



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

Feeding Lots Data to Language Models


CSCI 601-771 (NLP: Self-Supervised Models)

<https://self-supervised.cs.jhu.edu/sp2025/>


Long inputs/outputs in LM

- Books, scientific articles, government reports, videos, your daily experience, etc. they all are much longer than 10k tokens!!

Long Documents




Below is a paper. Memorize the paper and answer my question after the paper.
The paper begins. \n {paper} \n Now the paper ends.
Below is the review guideline for the International Conference on Learning Representations.
{Review guideline} Now the review guideline ends.
Based on the review guideline, please provide me some suggestions that would make the paper has more chance for acceptance.



Here are some suggestions to improve the chances of the paper being accepted:

- Clarify the contributions and novelty more concisely. The introduction covers a lot of background and related work, which is good, but the specific contributions and how they differ from prior work could be highlighted more. Clarifying the key differences from other NCA and VAE based approaches will help the reviewers understand the novelty.
- Provide more experimental details and results. Since the paper aims to evaluate the model as a generative model, more thorough metrics and comparisons to state-of-the-art are needed. Specifically:
 - Report log-likelihood numbers on test sets for both MNIST and CelebA using the same method (e.g. importance sampling)
 - Compare the log-likelihood to other relevant generative models for each dataset
 - Report Fr chet Inception Distance (FID) scores for the CelebA samples
 - Provide visual comparisons to other NCA and VAE based approaches on both MNIST and CelebA
- Discuss more potential applications and future directions. Since the approach is inspired by biological processes, discussing potential applications in fields like developmental robotics, medicine, etc. could be of interest. Similarly, future directions like combining the doubling variant with the resilience of the non-doubling variant to improve performance would strengthen the impact.

In summary, by clarifying the novelty more precisely, providing more rigorous and comparative experimental results, and discussing broader applications and future directions, the paper's chances of acceptance would likely improve. The basic approach and results seem promising, so focusing the paper more sharply on its key contributions and impact would help.



Inference-time Reasoning

of-preview

Assume laws of physics on Earth. A small strawberry is put into a normal cup and the cup is placed upside down on a table. Someone then takes the cup and puts it inside the microwave. Where is the strawberry now? Explain your reasoning step by step



Tracking the strawberry's path

To determine the current location of the strawberry, let's break down the events step by step:

1. **Initial Placement:** A small strawberry is placed into a normal cup. The strawberry is inside the cup, resting at the bottom.
2. **Cup Inversion:** The cup is then placed upside down on a table. This means the open end of the cup (where the strawberry was initially placed) is now facing down, in contact with the table surface. The strawberry is now sitting on the table, covered by the inverted cup.
3. **Cup Removal:** Someone picks up the cup to put it inside the microwave. When lifting the inverted cup off the table:

Nominal Context Length

Model	Creators	Year	Model Size	Context Size
GPT-2	OpenAI	2019	1.5 billion	1024
GPT-3	OpenAI	2020	175 billion	2048
PaLM	Google	2022	540 billion	2048
LLaMA	Meta	2023	65 billion	2048
LLaMA-2	Meta	2023	70 billion	4096
Claude-2	Anthropic	2023	? (130 billion)	100k
Claude-2.1	Anthropic	2023	? (130 billion)	200k
GPT-4	OpenAI	2023	? (1.76 trillion)	8192
Mistral	Mistral AI	2023	7 billion	8192 (32k)
Mixtral	Mistral AI	2023	47 billion	8192 (128k)
Gemini (Ultra)	Google	2023	? (1.5 trillion)	32k
LWModel	academia!	2023	7 billion	1 million
Gemini-1.5	Google	2024	? (1.5 trillion)	1 million
GPT-4o	OpenAI	2024	?	128k
LLaMA-3	Meta	2024	405 billion	8192
LLaMA-3.1	Meta	2024	405 billion	128k
Claude-3.5	Anthropic	2024	?	200k

Revisiting Encoding Positional Information

Recap: Self-Attention

- Given input \mathbf{x} :

$$Q = \mathbf{W}^q \mathbf{x}$$

$$K = \mathbf{W}^k \mathbf{x}$$

$$V = \mathbf{W}^v \mathbf{x}$$

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{QK^T}{\alpha}\right) V$$

- Remember that, without positional encoding the input is just bag of words for self-attention.

Positional Embeddings: The Flavors

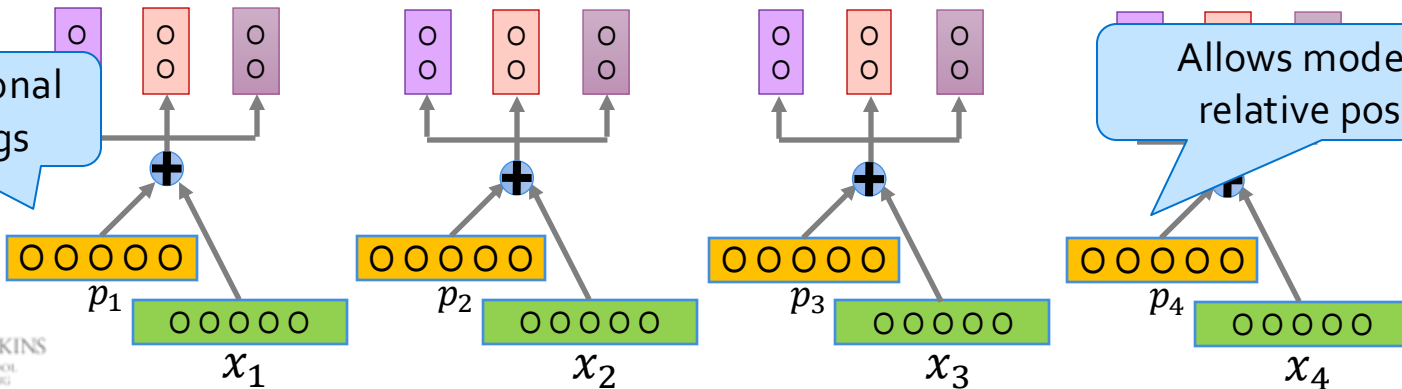
- **Sine embeddings:** add sines and cosines that enable localization

Notable models:
Original Transformer

$$Embed(x, i) = v_x + PE_{pos}$$

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

p_i are positional embeddings



Positional Embeddings: The Flavors

- **Sine embeddings:** add sines and cosines that enable localization

$$Embed(x, i) = v_x + PE_{pos}$$

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

Notable models:
Original Transformer

- **Absolute embeddings:** add a position vector to the embedding

$$Embed(x, i) = v_x + u_i$$

Notable models:
GPT1/2/3 - OPT

- **Limitations:**

- We can have fixed encoding for each index training position (e.g., 1, 2, 3, ... 1000).
 - What happens if we get a sequence with 5000 words at test time?
- We want something that can generalize to arbitrary sequence lengths.

Positional Embeddings: The Flavors

- **Sine embeddings:** add sines and cosines that enable localization

$$Embed(x, i) = v_x + PE_{pos}$$

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

Notable models:
Original Transformer

- **Absolute embeddings:** add a position vector to the embedding

$$Embed(x, i) = v_x + u_i$$

Notable models:
GPT1/2/3 - OPT

- **Relative embeddings:** add a vector to the attention computation

$$QK_{ij} = \mathbf{x}_i^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_j + \mathbf{P}_{ij}$$

Notable models:
T5, Gopher, Chinchilla

- Intuition: encoding the relative positions, for example based on the distance of the tokens in a local window to the current token.

Relative Positional Encoding: An Example

- *"I took the bus and then a cab after work"*
 - position of "bus": 4;
 - position of "cab": 8
- *"During the last two weeks, when my car was in the shop for repair, I took the bus and then a cab after work."*
 - position of "bus": 20;
 - position of "cab": 24
- You want a way of encoding the sentence so that the attention patterns depend on **the relative distances**.

Relative Positional Encoding

- There have been various choices:

- T5 models simplify this into learnable relative embeddings P_{ij} such that:

$$QK_{ij} = \mathbf{x}_i^T W_q^T W_k \mathbf{x}_j + P_{ij}$$

- DeBERTa learns relative positional embeddings \tilde{p}_{i-j} such that:

$$QK_{ij} = \mathbf{x}_i^T W_q^T W_k \mathbf{x}_j + \mathbf{x}_i^T W_q^T W_k \tilde{p}_{i-j} + \tilde{p}_{i-j}^T W_q^T W_k \mathbf{x}_j$$

- Transformer-XL learns relative positional embeddings \tilde{p}_{i-j} and trainable vectors \mathbf{u}, \mathbf{v} s.t.:

$$QK_{ij} = \mathbf{x}_i^T W_q^T W_k \mathbf{x}_j + \mathbf{x}_i^T W_q^T W_k \tilde{p}_{i-j} + \mathbf{u}^T W_q^T W_k \mathbf{x}_j + \mathbf{v}^T W_q^T W_k \tilde{p}_{i-j}$$

- ALiBi learns a scalar m such that:

$$QK_{ij} = \mathbf{x}_i^T W_q^T W_k \mathbf{x}_j - m |i - j|$$

Recap

- **Sine embeddings:** add sines and cosines that enable localization
- **Absolute embeddings:** add a position vector to the embedding
- **Relative embeddings:** add a vector to the attention computation
- **RoPE embeddings:** (forthcoming)

Notable models:
Original Transformer

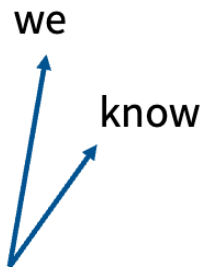
Notable models:
GPT1/2/3 - OPT

Notable models:
T5, Gopher,
Chinchilla, Deberta
Tranformer-XL,

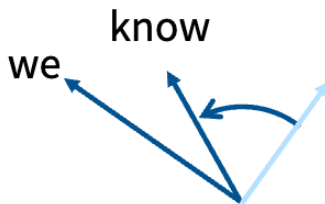
Notable models:
GPTJ, PaLM, LLaMA,
Gemma3, Qwen3,
GPL, Apple models.

Rotary Positional Encoding (RoPE)

- We want our embeddings to be invariant to absolute position.
- We know that inner products are invariant to arbitrary rotation.

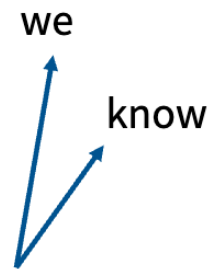


Position independent
embedding



Embedding
“of course we know”

Rotate by ‘2 positions’



Embedding
“we know that”

Rotate by ‘0 positions’

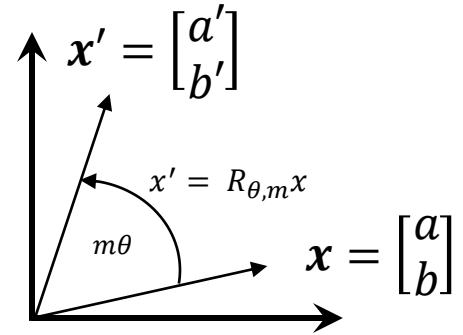
[Slide credit: Tatsu Hashimoto]

Thinking About Rotation Matrix

- In 2D, a rotation matrix can be defined in the following form:

$$R_{\theta,m} = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix}$$

- The rotation increases with increasing θ and m .



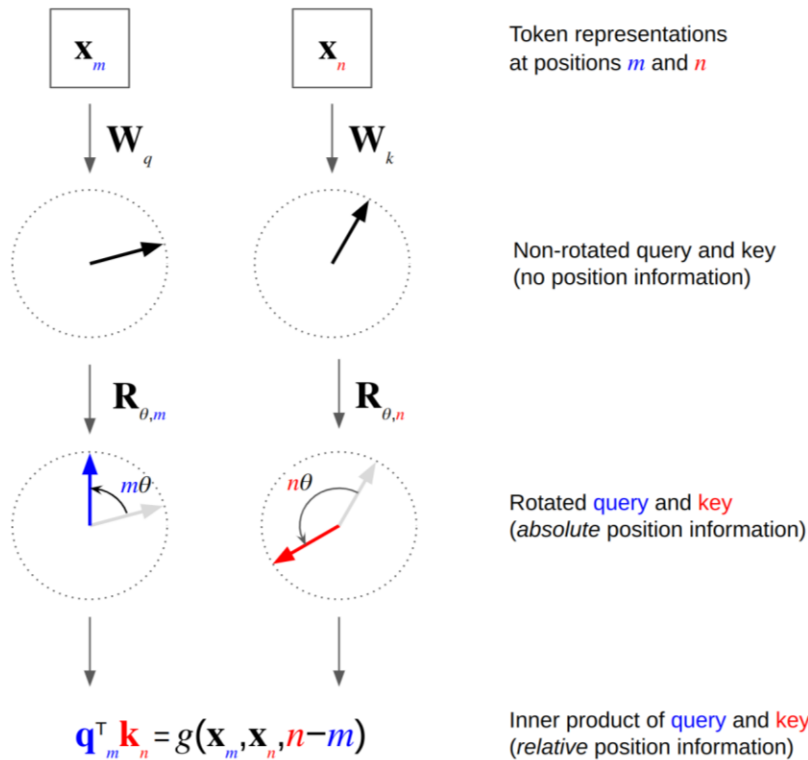
Rotary Positional Encoding (RoPE): 2D

- If you have two input representations x_n and x_m at positions n and m :

$$\begin{aligned} qk_{mn} &= (R_{\theta,m} W_q x_m)^T (R_{\theta,n} W_k x_n) \\ &= x_m^T W_q^T R_{\theta,m}^T R_{\theta,n} W_k x_n \\ &= x_m^T W_q^T R_{\theta,n-m} W_k x_n \end{aligned}$$

- $R_{\theta,m}$: rotation matrix, rotates a vector it gets multiplied to proportional to θ and the position indices m, n .

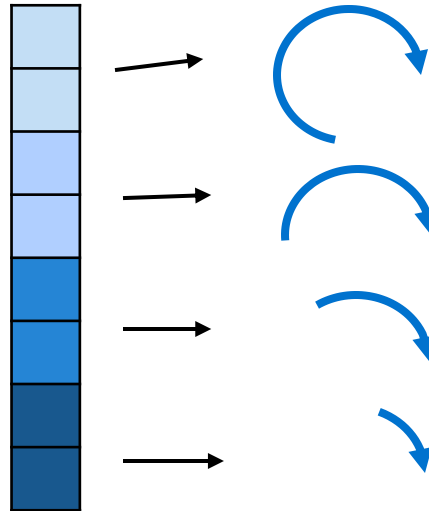
- Intuition: **nearby** words have **smaller relative rotation**.



[Figure source](#)

Thinking d Dimensions

- Create **pairs of** coordinates that would rotate together
- Each pair can have a different frequency of rotations



Thinking d Dimensions

- In practice, we are rotating d dimensional embedding matrices.
- Idea: rotate different dimensions with different angles:
 - $\Theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_{d/2}]$

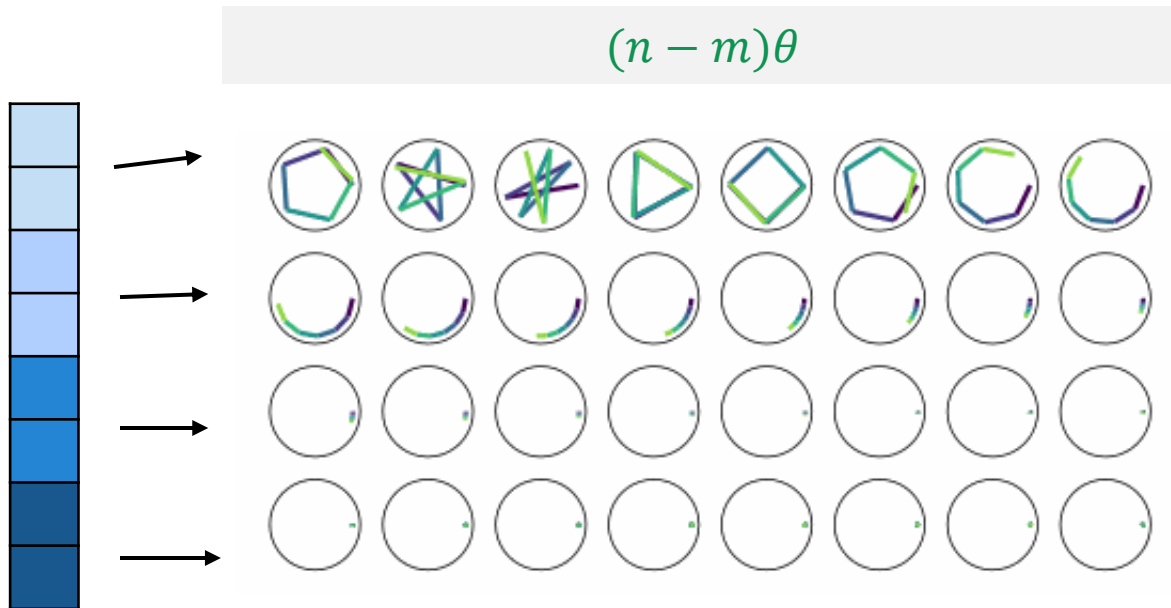
d = the dimension of the query and key space
 $m\theta$ = total rotation angle

$$R_{\Theta, m}^d = \begin{bmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \dots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \dots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \dots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{bmatrix}$$

- By using different frequencies for different dimensions of the positional encoding vector, a unique "fingerprint" can be created for each positional information.

RoPE, visualized.

- Each pair of coordinates are rotated by an angle proportional to the token position $n\theta$ or relative dist of words $(n - m)\theta$
- Top row: high-freq dimensions.
- Bottom row: low-freq dimensions.



Visualization by Sasha Rush:

<https://colab.research.google.com/drive/1JhVaH1Tr6BMLKX9mgljN83qFh9DOOyvC?usp=sharing#scrollTo=xeBojFzJRfu1>

RoPE via Faster (Dense) Operations

- While we can implement this matrix multiplication as is, we can make use of $R_{\Theta,m}^d$'s sparsity for faster multiplication via dense operations :

$$\mathbf{x}' = \mathbf{R}_{\Theta,m}^d \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_{d-1} \\ x_d \end{bmatrix} \otimes \begin{bmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{bmatrix} + \begin{bmatrix} -x_2 \\ x_1 \\ -x_4 \\ x_3 \\ \vdots \\ -x_d \\ x_{d-1} \end{bmatrix} \otimes \begin{bmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \\ \sin m\theta_{d/2} \end{bmatrix}$$

Integrating RoPE to Self-Attention

- **Drop the additive** positional encoding and make it **multiplicative to QK** matrices.
- Given input \mathbf{x} , the standard SA is computed as:

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$
$$\mathbf{Q} = \mathbf{W}^q\mathbf{x}, \quad \mathbf{K} = \mathbf{W}^k\mathbf{x}, \quad \mathbf{V} = \mathbf{W}^v\mathbf{x}$$

- Given input \mathbf{x} , SA with RoPE is computed as:

$$\text{Attention}(\mathbf{x}) = \text{softmax}\left(\frac{\mathbf{Q}_{\text{rope}}\mathbf{K}_{\text{rope}}^T}{\sqrt{d}}\right)\mathbf{V}$$

- Note since each the rotation depends on distance each word pair, it's not easy to write these in matrix form.

Setting $\Theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_{d/2}]$

- One common way is to set them using a “base” frequency:

$$\Theta = \left\{ \theta_i = B^{-2(i-1)/d}, i \in [1, 2, \dots, d/2] \right\}$$

- Example:
 - $\theta_1 = 1$, for any choice of B or d
 - $\theta_{d/2} = \frac{1}{B^{1-2/d}}$. If $d = 512$ and $B = 10K$, then $\theta_{256} = 1/10,000$
- The exponent $-2(i-1)/d$ scales the input so that:
 - **lower** dimensions correspond to **higher**-frequency sinusoids (rapid changes), and
 - **higher** dimensions correspond to **lower**-frequency sinusoids (slow changes)
- This **spread of frequencies** allows the model to represent both fine-grained (local) and coarse-grained (global) positional relationships.

Q: How should we set the “base” freq?

$$\Theta = \left\{ \theta_i = B^{-2(i-1)/d}, i \in [1, 2, \dots, d/2] \right\}$$

- You want the largest wavelength (lowest freq) be close to your longest context dist.
- Lowest frequency: $\theta_{d/2} = \frac{1}{B^{1-2/d}} \Rightarrow$ The largest wavelength: $T_{\max} = \frac{2\pi}{\theta_{d/2}} \approx 2\pi B$
- Examples:
 - If $B = 10,000 \Rightarrow T_{\max} \approx 62K$ tokens
 - If $B = 100,000 \Rightarrow T_{\max} \approx 628K$ tokens
- In both original RoPE paper (and the original Transformer paper), $B = 10,000$.

Examples of Base Frequencies

- Llama3 positional encoding: RoPE ($\theta = 500,000$)

We increase the RoPE base frequency hyperparameter to 500,000. This enables us to better support longer contexts; [Xiong et al. \(2023\)](#) showed this value to be effective for context lengths up to 32,768.

- Gemma3 positional encoding: A mix of RoPE ($\theta = 1,000,000$) + RoPE ($\theta = 10,000$)

Long context. Gemma 3 models support context length of 128K tokens, with the exception of the 1B model that has 32K. We increase RoPE base frequency from 10k to 1M on global self-attention layers, and keep the frequency of the local layers at 10k. We follow a process similar to the

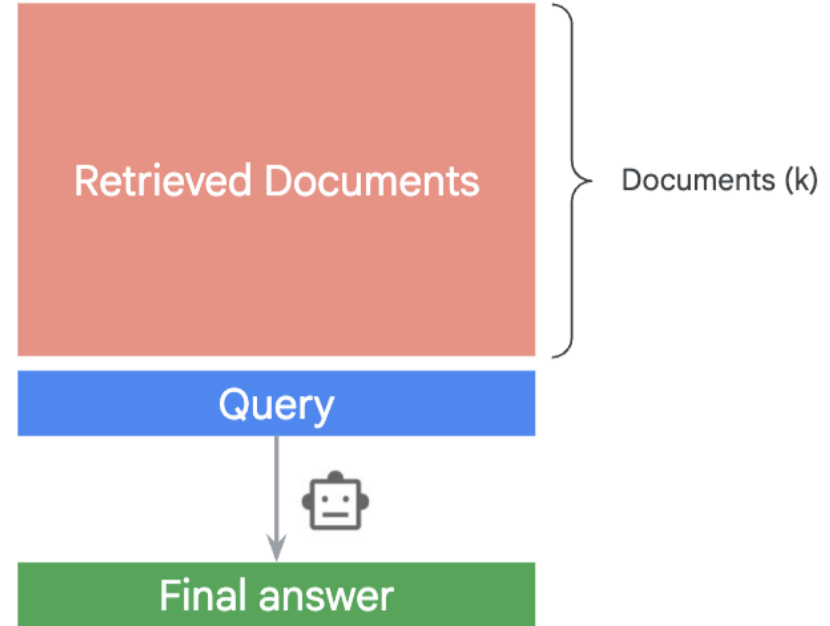
Summary

- Encoding positional information in language models is a non-trivial problem.
 - We discussed various proposal: learned, absolute, relative encoding, NoPos, etc.
- This is an important literature related to the length generalization of Transformers.
- This is an active research area and likely to change in the coming years.

Long Context Generalization

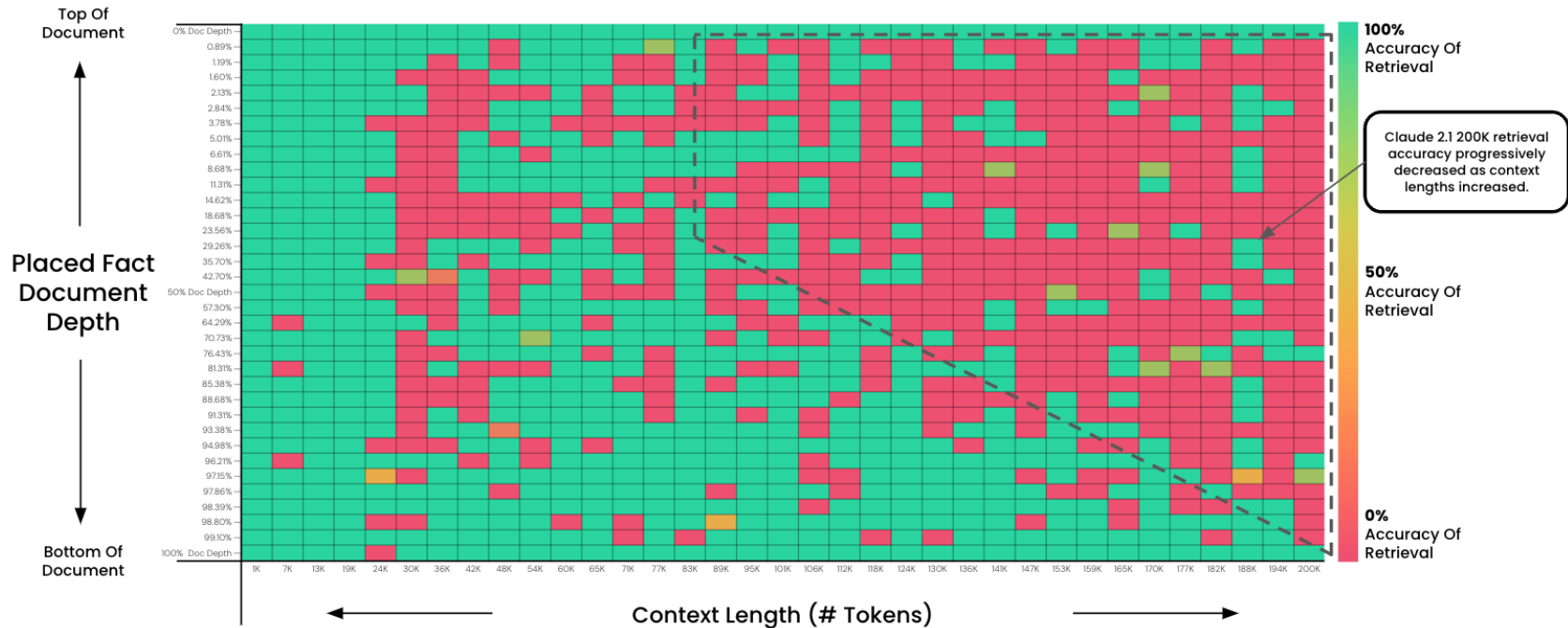
Needle-in-Haystack Problems

- You're given k documents in the input.
 - (e.g., the retrieval results)
- One of these docs may contain the answer to a given query which we want to answer.



Pressure Testing Claude-2.1 200K via "Needle In A HayStack"

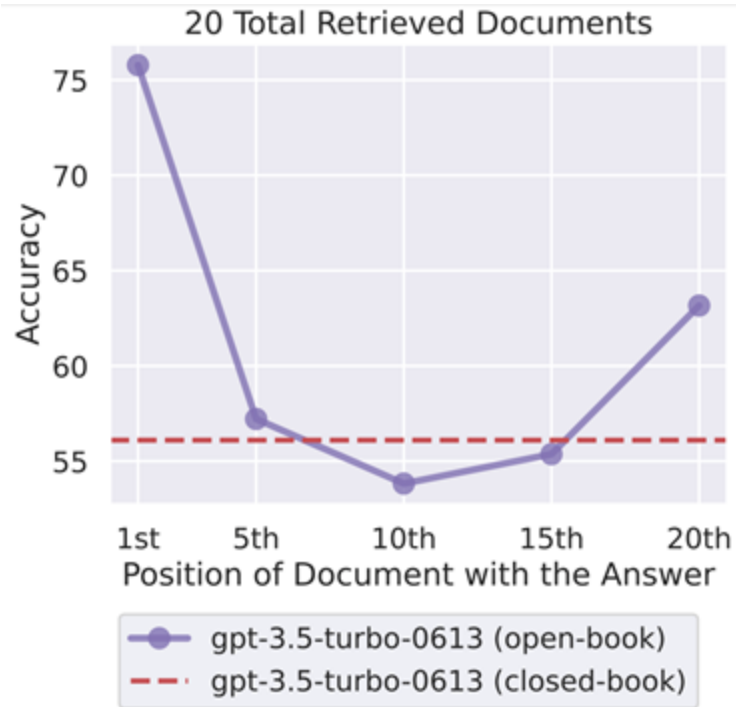
Asking Claude 2.1 To Do Fact Retrieval Across Context Lengths & Document Depth



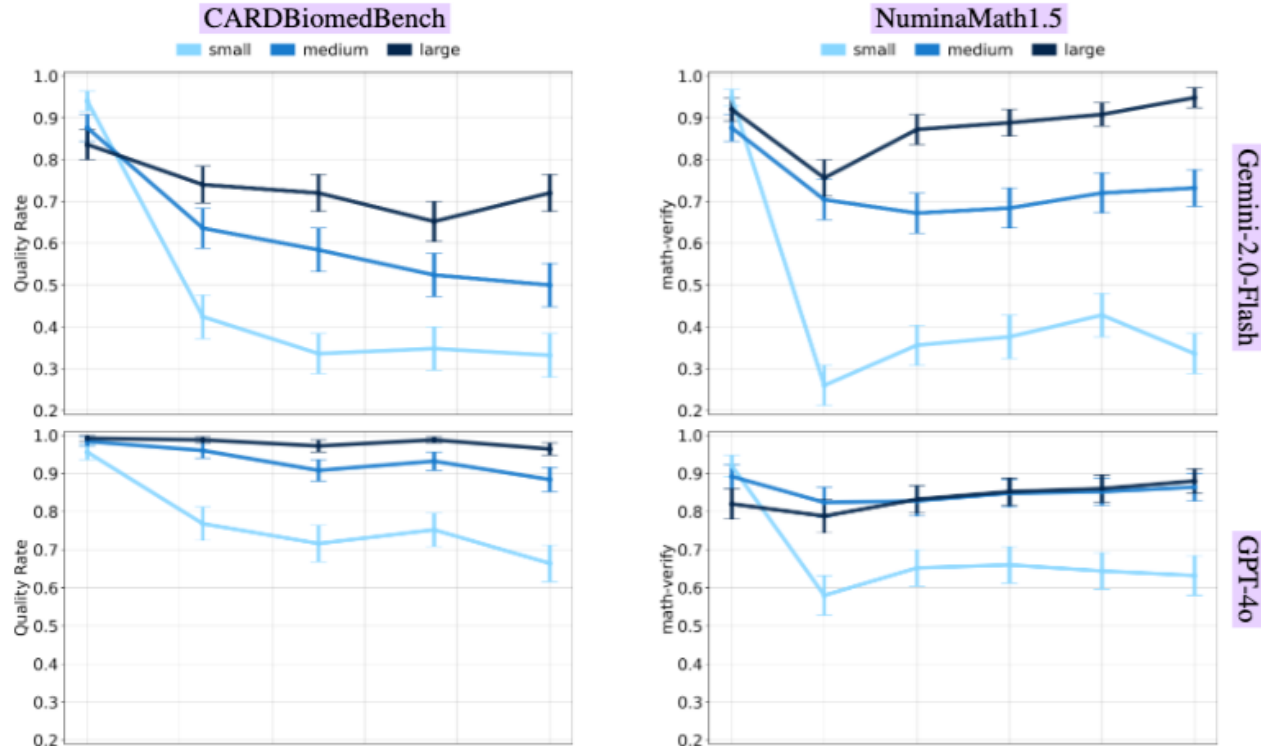
Goal: Test Claude 2.1 Ability To Retrieve Information From Large Context Windows

A fact was placed within a document. Claude 2.1 (200K) was then asked to retrieve it. The output was evaluated (with GPT-4) for accuracy. This test was run at 35 different document depths (top > bottom) and 35 different context lengths (1K > 200K tokens). Document Depths followed a sigmoid distribution

Recency vs Primacy bias



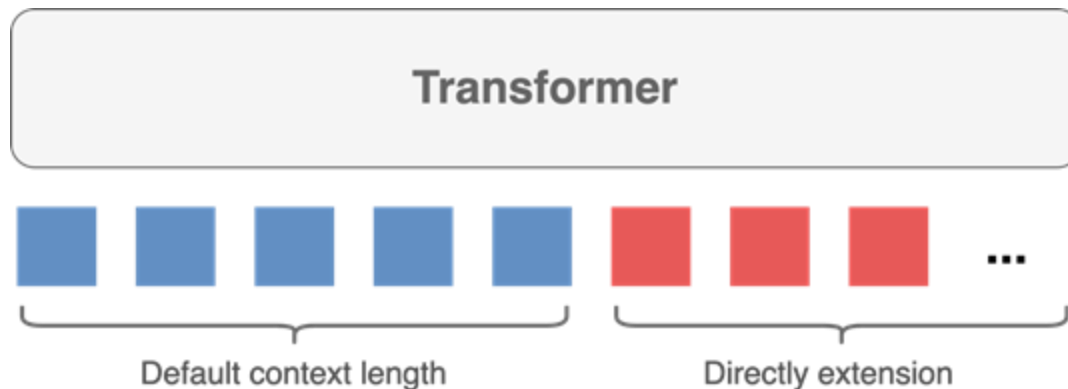
Recent vs Primacy Bias: Size of “Gold” doc



Long Context Extensions

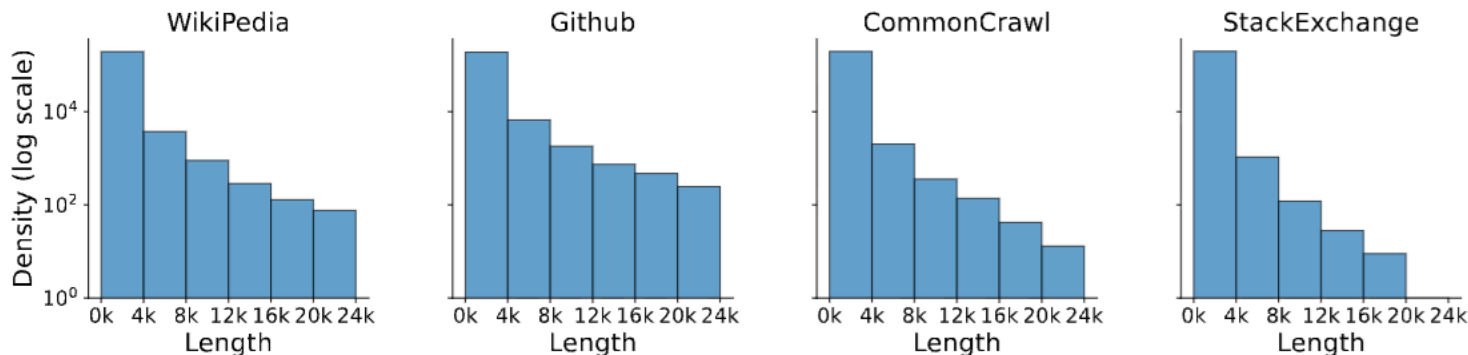
Context Extension: Challenges

- You're given a transformer trained on a default length.
- Now we want to extend its length.



Context Extension: Challenges

- **Challenge 1:** Adapting the positional encodings
- **Challenge 2:** Efficiency challenge since the FLOP/memory footprint grows quadratically.
- **Challenge 3:** Long data is very limited.



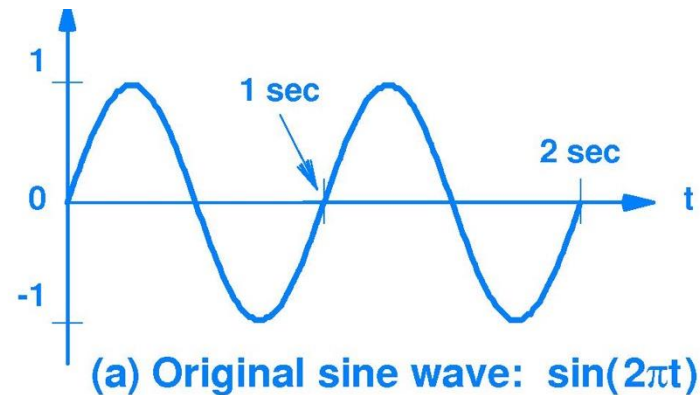
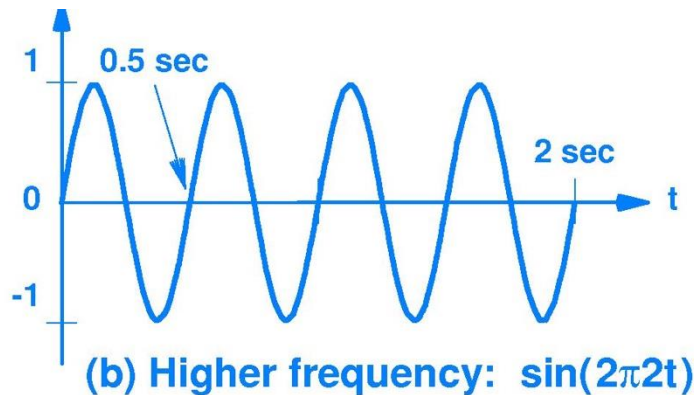
Context Extension during Pre-training

- Two ingredients:
 - Extensible positional encoding
 - Carefully selected long training data
- Recipe:
 - Pre-train a short context model(block size = 4k) on 1 trillion tokens of text
 - Adjust the hyperparameters of the positional embeddings
 - Continue pre-training but increase the block size to support long-contexts (block size = 80k). This may involve training on a smaller amount of data.

Approach: Changing ROPE Frequencies

- Assuming that we use RoPE encoding, **reduce** the base frequency by α so that more tokens fit within each period.

$$\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i, \boldsymbol{\theta}) = \mathbf{R}(\boldsymbol{\theta}, i) \mathbf{W}^{\top} \mathbf{x}_i \rightarrow \mathbf{f}(x_i, \alpha \odot \boldsymbol{\theta})$$



Uniform Scaling

- Linear Position Interpolation (PI):

- Scale all frequencies by the same scale:

$$\alpha_j^{\text{PI}} = \frac{c}{c'} = \frac{1}{t}$$

- This lowers every frequency *uniformly* so more tokens fit in each sinusoidal period.

NTK Scaling

$$\alpha_j = \kappa^{-\frac{2j}{d_k}}, \quad j = 1, \dots, \frac{d_k}{2}$$

- High frequencies scaled down less (\sim preserved)
 - Low frequencies are shrunk more (matching PI extension)
-
- Why is it called "NTK"?
 - Neural Tangent Theory is this idea that, high-frequency information are hard to learn.
 - So, it's better to not mess with them! (i.e., scale high-frequencies minimally)

YaRN Scaling

- Basically, piecewise scaling across dimensions.

$$\alpha_j^{\text{YaRN}} = ((1 - \gamma_j) \frac{1}{t} + \gamma_j) / \sqrt{T}$$

$$\gamma_j = \begin{cases} 0, & \text{if } \theta_j < p, \\ 1, & \text{if } \theta_j > q, \\ \frac{\theta_j - p}{q - p}, & \text{otherwise.} \end{cases}$$

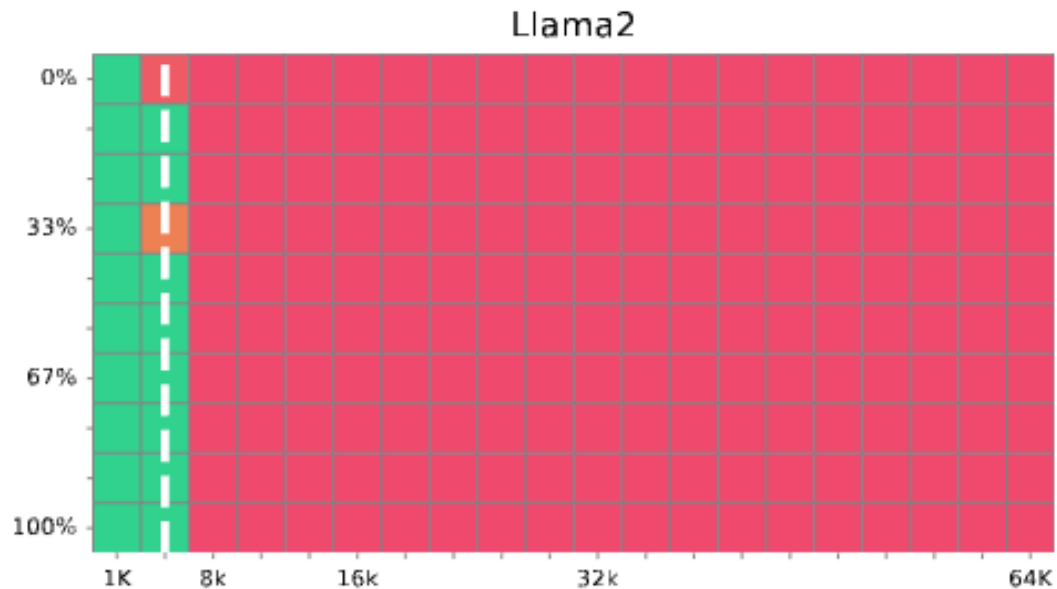
- Context length ratio: $t = C'/C$
- Temperature parameter: T to reduce all freq uniformly, if needed.
- Interpolation gate for dimension j : γ_j
- Hyperparameters: p, q, T which can be tuned as needed.
- Modes:
 - Low-frequency dimensions ($\theta_j < p$): fully PI-scaled ($1/t$).
 - High-frequency dimensions ($\theta_j > q$): no scaling (kept as original RoPE) if $T=1$.
 - Middle band ($q < \theta_j < p$): smoothly interpolated.

Empirical Comparison of RoPE extensions

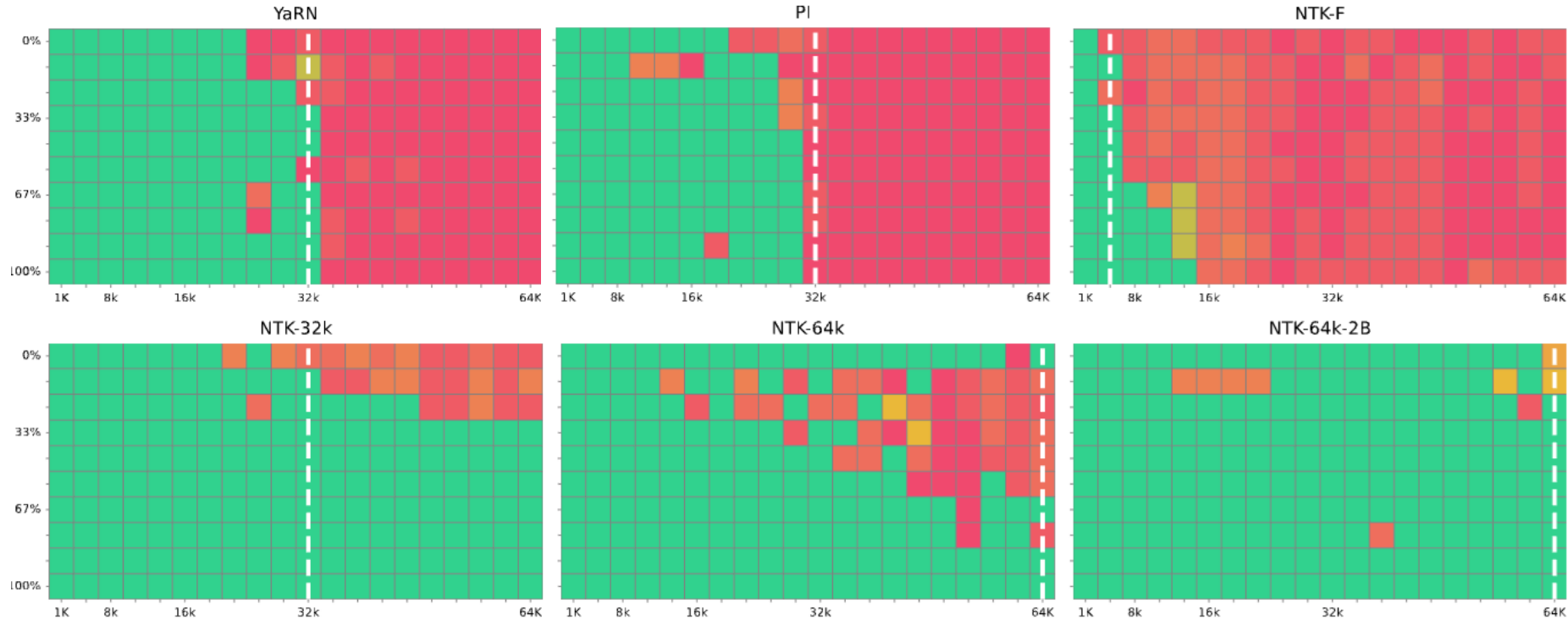
- Base model: llama2-7B which has a context of window of $\sim 4k$ tokens.
- Goal: extending the context window from 4k to 32k

Table 1: Overview of results across different extension types.

Attention Mechanisms		Model	PPL	Needle	Mshots	LongB	RULER
Exact Attention	Frozen	NTK-F	14.52	18.8	64.5	25.54	0.72
	Fine-Tuned	PI	5.85	42.1	75.5	33.48	57.66
		YaRN	5.85	46.7	75.0	33.45	36.95
		CLEX	5.82	71.1	74.0	33.48	52.17
		NTK-32K	5.79	83.7	71.0	35.32	59.42
		NTK-64K	5.93	69.1	73.0	34.30	60.03

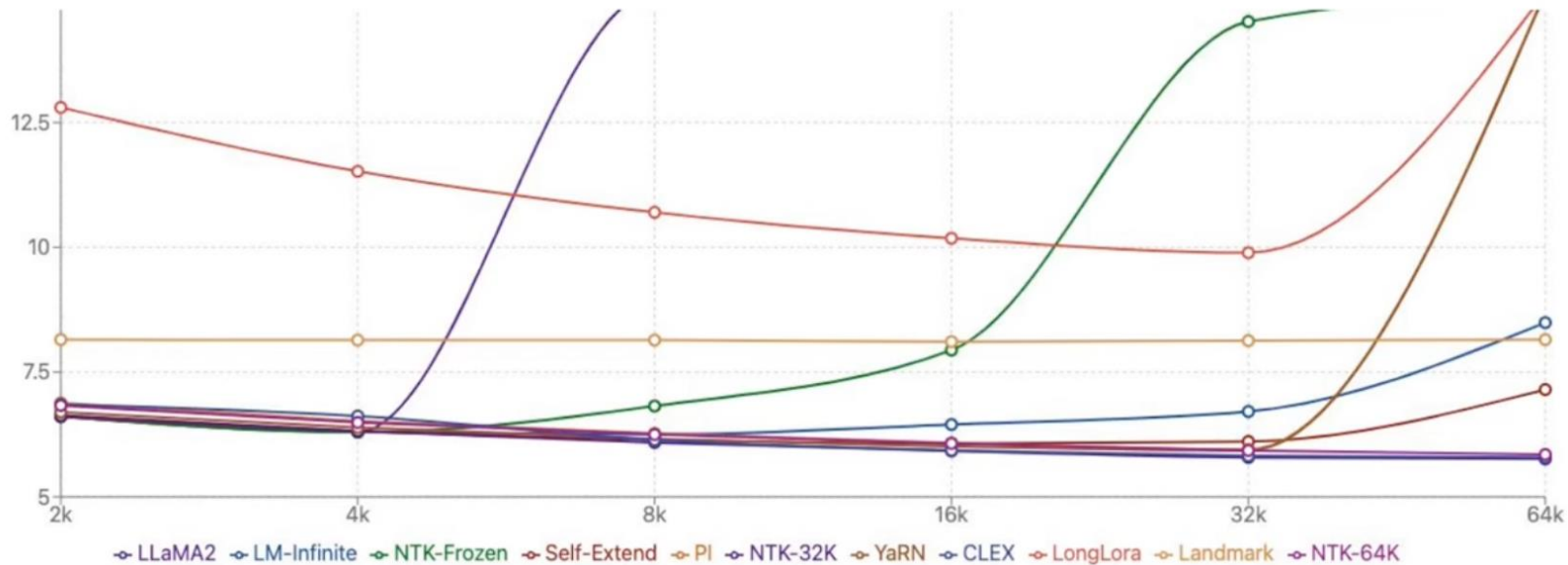


- X-axis shows length of the input.
- Y-axis shows where the “needle” was included in the context (first row: needle in the beginning)
- The white dashed line denotes the longest length examples seen at training or finetuning



- X-axis shows length of the input.
- Y-axis shows where the “needle” was included in the context (first row: needle in the beginning)
- The white dashed line denotes the longest length examples seen at training or finetuning

[A Controlled Study on Long Context Extension and Generalization in LLMs, 2025](#)



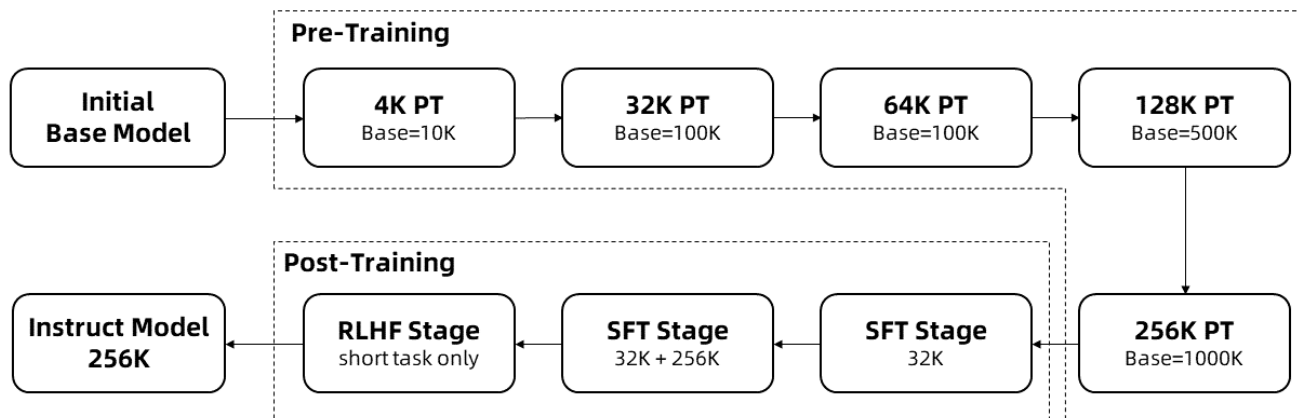
- Y-axis shows perplexity.
- The base model is trained up to 4k tokens. .
- The extensions are trained on data up to 32k.
- Note some generalize to 64k.

Model	Positional enc	Pre-training extension	Post-training extension
Qwen 3	RoPE	Increased base freq from 10K to 1M; used YaRN. Also uses Dual Chunk Attention.	Training on a dataset of long Chain-of-Thoughts.
Apple	RoPE (local attention) and no encoding for global attention	Trained on seq up to 65K tokens sampled from “natural” data e.g., licensed books, code repos, synthetic long-form data.	??
GLM	RoPE	Trained on 4K data (base=10K). Then did more training up to 128K (base=1M). Up-sample long documents from the pre-training corpus.	SFT with long reasoning traces (up to 128K)
Gemma3	RoPE	Pre-train on 32K seq (base=10K) and scale to 128K toward the end of training (base=1M).	??
Lamma3	RoPE	Pre-train on shorter sequences (base=500K) and then extend to 128K tokens.	SFT on some (0.1%) long data

[Qwen3 Technical Report](#), [Apple Intelligence Foundation Language Models](#),
[GLM-4.5: Agentic, Reasoning, and Coding \(ARC\) Foundation Models](#), [Gemma3 tech report](#)

Examples in the wild: Qwen3

- Qwen2.5-7B-Instruct-1M
 - **Changing RoPE base frequency** (to stretch positional embeddings further).
 - **Training recipe** to gradually expose the model to longer contexts.



Examples in the wild: Gemma3

- *"Instead of training with 128K sequences from scratch, we pre-train our models with 32K sequences and then scale the 4B, 12B, and 27B models up to 128K tokens at the end of pre-training while rescaling RoPE"*
- *"Our models generalize to 128K, but rapidly degrade as we continue to scale."*

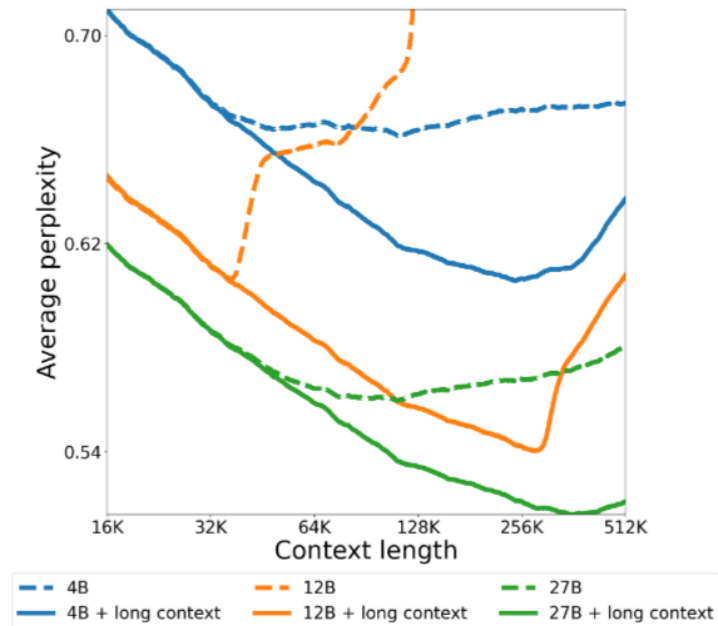


Figure 7 | **Long context** performance of pre-trained models before and after RoPE rescaling.

Examples in the wild: GLM-4.5

- We set the maximum sequence length to 4,096 during pre-training, and extended it to 32,768 and 131,072 during the mid-training stage ... When extending the sequence length to 32K, we also adjusted RoPE's base frequency from 10,000 to 1,000,000 for better long-context modeling ability."*

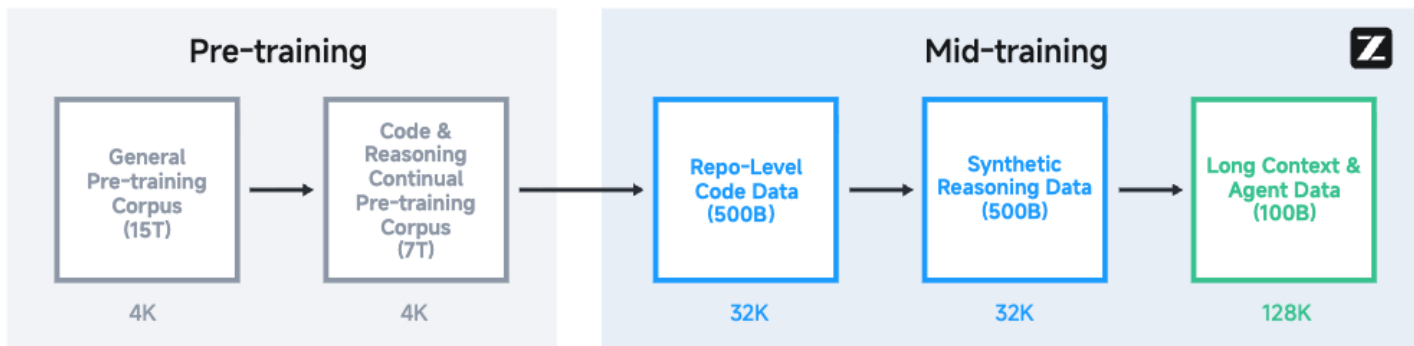


Figure 3: Pre-training and mid-training stages for GLM-4.5. We adapt a multi-stage training recipe and extend the sequence length from 4K to 128K.

Examples in the wild: Lamma3

3.4.2 Long Context Pre-Training

In the final stages of pre-training, we train on long sequences to support context windows of up to 128K tokens. We do not train on long sequences earlier because the compute in self-attention layers grows quadratically in the sequence length. We increase the supported context length in increments, pre-training until the model has successfully adapted to the increased context length. We assess successful adaptation by measuring whether **(1)** model performance on short-context evaluations has recovered completely and **(2)** the model perfectly solves “needle in a haystack” tasks up to that length. In Llama 3 405B pre-training, we increased context length gradually in six stages, starting from the original 8K context window and ending in the final 128K context window. This long-context pre-training stage was performed using approximately 800B training tokens.

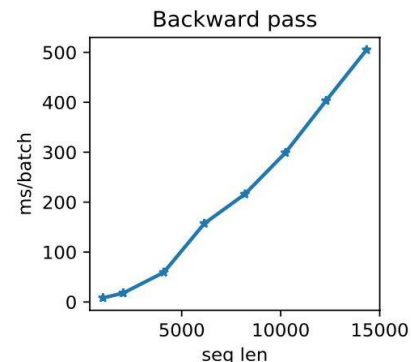
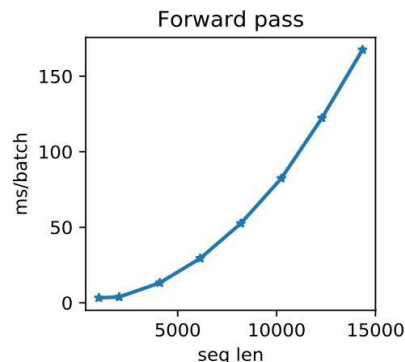
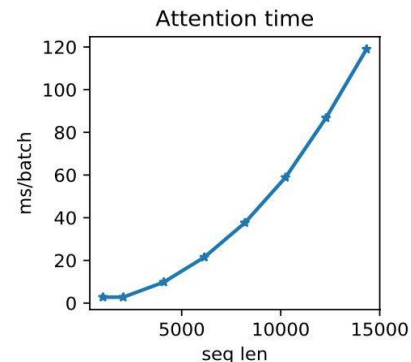
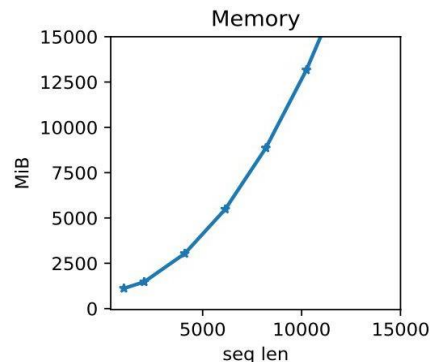
Summary

- RoPE the dominant positional encoding.
- People usually train short and extend longer.

Long Context: Efficiency Considerations

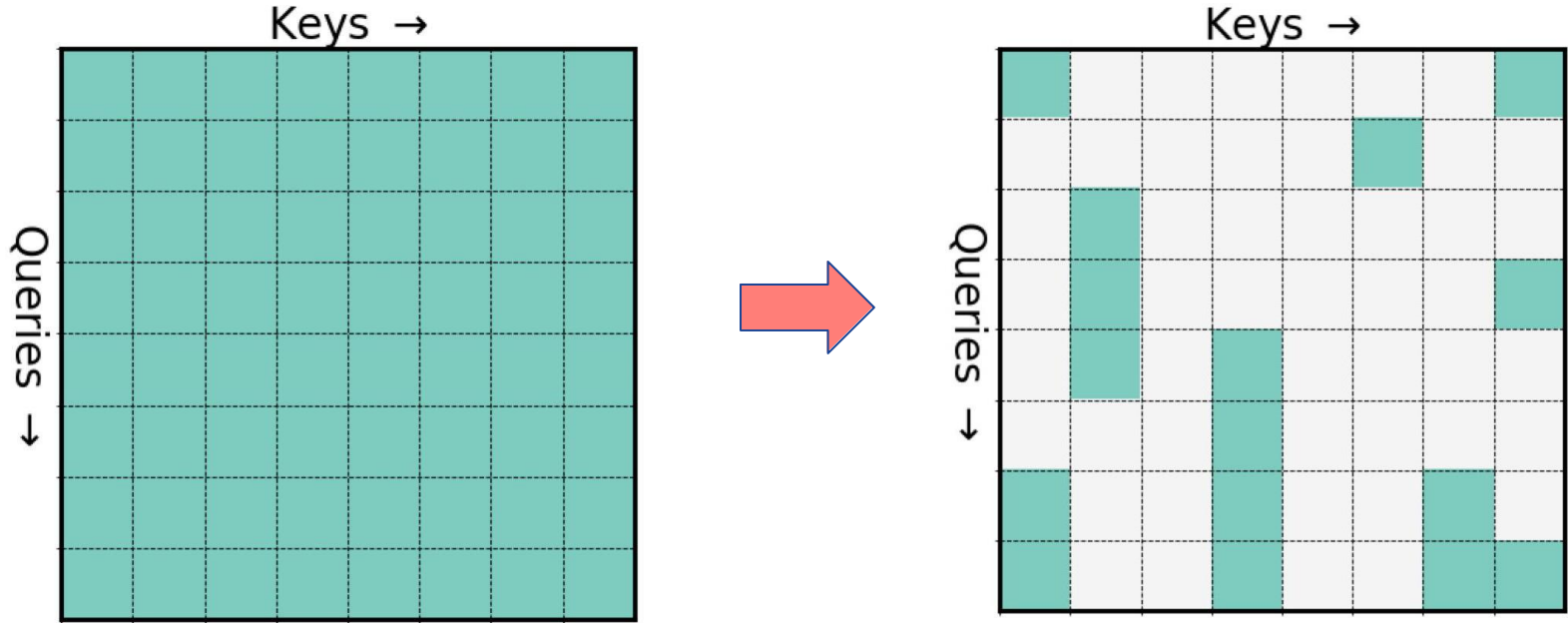
Transformer LMs and Long Inputs

- **Length generalization:** Do Transformers work accurately on long inputs?
- **Efficiency considerations:** How efficient are LMs on long inputs?



Sparse Attention Patterns

- The idea is to make the attention operation sparse

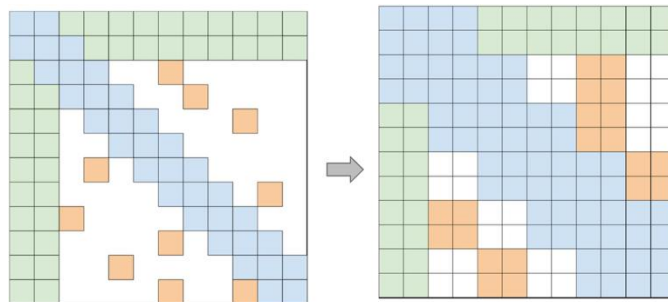


Sparse Attention Patterns: Challenge

- Ok sparsity is great, but how to efficiently implement this?
- **Challenge:** Arbitrary sparse matrix multiplication is not supported in DL libraries
- **A solution:** Perform computations in blocks

Pre-specified Sparsity Patterns: Computations

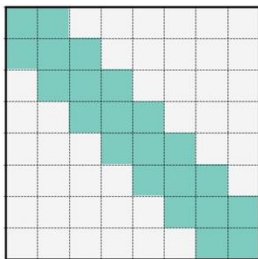
- Efficient blockified implementation
- There are libraries for implementing blockified sparse matrix multiplication.
 - Can be hardware specific
 - Block Sparse ([Gray et al., 2017](#))
 - TVM toolkit ([Chen et al., 2018](#))
 - cuSPARSE



Pre-specified Sparsity Patterns

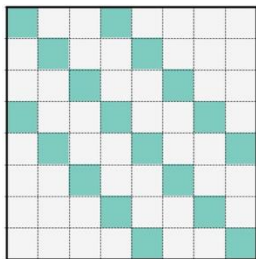
- A variety of patterns has been explored in the past work
 - Longformer ([Beltagy et al., 2020](#)), Sparse Transformer ([Child et al., 2019](#)), ...

Sliding window



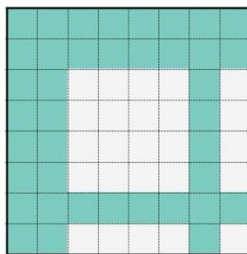
Sparse Transformer
Longformer

Dilated



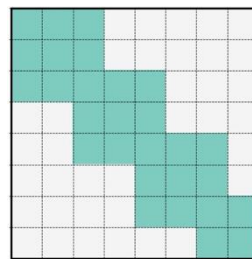
Longformer

Global



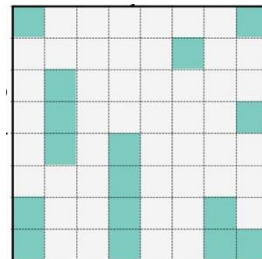
Big Bird

Blocked



Big Bird
Sinkhorn

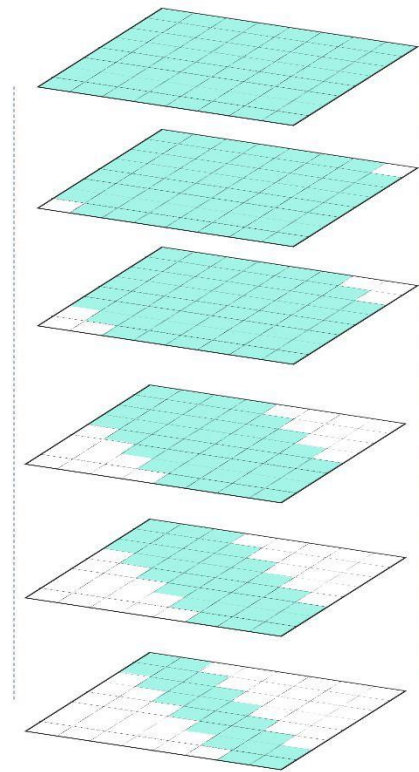
Random



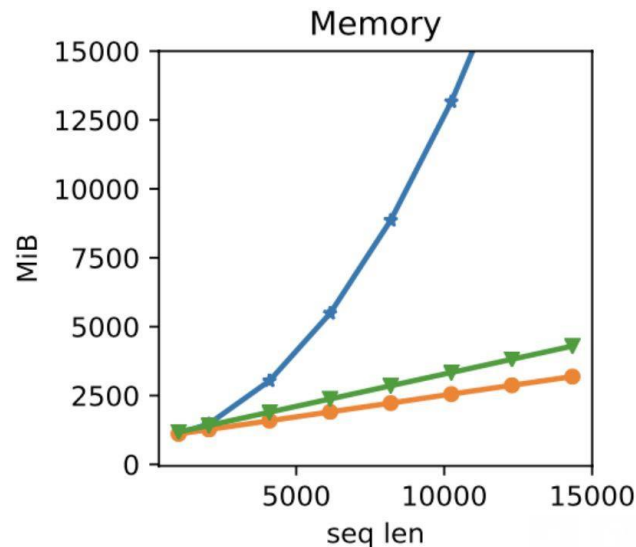
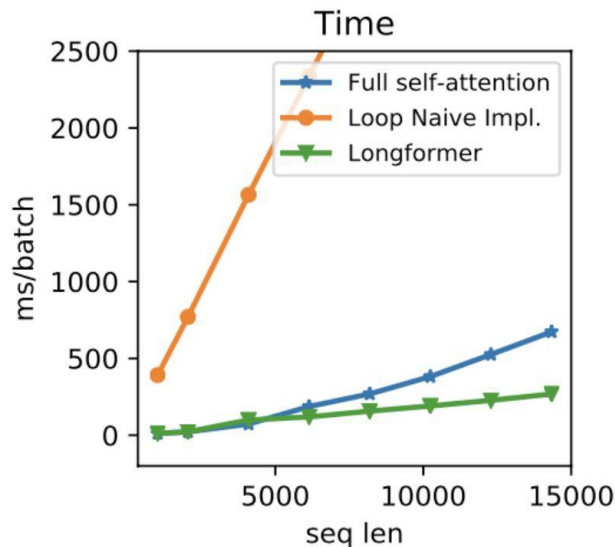
Big Bird

Pre-specified Sparsity Patterns

- Different layers and attention heads can follow different patterns
- A common setup is to have earlier layers with sparser attention pattern.
 - Longformer ([Beltagy et al, 2020](#))



Pre-specified Sparsity Patterns: Computations



A Notable Adoption: GPT-3

- Sparse patterns also used in GPT-3 ([Brown et al., 2020](#))

2.1 Model and Architectures

We use the same model and architecture as GPT-2 [RWC⁺19], including the modified initialization, pre-normalization, and reversible tokenization described therein, with the exception that we use alternating dense and locally banded sparse attention patterns in the layers of the transformer, similar to the Sparse Transformer [CGRS19]. To study the dependence of ML performance on model size, we train 8 different sizes of model, ranging over three orders of magnitude from 125 million parameters to 175 billion parameters, with the last being the model we call GPT-3. Previous work [KMH⁺20]

Summary

- How well do Transformers work on long sequences? Not so well.
- How can we make them more efficient? Induce sparsity.
- Many open questions we did not get into:
 - Limitations of sportifying Transformers
 - Alternatives to Transformers (e.g., state-space models)
 - ...

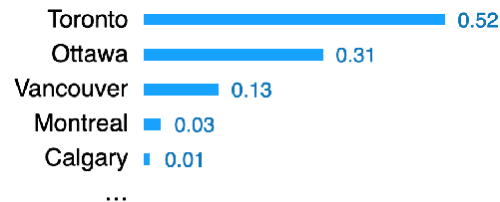
Retrieval-Augmented LMs

Retrieval-based Language Models

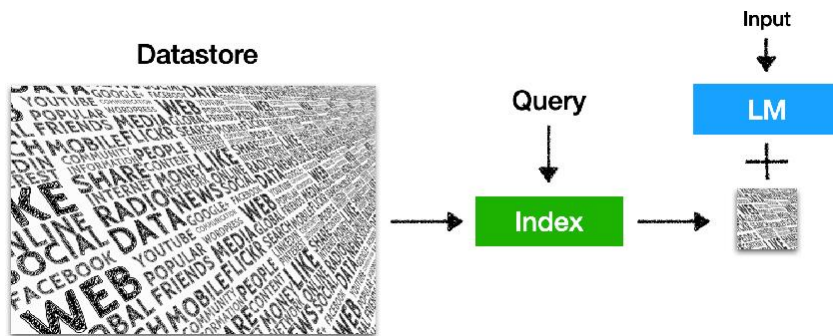
- It is a **language model** $P(x_n | x_1, x_2, \dots, x_{n-1})$

The capital city of Ontario is ____

(can be broadly extended to masked language models or encoder-decoder models)



- It retrieves from an **external datastore** (at least during inference time)



(Also referred to semiparametric and non-parametric models)

Why retrieval-based LMs?

LLMs can't memorize all (long-tail) knowledge in their parameters



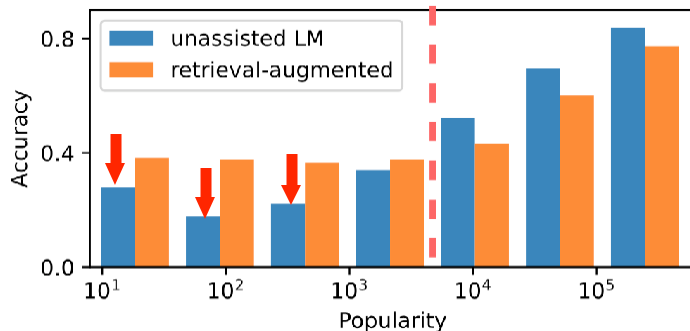
List 5 important papers authored by Geoffrey Hinton



Geoffrey Hinton is a renowned computer scientist ...
Here are five important papers authored by him:

1. **"Learning Internal Representations by Error Propagation"** (with **D. E. Rumelhart** and **R. J. Williams**) - This paper, published in **1986**, .. ✓
2. **"Deep Boltzmann Machines"** (with **R. Salakhutdinov**) - Published in **2009**, .. ✓
- ...
4. **"Deep Learning"** (with Y. Bengio and A. Courville) - Published as a book in **2016**,... ✗
5. **"Attention Is All You Need"** (with V. Vaswani, N. Shazeer, et al.) - Published in **2017**, this paper introduced the Transformer model,... ✗

What is Kathy Saltzman's occupation?



(Mallen et al., 2023)

GPT-3 davinci-003: 20%-30% accuracy

Why retrieval-based LMs?

LLMs' knowledge is easily outdated and hard to update



Who is the CEO of Twitter?



ChatGPT

As of my **knowledge cutoff in September 2021**, the CEO of Twitter is **Jack Dorsey**....



Who is the CEO of Twitter?



About 1,090,000,000 results (0.45 seconds)

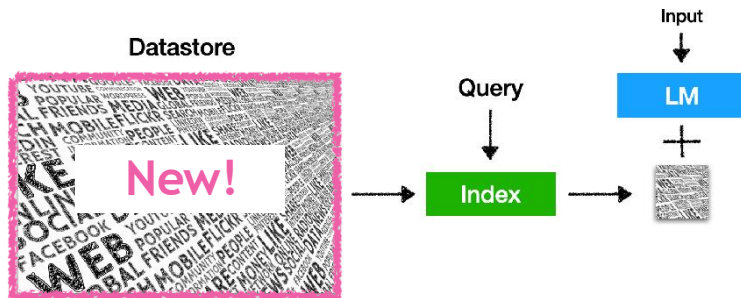
Twitter / CEO

Linda Yaccarino

Jun 5, 2023–



- Existing **knowledge editing** methods are still NOT scalable (**active research!**)
- The datastore can be easily **updated** and **expanded** - even without retraining!



Why retrieval-based LMs?

LLMs' output is challenging to interpret and verify



Create an itinerary for exploring Toronto over a span of two days.



Toronto is a vibrant and diverse city with many attractions and activities to suit different tastes and preferences. Here are some suggestions for creating an itinerary for exploring Toronto over two days:

- **Day 1:** Start your day with a visit to the **CN Tower**, the iconic landmark that offers panoramic views of the city and beyond ¹. You can also experience the thrill of walking on the glass floor or the EdgeWalk, the world's highest hands-free walk ¹. Next, head to the nearby **St. Lawrence Market**, one of the world's best food markets, where you can sample a variety of cuisines and local specialties ². After lunch, take a stroll along **Queen West**, a trendy neighborhood with eclectic shops, galleries, cafes and street art ³. In the evening, enjoy a show at one of the many theaters or comedy clubs in the **Entertainment District**, or catch a game at the **Scotiabank Arena** if you're a sports fan.

- **Day 2:** Explore the history and culture of Toronto at the **Royal Ontario Museum**, one of the largest and most comprehensive museums in North America, featuring exhibits on art, natural history, world cultures and more ⁴. Then, hop on a ferry to the **Toronto Islands**, a group of islands that offer a relaxing escape from the city, with beaches, parks, trails and amusement rides ³ ⁵. You can also rent a bike or kayak to explore the islands at your own pace. For dinner, head to **Chinatown**, one of the largest and most vibrant in North America, where you can find a variety of Asian cuisines and shops ³.

I hope this helps you plan your trip to Toronto. Have fun! 😊

Learn more:

1. cntower.ca

2. travel.usnews.com

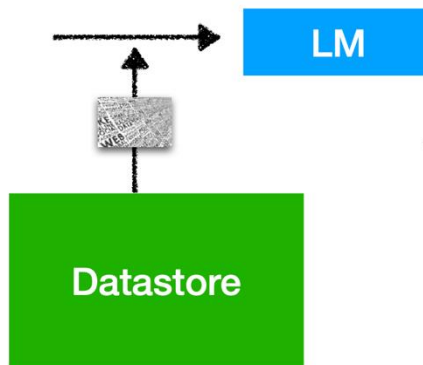
3. bing.com

4. rom.on.ca

5. tripadvisor.com

Why retrieval-based LMs?

LLMs are *large* and expensive to train and run



vs.



Long-term goal: can we possibly reduce the **training** and **inference costs**, and scale down the size of LLMs?

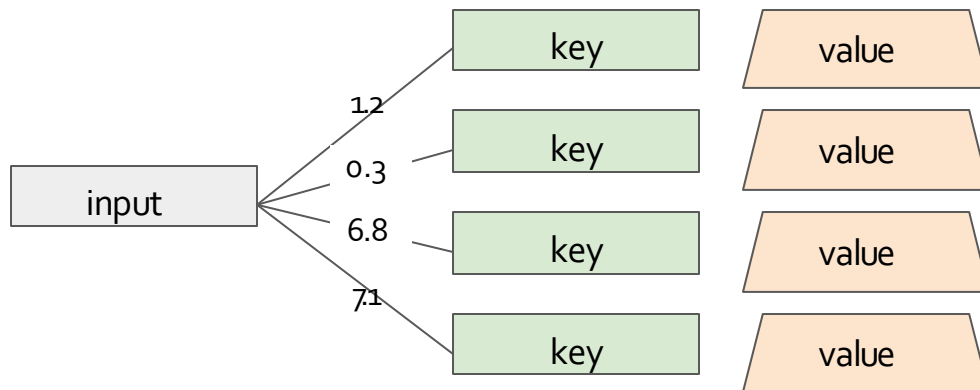
e.g., RETRO (Borgeaud et al., 2021): “obtains comparable performance to GPT-3 on the Pile, despite using **25x fewer parameters**”

What are the Key Design Questions?

- What are your memories?
 - Documents, database records, training examples, etc.
- How to retrieve memories?
 - Use an off-the-shelf search engine (e.g. Google, StackOverflow).
 - How to train your own memory retriever.
- How to use retrieved memories?
 - "Text fusion"
 - Common failure modes:
 - Underutilization: model ignores retrieved memories.
 - Overreliance: model depends too much on memories!

Anatomy of a Neural Retriever

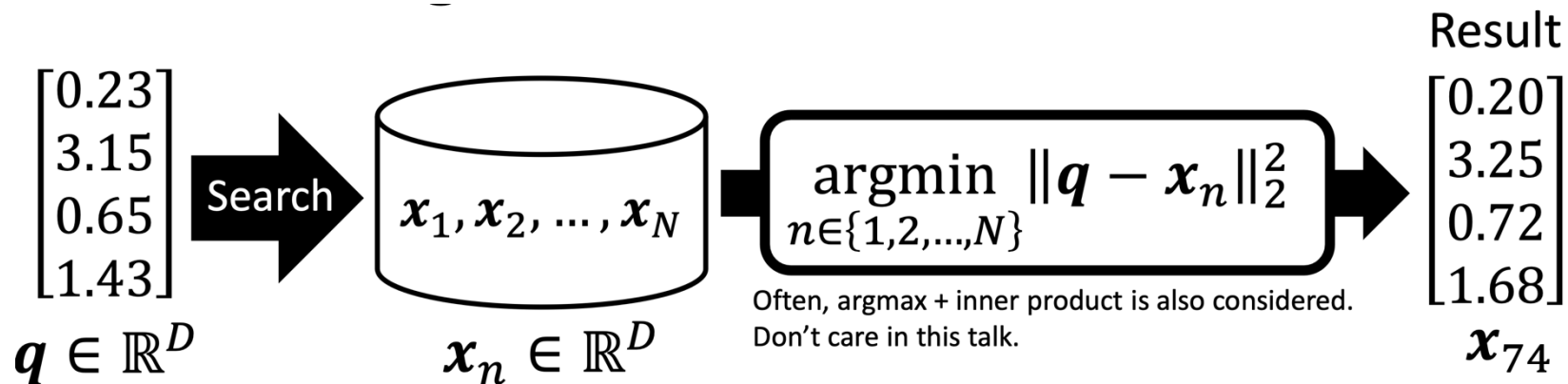
1. Score the input against each key.
2. Return the value for the highest scoring key.



A similarity function:

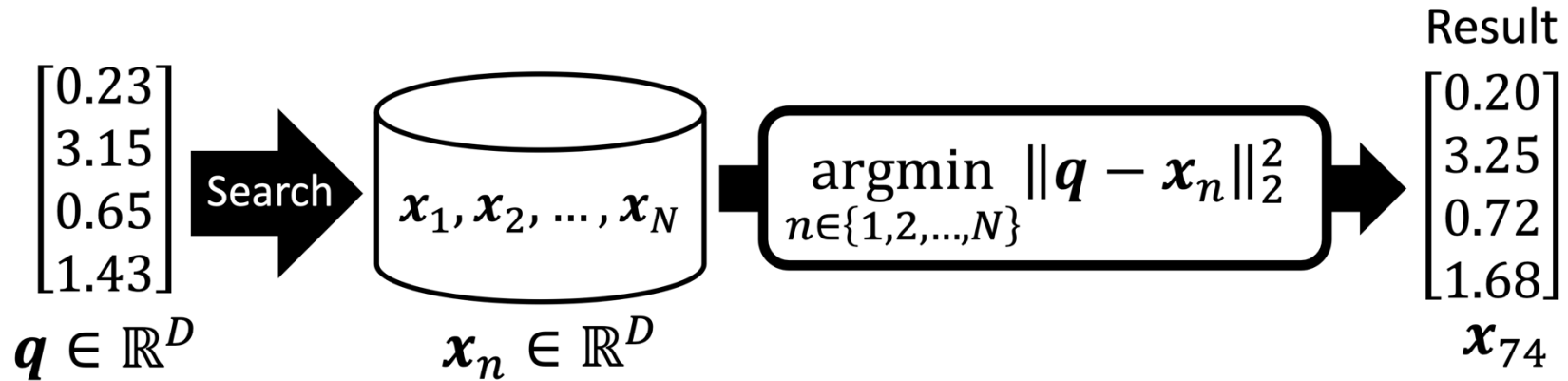
$$\text{sim}(\text{input}, \text{key}) \rightarrow \text{score}$$

Finding Nearest Neighbors



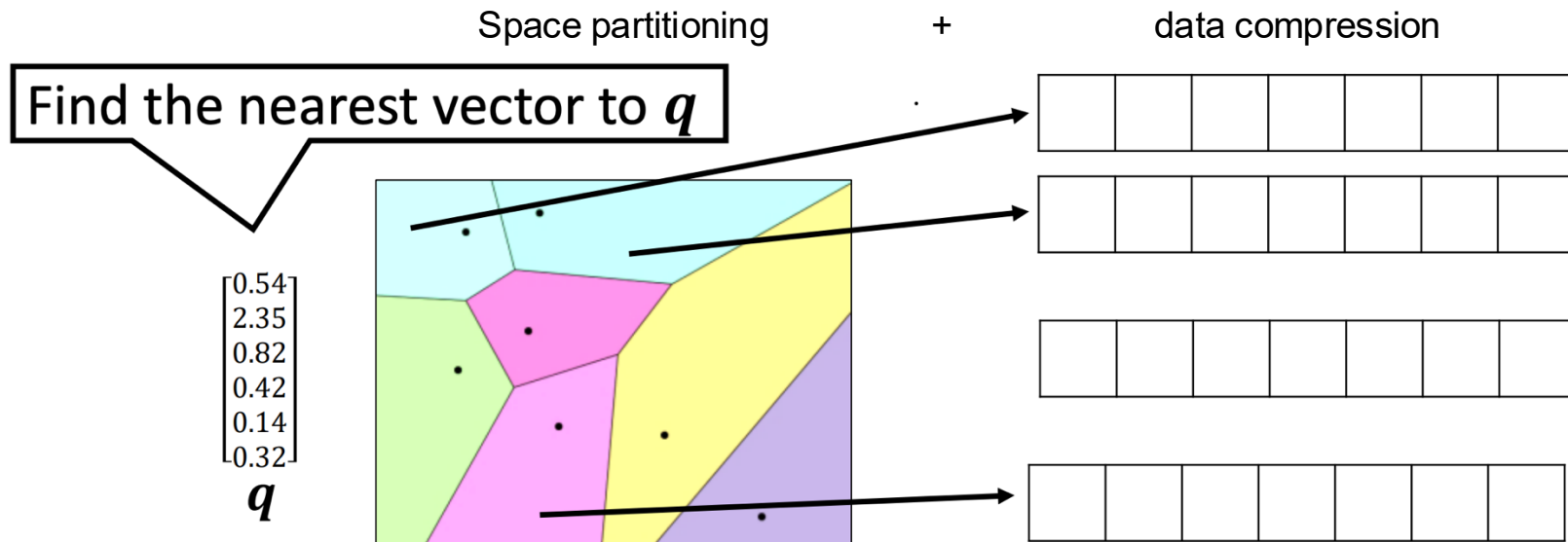
- N D -dim database vectors: $\{x_n\}_{n=1}^N$
- Given a query q , find the closest vector from the database
- One of the fundamental problems in computer science
- Solution: linear scan, $O(ND)$, slow ☹

Approximate Finding Nearest Neighbors



- Faster search
- Don't necessarily have to be exact neighbors
- Trade off: runtime, accuracy, and memory-consumption

Approximate NNs: Algorithms, Libraries, Services



- More information: <https://github.com/facebookresearch/faiss/wiki>

Approximate NNs: Algorithms, Libraries, Services

Algorithm

- Scientific paper
- Math
- Often, by researchers

Product Quantization +
Inverted Index (PQ, IVFPQ)
[Jégou+, TPAMI 2011]

Hierarchical Navigable
Small World (HNSW)
[Malkov+, TPAMI 2019]

ScaNN (4-bit PQ)
[Guo+, ICML 2020]

Library

- Implementations of algorithms
- Usually, a search function only
- By researchers, developers, etc

faiss

NMSLIB

hnswlib

ScaNN

Service (e.g., vector DB)

- Library + (handling metadata, serving, scaling, IO, CRUD, etc)
- Usually, by companies

Pinecone

Qdrant

Milvus

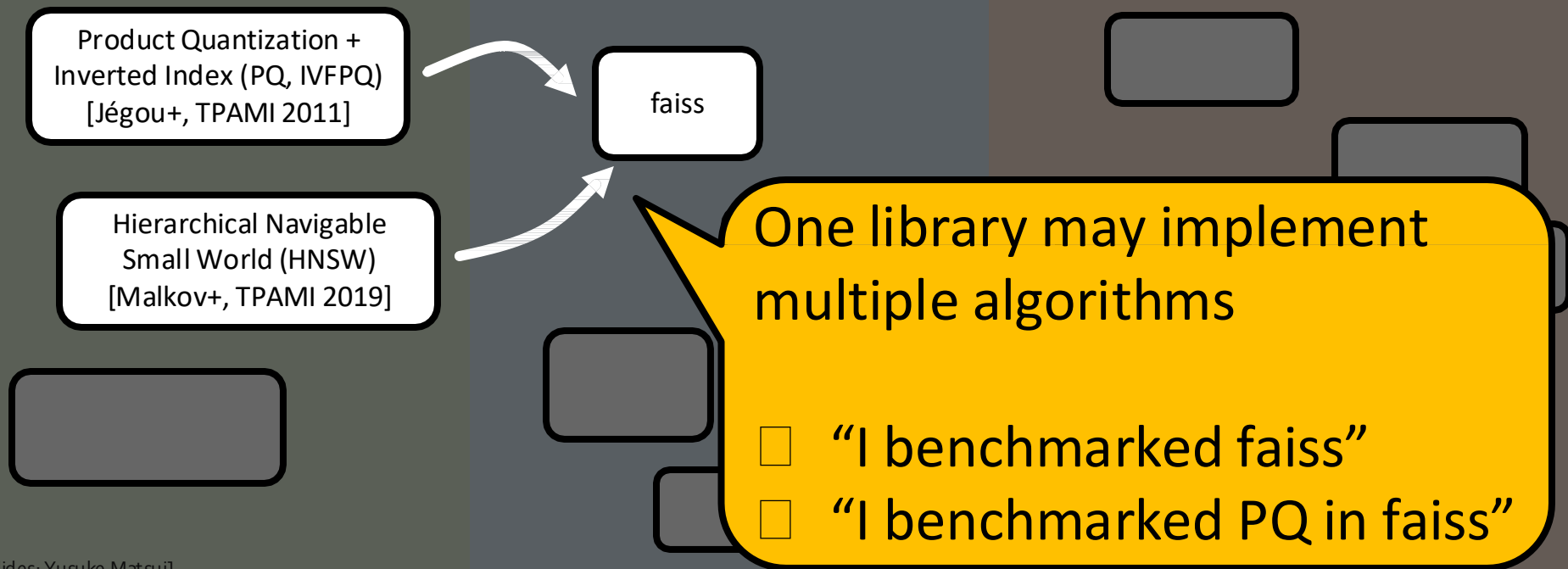
jina

Vald

Vertex AI
Matching Engine

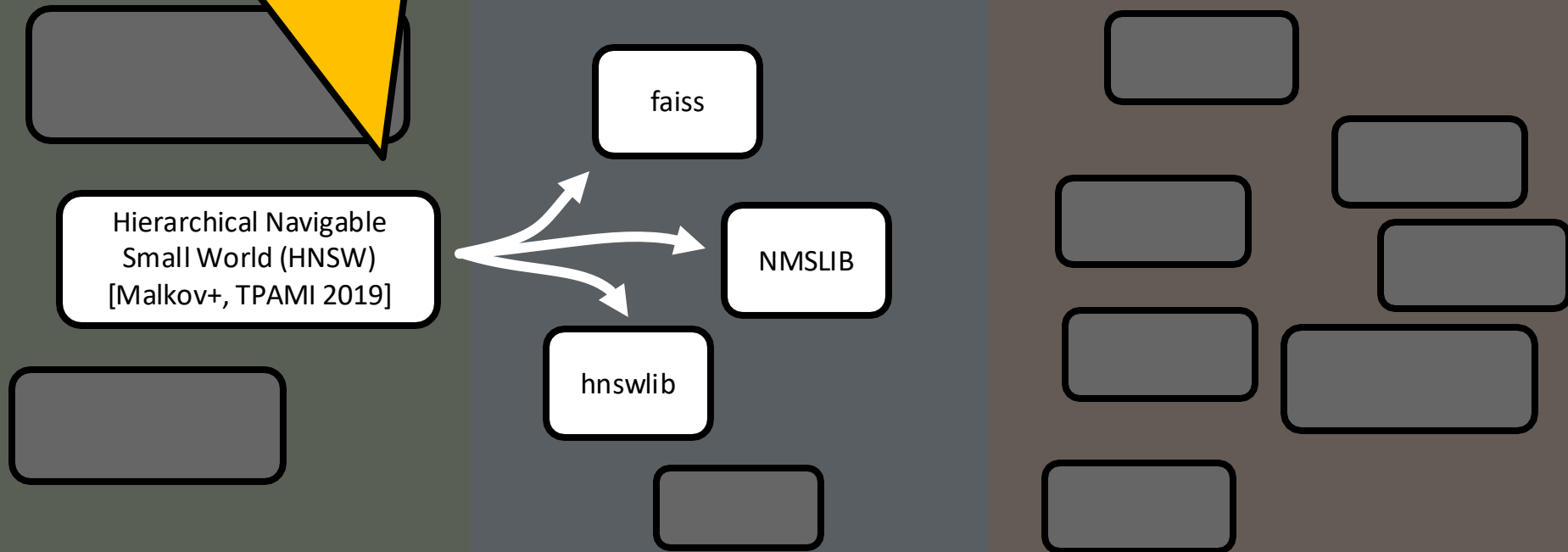
Weaviate

Three levels of technology



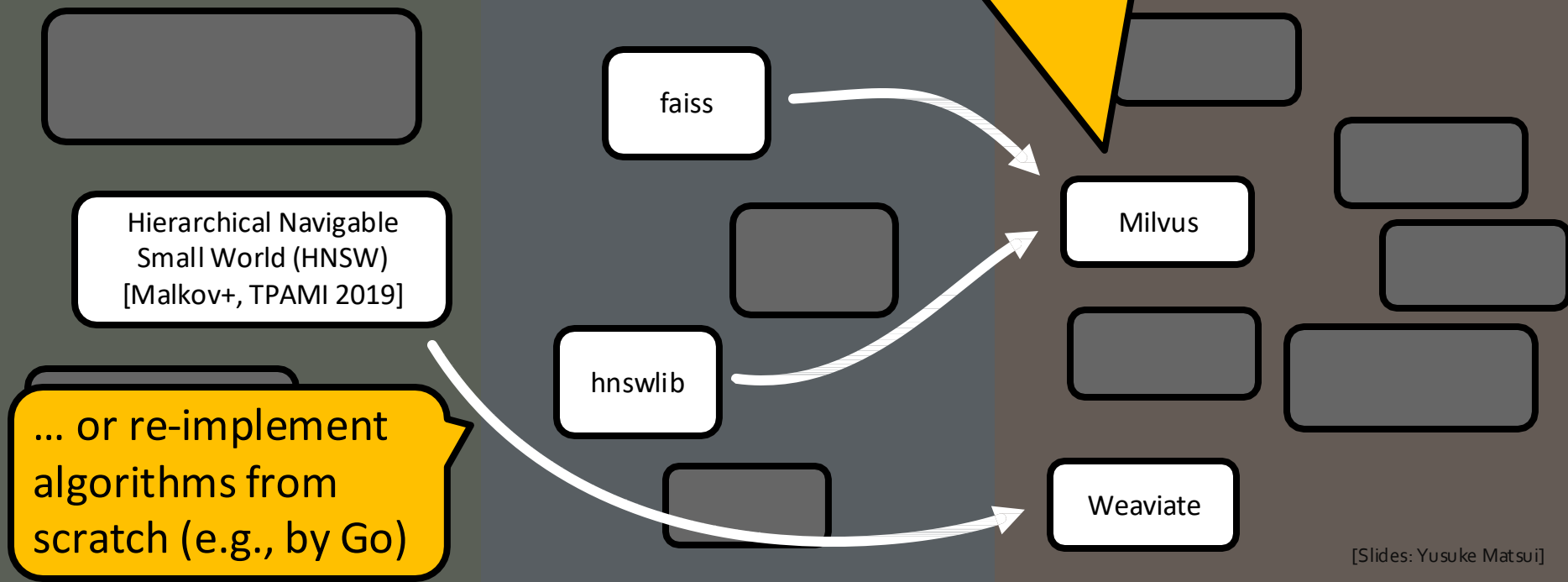
Three levels of technology


One algorithm may be implemented in multiple libraries




Three levels of technology

One service may use some libraries





Let's assume that we have our
retrieval engine and data ready



Defining Similarity Metrics: Approach 1

$$\text{sim}(I, M) = \text{tf}(I, M) \times \log N / d_I$$

- $\text{tf}(I, M)$: # of occurrences of I in M .
- N : # of documents
- d_I : # of documents that contain I .

- **Advantages:**

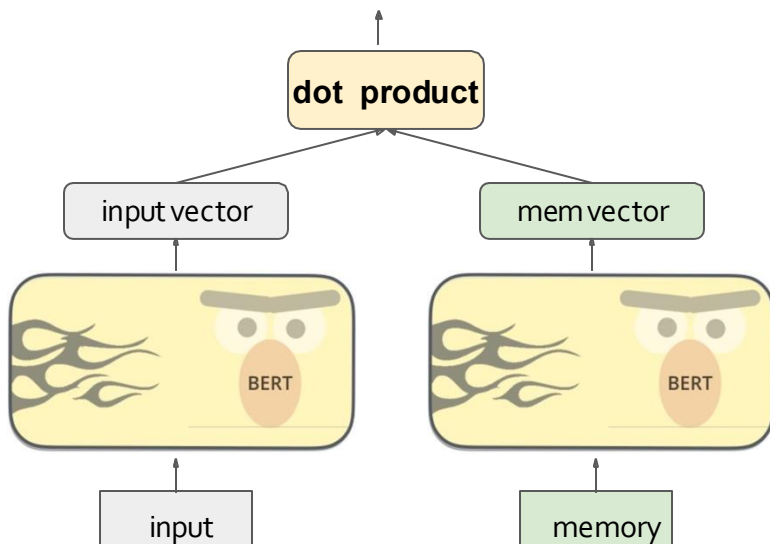
- Not differentiable; you can't easily optimize it.

- **Disadvantages:**

- Simple and fast. Often the first idea that you should try!

Defining Similarity Metrics: Approach 2

$$\text{sim}(\mathbf{I}, \mathbf{M}) = \text{Encoder}(\mathbf{I}) \times \text{Encoder}(\mathbf{M})$$



- **Advantages:**

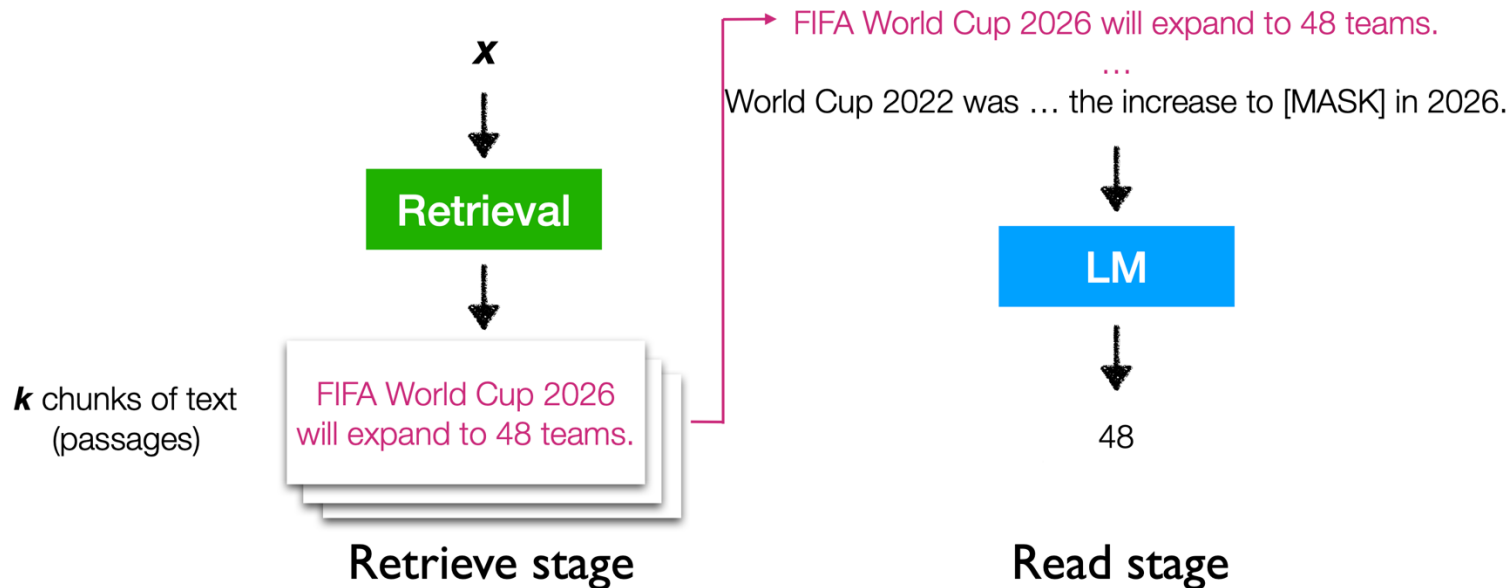
- Differentiable -- can optimize with gradient descent.

- **Disadvantages:**

- Works well for data on which your LM is pre-trained on.

Retrieval-Augmented LM

- x = World Cup 2022 was the last before the increase to [MASK] in the 2026 tournament.



Retrieval-Augmented LM

FIFA World Cup 2026
will expand to 48 teams.

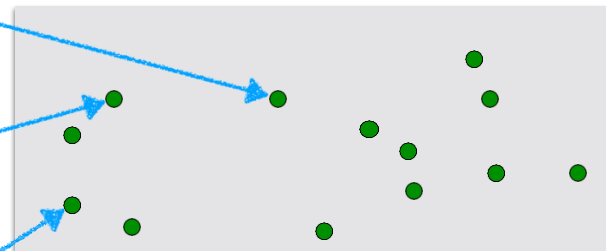
In 2022, the 32
national teams involved
in the tournament.

Team USA celebrated
after winning its match
against Iran ...

Encoder

Encoder

Encoder



$$\mathbf{z} = \text{Encoder}(z)$$

Wikipedia

13M chunks (passages)
(called *documents* in the paper)

Retrieval-Augmented LM

x = World Cup 2022 was ... the increase to [MASK] in 2026.

FIFA World Cup 2026
will expand to 48 teams.

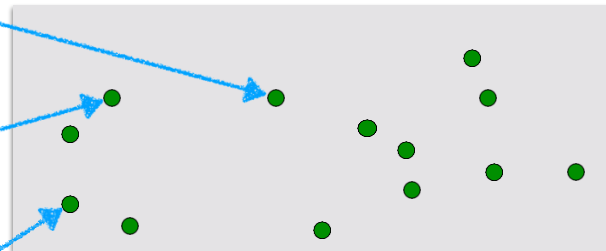
In 2022, the 32
national teams involved
in the tournament.

Team USA celebrated
after winning its match
against Iran ...

Encoder

Encoder

Encoder



$z = \text{Encoder}(z)$

Wikipedia

13M chunks (passages)
(called *documents* in the paper)

Retrieval-Augmented LM

\mathbf{x} = World Cup 2022 was ... the increase to [MASK] in 2026.

FIFA World Cup 2026
will expand to 48 teams.

In 2022, the 32
national teams involved
in the tournament.

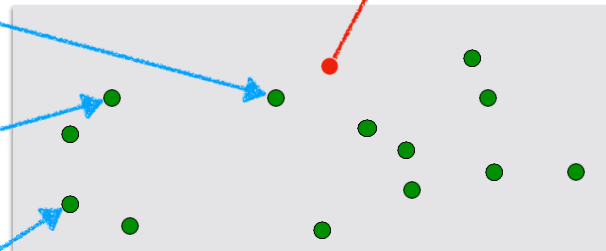
Team USA celebrated
after winning its match
against Iran ...

Encoder

Encoder

Encoder

Encoder



$$\mathbf{z} = \text{Encoder}(\mathbf{z})$$

$$\mathbf{x} = \text{Encoder}(\mathbf{x})$$

Wikipedia

13M chunks (passages)
(called *documents* in the paper)

Retrieval-Augmented LM

\mathbf{x} = World Cup 2022 was ... the increase to [MASK] in 2026.

FIFA World Cup 2026
will expand to 48 teams.

In 2022, the 32
national teams involved
in the tournament.

Team USA celebrated
after winning its match
against Iran ...

Encoder

Encoder

Encoder

Encoder

Fast nearest neighbor
search

$\mathbf{z} = \text{Encoder}(\mathbf{z})$

$\mathbf{x} = \text{Encoder}(\mathbf{x})$

Wikipedia

13M chunks (passages)
(called *documents* in the paper)

Retrieval-Augmented LM

\mathbf{x} = World Cup 2022 was ... the increase to [MASK] in 2026.

FIFA World Cup 2026
will expand to 48 teams.

In 2022, the 32
national teams involved
in the tournament.

Team USA celebrated
after winning its match
against Iran ...

Encoder

Encoder

Encoder

Encoder

Fast nearest neighbor
search

$$\mathbf{z} = \text{Encoder}(\mathbf{z})$$

$$\mathbf{x} = \text{Encoder}(\mathbf{x})$$

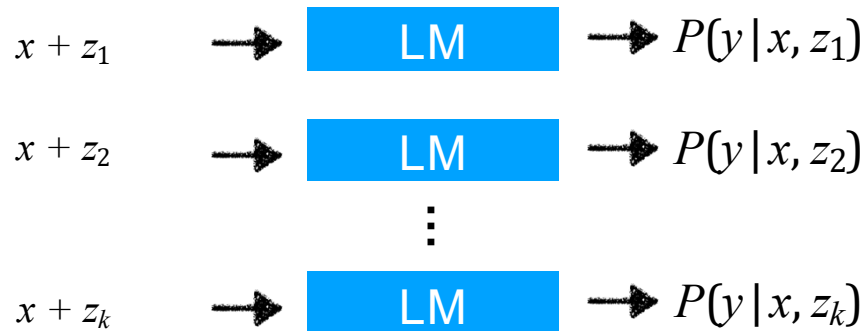
$$\mathbf{z}_1, \dots, \mathbf{z}_k = \text{argTop-}k(\mathbf{x} \cdot \mathbf{z})$$

k retrieved chunks

Wikipedia

13M chunks (passages)
(called *documents* in the paper)

Retrieval-Augmented LM

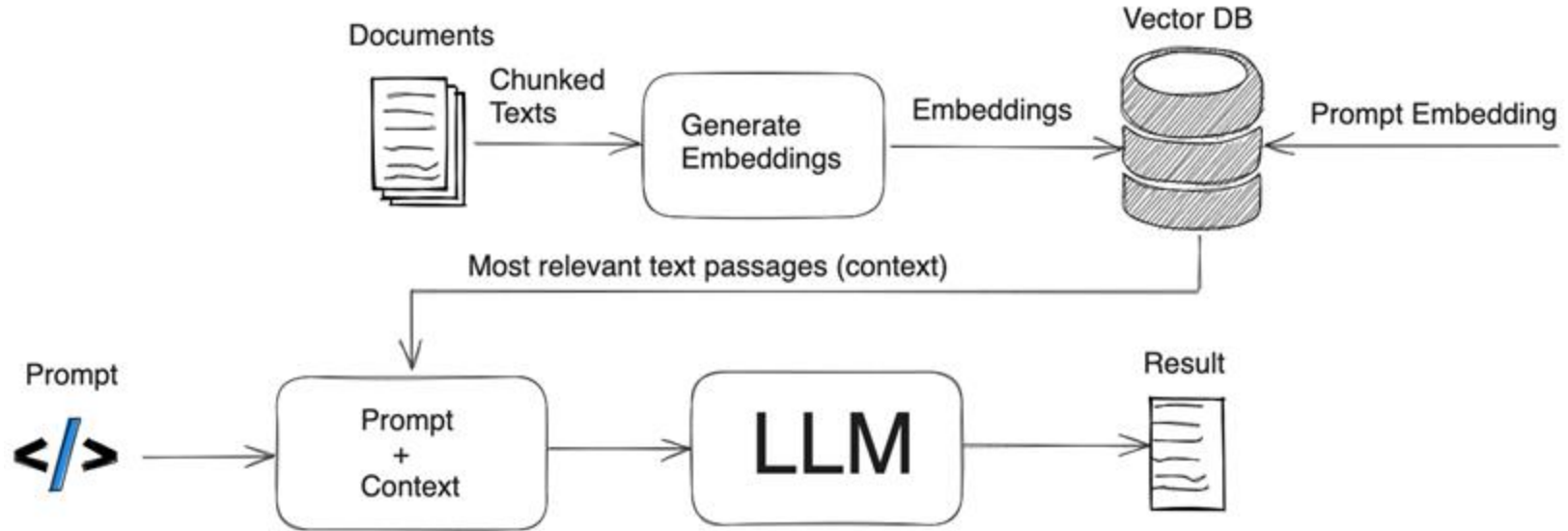


Need to approximate
→ Consider top k chunks only

$$\sum_{z \in \mathcal{D}} \underbrace{P(z | x)}_{\text{from the retrieve stage}} \underbrace{P(y | x, z)}_{\text{from the read stage}}$$

0 if not one of top k

RAG: A more simplified pipeline



Retrieval-Augmented LM: Variants

What to retrieve?

- Chunks
- Tokens
- Others

How to use retrieval?

- Input layer
- Intermediate layers
- Output layer

When to retrieve?

- Once
- Every n tokens ($n > 1$)
- Every token

Retrieval-Augmented LM: Common Variant

What to retrieve?

- **Chunks**
- Tokens
- Others

How to use retrieval?

- **Input layer**
- Intermediate layers
- Output layer

When to retrieve?

- **Once**
- Every n tokens ($n > 1$)
- Every token

IR in the Middle of LM

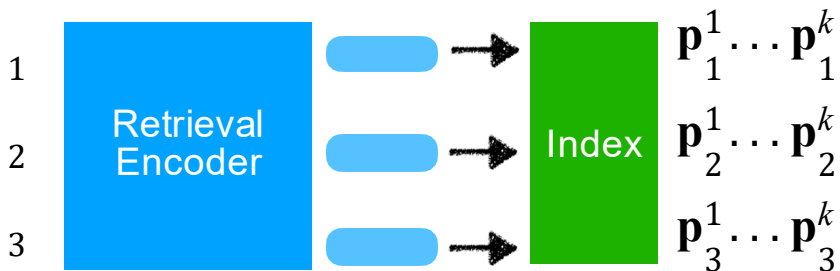
x = World Cup 2022 was the last with 32 teams, before the increase to

1

2

3

(k chunks of text per split)



IR in the Middle of LM

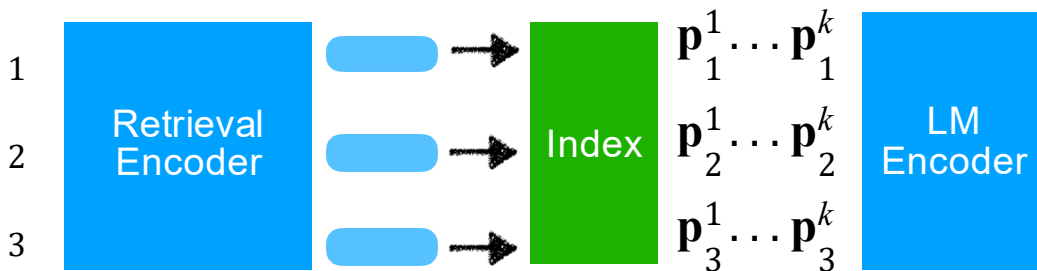
x = World Cup 2022 was the last with 32 teams, before the increase to

1

2

3

(k chunks of text per split)



IR in the Middle of LM

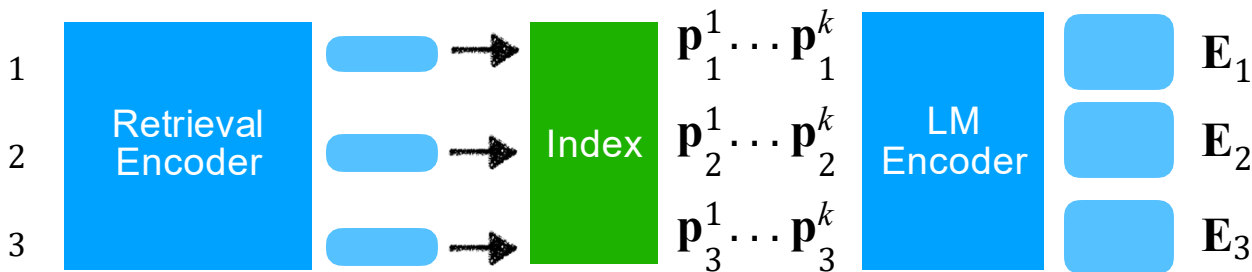
x = World Cup 2022 was the last with 32 teams, before the increase to

1

2

3

(k chunks of text per split)



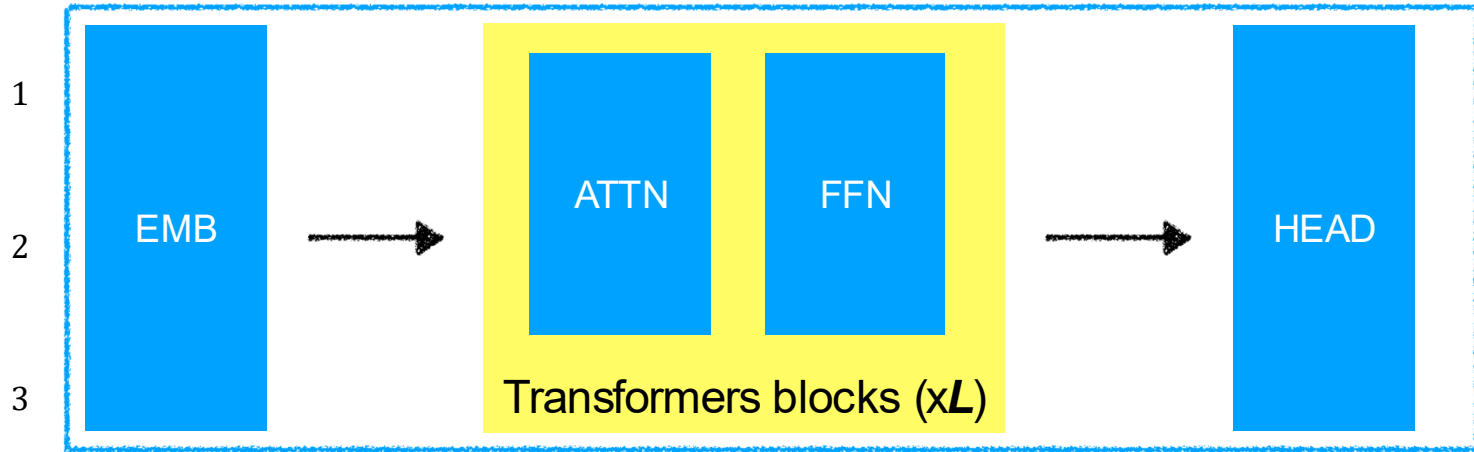
(A $r \times k \times d$ matrix)

(r = # tokens per text chunk)

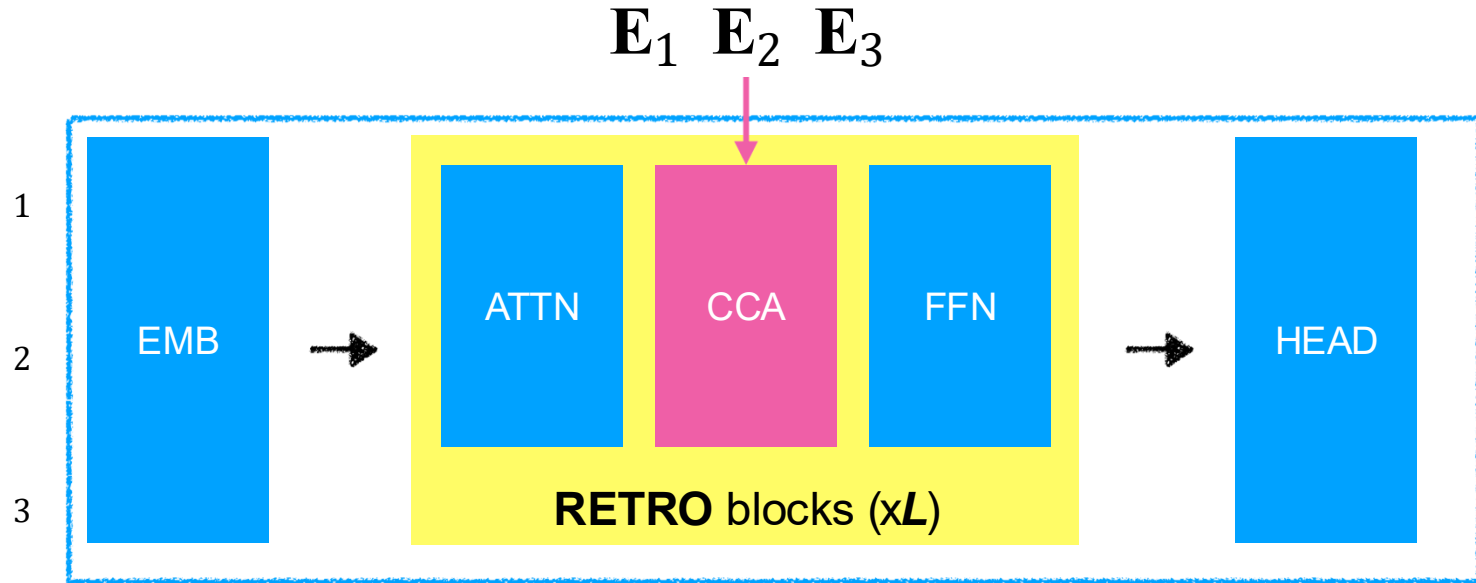
(d = hidden dimension)

(k = # retrieved chunks per split)

Regular Decoder



Regular Decoder with IR Embeddings



Chunked CrossAttention (CCA)

Results

Perplexity: The lower the better

Model	Retrieval Set	#Database tokens	#Database keys	Valid	Test
Adaptive Inputs (Baevski and Auli, 2019)	-	-	-	17.96	18.65
SPALM (Yogatama et al., 2021)	Wikipedia	3B	3B	17.20	17.60
kNN-LM (Khandelwal et al., 2020)	Wikipedia	3B	3B	16.06	16.12
Megatron (Shoeybi et al., 2019)	-	-	-	-	10.81
Baseline transformer (ours)	-	-	-	21.53	22.96
kNN-LM (ours)	Wikipedia	4B	4B	18.52	19.54
RETRO	Wikipedia	4B	0.06B	18.46	18.97
RETRO	C4	174B	2.9B	12.87	10.23
RETRO	MassiveText (1%)	18B	0.8B	18.92	20.33
RETRO	MassiveText (10%)	179B	4B	13.54	14.95
RETRO	MassiveText (100%)	1792B	28B	3.21	3.92

Significant improvements by retrieving from 1.8 trillion tokens

Results

Perplexity: The lower the better



Model	Retrieval Set	#Database tokens	#Database keys	Valid	Test
Adaptive Inputs (Baevski and Auli, 2019)	-	-	-	17.96	18.65
SPALM (Yogatama et al., 2021)	Wikipedia	3B	3B	17.20	17.60
kNN-LM (Khandelwal et al., 2020)	Wikipedia	3B	3B	16.06	16.12
Megatron (Shoeybi et al., 2019)	-	-	-	-	10.81
Baseline transformer (ours)	-	-	-	21.53	22.96
kNN-LM (ours)	Wikipedia	4B	4B	18.52	19.54
RETRO	Wikipedia	4B	0.06B	18.46	18.97
RETRO	C4	174B	2.9B	12.87	10.23
RETRO	MassiveText (1%)	18B	0.8B	18.92	20.33
RETRO	MassiveText (10%)	179B	4B	13.54	14.95
RETRO	MassiveText (100%)	1792B	28B	3.21	3.92

Significant improvements by retrieving from 1.8 trillion tokens

Retrieval-Augmented LM

What to retrieve?


- **Chunks**
- Tokens
- Others

How to use retrieval?


- Input layer
- **Intermediate layers**
- Output layer

When to retrieve?

- Once
- **Every n tokens ($n > 1$)**
- Every token

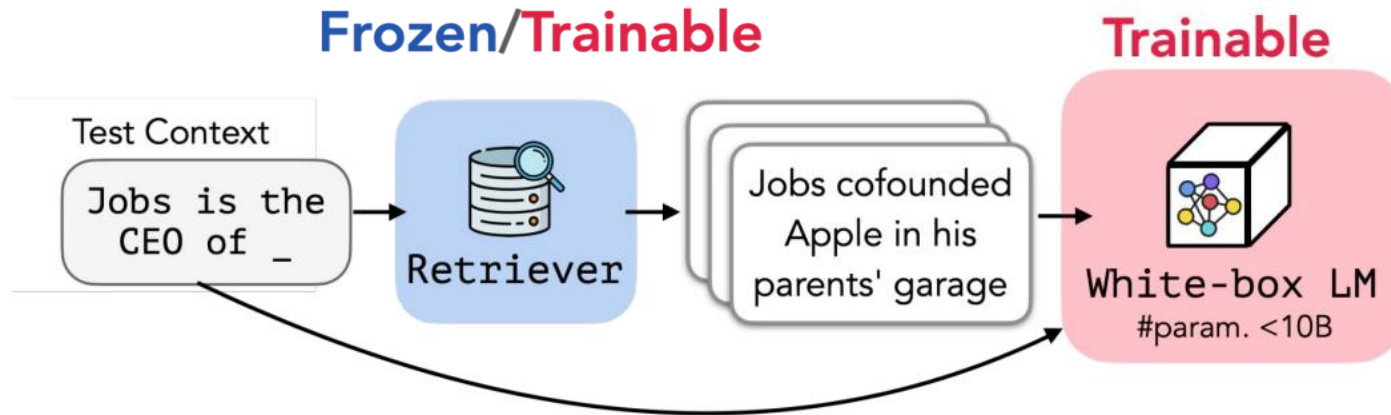


How do you train
these models?



End-to-end Training

- There are various ideas in the literature for how to train these models efficiently and in an end-to-end fashion.



Main takeaways

- How do we enable LMs to utilize external knowledge?
 - Retrieval-augmented language models
- A retriever is a function, $f(\text{input}, \text{memory}) \rightarrow \text{score}$
- What we did not discuss:
 - Attribution: Tracing decisions to the source knowledge
 - How to modify the knowledge
 - Conflicting knowledge
 - Editing knowledge
 - More efficient scaling
 -

Practical tips for building IR systems

- Try the easiest approach.
- For example, you can fire up an Elastic Search IR with a single line of command:
- ...
- And here is how you query.

Practical tips for building IR systems

- Try the easiest approach.
- For example, you can fire up an Elastic Search IR with a single line of command:
- ...
- And here is how you query.

Always try the easiest test

- Sample sentences from indexed corpus and see if the IR engine is able to retrieve them.
 - (may be useful to define a few retrieval metrics here)
- If it's robust, try queries that are perturbations of the indexed content.

Retrieval-Augmented Generation Limitations

- **However, RAG has several limitations:**
- **Extra complexity**
 - It needs to define what to store in the database.
- **Lacks contextual understanding.**
 - Summarization / Complex Reading Comprehension
- **Stateless.**
 - It does not hold memory of what is previously retrieved.



How were retrieval engines
designed before NNs?



Classic Information Retrieval

```
t = token
d = document (movie review)
term_frequency(t,d) = number of times t occurs in d
document_frequency(t) = # documents t occurs in
N = number of documents
inverse_document_frequency(t) = N / document_frequency(t)
tf-idf(t,d) = tf(t,d) x inverse_document_frequency(t)
```

- A vector for d is produced by stacking $\text{tf-idf}(t,d)$ in an order determined by the token vocab IDs
- **BM25** adds two parameters: k , a knob that adjust the balance between term frequency and IDF, and b , which controls the importance of document length normalization [[Kamphuis et al. \(2020\)](#)]
- We can represent a question and each document in the external databases like this, and use vector similarity measurements to find the most similar documents for a given question

Inverted Index

t = token

d = document (movie review)

$\text{term_frequency}(t, d)$ = number of times t occurs in d

$\text{document_frequency}(t)$ = # documents t occurs in

N = number of documents

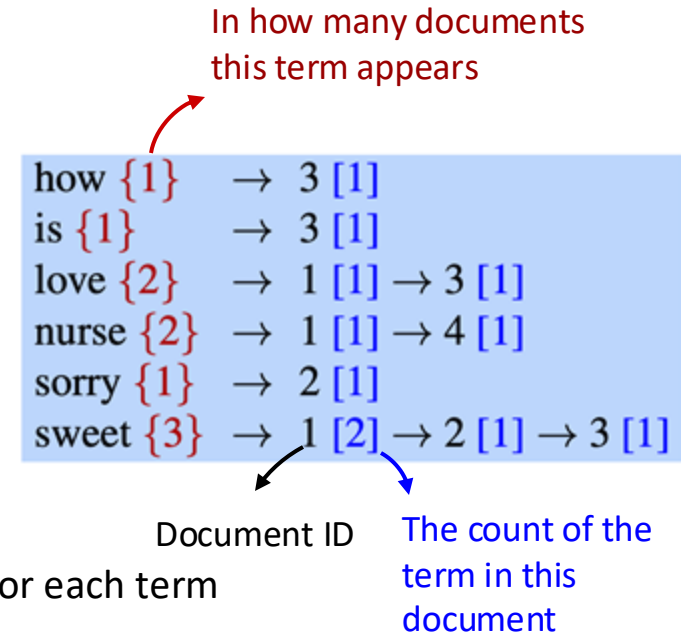
$\text{inverse_document_frequency}(t) = N / \text{document_frequency}(t)$

$\text{tf-idf}(t, d) = \text{tf}(t, d) \times \text{inverse_document_frequency}(t)$

- We need to efficiently find documents that contain words in the question
- The basic search problem in IR is thus to find all documents that contain a term
- The data structure for this task is the **inverted index**, which we use for making this search efficient, and also conveniently storing useful information like the document frequency and the count of each term in each document

Inverted Index (cont.)

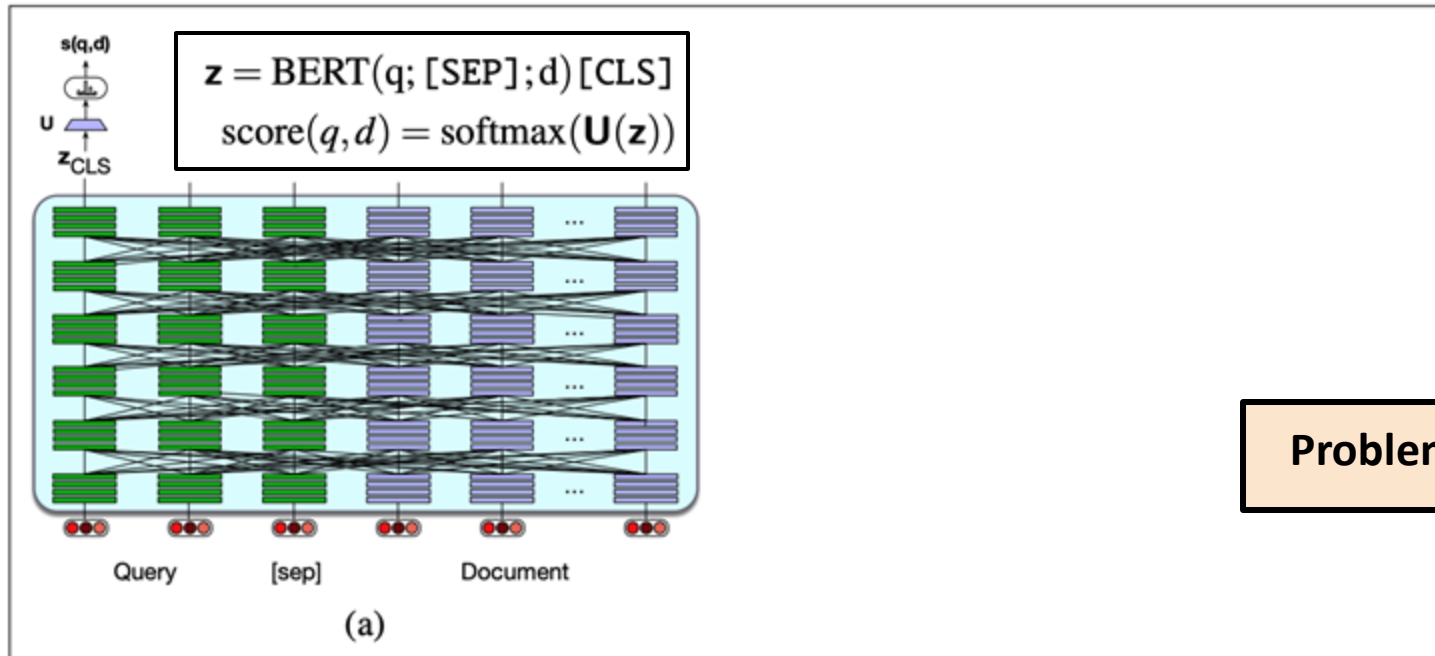
- An **inverted index**, given a query term, gives a list of documents that contain the term
- It consists of two parts:
 - Dictionary
 - Postings
- **Dictionary** is a list of terms (designed to be efficiently accessed), each pointing to a postings list for the term.
 - The dictionary can also store the document frequency for each term
- A **postings list** is the list of document IDs associated with each term
 - It can also contain information like the term frequency or even the exact positions of terms in the document



Why could this classic IR be limited?

- tf-idf/BM25 algorithms work only if there is exact overlap of words between the question and document
- What if a question contains a lot of synonyms of tokens in a relevant document?
- **Vocabulary mismatch problem:** The user posing a query (or asking a question) needs to guess exactly what words the writer of the answer might have used
 - Instead of (sparse) word-count vectors, using (dense) embeddings

Two ways to dense retrieval



Problem?

Figure 14.7 Two ways to do dense retrieval, illustrated by using lines between layers to schematically represent self-attention: (a) Use a single encoder to jointly encode query and document and finetune to produce a relevance score with a linear layer over the CLS token. This is too compute-expensive to use except in rescoring (b) Use separate encoders for query and document, and use the dot product between CLS token outputs for the query and document as the score. This is less compute-expensive, but not as accurate.

Two ways to dense retrieval

Figure: [Jurafsky and Martin]

$$\begin{aligned} \mathbf{z}_q &= \text{BERT}_Q(q) [\text{CLS}] \\ \mathbf{z}_d &= \text{BERT}_D(d) [\text{CLS}] \\ \text{score}(q, d) &= \mathbf{z}_q \cdot \mathbf{z}_d \end{aligned}$$

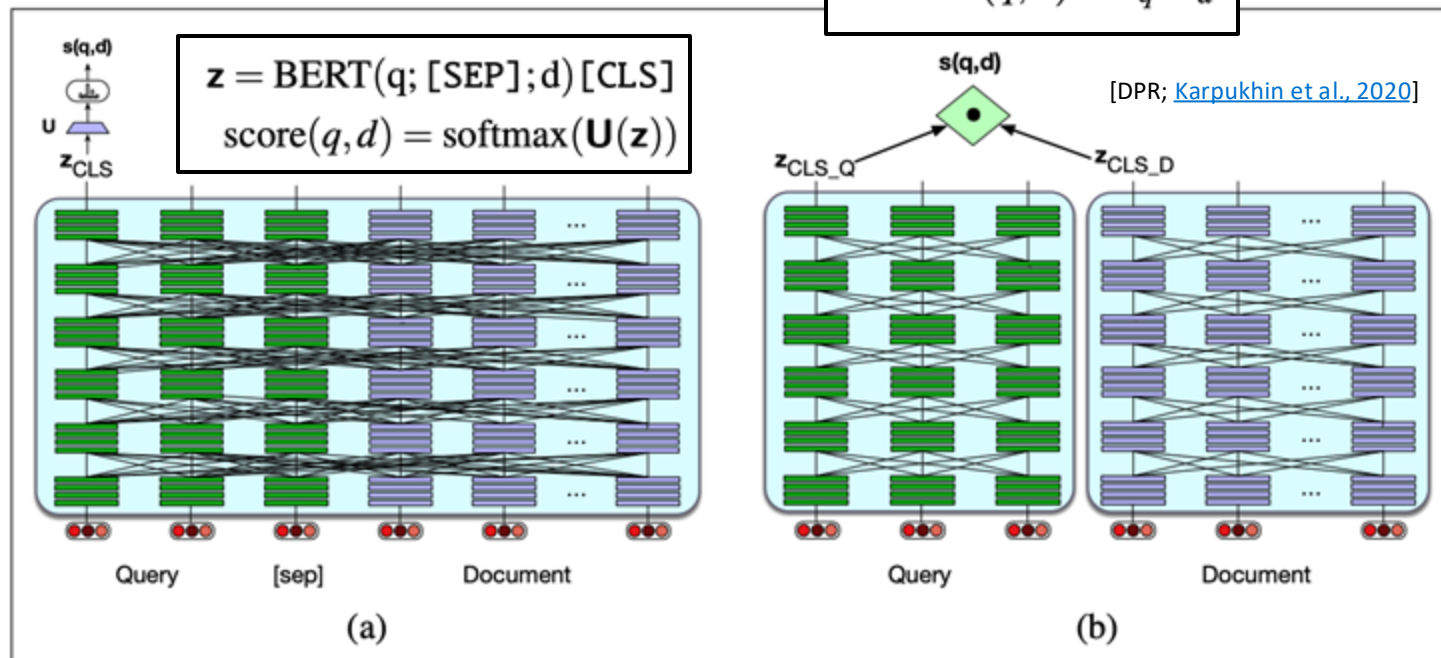


Figure 14.7 Two ways to do dense retrieval, illustrated by using lines between layers to schematically represent self-attention: (a) Use a single encoder to jointly encode query and document and finetune to produce a relevance score with a linear layer over the CLS token. This is too compute-expensive to use except in rescoring (b) Use separate encoders for query and document, and use the dot product between CLS token outputs for the query and document as the score. This is less compute-expensive, but not as accurate.

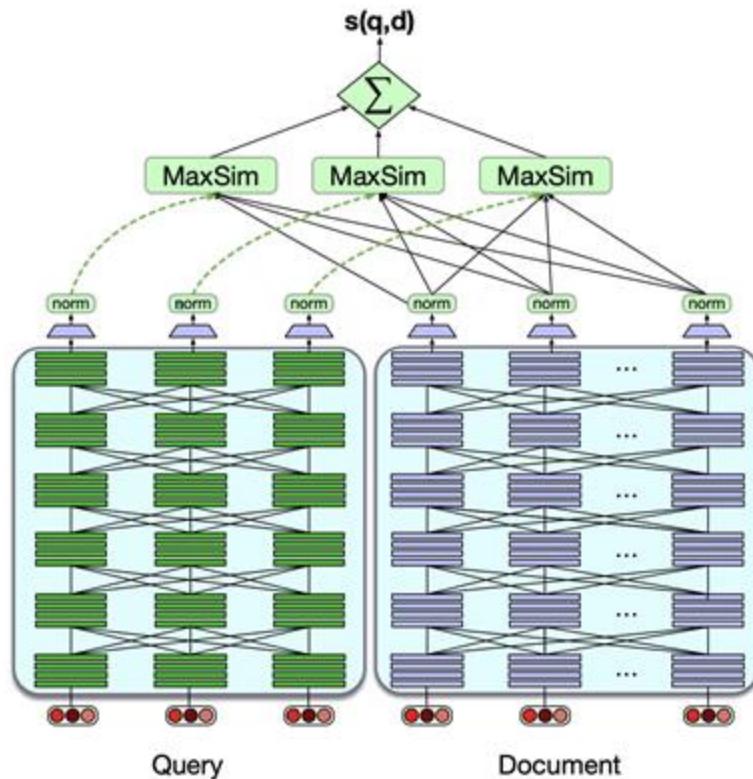
Another way to the second approach – ColBERT

[Khattab et al., 2021]

$$\text{score}(q, d) = \sum_{i=1}^N \max_{j=1}^m \mathbf{E}_{q_i} \cdot \mathbf{E}_{d_j}$$

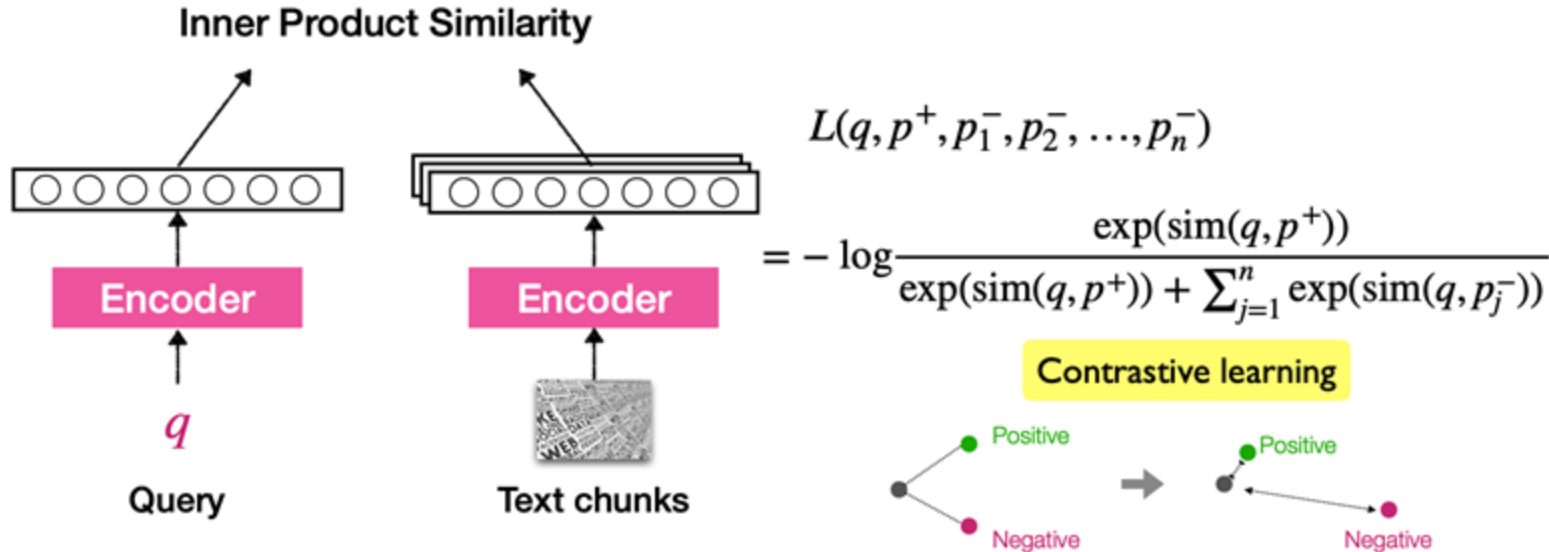
BERT output vectors
rescaled to unit length

Essentially, for each token
in q , ColBERT finds the
most contextually similar
token in d , & then sums up
these similarities



Training dense retrievals

- Train the model to maximize the score for positive documents and minimize the score for negative ones:



No “Free Lunch” in Long Context



- No single approach is universally optimal for long context learn
- There is always a tradeoff between performance and cost.
- It makes long context learning a challenging problem.