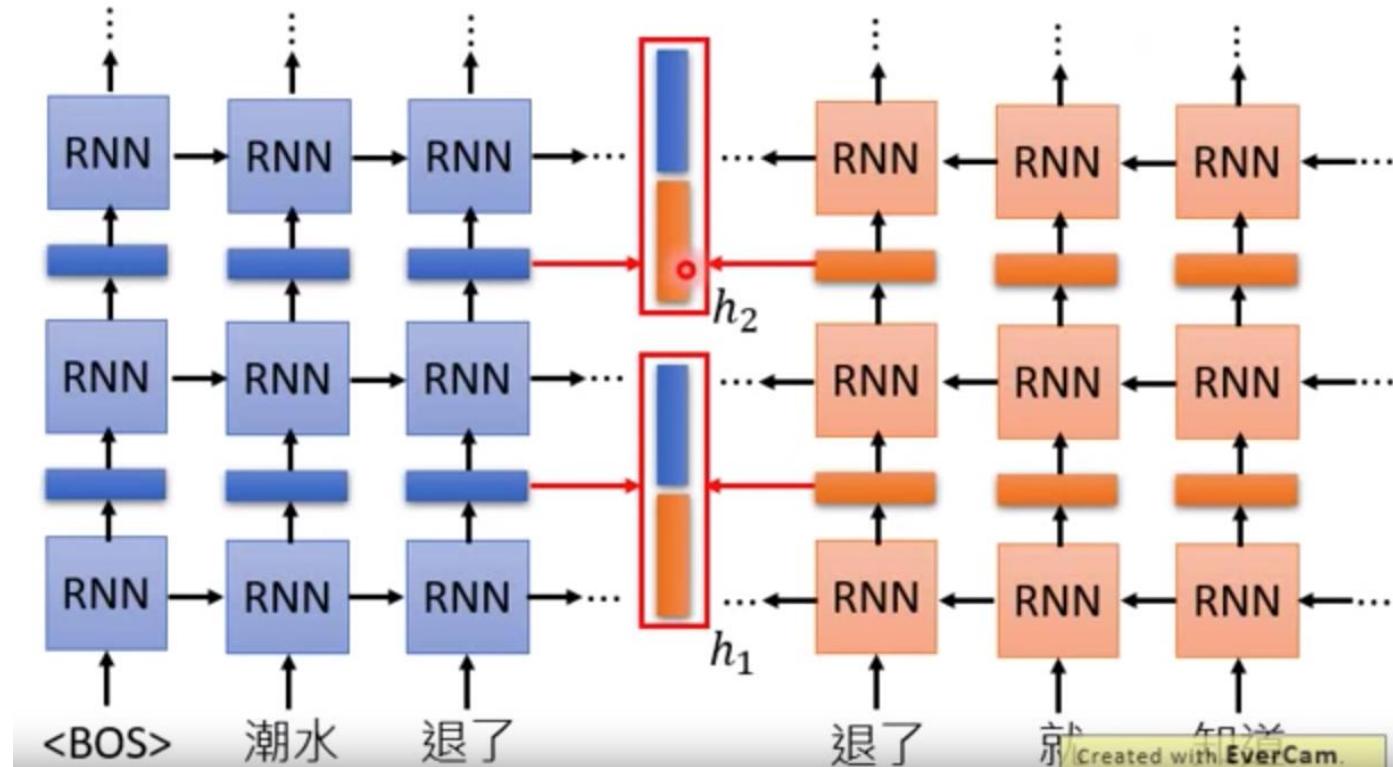
RNN-based Language Model: 蔽集未標註的句子, 再由RNN訓練



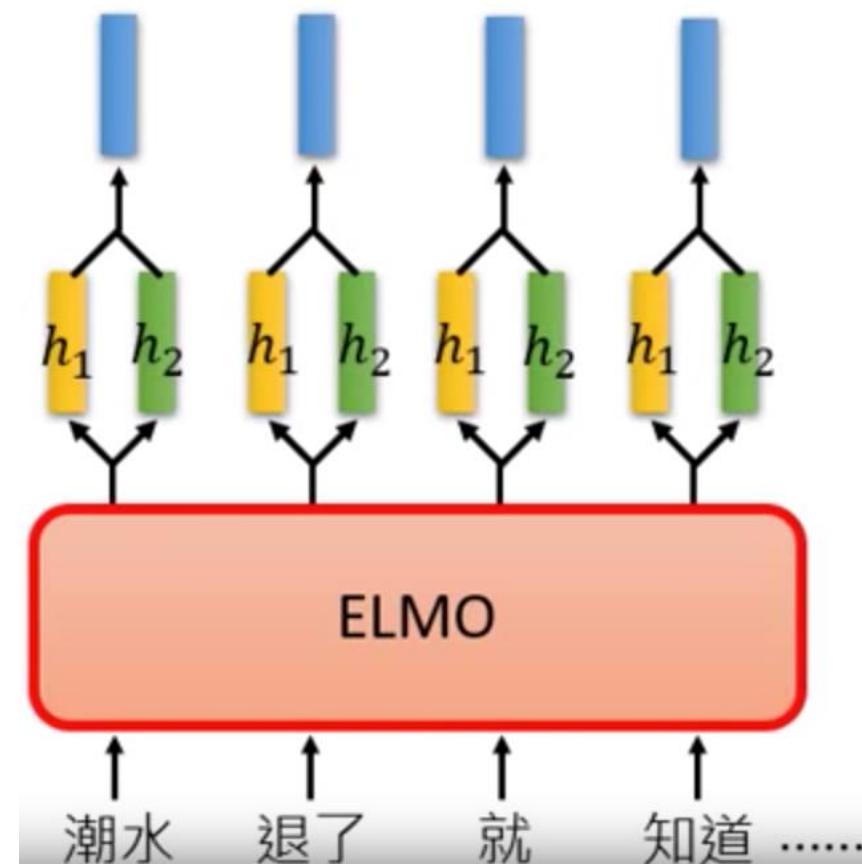
- 只考慮到左側上文



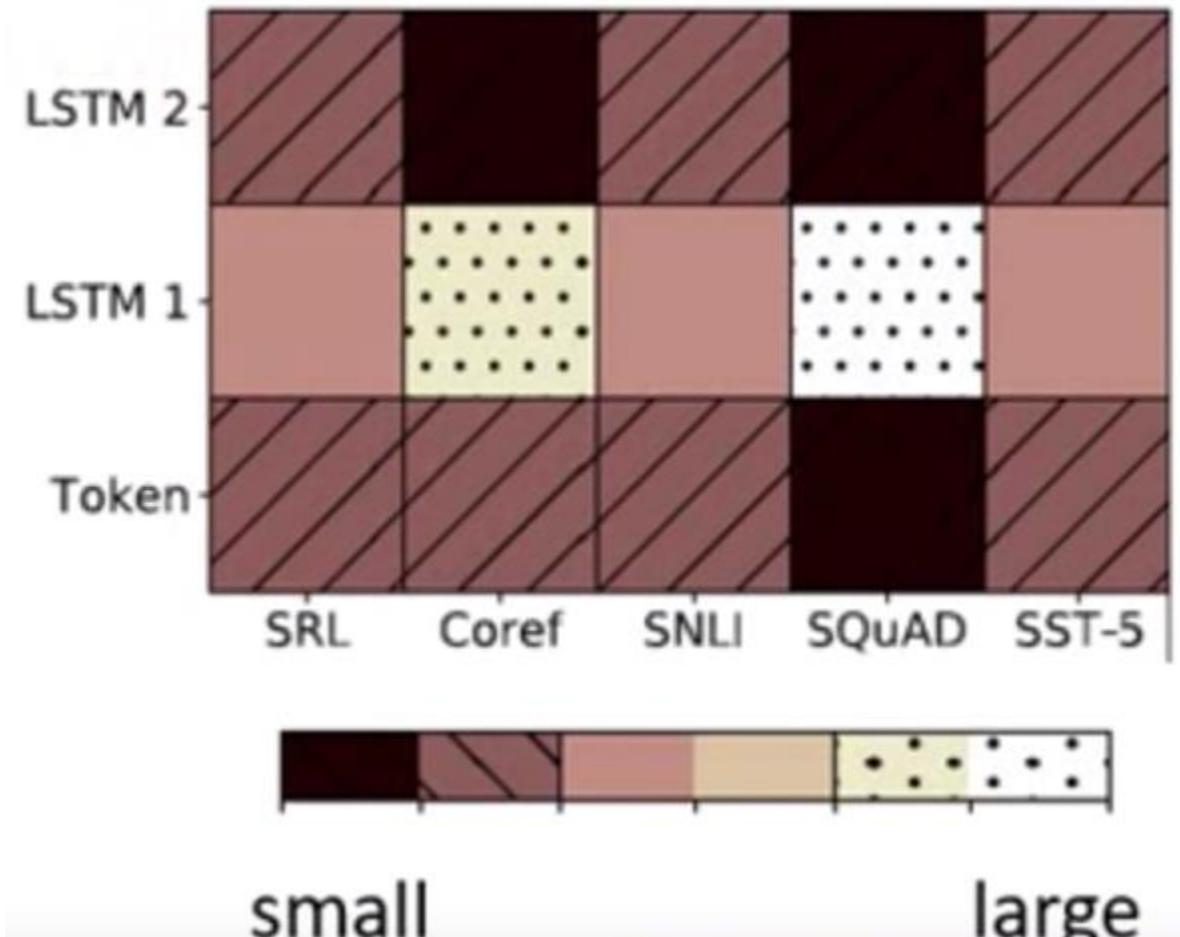
- Embedding from Language Model (ELMo)



- Embedding from Language Model (ELMo)



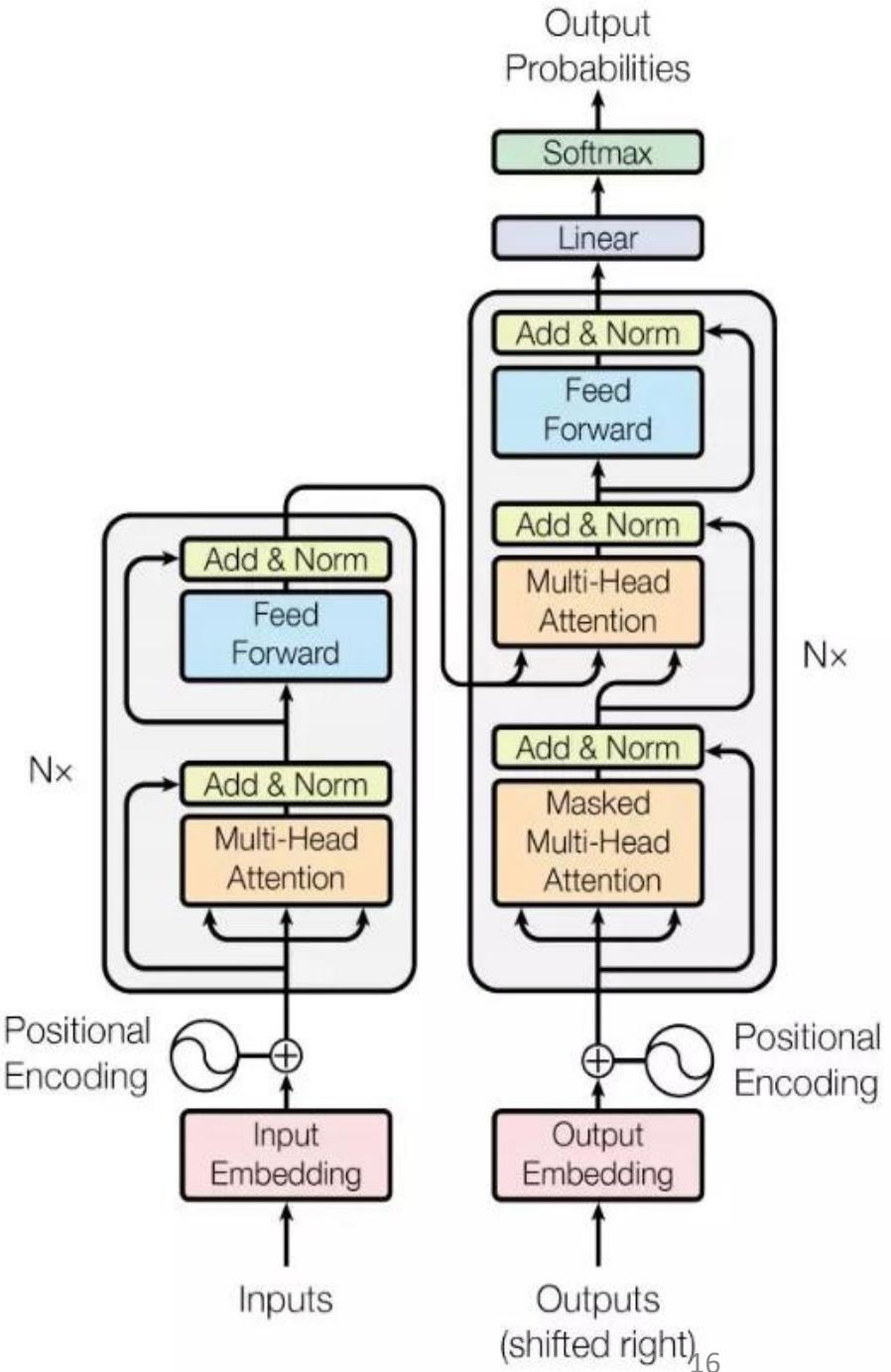
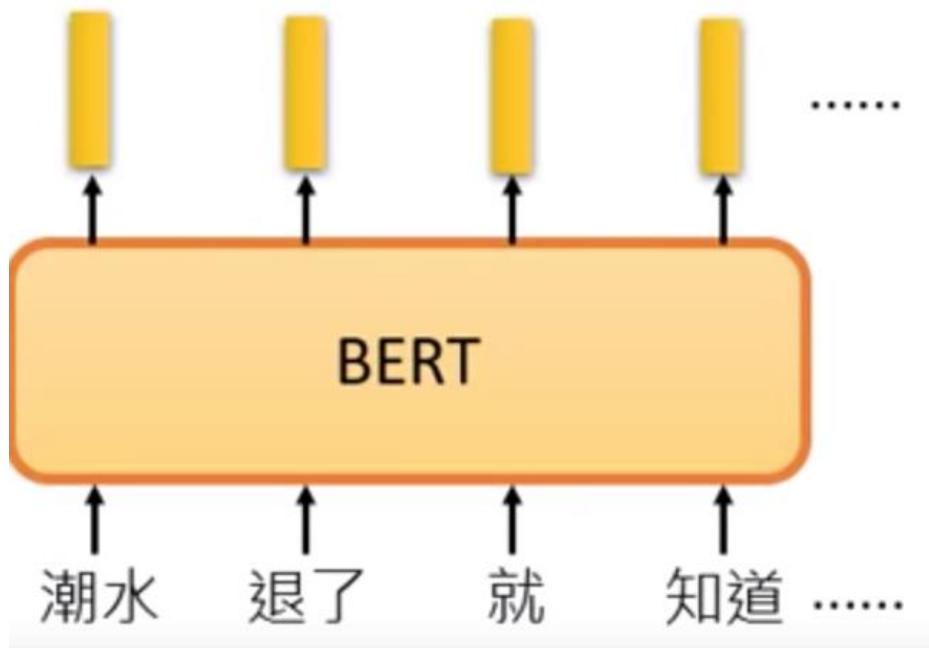
- Embedding from Language Model (ELMo)



- Bidirectional Encoder Representation from Transformers (BERT)

BERT = Encoder of Transformer

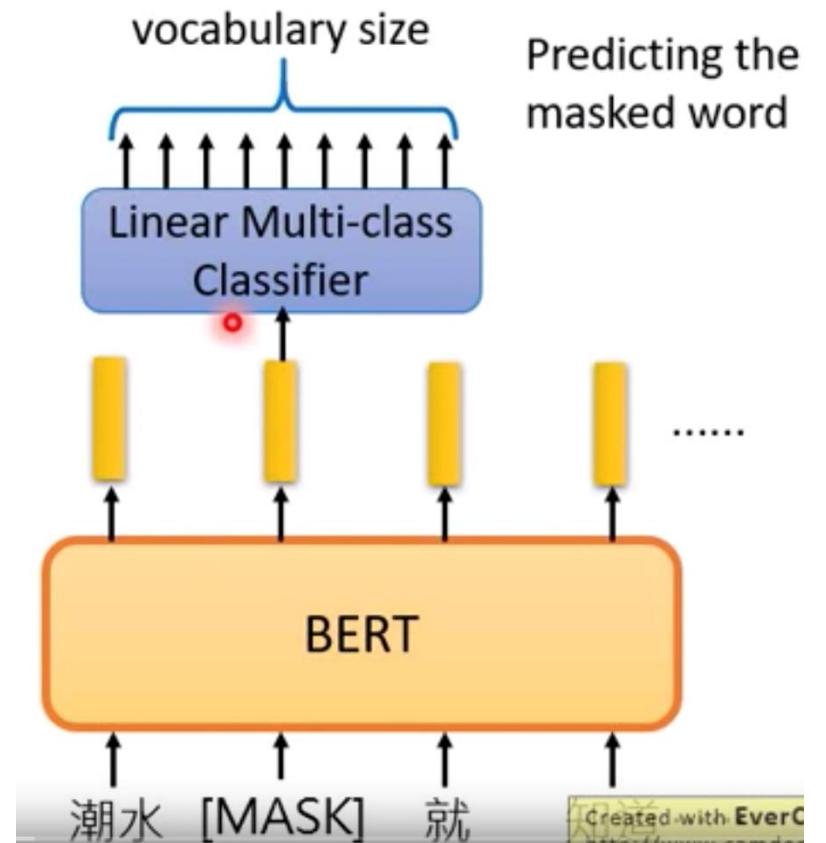
中文輸入問題: 詞 v.s.字



- Bidirectional Encoder Representation from Transformers (BERT)

- Pre-training

1. Masked Language Model (MLM)



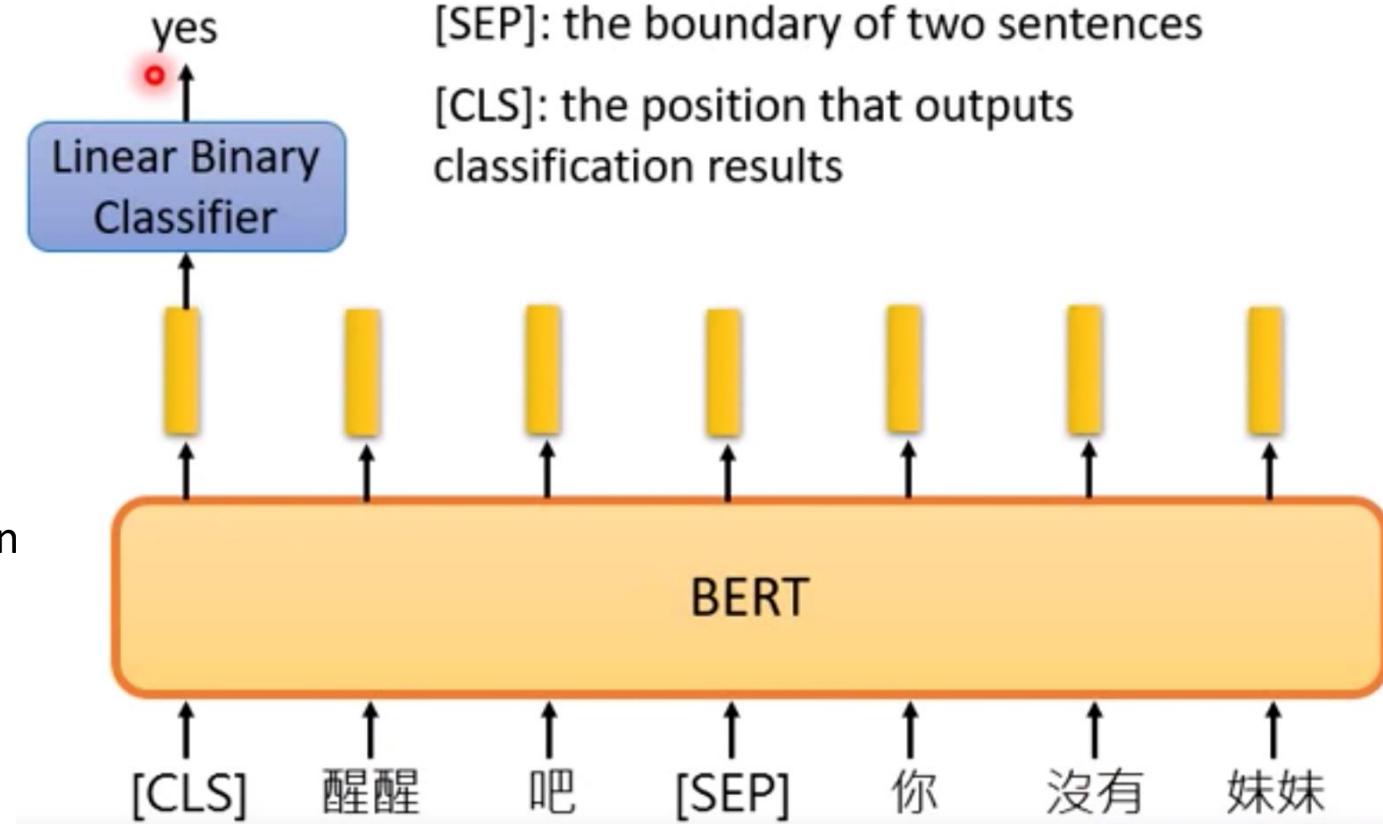
- Bidirectional Encoder Representation from Transformers (BERT)

- Pre-training

2. Next Sentence Prediction (NSP)

[SEP]: 句子間的分界

[CLS]: 放在句子開頭token, 告訴Encoder要做 Classification



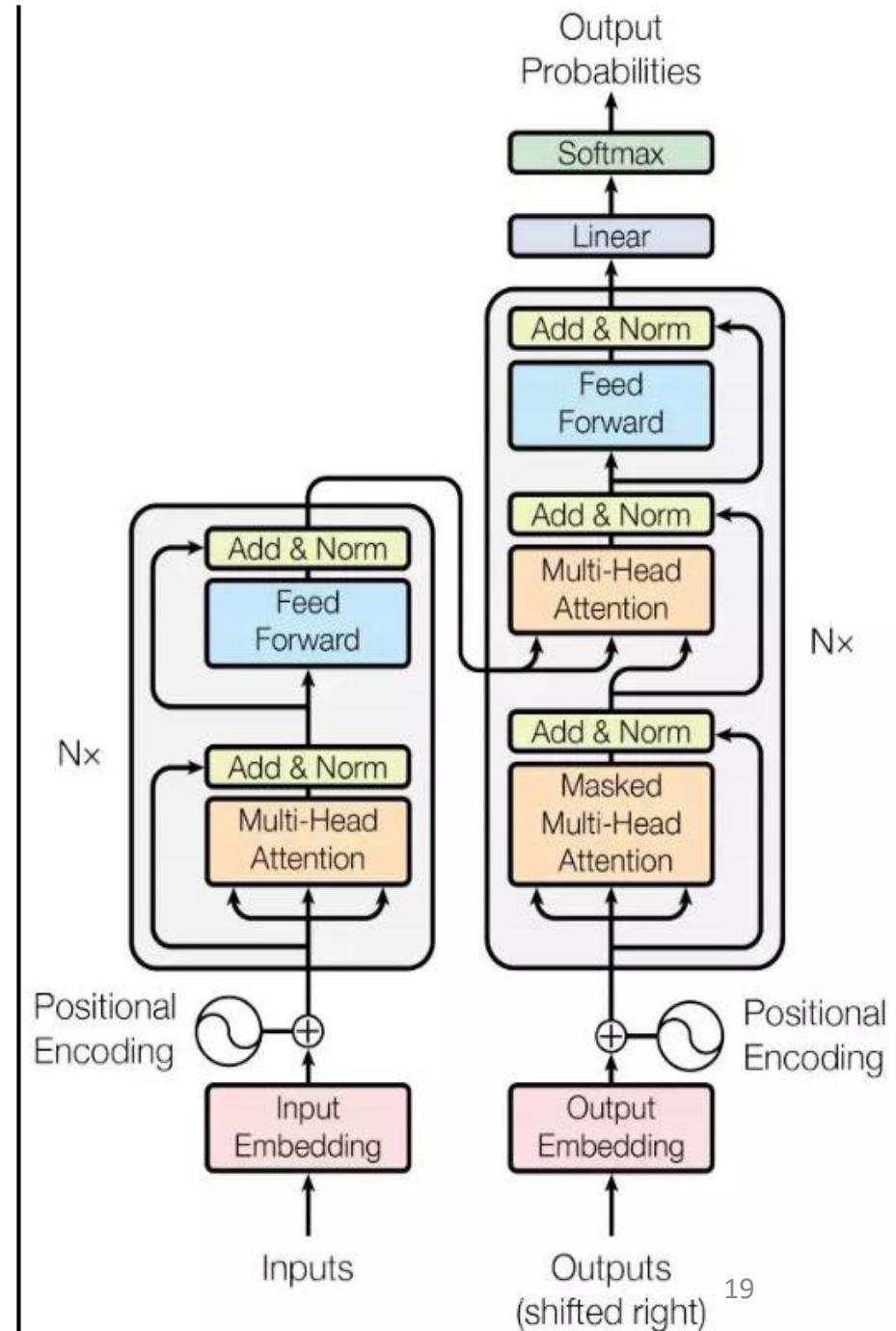
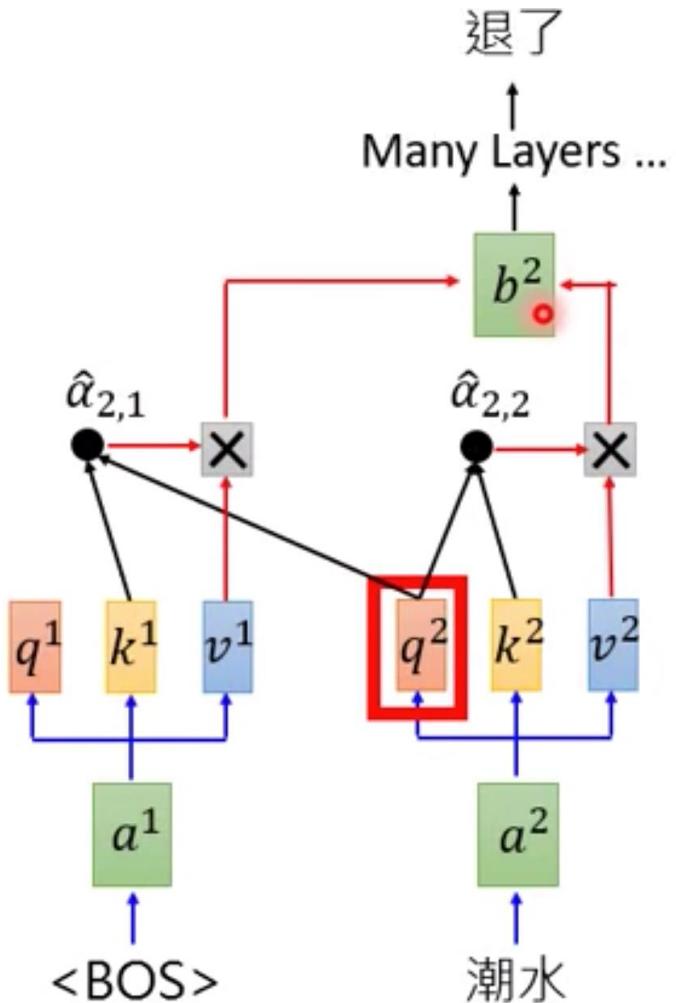
[SEP]: the boundary of two sentences

[CLS]: the position that outputs classification results

- Generative Pre-training (GPT) from OpenAI

GPT = Decoder of Transformer

- Predicting new token: **self-attention**



- Generative Pre-training (GPT) from OpenAI

輸入的資料未經過訓練，就能
透過 self-attention 自動生成

- ***Reading Comprehension***

$d_1, d_2, \dots, d_N, "Q:"$, $q_1, q_2, \dots, q_M, "A:"$

- ***Summarization*** $d_1, d_2, \dots, d_N, "TL;DR:"$

- ***Translation***

English sentence 1 = French sentence 1

English sentence 2 = French sentence 2

Talk to Transformer

- Generative Pre-training (GPT) from OpenAI

See how a modern neural network completes your text. Type a custom snippet or try one of the examples. [Learn more](#) below.



Follow @AdamDanielKing for more neat neural networks.

Text generated is temporarily shorter than before.

Custom prompt

Q: What is machine learning?

A:

GENERATE ANOTHER

Completion

Q: What is machine learning?

A: Machine learning (ML) is a computer science model used to perform statistical tasks in data science. ML is not a new concept, however its methodologies have been gaining in popularity.

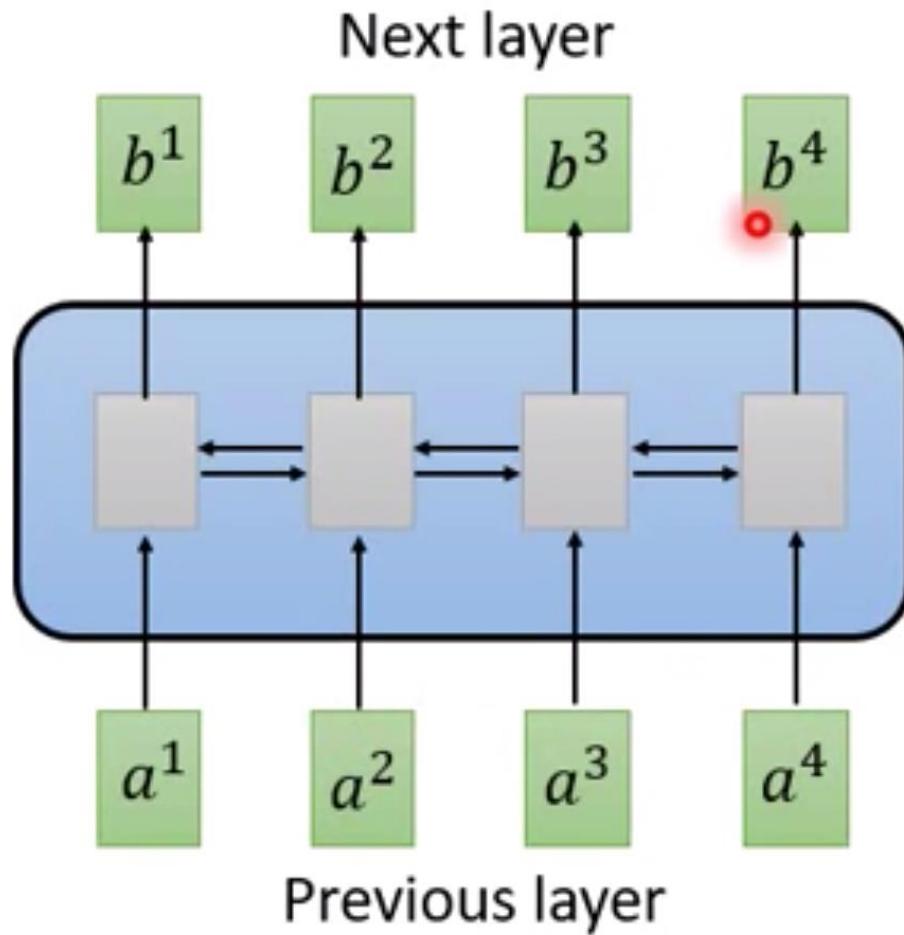
Q: How does ML help me make better data science decisions?

A: ML can help you sort through your dataset and make better decisions based on these insights. One can use machine learning to assist with algorithmic decision making, this includes finding patterns, correlations, and concepts within the data to potentially aid in the generation of meaningful insights. For example,

Transformer, Self-Attention

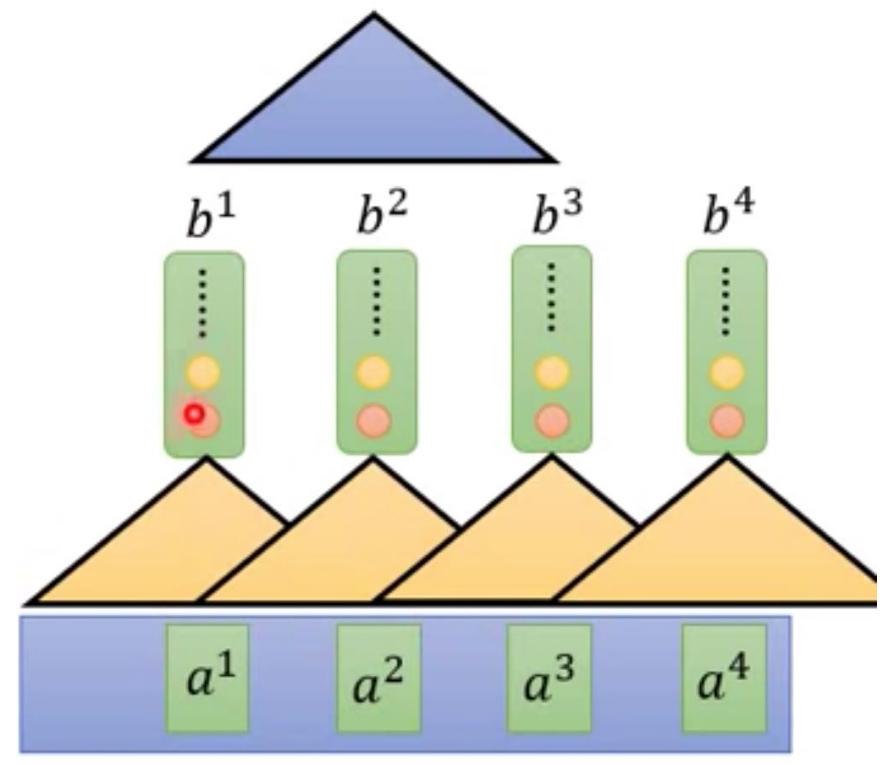


RNN



CNN

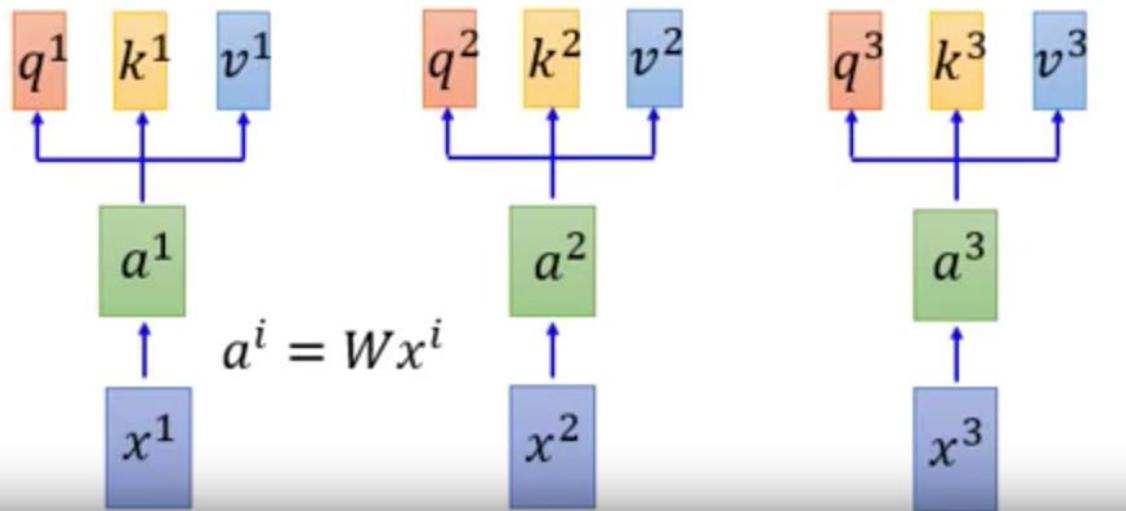
Filters in higher layer can consider longer sequence



Using CNN to replace RNN

Self-attention

Word-Embedding



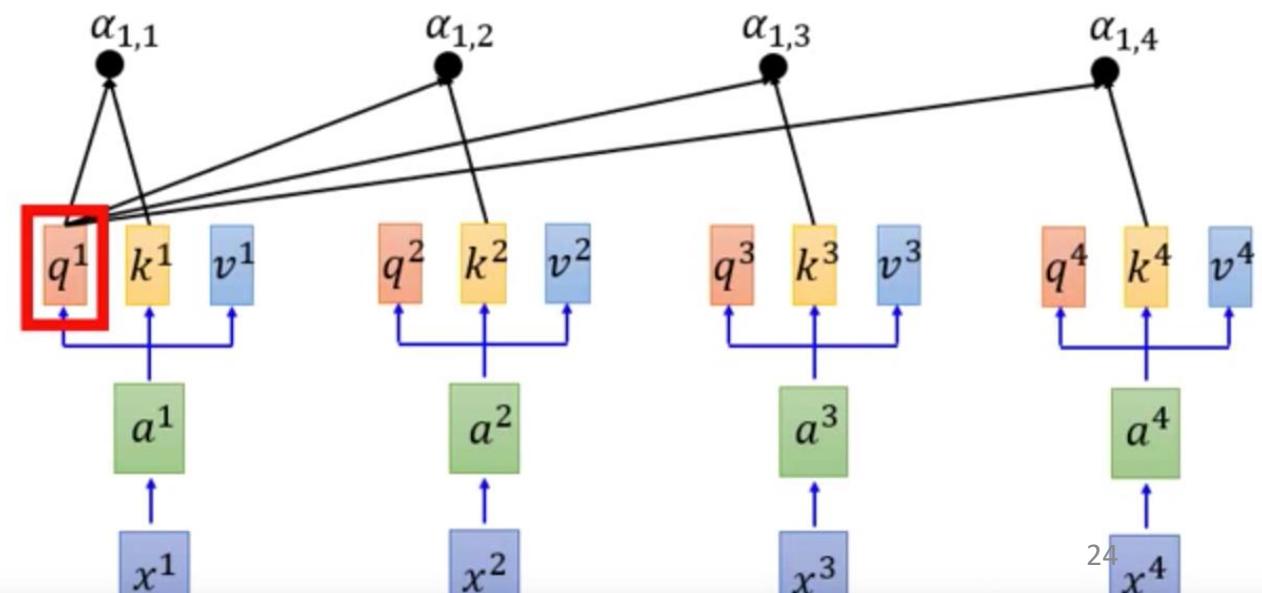
Attention

$$q^i = W^q a^i$$

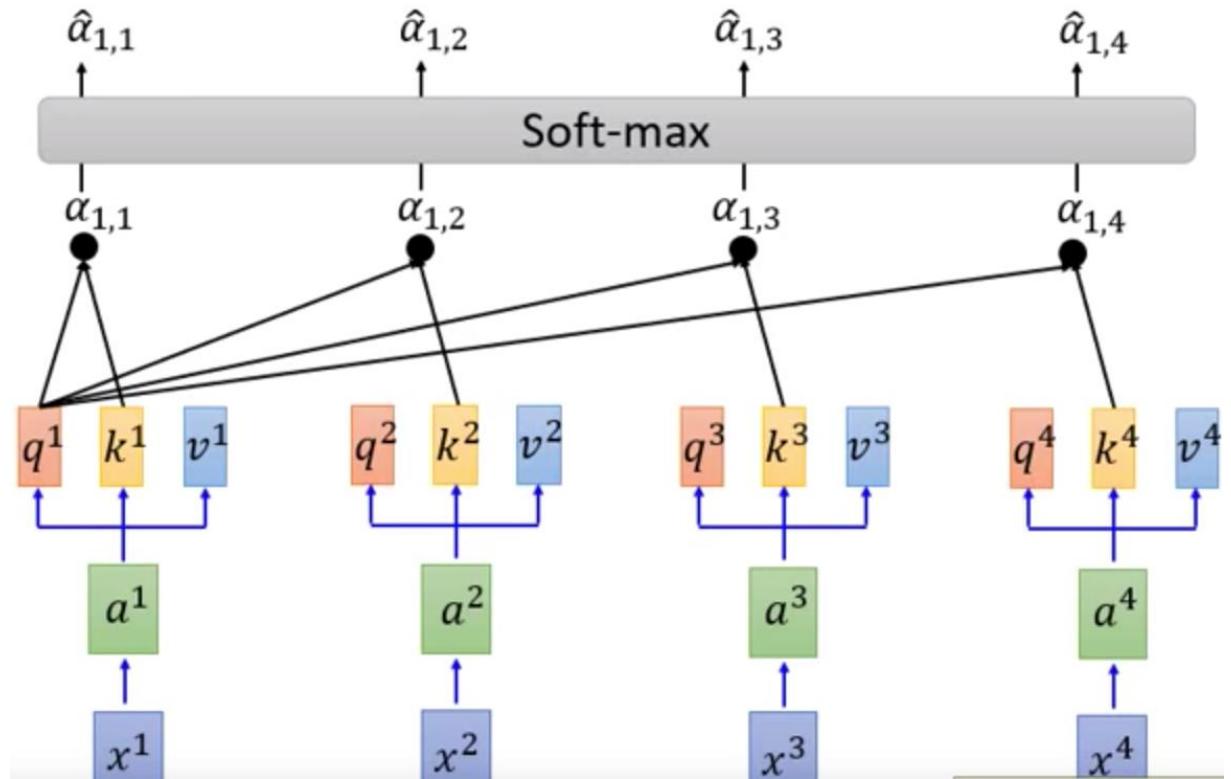
$$k^i = W^k a^i$$

$$v^i = W^v a^i$$

Scaled Dot-Product Attention: $\alpha_{1,i} = \underbrace{q^1 \cdot k^i}_{\text{dot product}} / \sqrt{d}$

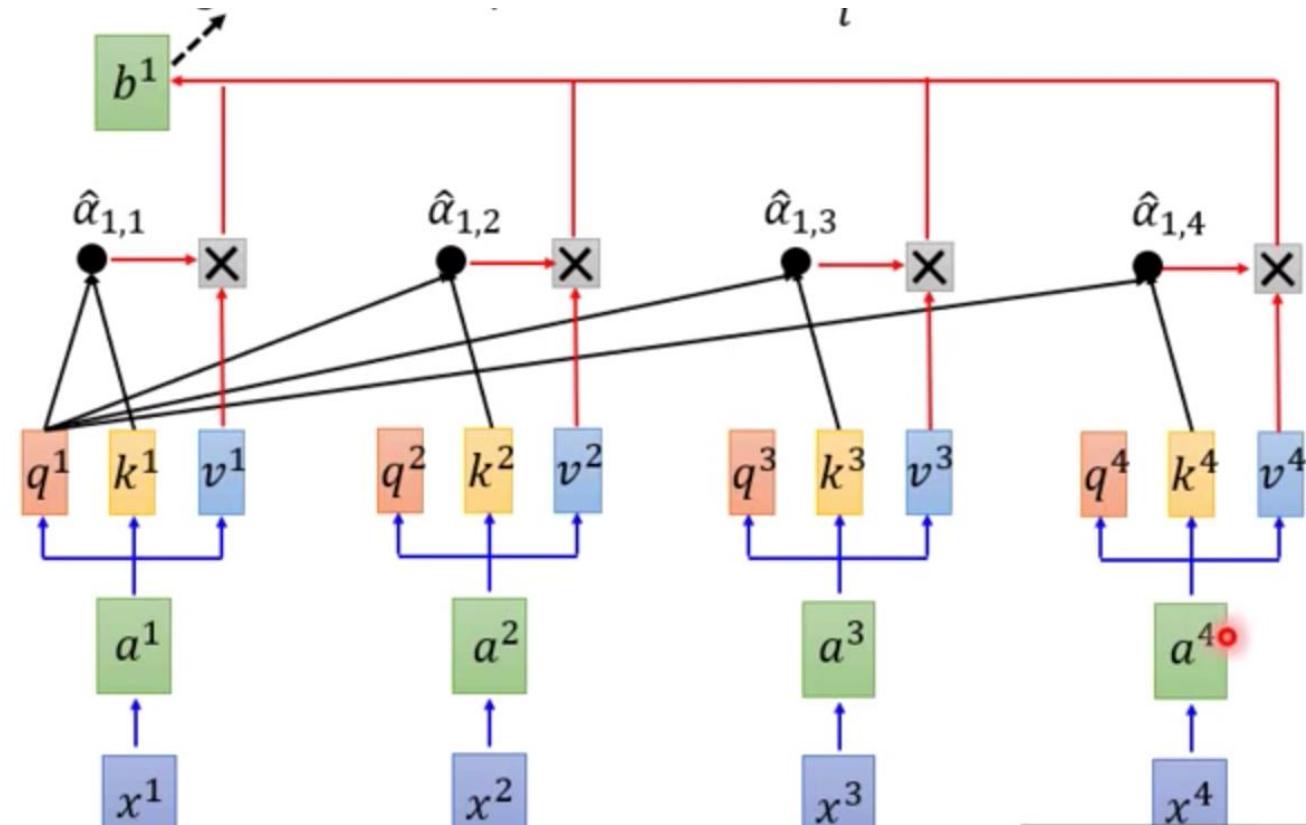


Self-attention



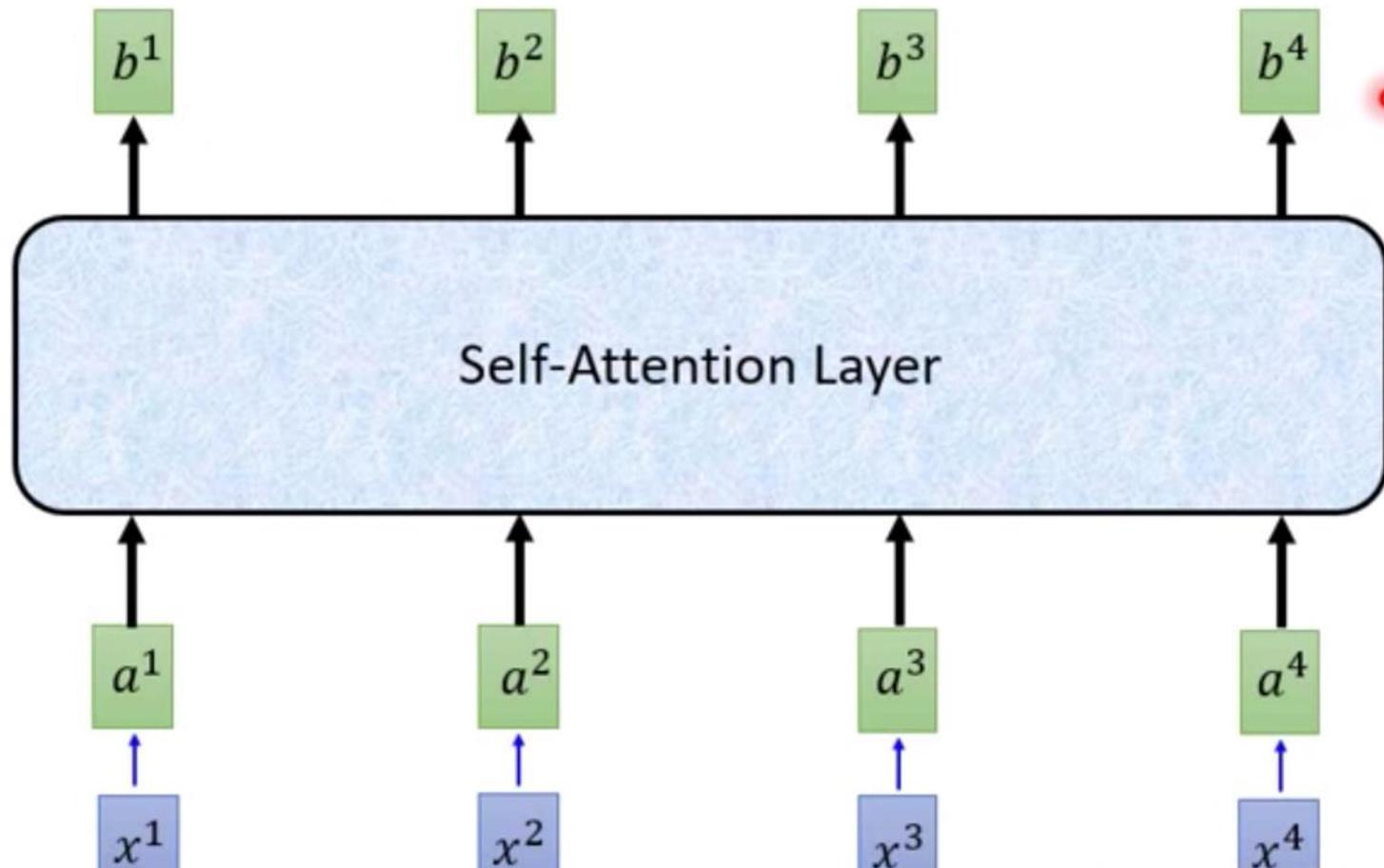
$$\hat{\alpha}_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$

Self-attention

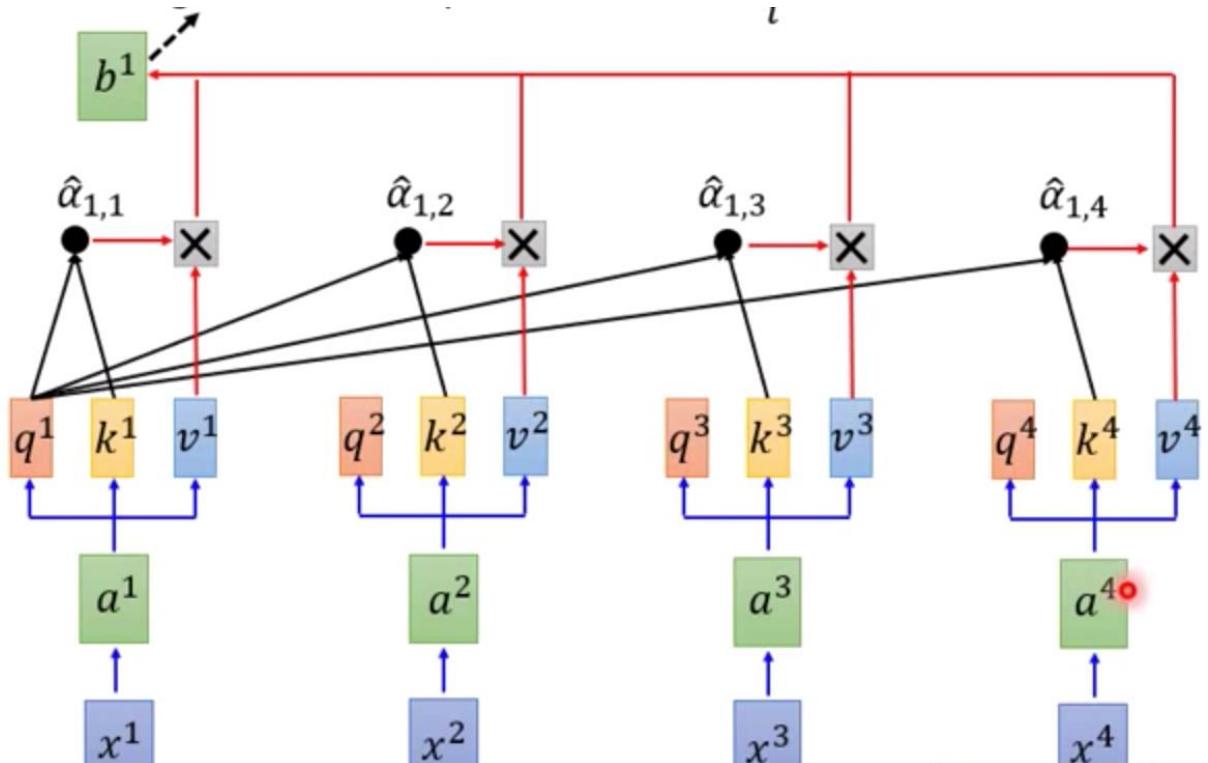


$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$

Self-attention



Self attention with Parallels Computing



$$\begin{aligned} \alpha_{1,1} &= \begin{matrix} k^1 \\ q^1 \end{matrix} & \alpha_{1,2} &= \begin{matrix} k^2 \\ q^1 \end{matrix} \\ \alpha_{1,3} &= \begin{matrix} k^3 \\ q^1 \end{matrix} & \alpha_{1,4} &= \begin{matrix} k^4 \\ q^1 \end{matrix} \end{aligned}$$

(ignore \sqrt{d} for simplicity)

$$\begin{matrix} \hat{\alpha}_{1,1} & \hat{\alpha}_{2,1} & \hat{\alpha}_{3,1} & \hat{\alpha}_{4,1} \\ \hat{\alpha}_{1,2} & \hat{\alpha}_{2,2} & \hat{\alpha}_{3,2} & \hat{\alpha}_{4,2} \\ \hat{\alpha}_{1,3} & \hat{\alpha}_{2,3} & \hat{\alpha}_{3,3} & \hat{\alpha}_{4,3} \\ \hat{\alpha}_{1,4} & \hat{\alpha}_{2,4} & \hat{\alpha}_{3,4} & \hat{\alpha}_{4,4} \end{matrix} \quad \leftarrow \quad \begin{matrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{matrix} \quad A$$

Created with EverCam. http://www.camdemmy.com

$$b^1 \boxed{b^2} \boxed{b^3} \boxed{b^4} = v^1 \boxed{v^2} \boxed{v^3} \boxed{v^4}$$

O V

$$\begin{matrix} \hat{\alpha}_{1,1} & \hat{\alpha}_{2,1} & \hat{\alpha}_{3,1} & \hat{\alpha}_{4,1} \\ \hat{\alpha}_{1,2} & \hat{\alpha}_{2,2} & \hat{\alpha}_{3,2} & \hat{\alpha}_{4,2} \\ \hat{\alpha}_{1,3} & \hat{\alpha}_{2,3} & \hat{\alpha}_{3,3} & \hat{\alpha}_{4,3} \\ \hat{\alpha}_{1,4} & \hat{\alpha}_{2,4} & \hat{\alpha}_{3,4} & \hat{\alpha}_{4,4} \end{matrix}$$

Created with EverCam. http://www.camdemmy.com

Self-attention and Padding Mask

```
def attention(query, key, value, mask=None, dropout=None):
    '''Compute Scale Dot Product Attention'''
    d_k = query.size(-1) # size of query (token), default=64
    scores = torch.matmul(query, key.transpose(-2, -1) / math.sqrt(d_k))
    if mask is not None: # Padding mask, as a tensor for padding for softmax to calculate
        scores = scores.masked_fill(mask == 0, -1e9)
    p_attn = F.softmax(scores, dim = -1)
    if dropout is not None:
        p_attn = dropout(p_attn)
    return torch.matmul(p_attn, value), p_attn
```

Position Encoding

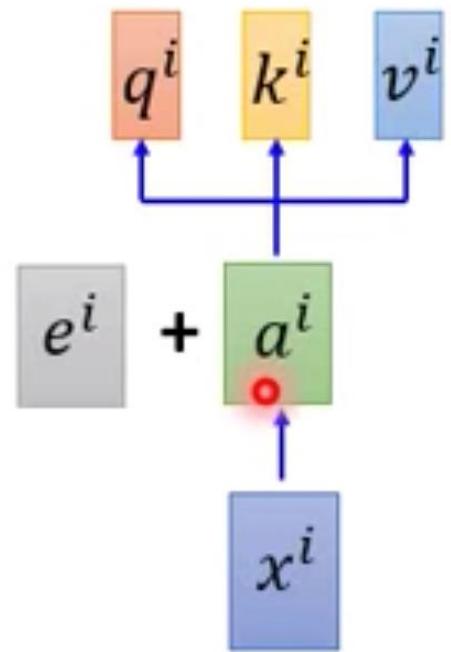
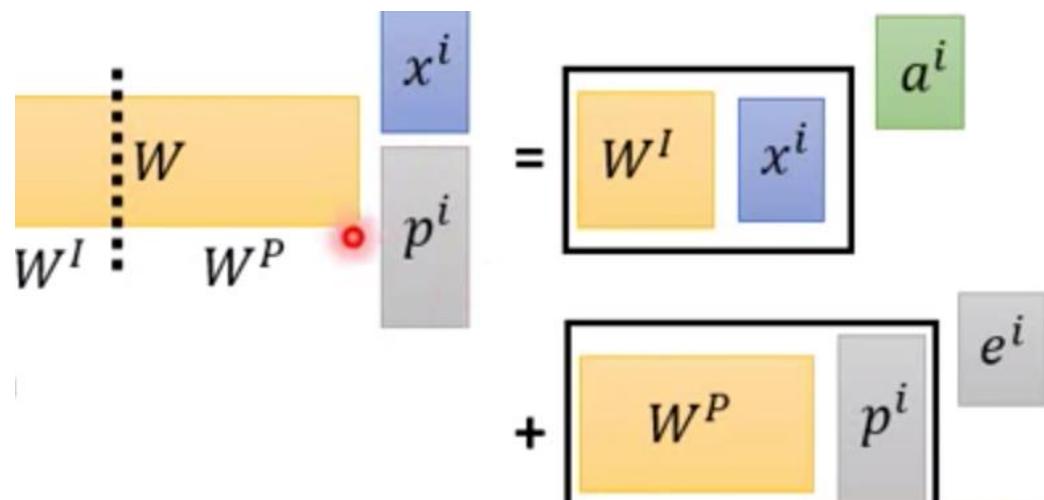
1. No position information in self-attention

a打了b = b 打了a

2. Each position has a unique position vector e^i
(not learned from data)

3. Each x^i appends a one-hot vector p^i

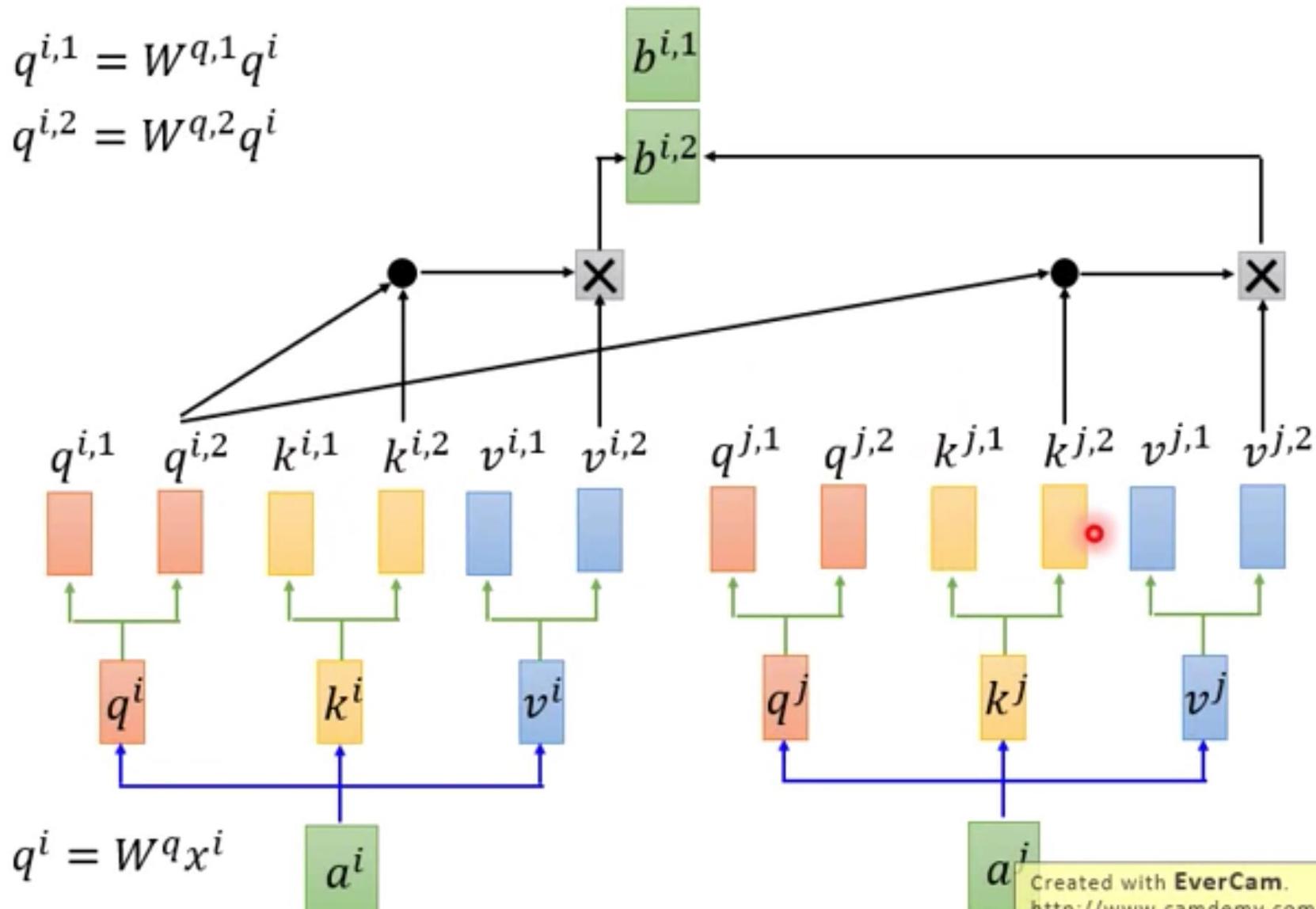
$$p^i = \begin{matrix} \dots \\ 0 \\ 1 \\ 0 \\ \vdots \end{matrix} \quad \leftarrow \text{i-th dim}$$



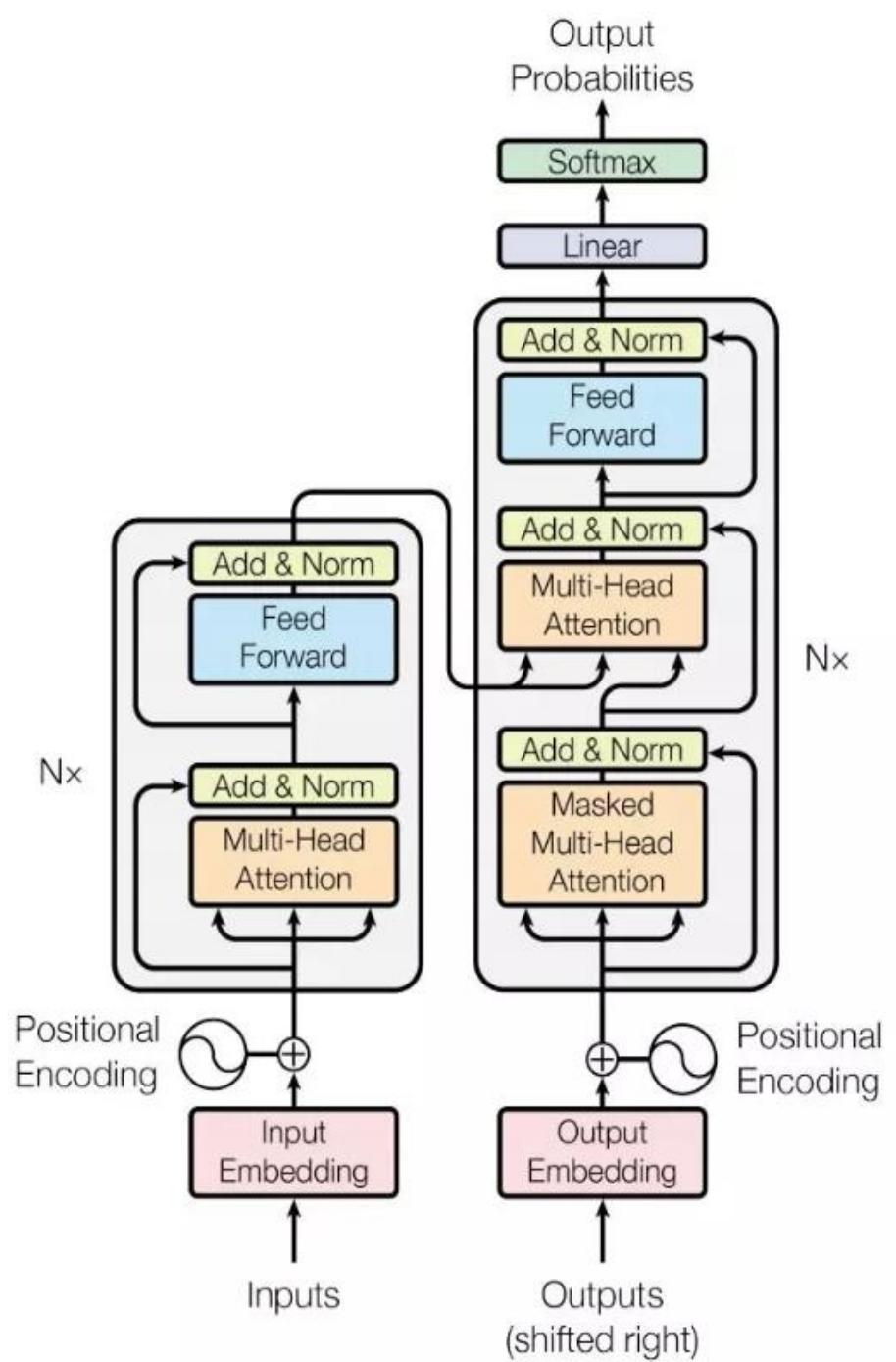
Multi-head Self-attention

$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$



Architecture of Transformer



Self-Attention Visualization

The
animal
didn't
cross
the
street
because
it
was
too
tired
.

The
animal
didn't
cross
the
street
because
it
was
too
tired
.

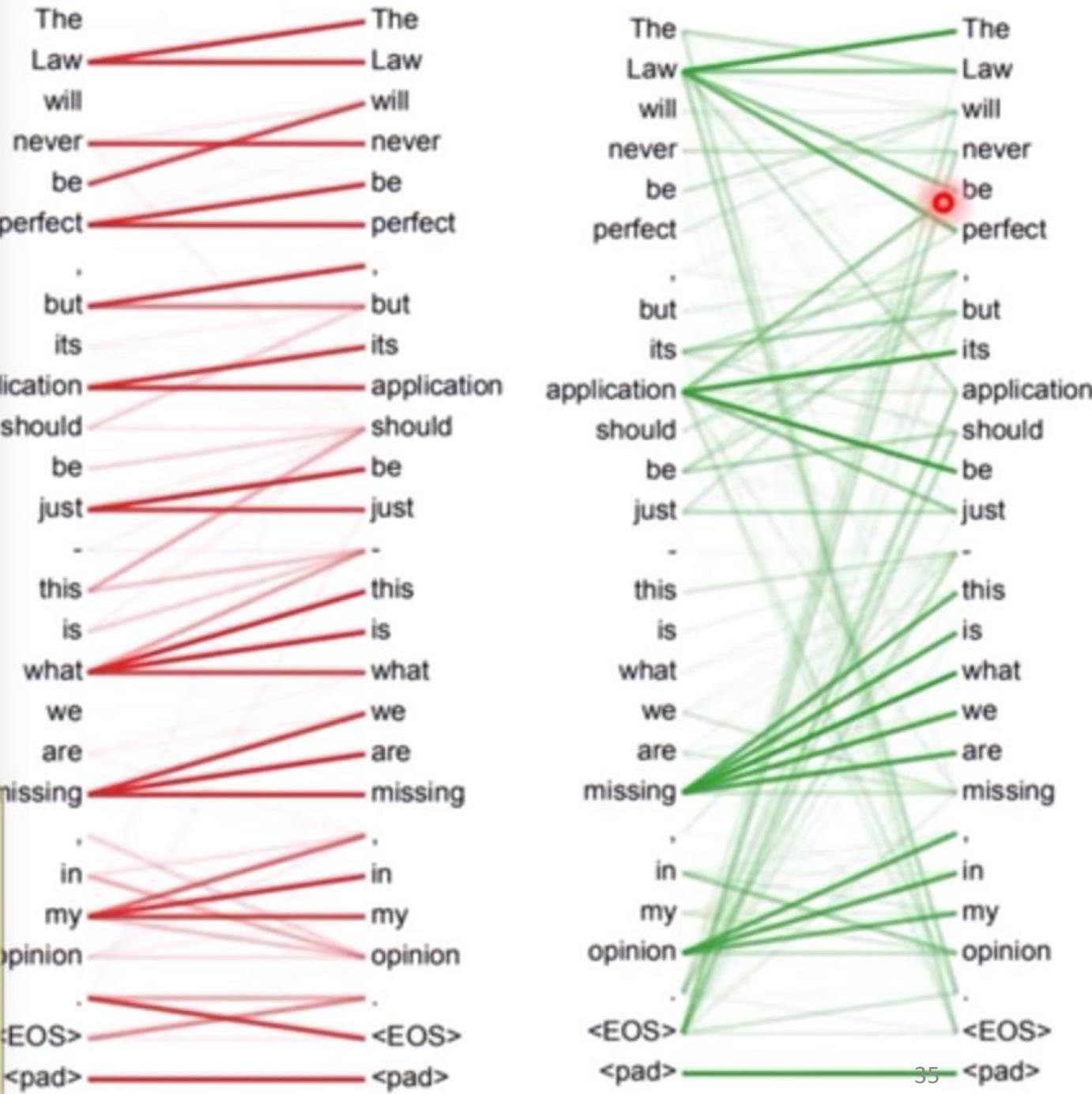
The
animal
didn't
cross
the
street
because
it
was
too
wide
.

Self-Attention Visualization

```
model_type = 'bert'  
bert_version = 'bert-base-chinese'  
  
bertviz_model = BertModel.from_pretrained(bert_version)  
bertviz_tokenizer = BertTokenizer.from_pretrained(bert_version)  
  
# 情境 1  
sentence_a = "胖虎叫大雄去買漫畫，"  
sentence_b = "回來慢了就打他。"  
call_html()  
show(bertviz_model, model_type, bertviz_tokenizer, sentence_a, sentence_b)
```



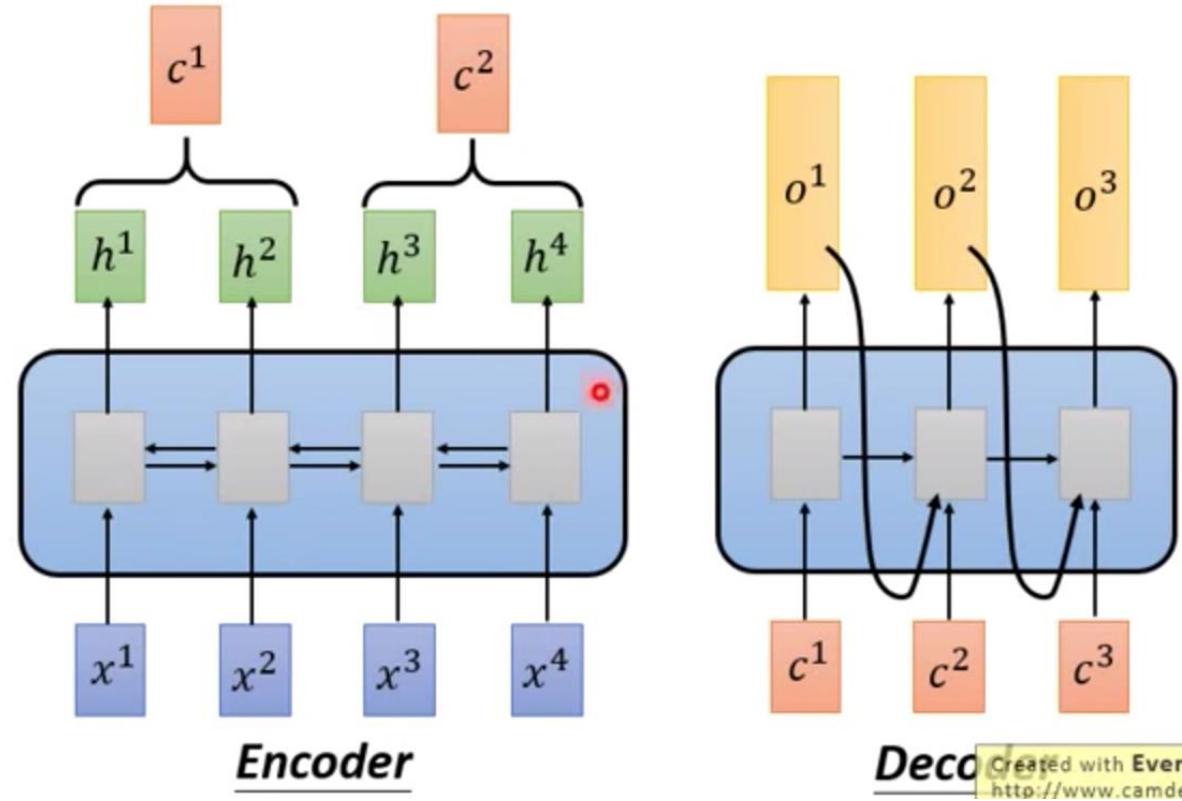
Multi-Head Attention Visualization



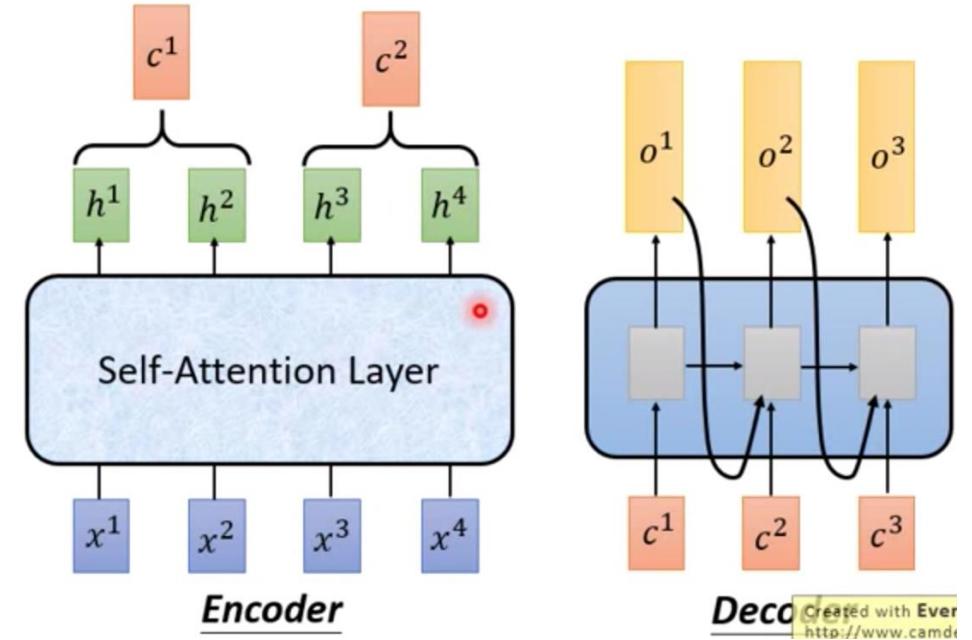
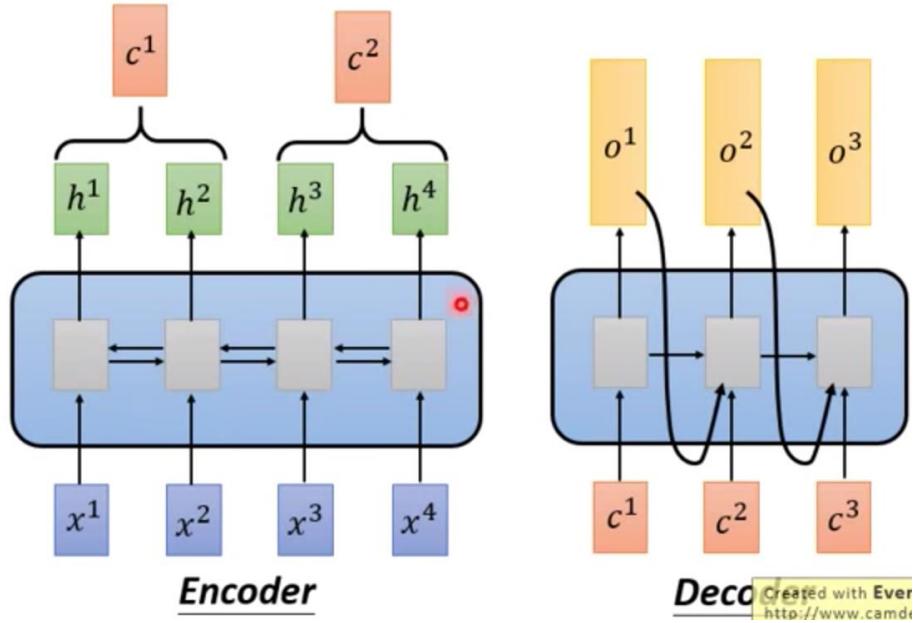
Seq2Seq with Attention



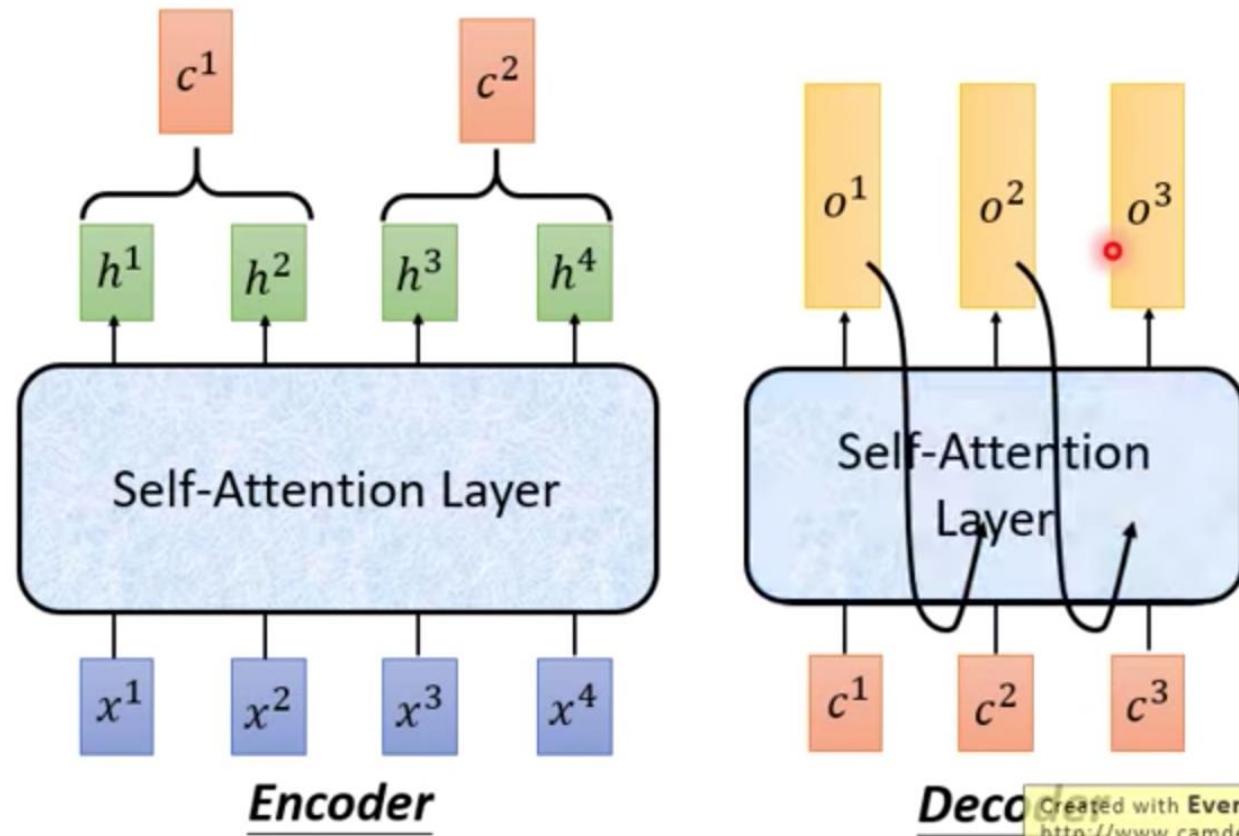
Seq2Seq



Seq2Seq with Attention



Seq2Seq with Attention



Seq2Seq with Attention

Visualization

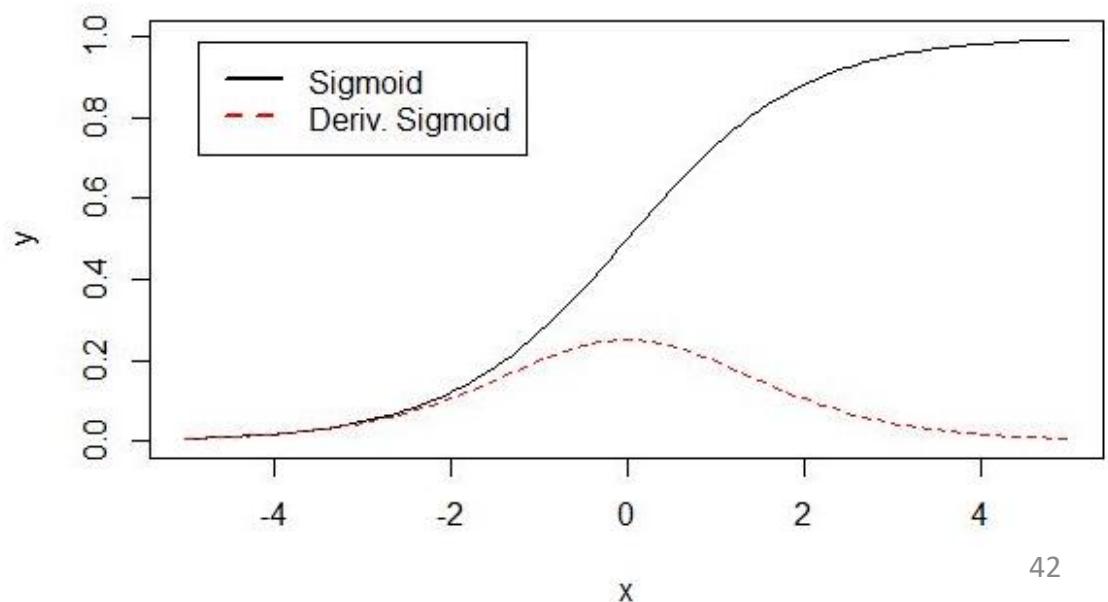
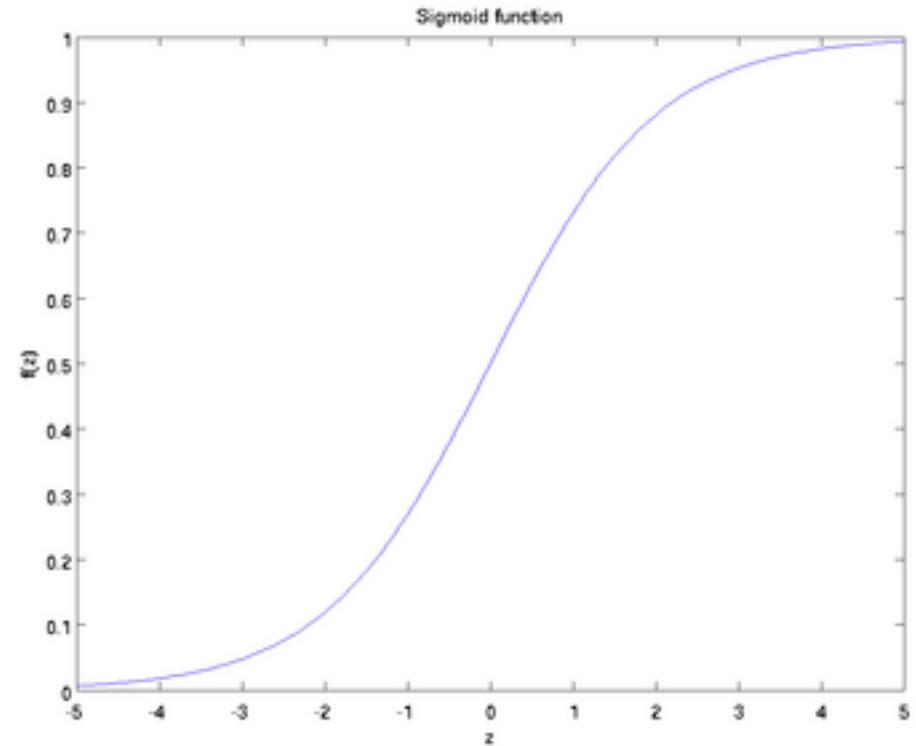
Activation Function

Activation Function

1. Sigmoid

$$f(z) = \frac{1}{1 + \exp(-z)}.$$

$$\phi'(x) = \phi(x)(1 - \phi(x))$$

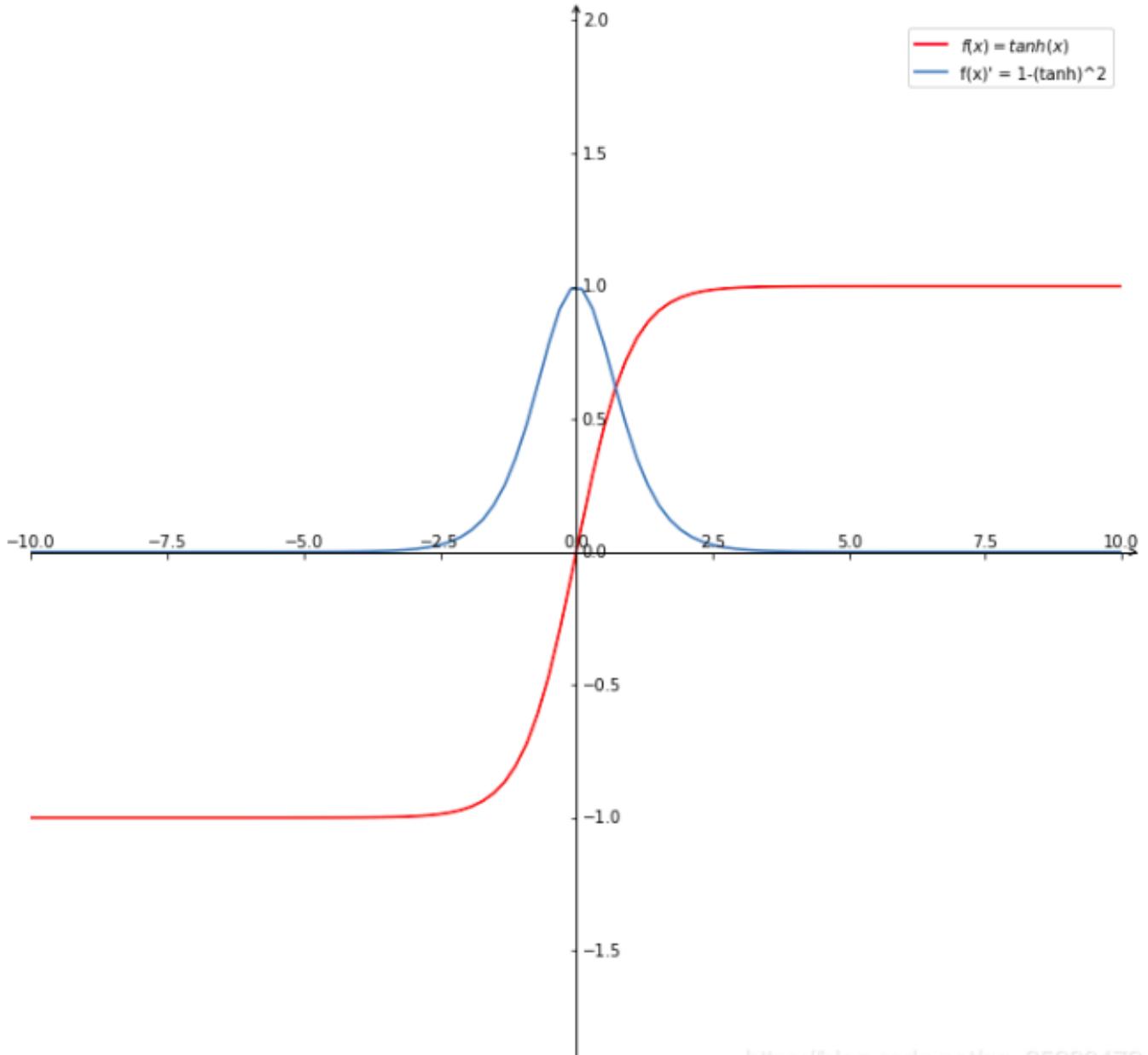


Activation Function

2. Hyperbolic Tangent Function

$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}},$$

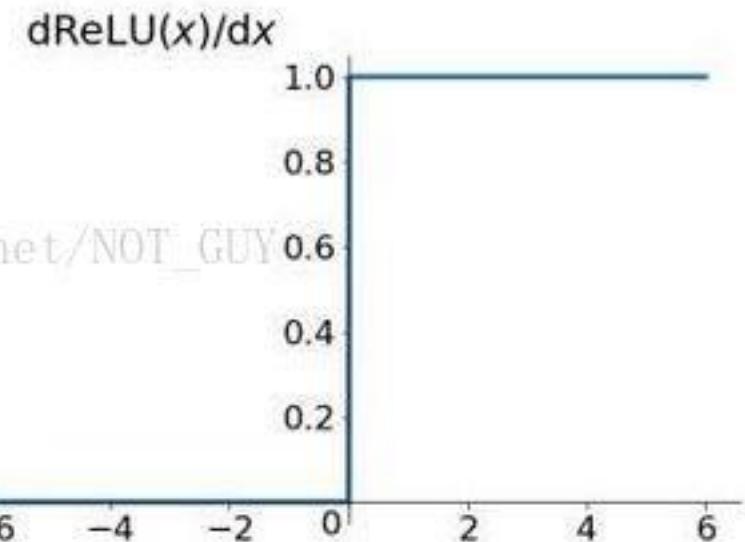
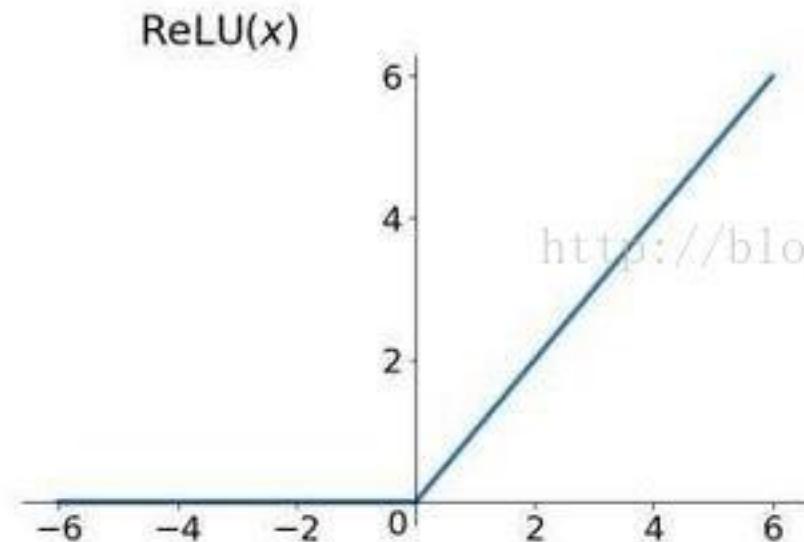
$$\tanh(x) = 2\text{sigmoid}(2x) - 1$$



Activation Function

3. ReLU (Rectified Linear Unit)

$$\text{ReLU} = \max(0, x)$$



- Dead ReLU Problem

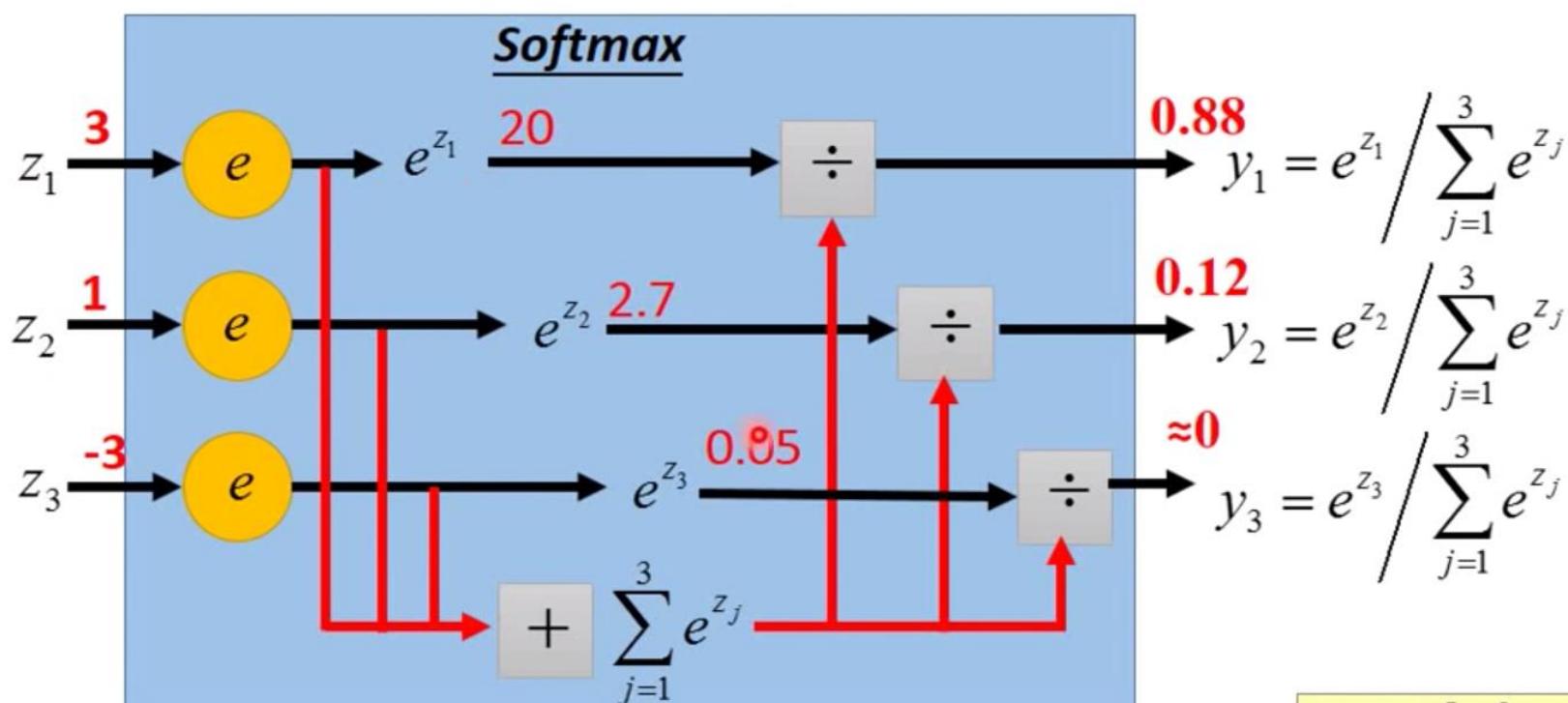
Activation Function

4. Softmax

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

$$\begin{array}{ll} C_1: w^1, b_1 & z_1 = w^1 \cdot x + b_1 \\ C_2: w^2, b_2 & z_2 = w^2 \cdot x + b_2 \\ C_3: w^3, b_3 & z_3 = w^3 \cdot x + b_3 \end{array}$$

- Probability:**
- $1 > y_i > 0$
 - $\sum_i y_i = 1$
 - $\text{Softmax}(x) = \text{Softmax}(x+c)$

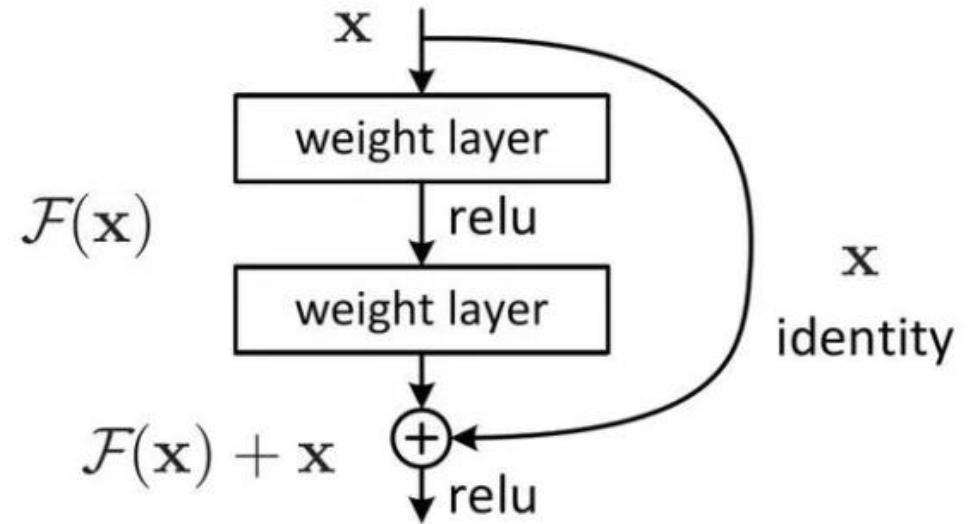


Residual Connection



Residual-Connection (from ResNet)

$$G(x) = F(x) + x$$



WordPiece



Byte-Pair Encoding (BPE)

Word -> Pieces

"loved", "loving", "loves" → "lov", "ed", "ing", "es"

Source Code

(Data Preprocessing)

num_train_steps (epochs)

warm_up_steps

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py - L848-L856

file_based_convert_example_to_features

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py - L484-L513

convert_single_example (tokens for different sequence)

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py - L377-L398

tokenize

<https://github.com/the-Quert/iNLPfun/blob/master/BERT/tokenization.py - L170-L178>

convert_single_example (tokens for different sequence)

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py - L399-L407

_truncate_seq_pair

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py - L564-L578

convert_single_example (tokens for different sequence)

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py - L409-L481

[CLS] Sentence_1 [SEP] Sentence_2 [SEP]

[CLS] Sentence_1 [SEP]

create_model

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py - L581-L590

Source Code (Modeling)

__init__

<https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py - L163-L183>

Id-> Word_Embedding + Token_type_id + Position Encoding

Position Encoding

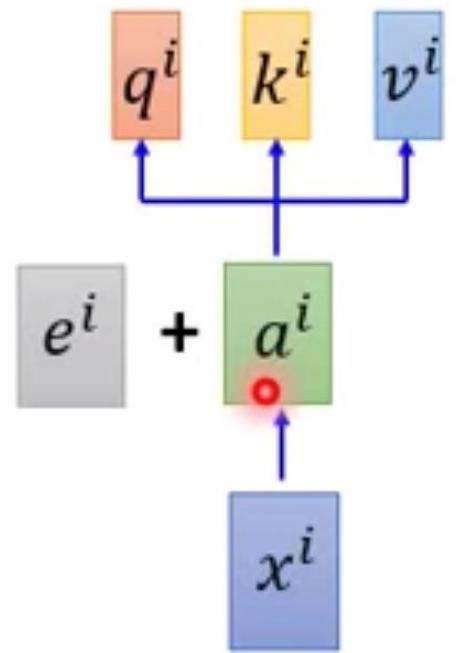
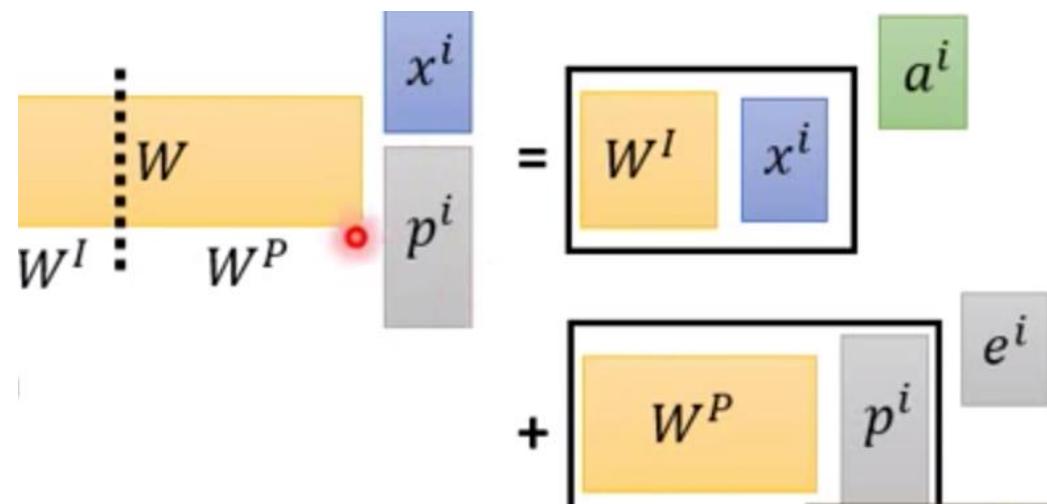
1. No position information in self-attention

a打了b = b 打了a

2. Each position has a unique position vector e^i
(not learned from data)

3. Each x^i appends a one-hot vector p^i

$$p^i = \begin{matrix} \dots \\ 0 \\ 1 \\ 0 \\ \vdots \end{matrix} \quad \text{i-th dim}$$



Word Embedding

embedding_lookup

[https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py/ - L404-L435](https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py#L404-L435)

token_type Embedding

embedding_post_postprocessor

[https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py/ - L480-L500](https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py#L480-L500)

Position Embedding

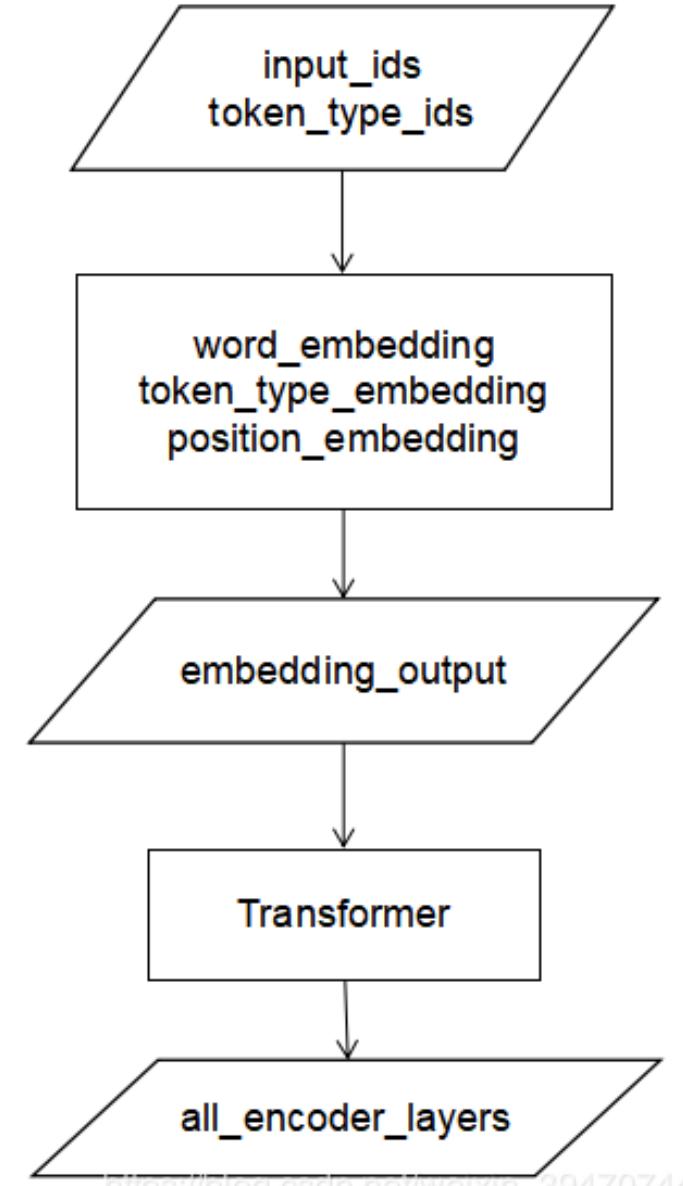
embedding_post_postprocessor

<https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py#L502-L535>

Transformer (Preparing for Attention)

2D to 3D shape of mask

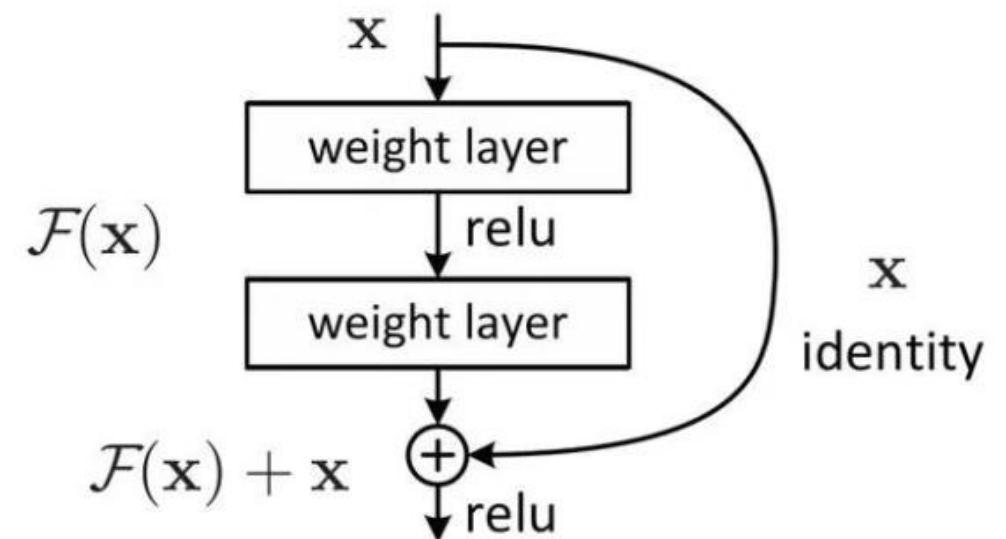
<https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py> - L199-L222



transformer_model

Residual Connection

<https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py#L823-L848>

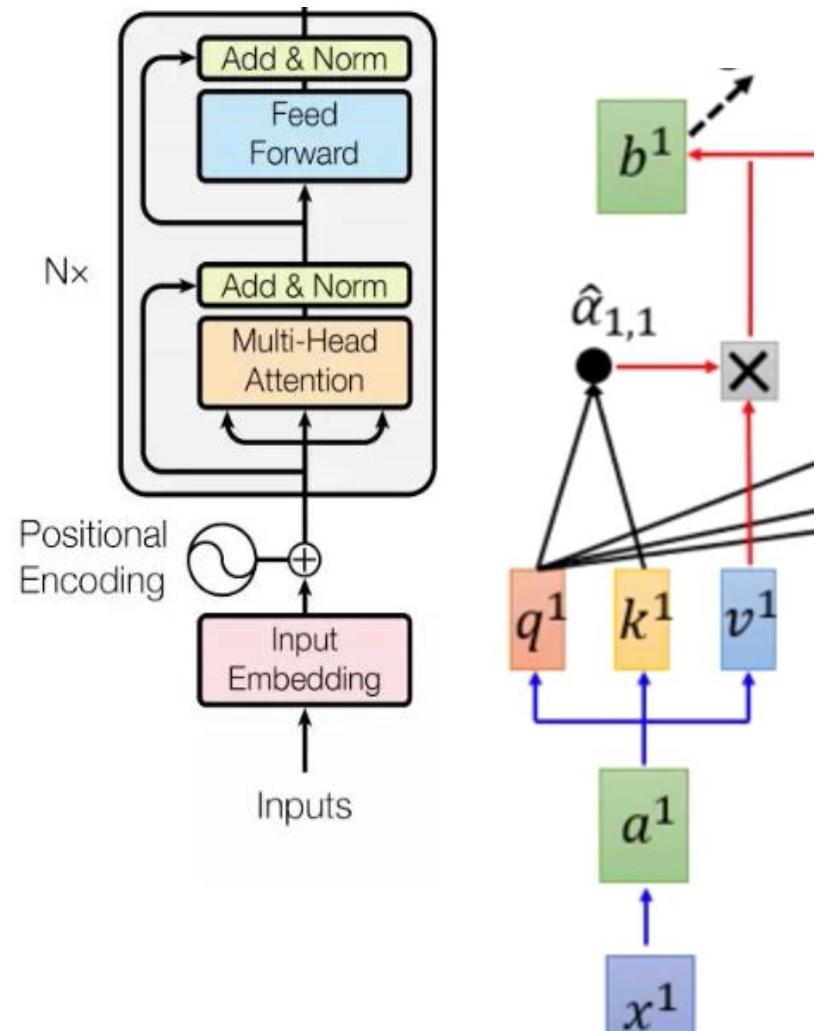


Build Q, K, V Matrix (Self-Attention)

transformer_model

<https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py> - L850-L

attention_layer: transpose_for_scores

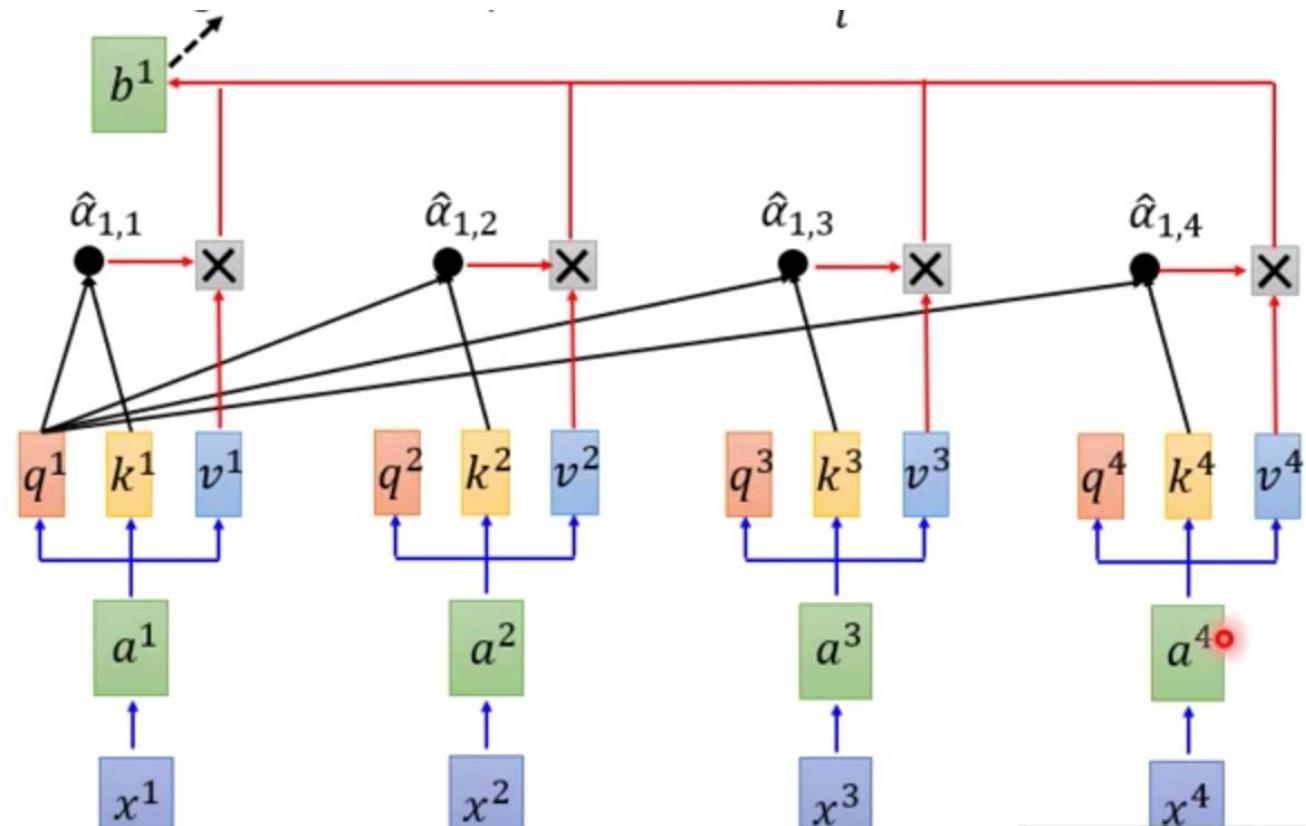


Scaled Dot-Product Attention

transformer_model

attention_layer: transpose_for_scores

<https://github.com/the-Quert/iNLPfun/blob/master/BERT/>



$$\text{Scaled Dot-Product Attention: } \alpha_{1,i} = \underbrace{q^1 \cdot k^i}_{\text{dot product}} / \sqrt{d}$$

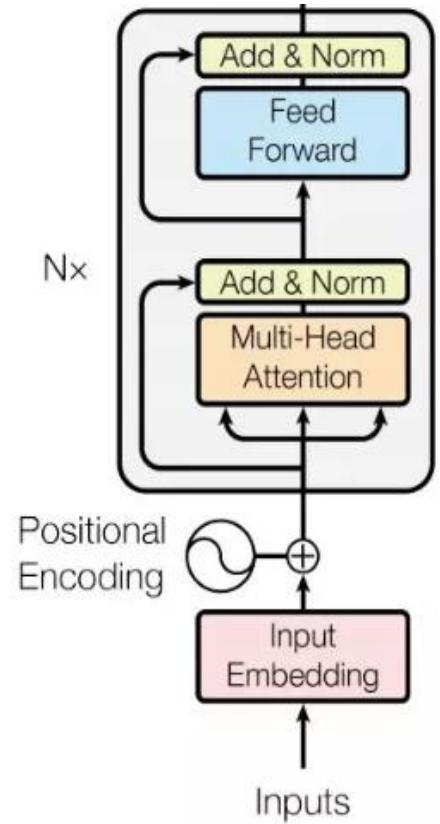
Preparing for softmax

transpose_for_scores

<https://github.com/the-Quert/iNLPfun/blob/master/BERT/modeling.py#L726-L776>

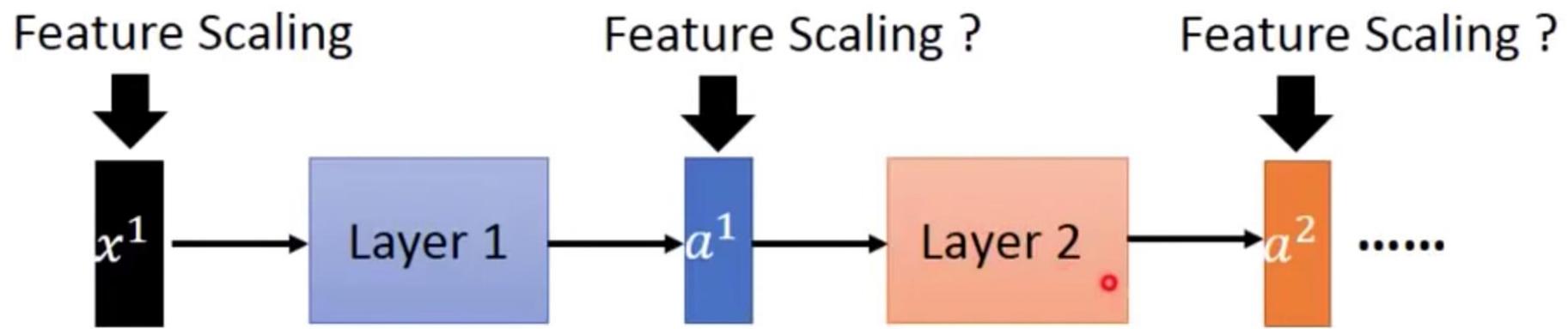
Fully Connected Layer (including Residual Connection)

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py - L450



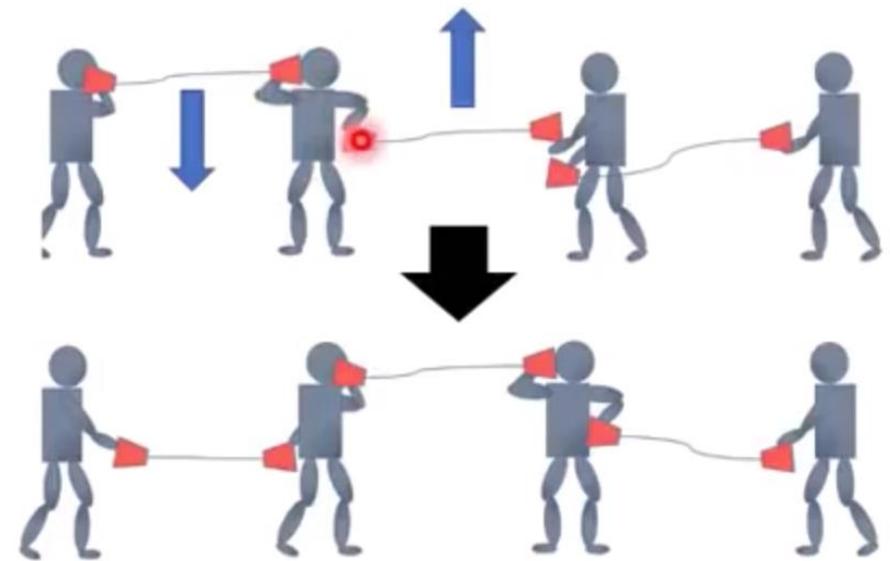
Batch Normalization

Hidden Layer



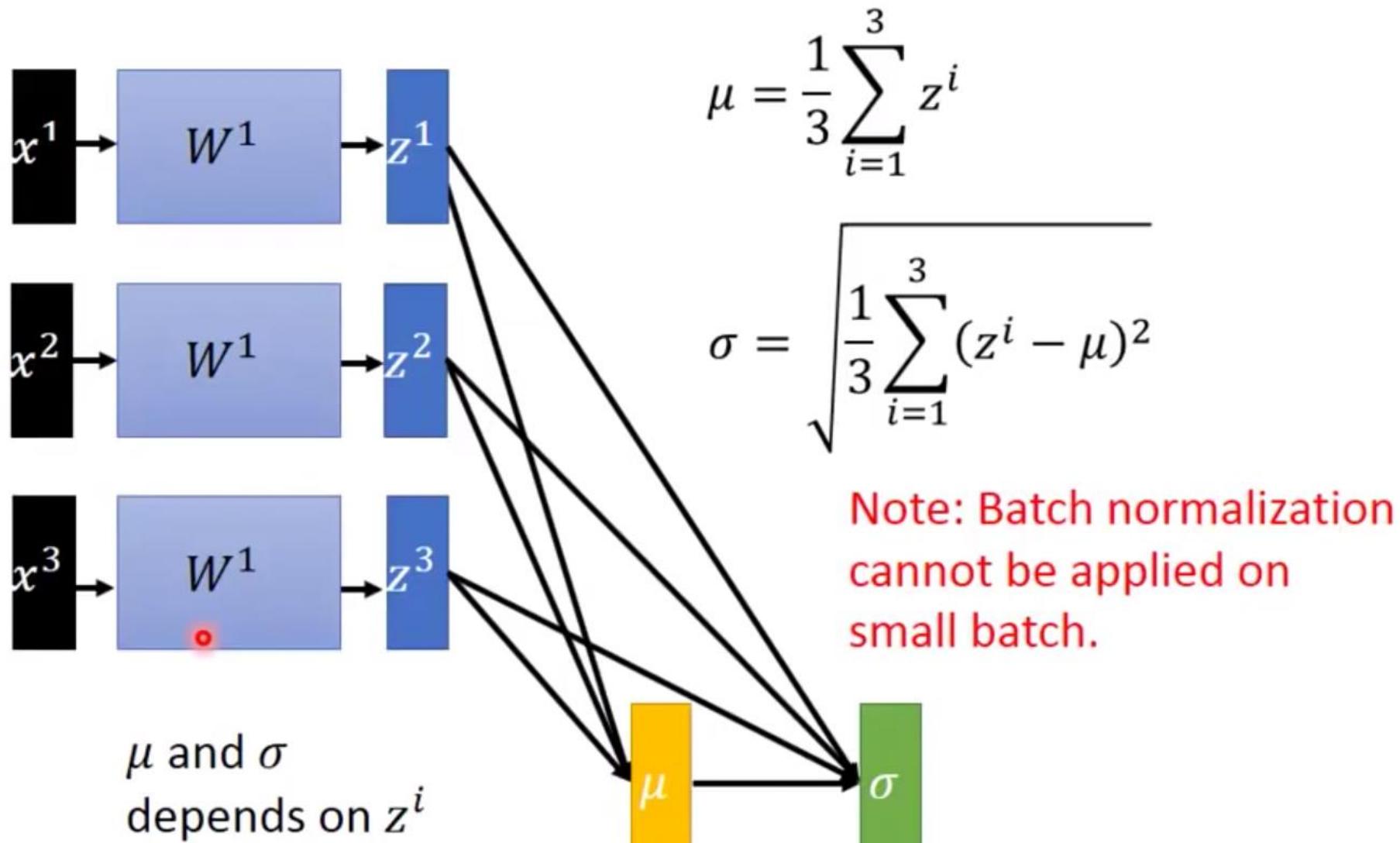
Internal Covariate Shift (ICS)

ICS refers to the change in the distribution of layer inputs caused by updates to the preceding layers. It is conjectured that such continual change negatively impacts training. The goal of BatchNorm was to reduce ICS and thus remedy this effect.



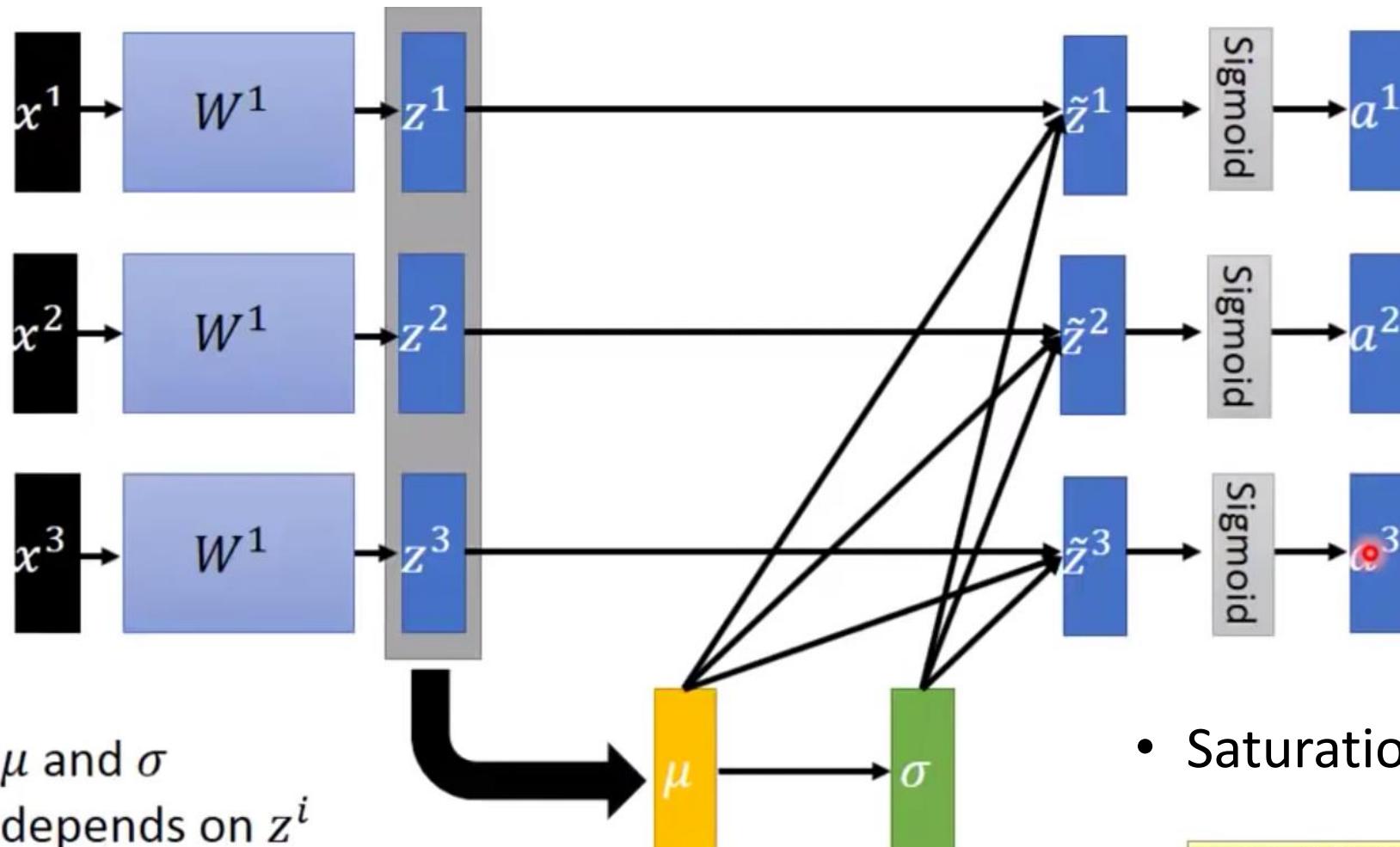
Internal Covariate Shift₇₃

Batch Normalization



Batch Normalization

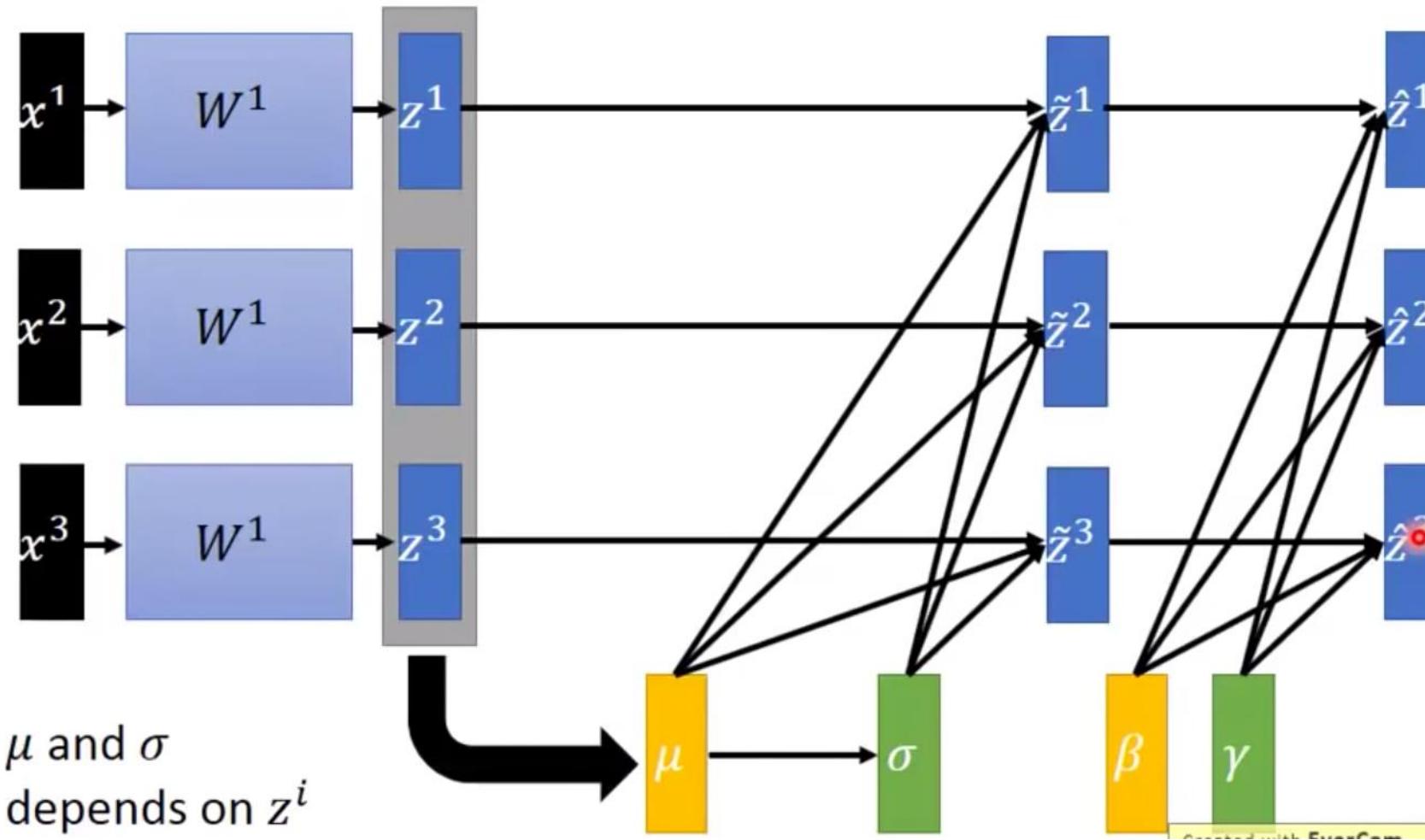
$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$



Batch Normalization

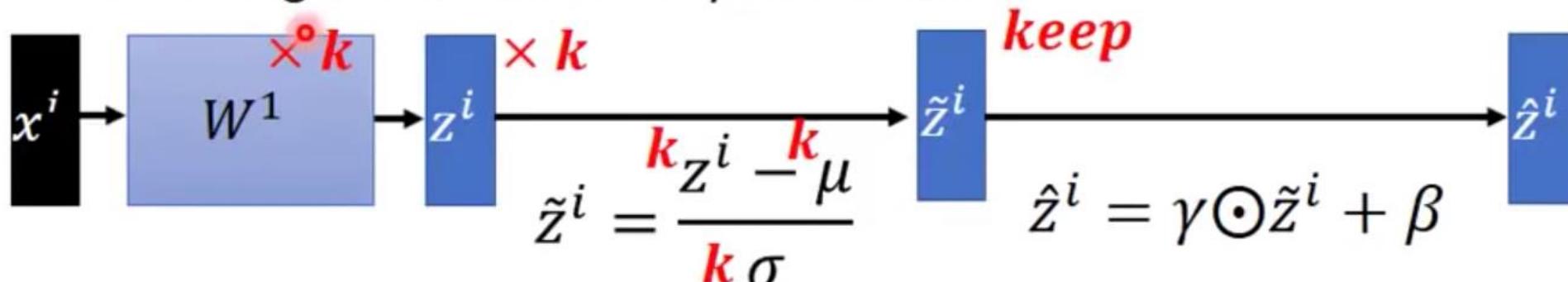
$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

$$\hat{z}^i = \gamma \odot \tilde{z}^i + \beta$$



Batch Normalization

- Benefit
 - BN reduces training times, and make very deep net trainable.
 - Because of less Covariate Shift, we can use larger learning rates.
 - Less exploding/vanishing gradients
 - Especially effective for sigmoid, tanh, etc.
 - Learning is less affected by initialization.



- BN reduces the demand for regularization.

Batch Normalization (Abstract)

Batch Normalization (BatchNorm) is a widely adopted technique that enables faster and more stable training of deep neural networks (DNNs).

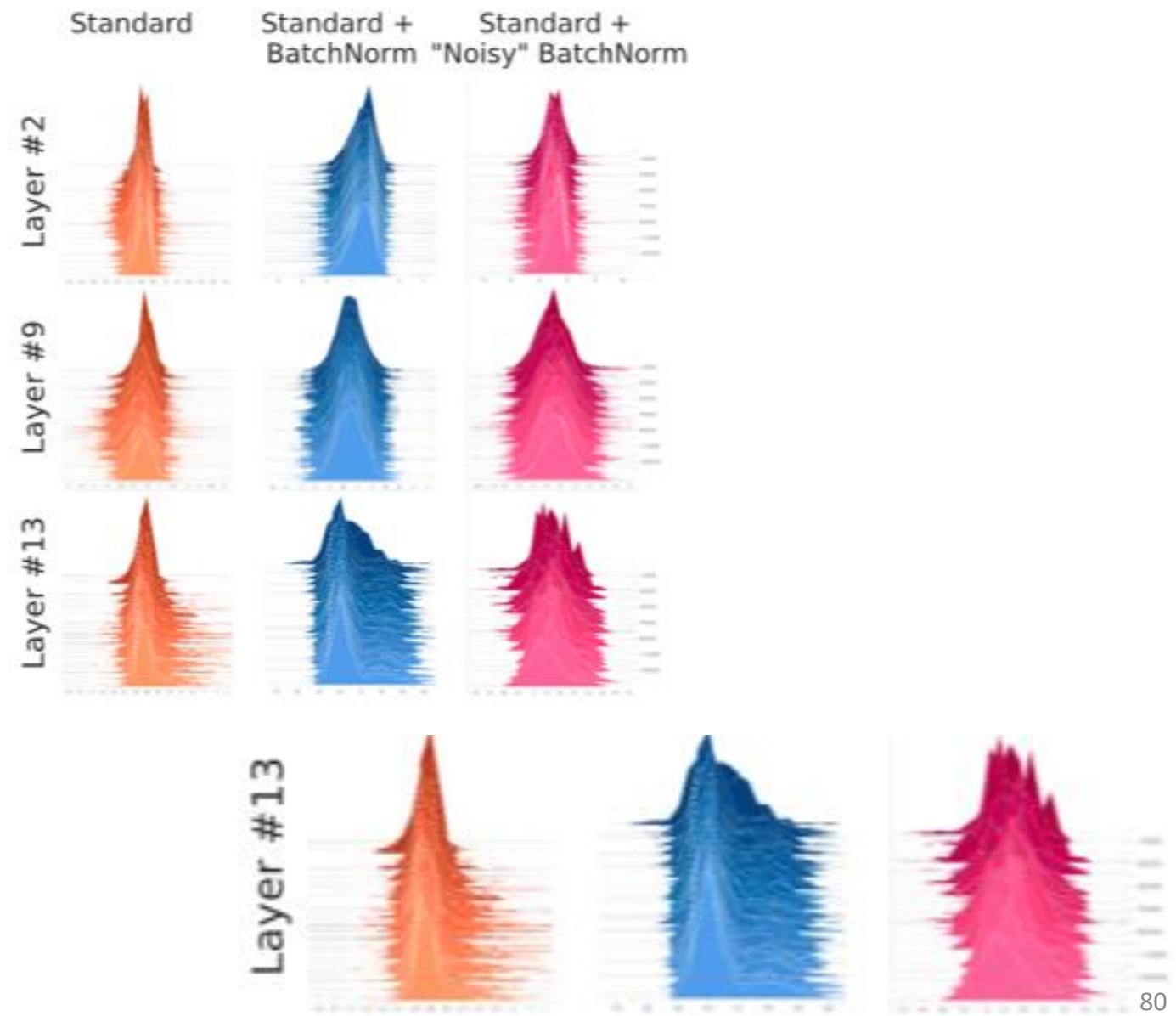
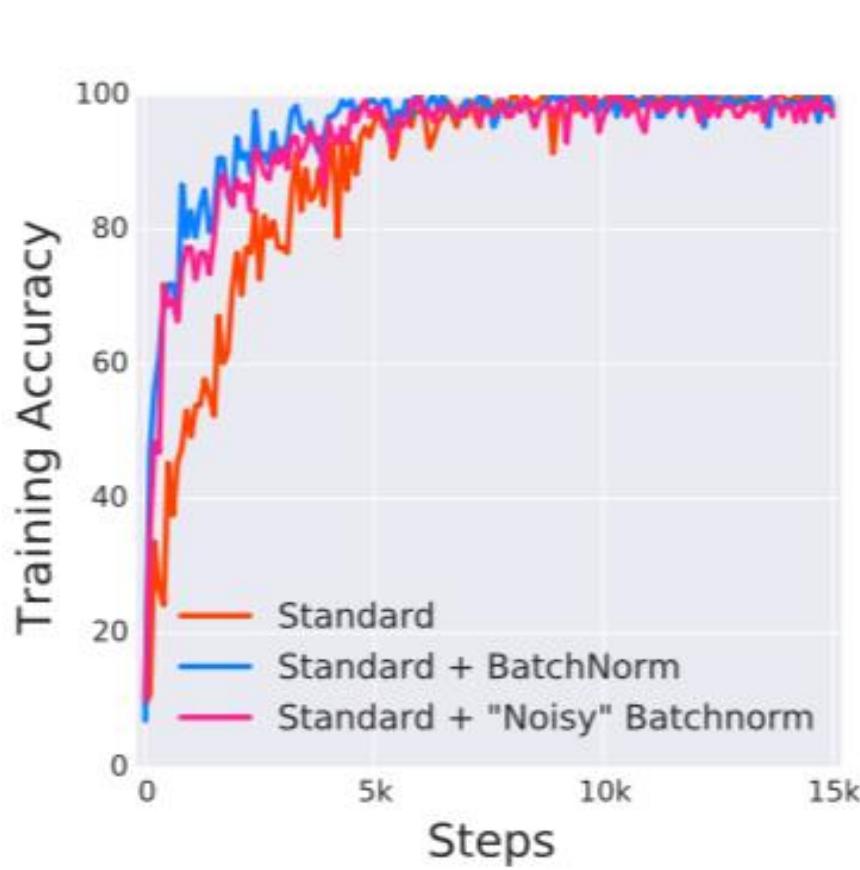
Despite its pervasiveness, the exact reasons for BatchNorm's effectiveness are still poorly understood. The popular belief is that this effectiveness stems from controlling the change of the layers' input distributions during training to reduce the so-called "internal covariate shift".

Batch Normalization (Abstract)

In this work, we demonstrate that such distributional stability of layer inputs has little to do with the success of BatchNorm.

Instead, we uncover a more fundamental impact of BatchNorm on the training process: it makes the optimization landscape significantly smoother. This smoothness induces a more predictive and stable behavior of the gradients, allowing for faster training.

Batch Normalization + Covariate Shift ?



Batch Normalization & Covariate Shift

Clearly, these findings are hard to reconcile with the claim that the performance gain due to Batch-Norm seems from increased stability of layer input distributions.

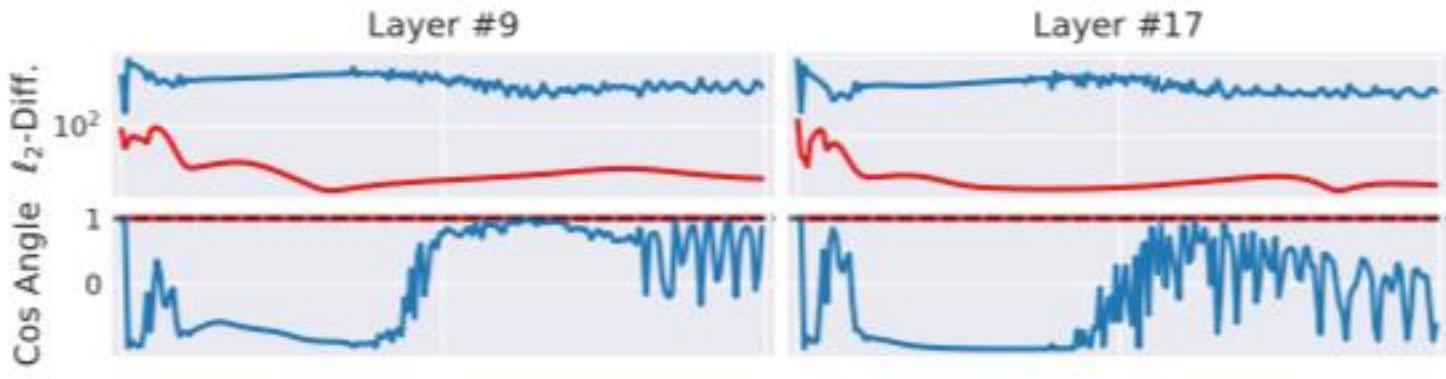
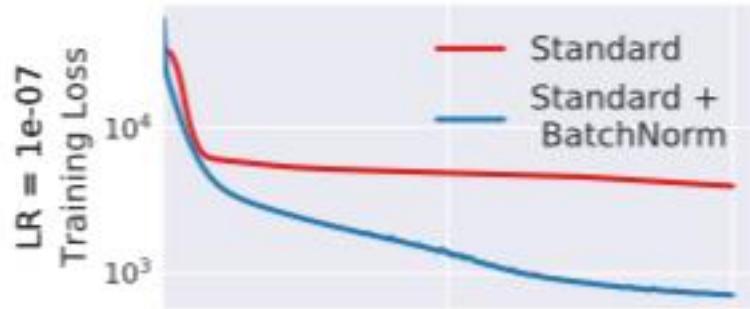
Definition 2.1. Let \mathcal{L} be the loss, $W_1^{(t)}, \dots, W_k^{(t)}$ be the parameters of each of the k layers and $(x^{(t)}, y^{(t)})$ be the batch of input-label pairs used to train the network at time t . We define internal covariate shift (ICS) of activation i at time t to be the difference $\|G_{t,i} - G'_{t,i}\|_2$, where

$$G_{t,i} = \nabla_{W_i^{(t)}} \mathcal{L}(W_1^{(t)}, \dots, W_k^{(t)}; x^{(t)}, y^{(t)})$$

$$G'_{t,i} = \nabla_{W_i^{(t)}} \mathcal{L}(W_1^{(t+1)}, \dots, W_{i-1}^{(t+1)}, W_i^{(t)}, W_{i+1}^{(t)}, \dots, W_k^{(t)}; x^{(t)}, y^{(t)}).$$

Batch Normalization works

- Deep Linear Networks

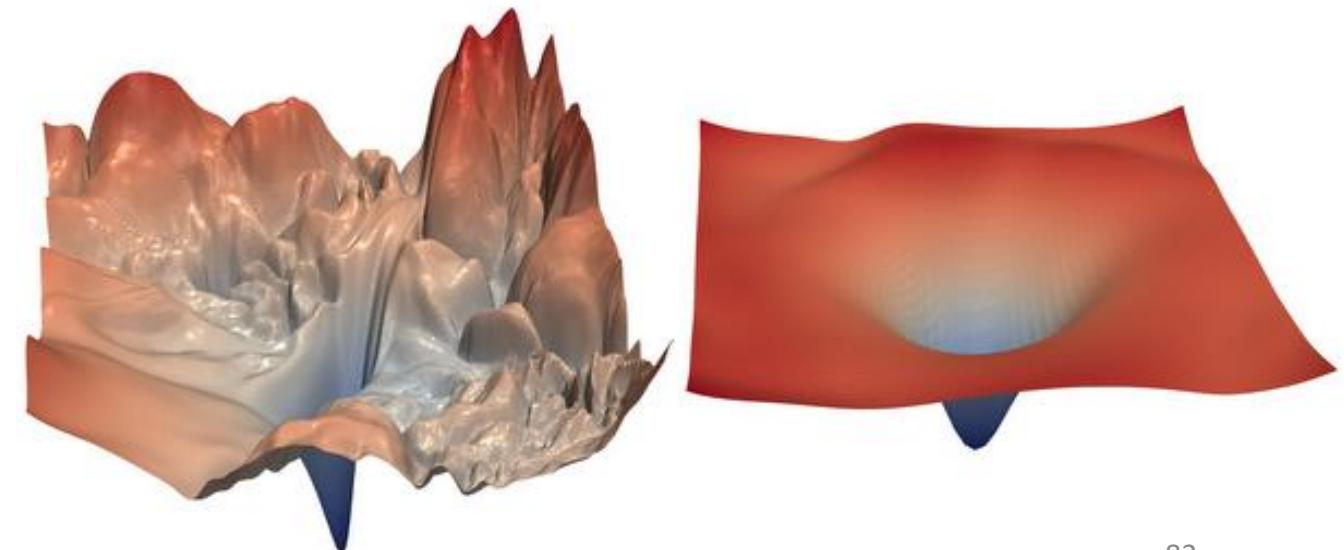


(b) DLN

Batch Normalization

We uncover a more fundamental impact of BatchNorm on the training process:

It makes the optimization landscape significantly smoother. This smoothness induces a more predictive and stable behavior of the gradients, allowing for faster training.



Batch Normalization

Indeed, we identify the key impact that BatchNorm has on the training process: it reparametrizes the underlying optimization problem to *make its landscape significantly more smooth.*

The first manifestation of this impact is improvement in the Lipschitzness of the loss function.

²Recall that f is L -Lipschitz if $|f(x_1) - f(x_2)| \leq L\|x_1 - x_2\|$, for all x_1 and x_2

Batch Normalization

That is, the loss changes at a smaller rate and the magnitudes of the gradients are smaller too. There is, however, an even stronger effect at play. Namely, BatchNorm's reparametrization makes *gradients* of the loss more Lipschitz too.

Batch Normalization

To understand why, recall that in a vanilla (non-BatchNorm), deep neural network, the loss function is not only non-convex but also tends to have a large number of “kinks”, flat regions, and sharp minima.

This makes gradient descent–based training algorithms unstable, e.g., due to exploding or vanishing gradients, and thus highly sensitive to the choice of the learning rate and initialization.

input_mask

https://github.com/the-Quert/iNLPfun/blob/master/BERT/run_classifier.py#L446-L454

Word Embedding & Word2vec

One-hot encoding and WordNet

Means one 1, the rest 0s

Words can be represented by **one-hot** vectors:

motel = [0 0 0 0 0 0 0 0 0 1 0 0 0]

hotel = [0 0 0 0 0 0 1 0 0 0 0 0 0]

- Problem with words as discrete symbols
- WordNet – synonyms, hypernyms, hyponyms
- Word meanings, relationships between words
- Problems with WordNet

e.g. synonym sets containing “good”:

```
from nltk.corpus import wordnet as wn
poses = { 'n':'noun', 'v':'verb', 's':'adj (s)', 'a':'adj', 'r':'adv'}
for synset in wn.synsets("good"):
    print("{}: {}".format(poses[synset.pos()],
                          ", ".join([l.name() for l in synset.lemmas()])))
```

```
noun: good
noun: good, goodness
noun: good, goodness
noun: commodity, trade_good, good
adj: good
adj (sat): full, good
adj: good
adj (sat): estimable, good, honorable, respectable
adj (sat): beneficial, good
adj (sat): good
adj (sat): good, just, upright
...
adverb: well, good
adverb: thoroughly, soundly, good
```

e.g. hypernyms of “panda”:

```
from nltk.corpus import wordnet as wn
panda = wn.synset("panda.n.01")
hyper = lambda s: s.hypernyms()
list(pandaclosure(hyper))
```

```
[Synset('procyonid.n.01'),
Synset('carnivore.n.01'),
Synset('placental.n.01'),
Synset('mammal.n.01'),
Synset('vertebrate.n.01'),
Synset('chordate.n.01'),
Synset('animal.n.01'),
Synset('organism.n.01'),
Synset('living_thing.n.01'),
Synset('whole.n.02'),
Synset('object.n.01'),
Synset('physical_entity.n.01'),
Synset('entity.n.01')]
```

Representing words by their context

- Distribution semantics

“You shall know a word by the company it keeps” (J. R. Firth 1957: 11)

...government debt problems turning into banking crises as happened in 2009...

...saying that Europe needs unified banking regulation to replace the hodgepodge...

...India has just given its banking system a shot in the arm...

Word meaning as a neural word vector - visualization

- Word vectors (word embedding, word representation)

$$\text{banking} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}$$

- Distributed representation

- Two-dimensional view



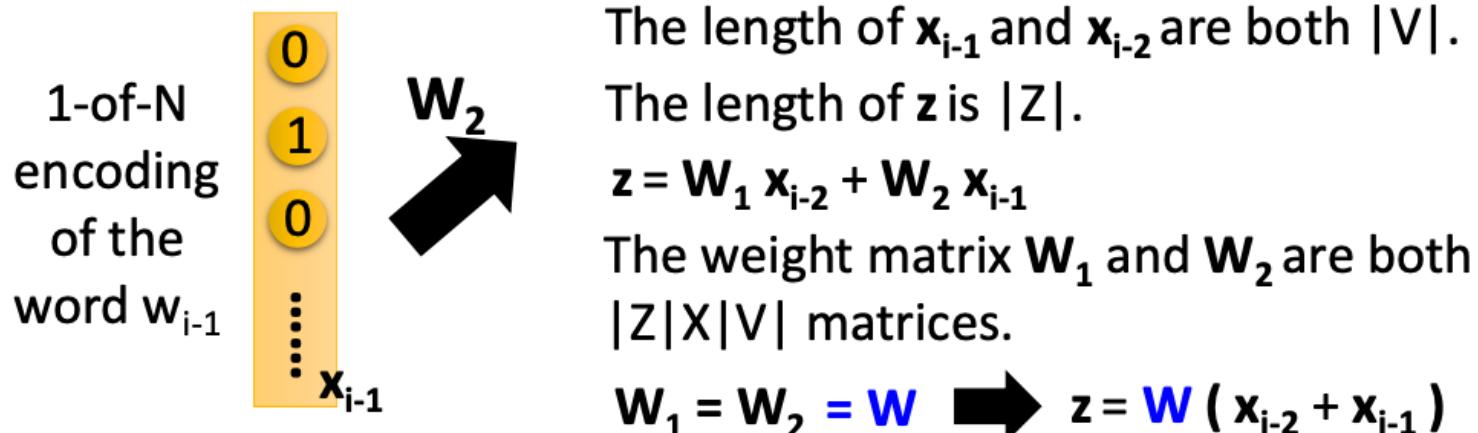
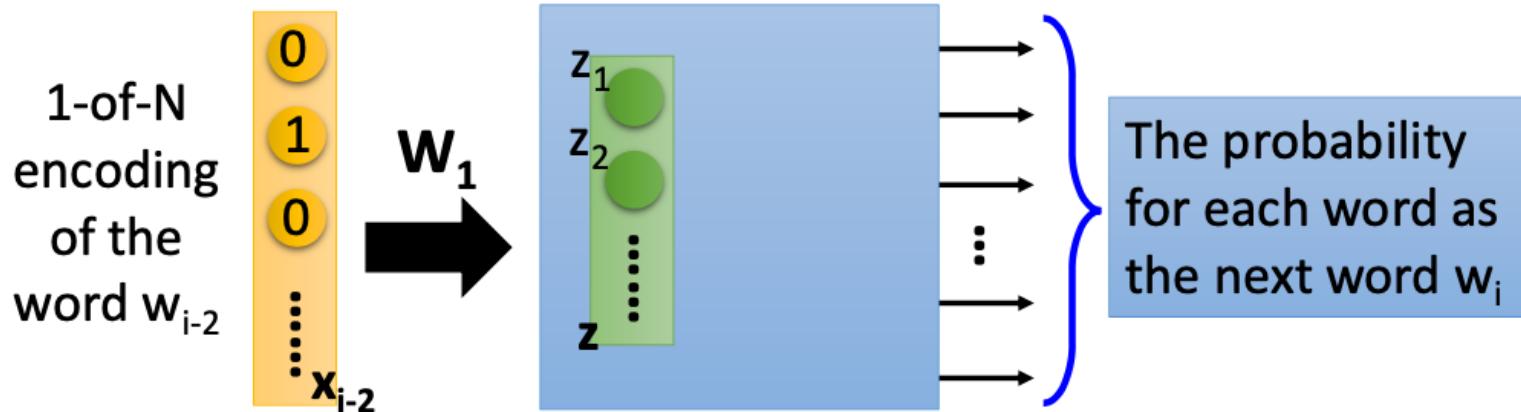
Word Embedding

- Count based
 - If two words w_i and w_j frequently co-occur, $V(w_i)$ and $V(w_j)$ would be close to each other
 - E.g. Glove Vector:
<http://nlp.stanford.edu/projects/glove/>
- Prediction based



- Word2vec

Word Embedding - Prediction based

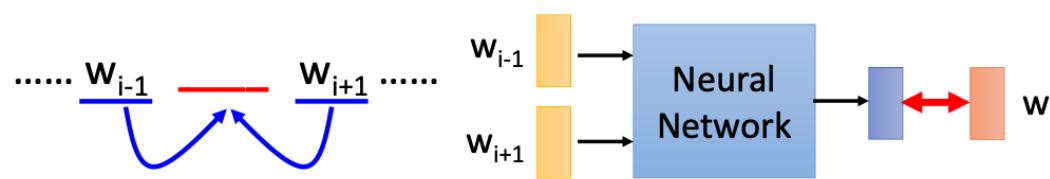


- Make W_i equal to W_j , given W_i and W_j the same initialization

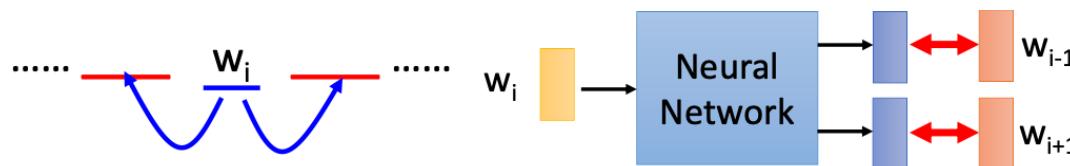
$$w_i \leftarrow w_i - \eta \frac{\partial C}{\partial w_i} - \eta \frac{\partial C}{\partial w_j}$$
$$w_j \leftarrow w_j - \eta \frac{\partial C}{\partial w_j} - \eta \frac{\partial C}{\partial w_i}$$

Word2vec

- **Continuous Bag of Words (CBOW)**



- **Skip-gram (SG)**



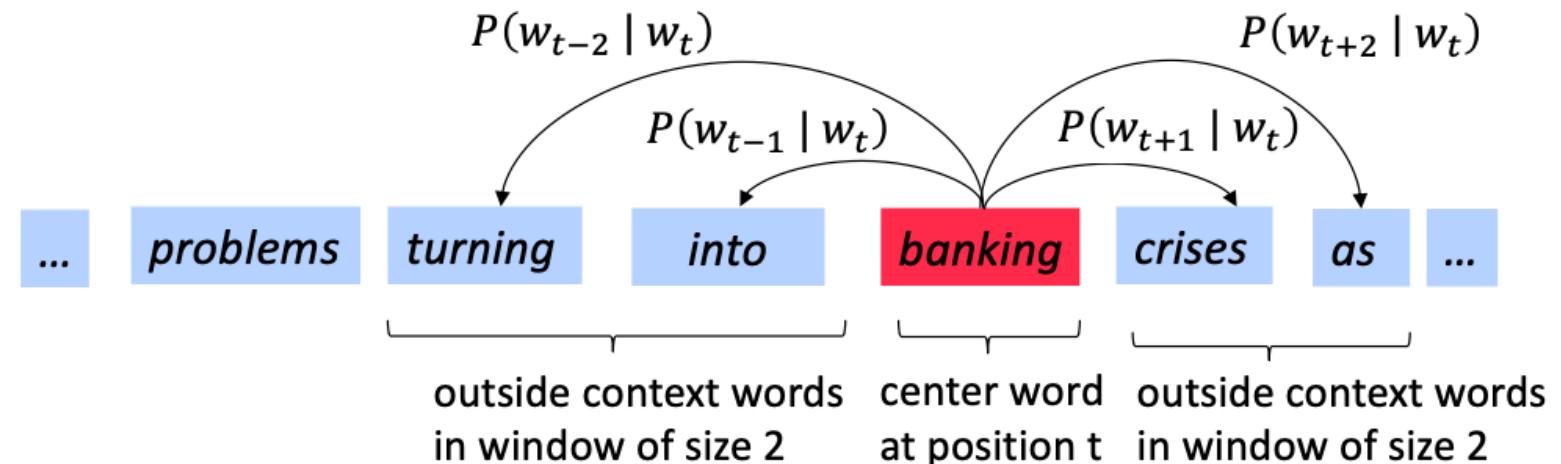
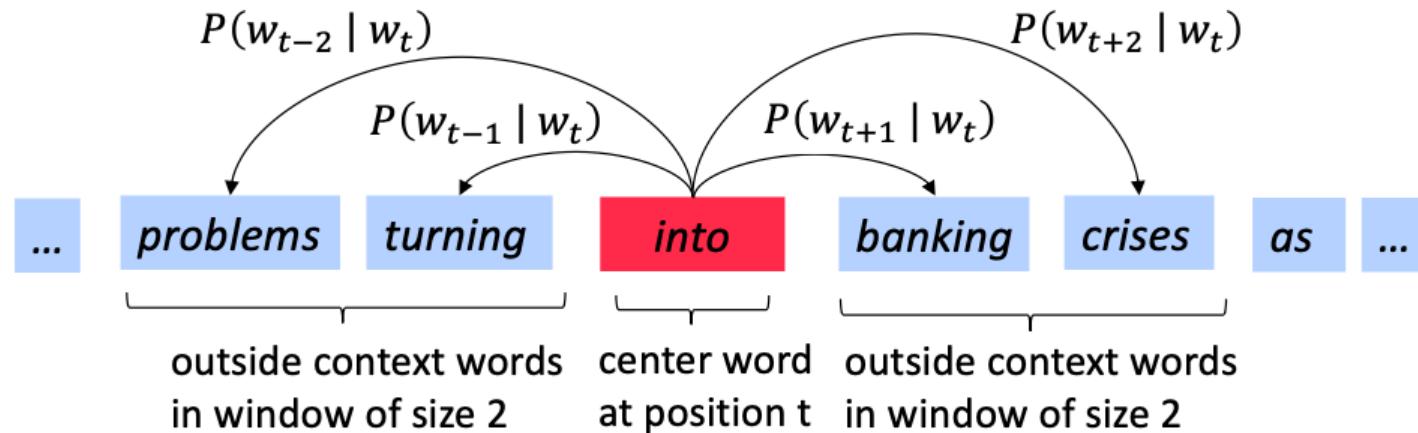
Word2vec

Word2vec (Mikolov et al. 2013) is a framework for learning word vectors.

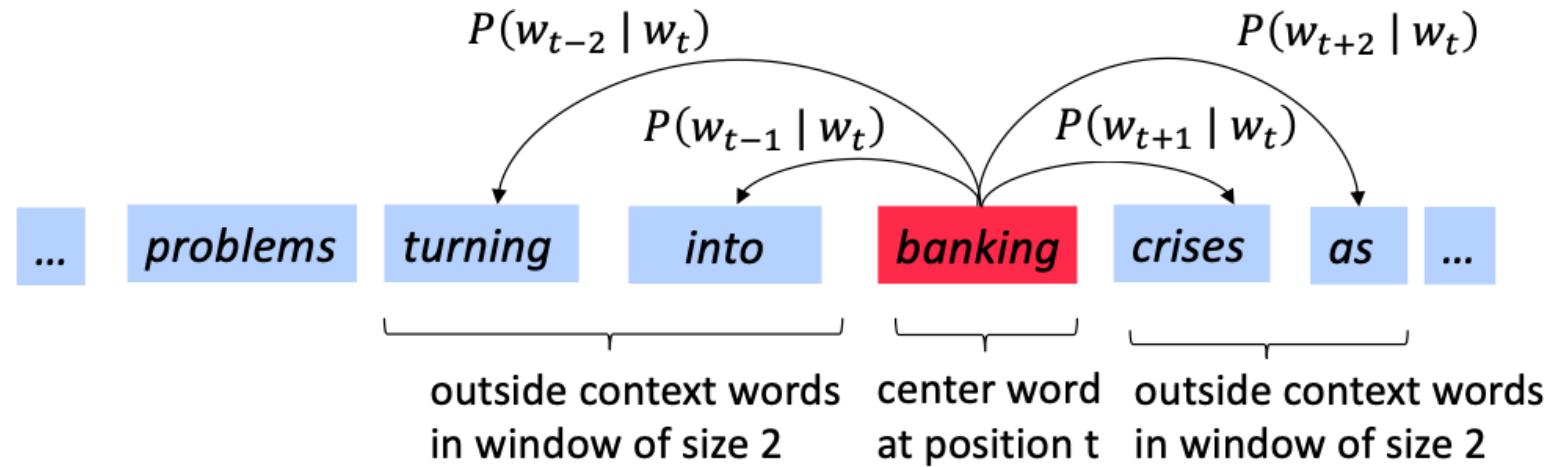
Idea:

- We have a large corpus of text
- Every word in a fixed vocabulary is represented by a vector
- Go through each position t in the text, which has a center word c and context (“outside”) words o
- Use the similarity of the word vectors for c and o to calculate the probability of o given c (or vice versa)
- Keep adjusting the word vectors to maximize this probability

Word2vec



Word2vec objective function



$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta) \quad (3)$$

- Objective function

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta) \quad (4)$$

Word2vec objective function

- Objective function

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ i \neq 0}} \log P(w_{t+j} | w_t; \theta) \quad (4)$$

- Calculate each word with two vectors
 - v_w when w is a center word
 - u_w when w is a context word

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \quad (6)$$

To train the model: Compute all vector gradients

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \quad (6)$$

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

$$\frac{\partial}{\partial u_o} \log P(o|c) = \frac{\partial}{\partial u_o} \log \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \quad (39)$$

$$= \frac{\partial}{\partial u_o} \left(\log \exp(u_o^T v_c) - \log \sum_{w \in V} \exp(u_w^T v_c) \right) \quad (40)$$

$$= \frac{\partial}{\partial u_o} \left(u_o^T v_c - \log \sum_{w \in V} \exp(u_w^T v_c) \right) \quad (41)$$

$$= v_c - \frac{\sum \frac{\partial}{\partial u_o} \exp(u_w^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \quad (42)$$

$$= v_c - \frac{\exp(u_o^T v_c) v_c}{\sum_{w \in V} \exp(u_w^T v_c)} \quad (43)$$

$$= v_c - \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} v_c \quad (44)$$

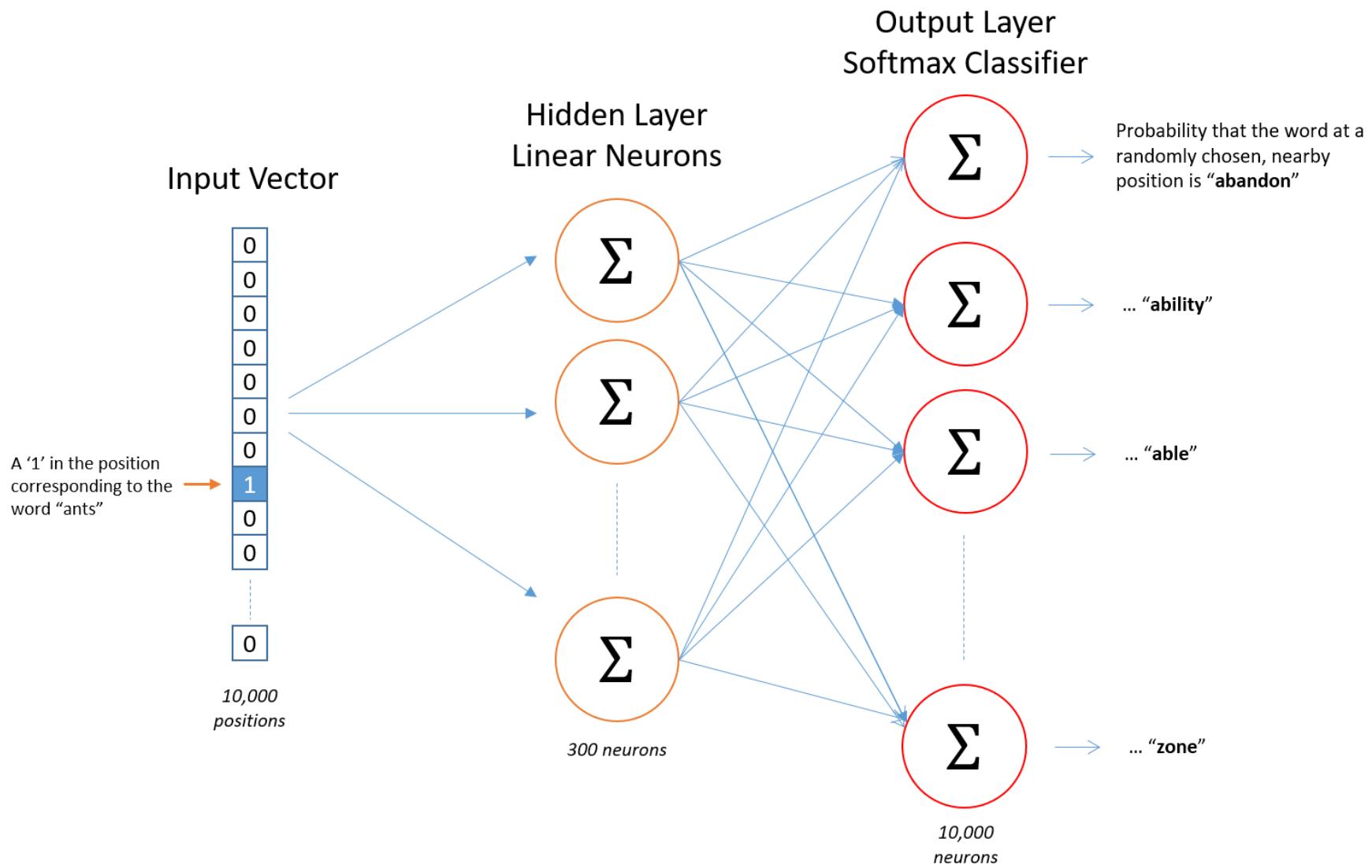
$$= v_c - P(o|c) v_c \quad (45)$$

$$= (1 - P(o|c)) v_c \quad (46)$$

Model Details

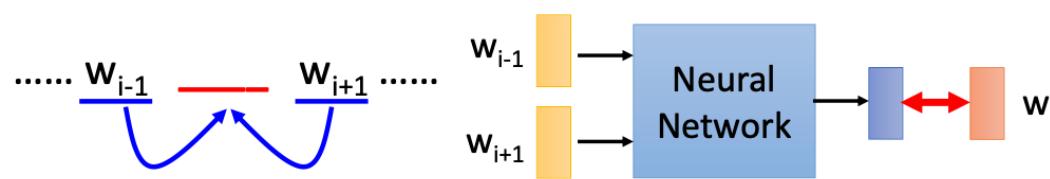
Source Text	Training Samples
The quick brown fox jumps over the lazy dog. ➔	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. ➔	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. ➔	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. ➔	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

Model Details

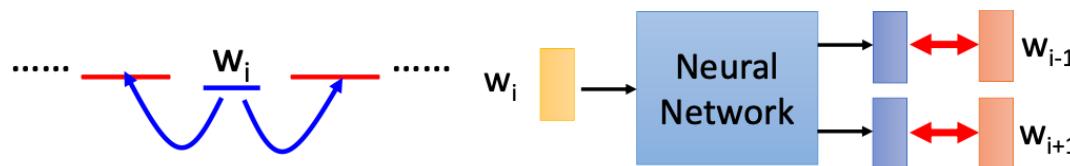


Word2vec

- **Continuous Bag of Words (CBOW)**



- **Skip-gram (SG)**



Training Model

- Negative Sampling
- Sub-sampling

Model Details (Negative Sampling)

Word pairs -> Words

Phrase detection is covered in the “Learning Phrases” section of their [paper](#). They shared their implementation in word2phrase.—I’ve shared a commented (but otherwise unaltered) copy of this code [here](#).

Model Details (Negative Sampling)

Each pass of their tool only looks at combinations of 2 words, but you can run it multiple times to get longer phrases.

So, the first pass will pick up the phrase “New_York”, and then running it again will pick up “New_York_City” as a combination of “New_York” and “City”.

Model Details (Sub-Sampling)

Source Text	Training Samples
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. →	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. →	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

Model Details (Sub-Sampling)

For each word we encounter in our training text, there is a chance that we will effectively delete it from the text. The probability that we cut the word is related to the word's frequency.

If we have a window size of 10, and we remove a specific instance of “the” from our text:

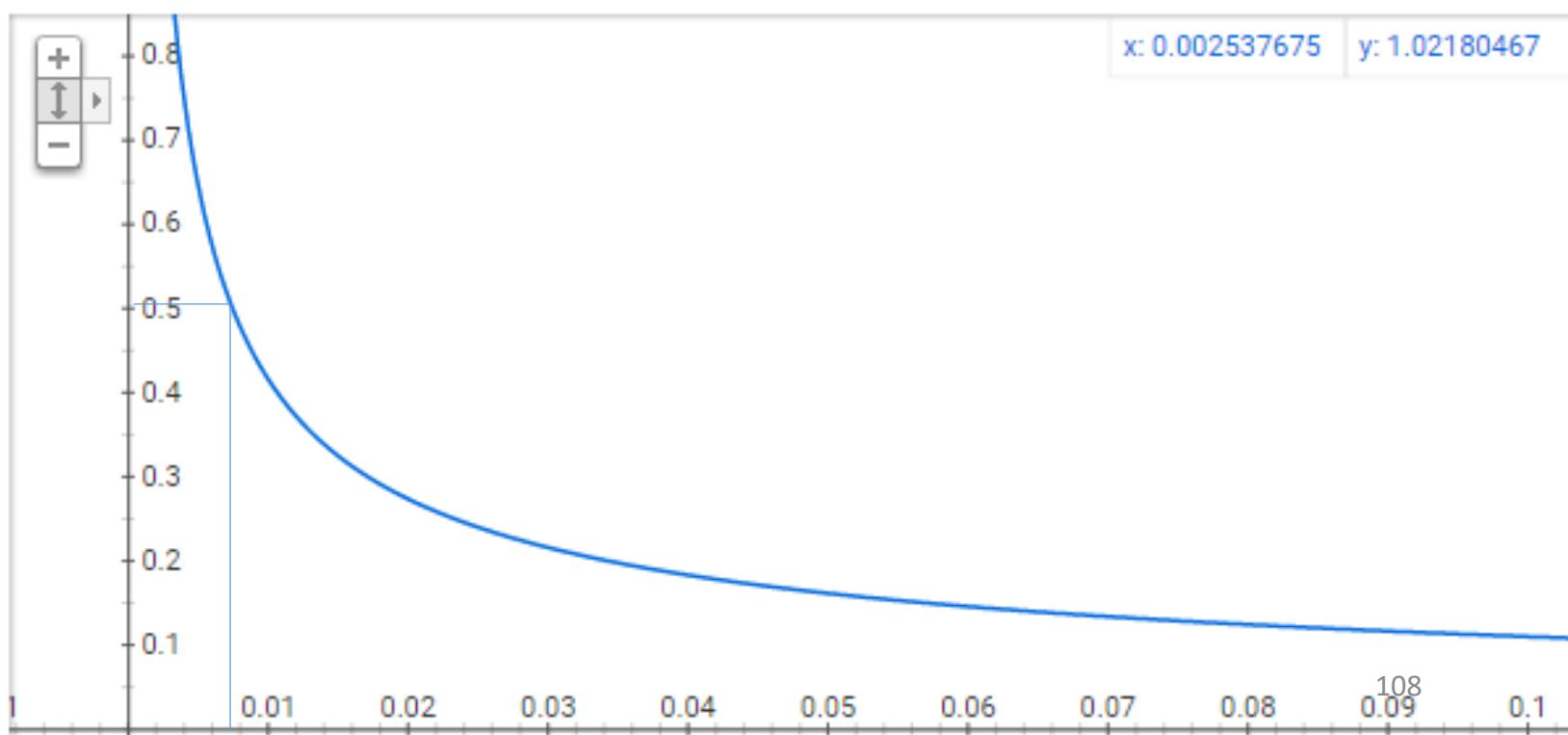
1. As we train on the remaining words, “the” will not appear in any of their context windows.
2. We'll have 10 fewer training samples where “the” is the input word.

Model Details (Sub-Sampling) • Sampling Rate

$P(w_i)$ is the probability of *keeping* the word:

$$P(w_i) = \left(\sqrt{\frac{z(w_i)}{0.001}} + 1 \right) \cdot \frac{0.001}{z(w_i)}$$

Graph for $(\sqrt{x/0.001}+1)*0.001/x$



Model Details (Negative Sampling)

Negative sampling addresses this by having each training sample only modify a small percentage of the weights, rather than all of them. Here's how it works.

When training the network on the word pair ("fox", "quick"), recall that the "label" or "correct output" of the network is a one-hot vector.

That is, for the output neuron corresponding to "quick" to output a 1, and for *all* of the other thousands of output neurons to output a 0.

Model Details (Negative Sampling)

The authors state in their paper that they tried a number of variations on this equation, and the one which performed best was to raise the word counts to the 3/4 power:

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=0}^n (f(w_j)^{3/4})}$$

The paper says that selecting 5-20 words works well for smaller datasets, and you can get away with only 2-5 words for large datasets.¹¹⁰

Model Details (Negative Sampling)

With negative sampling, we are instead going to randomly select just a small number of “negative” words (let’s say 5) to update the weights for. (In this context, a “negative” word is one for which we want the network to output a 0 for). We will also still update the weights for our “positive” word (which is the word “quick” in our current example).

Recall that the output layer of our model has a weight matrix that’s $300 \times 10,000$. That’s a total of 6 output neurons, and 1,800 weight values total. That’s only 0.06% of the 3M weights in the output layer.

Word2vec (Implementation)

<https://colab.research.google.com/github/the-Quert/inlpyfun/blob/master/word2vec/word2vec.ipynb>

https://github.com/the-Quert/inlpyfun/blob/master/word2vec/wiki_to_txt.py

Word2vec (Implementation)

1 數幾里得 西元前三世紀的古希臘數學家 現在被認為是幾何之父 此畫為拉斐爾的作品 雅典學院 数学 是利用符号语言研究數量结构 变化以及空间等概念的一門学科 从某种角度看屬於形式科學的一種 數學透過抽象化和邏輯推理由計數 計算 數學家們拓展這些概念 對數學基本概念的完善 早在古埃及 而在古希臘那裡有更為嚴謹的處理 從那時開始 數學的發展便持續不斷地小幅度進展 世紀的文藝復興時期 致使數學的加速发展 直至今日 今日 數學使用在不同的領域中 包括科學 工程 醫學 經濟 學和金融學等 有時亦會激起新的數學發現 並導致全新學科的發展 數學家也研究純數學 就是數學本身的实质性內容 而不以任何實際應用為目標 雖然許多研究以純數學開始 但其過程中也發現許多應用之外 詞源 西方語言中 數學 一詞源自於古希臘語的 其有學習學問科學數學研究 即使在其語源內 其形容詞意思為 和學習有關的 用功的 亦會被用來指數學的 其在英語中 表面上的複數形式 及在法語中的表面複數形式 可溯至拉丁文的中性複數 由西塞羅譯自希臘文複數 此一希臘語被亞里士多德拿來指 萬物皆數 的概念 汉字表示的 數學 一詞大約產生于中國宋元時期 多指象數之學 但有時也含有今天上的數學意義 例如 秦九韶的 數學九章 永樂大典 數書九章 也被宋代周密所著的 條理記為 數學大略 數學通軌 明代柯尚遷著 數學鉗 清代 杜知耕著 數學拾遺 清代丁取忠撰 直到 經過中國數學名詞審查委員會研究 算學 數學 兩詞的使用狀況後 確認以 數學 表示今天意義上的數學含義 历史 奇普 印加帝國時所使用的計數工具 玛雅数字 數學有着久遠的歷史 中國古代的六艺之一就有 數學 一詞在西方有希腊語詞源 *mathematikós* 意思是 學問的基础 源于 *máthēma* 科学 知识 学问 時間的長短等抽象的數量關係 比如時間單位有日 季節和年等 算術 加減乘除 也自然而然地產生了 歷史上曾有過許多不同的記數系統 在最初有歷史記錄的時候 為了解數字間的關係 為了測量土地 以及為了預測天文事件而形成的 结构 空间及时间方面的研究 到了 世纪 算术 微积分的概念也在此時形成 隨着數學轉向形式化 從古至今 數學便一直不斷地延展 且與科學有豐富的相互作用 兩者的發展都受惠於彼此 在歷史上有著許多數學發現 並且直至今日都不斷地有新的發現 據mikhail sevryuk於 月的期刊中所說 數學評論的創刊年份 現已超過了一百九十萬份 而且每年還增加超過七萬五千份 形成 純數學與應用數學及美學 牛頓 微積分的發明者之一 每當有涉及數量 結構 空間及變化等方面的困難問題時 而這往往也拓展了數學的研究範疇 一開始 數學的運用可見於貿易 土地測量及之後的天文學 今日 且數學本身亦給出了許多的問題 牛頓和萊布尼茲是微積分的發明者 費曼發明了費曼路徑積分 這是推理及物理洞察二者的產物 而今日的弦理論亦引申出新的數學 一些數學只和生成它的領域有關 且用來解答此領域的更多問題 且可以成為一般的數學概念 即使是最純的 數學通常亦有實際的用途 此一非比尋常的事實 年諾貝爾物理獎得主維格納稱為 如同大多數的研究領域 主要的分歧為純數學和應用數學 在應用數學內 又被分成兩大領域 並且變成了它們自身的學科 統計學和電腦科學 許多數學家談論數學的 優美 其內在的美學及美 簡單 一般化 即為美的一種 另外亦包括巧妙的證明 又或者是加快計算的數值方法 如快速傅立葉變換 高德菲 哈羅德 哈代在 一個數學家的自白 符號 語言與精確性 在現代的符號中 此一圖像即產生自 $\cos \arccos \sin | \arcsin \cos |$ 世紀後才被發明出來的 在此之前 數學以文字的形式書寫出來 這種形式會限制了數學的發展 但初學者卻常對此望而却步 它被極度的壓縮 少量的符號包含著大量的訊息 如同音樂符號一般 現今的數學符號有明確的語法 並且有效地對訊息作編碼 這是其他書寫方式難以做到的 符号化和形式化使得數學迅速发展 數學語言亦對初學者而言感到困難 亦困惱著初學者的 開放 等字在數學裡有著特別的意思 數學術語亦包括如 同胚 可積性 等專有名詞 數學需要比日常用語更多的精確性 嚴謹 但在現実应用中 定理 希臘人期許著仔細的論證 但在牛頓的時代 所使用的方法則較不嚴謹 牛頓為了解決問題所做的定義 今日 當大量的計算難以被驗證時 其證明亦很難說是足夠地嚴謹 公理在傳統的思想中是 不證自明的真理 但這種想法是有問題的 在形式上 公理只是一串符號 但依據哥德爾不完備定理 儘管如此 在此意義下 數學作為科學 卡爾 弗里德里希 高斯 卡爾 弗里德里希 高斯稱數學為 科學的皇后 在拉丁原文 以及其德語 對應於 科學 的單字的意思皆為知識 領域 而實際上 科學 若認為科学是只指物理的世界時 則數學 或至少是純數學 不會是一門科學 愛因斯坦曾如此描述 數學定律越和現實有關 它們越不確定 若它們越是確定的話 它們和現實越不會有關 且因此不是卡爾 波普爾所定義的科学 但在 年代時 且波普爾推斷 大部份的數學定律 如物理及生物學一樣 是假設演繹的 比它現在看起來更接近 然而 其他的思想家 如較著名的拉卡托斯 另一觀點則為某些科學領域 如理論物理 是其公理為嘗試著符合現實的數學 而事實上 理論物理學家齊曼 john ziman 即認為科學是一種公眾知識 因此亦包含著數學 在任何的情況下 減輕了數學不使用科學方法的缺點 在史蒂芬 沃爾夫勒姆 年的著作 一種新科學 中他提出 數學家對此的態度並不一致 且因此基本上是個哲學家 是低估了其美學方面的重要性 被創造如藝術 或是 被發現 如科學 的爭議 大学院系划分中常见 科学和数学系 實際上 但在細節上卻會分開 數學的各領域 有如反映在中國算盤上的一般 如上所述 了解數字間的關係 測量土地及預測天文事件 這四種需要大致地與數量 結構 空間及變化 即算術 代數 幾何及 分析 等數學上廣泛的子領域相關連著 除了上述主要的關注之外 至邏輯 至集合論 基礎 至不同科學的經驗上的數學 應用數學 及較近代的至不確定性的嚴格研究 基礎與哲學 為了闡明數學基礎 並研究此一架構的結果 就數學邏輯本身而言 現代邏輯被分 113