

Generalising IR Operations

Many operations in an information retrieval pipeline can be thought of as “transformer” functions.

Example: BM25 Retrieval

Assignment Project Exam Help



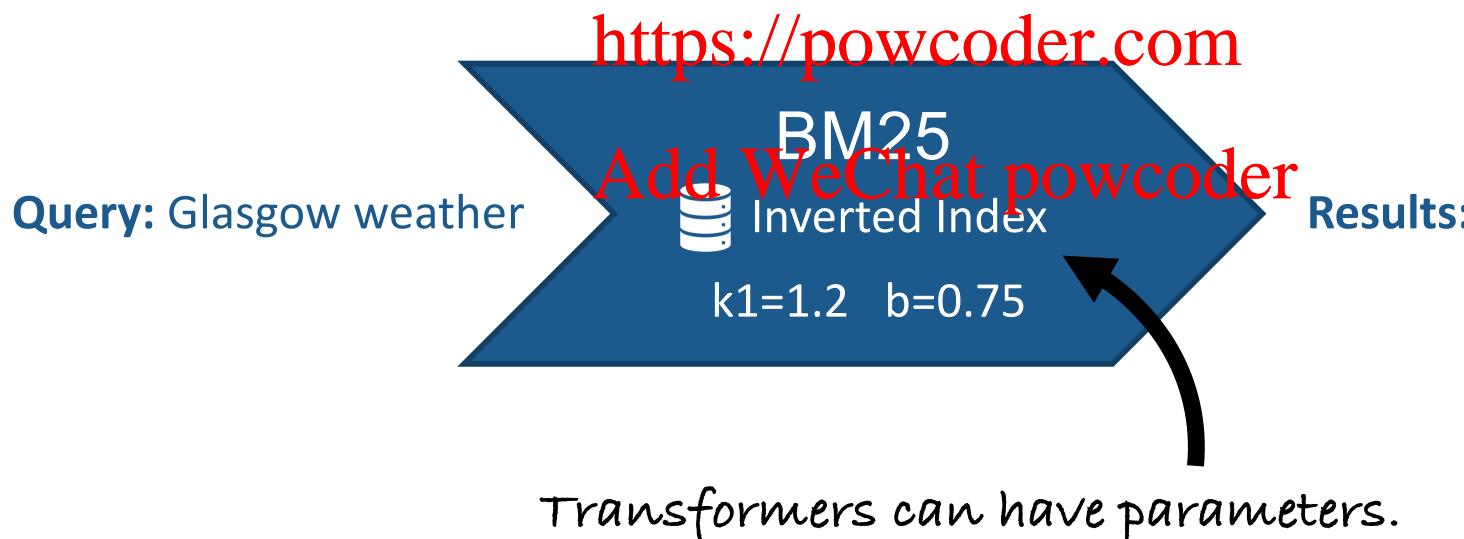
docno	score
15213	0.5134
42635	0.3742
26340	0.3223
...	...

Generalising IR Operations

Many operations in an information retrieval pipeline can be thought of a “transformer” functions.

Example: BM25 Retrieval

Assignment Project Exam Help



docno	score
15213	0.5134
42635	0.3742
26340	0.3223
...	...

Generalising IR Operations

Many operations in an information retrieval pipeline can be thought of a “transformer” functions.

Example: BM25 Retrieval

Assignment Project Exam Help

<https://powcoder.com>

BM25
Add WeChat powcoder



Inverted Index

$k_1=1.2 \ b=0.75$

qid	query
0	glasgow weather
1	flights to glasgow
...	...

qid	query	docno	score
0	glasgow...	15213	0.5134
0	glasgow...	42635	0.3742
0	glasgow...	26340	0.3223
...

More generally: a batch of queries to a batch results

Generalising IR Operations

Many operations in an information retrieval pipeline can be thought of a “transformer” functions.

Example: BM25 Retrieval

Assignment Project Exam Help

<https://powcoder.com>

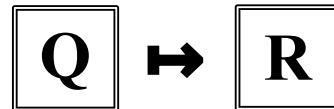
BM25
Add WeChat powcoder



Inverted Index
 $k_1=1.2 \quad b=0.75$

qid	query
0	glasgow weather
1	flights to glasgow
...	...

qid	query	docno	score
0	glasgow...	15213	0.5134
0	glasgow...	42635	0.3742
0	glasgow...	26340	0.3223
...



Shorthand: BM25 maps a query (Q) frame to a Result (R) frame.

Generalising IR Operations

The PyTerrier Data Model:



qid	query
0	glasgow weather
1	flights to glasgow
...	...

Assignment Project Exam Help

<https://powcoder.com>

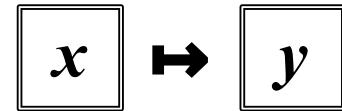


docno	text	title
0	Science & Mathematics Physics	The hot glowing...
1	School-Age Kids Growth &...	Developmental...
...



qid	query	docno	score
0	glasgow weather	15213	0.5134
0	glasgow weather	42635	0.3742
0	glasgow weather	26340	0.3223
...

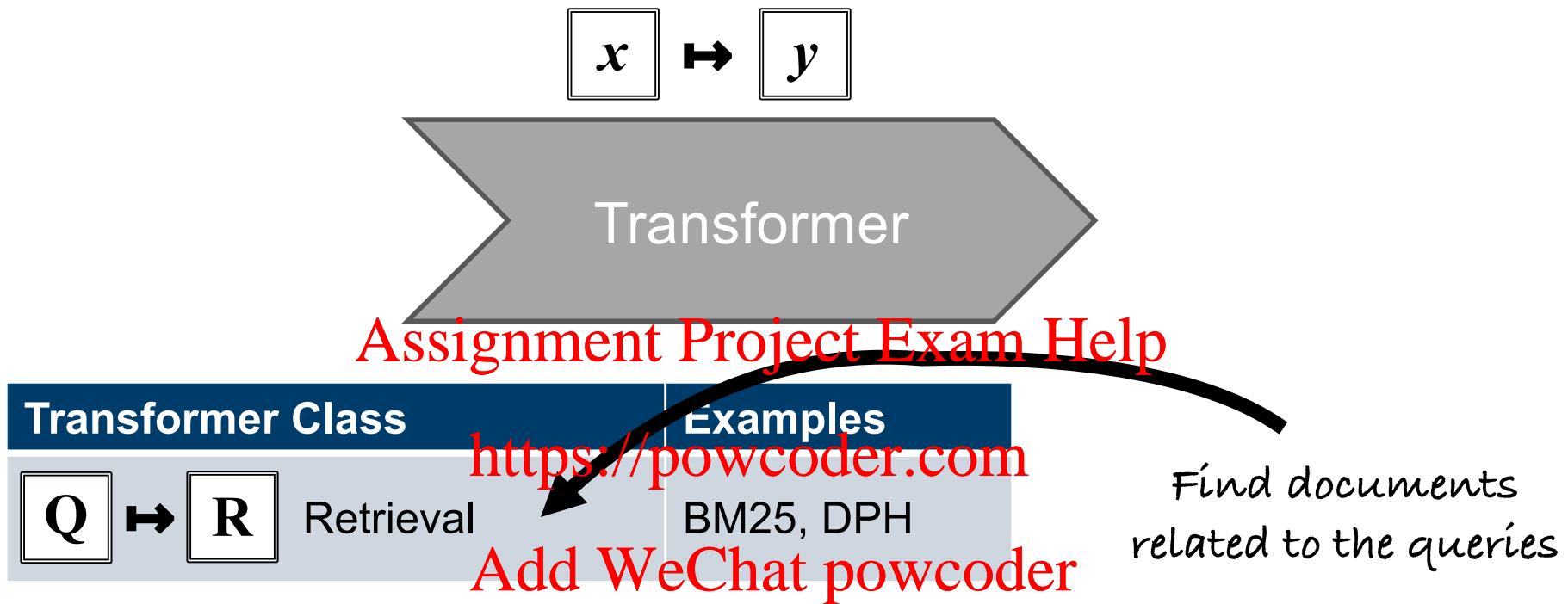
Generalising IR Operations



<https://powcoder.com>

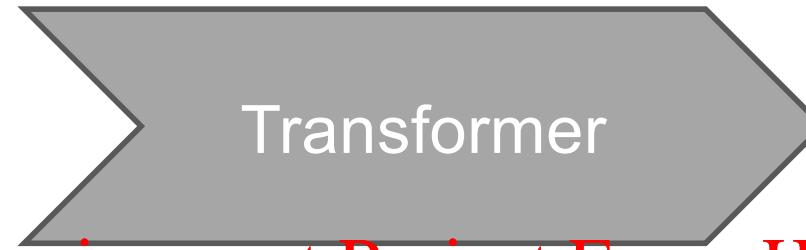
Add WeChat powcoder

Generalising IR Operations



Generalising IR Operations

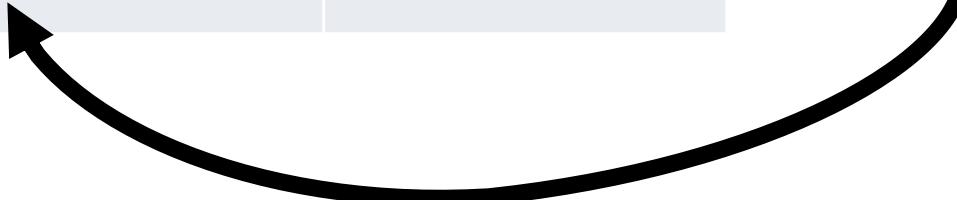
$$\boxed{x} \rightarrow \boxed{y}$$



Assignment Project Exam Help

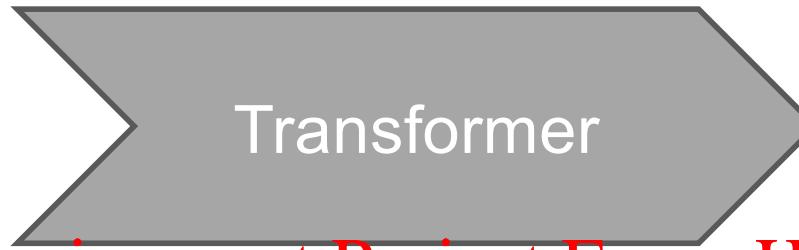
Transformer Class	Examples
$\boxed{Q} \rightarrow \boxed{R}$ Retrieval	https://powcoder.com BM25, DPH Add WeChat powcoder
$\boxed{R} \rightarrow \boxed{Q}$ PRF	RM3, Bo1

Rewrite queries based on
the returned documents.



Generalising IR Operations

$$x \rightarrow y$$



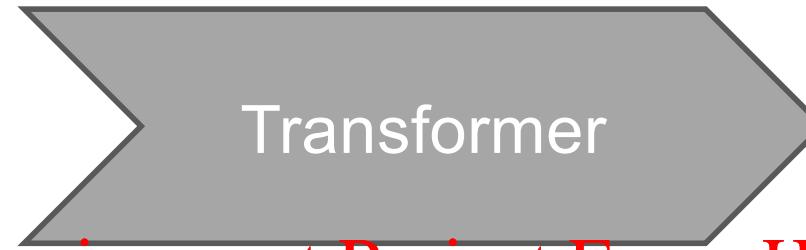
Assignment Project Exam Help

Transformer Class	Examples
$Q \rightarrow R$ Retrieval	https://powcoder.com BM25, DPH <i>Add WeChat powcoder</i>
$R \rightarrow Q$ PRF	RM3, Bo1
$R \rightarrow R$ Re-ranking	LambdaMART

Find a better order
for the results.
(usually a more expensive
method than retrieval.)

Generalising IR Operations

$$x \rightarrow y$$



Assignment Project Exam Help

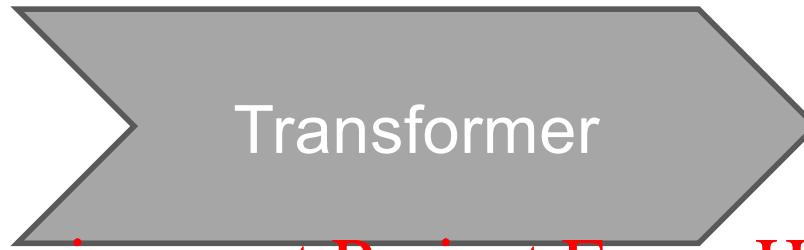
Transformer Class	Examples https://powcoder.com
$Q \rightarrow R$ Retrieval	BM25, DPH Add WeChat powcoder
$R \rightarrow Q$ PRF	RM3, Bo1
$R \rightarrow R$ Re-ranking	LambdaMART
$Q \rightarrow Q$ Query Re-writing	SDM

Build a better version
of the user's query



Generalising IR Operations

$$x \rightarrow y$$



Assignment Project Exam Help

Transformer Class	Examples
$Q \rightarrow R$ Retrieval	https://powcoder.com BM25, DPH Add WeChat powcoder
$R \rightarrow Q$ PRF	RM3, Bo1
$R \rightarrow R$ Re-ranking	LambdaMART
$Q \rightarrow Q$ Query Re-writing	SDM
$D \rightarrow D$ Doc. Re-writing	Passaging

Build a better version of documents to index.

Generalising IR Operations

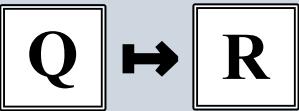
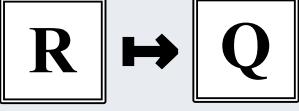
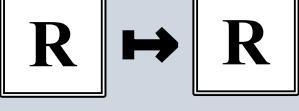
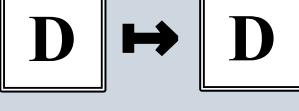
These operations are composable!

Example: Perform SDM, then BM25, then RM3, then BM25 again.

Assignment Project Exam Help

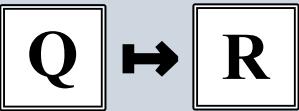
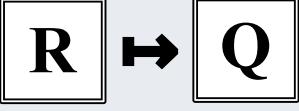
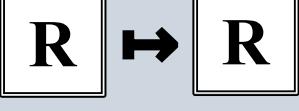
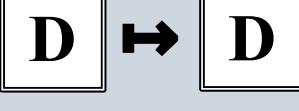


What's the point of all this?

Transformer Class	Examples
 Retrieval	Assignment Project Exam Help BM25, DPH
 PRF	https://powcoder.com RM3, Bo1
 Re-ranking	Add WeChat powcoder LambdaMART
 Query Re-writing	SDM
 Doc. Re-writing	Passaging

What's the point of all this?

All of these transformations can be replaced with a Neural Network!

Transformer Class	Examples	Neural Examples
 Retrieval	Assignment Project Exam Help BM25, DPH	ANCE, CoBERT
 PRF	https://powcoder.com Add WeChat powcoder	CoBERT-PRF
 Re-ranking	LambdaMART	CEDR, monoT5
 Query Re-writing	SDM	T5-QE, IntenT5
 Doc. Re-writing	Passaging	Doc2Query, DeepImpact

Why Neural Networks?

- Neural Natural Language Processing (NLP) techniques are highly effective: state-of-the-art at most NLP tasks.
 - Particularly: Using models trained on a “language modeling” objective transfer the knowledge well to other tasks.

Assignment Project Exam Help

<https://powcoder.com>

- Able to learn complex & subtle patterns from training data
 - Automatically learns to overcome the lexical gap, proximity relations, term salience (importance), coreference, etc.
 - This means less manual “feature engineering”

Building comprehensive rules/heuristics for language is challenging.

What about other usages of “where”?

Where was iodine discovered?

BM25=12.05

Iodine was discovered by French chemist Bernard Courtois in 1811 in France. He was attempting to extract potassium chloride from seaweed. After crystallizing the potassium chloride, he added sulfuric acid to the remaining liquid, which rather surprisingly, produced a purple vapor, which condensed into dark crystals.

BM25=12.18

Should “where” queries match “in”? <https://powcoder.com>

What about other usages of “in”?

How close does it need to be to other query terms?

Building comprehensive rules/heuristics for language is challenging.

iodine discovered place

BM25=12.05

Iodine was discovered by Bernard Courtois in 1811 in France. Courtois was trying to extract potassium chloride from sea kelp. Would this document be less relevant if this phrase was later in the document? rather surprisingly, produced a purple vapor, which condensed into dark crystals.

Assignment Project Exam Help

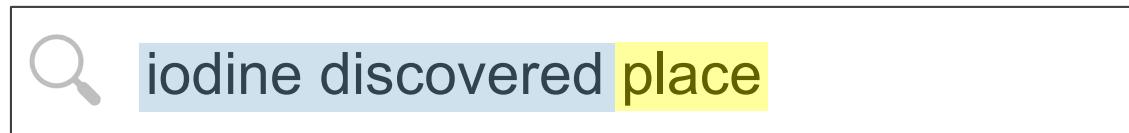
https://powcoder.com Add WeChat powcoder

What about other formulations of the query?

Or other formulations of the document?

18

Building comprehensive rules/heuristics for language is challenging.



BM25=12.05

Iodine was discovered by Bernard Courtois in 1811 in France. Courtois was trying to extract potassium chloride from seaweed. After crystallization, he found a purple vapor, which condensed into dark crystals.

Is "France" really enough? Would the city be better?

Is the fact that the discoverer was French close enough? Iodine's property may be the ability to kill germs.

BM25=12.18

Project Exam Help

The element was discovered in 1811 by French chemist Bernard Courtois (1777-1838). The element occurs primarily in seawater. Is the fact that the discoverer was French close enough? Property may be the ability to kill germs.

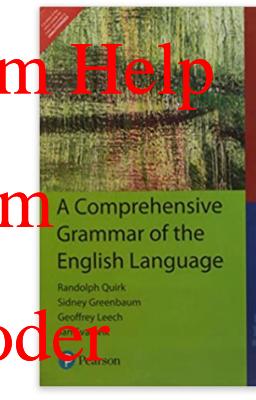
Natural language is messy

The same idea can be described by many sequences of words.

Meanwhile... Assignment Project Exam Help
words can describe
different ideas (ambiguity).

Add WeChat powcoder
Even the grammar of languages
themselves are challenging to
fully describe.

We can instead let a neural
network learn how to deal with
this mess.



A Comprehensive Grammar o...

Paperback – 1 Jan. 2011

by Geoffrey Leech & Jan Svartvik Randolph Quirk, Sidney G
★★★★★ 126 ratings

See all formats and editions

Hardcover
from £96.95

Paperback
£39.95

1 Used from £96.95

4 Used from £26.98
6 New from £26.99

BRAND NEW, perfect condition.

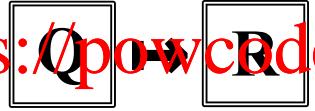
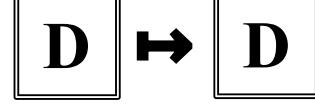
Report incorrect product information.

Print length
1
1779 pages

Language
English

Publisher
Pearson

Today

1. Review of LTR & Basics of Neural Networks for NLP
2. Neural Re-ranking  Assignment Project Exam Help
3. Neural Retrieval  <https://powcoder.com>
4. Neural Query Rewriting & PRF  Add WeChat powcoder
5. Neural Document Rewriting 
6. Neural IR in PyTerrier

Review of Supervised Learning

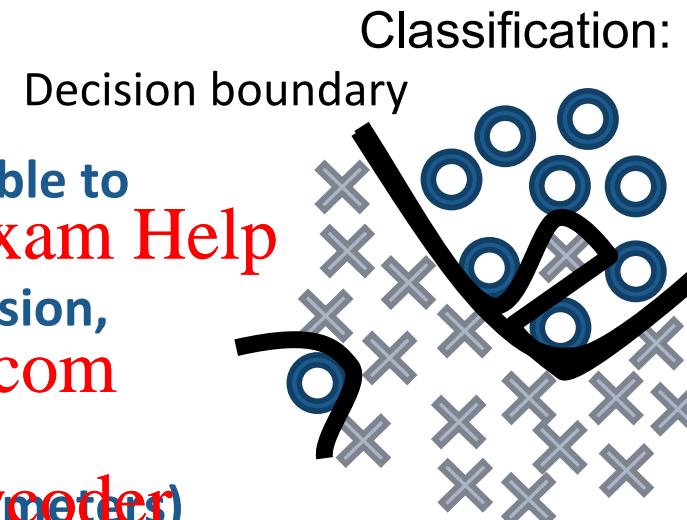
Use **training data** to make a model that is able to predict something about new inputs.

Many algorithms exist (e.g., Logistic Regression, LambdaMART, etc.) <https://powcoder.com>

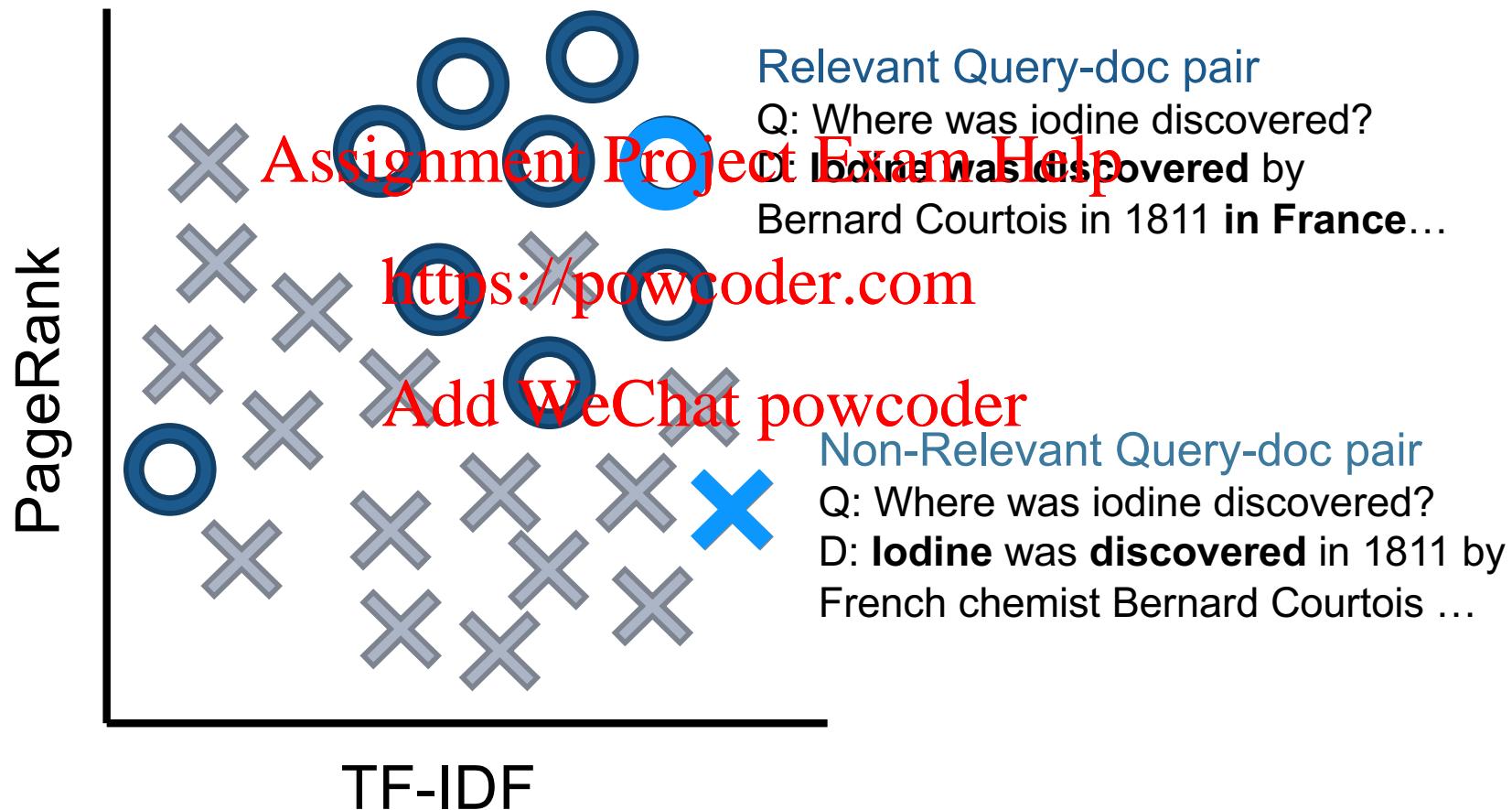
- Today we're focusing on neural networks

These algorithms have ~~the settings (hyperparameters)~~ **Add WeChat powcoder** that affect how the model is built

In reality: these algorithms operate over many features (figure shows 2: x and y axis)



Learning to Rank



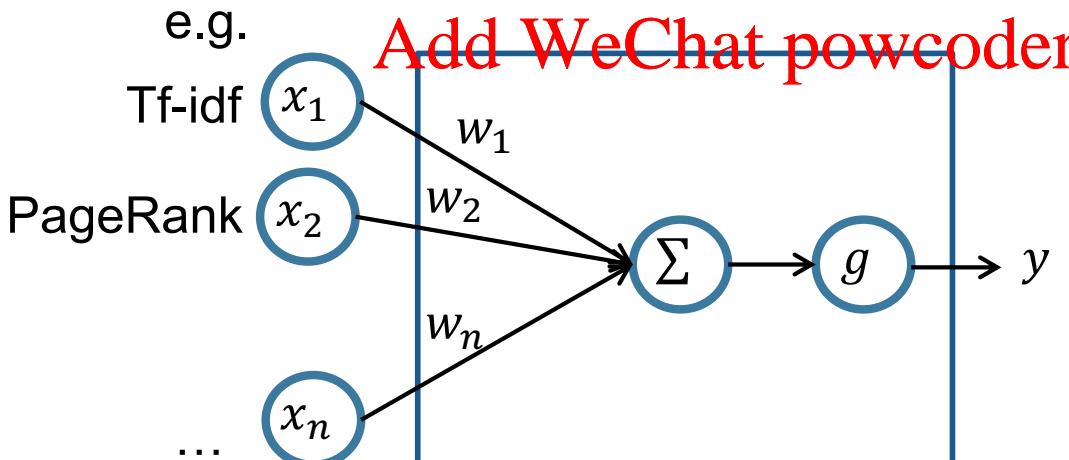
Neural Networks (Intro.)

Consist of computational elements (neurons)

A neuron receives INPUT from other nodes and each INPUT is associated with a weight w (learned)

The unit computes some function f of the weighted sum of its inputs: $y = g(\sum_{i=1}^n w_i x_i)$

Also sometimes referred to as a perceptron
<https://powcoder.com>



g is called the “activation function”. It transforms the OUTPUT shape of the weighted sum.

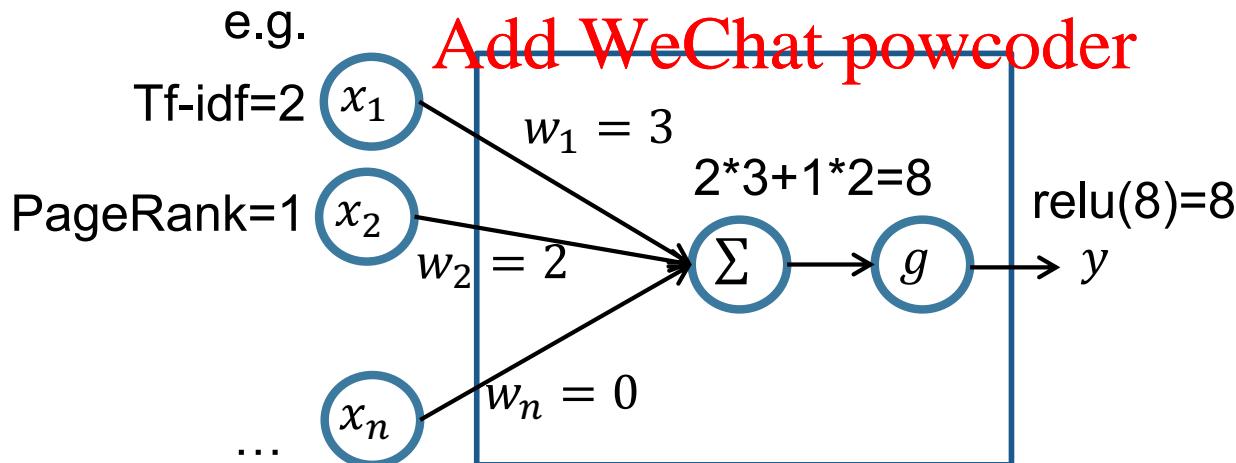
Neural Networks (Intro.)

Consist of computational elements (neurons)

A neuron receives INPUT from other nodes and each INPUT is associated with a weight w (learned)

The unit computes some function f of the weighted sum of its inputs: $y = g(\sum_{i=1}^n w_i x_i)$

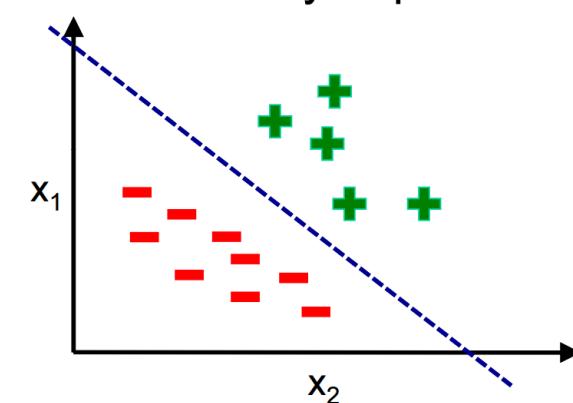
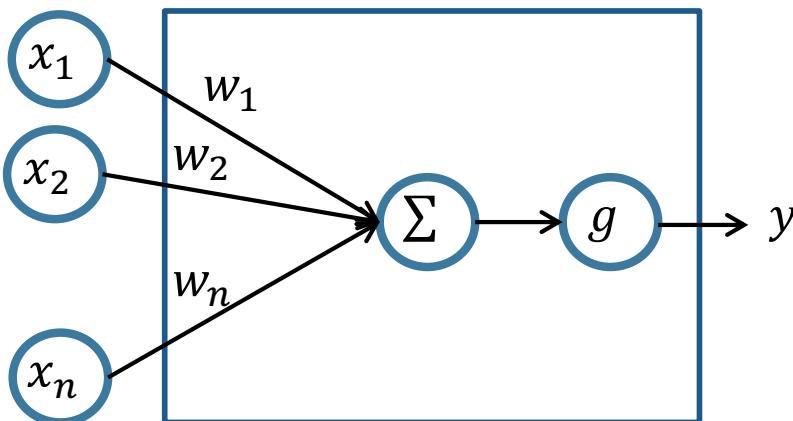
Also sometimes referred to as a perceptron
<https://powcoder.com>



Neural Networks (Intro.)

Training

- Feeding in learning examples $D = (x_i, y_i)$
- The network adjusts the weights based on training samples
Assignment Project Exam Help
- Uses a loss function to determine how “wrong” the predicated value is <https://powcoder.com>
- Propagates loss signals obtained by gradients of OUTPUT with respect to INPUT



Neural Networks (Intro.)

A neural network consists of multiple layers of neurons

Helps in capture non-linear target functions

- Many interesting problems are non-linear in nature!

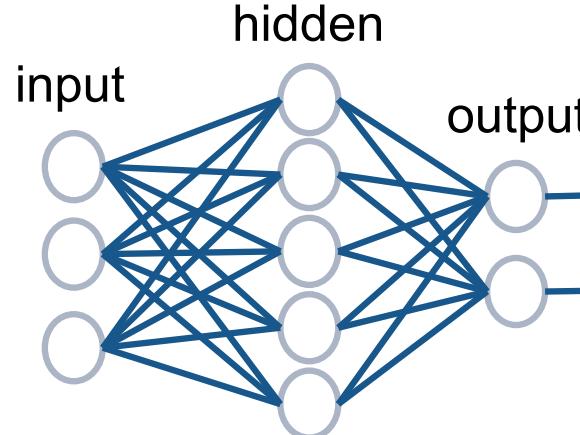
Assignment Project Exam Help

Several special structures allow functionality on certain types of data: e.g., convolutional, recurrent, transformer

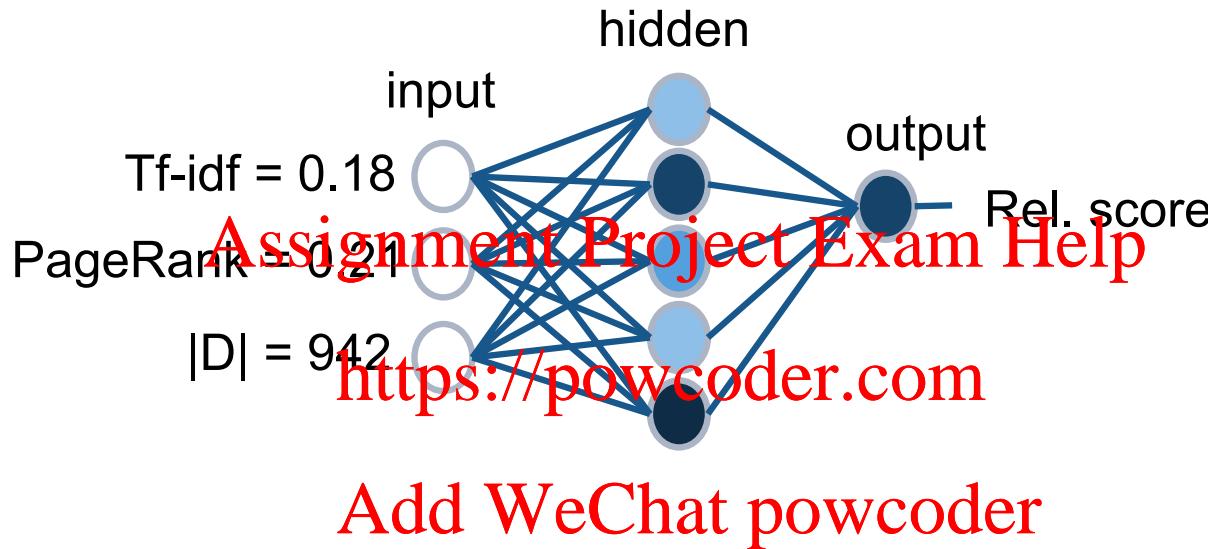
<https://powcoder.com>

This is the basic “feed forward” neural network:

Add WeChat powcoder



Learning to rank with neural networks



This could work, but doesn't provide much beyond what more basic methods can do (e.g., LambdaMART).

Learning to rank with neural networks

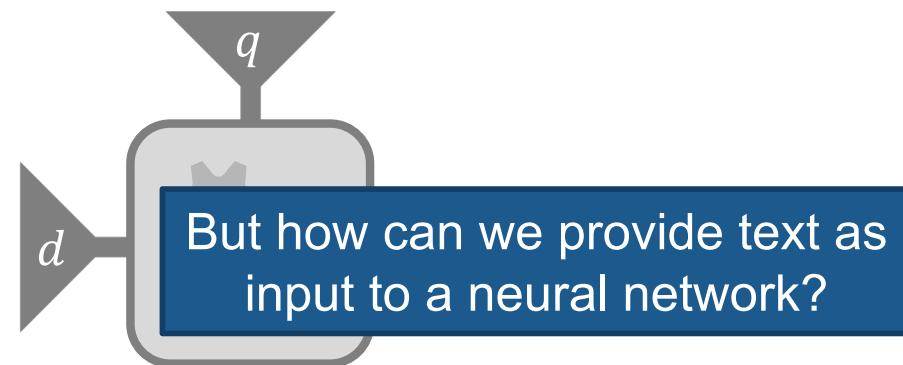
We explore feeding the text of the query and document as features directly into the neural networks.

Goal: Determine relevance scores based on the **query and document text itself**.

Let the neural network learn rules for vocabulary mismatch, proximity, etc.

Q: Where was iodine discovered?

D2: Iodine was discovered by Bernard Courtois in 1811 in France...



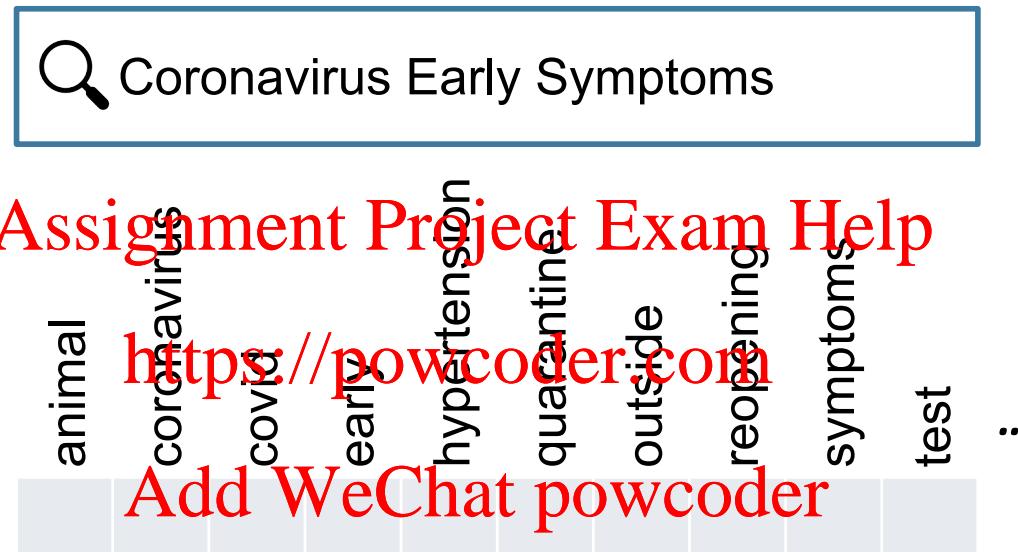
Representing Text

One option: One-hot encoding

Assignment Project Exam Help
animal coronavirus covid early hypertension quarantine outside reopening symptoms test :
<https://powcoder.com>
Add WeChat powcoder

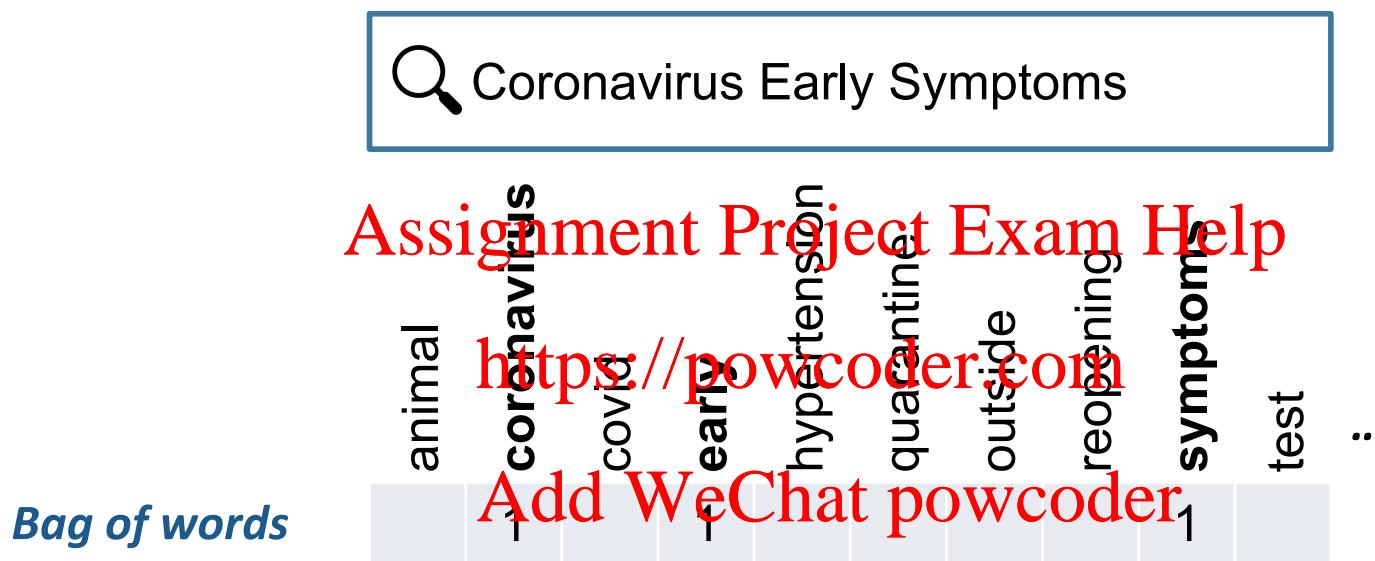
Representing Text

One option: One-hot encoding



Representing Text

One option: One-hot encoding



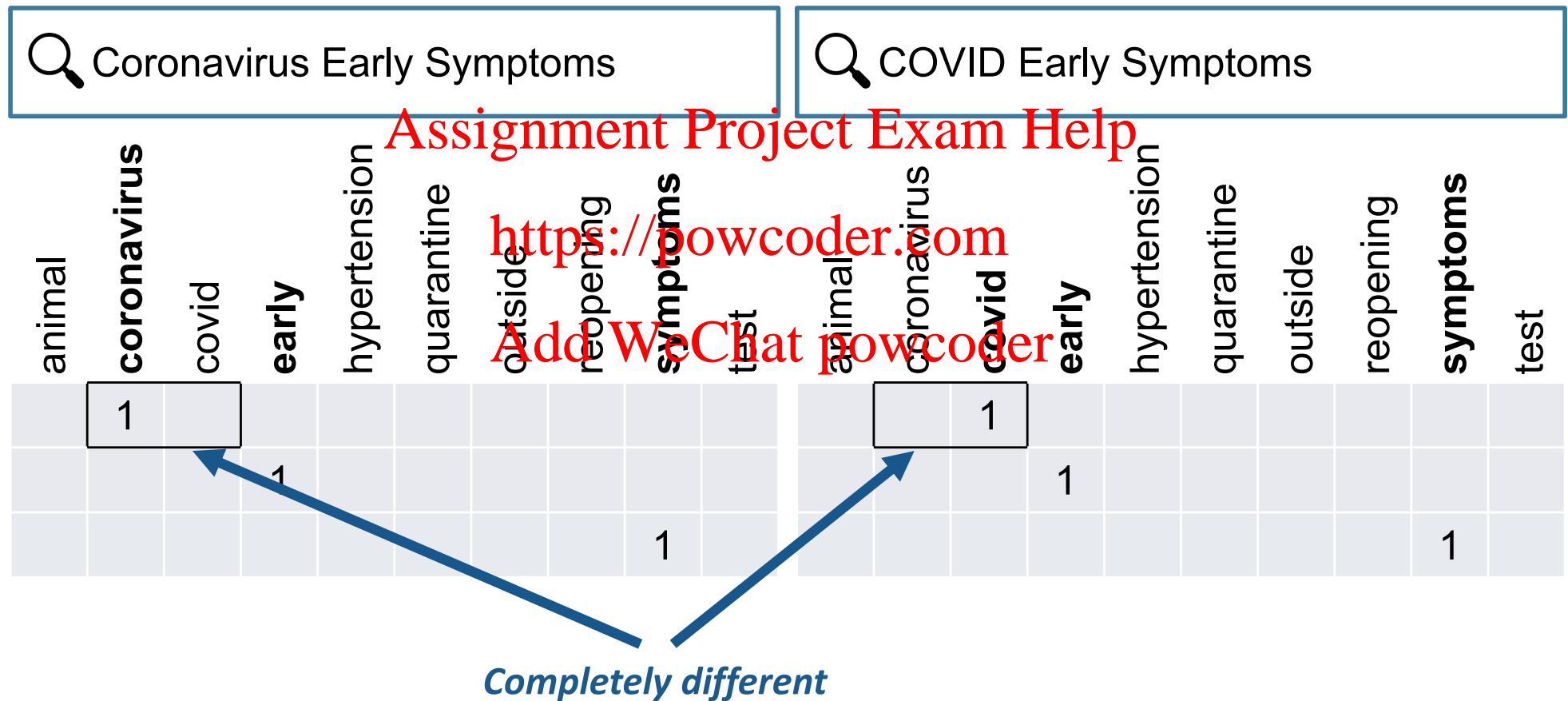
Representing Text

One option: One-hot encoding

Sequence	animal	coronavirus	https://powcoder.com	covid	early	hypertension	quarantine	outside	reopening	symptoms	test	:
				1						1		

Representing Text

Problem: no relationship between words



Representing Text

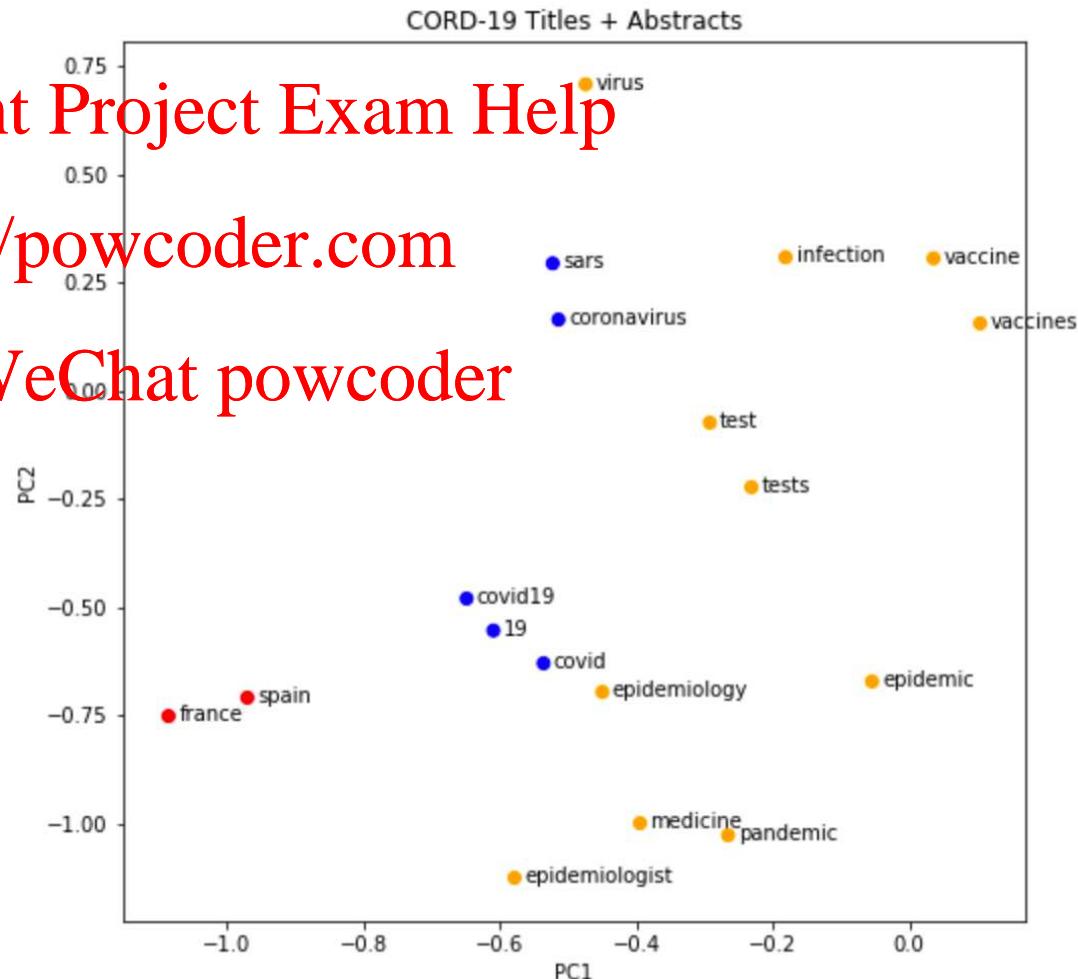
Word Vectors: Map each word to a dense vector

Also known as a word “embedding”

```
covid =  
[0.60586,  
 0.04596,  
 0.12191,  
 -0.18414,  
 -0.04422,  
 0.13495,  
 0.31471,  
 0.33992,  
 0.01285,  
 -0.18592,  
 -0.43352,  
 -0.62741,  
 0.24341,  
 0.07149,  
 ...]
```

```
coronavirus =  
[0.5583,  
 -0.27798,  
 0.08317,  
 -0.19729,  
 -0.49235,  
 0.26514,  
 0.03004,  
 0.25704,  
 -0.38031,  
 -0.32722,  
 -0.47273,  
 -0.01596,  
 0.32322,  
 -0.04947,  
 ...]
```

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder



Representing Text

Building word vectors

- Based on word co-occurrences
- Trained using neural network
- e.g. word2vec, glove, etc.
- Important: these vectors are the same, regardless of context that the word appears in (i.e., “static”).

Assignment Project Exam Help



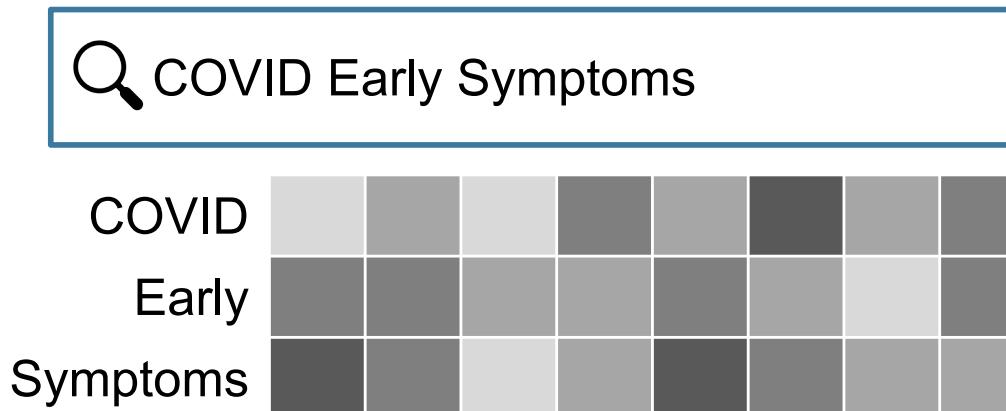
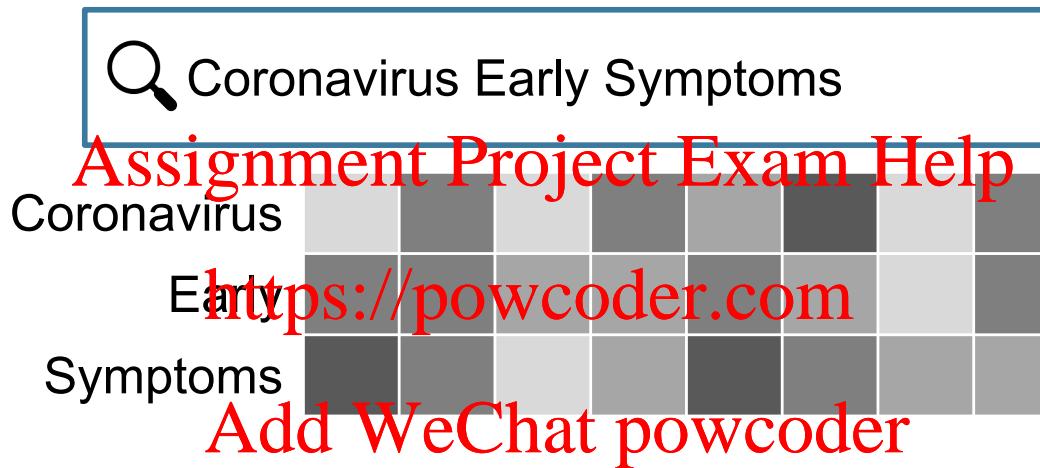
<https://powcoder.com>

Add WeChat to powcoder

```
covid = [0.60586, 0.04596, 0.12191, -0.18414, -0.04422, 0.13495, 0.31471, 0.33992, 0.01285, -0.18592, -0.43352, -0.62741, 0.24341, 0.07149, ...]  
coronavirus = [0.33853, -0.27798, 0.08317, -0.19729, -0.49235, 0.26514, 0.03004, 0.25704, -0.38031, -0.32722, -0.47273, -0.01596, 0.32322, -0.04947, ...]
```

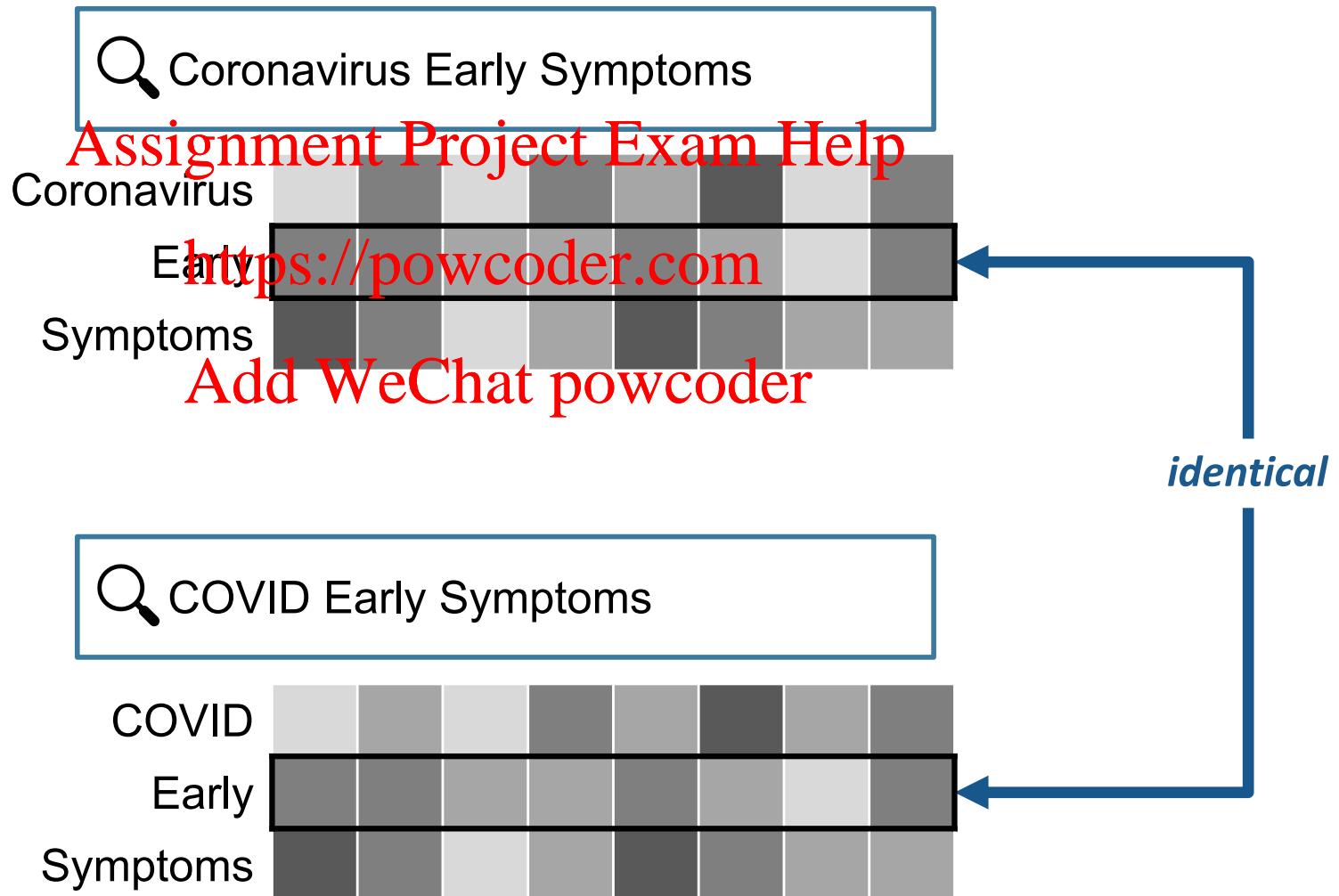
Representing Text

Word Vectors: Map each word to a dense vector



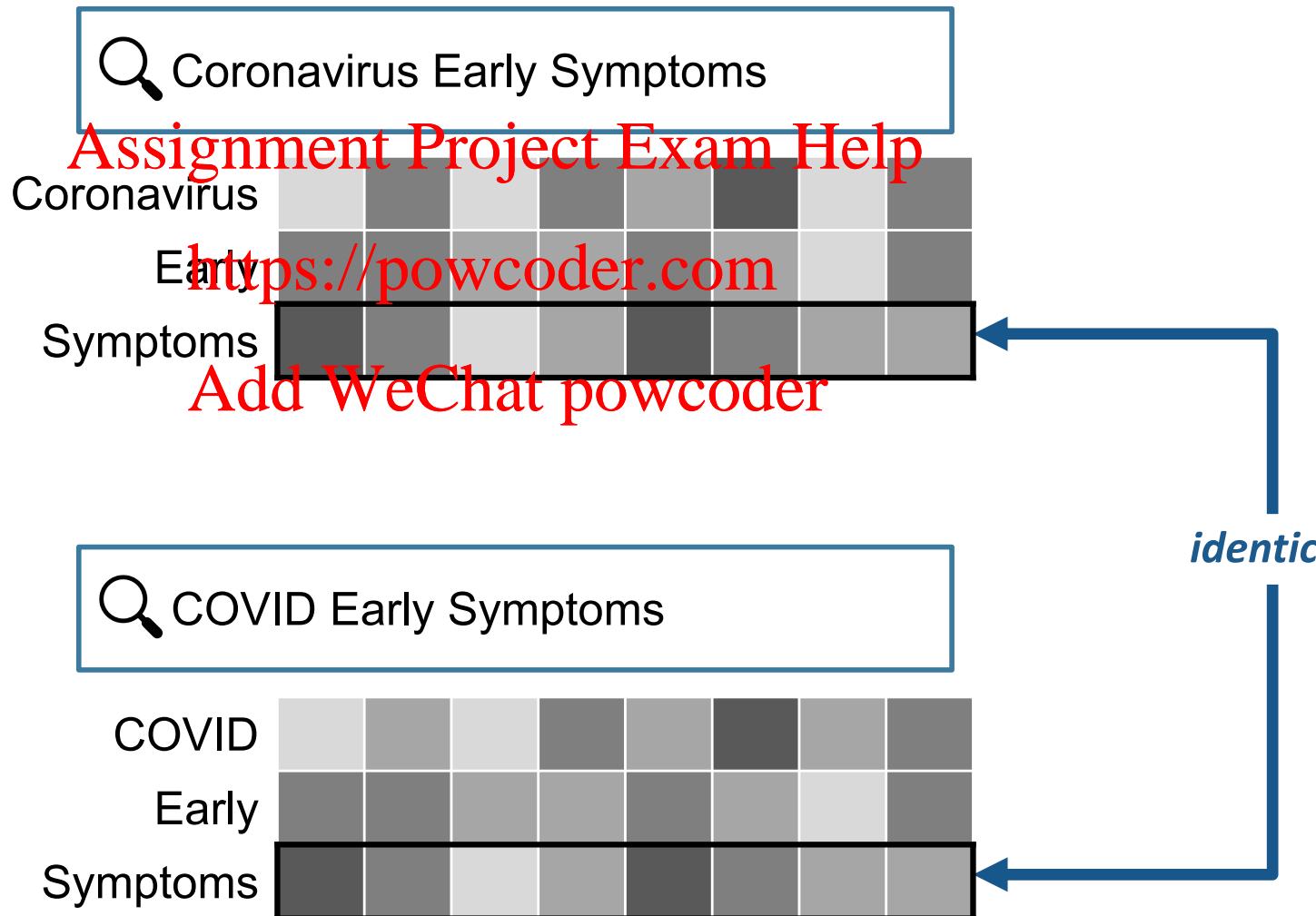
Representing Text

Word Vectors: Map each word to a dense vector



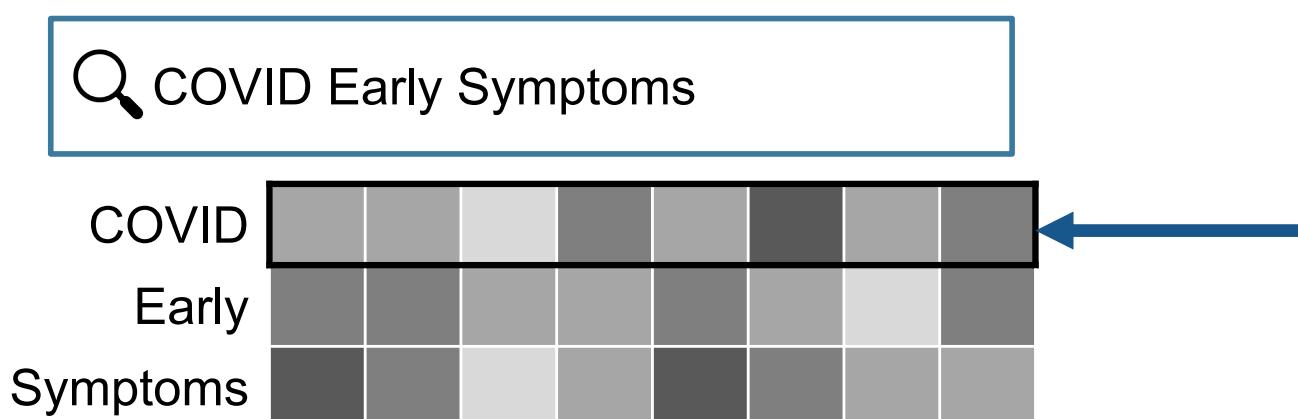
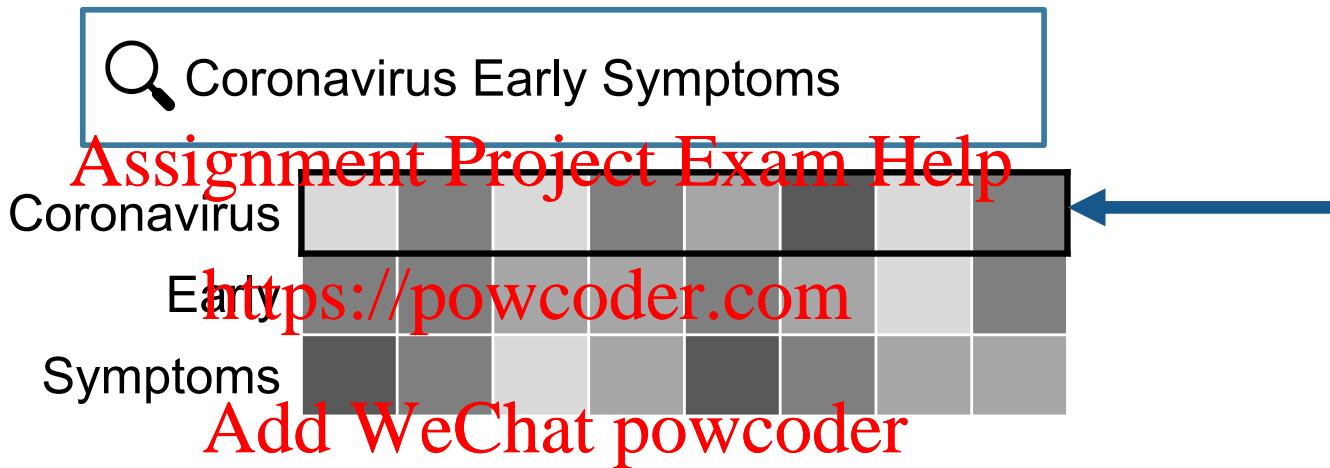
Representing Text

Word Vectors: Map each word to a dense vector



Representing Text

Word Vectors: Map each word to a dense vector



*Not identical
vectors, but close*

Problem: Not Context-Aware

- ✓ Handles different tokens with similar meanings

“Coronavirus” has a similar vector than “COVID”



Assignment Project Exam Help

- ✗ Doesn't handle a single token with multiple possible meanings

<https://powcoder.com>

“A bear is raiding homes in California.”



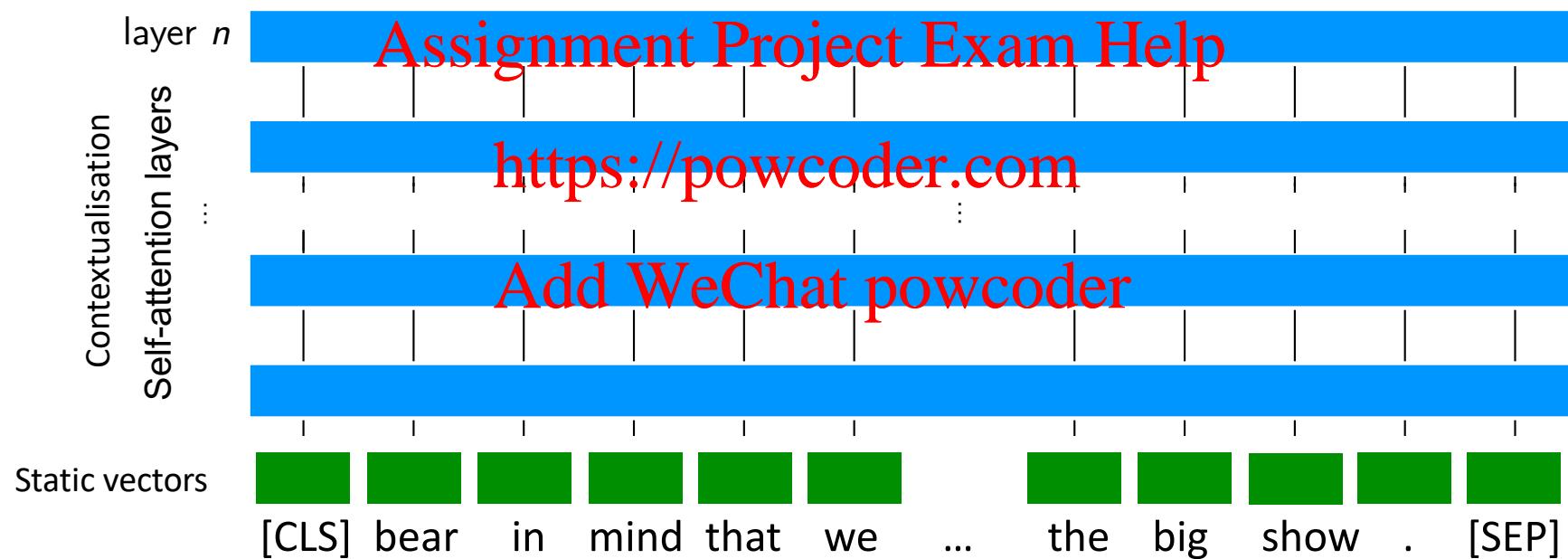
Add WeChat powcoder

has the same vector than

“Bear in mind that we do not have much time before the big show.”

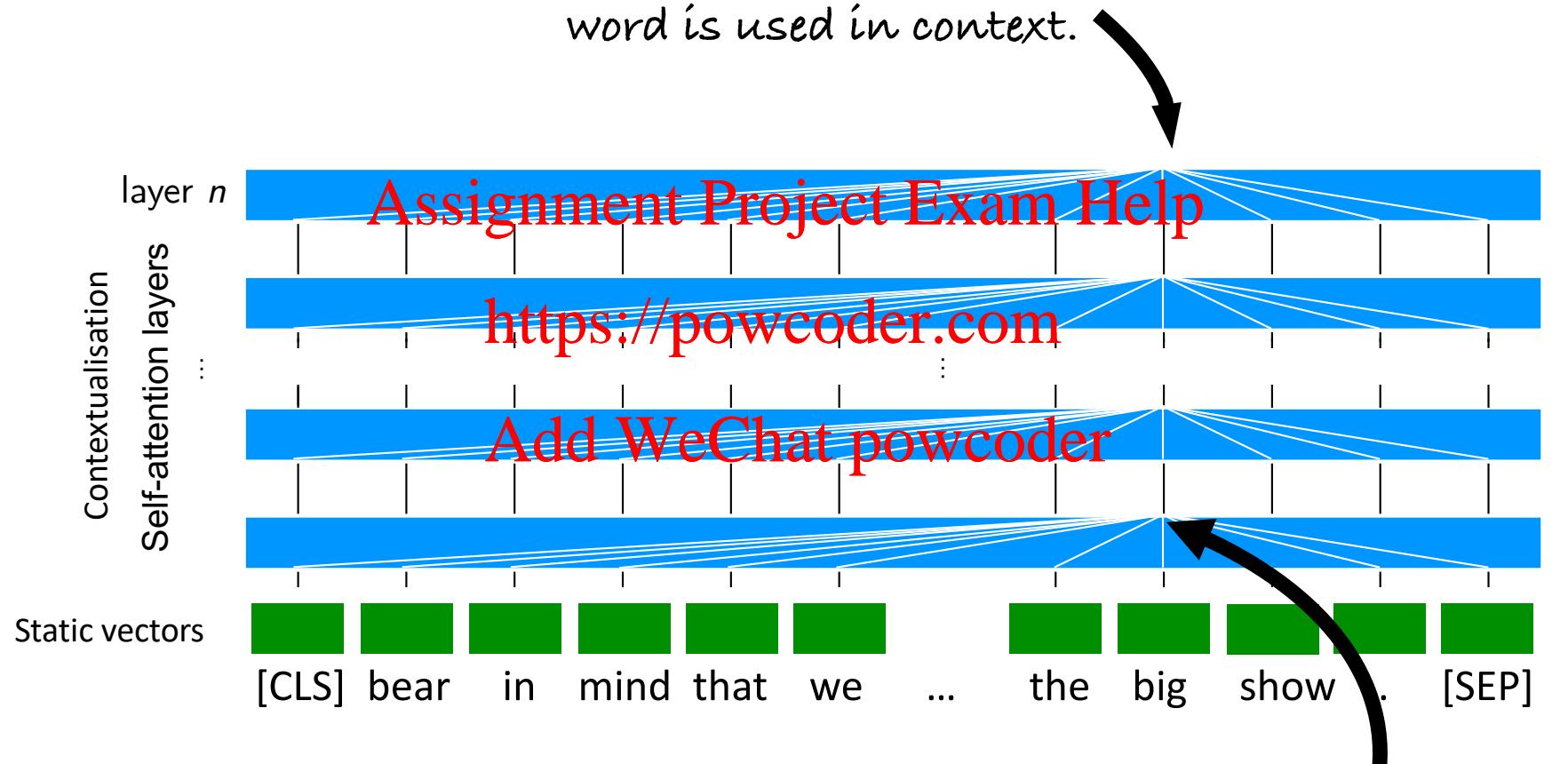


BERT & Transformer Networks



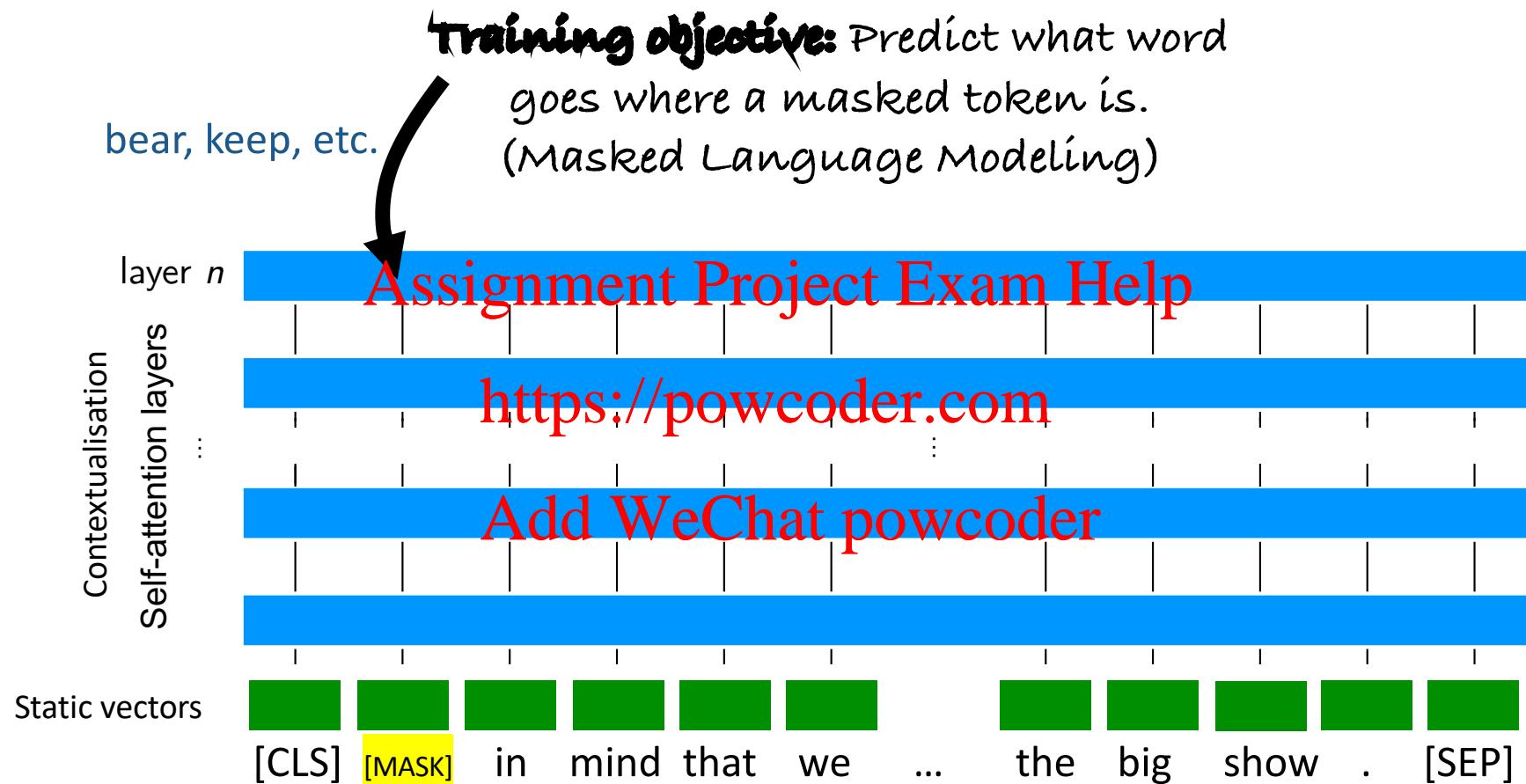
BERT & Transformer Networks

By the end, a “contextualised” vector is produced – one that knows how the word is used in context.

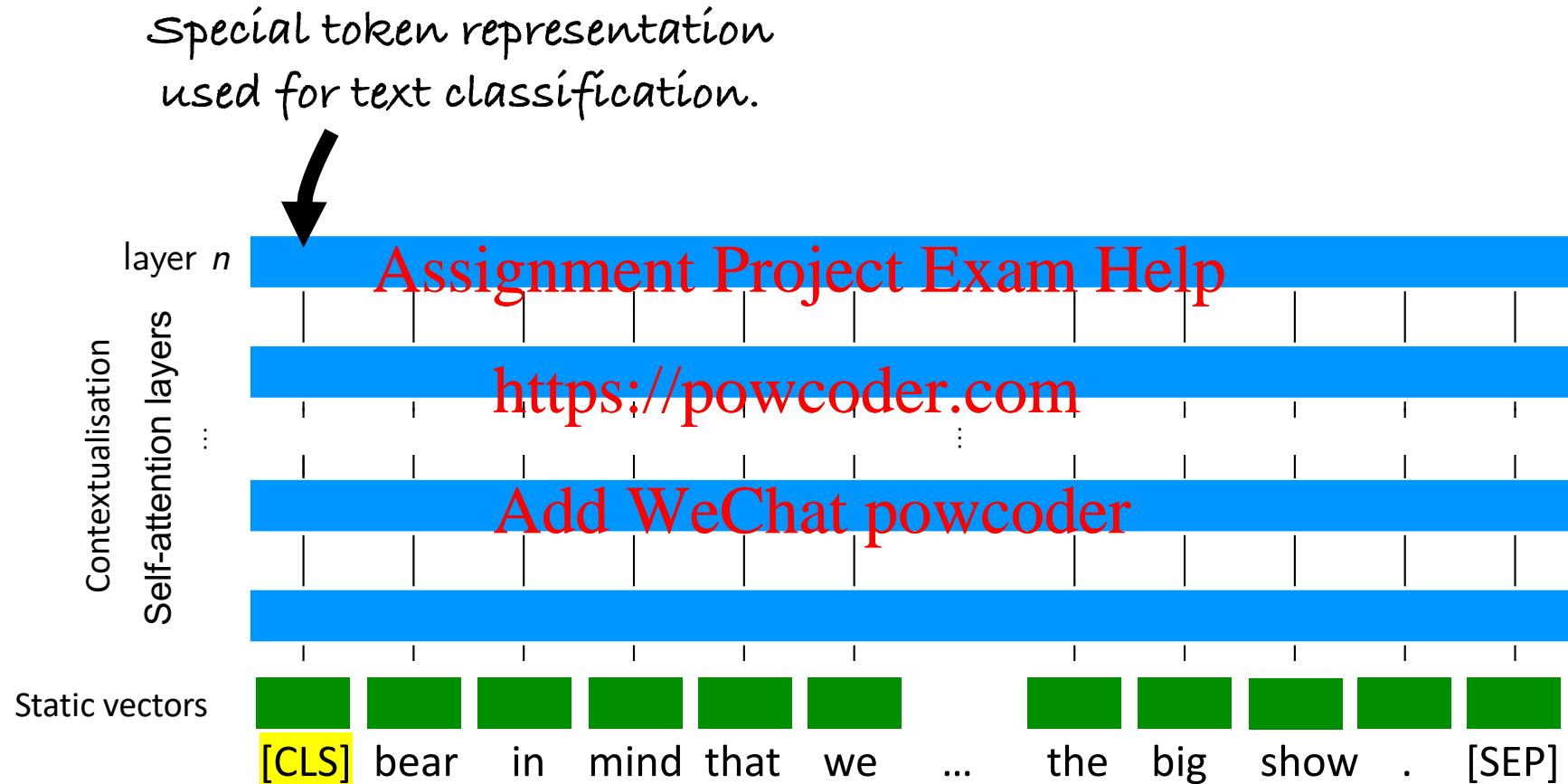


At each layer, a token's vector is a learned combination of all the other vectors in the text.

BERT & Transformer Networks



BERT & Transformer Networks



Problem: Not Context-Aware

- ✓ Handles different tokens with similar meanings

“Coronavirus” has a similar vector than “COVID”



Assignment Project Exam Help

- ✓ Doesn't handle a single token with multiple possible meanings

<https://powcoder.com>

“A bear is raiding homes in California.”



Add WeChat powcoder

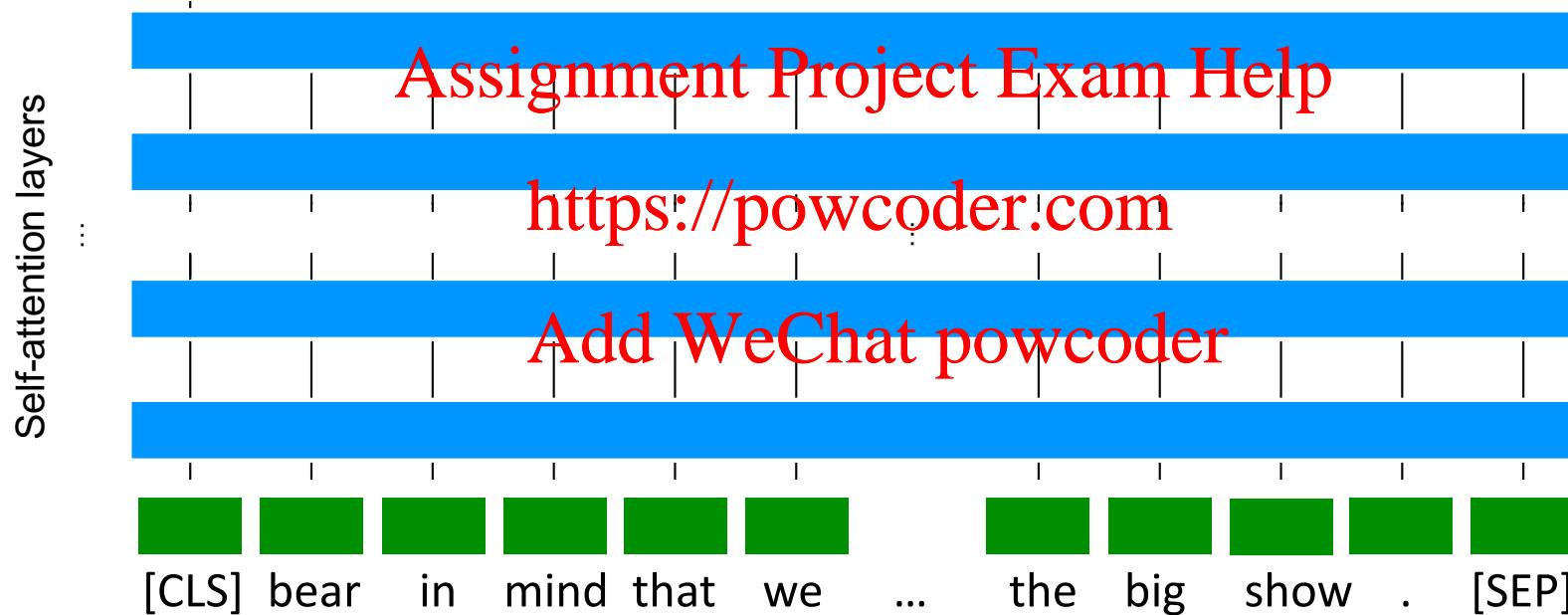
has a different vector than

“Bear in mind that we do not have much time before the big show.”



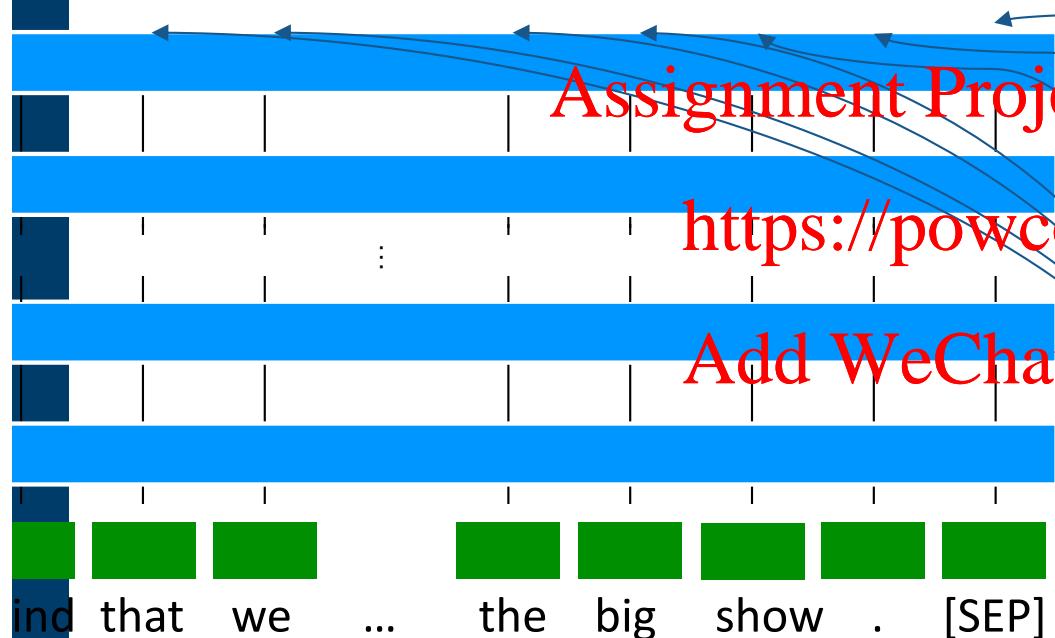
T5 & Text Generation

Transformer “Encoder”



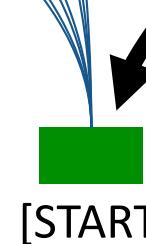
T5 & Text Generation

Transformer “Encoder”



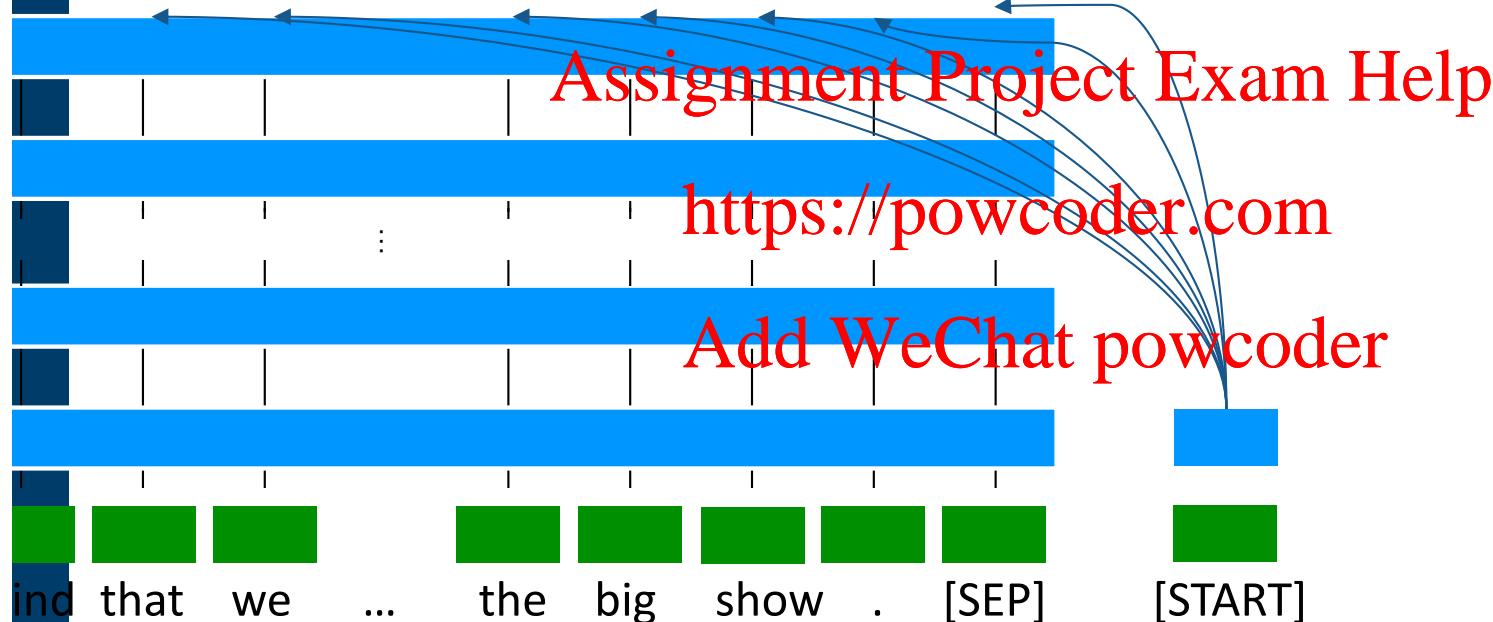
Transformer “Decoder”

Build a new representation based on the encoder outputs.



T5 & Text Generation

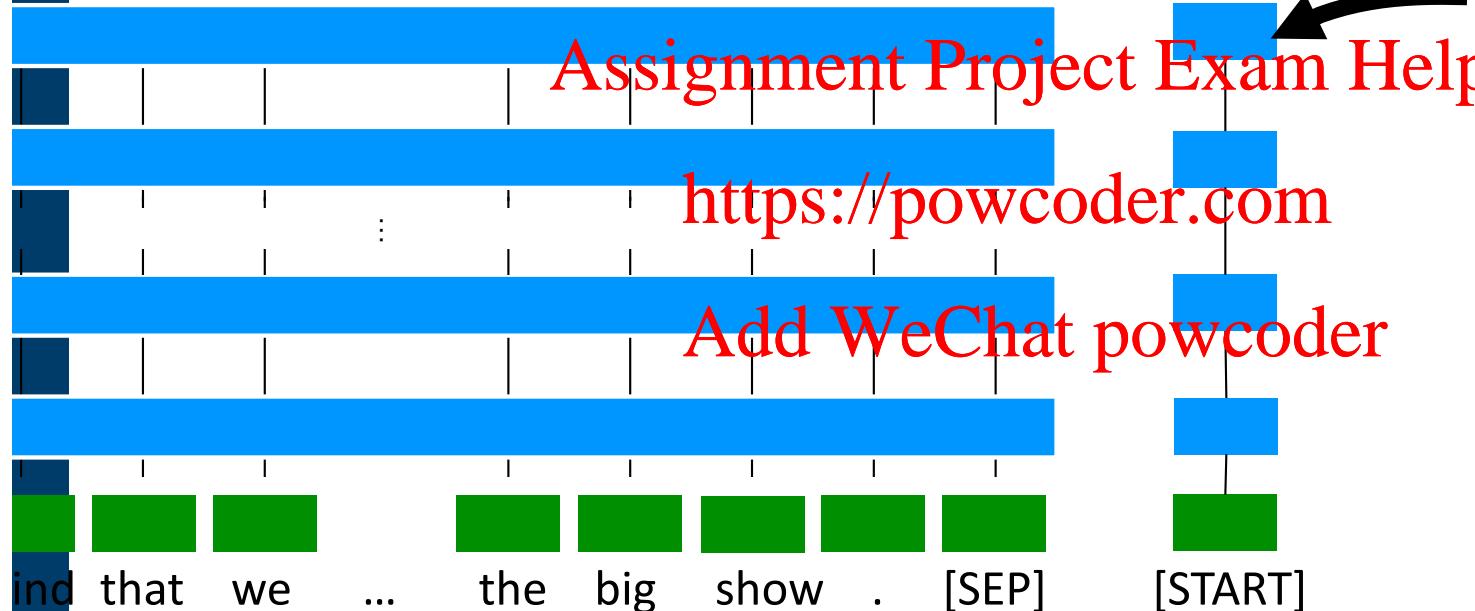
Transformer “Encoder”



Transformer “Decoder”

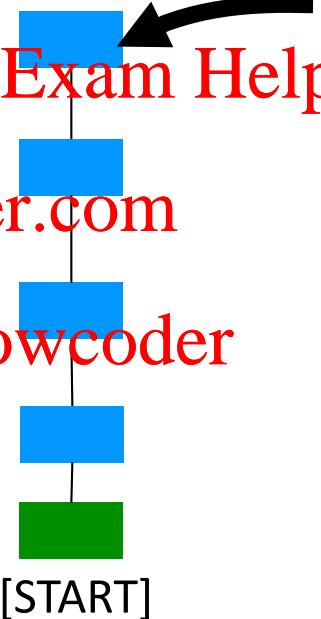
T5 & Text Generation

Transformer “Encoder”



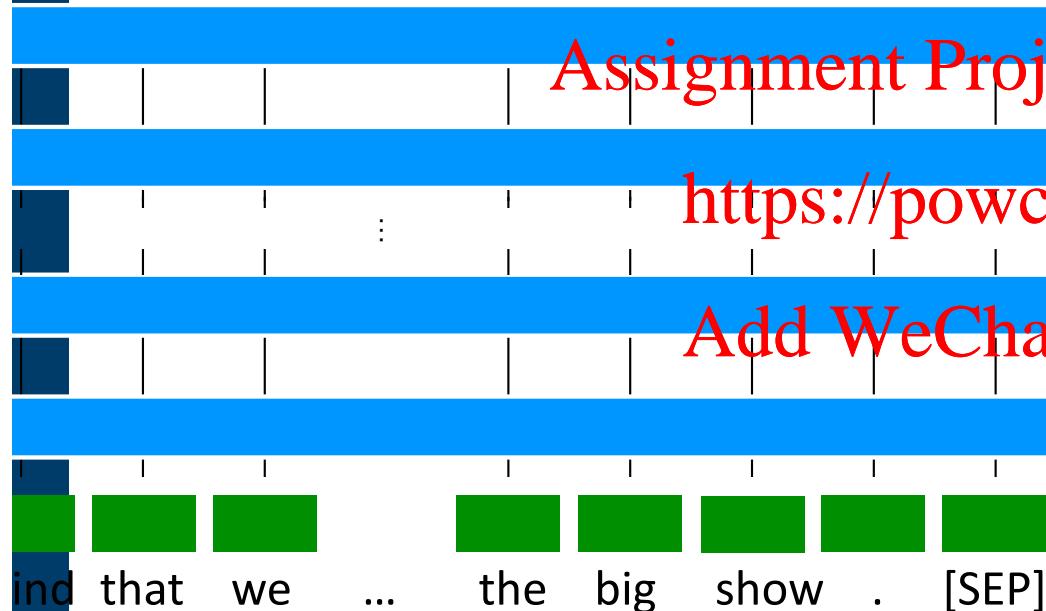
Transformer “Decoder”

It
Predict the next word
in the sequence.

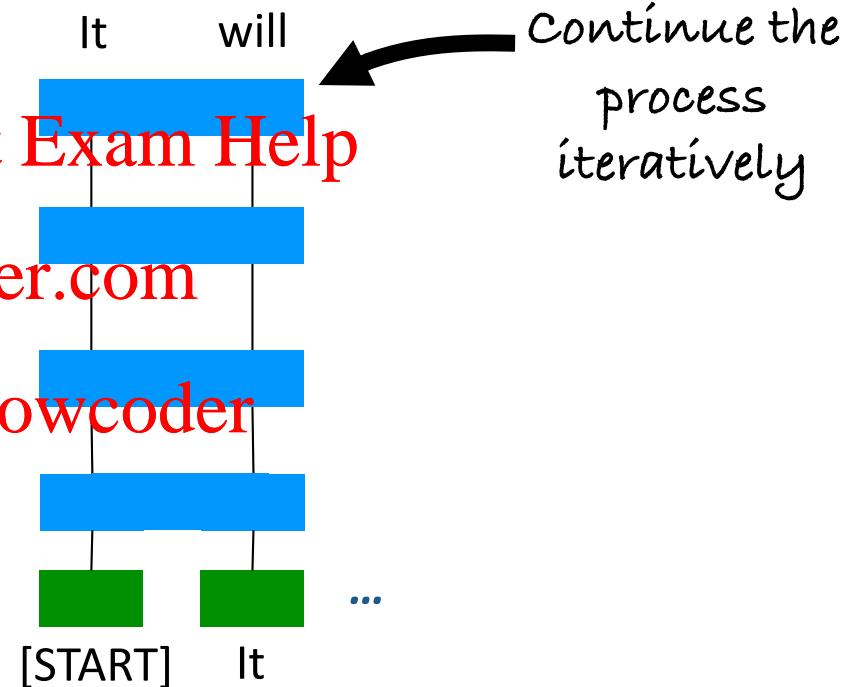


T5 & Text Generation

Transformer “Encoder”



Transformer “Decoder”



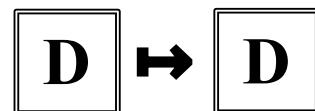
Summary

Neural NLP techniques can:

- Build representations that place similar words near each other
Assignment Project Exam Help
- Build representations that can distinguish different meanings a single word may have
https://powcoder.com
- Generate text sequences
Add WeChat powcoder

We will now use these tools to perform IR operations!

Today

1. Review of LTR & Basics of Neural Networks for NLP
2. Neural Re-ranking  Assignment Project Exam Help
3. Neural Retrieval 
4. Neural Query Rewriting & PRF 
5. Neural Document Rewriting 
6. Neural IR in PyTerrier

Neural Re-Ranking

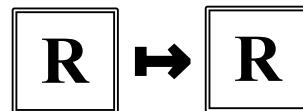
Q → R → R
Assignment Project Exam Help



Add WeChat powcoder

Why Neural Re-Ranking?

Simple formulation: Given a query and document, assign a new score.



qid	query	docno	score	qid	query	docno	score
0	glasgow weather	15213	0.5134	0	glasgow weather	26340	0.9312
0	glasgow weather	42613	0.3742	0	glasgow weather	15213	0.1151
0	glasgow weather	26340	0.3223	0	glasgow weather	52363	0.0215
...

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Computational cost: Lexical retrieval method like BM25 are simple and achieve reasonably recall already; you just need to re-score a set of 100-1000 query-document pairs using an expensive NN.

Test Collection Suitability: Most benchmarks are based on pooling of lexical systems like BM25, so there is a higher likelihood that documents in this set have relevance assessments.

Formulation

$(query, document) \rightarrow score$



Document:

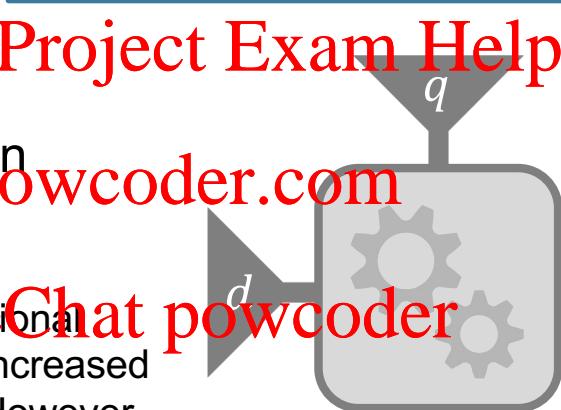
Assignment Project Exam Help

Title: How can we evaluate an interrelation of symptoms?

Abstract: A pandemic of 2019 novel coronavirus (COVID-19) is an international problem and factors associated with increased risk of mortality have been reported. However, there exists limited statistical method to estimate a comprehensive risk for a case in which a patient has several characteristics...

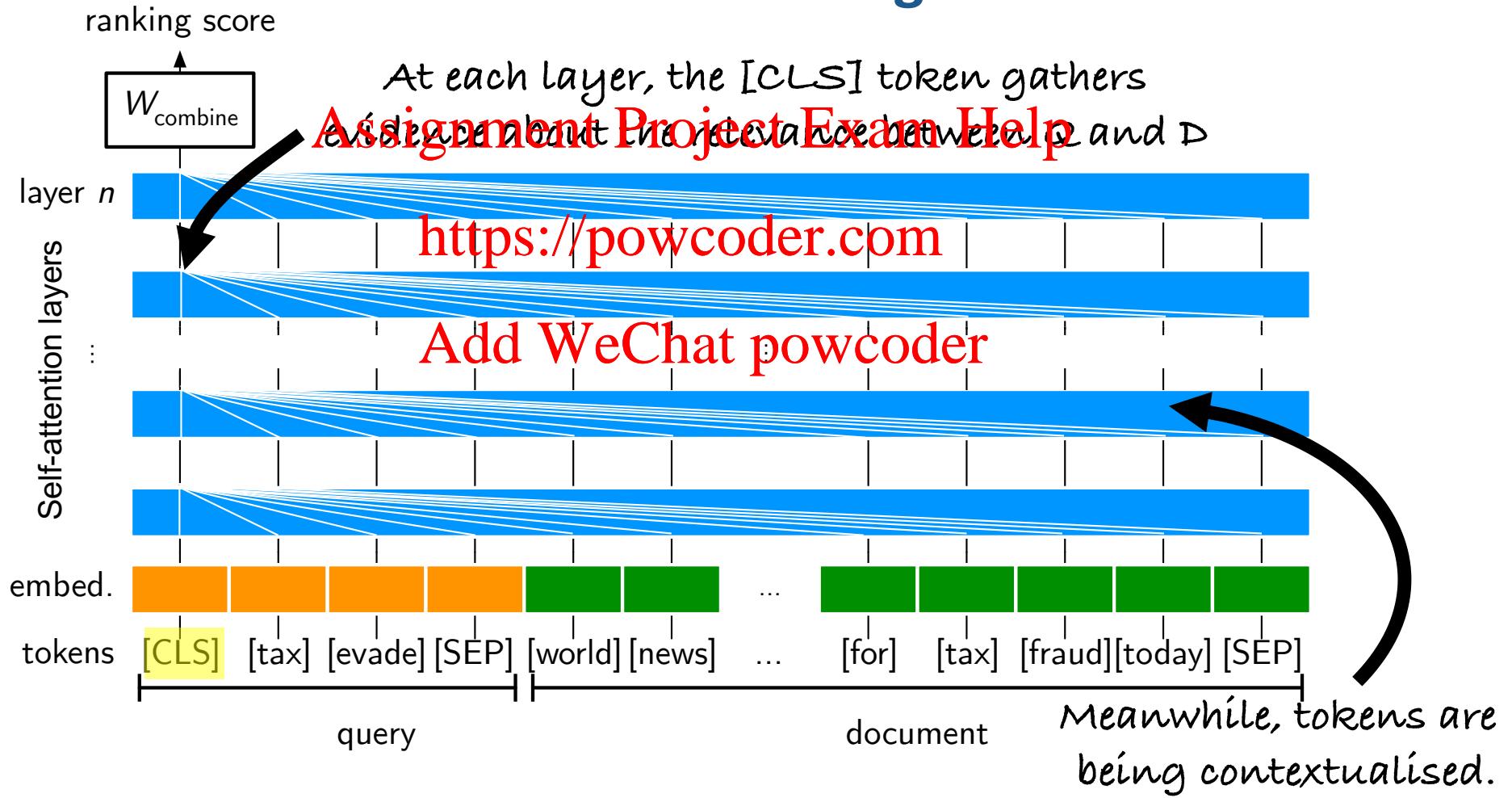
<https://powcoder.com>

Add WeChat powcoder



“Vanilla BERT” – The simplest approach that works (really well)!

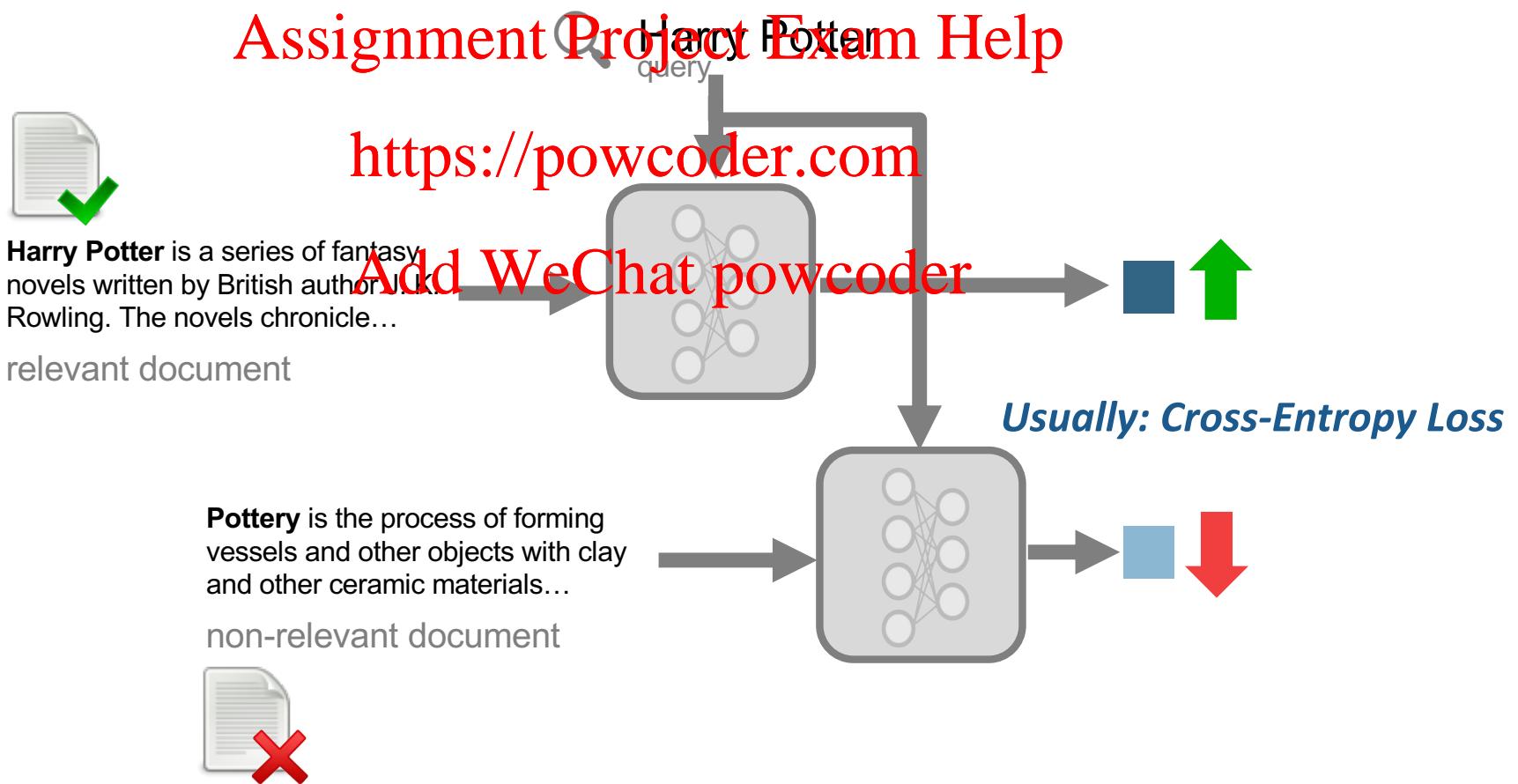
Idea: Concatenate query and document; let the model learn how to combine into a ranking score.



Training

Usually: Pairwise training.

With a query, a relevant doc and non-relevant doc, maximize the relative score between them.



Handling Long Documents

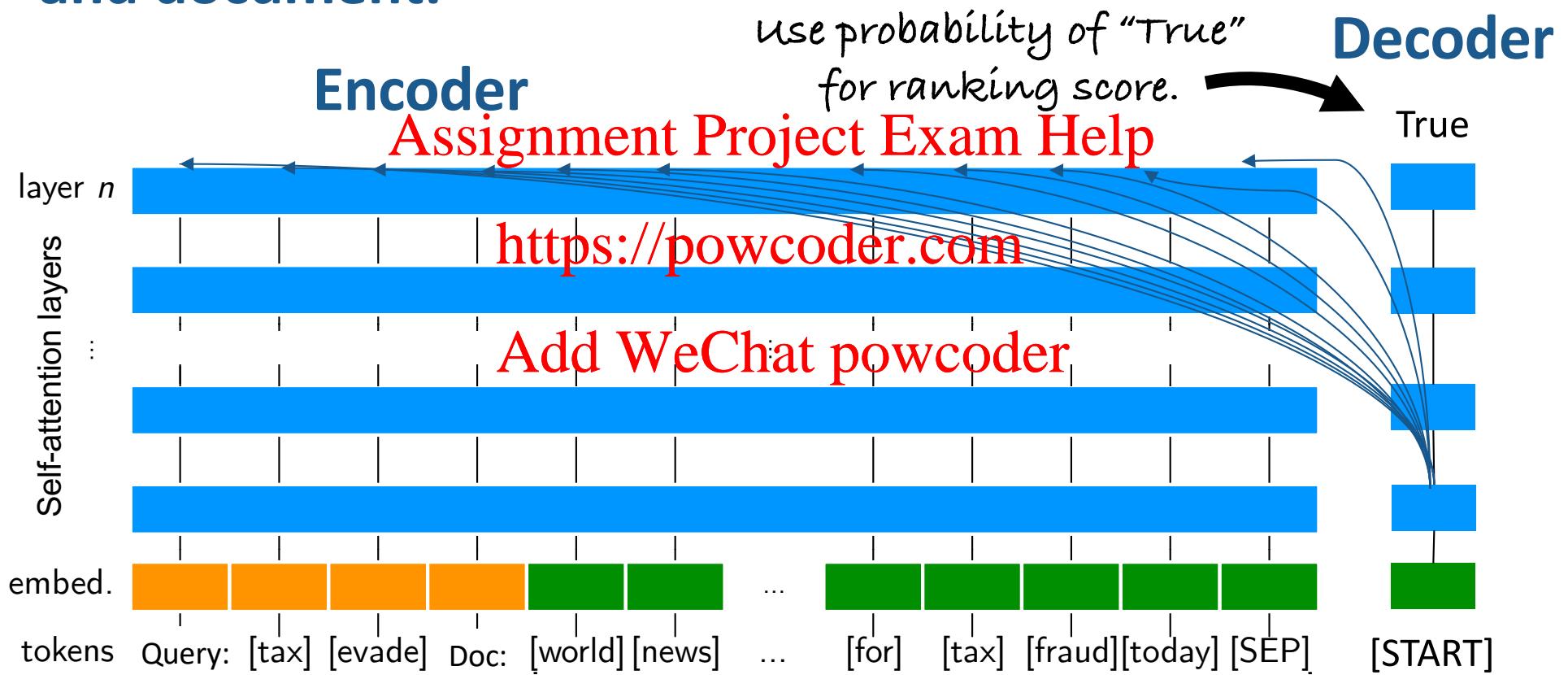
The cost of models like BERT are quadratic over the number of tokens. There is also an inherit maximum length (usually 512).

Strategies for long documents:

- **FirstP** – Only take the first 512 tokens
- **MaxP** – Split the document into passages; take the maximum score over all the passages
- **MeanCLS** – Take the average over [CLS] representations of each passage before computing a relevance score
- **PARADE** – Use another transformer to combine the [CLS] representations of each passage

Re-Ranking with a Text Generation Model (e.g., monoT5)

Idea: Generate “True” or “False”, prompted by query and document.



Approaches for Re-Ranking

Technique	Example
Use [CLS] Token	Vanilla BERT [1]
Generate True/False Token	monoT5 [2] https://powcoder.com

Add WeChat powcoder

[1] MacAvaney et al. CEDR: Contextualized Embeddings for Document Ranking. SIGIR 2019.

[2] Nogueira et al. Document Ranking with a Pretrained Sequence-to-Sequence Model. xxiv 2020.

Approaches for Re-Ranking

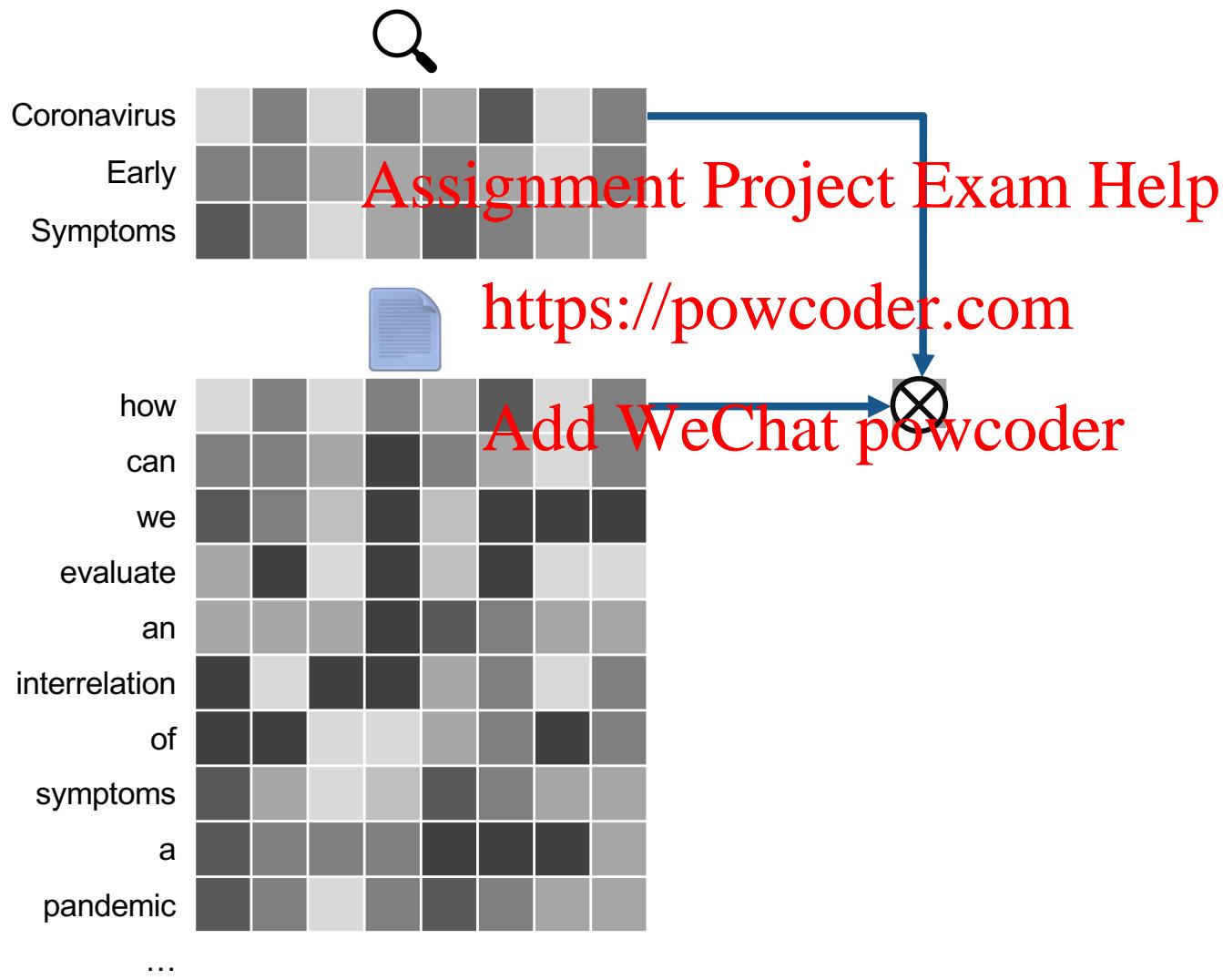
Technique	Example
Use [CLS] Token	Vanilla BERT [1]
Generate True/False Token	monoT5 [2]
Compare token representations	CEDR [1]

Add WeChat powcoder

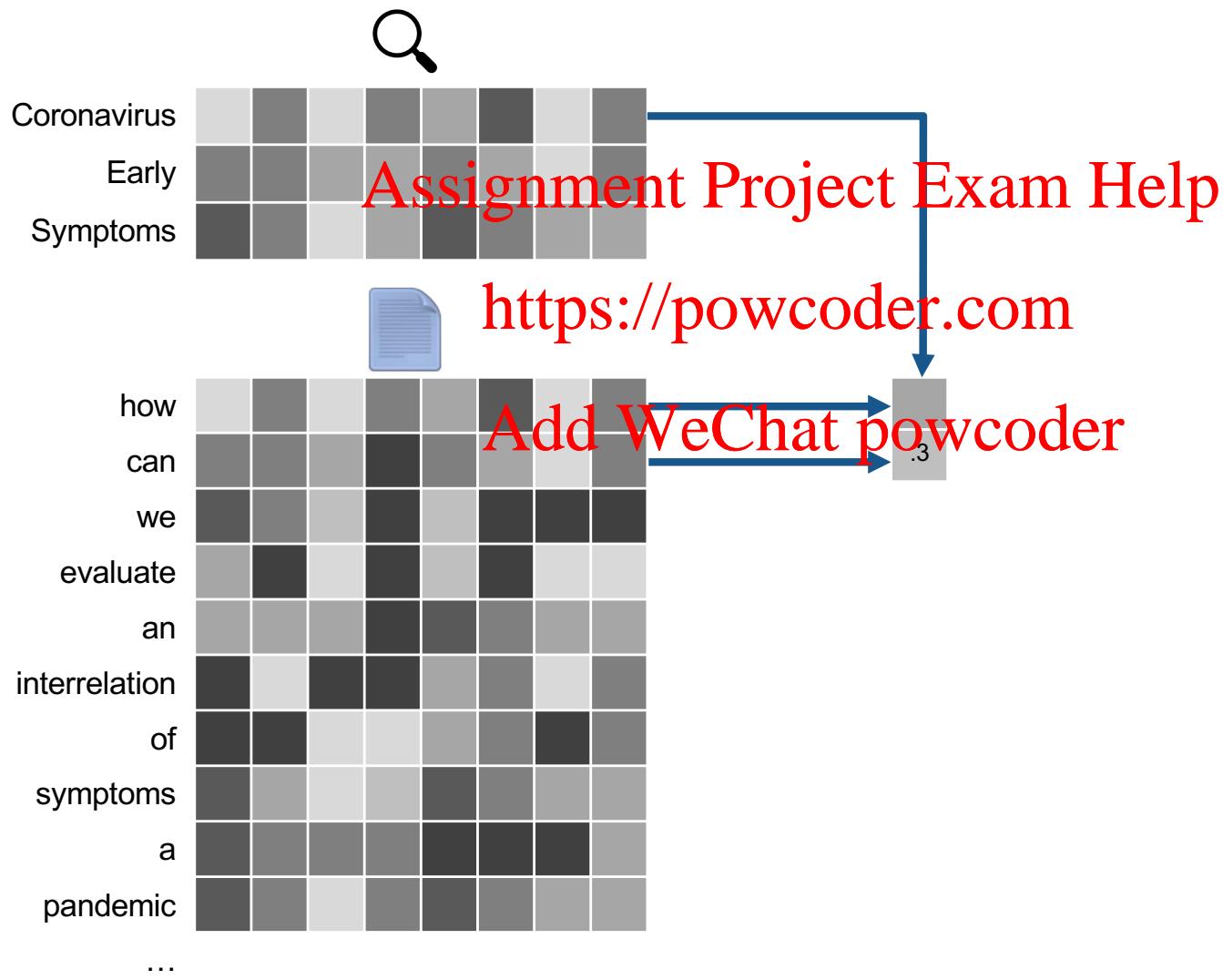
[1] MacAvaney et al. CEDR: Contextualized Embeddings for Document Ranking. SIGIR 2019.

[2] Nogueira et al. Document Ranking with a Pretrained Sequence-to-Sequence Model. xxiv 2020.

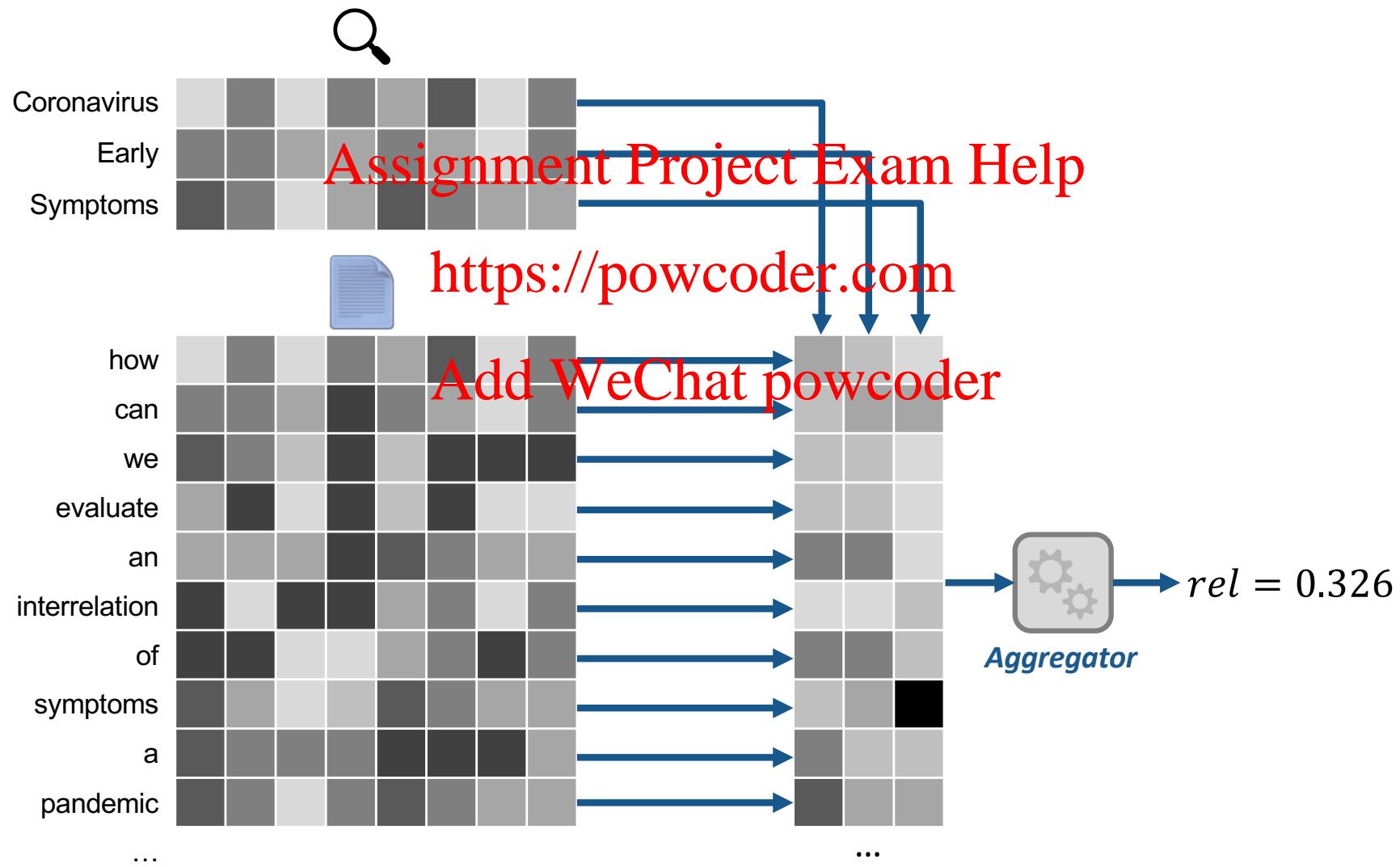
Comparing Token Representations



Comparing Token Representations

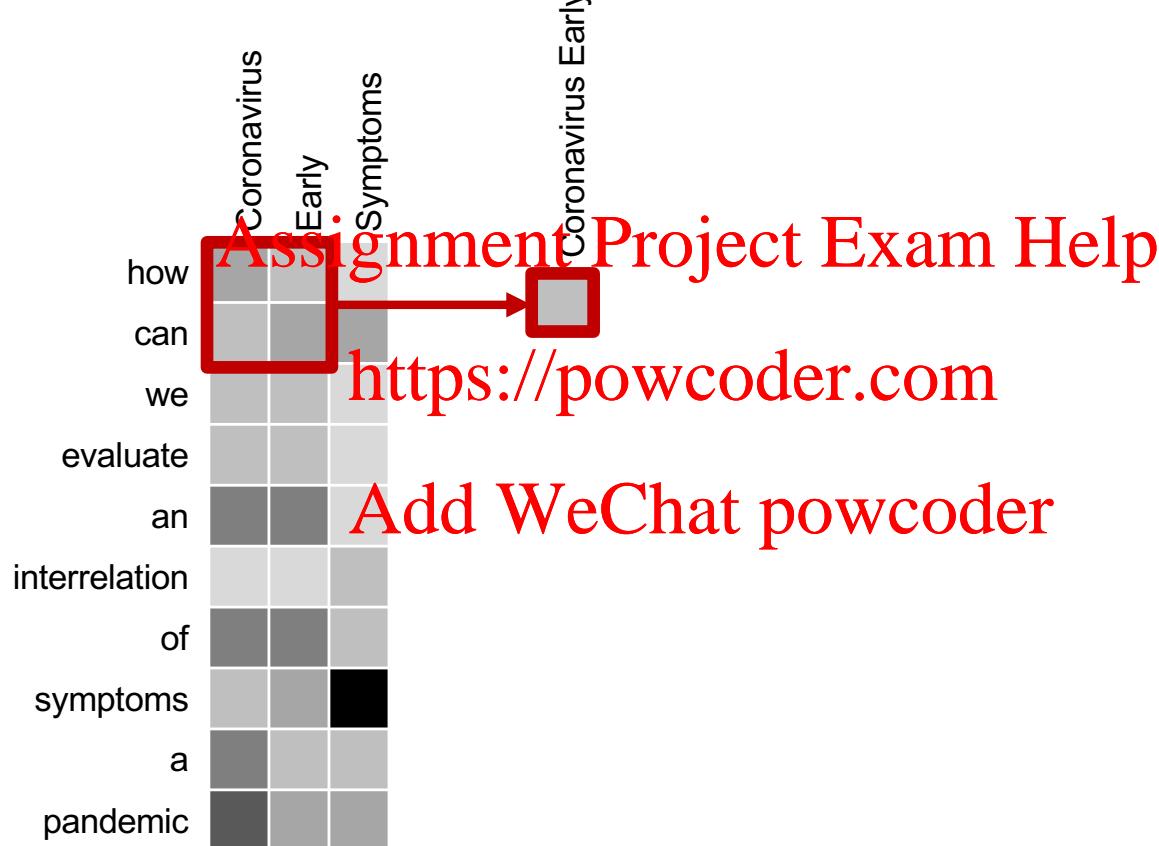


Comparing Token Representations



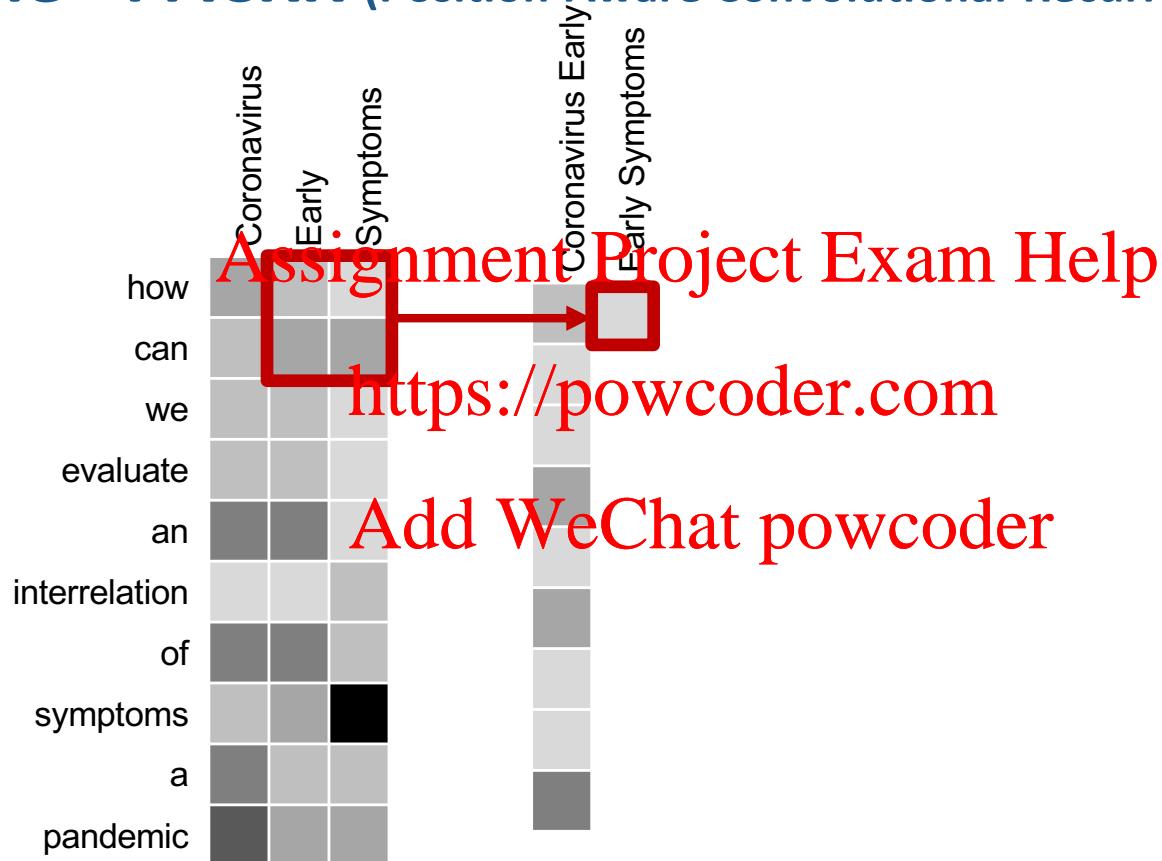
Interaction Aggregation

Example - PACRR (Position-Aware Convolutional-Recurrent Relevance)



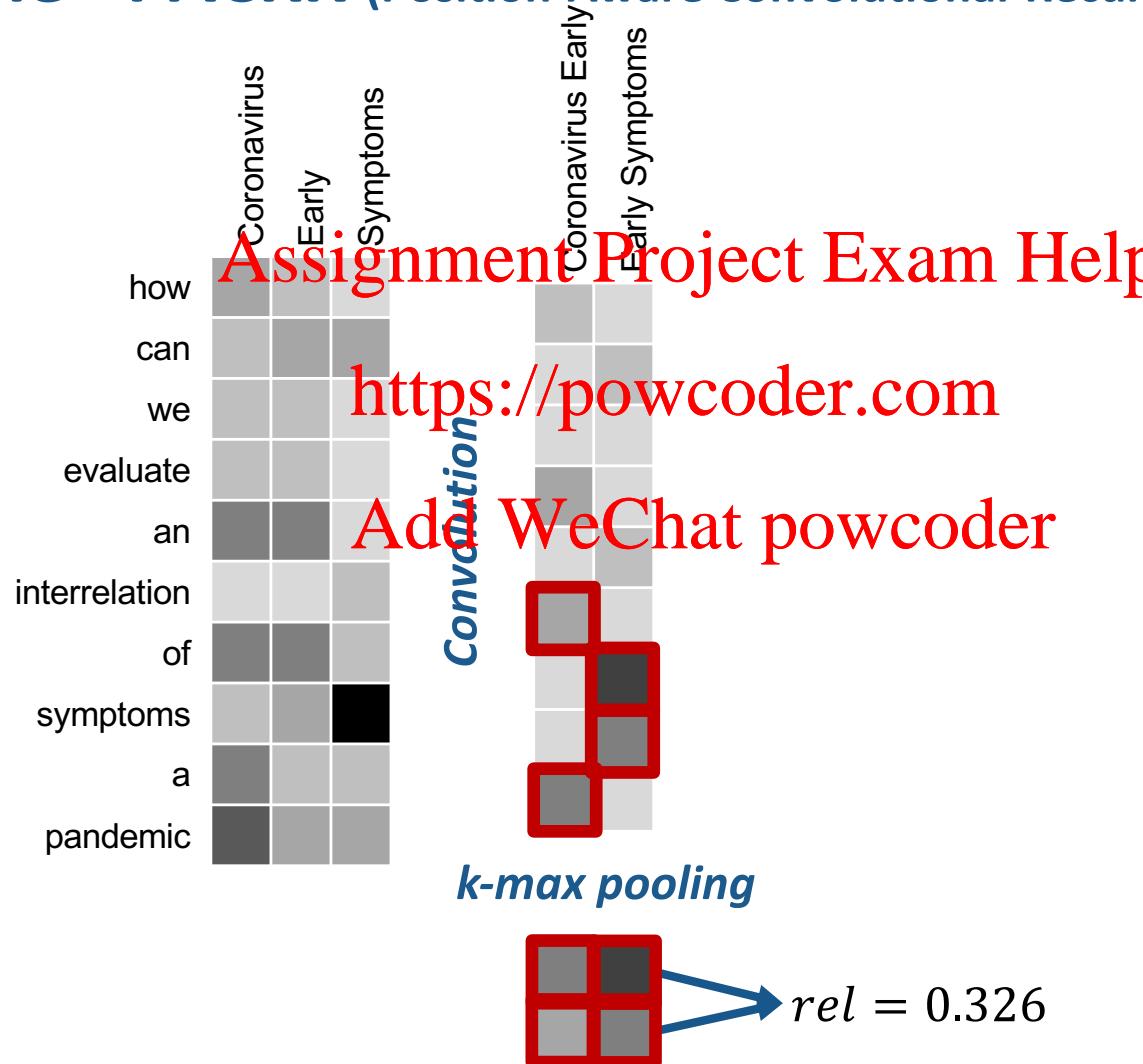
Interaction Aggregation

Example - PACRR (Position-Aware Convolutional-Recurrent Relevance)



Interaction Aggregation

Example - PACRR (Position-Aware Convolutional-Recurrent Relevance)



Problem: Efficiency

Using methods like BERT and T5 for ranking is effective, but also very slow.

E.g., BERT takes 45x longer than BM25: (1.7 seconds)

<https://powcoder.com>

	name	map	ndcg	ndcg_cut_10	mrt
0	BM25	0.075880	0.177728	0.644374	38.557969
1	BM25 >> BERT	0.085371	0.185821	0.740331	1748.576758

Add WeChat powcoder

Approaches for Re-Ranking

Technique	Example
Use [CLS] Token	Vanilla BERT [1]
Generate True/False Token	monoT5 [2]
Compare token representations	CEDR [1]
Predict token importance	EPIC [3]

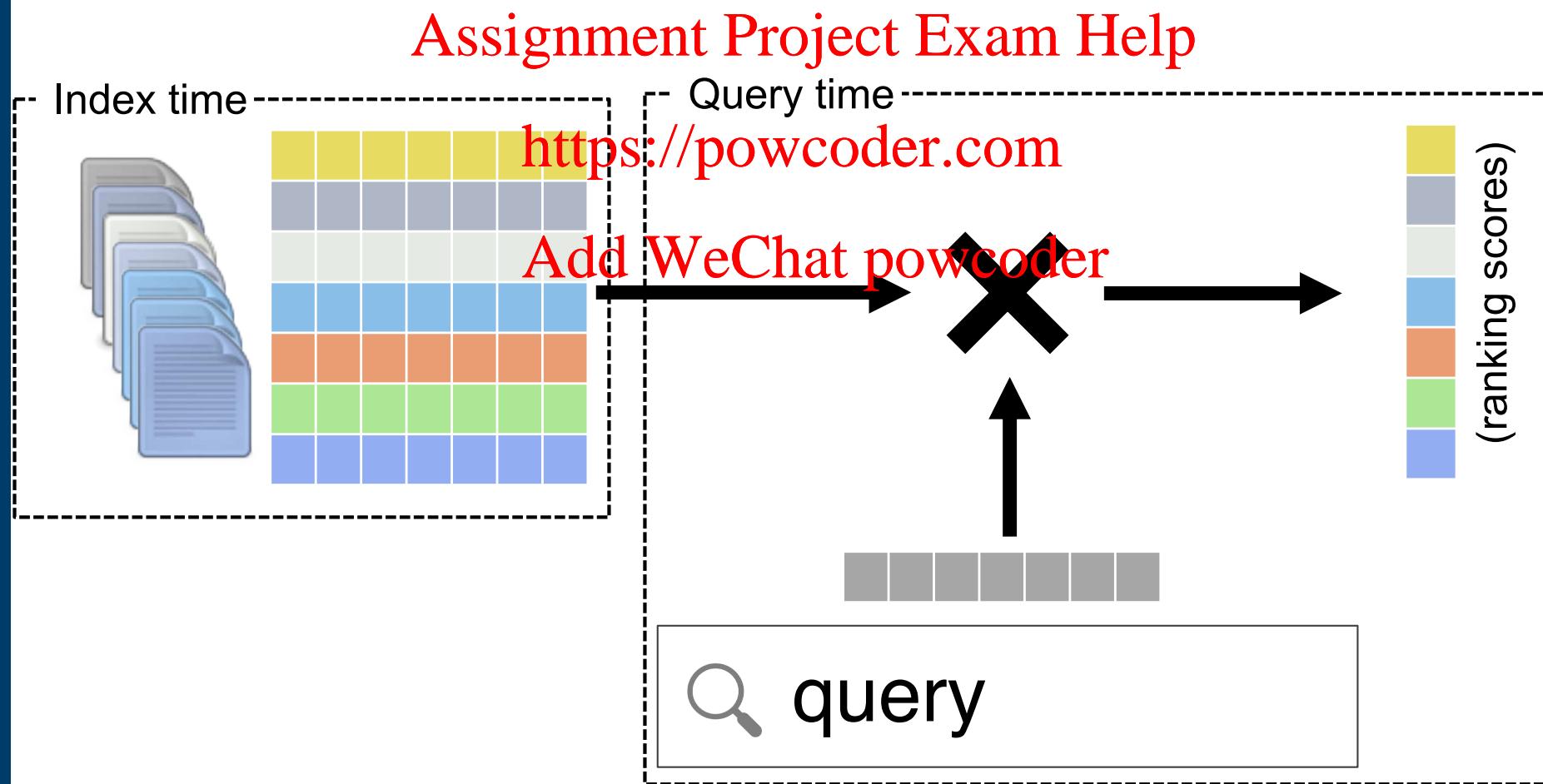
For inexpensive re-ranking

[1] MacAvaney et al. CEDR: Contextualized Embeddings for Document Ranking. SIGIR 2019.

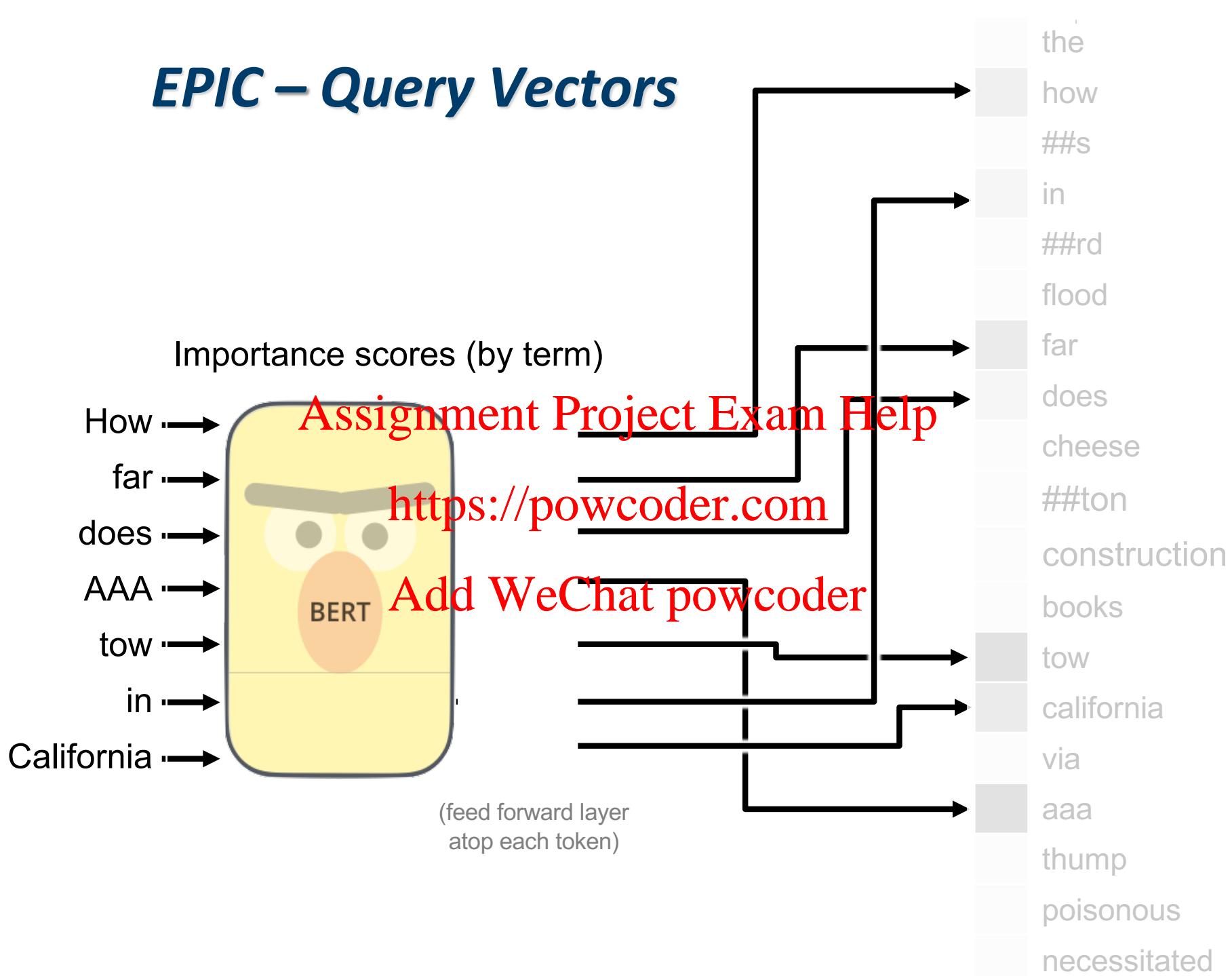
[2] Nogueira et al. Document Ranking with a Pretrained Sequence-to-Sequence Model. xxiv 2020.

[3] MacAvaney et al. Expansion via Prediction of Importance with Contextualization. SIGIR 2020.

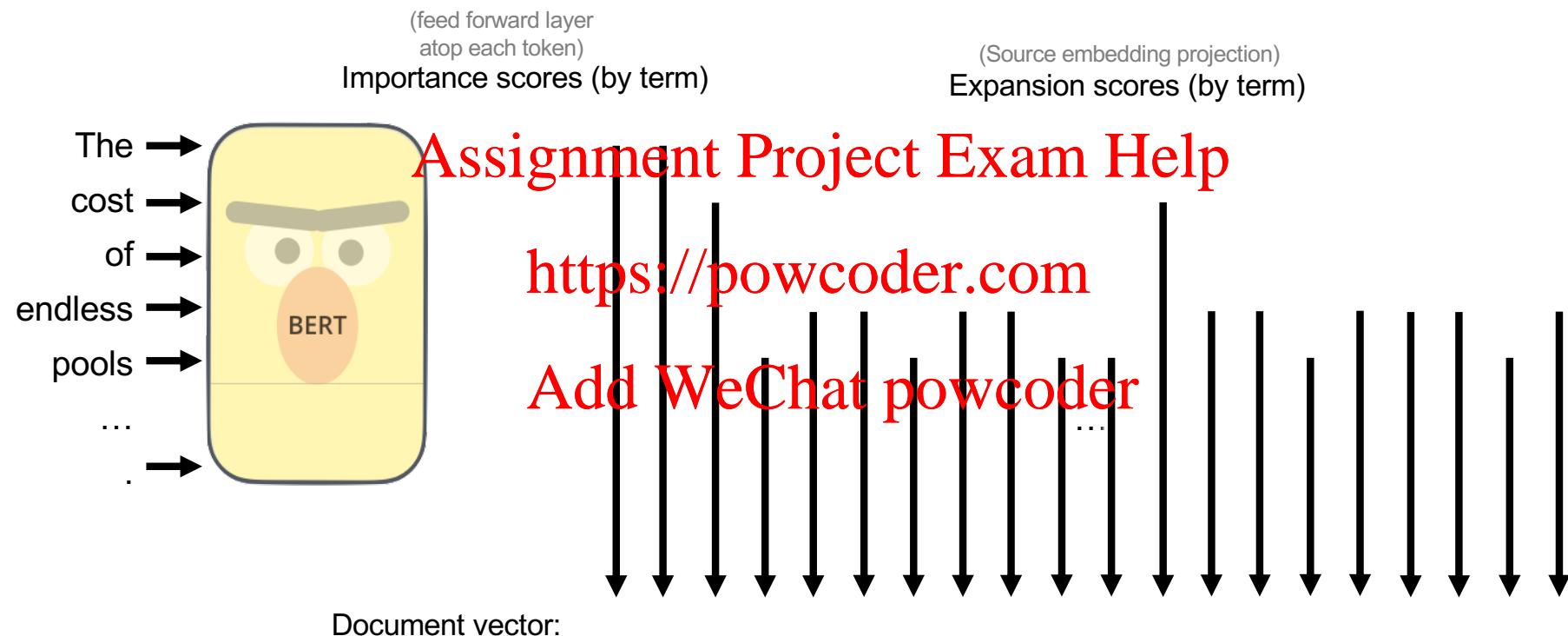
Main idea: Learn representations that are fast to score.



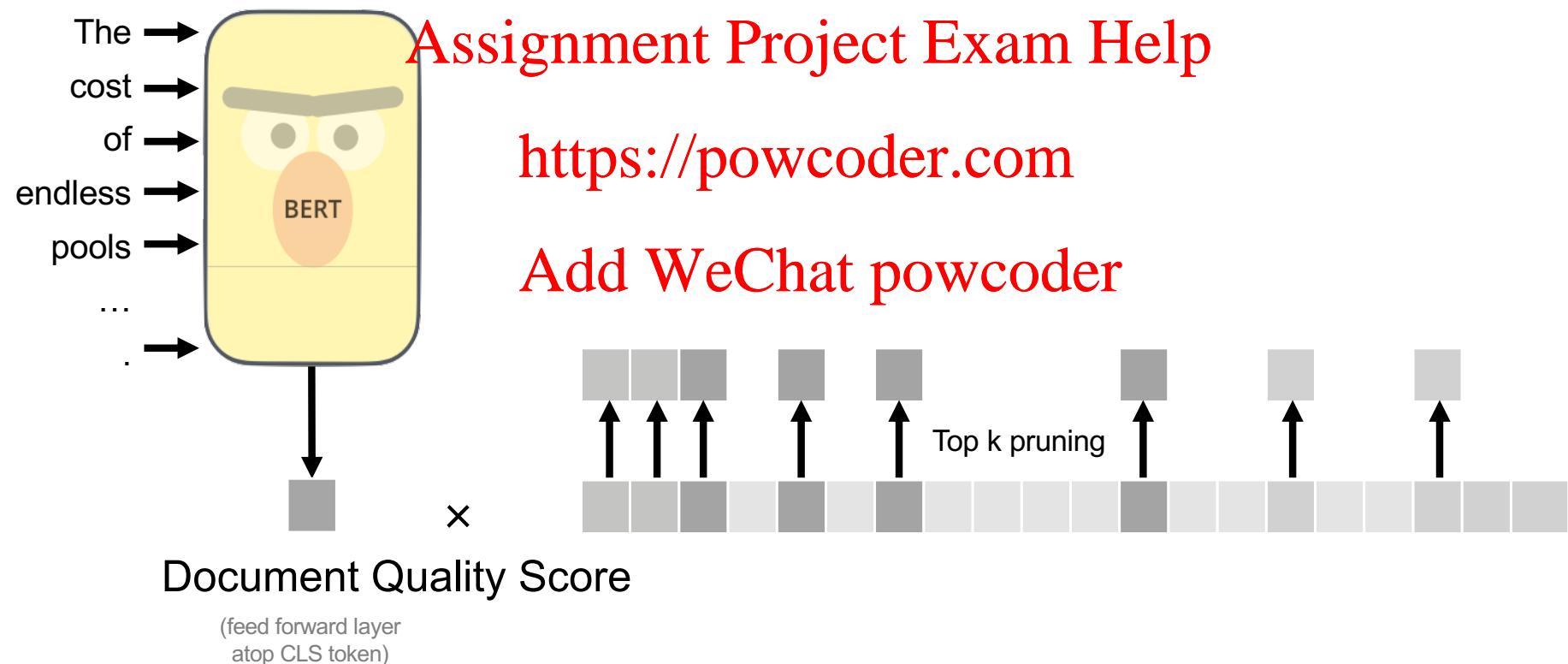
EPIC – Query Vectors



EPIC – Document Vectors



EPIC – Document Vectors



Today

1. Review of LTR & Basics of Neural Networks for NLP
2. Neural Re-ranking
Assignment Project Exam Help
3. Neural Retrieval
<https://powcoder.com>
4. Neural Query Rewriting & PRF
Add WeChat powcoder
5. Neural Document Rewriting
6. Neural IR in PyTerrier

Neural Ranking



BM25
E.g., CEDR
Assignment Project Exam Help

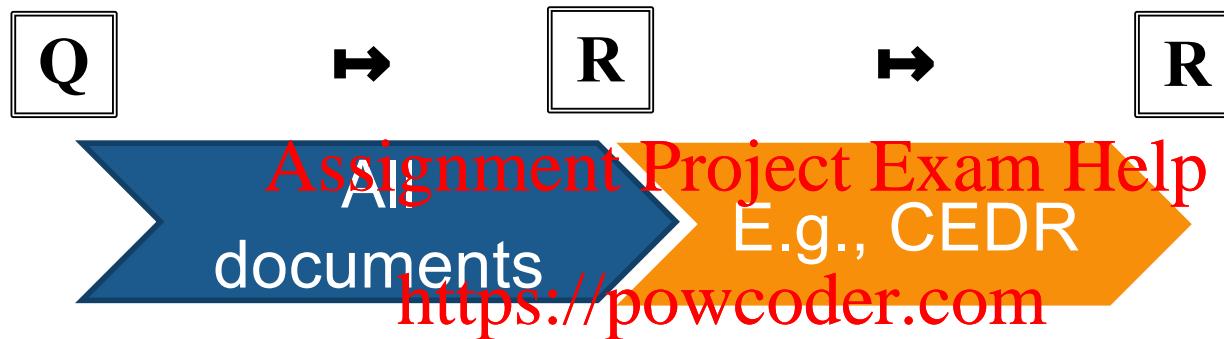
Can we replace the lexical scoring function with a neural network
~~Add WeChat powcoder what's important for the neural re-ranker?~~



Neural Network
E.g., CEDR

Neural Ranking

One solution: Score all documents in the corpus:

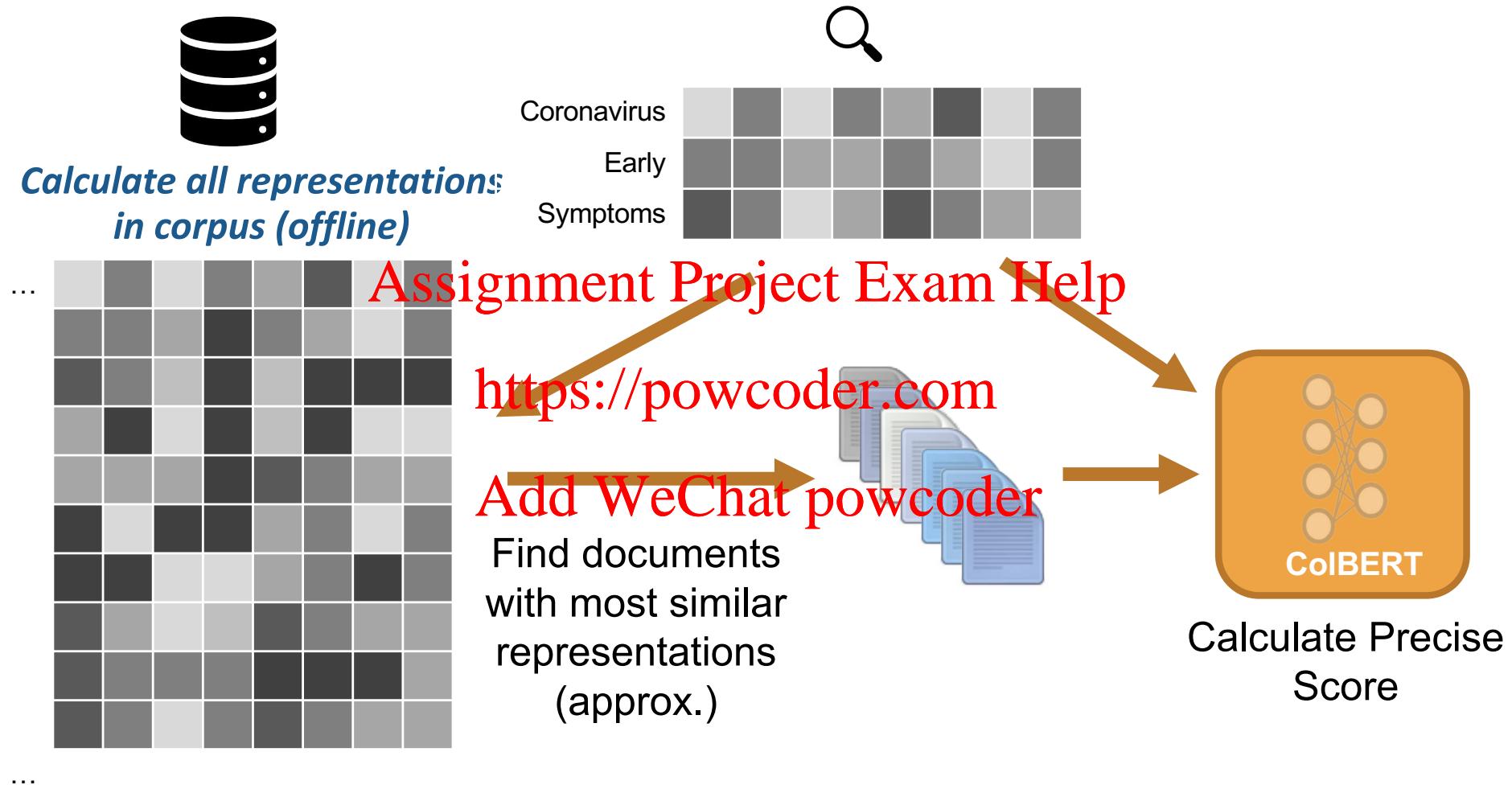


Add WeChat powcoder

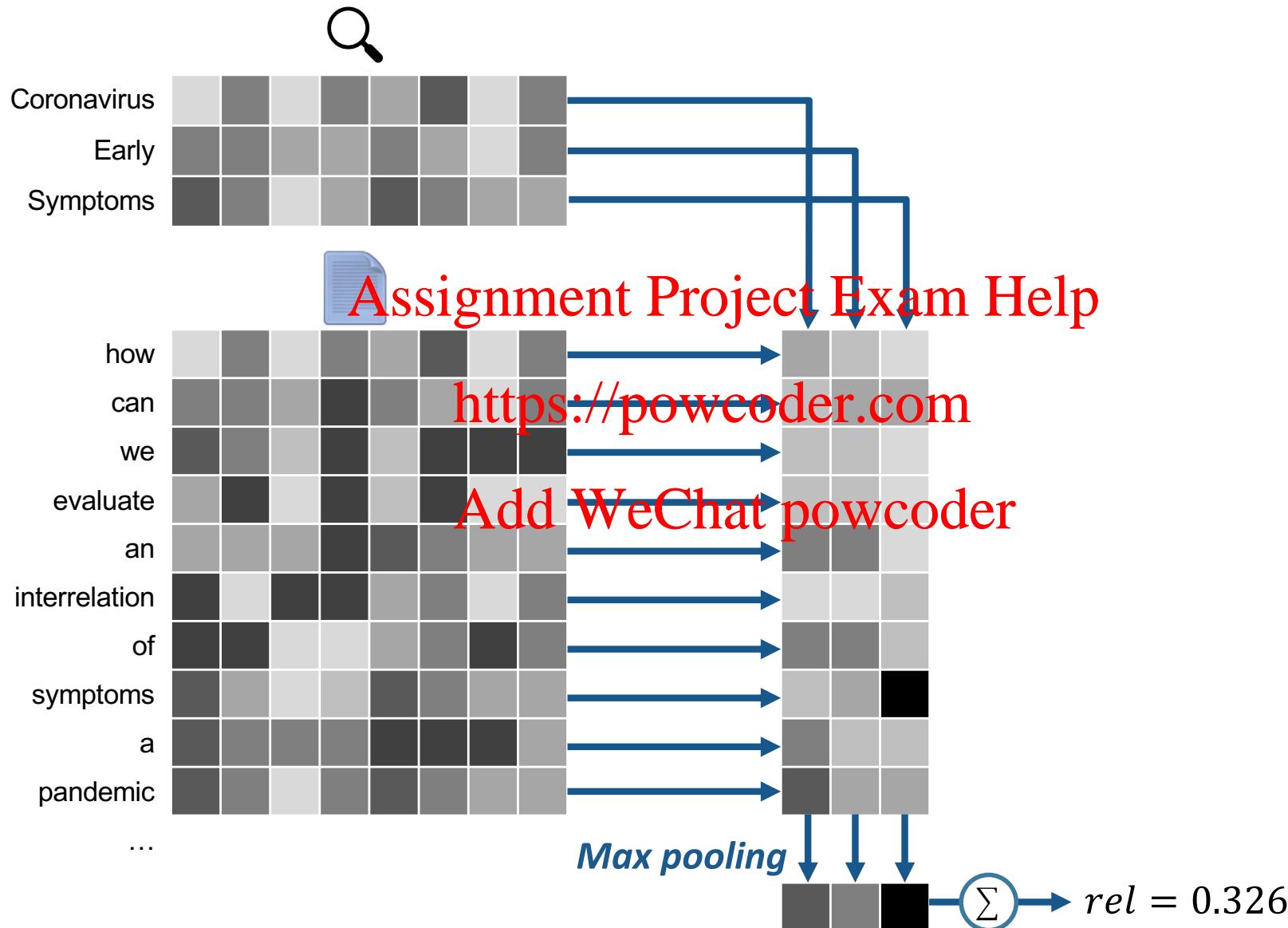
**Major Problem: Remember that these approaches
are very expensive – time-prohibitive to score all**

**Solution: Pre-compute representations
that are fast to score.**

ColBERT: Approx. Nearest Neighbour



Contextualized Late Interaction over BERT (ColBERT)



Problem with ColBERT

The token representations are very large (128 floats per token in each document)

This means a lot of storage is required.

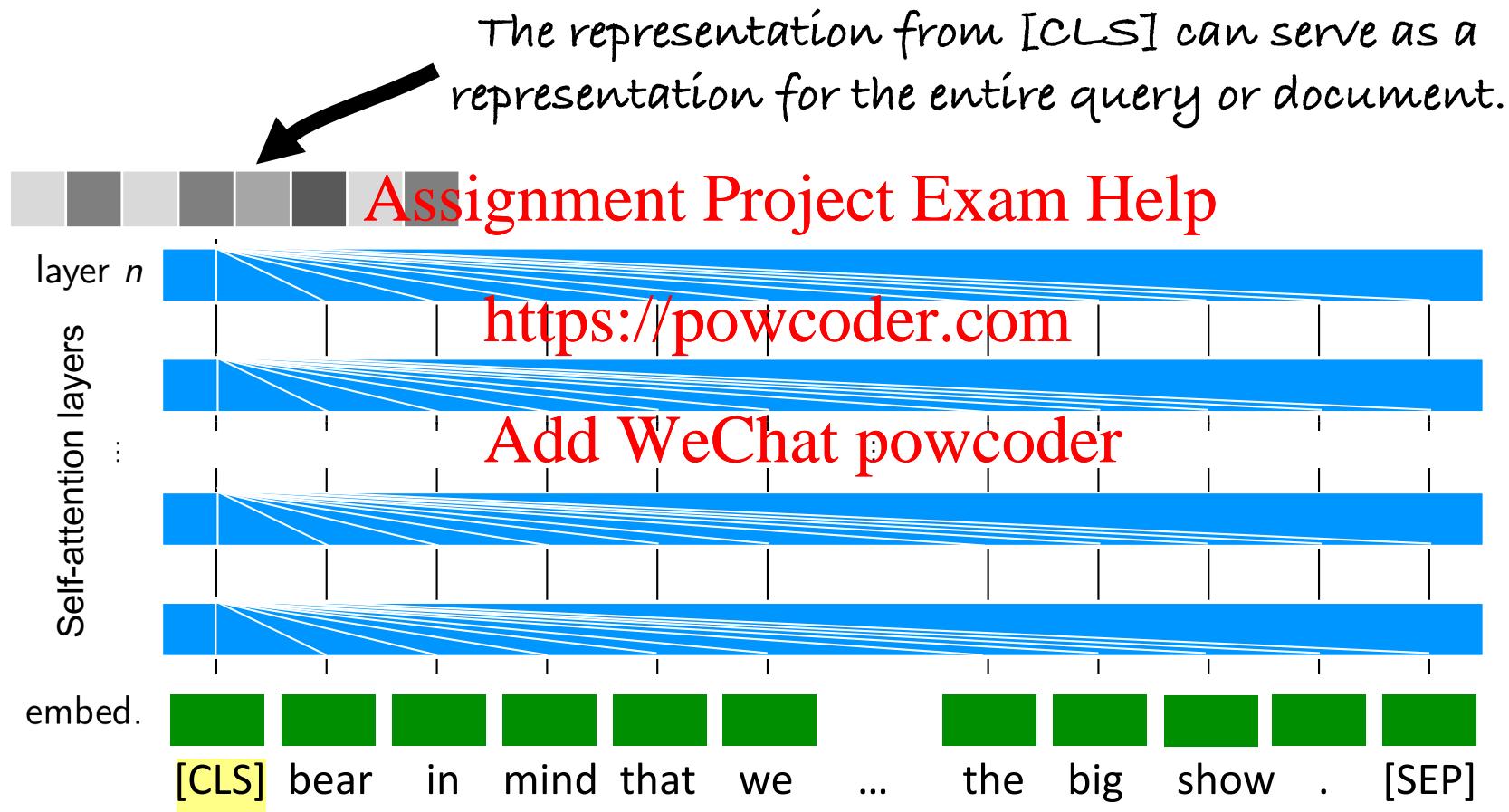
Assignment Project Exam Help

Alternative: build just a single representation
for each document

<https://powcoder.com>

Add WeChat powcoder

Single-Representation BERT (Example: ANCE)



Overview: Dense Retrieval

Two main strategies:

1. Multiple-Representation

- One vector for each term
[Assignment](#) [Project](#) [Exam](#) [Help](#)
- Find similar vectors to those in the query
<https://powcoder.com>
- Second stage required: score entire document
[Add WeChat](#) [powcoder](#)
- Tends to be more robust across datasets than single-rep

2. Single Representation

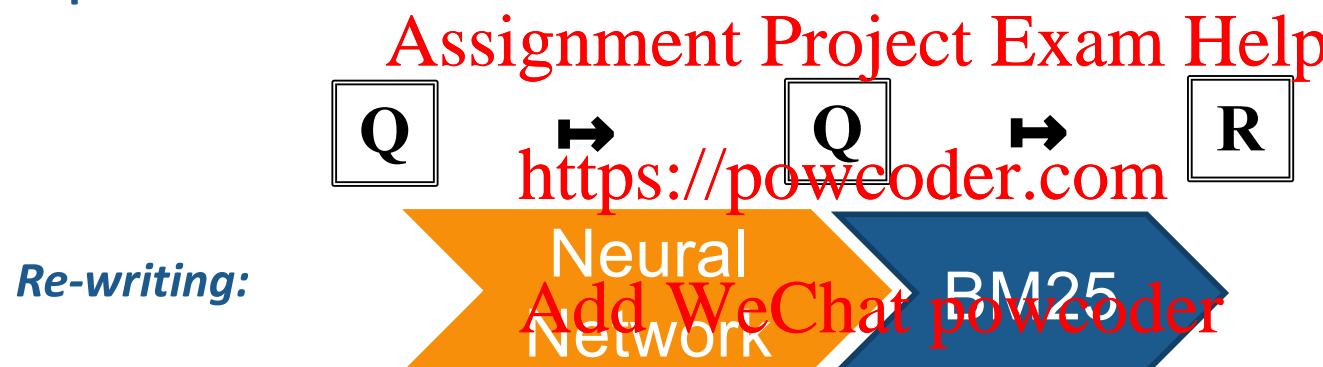
- One vector for document
- Find similar vectors to a representation for the query
- Tends to be faster & smaller than multi-rep

Today

1. Review of LTR & Basics of Neural Networks for NLP
2. Neural Re-ranking
Assignment Project Exam Help
3. Neural Retrieval
<https://powcoder.com>
4. Neural Query Rewriting & PPE
Add WeChat powcoder
5. Neural Document Rewriting
6. Neural IR in PyTerrier

Neural Query Rewriting & PRF

Can we use a neural network to build better queries?



T5-QR & T5-PRF

Main idea: use a text generation model to generate better queries either based on:

(1) The query text alone

Assignment Project Exam Help

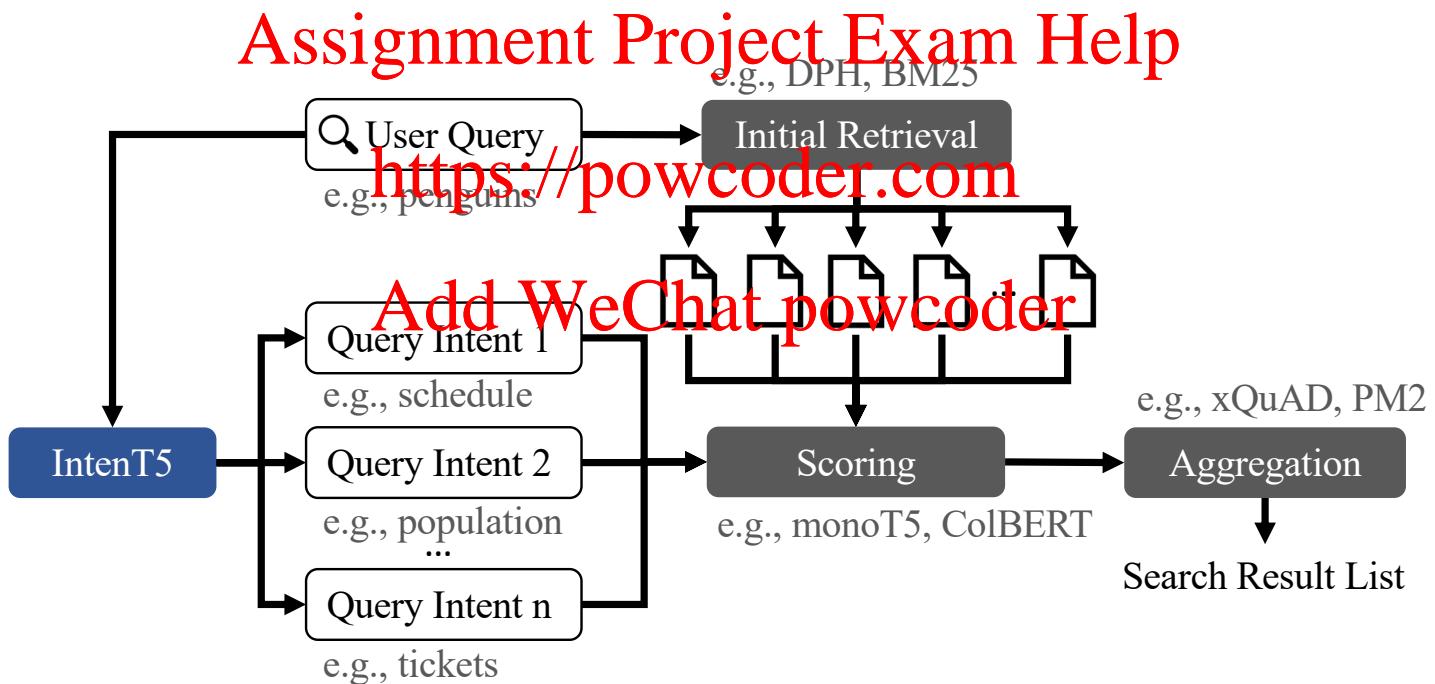
*Example: Do goldfish Grow -> How long do goldfish grow
<https://powcoder.com>*

(2) The query text and the top-retrieved results

*Example: do goldfish grow. Context: A: The conditions goldfish are kept in plus their diet determine how large they will grow. I have seen goldfish grow ridiculously large in very small containers when their water was changed frequently. Goldfish will not grow if water conditions are poor. Fancy goldfish don't grow as large as Common goldfish. A good size would be around 5 inches body length for most fancy varieties, 8 inches for Comets and 12 inches for Common Goldfish. These sizes are usually only attained by pond grown fish.
-> large kept container water conditions*

IntenT5

Generating multiple queries and aggregating them can help improve search result diversity.

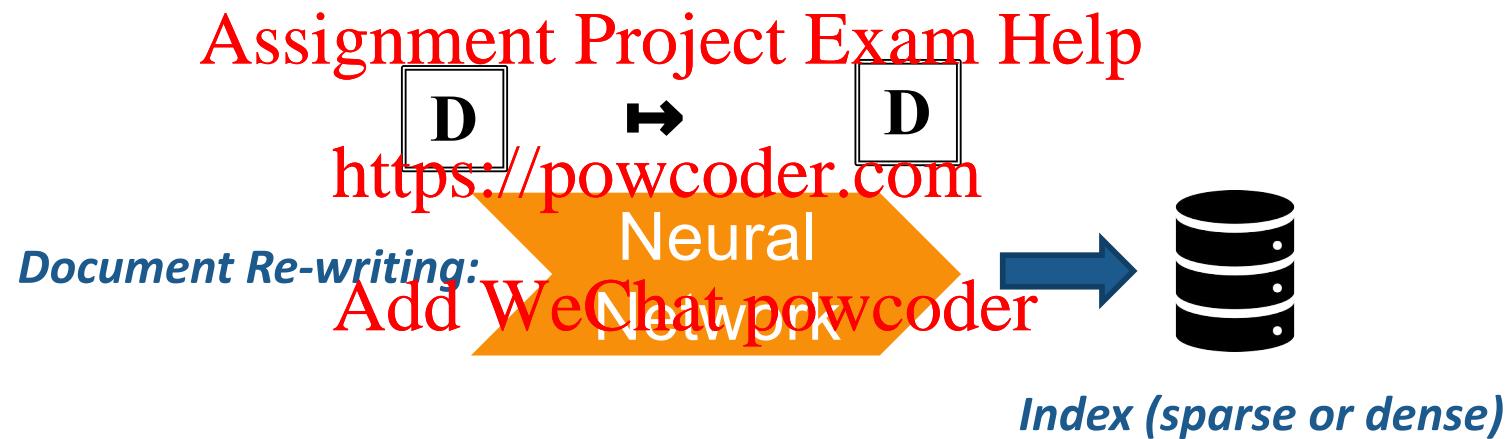


Today

1. Review of LTR & Basics of Neural Networks for NLP
2. Neural Re-ranking
Assignment Project Exam Help
3. Neural Retrieval
<https://powcoder.com>
4. Neural Query Rewriting & PRF
Add WeChat powcoder
5. Neural Document Rewriting
6. Neural IR in PyTerrier

Document Re-Writing

Can we find better representations of documents to index?



DeepCT

Main idea: Boost the term frequency of important words by repeating them in the text.

Abstract:

Nidovirus subgenomic mRNAs contain a leader sequence derived from the 5' end of the genome fused to different sequences ('bodies') derived from the 3' end. Their generation involves a unique mechanism of discontinuous subgenomic RNA synthesis that resembles copy-choice RNA recombination. During this process, the nascent RNA strand is transferred from one site in the template to another, during either plus or minus...

Abstract:

Nidovirus Nidovirus Nidovirus subgenomic subgenomic subgenomic mRNAs mRNAs mRNAs contain a leader sequence derived from the 5' end of the genome fused to different sequences ('bodies') derived from the 3' end. Their generation involves a unique mechanism of discontinuous subgenomic subgenomic RNA RNA synthesis that resembles copy-choice RNA RNA recombination. During this process, the nascent RNA RNA RNA strand is transferred from one site in the template to another, during either plus or minus...

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

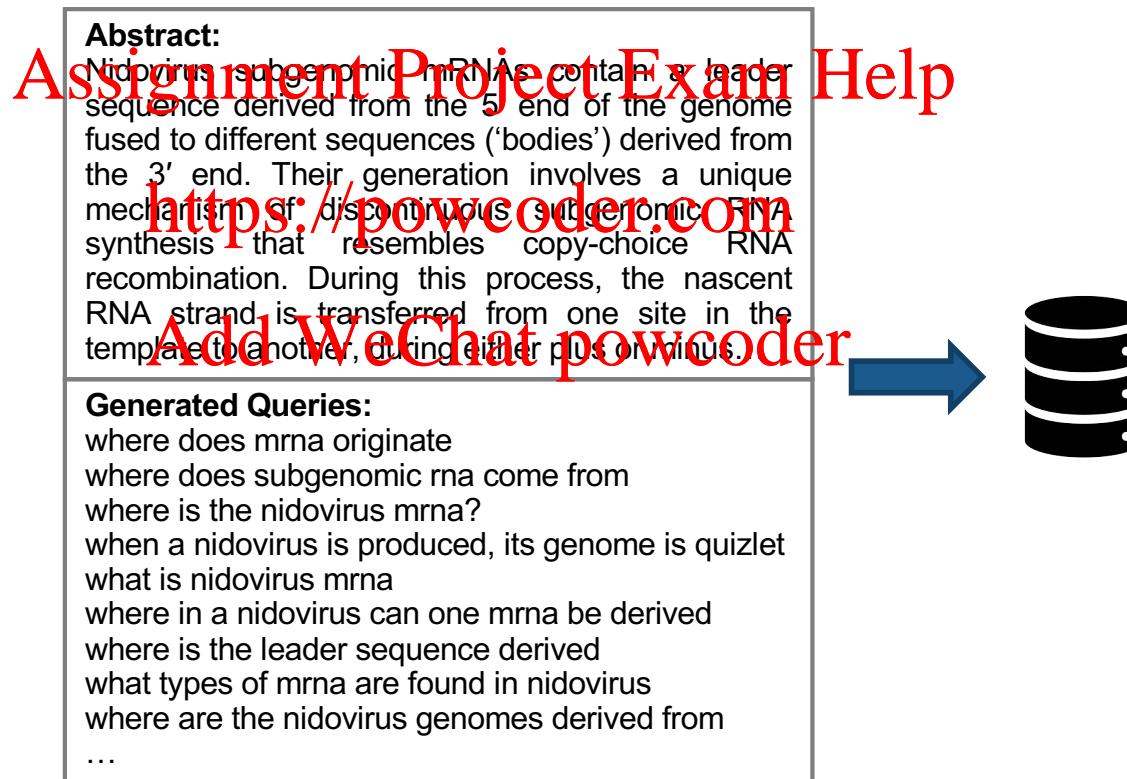


Training: Predict the terms that match relevant query terms.

Doc2Query

Main idea: Use a causal language model to generate additional text to add to documents when indexing.

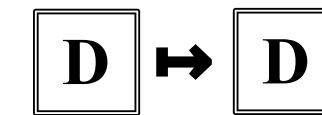
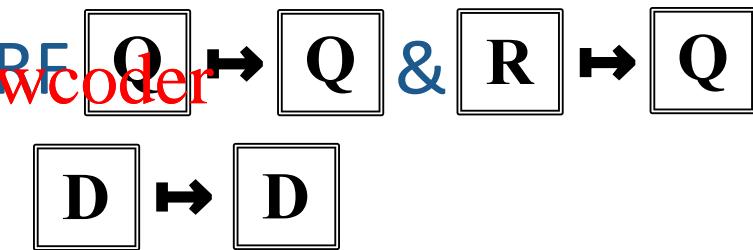
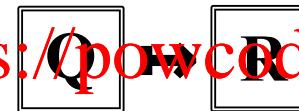
At retrieval time, use standard retrieval model (e.g., BM25).



Training: optimize to generate queries from passages on a collection with many queries (e.g., MS-MARCO)

Today

1. Review of LTR & Basics of Neural Networks for NLP
2. Neural Re-ranking
Assignment Project Exam Help
3. Neural Retrieval
<https://powcoder.com>
4. Neural Query Rewriting & PRF
Add WeChat powcoder
5. Neural Document Rewriting
6. Neural IR in PyTerrier



Neural IR In PyTerrier

Ranker: Which model to use?

```
drmm = onir_pt.reranker('drmm', 'wordvec_hash', text_field='abstract')
knrm = onir_pt.reranker('knrm', 'wordvec_hash', text_field='abstract')
pacrr = onir_pt.reranker('pacrr', 'wordvec_hash', text_field='abstract')
```

Which doc text field to use?

Vocab: Which word embeddings to use?

Assignment Project Exam Help

```
br = pt.BatchRetrieve(index) # 100
pipeline = br >> pt.text.get_text(dataset, 'abstract') >> reranker
```

These models need the document text. You can fetch it using `pt.text.get_text` or `BatchRetrieve`

```
pt.Experiment(
    [br, drmm_pipeline, knrm_pipeline, pacrr_pipeline],
    dataset.get_topics('title'),
    dataset.get_qrels(),
    names=['DPH', 'DPH >> DRMM', 'DPH >> KNRM', 'DPH >> PACRR'],
    eval_metrics=["recip_rank", "P.5", 'ndcg_cut.10', 'mrt']
)
```

Add WeChat powcoder

	name	recip_rank	P_5	ndcg_cut_10	mrt
0	DPH	0.767259	0.684	0.584309	33.919908
1	DPH >> DRMM	0.515536	0.420	0.377125	88.168377
2	DPH >> KNRM	0.488852	0.340	0.329785	79.464181
3	DPH >> PACRR	0.623139	0.532	0.457561	89.766008

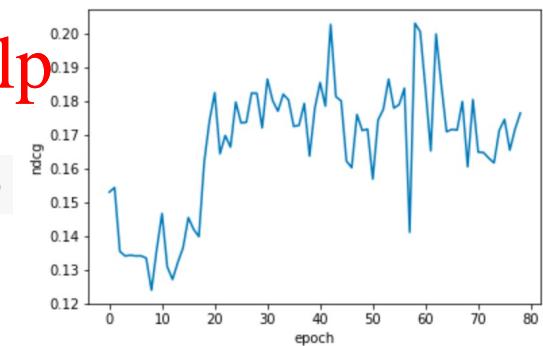
Neural IR In PyTerrier - Training

Ranker: Train model using medical subset of MS MARCO [1,2]

```
train_ds = pt.datasets.get_dataset('irds:msmarco-passage/train/medical')
```

```
fit_res = knrm.fit(  
    train_ds.get_topics(),  
    train_ds.get_qrels(),  
    valid_ds.get_topics(),  
    valid_ds.get_qrels())
```

Assignment Project Exam Help
<https://powcoder.com>



Add WeChat powcoder

	name	map	recip_rank	ndcg	ndcg_cut_10	mrt
0	DPH	0.075329	0.767259	0.164584	0.584309	30.752108
1	DPH >> KNRM	0.075105	0.810732	0.164964	0.619237	77.134939

[1] Bajaj et al. MS MARCO: A Human Generated Mchine Reading COnprehension Dataset}. InCoCo@NeurIPS 2016.
[2] MacAvaney et al. SLEDGE-Z: A Zero-Shot Baseline for COVID-19 Literature Search. EMNLP 2020.

Neural IR In PyTerrier - Pretrained

Load trained model from URL

```
bert = onir_pt.reranker.from_checkpoint('https://macavaney.us/scibert-medmarco.tar.gz',
                                         text_field='title_abstract',
                                         expected_md5='854966d0b61543fffffa44cea627ab63b')
```

(optional) verification of download's hash

```
def cat_title_abstract(df):
    df['title_abstract'] = df[['title', 'abstract']].apply(lambda row: row[0] + ' ' + row[1], axis=1)
    return df
```

```
bert_pipeline = (br >> https://powcoder.com
                  pt.text.get_text(dataset, ['title', 'abstract']) >>
                  pt.apply.generic(cat_title_abstract) >>
                  bert)
```

Include both the title and abstract

```
pt.Experiment(
    [br, bert_pipeline],
    dataset.get_topics('title'),
    dataset.get_qrels(),
    names=['BM25', 'BM25 >> BERT'],
    eval_metrics=['map', 'ndcg', 'ndcg_cut.10', 'mrt']
)
```

TREC COVID results

Title queries

	name	map	ndcg	ndcg_cut_10	mrt
0	BM25	0.073623	0.162657	0.583665	29.487896
1	BM25 >> BERT	0.077678	0.168394	0.650975	1637.205799

Description queries

	name	map	ndcg	ndcg_cut_10	mrt
0	BM25	0.077880	0.177728	0.644374	38.557969
1	BM25 >> BERT	0.085371	0.185821	0.740331	1748.576758

Neural IR In PyTerrier - EPIC

EPIC (indexed)

```
indexed_epic = onir_pt.indexed_epic.from_checkpoint(  
    'https://macavaney.us/epic.msmarco.tar.gz',  
    index_path='./epic_cord19')
```

Assignment Project Exam Help
indexed_epic.index(dataset.get_corpus_iter(), fields=['abstract',])

<https://powcoder.com>

```
pipeline = br >> indexed_epic.reranker()  
pt.Experiment(  
    [br, pipeline], Add WeChat powcoder  
    dataset.get_topics('title'),  
    dataset.get_qrels(),  
    names=[ "DPH", "DPH >> EPIC (indexed)" ],  
    eval_metrics=[ "recip_rank", "P.5", "mrt" ]  
)
```

	name	recip_rank	P_5	mrt
0	DPH	0.766833	0.684	30.500175
1	DPH >> EPIC (indexed)	0.821500	0.700	53.264584

Neural IR In PyTerrier – Indexing Pipelines

```
deepct = pyterrier_deepct.DeepCTTransformer(  
    "bert-base-uncased/bert_config.json",  
    "marco/model.ckpt-65816")
```

Assignment Project Exam Help

Example DeepCT outputs: <https://powcoder.com>

```
df.iloc[0]['text']
```

'OBJECTIVE: This retrospective chart review describes the epidemiology and clinical features of 40 patients with culture-proven *Mycoplasma pneumoniae* infections at King Abdulaziz University Hospital, Jeddah, Saudi Arabia. **METHODS:** Patients with positive *M. pneumoniae* cultures from respiratory specimens from January 1997 through December

```
deepct.transform(df).iloc[0][ "text" ]
```

'objective objective retrospective retrospective chart chart chart chart chart chart review describes epidemiology clinical features features features features patients patient

Neural IR In PyTerrier – Indexing Pipelines

```
dataset = pt.get_dataset("irds:cord19/trec-covid")
index_loc = "./deepct_index_path"
indexer = deepct>>>pt.terrier.Indexer(index_loc)
```

Assignment Project Exam Help

Build an *indexing* pipeline

<https://powcoder.com>

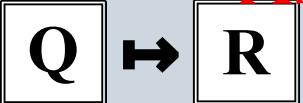
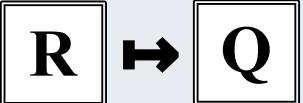
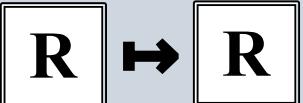
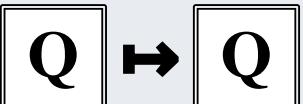
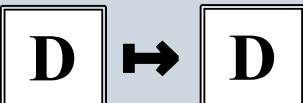
```
indexref = indexer.index(dataset.get_corpus_iter())
```

Index the collection
Add WeChat powcoder

	name	map	ndcg	ndcg_cut_10
0	BM25	0.181478	0.373328	0.583665
1	BM25_deepct	0.132069	0.304521	0.538936

Overview

We covered techniques for replacing various transformations in IR to use neural networks!

Transformer Class	Neural Examples
 Retrieval	ANCE, ColBERT https://powcoder.com
 PRF	ColBERT-PRF Add WeChat powcoder
 Re-ranking	CEDR, monoT5
 Query Re-writing	T5-QE, IntenT5
 Doc. Re-writing	Doc2Query, DeepImpact

Closing Remarks / Future Directions

This was just a survey of recent advances.

1. Are there inherent trade-offs among:

- Effectiveness
- Efficiency
- Interpretability

<https://powcoder.com>

2. Can we prevent unintended biases/side-effects?

3. NLP technique

- In many models (particularly long-form seq2seq), models have the tendency to produce coherent text
- But these often contain factual or common-sense errors
- Can IR systems be incorporated into these models?
 - Example: K. Guu, K. Lee, Z. Tung, P. Pasupat, and M.-W. Chang. REALM: Retrieval-augmented language model pre-training. arXiv:2002.08909, 2020.

Bibliography

- Z. Dai, C. Xiong, J. Callan, and Z. Liu. Convolutional neural networks for soft-matching n-grams in ad-hoc search. In Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining (WSDM 2018), pages 126–134, 2018.
- J. Guo, Y. Fan, Q. Ai, and W. B. Croft. A deep relevance matching model for ad-hoc retrieval. In Proceedings of the 25th ACM International on Conference on Information and Knowledge Management (CIKM 2016), pages 55–64, Indianapolis, Indiana, 2016.
- K. Hui, A. Yates, K. Berberich, and G. de Melo. PACRR: A position-aware neural IR model fore levance matching. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 1049–1058, Copenhagen, Denmark, 2017.
- S. MacAvaney, A. Yates, A. Cohan, and N. Goharian. CEDR: Contextualized embeddings for document ranking. In Proceedings of the 42nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2019), pages 1101–1104, Paris, France, 2019a.
- S. MacAvaney, A. Cohan, and N. Goharian. SLDC-E: A simple yet effective baseline for coronavirus scientific knowledge search.arXiv:2005.02365, 2020a.
- S. MacAvaney, F. M. Nardini, R. Perego, N. Tonellootto, N. Goharian, and O. Frieder. Expansion via prediction of importance with contextualization. In Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2020), pages 1573–1576, 2020c.
- S. MacAvaney, S. Feldman, N. Goharian, D. Downey, and A. Cohan. ABNIRML: Analyzing the Behavior of Neural IR Models. In TACL 2022.
- K. Guu, K. Lee, Z. Tung, P. Pasupat, and M.-Y. Chang. TFAE: Retrieval-augmented language model pre-training.arXiv:2002.08909, 2020.
- R. Nogueira and J. Lin. From doc2query to docTTTTquery, 2019.
- R. Nogueira and K. Cho. Passage re-ranking with BERT.arXiv:1901.04085, 2019.
- C. Xiong, Z. Dai, J. Callan, Z. Liu, and R. Power. End-to-end neural ad-hoc ranking with kernel pooling. In Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2017), pages 55–64, Tokyo, Japan, 2017.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder