

CSCSE 638 Natural Language Processing Foundation and Techniques

Lecture 15: Text Similarity, Retrieval-Augmented Generation

Kuan-Hao Huang

Spring 2026



Quiz 2

- Average: 90.67
- Std: 8.67
- Q1: 84
- Median: 93
- Q3: 97

Check [Gradescope](#) for details. For questions, send emails to csce638-ta-26s@lists.tamu.edu with “[CSCE 638] Subject ...” or check with TA in office hours

Quiz 3

- Mar 18 (Wednesday)
- Coverage: mainly Lecture 11 to 15
 - Naturally include some concepts in Lecture 1 to 10
- In-class, 20 minutes, closed book, no cheat sheet
- Written quiz, 5 questions
 - Please bring a pen
- Tips
 - Get familiar with formula
 - Understand the intuition behind the formula and the design
 - Know the pros and cons of different approaches

Assignment 3

Assignment 3

RELEASE DATE: 03/15/2026

DUE DATE: 04/02/2026 11:59pm on [Gradescope](#)

L^AT_EX Template: <https://www.overleaf.com/read/prmshxybxqgh#4fa8a9>

Name: First-Name Last-Name UIN: 000000000

This assignment consists of two parts: a writing section and a programming section. For the writing section, please use the provided L^AT_EX template to prepare your solutions and remember to fill in your name and UIN. For the programming section, please follow the instructions carefully.

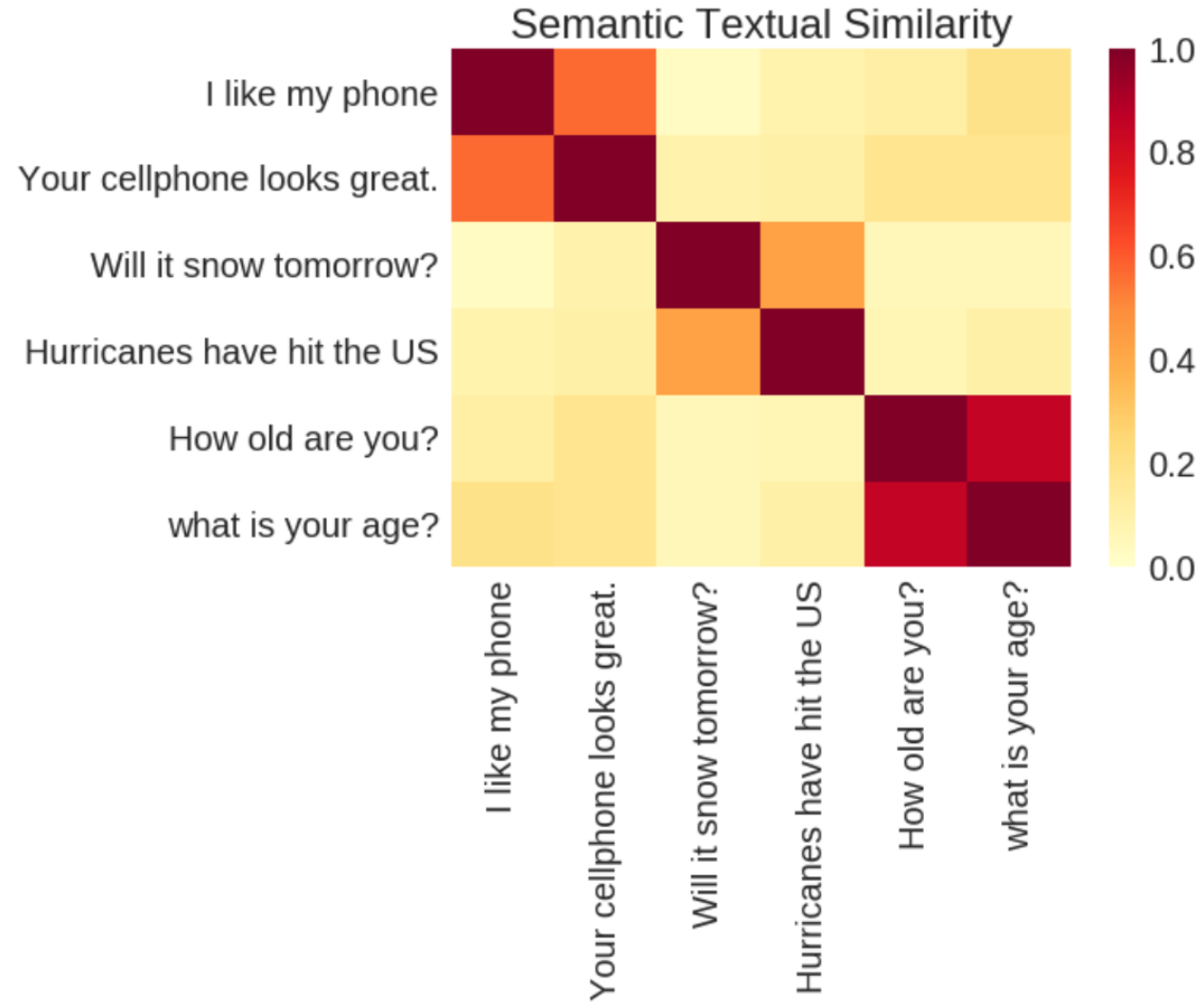
*Discussions with others on course materials and assignment solutions are encouraged, and the use of AI tools as assistance is permitted. However, you must ensure that **the final solutions are written in your own words**. It is your responsibility to avoid excessive similarity to others' work. Additionally, please clearly **indicate any parts where AI tools were used as assistance**.*

If you have any question, please send an email to csce638-ta-26s@list.tamu.edu

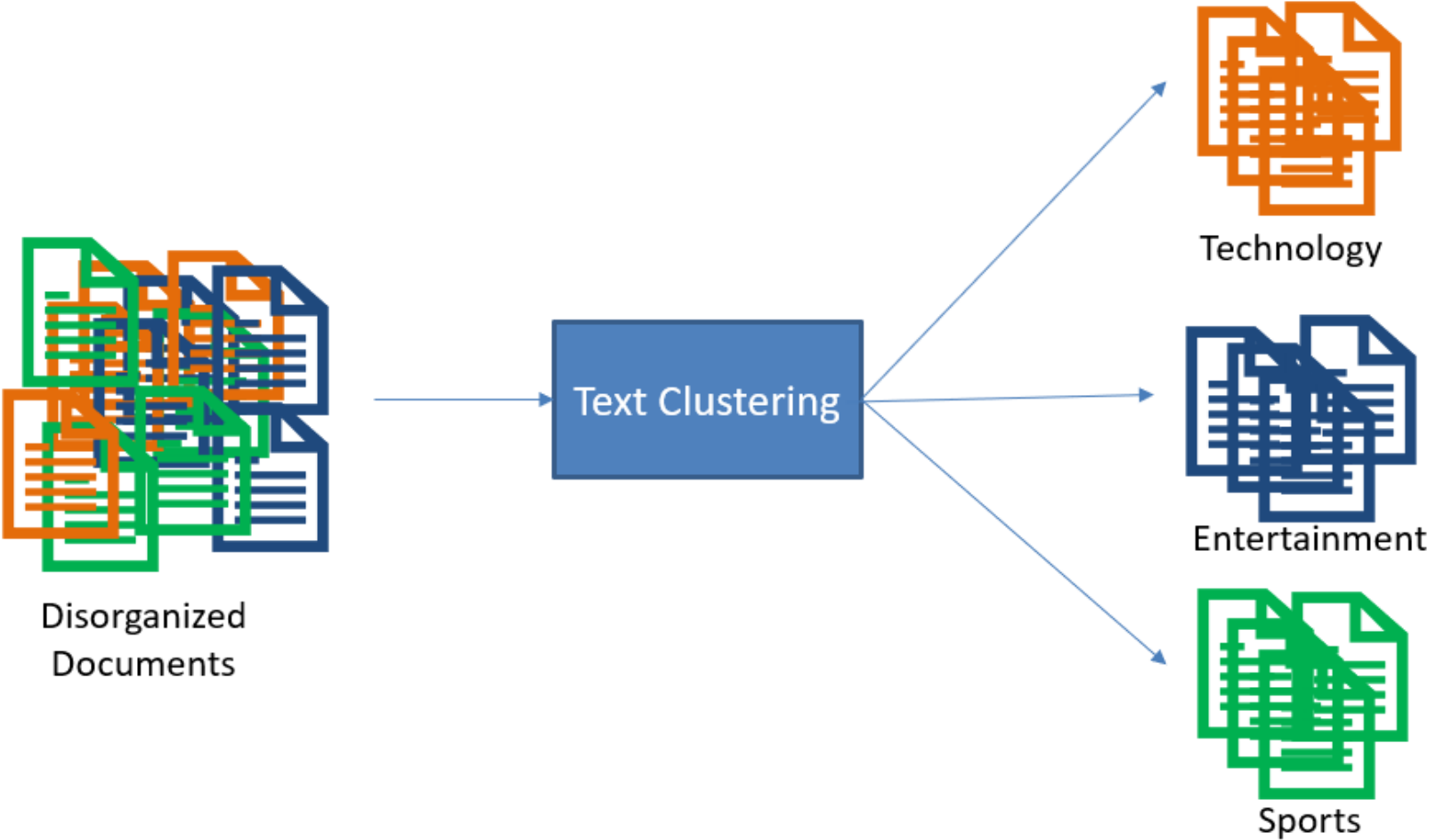
Lecture Plan

- Text Similarity
- Retrieval-Augmented Generation

Text Similarity



Document Clustering




Information Retrieval

The image shows a Google search interface with the search term "texas a&m". The search results are displayed in a list format. The first result is for "Texas A&M" with the URL "https://www.tamu.edu". The second result is for "Texas A&M Athletics" with the URL "https://12thman.com" and includes a small image of a person in a red hat. The third result is for "Texas A&M University-Corpus Christi" with the URL "https://www.tamucc.edu". The fourth result is for "Texas A&M Athletics" with the URL "https://12thman.com" and includes a link to "2024 Football Schedule".


Google

texas a&m


All News Images Maps Videos Shopping Forums More Tools

 Texas A&M
https://www.tamu.edu


Texas A&M University
Howdy from Texas A&M University. Texas A&M University is an engine of imagination, learning, discovery and innovation. Here, you'll learn essential career ...

 Texas A&M Athletics
https://12thman.com


Texas A&M Athletics - 12thMan.com
The official athletics website for the Texas A&M Aggies.
[Football](#) · [Staff Directory](#) · [2024 Football Schedule](#) · [Composite Calendar](#)

 Texas A&M University-Corpus Christi
https://www.tamucc.edu

Texas A&M University-Corpus Christi: Welcome Home
Welcome to THE ISLAND! Discover the Island University, the only university in the nation located on its own island, at the heart of the Texas Gulf Coast.

 Texas A&M Athletics
https://12thman.com › sports › football › schedule

2024 Football Schedule
2024 Football Schedule · Early: Game will have a start time between 11AM-Noon CT · Afternoon: Game will have a start time between 2:30PM – 3:30PM CT · Night: ...



Recommendation Systems

Your recently viewed items and featured recommendations

Sponsored products related to this search [What's this?](#)

									
	<p>All-new Echo Show (2nd Gen) + Ring Video Doorbell 2- Charcoal 1 offer from \$428.99</p>	<p>AmazonBasics Microwave, Small, 0.7 Cu. Ft, 700W, Works with Alexa ★★★★☆ 1,375 \$59.99 ✓prime</p>	<p>Echo Look Hands-Free Camera and Style Assistant with Alexa— includes Style Check to... ★★★★☆ 413 \$99.99 ✓prime</p>	<p>Sonos Beam - Smart TV Sound Bar with Amazon Alexa Built-in - Black ★★★★☆ 474 \$399.00 ✓prime</p>	<p>Echo Wall Clock - see timers at a glance - requires compatible Echo device ★★★★☆ 1,231 \$29.99 ✓prime</p>	<p>Echo Spot Adjustable Stand - Black ★★★★☆ 933 \$19.99 ✓prime</p>	<p>AHASTYLE Wall Mount Hanger Holder ABS for New Dot 3rd Generation Smart Home Speakers... ★★★★☆ 12 \$10.99 ✓prime</p>	<p>Angel Statue Crafted Stand Holder for Amazon Echo Dot 3rd Generation, Alexa Smart... ★★★★☆ 57 \$25.99 ✓prime</p>	

Page 1 of 3

Explore more from across the store

									
	<p>Actionable Gamification: Beyond Points, Badges, and Leaderboards › Yu-kai Chou</p>	<p>The Model Thinker: What You Need to Know to... › Scott E. Page</p>	<p>Don't Make Me Think, Revisited: A Common... › Steve Krug</p>	<p>Hooked: How to Build Habit-Forming Products › Nir Eyal</p>	<p>Microservices Patterns: With examples in Java › Chris Richardson</p>	<p>Solving Product Design Exercises: Questions &... › Artiom Dashinsky</p>	<p>100 Things Every Designer Needs to Know About... Susan Weinschenk</p>	<p>Infinity › Jonathan Hickman ★★★★☆ 182</p>	

Page 1 of 6

Semantic Quality Control

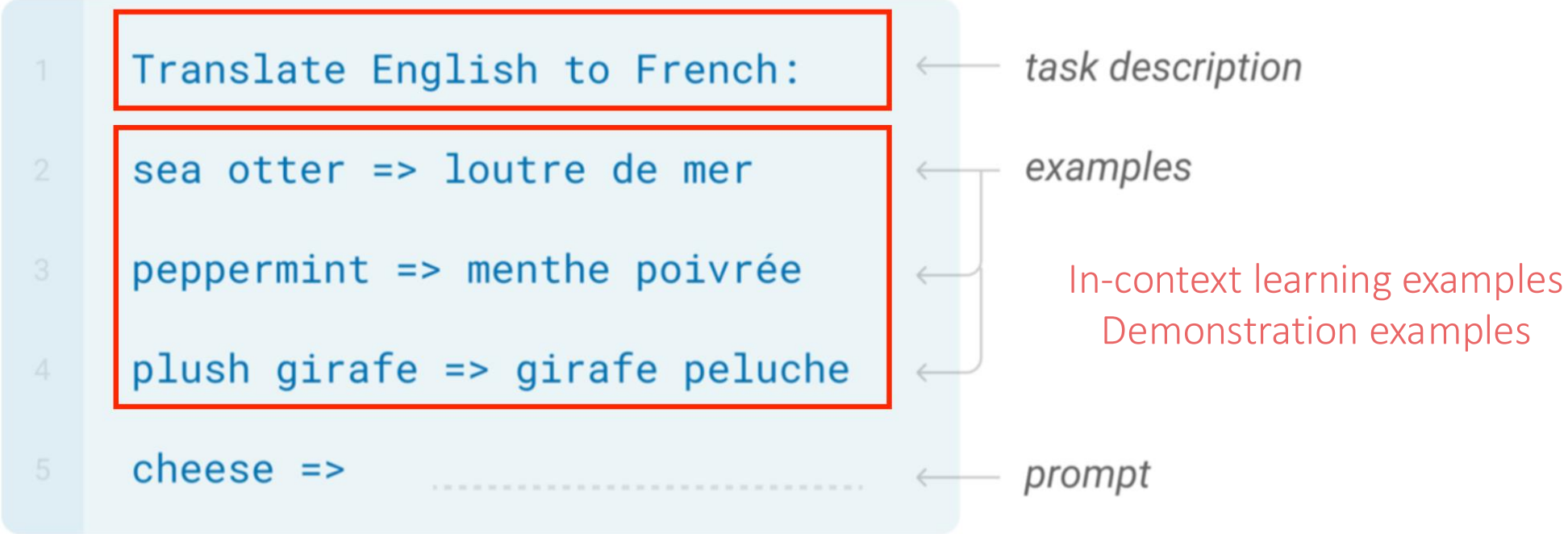
- Paraphrase generation

We will go hiking if tomorrow is a sunny day.

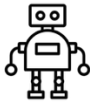
If it is sunny tomorrow, we will go hiking.

- Style transfer
- Plagiarism detection
- ...

In-Context Example Selection



Semantic Textual Similarity Benchmark



A soccer player is kicking the soccer ball into the goal from a long way down the field.

A soccer player kicks the ball into the goal.

3.25

3.94

Earlier this month, RIM had said it expected to report second-quarter earnings of between 7 cents and 11 cents a share.

Excluding legal fees and other charges it expected a loss of between 1 and 4 cents a share.

1.2

0.5

...

...

...

...

David Beckham Announces Retirement From Soccer.

David Beckham retires from football.

4.4

3.8

Pearson's Correlation Coefficient

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

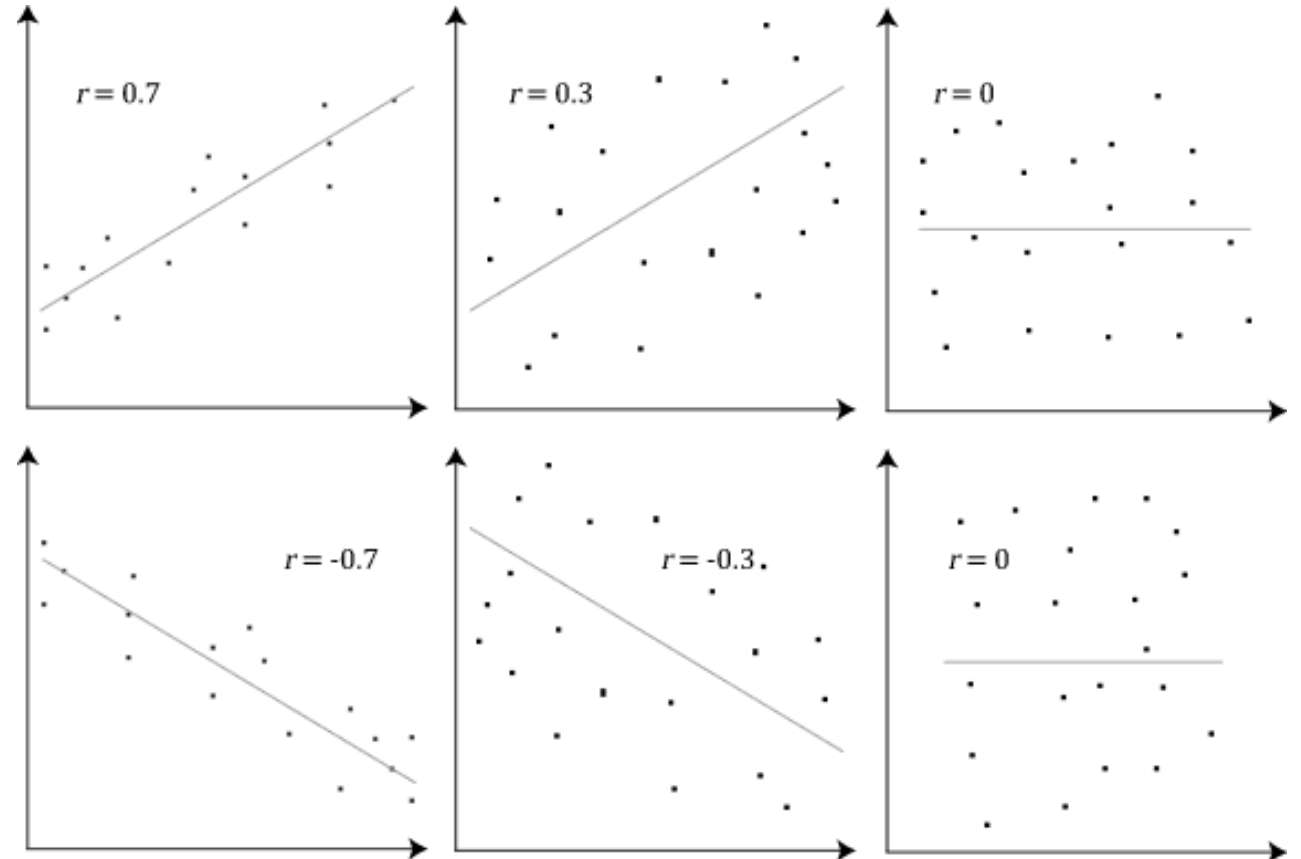
r = correlation coefficient

x_i = values of the x-variable in a sample

\bar{x} = mean of the values of the x-variable

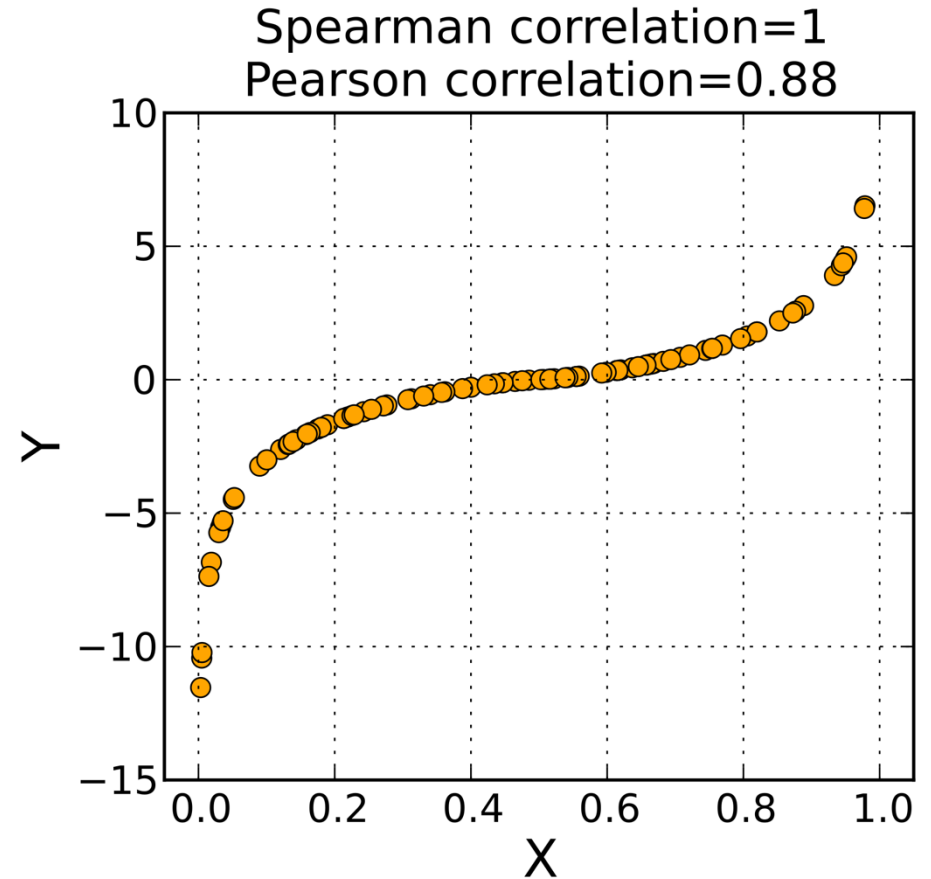
y_i = values of the y-variable in a sample

\bar{y} = mean of the values of the y-variable

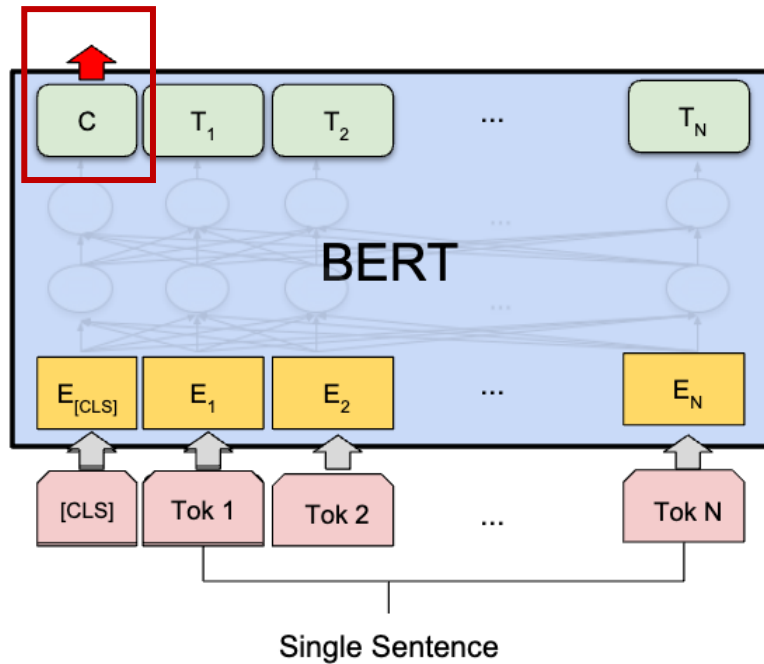


Spearman's Correlation Coefficient

- Pearson's correlation coefficient on **rank**
- Score
 - Human: [1.2, 3.4, 2.5, 0.7, 4.0]
 - Machine: [0.5, 3.3, 1.0, 1.2, 3.4]
- Rank
 - Human: [4, 2, 3, 5, 1]
 - Machine: [5, 2, 4, 3, 1]
- Assesses monotonic relationships
 - whether linear or not

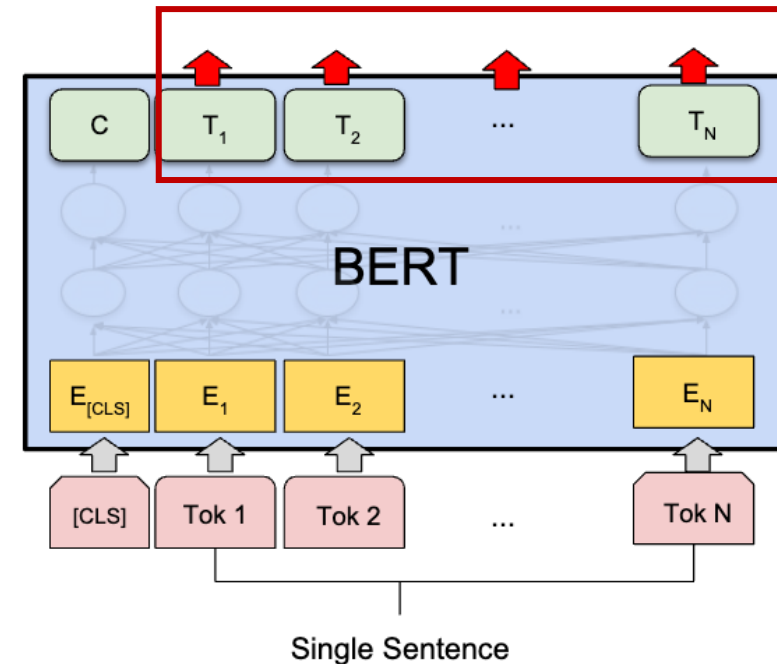


A Simple Approach: Text Encoder + Cosine Similarity



$$E_1 = \text{Encoder}(S_1)$$

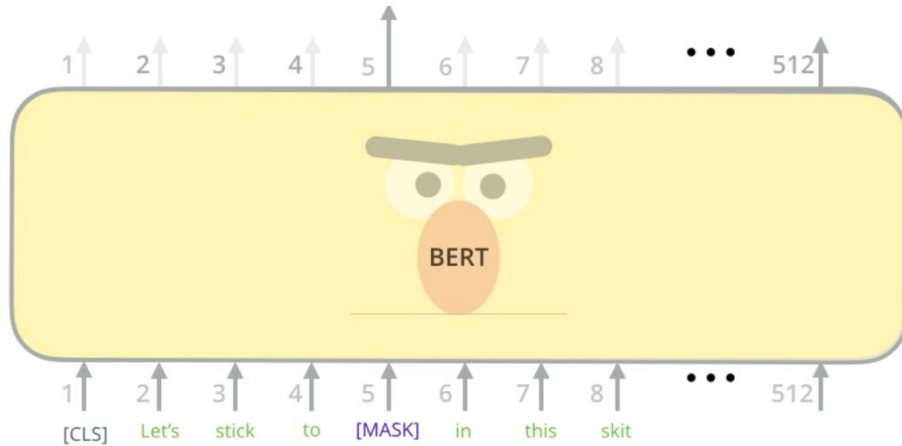
$$E_2 = \text{Encoder}(S_2)$$



$$\text{Similarity}(S_1, S_2) = \frac{E_1 \cdot E_2}{\|E_1\| \|E_2\|}$$

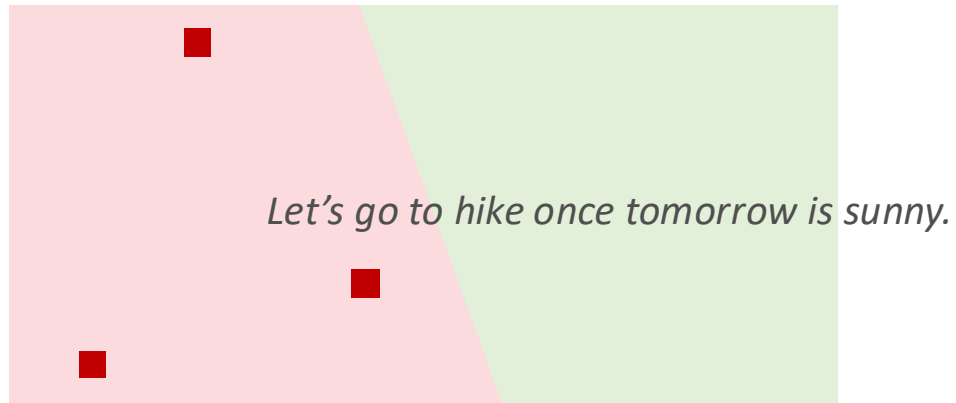
Unfortunately, the performance is bad (why?)

A Simple Approach: Text Encoder + Cosine Similarity



Pre-trained BERT embeddings are more about lexical information

If it is sunny tomorrow, we will go hiking.



Good classification performance \neq Good similarity

We will go hiking if tomorrow is a sunny day.

Sentence-BERT

- Consider SNLI dataset
 - Stanford Natural Language Inference

A boy is jumping on skateboard in the middle of a red bridge.

The boy skates down the sidewalk.

Contradiction

A boy is jumping on skateboard in the middle of a red bridge.

The boy is wearing safety equipment.

Neutral

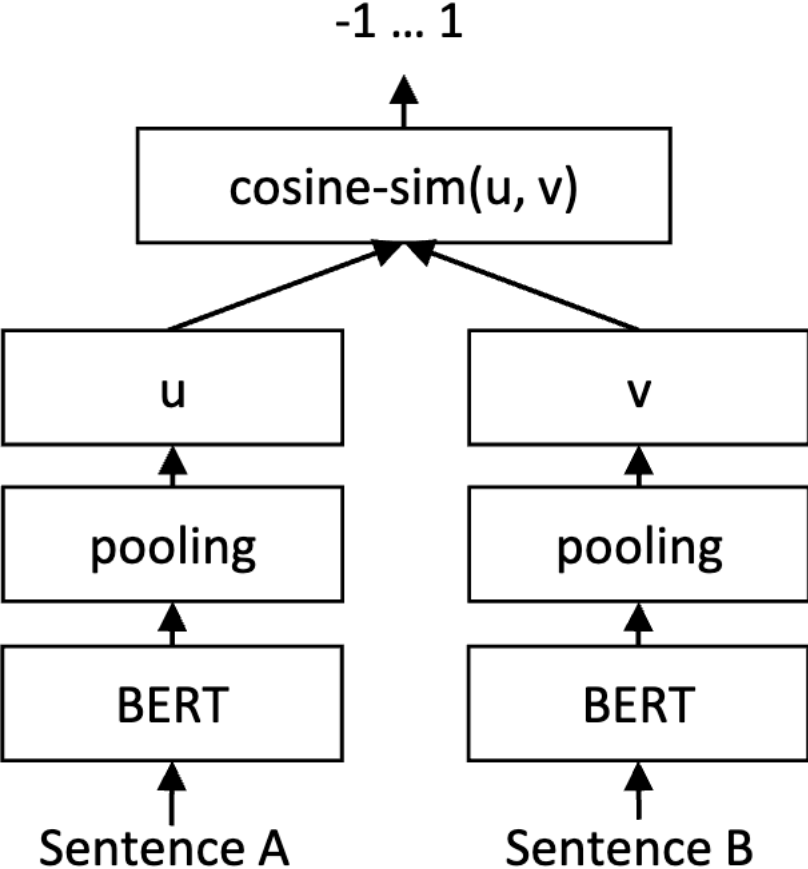
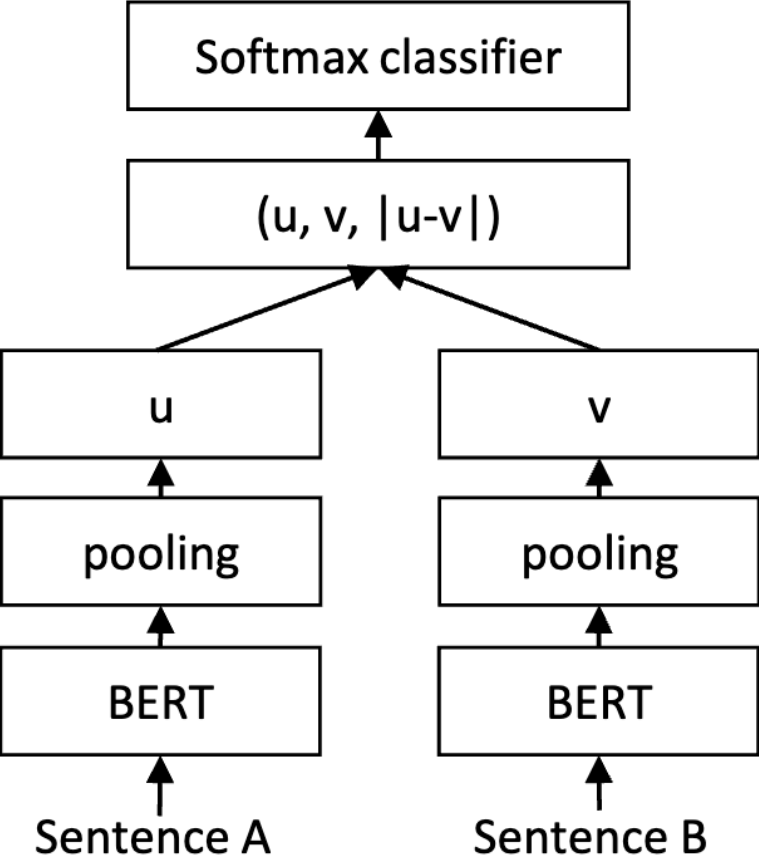
A boy is jumping on skateboard in the middle of a red bridge.

The boy does a skateboarding trick.

Entailment

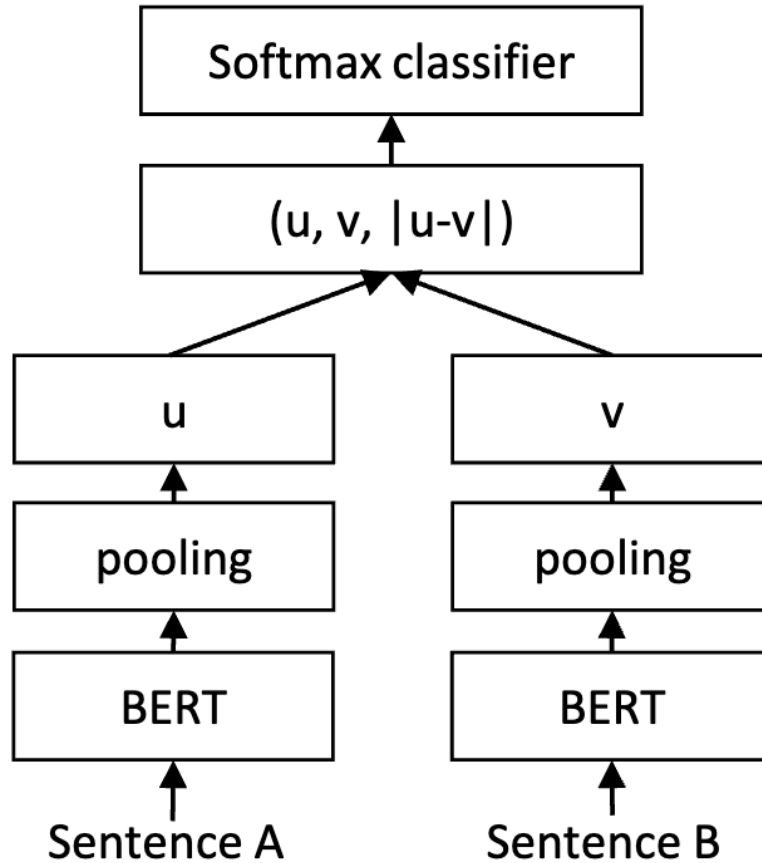
Sentence-BERT

Contradiction Neutral Entailment



Sentence-BERT

Contradiction Neutral Entailment



Cross Entropy Loss

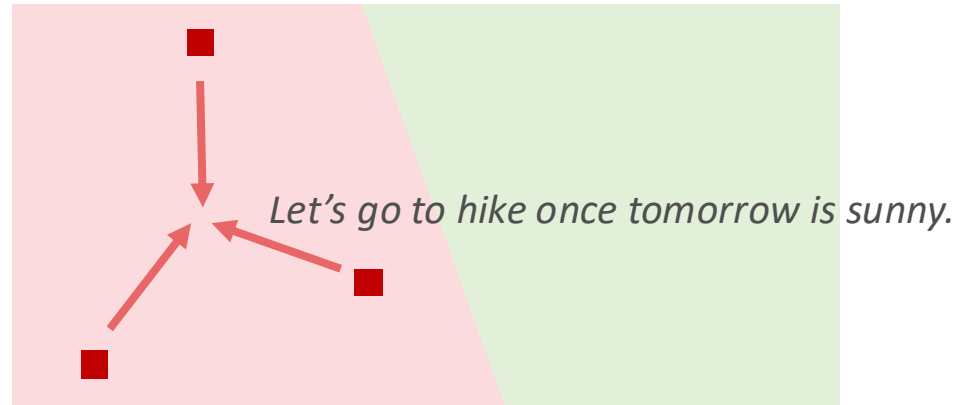
$$o = \text{softmax}(W_t(u, v, |u - v|))$$

Triplet Loss

$$\max(\|s_a - s_p\| - \|s_a - s_n\| + \epsilon, 0)$$

Sentence-BERT

If it is sunny tomorrow, we will go hiking.



We will go hiking if tomorrow is a sunny day.

Sentence-BERT: Performance

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICK-R	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	76.69	71.22
SBERT-NLI-base	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT-NLI-large	72.27	78.46	74.90	80.99	76.25	79.23	73.75	76.55
SRoBERTa-NLI-base	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa-NLI-large	74.53	77.00	73.18	81.85	76.82	79.10	74.29	76.68

SentenceTransformers

- <https://sbert.net/>

```
from sentence_transformers import SentenceTransformer

# 1. Load a pretrained Sentence Transformer model
model = SentenceTransformer("all-MiniLM-L6-v2")

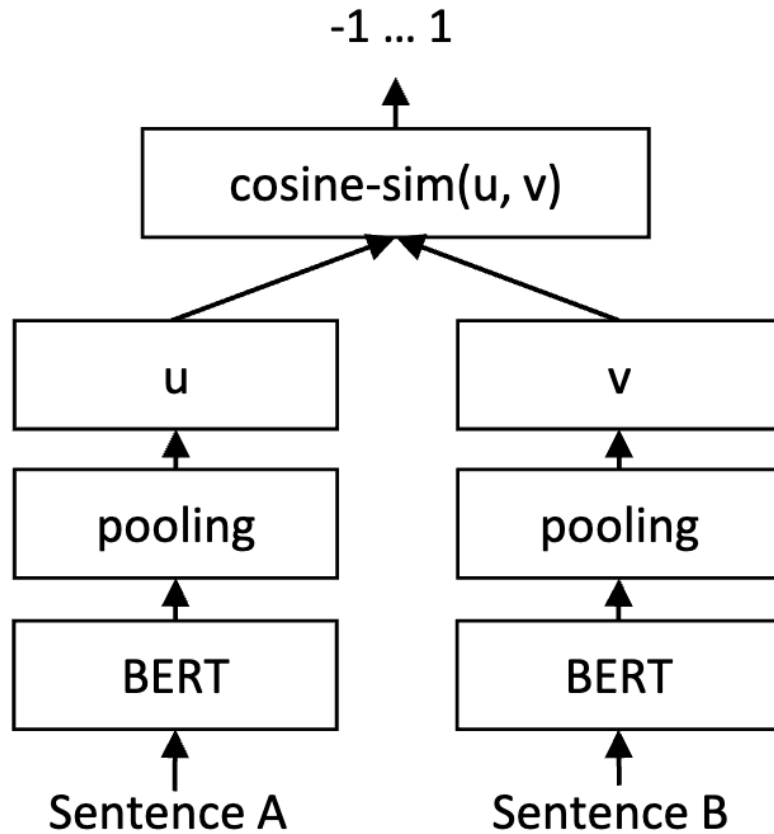
# The sentences to encode
sentences = [
    "The weather is lovely today.",
    "It's so sunny outside!",
    "He drove to the stadium.",
]

# 2. Calculate embeddings by calling model.encode()
embeddings = model.encode(sentences)
print(embeddings.shape)
# [3, 384]

# 3. Calculate the embedding similarities
similarities = model.similarity(embeddings, embeddings)
print(similarities)
# tensor([[1.0000, 0.6660, 0.1046],
#         [0.6660, 1.0000, 0.1411],
#         [0.1046, 0.1411, 1.0000]])
```

SimCSE

- Simple Contrastive Learning of Sentence Embeddings



Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+) / \tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+) / \tau}}$$

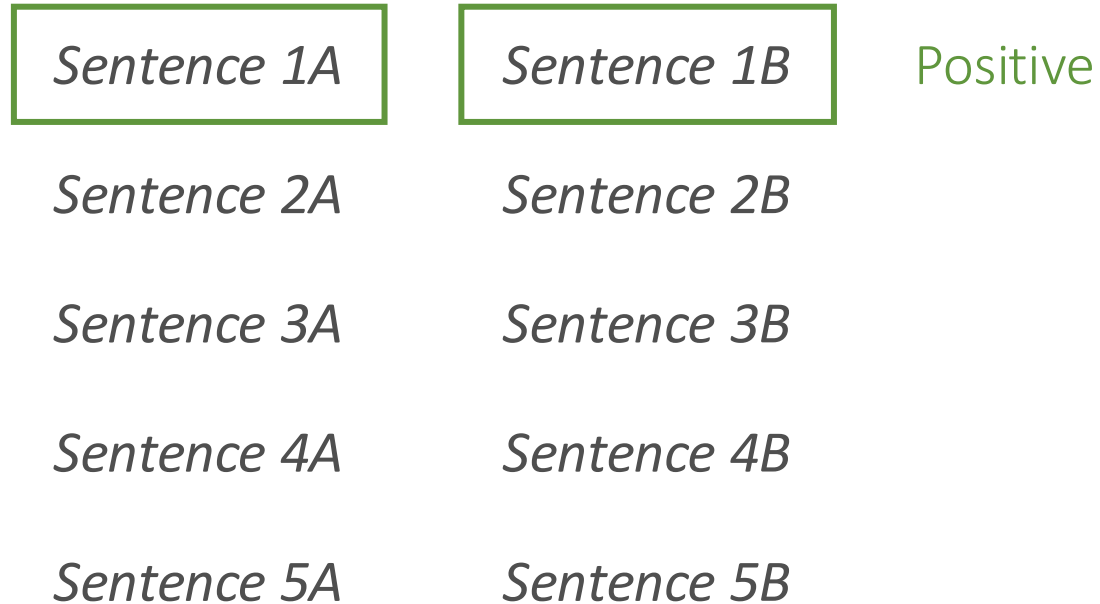
SimCSE: Supervised Contrastive Learning

<i>Sentence 1A</i>	<i>Sentence 1B</i>	Positive
<i>Sentence 2A</i>	<i>Sentence 2B</i>	Negative
<i>Sentence 3A</i>	<i>Sentence 3B</i>	Negative
<i>Sentence 4A</i>	<i>Sentence 4B</i>	Negative
<i>Sentence 5A</i>	<i>Sentence 5B</i>	Negative

Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

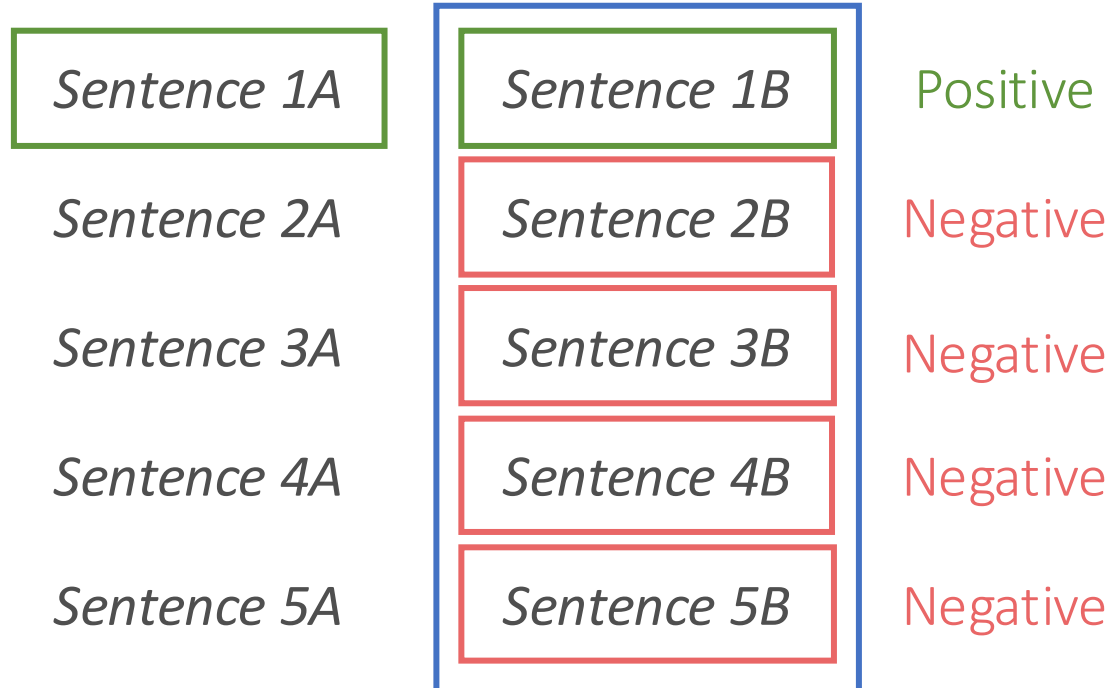
SimCSE: Supervised Contrastive Learning



Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

SimCSE: Supervised Contrastive Learning



Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

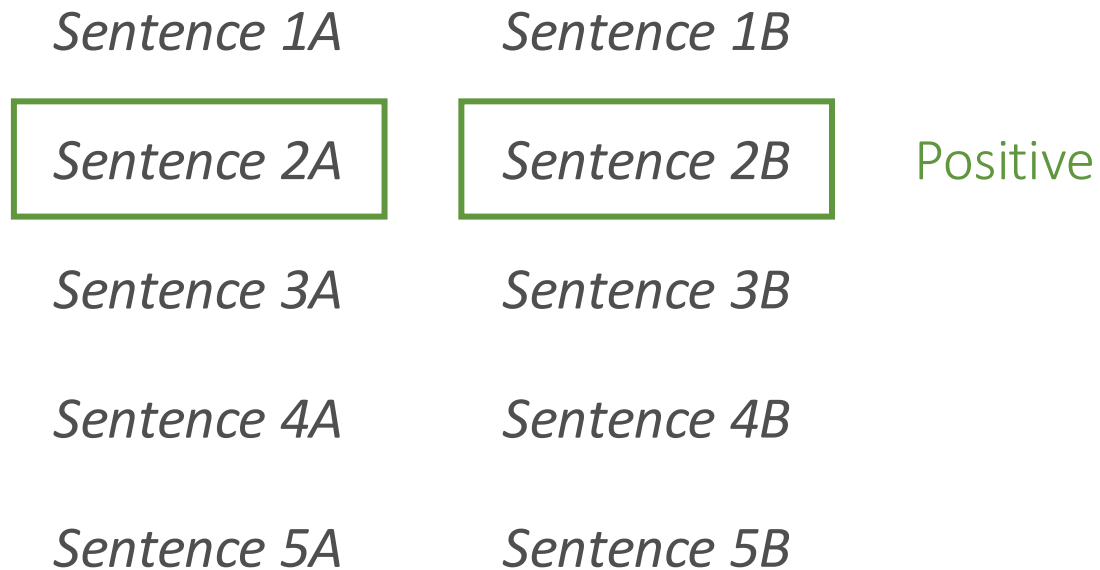
SimCSE: Supervised Contrastive Learning

Sentence 1A	Sentence 1B	Negative
Sentence 2A	Sentence 2B	Positive
Sentence 3A	Sentence 3B	Negative
Sentence 4A	Sentence 4B	Negative
Sentence 5A	Sentence 5B	Negative

Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

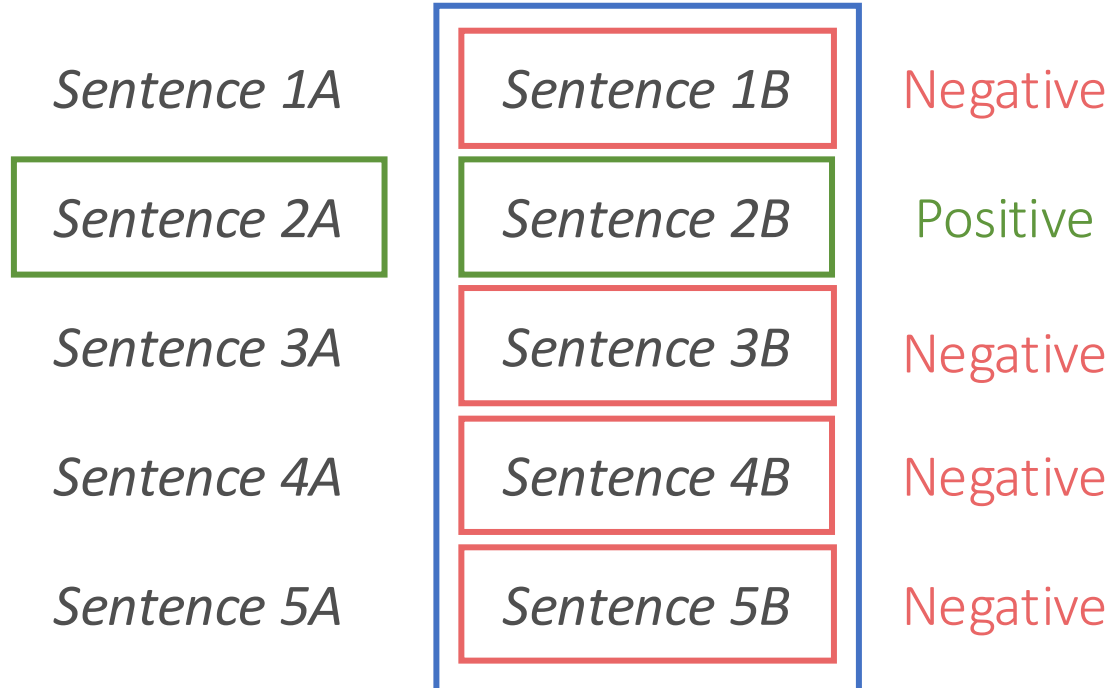
SimCSE: Supervised Contrastive Learning



Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

SimCSE: Supervised Contrastive Learning

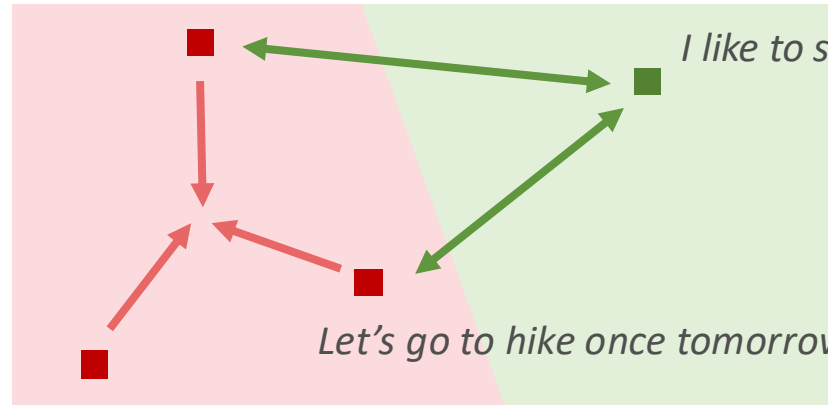


Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

SimCSE: Supervised Contrastive Learning

If it is sunny tomorrow, we will go hiking.



I like to study large language models.

Let's go to hike once tomorrow is sunny.

We will go hiking if tomorrow is a sunny day.

SimCSE: Performance

<i>Supervised models</i>								
InferSent-GloVe [♣]	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder [♣]	64.49	67.80	64.61	76.83	73.18	74.92	76.69	71.22
SBERT _{base} [♣]	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT _{base} -flow	69.78	77.27	74.35	82.01	77.46	79.12	76.21	76.60
SBERT _{base} -whitening	69.65	77.57	74.66	82.27	78.39	79.52	76.91	77.00
CT-SBERT _{base}	74.84	83.20	78.07	83.84	77.93	81.46	76.42	79.39
* SimCSE-BERT _{base}	75.30	84.67	80.19	85.40	80.82	84.25	80.39	81.57
SRoBERTa _{base} [♣]	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa _{base} -whitening	70.46	77.07	74.46	81.64	76.43	79.49	76.65	76.60
* SimCSE-RoBERTa _{base}	76.53	85.21	80.95	86.03	82.57	85.83	80.50	82.52
* SimCSE-RoBERTa _{large}	77.46	87.27	82.36	86.66	83.93	86.70	81.95	83.76

Supervised Data

- Consider SNLI dataset
 - Stanford Natural Language Inference

A boy is jumping on skateboard in the middle of a red bridge.

The boy skates down the sidewalk.

Contradiction

A boy is jumping on skateboard in the middle of a red bridge.

The boy is wearing safety equipment.

Neutral

A boy is jumping on skateboard in the middle of a red bridge.

The boy does a skateboarding trick.

Entailment

What if we don't have supervised data?

SimCSE: Unsupervised Contrastive Learning

<i>Sentence 1</i>	<i>Sentence 1'</i>	Positive
<i>Sentence 2</i>	Negative	
<i>Sentence 3</i>	Negative	
<i>Sentence 4</i>	Negative	
<i>Sentence 5</i>	Negative	

Contrastive Loss

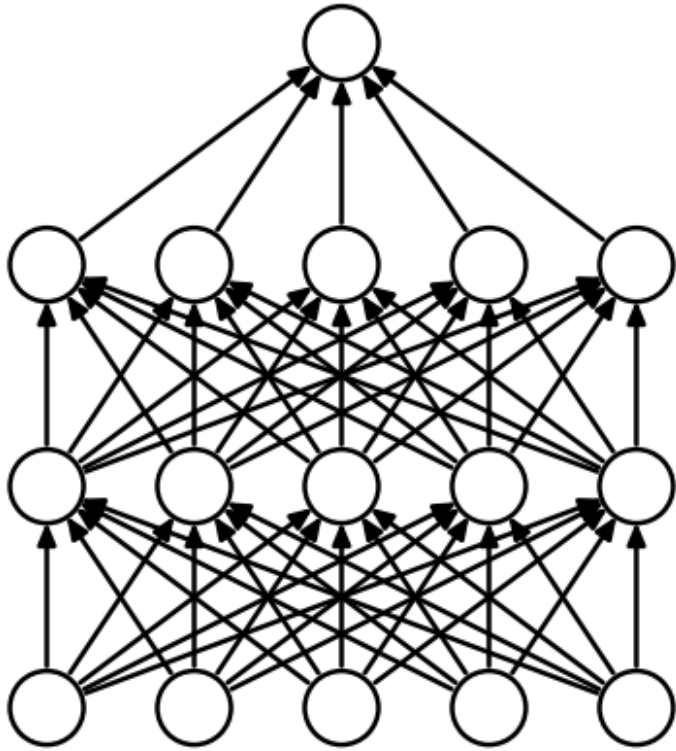
$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

Generate positive example with masking

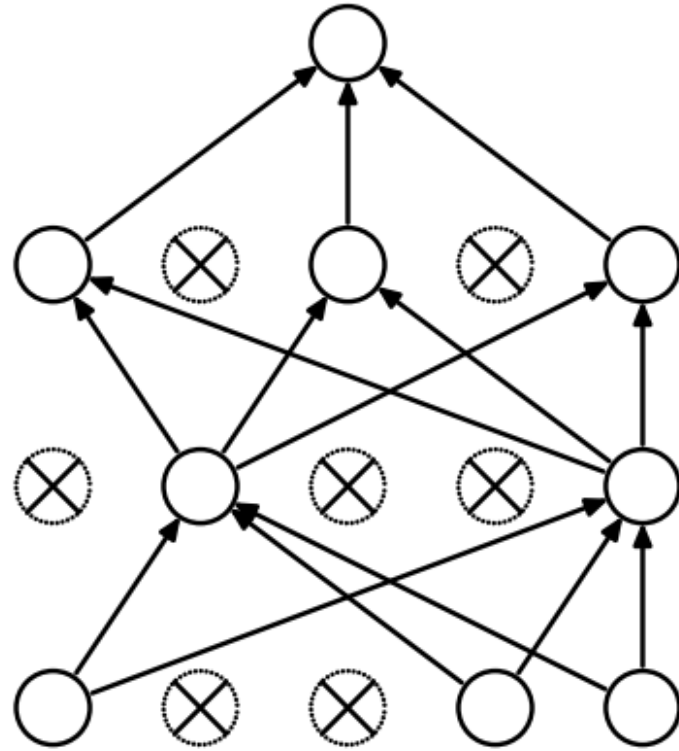
If it is sunny tomorrow, we will go hiking.

If [mask] is sunny tomorrow, we [mask] go hiking.

Dropout



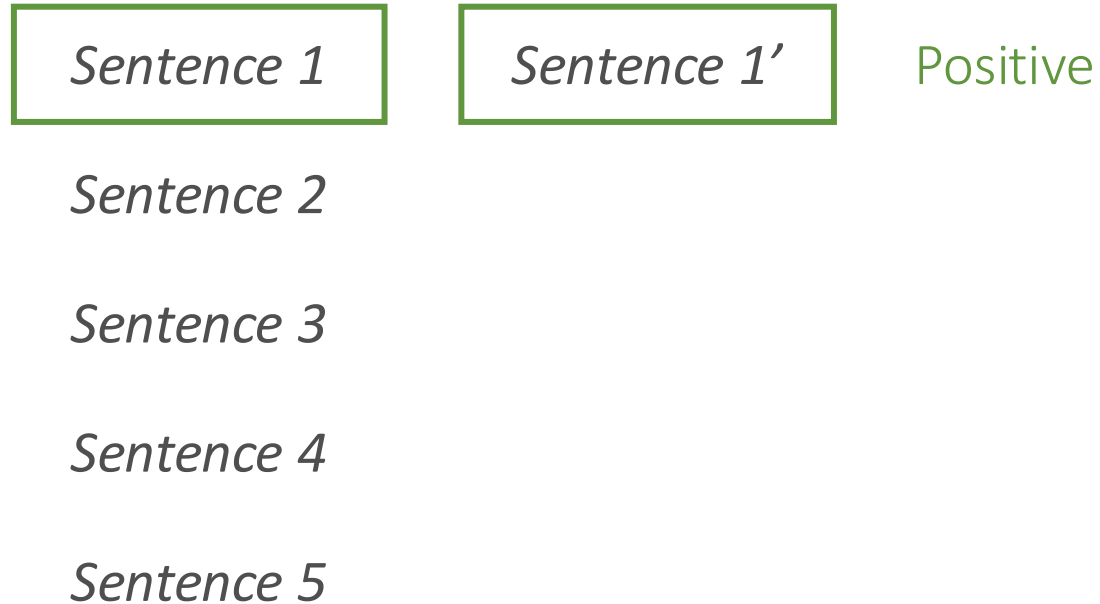
(a) Standard Neural Net



(b) After applying dropout.

Generate positive example with neuron masking

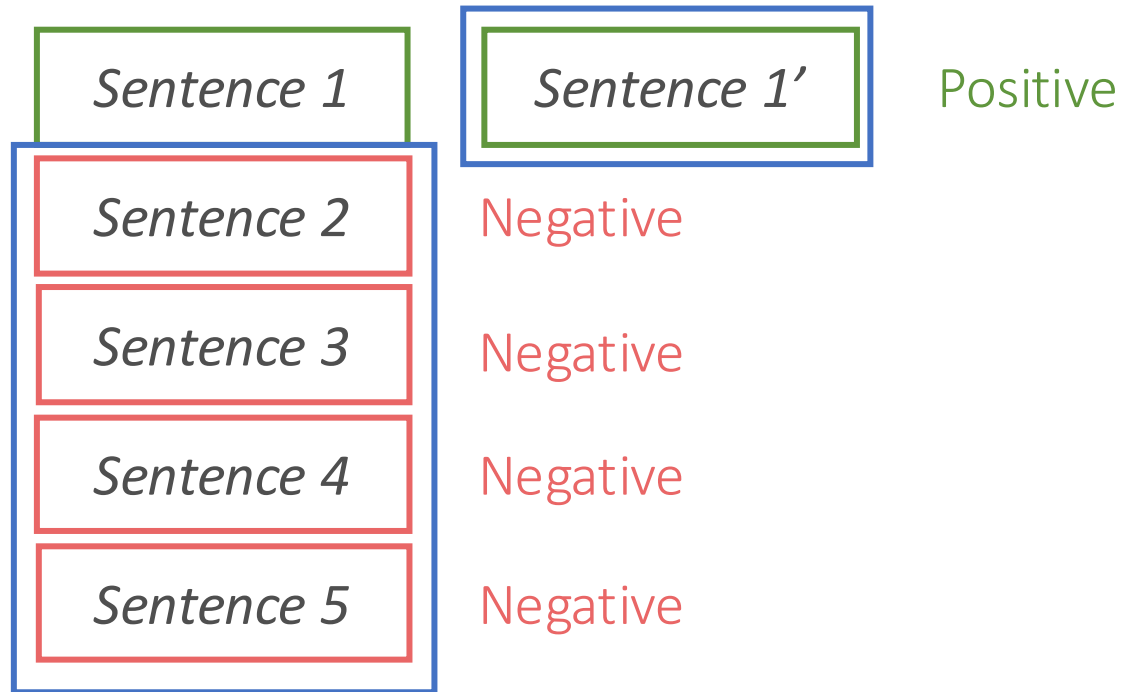
SimCSE: Unsupervised Contrastive Learning



Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

SimCSE: Unsupervised Contrastive Learning



Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

SimCSE: Unsupervised Contrastive Learning

Sentence 1

Sentence 2

Sentence 2'

Positive

Sentence 3

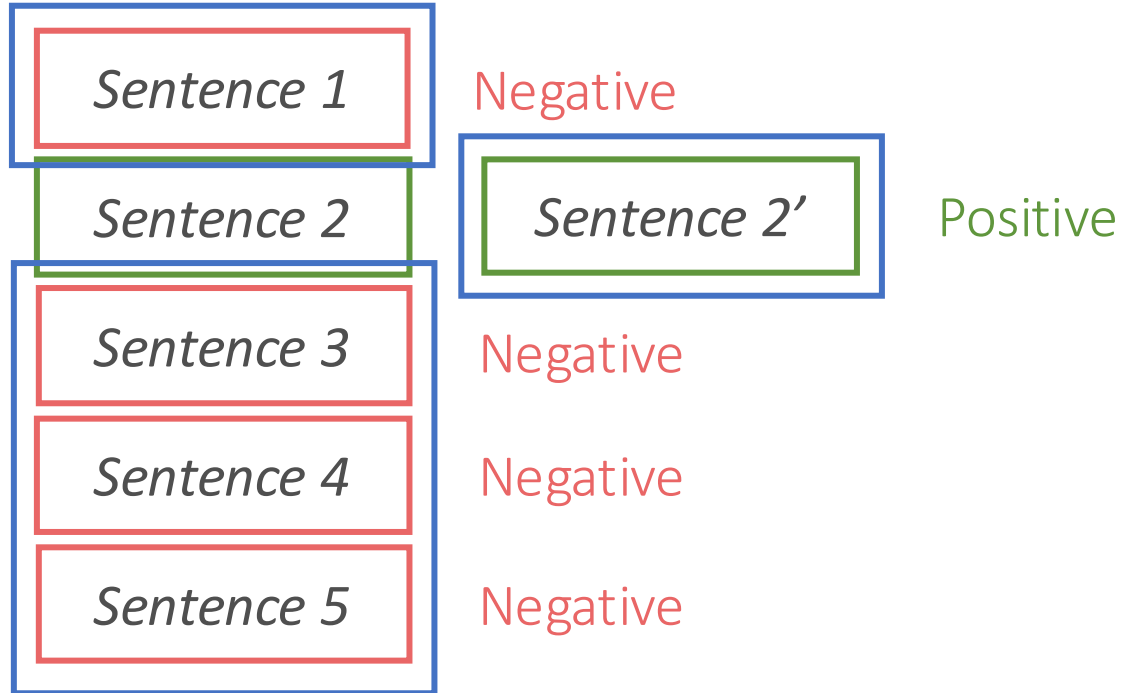
Sentence 4

Sentence 5

Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

SimCSE: Unsupervised Contrastive Learning



Contrastive Loss

$$l_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}$$

SimCSE: Performance

Model	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
<i>Unsupervised models</i>								
GloVe embeddings (avg.) [♣]	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
BERT _{base} (first-last avg.)	39.70	59.38	49.67	66.03	66.19	53.87	62.06	56.70
BERT _{base} -flow	58.40	67.10	60.85	75.16	71.22	68.66	64.47	66.55
BERT _{base} -whitening	57.83	66.90	60.90	75.08	71.31	68.24	63.73	66.28
IS-BERT _{base} [♡]	56.77	69.24	61.21	75.23	70.16	69.21	64.25	66.58
CT-BERT _{base}	61.63	76.80	68.47	77.50	76.48	74.31	69.19	72.05
* SimCSE-BERT _{base}	68.40	82.41	74.38	80.91	78.56	76.85	72.23	76.25
RoBERTa _{base} (first-last avg.)	40.88	58.74	49.07	65.63	61.48	58.55	61.63	56.57
RoBERTa _{base} -whitening	46.99	63.24	57.23	71.36	68.99	61.36	62.91	61.73
DeCLUTR-RoBERTa _{base}	52.41	75.19	65.52	77.12	78.63	72.41	68.62	69.99
* SimCSE-RoBERTa _{base}	70.16	81.77	73.24	81.36	80.65	80.22	68.56	76.57
* SimCSE-RoBERTa _{large}	72.86	83.99	75.62	84.77	81.80	81.98	71.26	78.90
<i>Supervised models</i>								
InferSent-GloVe [♣]	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder [♣]	64.49	67.80	64.61	76.83	73.18	74.92	76.69	71.22
SBERT _{base} [♣]	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT _{base} -flow	69.78	77.27	74.35	82.01	77.46	79.12	76.21	76.60
SBERT _{base} -whitening	69.65	77.57	74.66	82.27	78.39	79.52	76.91	77.00
CT-SBERT _{base}	74.84	83.20	78.07	83.84	77.93	81.46	76.42	79.39
* SimCSE-BERT _{base}	75.30	84.67	80.19	85.40	80.82	84.25	80.39	81.57
SRoBERTa _{base} [♣]	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa _{base} -whitening	70.46	77.07	74.46	81.64	76.43	79.49	76.65	76.60
* SimCSE-RoBERTa _{base}	76.53	85.21	80.95	86.03	82.57	85.83	80.50	82.52
* SimCSE-RoBERTa _{large}	77.46	87.27	82.36	86.66	83.93	86.70	81.95	83.76

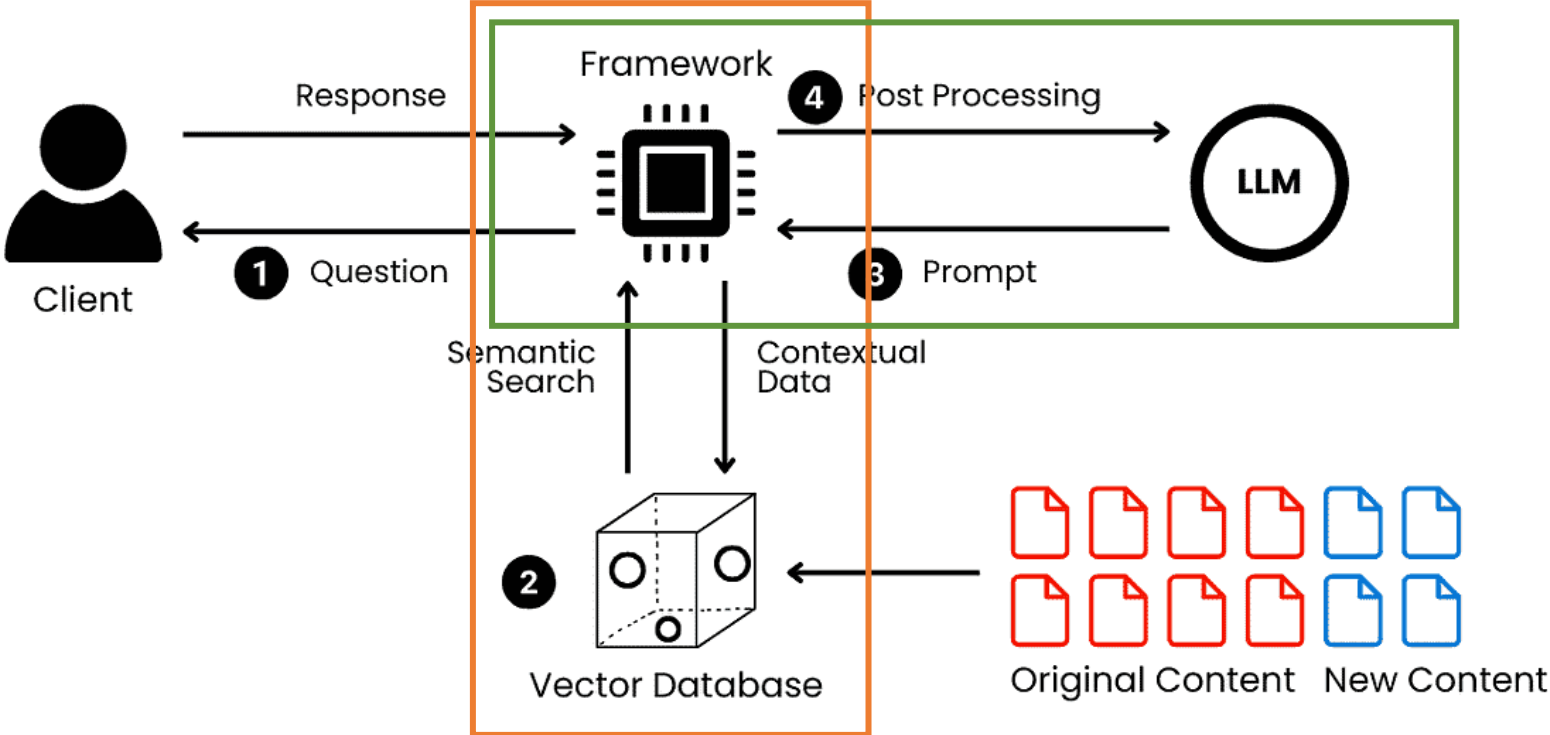
Lecture Plan

- Text Similarity
- Retrieval-Augmented Generation

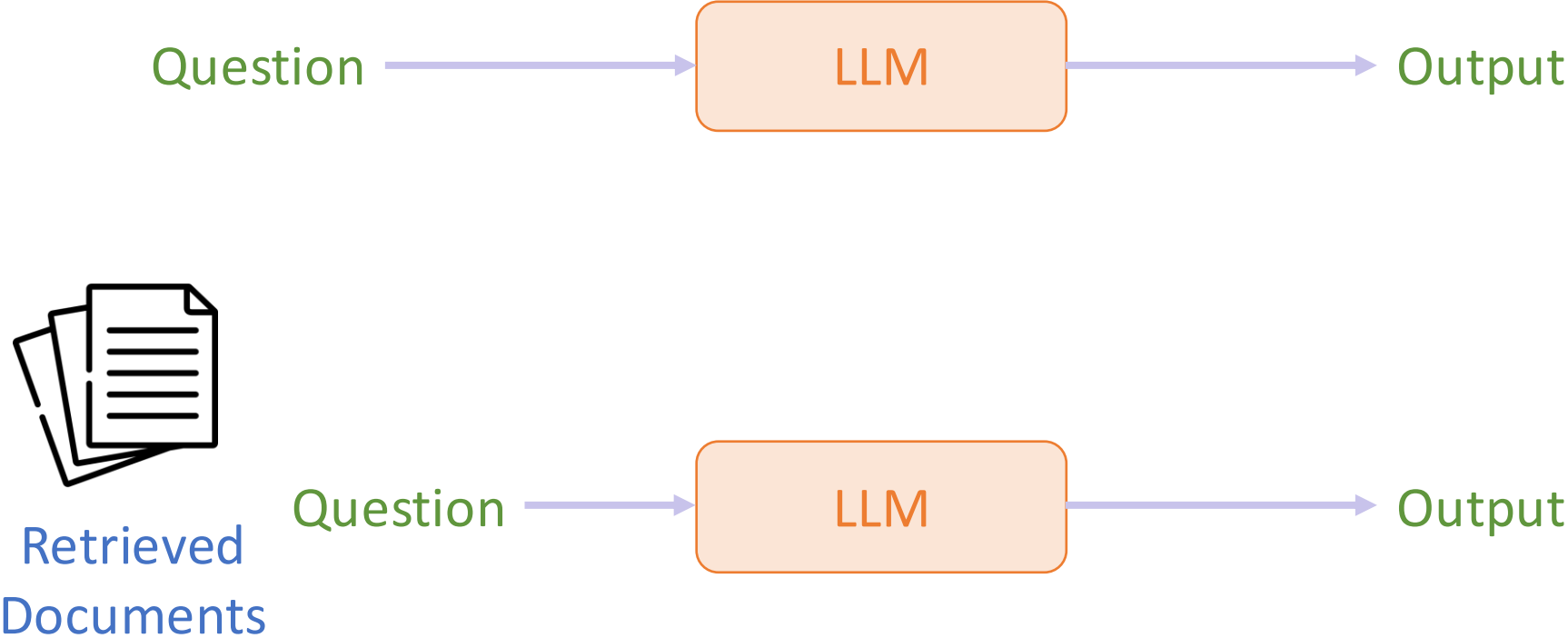
Retrieval-Augmented Generation (RAG)

- Information Retrieval + Large Language Models

RAG Architecture Model



Retrieval-Augmented Generation (RAG)



Why RAG?

- Knowledge Cutoff

What is 2026 WBC US vs Canada result?

Assistant B



There is **no result** yet for a 2026 WBC game between the USA and Canada because **the tournament has not happened yet.**

The 2026 World Baseball Classic is scheduled to take place in **March 2026.**

Here is what we know so far:

- **Status:** The brackets and schedule have not been announced.
- **Qualification:** Both Team USA and Team Canada automatically qualified for the 2026 tournament based on their performance in the 2023 WBC.
- **Likelihood:** While they are often placed in the same pool due to geography, it is not guaranteed they will play each other in the group stage until the official schedule is released (likely in late 2025).

Team USA beats Canada to set up WBC semis showdown vs. D.R.



Jeff Passan
Mar 13, 2026, 11:21 PM ET

Follow

Share Like

3.3K

HOUSTON -- In the aftermath of the Dominican Republic's mercy-rule win over Korea to advance to the World Baseball Classic semifinals, Dominican fans at LoanDepot Park in Miami made clear what they wanted to see next, chanting: "Queremos a USA" -- we want the USA.

They got what they asked for.

The United States booked its spot in the WBC semifinals with a 5-3 victory against Canada on Friday night, setting up a dream clash between the two most talented teams in the tournament. National League Cy Young winner [Paul Skenes](#) will start for Team USA against the formidable Dominican lineup that has wrecked opposing pitchers, while [Athletics](#) right-hander [Luis Severino](#) is expected to face an American lineup big on names but not so much on production thus far.

Why RAG?

- Domain Knowledge

what is the name of this molecule:

CCOCC(=O)CN[C@H]

(C(=O)N1CC[C@H]1C(=O)NCc1ccc(/C(N)=N/O)cc1)C1CCCCC1

The SMILES you gave:

CCOCC(=O)CN[C@H](C(=O)N1CC[C@H]1C(=O)NCc1ccc(/C(N)=N/O)cc1)C1CCCCC1



corresponds to the **direct thrombin inhibitor drug**:

Dabigatran

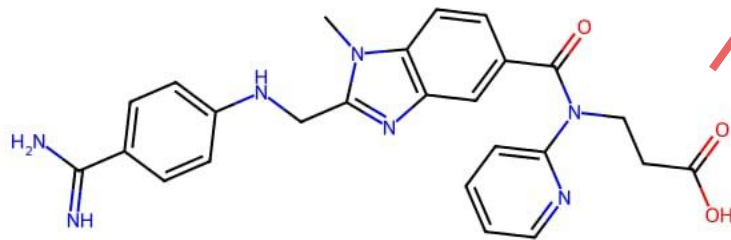
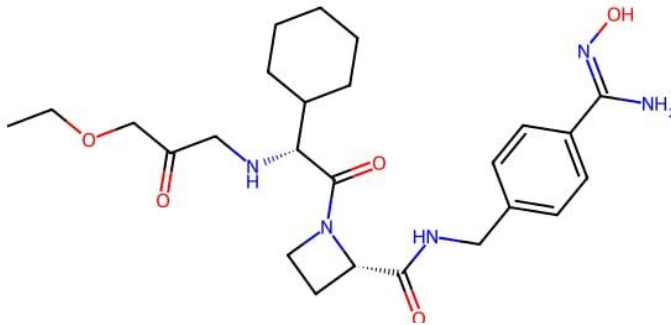
More precisely, the structure matches **dabigatran (the active drug)** rather than the orally used prodrug

Dabigatran etexilate.

Key structural features that identify it


From the SMILES:

- **Amidoxime group:** /C(N)=N/O
 - characteristic benzamidine bioisostere used in thrombin inhibitors
- **p-substituted benzyl amide**
- **proline-like tertiary amide ring:** N1CC[C@H]1 ↓
- **cyclohexyl substituent:** C1CCCCC1




Why RAG?


- Hallucination



Who was the first person to walk on the moon?

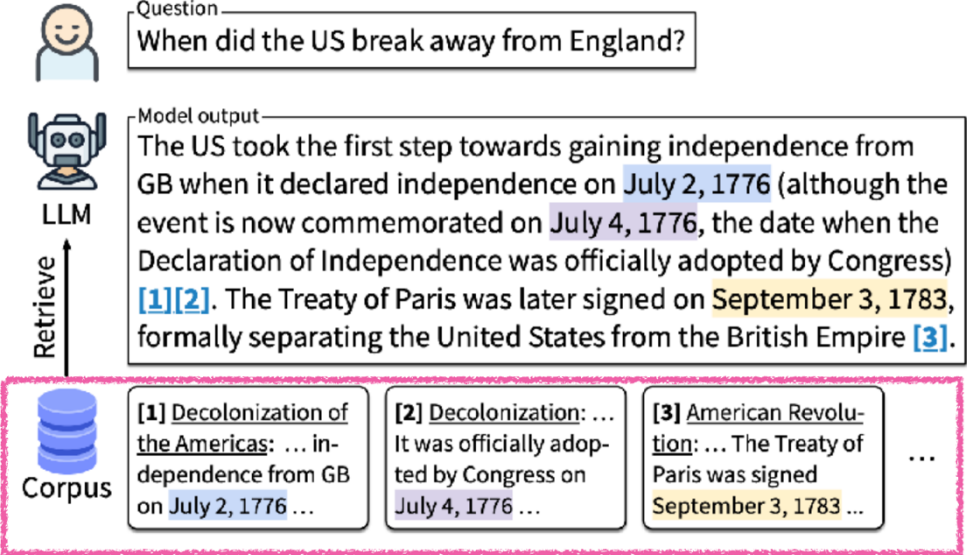


Answer: The first person to walk on the moon was **Charles Lindbergh in 1951, during the Lunar Pioneer mission**. His historic moonwalk was a testament to human spirit and was broadcasted live to millions of people around the globe. ❌



Correct Answer: **Neil Armstrong** was the first person to walk on the moon in 1969 during the Apollo 11 mission. ✅

Generating text with citations



(Nakano et al. 2021; Menick et al., 2022; Gao et al., 2023)

Why RAG?

- Private data

Extracting Training Data from Large Language Models

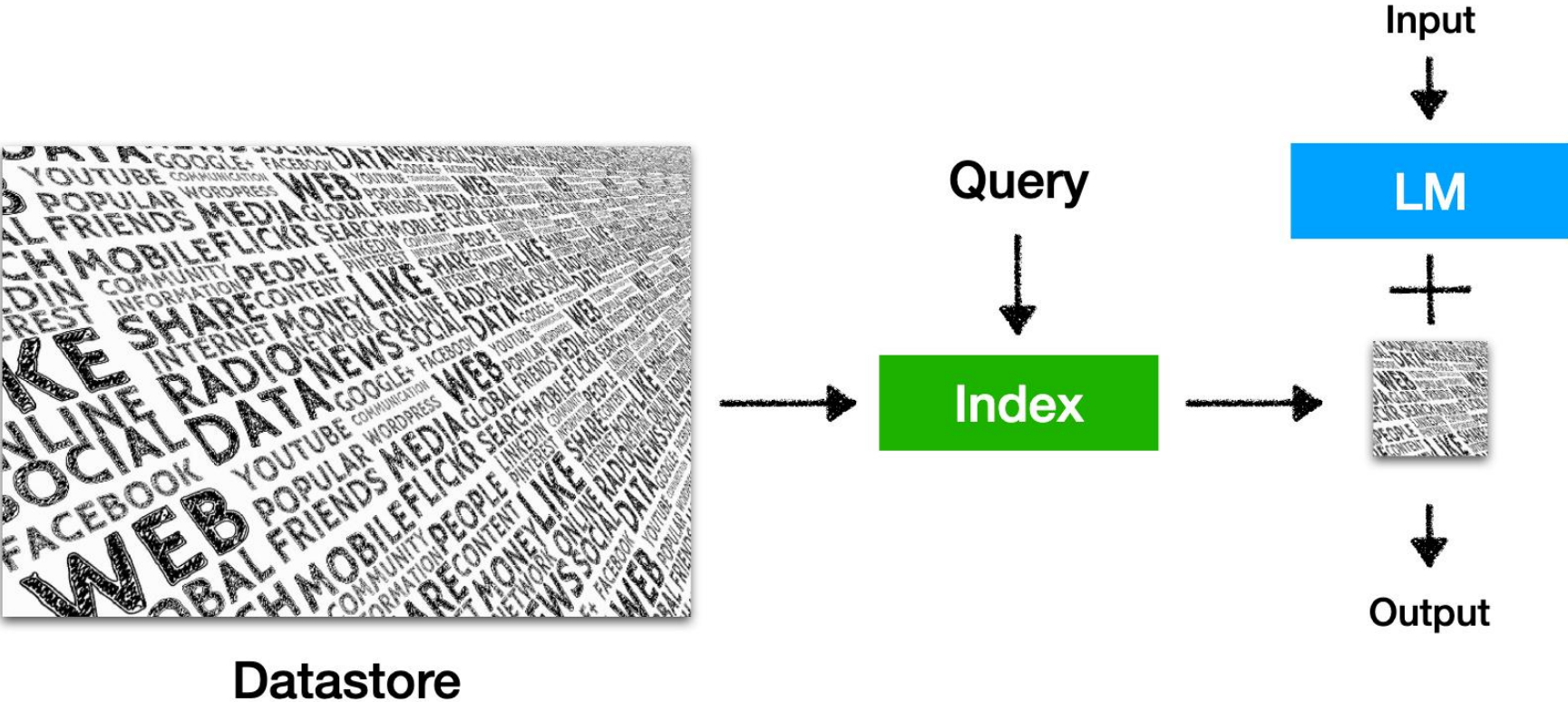
Nicholas Carlini¹ Florian Tramèr² Eric Wallace³ Matthew Jagielski⁴
Ariel Herbert-Voss^{5,6} Katherine Lee¹ Adam Roberts¹ Tom Brown⁵
Dawn Song³ Úlfar Erlingsson⁷ Alina Oprea⁴ Colin Raffel¹
¹Google ²Stanford ³UC Berkeley ⁴Northeastern University ⁵OpenAI ⁶Harvard ⁷Apple

Category	Count
US and international news	109
Log files and error reports	79
License, terms of use, copyright notices	54
Lists of named items (games, countries, etc.)	54
Forum or Wiki entry	53
Valid URLs	50
Named individuals (non-news samples only)	46
Promotional content (products, subscriptions, etc.)	45
High entropy (UUIDs, base64 data)	35
Contact info (address, email, phone, twitter, etc.)	32
Code	31
Configuration files	30
Religious texts	25
Pseudonyms	15
Donald Trump tweets and quotes	12
Web forms (menu items, instructions, etc.)	11
Tech news	11
Lists of numbers (dates, sequences, etc.)	10

Why RAG?

- Potentially leverage other modalities
 - Knowledge base
 - Tabular data
 - Graphs
 - ...

Retrieval-Augmented Generation (RAG)



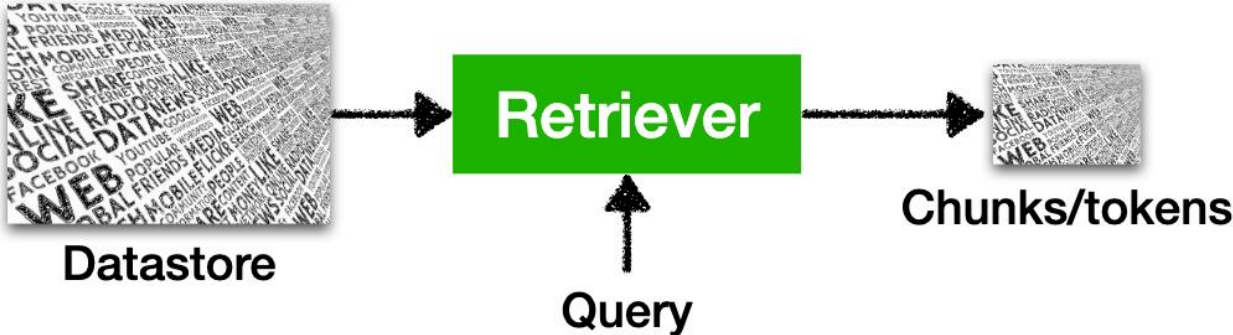
Retrieval-Augmented Generation (RAG)

Retrieval models and **language models** are trained **independently**

- Training language models

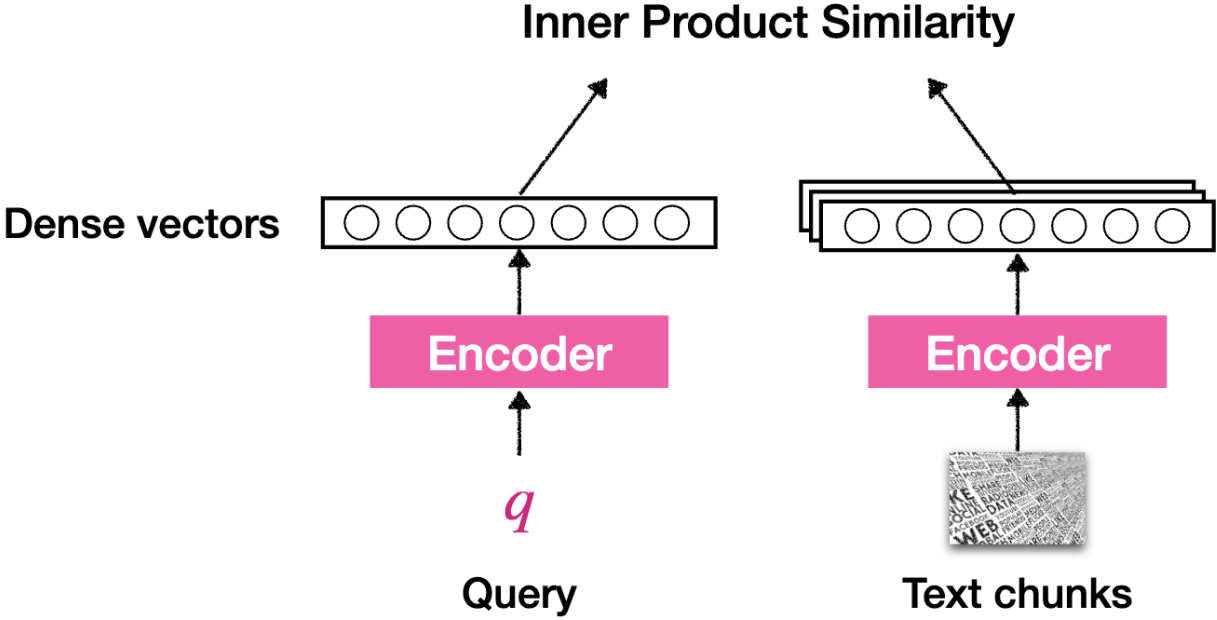


- Training retrieval models



How to Train A Retriever?

Dense retrieval models: DPR (Karpukhin et al. 2020)

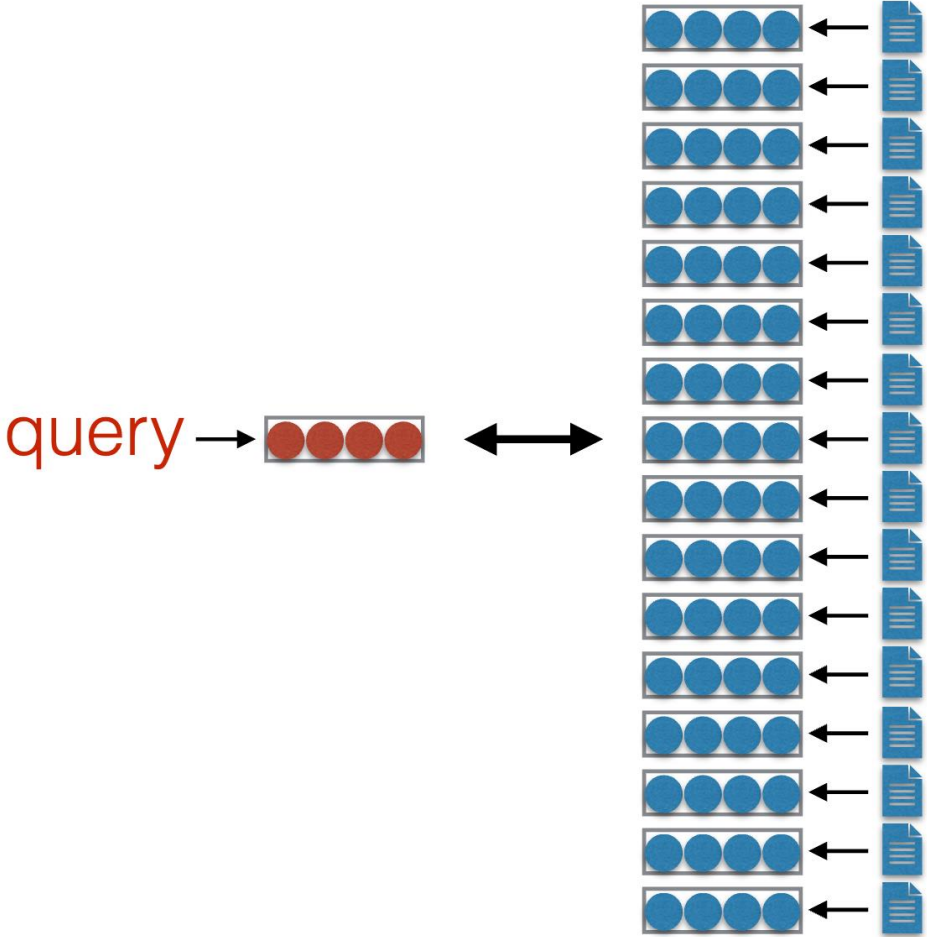


Dense Passage Retrieval (DPR)

Similarity between query and documents

Similarity between two sentences

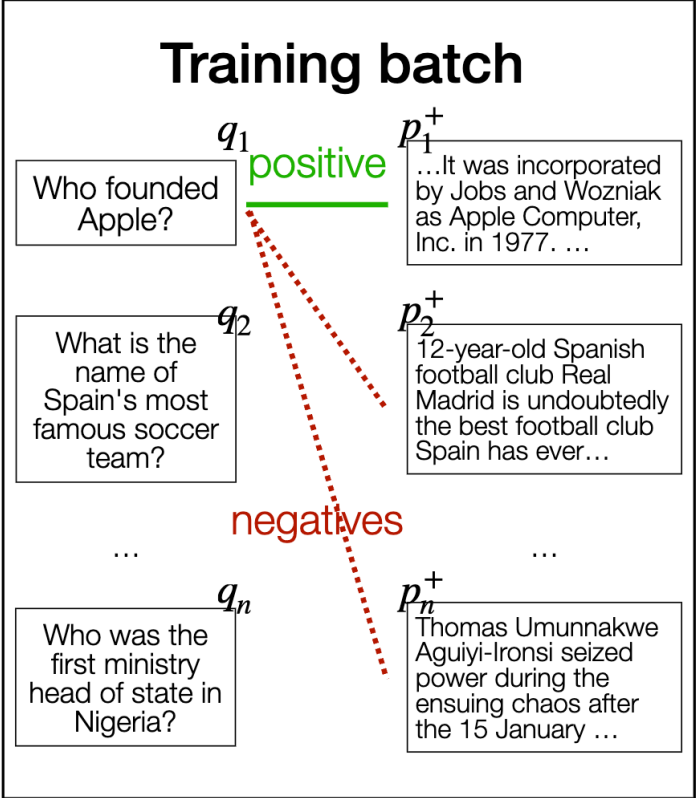
We will go hiking if tomorrow is a sunny day.
If it is sunny tomorrow, we will go hiking.



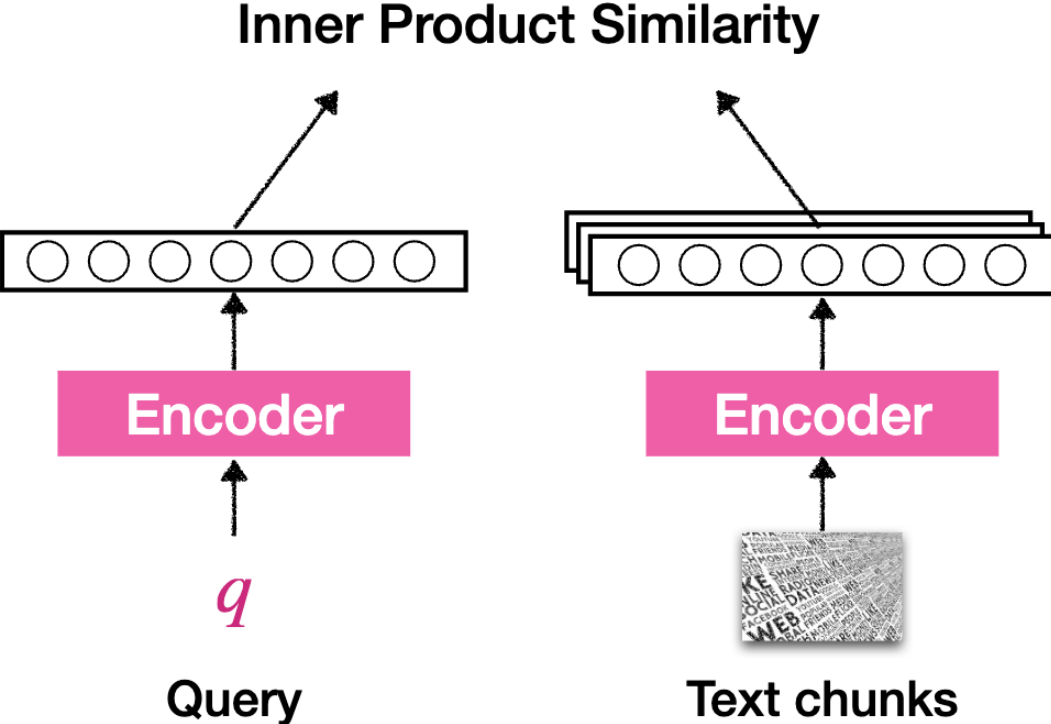
Dense Passage Retrieval (DPR)

$$L(q, p^+, p_1^-, p_2^-, \dots, p_n^-)$$

$$= -\log \frac{\exp(\text{sim}(q, p^+))}{\exp(\text{sim}(q, p^+)) + \sum_{j=1}^n \exp(\text{sim}(q, p_j^-))}$$



Dense Passage Retrieval (DPR)

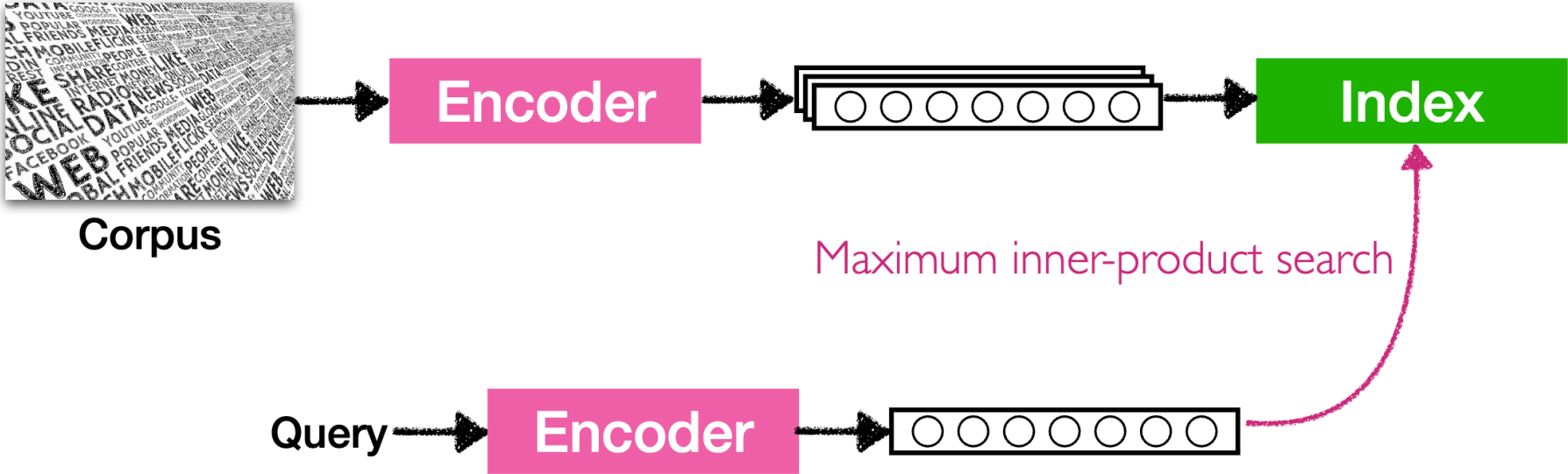


$$L(q, p^+, p_1^-, p_2^-, \dots, p_n^-) = -\log \frac{\exp(\text{sim}(q, p^+))}{\exp(\text{sim}(q, p^+)) + \sum_{j=1}^n \exp(\text{sim}(q, p_j^-))}$$

Contrastive learning



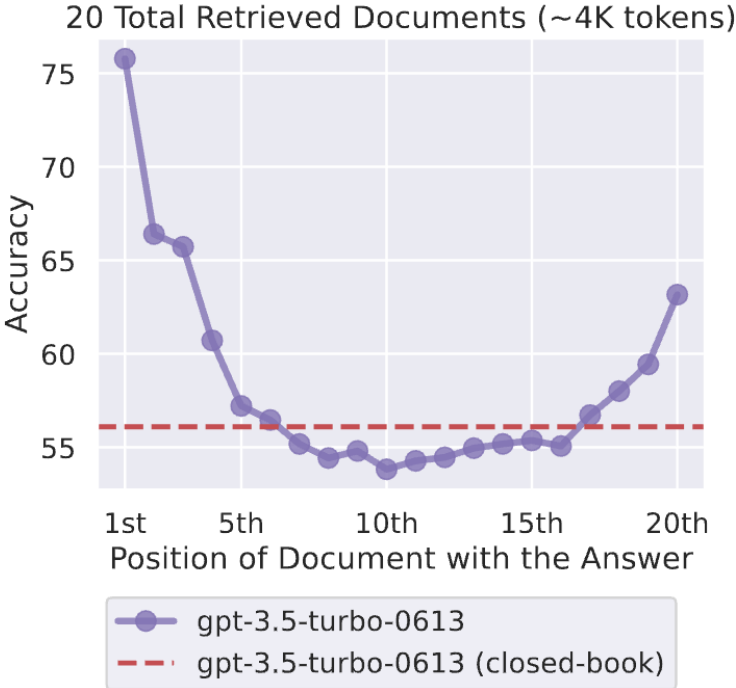
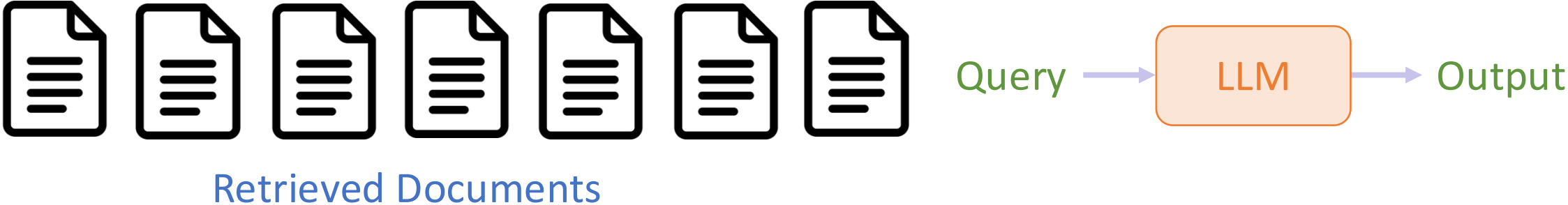
Dense Passage Retrieval (DPR)



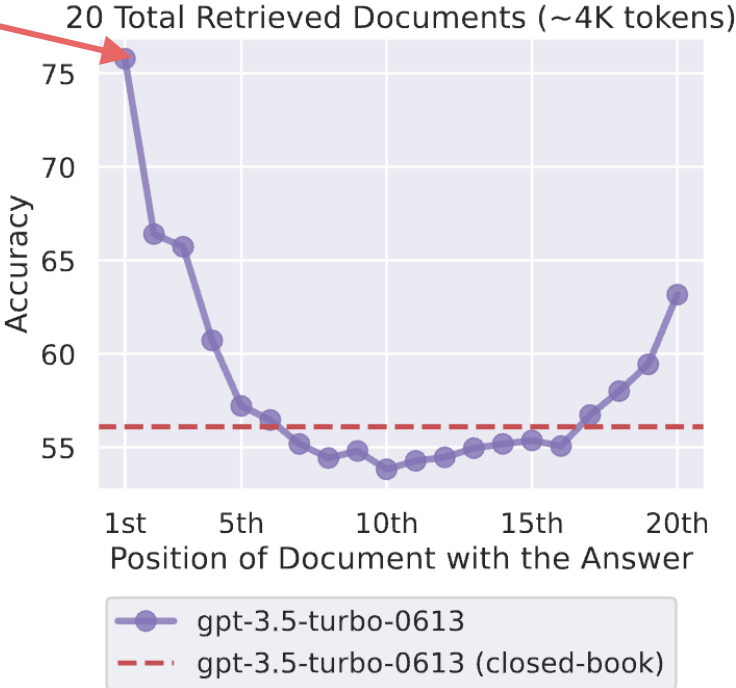
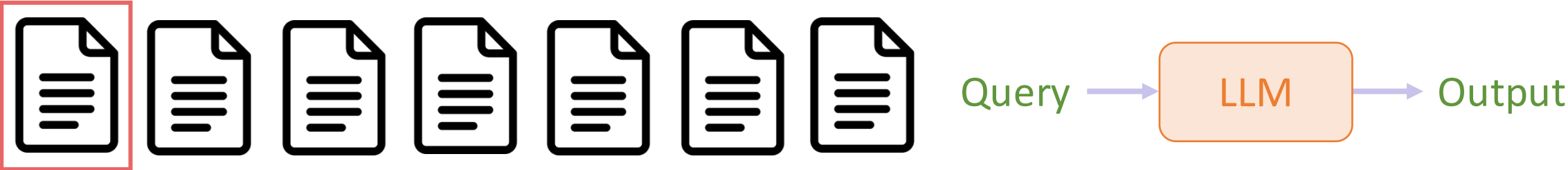
Challenges with RAG

- Longer input text
 - Length generalization
 - KV cache
- The lost-in-the-middle problem

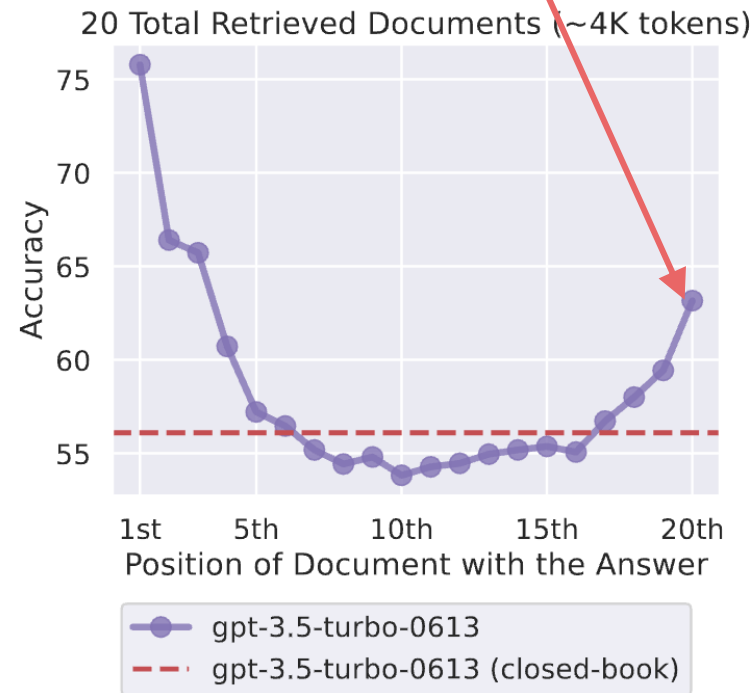
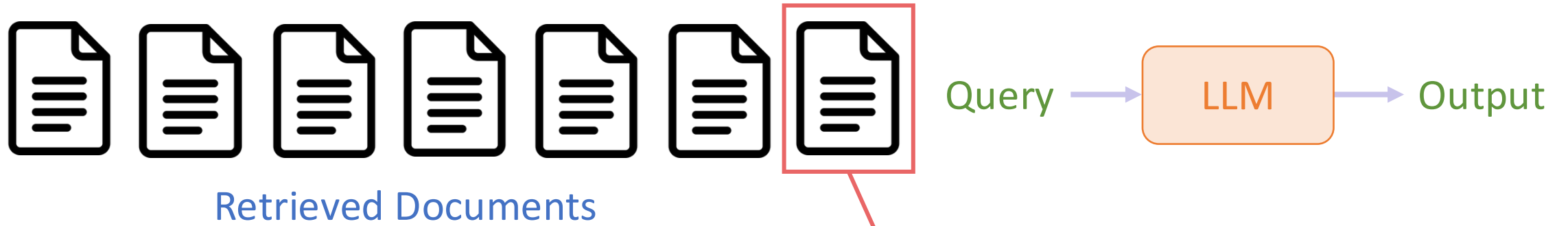
The Lost-in-the-Middle Problem



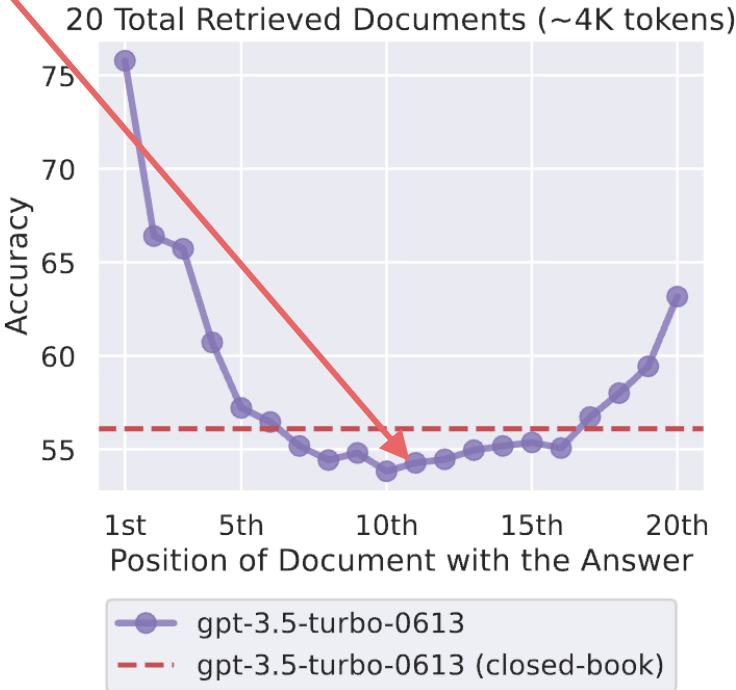
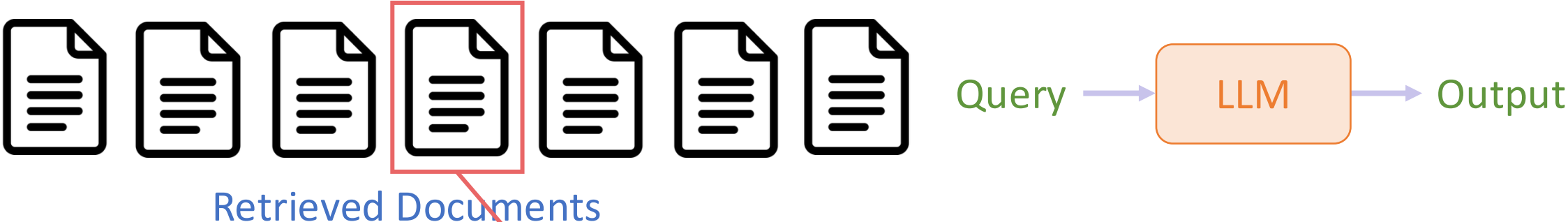
The Lost-in-the-Middle Problem



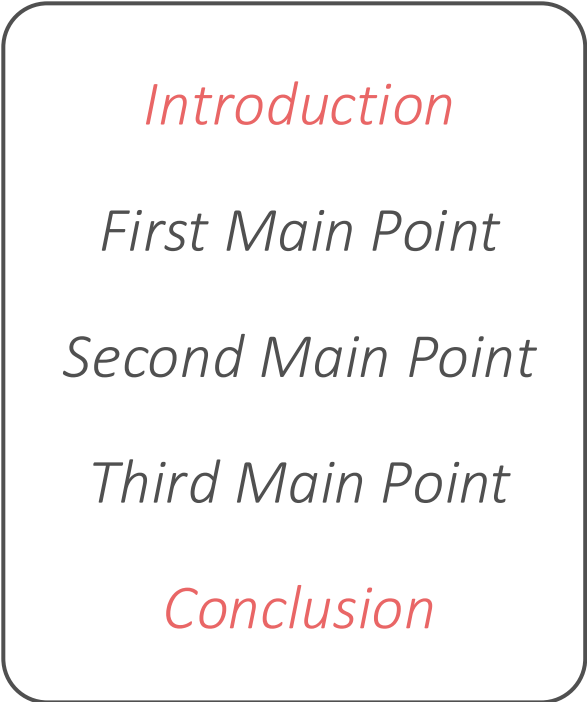
The Lost-in-the-Middle Problem



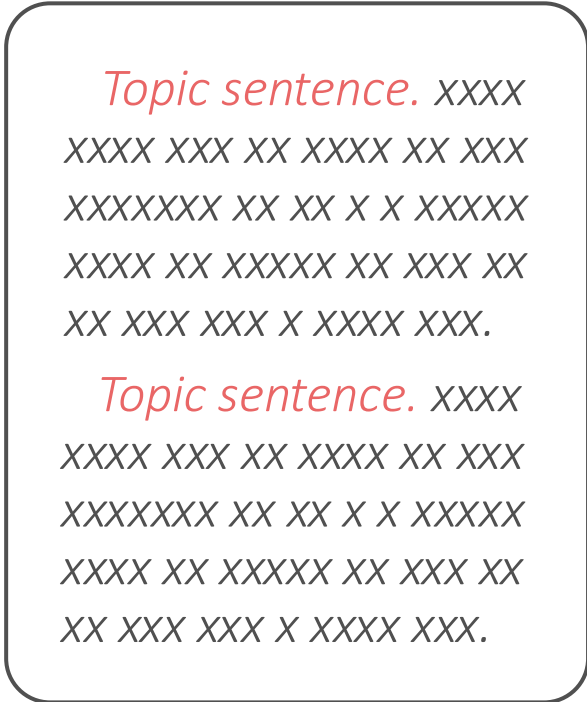
The Lost-in-the-Middle Problem



Reasons for Positional Bias: Pre-Training Data

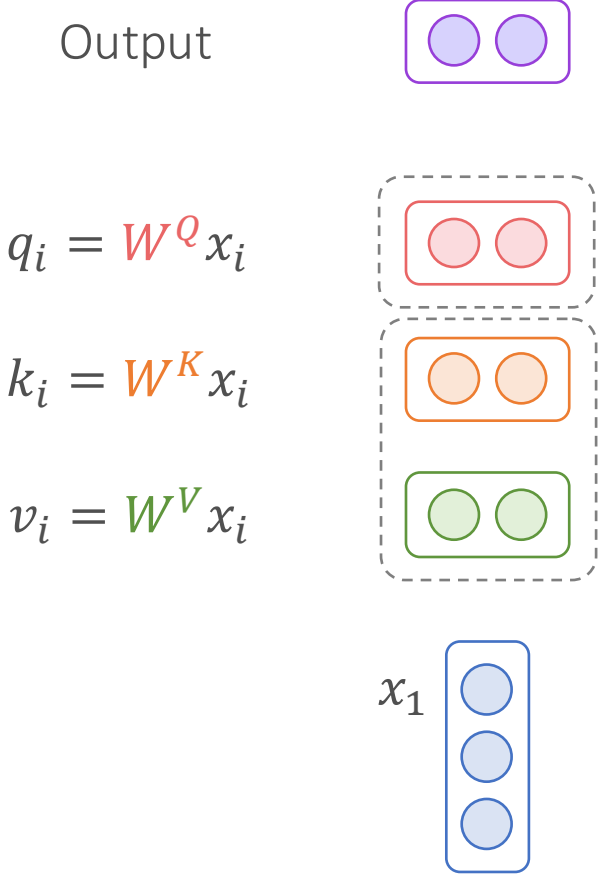


The 5 Paragraph Essay Outline

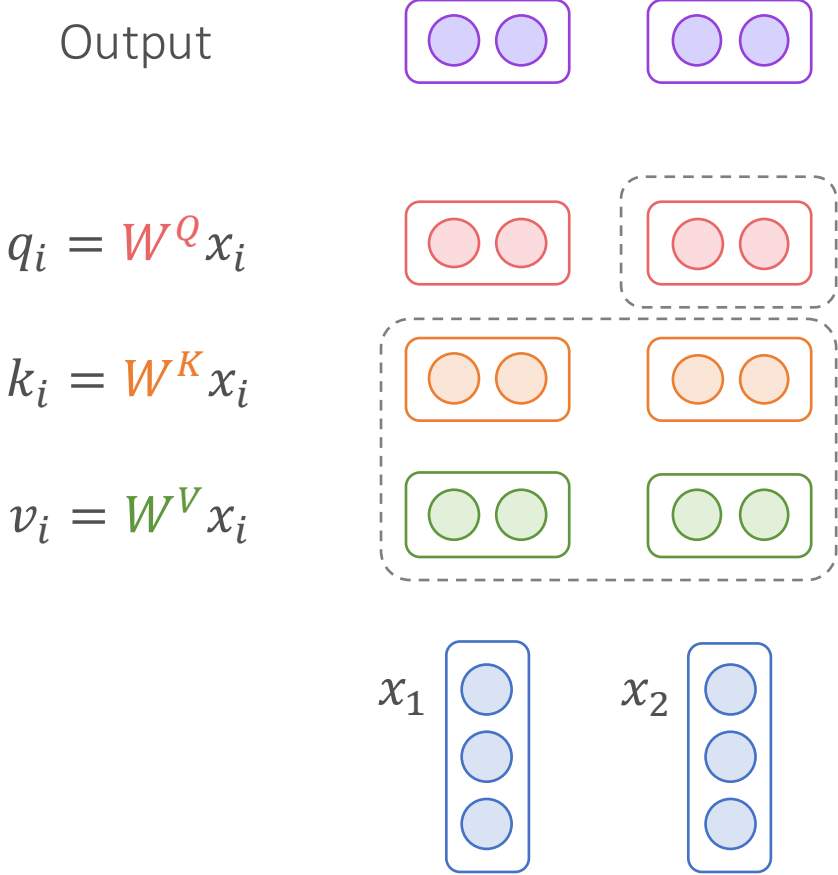


Topic Sentence

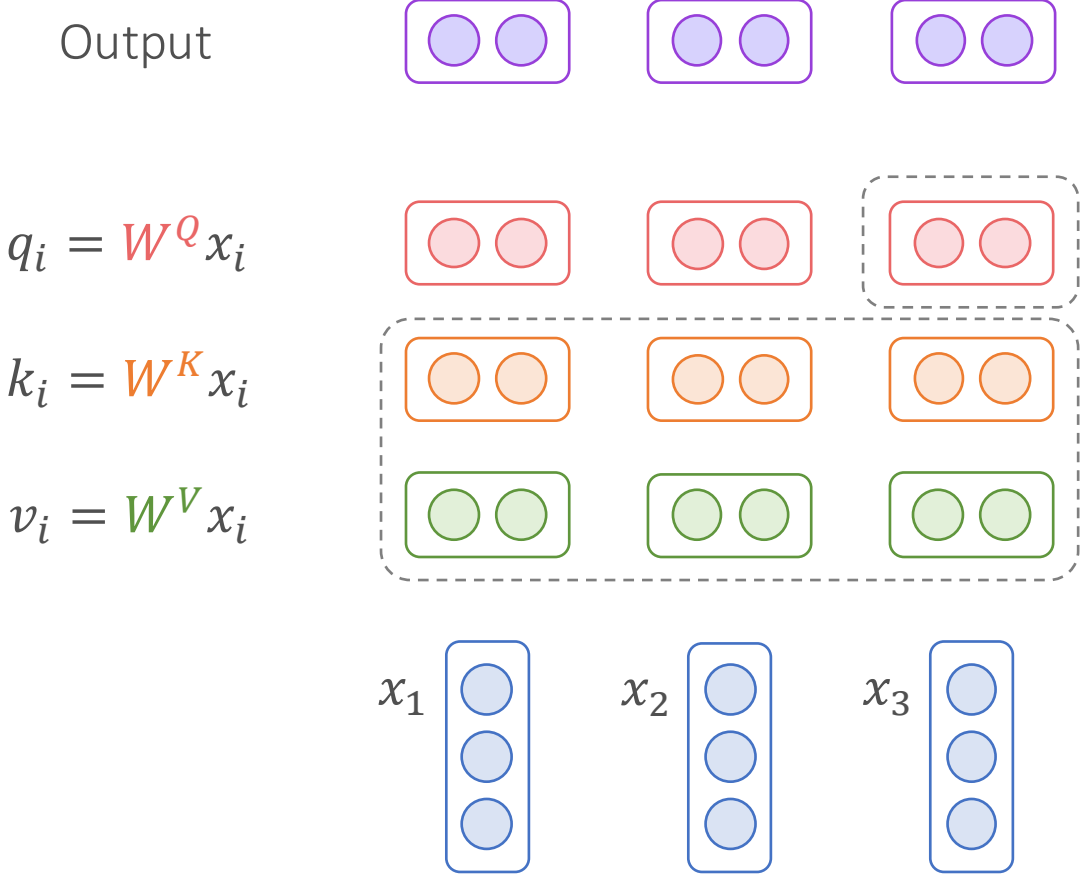
Reasons for Positional Bias: Attention Mechanism



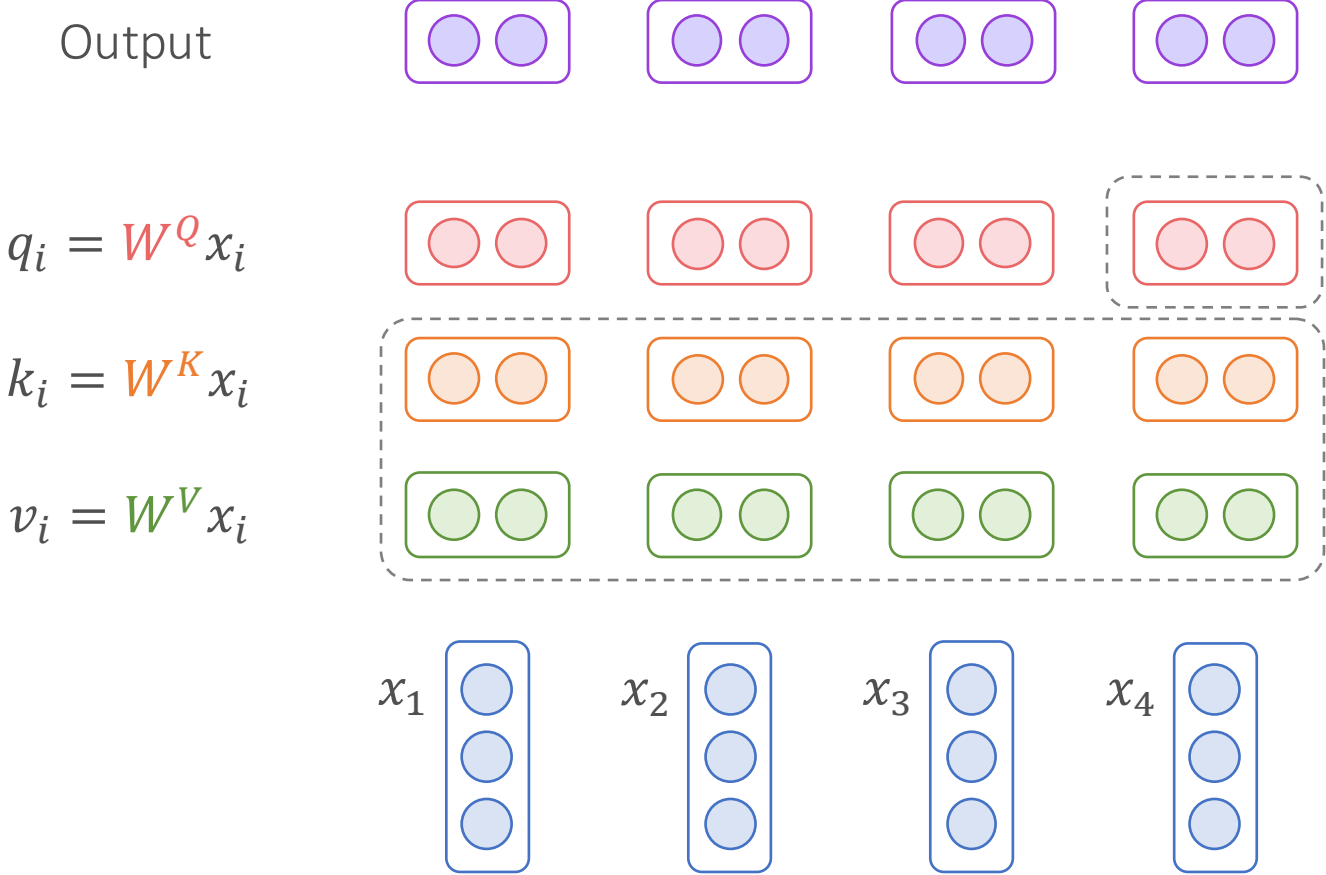
Reasons for Positional Bias: Attention Mechanism



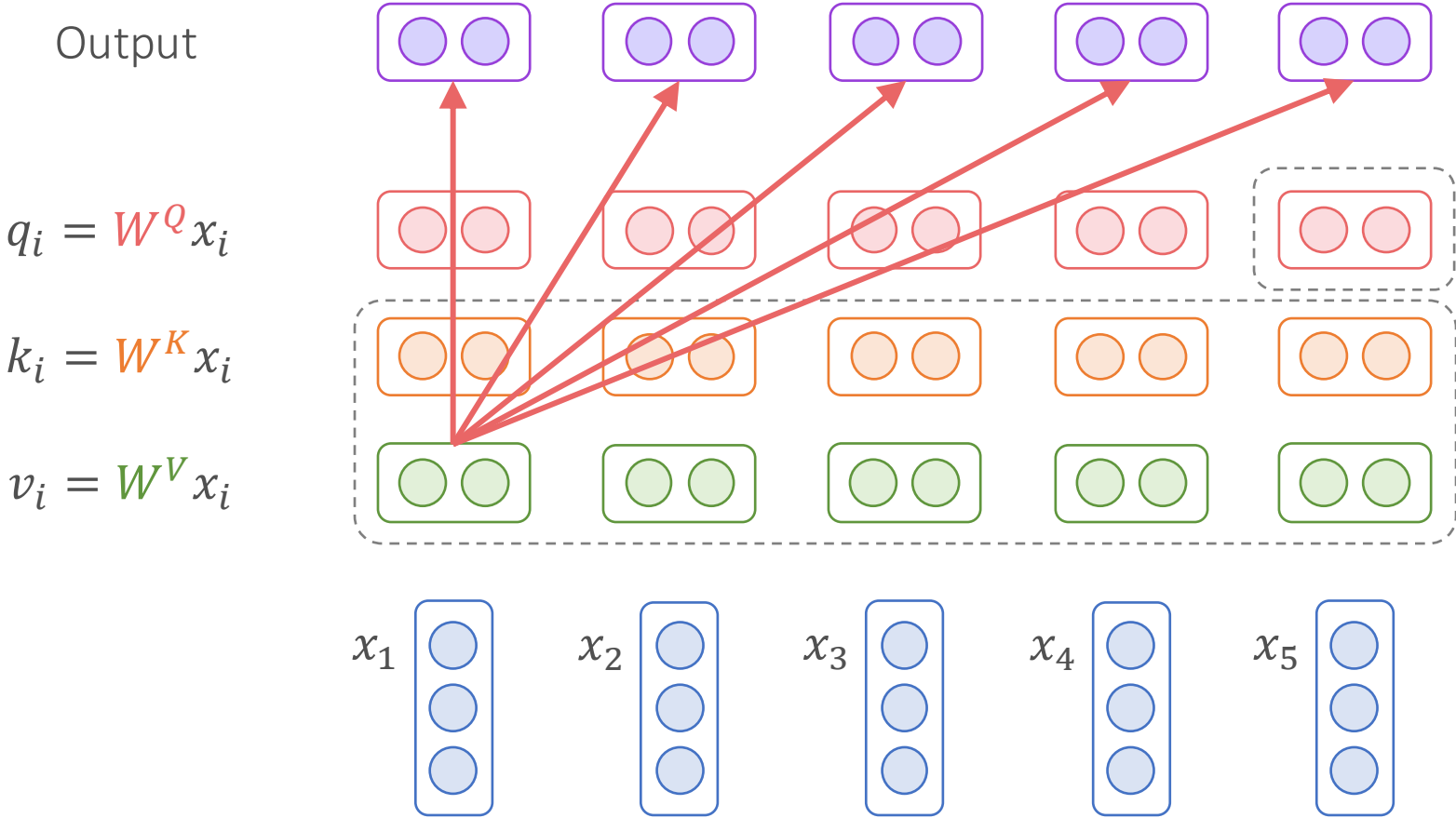
Reasons for Positional Bias: Attention Mechanism



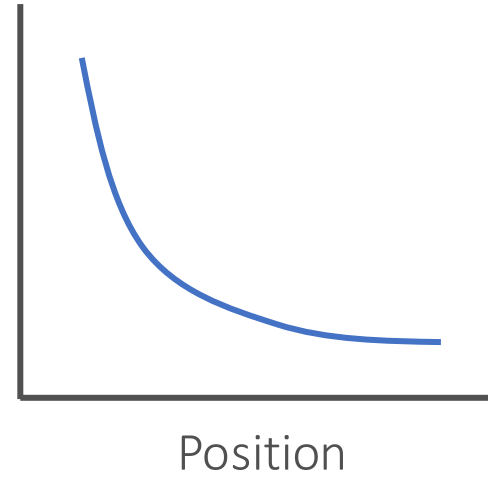
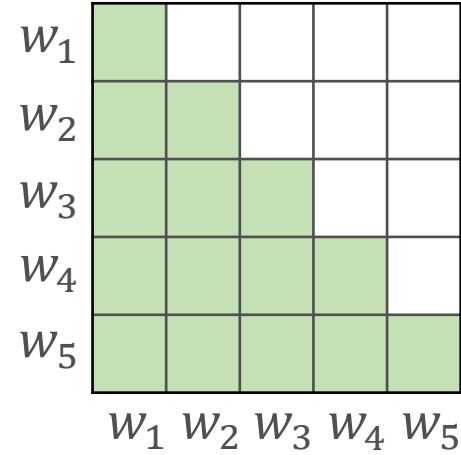
Reasons for Positional Bias: Attention Mechanism



Reasons for Positional Bias: Attention Mechanism



Causal Attention Mask



Reasons for Positional Bias: Positional Encoding

Rotary Position Embedding
(RoPE)

$$\mathbf{q}_m = f_q(\mathbf{x}_m, m)$$

$$\mathbf{k}_n = f_k(\mathbf{x}_n, n)$$

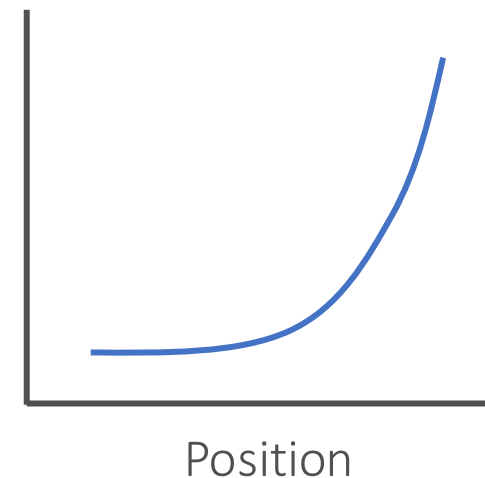
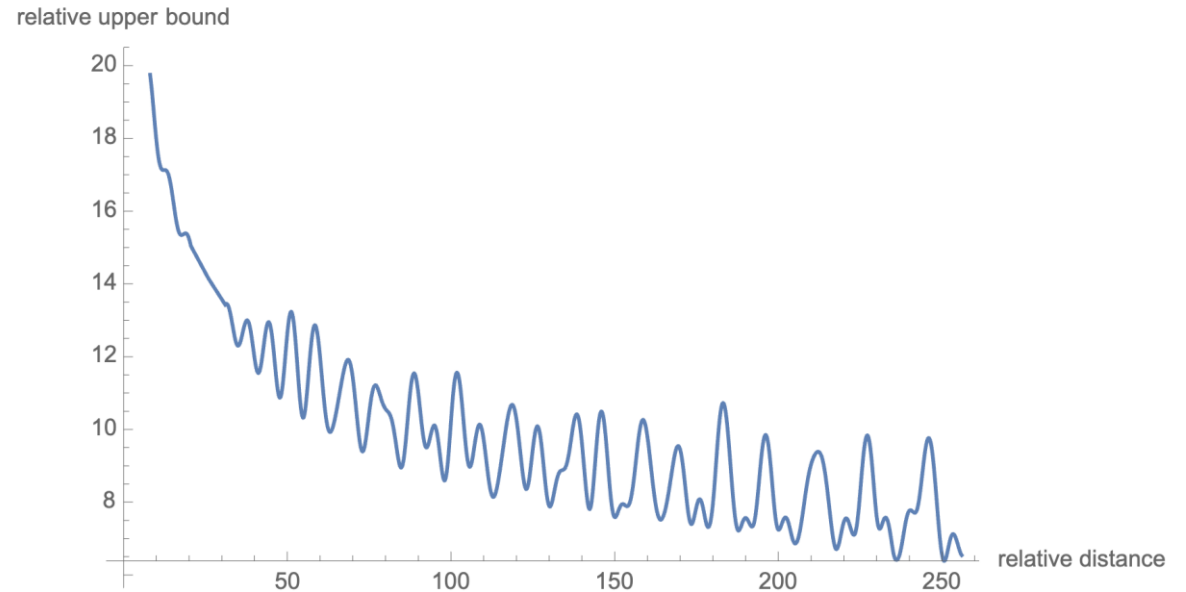
$$\mathbf{v}_n = f_v(\mathbf{x}_n, n)$$

$$f_q(\mathbf{x}_m, m) = (\mathbf{W}_q \mathbf{x}_m) e^{im\theta}$$

$$f_k(\mathbf{x}_n, n) = (\mathbf{W}_k \mathbf{x}_n) e^{in\theta}$$

$$\langle f_q(\mathbf{x}_m, m), f_k(\mathbf{x}_n, n) \rangle =$$

$$\text{Re}[(\mathbf{W}_q \mathbf{x}_m)(\mathbf{W}_k \mathbf{x}_n)^* e^{i(m-n)\theta}]$$

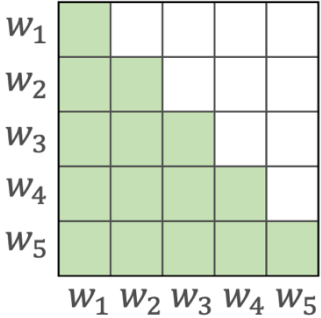


Combine All Together

Introduction
 First Main Point
 Second Main Point
 Third Main Point
Conclusion



Causal Attention Mask



Position

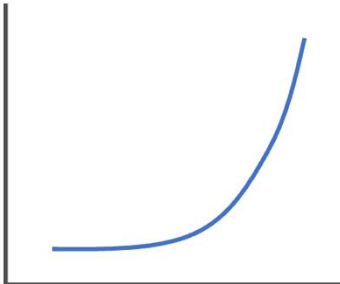
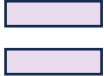


Rotary Position Embedding (RoPE)

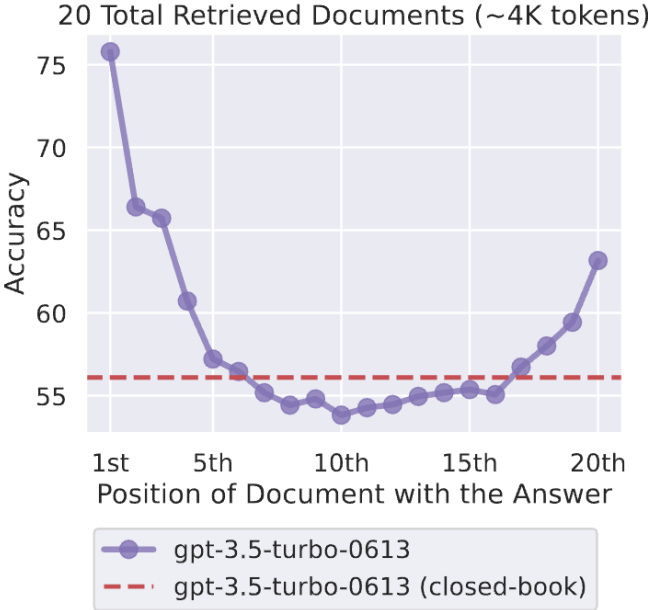
$$\mathbf{q}_m = f_q(\mathbf{x}_m, m)$$

$$\mathbf{k}_n = f_k(\mathbf{x}_n, n)$$

$$\mathbf{v}_n = f_v(\mathbf{x}_n, n)$$



Position



Lecture Plan

- Text Similarity
- Retrieval-Augmented Generation