

# Adaptability

Long context and retrieval-augmented code generation

---

Sean Welleck

Neural Code Generation  
Carnegie Mellon University  
March 19, 2024

# Code generation in 2024



Figure 1: <https://www.youtube.com/watch?v=SSnsmqIj1MI>

Adapt to:

- New code
- New repositories
- New libraries
- New languages
- ...

$$p_{\theta}(y|x, C)$$

Example contexts  $C$  :

- Preceding file contents
- Repository
- New code developed after training
- Documentation, examples, tutorials, etc...
- Discuss: any others?

$$p_{\theta}(y|x, C)$$

Example contexts  $C$  :

- **Preceding file contents**
  - E.g. HuggingFace modeling\_utils.py :  $\approx$  48k tokens
- **Repository**
  - E.g. *three.js*:  $\approx$  800k tokens

$$p_{\theta}(y|x, C)$$

Example contexts  $C$  :

- **Preceding file contents**
  - E.g. HuggingFace modeling\_utils.py :  $\approx$  48k tokens
- **Repository**
  - E.g. *three.js*:  $\approx$  800k tokens
- **New code developed after training**
  - Which data is relevant?
- **Documentation, examples, tutorials, etc...**
  - Which data is relevant?

# Today's lecture

1. Part I: long-context generation
2. Part II: intro to retrieval-augmented generation
  - Case study: retrieving documentation

# Today's lecture

1. Part I: long-context generation
2. Part II: intro to retrieval-augmented generation
  - Case study: retrieving documentation

**Next week:** *Student discussion* on repository-level generation, retrieval augmented generation for code



# I. Long-context generation

---

# Long-context generation

- **Standard transformers** → 2048 tokens
- **Efficient attention** → millions of tokens
  - Example recipe: Long-context fine-tuning → 100k tokens

Key problem in standard transformers:

- Memory bottlenecks (e.g., in attention)

A transformer contains a sequence of 'transformer blocks'.

- Each block has an **attention** layer and a **feedforward layer**

# Standard transformers | recap

```
73  class Block(nn.Module):
74      """ an unassuming Transformer block """
75
76  def __init__(self, config):
77      super().__init__()
78      self.ln_1 = nn.LayerNorm(config.n_embd)
79      self.attn = CausalSelfAttention(config)
80      self.ln_2 = nn.LayerNorm(config.n_embd)
81      self.mlp = nn.ModuleDict(dict(
82          c_fc    = nn.Linear(config.n_embd, 4 * config.n_embd),
83          c_proj  = nn.Linear(4 * config.n_embd, config.n_embd),
84          act     = NewGELU(),
85          dropout = nn.Dropout(config.resid_pdrop),
86      ))
87      m = self.mlp
88      self.mlpf = lambda x: m.dropout(m.c_proj(m.act(m.c_fc(x)))) # MLP forward
89
90  def forward(self, x):
91      x = x + self.attn(self.ln_1(x))
92      x = x + self.mlpf(self.ln_2(x))
93      return x
```

Figure 2: Transformer block from *minGPT*

<https://github.com/karpathy/minGPT/blob/master/mingpt/model.py>

- Let  $Q, K, V \in \mathbb{R}^{s \times d}$ ,  $s$  is sequence length

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V$$

## Standard transformers | Problem : quadratic attention

- Let  $Q, K, V \in \mathbb{R}^{s \times d}$ ,  $s$  is sequence length

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V$$

For a single query  $q$ , write as  $\text{Attention}(q, k_1, \dots, k_s, v_1, \dots, v_s)$

$$\text{Attention}(q) = \sum_{i=1}^s v_i w_i$$

where  $w_i = \frac{\exp(w'_i)}{\sum_{j=1}^s \exp(w'_j)}$ ,

$$w'_i = \text{dot}(q, k_i)$$

Intuition: weighted sum of values  $v$ , with weights determined by similarity of query and keys

# Standard transformers | Problem : quadratic attention

- Let  $Q, K, V \in \mathbb{R}^{s \times d}$ ,  $s$  is sequence length

$$\underbrace{\text{Attention}(Q, K, V)}_{\mathbb{R}^{s \times d}} = \text{softmax} \left( \underbrace{\begin{pmatrix} \underbrace{QK^T}_{\mathbb{R}^{s \times s}} \\ \underbrace{\sqrt{d}}_{\mathbb{R}} \end{pmatrix}}_{\mathbb{R}^{s \times s}} \right) V$$



# Standard transformers | Problem : quadratic attention

- Let  $Q, K, V \in \mathbb{R}^{s \times d}$ ,  $s$  is sequence length

$$\underbrace{\text{Attention}(Q, K, V)}_{\mathbb{R}^{s \times d}} = \text{softmax} \left( \underbrace{\begin{pmatrix} \underbrace{QK^T}_{\mathbb{R}^{s \times s}} \\ \underbrace{\sqrt{d}}_{\mathbb{R}} \end{pmatrix}}_{\mathbb{R}^{s \times s}} \right) V$$

$\mathbb{R}^{s \times s}$ :

- **memory requirement** grows quadratically in sequence length  $s$

# Standard transformers | Problem : quadratic attention

Accelerators (e.g., CUDA GPU) have *limited memory capacity and bandwidth*:

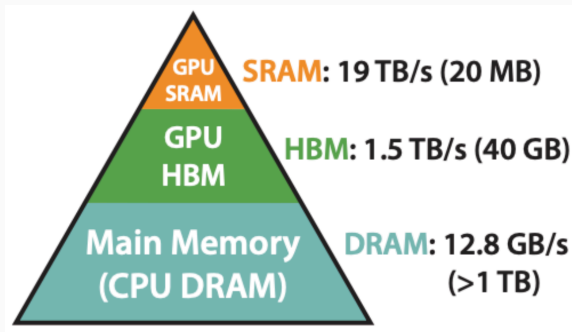


Figure 3: From *FlashAttention* [2]

# Standard transformers | Problem : quadratic attention

Accelerators (e.g., CUDA GPU) have *limited memory capacity and bandwidth*:

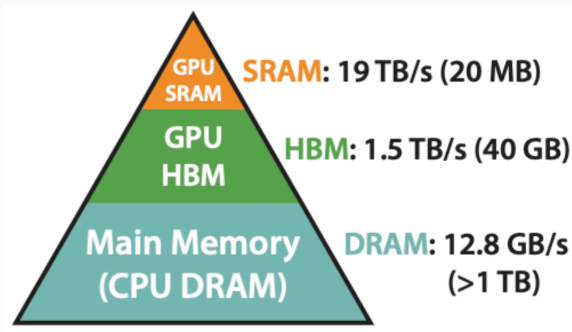


Figure 3: From *FlashAttention* [2]

$\mathbb{R}^{S \times S}$  memory requirement: many slow SRAM  $\leftrightarrow$  HBM transfers

# Standard transformers | Problem : quadratic attention

Accelerators (e.g., CUDA GPU) have *limited memory capacity and bandwidth*:

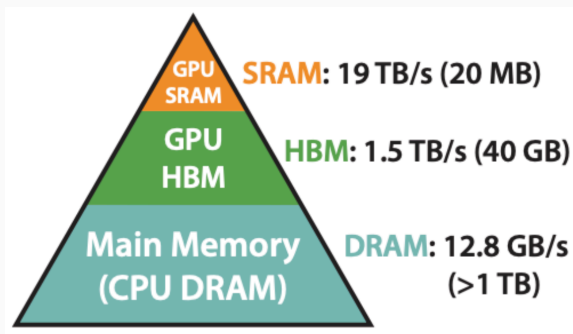


Figure 3: From *FlashAttention* [2]

$\mathbb{R}^{S \times S}$  memory requirement: many slow SRAM  $\leftrightarrow$  HBM transfers

- $\therefore$  memory requirement is the practical bottleneck!

Implications:

- Expensive to (pre-)train with a long context length.
  - Example: Code Llama: used 4,096 tokens during pre-training
- Expensive to generate (inference-time)

Implications:

- Expensive to (pre-)train with a long context length.
  - Example: Code Llama: used 4,096 tokens during pre-training
- Expensive to generate (inference-time)

Expensive can mean:

- Slow: bandwidth leads to transfers
- Infeasible: simply run out of memory
  - Example: Inference with batch size 1, hidden size 1024, 100M tokens  
⇒ 1000GB memory [8]

# Efficient attention

- Standard transformers → 2048 tokens
- **Efficient attention** → **millions of tokens**
  - Example: Long-context fine-tuning → 100k tokens

---

## SELF-ATTENTION DOES NOT NEED $O(n^2)$ MEMORY

---

A PREPRINT

**Markus N. Rabe and Charles Staats**  
Google Research  
{mrabe, cstaats}@google.com

**Memory-efficient attention** [Rabe & Staats arXiv 2021] [10]:

- Compute attention in chunks via clever softmax trick
- Avoids full  $\mathbb{R}^{S \times S}$  matrix: **better memory requirement!**



For a single query  $q$ , write as  $\text{Attention}(q, k_1, \dots, k_s, v_1, \dots, v_s)$

$$\text{Attention}(q) = \sum_{i=1}^s v_i w_i$$

$$\text{where } w_i = \frac{\exp(w'_i)}{\sum_{j=1}^s \exp(w'_j)},$$

$$w'_i = \text{dot}(q, k_i)$$

For a single query  $q$ , write as  $\text{Attention}(q, k_1, \dots, k_s, v_1, \dots, v_s)$

$$\text{Attention}(q) = \sum_{i=1}^s v_i w_i$$

$$\text{where } w_i = \frac{\exp(w'_i)}{\sum_{j=1}^s \exp(w'_j)},$$

$$w'_i = \text{dot}(q, k_i)$$

Requires storing  $s$  values of  $w'_i$ , implying  $O(s^2)$  for  $s$  queries.

For a single query  $q$ , write as  $\text{Attention}(q, k_1, \dots, k_s, v_1, \dots, v_s)$

$$\text{Attention}(q) = \sum_{i=1}^s v_i w_i$$

$$\text{where } w_i = \frac{\exp(w'_i)}{\sum_{j=1}^s \exp(w'_j)},$$

$$w'_i = \text{dot}(q, k_i)$$

Requires storing  $s$  values of  $w'_i$ , implying  $O(s^2)$  for  $s$  queries.

Insight 1:

$$\text{Attention}(q) \equiv \frac{\sum_{i=1}^s v_i \exp(w'_i)}{\sum_{j=1}^s \exp(w'_j)}$$

For a single query  $q$ , write as  $\text{Attention}(q, k_1, \dots, k_s, v_1, \dots, v_s)$

$$\text{Attention}(q) = \sum_{i=1}^s v_i w_i$$

$$\text{where } w_i = \frac{\exp(w'_i)}{\sum_{j=1}^s \exp(w'_j)},$$

$$w'_i = \text{dot}(q, k_i)$$

Requires storing  $s$  values of  $w'_i$ , implying  $O(s^2)$  for  $s$  queries.

Insight 1:

$$\text{Attention}(q) \equiv \frac{\sum_{i=1}^s v_i \exp(w'_i)}{\sum_{j=1}^s \exp(w'_j)}$$

Insight 2:

- Compute by iteratively adding to a vector  $v^* \in \mathbb{R}^d$  and  $w^* \in \mathbb{R}$ .  
Less memory!

Specifically, compute softmax normalization online [9]:

For each query  $q$ :

For  $i = 1, \dots, s$ :<sup>1</sup>

$$w'_i = \text{dot}(q, k_i)$$

$$v^* \leftarrow v^* + v_i \exp(w'_i)$$

$$s^* \leftarrow s^* + \exp(w'_i)$$

At the end,  $\text{Attention}(q) = v^*/s^*$ .

---

<sup>1</sup>In general, segment into “chunks” (aka “blocks” / “tiles”). See paper for full version.

Sequence length	$n = 2^8$	$2^{10}$	$2^{12}$	$2^{14}$	$2^{16}$	$2^{18}$	$2^{20}$
Size of inputs and outputs	160KB	640KB	2.5MB	10MB	40MB	160MB	640MB
Memory overhead of <b>standard attention</b>	270KB	4.0MB	64MB	<b>1GB</b>	<b>OOM</b>	<b>OOM</b>	<b>OOM</b>
Memory overhead of <b>memory-eff. attn.</b>	270KB	4.0MB	16MB	<b>17MB</b>	<b>21MB</b>	<b>64MB</b>	<b>256MB</b>
Compute time on TPUv3	0.06ms	0.11ms	0.7ms	11.3ms	177ms	2.82s	45.2s
Relative compute speed	$\pm 5\%$	$\pm 5\%$	$-8 \pm 2\%$	<b><math>-13 \pm 2\%</math></b>	-	-	-


~ 1 million tokens 

Table 2: Memory and time requirements of self-attention during **inference**.

Figure 4: Inference benchmarking

Sequence length	$n = 2^8$	$2^{10}$	$2^{12}$	$2^{14}$	$2^{16}$	$2^{18}$	$2^{20}$
Size of inputs and outputs	192KB	768KB	2.9MB	12MB	47MB	188MB	750MB
Memory overhead of standard attention	532KB	8.0MB	128MB	2.0GB	OOM	OOM	OOM
Memory overhead of memory-eff. attn.	532KB	8.0MB	41MB	64MB	257MB	1.0GB	4.0GB
Compute time on TPUv3	0.1ms	0.18ms	1.4ms	21ms	336ms	5.3s	85s
Relative compute speed	$\pm 5\%$	$\pm 5\%$	$-30 \pm 5\%$	$-35 \pm 5\%$	-	-	-

Table 3: Memory and time requirements of self-attention during **differentiation**. Note that the slowdown in compute speed is expected due to the use of checkpointing in memory-efficient attention.

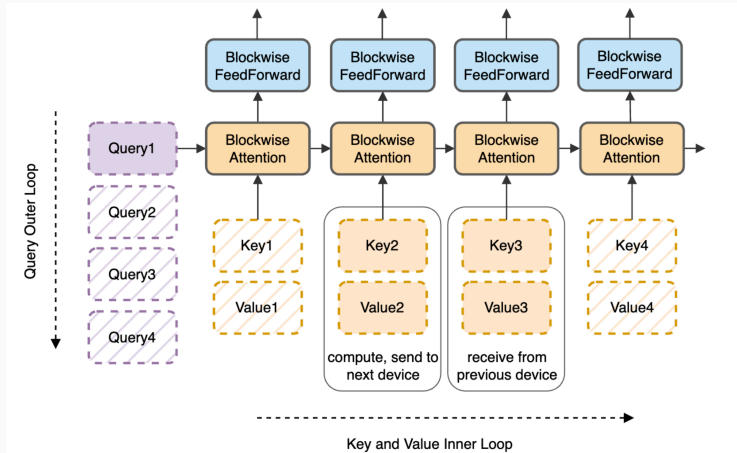
Figure 5: Differentiation (proxy for training) benchmarking

Memory-efficient attention [Rabe & Staats 2021 [10]]:

- *FlashAttention*: Analogy via new CUDA kernel [Dao et al 2022 [2]]
- *Blockwise Transformers*: also chunk feedforward layer [Liu et al 2023 [7]]
- **Ring Attention**: Distribute chunks across devices [Liu et al 2023 [8]]



# Efficient attention | Ring Attention [Liu, Zaharia, Abbeel arXiv 2023] [8]



**Figure 6:** From [8]: each host holds one query block, and key-value blocks traverse through a ring of hosts in a block-by-block fashion.

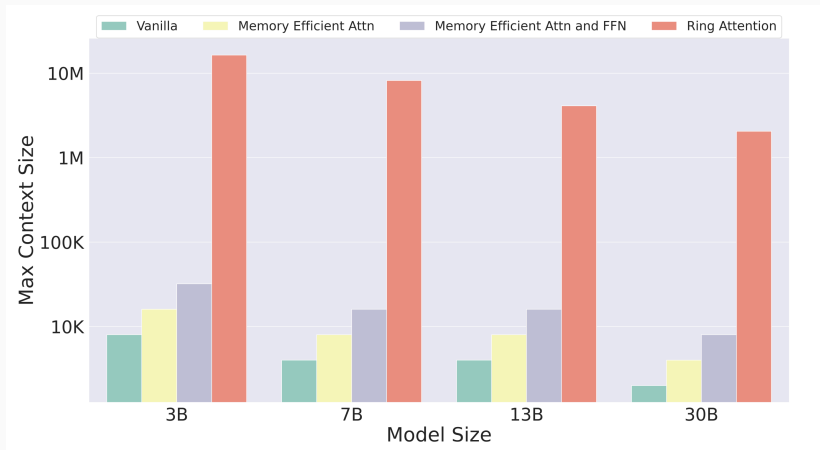


Figure 7: From [8]. Maximum *training* context size on TPUv4-1024 with vanilla, efficient attention variants, and Ring Attention.

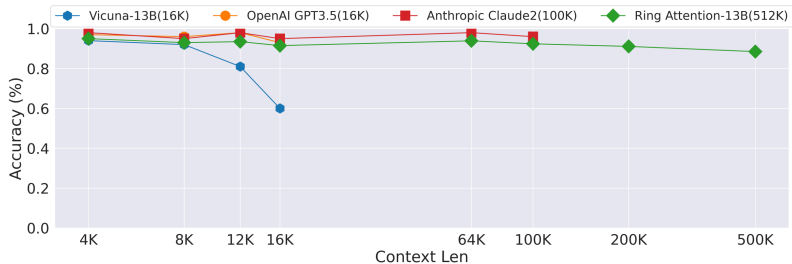


Figure 3: Comparison of different models on the long-range line retrieval task.

Figure 8: Long-range line retrieval task

# Long-context generation

- Standard transformers → 2048 tokens
- Efficient attention → millions of tokens
  - **Example adaptation recipe:** → 100k tokens

## *Effective Long-Context Scaling of Foundation Models (Meta 2023)*

- Llama / analysis
- + ideas used in CodeLlama (Meta 2023)

# Long context fine-tuning of [code-]llama

Pretrain with short context (e.g. 4,096), adapt to long contexts

- Positional embeddings
- Data

# Long context fine-tuning | positional embeddings

RoPe embeddings map inputs to a sphere, with a periodic structure

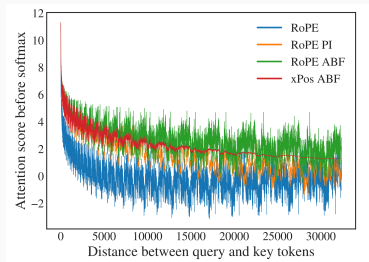


Figure 9: Vanilla RoPE naturally decays as  $s$  increases

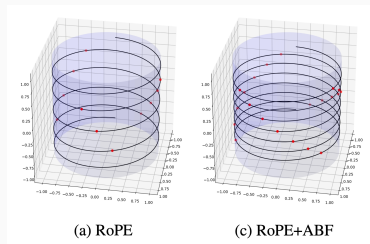


Figure 10: Solution: increase frequency (RoPE ABF)

# Long context fine-tuning Llama | positional embeddings

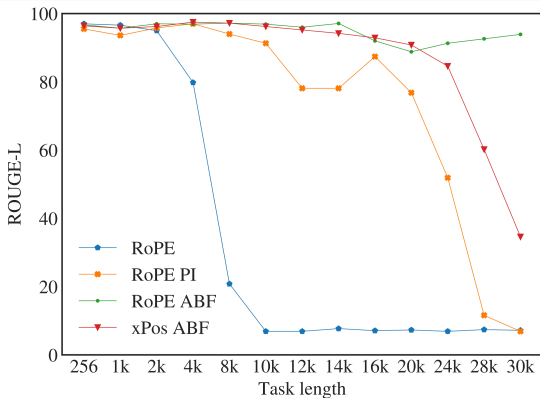


Figure 11: Fixes degradation on diagnostic task



Data:

- Continue pretraining on special mix of 400B tokens
- Finetune on long instruction-tuning data

Special mix: “new long text data” and up-weight long data samples

<b>Continual Pretrain Data</b>	NarrativeQA $\Delta$ F1	Qasper $\Delta$ F1	Quality $\Delta$ EM	QMSum $\Delta$ ROUGE-geo
LLAMA 2 LONG data mix	<b>23.70%</b>	<b>43.64%</b>	<b>75.5%</b>	<b>45.70%</b>
LLAMA 2 data mix	18.23%	38.12%	60.3%	44.87%
- remove long text data	19.48%	39.14%	67.1%	36.60%
- upsample existing long text data	22.15%	36.82%	65.0%	42.83%

Table 7: Comparison of different pretraining data mix on long-context tasks. Instead of showing the absolute performance, we report relative improvements over the 7B LLAMA 2 which has a 4,096-token context window. All models are evaluated with prompts truncated at 16,384 tokens.

Synthetic long + concatenated short instruction data:

Settings	Qasper	NarrativeQA	QuALITY	SummScreenFD	QMSum
LLAMA 2 CHAT baseline	12.2	9.13	56.7	10.5	14.4
LLAMA 2 LONG <i>finetuned</i> with:					
"RLHF V5"	22.3	13.2	71.4	14.8	16.9
"RLHF V5" mix pretrain	23.7	16.6	76.2	<b>15.7</b>	17.8
"RLHF V5" mix self-inst w/o LM loss	35.7	22.3	59.3	12.2	13.4
"RLHF V5" mix self-inst with LM loss	<b>38.9</b>	<b>23.3</b>	<b>77.3</b>	14.5	<b>18.5</b>

Table 9: Comparison of different instruction finetuning data mixes.

# Long context fine-tuning CodeLlama

## Mixed result for CodeLlama:

**Synthetic Key Retrieval Task.** We prompt the model with a variable number of tokens by concatenating Python solutions from the CodeContest dataset (Li et al., 2022), which results in syntactically valid source code. At a specified relative position within the prompt, we insert the following key, where <VALUE> is a two-digit number that is randomly sampled based on the overall number of tokens in the prompt:

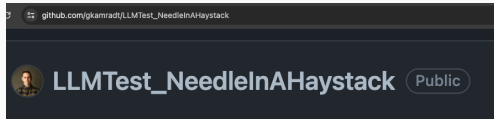
```
def my_function() -> int:
    """Note that this function is used at the end
    """
    return <VALUE>
```

Model	Size	Context Length / Key Position								
		8,000			16,000			24,000		
		0	0.2	0.4	0	0.2	0.4	0	0.2	0.4
CODE LLAMA	7B	100.0	95.3	100.0	54.7	100.0	98.4	3.1	85.9	85.9
CODE LLAMA	13B	100.0	100.0	100.0	100.0	100.0	100.0	100.0	89.1	6.3
CODE LLAMA	34B	76.6	100.0	100.0	95.3	96.9	100.0	81.3	0.0	81.3
CODE LLAMA - INSTRUCT	7B	100.0	97.7	100.0	7.0	96.9	96.1	0.0	62.5	54.7
CODE LLAMA - INSTRUCT	13B	100.0	100.0	100.0	100.0	100.0	93.8	4.7	84.4	100.0
CODE LLAMA - INSTRUCT	34B	92.2	100.0	100.0	68.8	95.3	100.0	46.9	0.0	85.9
gpt-3.5-turbo-16k-0630	-	100.0	100.0	95.3	95.3	90.6	98.4	-	-	-

Table 17: Function Key Retrieval Accuracy (%) for Code Llama models.

Long sequence modeling and generation is an active research area:

- Alternative models: state-space models, ...
- **Proprietary approaches**
- ...



## The Test

1. Place a random fact or statement (the 'needle') in the middle of a long context window (the 'haystack')
2. Ask the model to retrieve this statement
3. Iterate over various document depths (where the needle is placed) and context lengths to measure performance

Figure 12: [https://github.com/gkamradt/LLMTest\\_NeedleInAHaystack](https://github.com/gkamradt/LLMTest_NeedleInAHaystack)

# Black box models

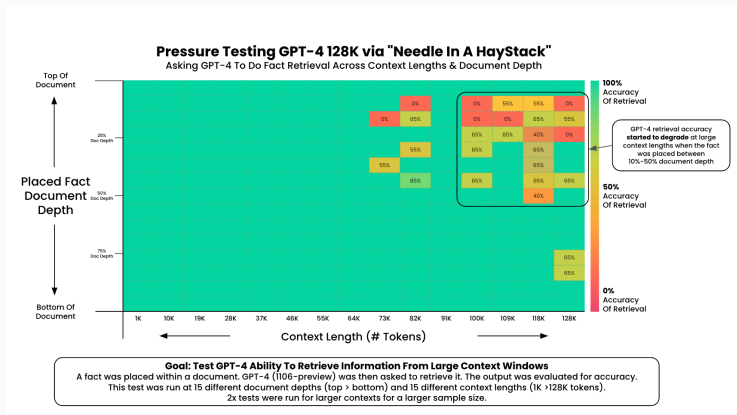


Figure 13: OpenAI's GPT-4-128K (Run 11/8/2023); from [https://github.com/gkamradt/LLMTest\\_NeedleInAHaystack](https://github.com/gkamradt/LLMTest_NeedleInAHaystack)

# Black box models

## Claude 3 Opus

### Recall accuracy over 200K

(averaged over many diverse document sources and 'needle' sentences)

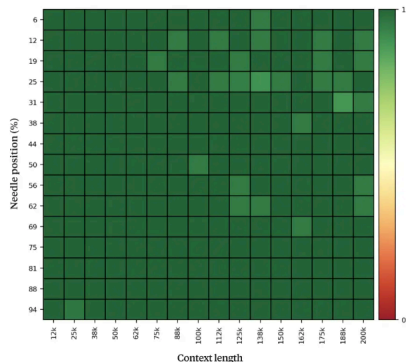


Figure 14: Claude 3 (March 4 2024); from <https://www.anthropic.com/news/claude-3-family>



# Black box models

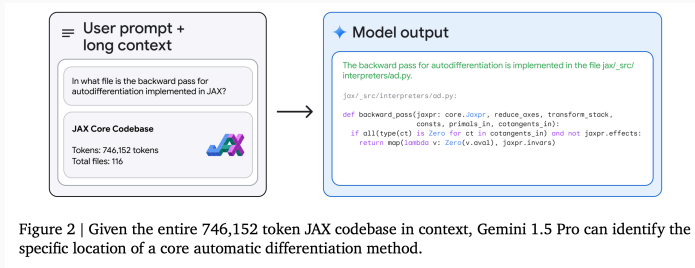


Figure 15: Gemini 1.5 (Feb 2024)

# Black box models

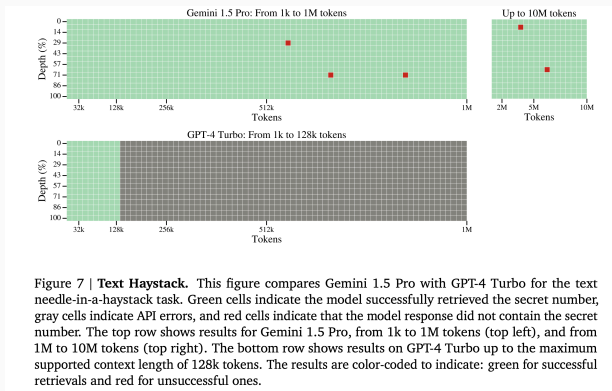


Figure 16: Gemini 1.5 (Feb 2024)

# Black box models

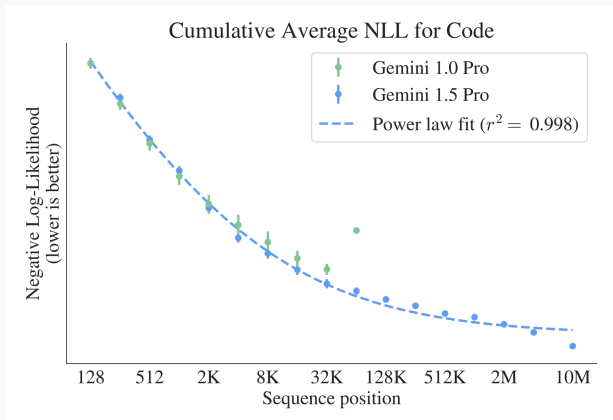


Figure 17: Gemini 1.5 (Feb 2024)

$$p(y|x, C)$$

- Approach 1: place a long context  $C$  in the input of a LM
  - Use techniques to extend the context size of transformers

*Observation:* how do we select a context?

## II. Retrieval-augmented generation

---

# Retrieval-augmented generation

- $x$ : input (e.g. query),  $y$ : output (e.g. code)
- $\mathcal{Z}$ : data store of items (e.g. code snippets)

$$\begin{aligned} p(y|x) &= \sum_{Z_k \subset \mathcal{Z}} p(y|x, Z_k) p(Z_k|x) \\ &\approx \arg \max_{Z_k} p(y|x, Z_k) p(Z_k|x) \\ &\approx p(y|x, \hat{Z}_k) p(\hat{Z}_k|x), \end{aligned}$$

# Retrieval-augmented generation

- $x$ : input (e.g. query),  $y$ : output (e.g. code)
- $\mathcal{Z}$ : data store of items (e.g. code snippets)

$$\begin{aligned} p(y|x) &= \sum_{Z_k \subset \mathcal{Z}} p(y|x, Z_k) p(Z_k|x) \\ &\approx \arg \max_{Z_k} p(y|x, Z_k) p(Z_k|x) \\ &\approx p(y|x, \hat{Z}_k) p(\hat{Z}_k|x), \end{aligned}$$

where

$$\hat{Z}_k = \text{topk}_{z_i \in \mathcal{Z}} s_\phi(x, z_i)$$

is a set of *retrieved* items using scorer  $s_\phi(x, z)$ .

# Retrieval-augmented generation

Concretely:

1. Specify a **datastore**  $\mathcal{Z}$  of items to retrieve  $z$
2. **Retriever**: map a query  $x$  to ranked list  $z_1, z_2, z_3, \dots$
3. **Generator**: use top-ranked items,  $p_\theta(y|x, \{z_k\}_{k=1}^K)$



# Retrieval-augmented generation

Concretely:

1. Specify a **datastore**  $\mathcal{Z}$  of items to retrieve  $z$
2. **Retriever**: map a query  $x$  to ranked list  $z_1, z_2, z_3, \dots$
3. **Generator**: use top-ranked items,  $p_\theta(y|x, \{z_k\}_{k=1}^K)$

Key feature:  $\mathcal{Z}$  can have new code, relevant documents, etc...

Example retriever: *BM-25*

- Score text based on term (n-gram) frequencies and length
- Doesn't use a neural model

# Retrieval-augmented generation | retrievers

First thing to do is create an instance of the BM25 class, which reads in a corpus of text and does some indexing on it:

```
from rank_bm25 import BM25Okapi

corpus = [
    "Hello there good man!",
    "It is quite windy in London",
    "How is the weather today?"
]

tokenized_corpus = [doc.split(" ") for doc in corpus]

bm25 = BM25Okapi(tokenized_corpus)
# <rank_bm25.BM25Okapi at 0x1047881d0>
```

Now that we've created our document indexes, we can give it queries and see which documents are the most relevant:

```
query = "windy London"
tokenized_query = query.split(" ")

doc_scores = bm25.get_scores(tokenized_query)
# array([0.          , 0.93729472, 0.          ])
```

Instead of getting the document scores, you can also just retrieve the best documents with

```
bm25.get_top_n(tokenized_query, corpus, n=1)
# ['It is quite windy in London']
```

Figure 18: Usage example from [https://github.com/dorianbrown/rank\\_bm25](https://github.com/dorianbrown/rank_bm25)

Example retriever: *dense passage retriever* (DPR) [6]

- Encoder

$$s_{\phi}(x) \rightarrow \mathbb{R}^d$$

- Score

$$s(x, z_i) = s_{\phi}(x)^{\top} s_{\phi}(z_i) \in \mathbb{R}$$

- Retrieval

$$\mathit{top-K}(s(x, z_1), s(x, z_2), \dots, s(x, z_{|Z|}))$$

Example retriever: *dense passage retriever* (DPR) [6]

- Encoder

$$s_{\phi}(x) \rightarrow \mathbb{R}^d$$

- Score

$$s(x, z_i) = s_{\phi}(x)^{\top} s_{\phi}(z_i) \in \mathbb{R}$$

- Retrieval

$$\mathit{top}\text{-}K(s(x, z_1), s(x, z_2), \dots, s(x, z_{|Z|}))$$

In practice, precompute all  $s(z_i)$ 's and use a fast maximum inner-product search (MIPS) library (e.g., *faiss* (<https://github.com/facebookresearch/faiss>)).

Example retriever: *dense passage retriever* (DPR)

- Neural encoder  $s_\phi(x) \rightarrow \mathbb{R}^d$  (e.g., BERT)
- Train with *InfoNCE loss* [11]

$$L(x, z^+, Z^-) = -\log \frac{\exp(s(x, z^+))}{\exp(s(x, z^+)) + \sum_{z^- \in Z^-} \exp(s(x, z^-))}$$

As negatives, use other examples in the current minibatch [4, 3]

# Retrieval-augmented generation | retrievers

Task-aware retrieval with instructions (TART) [1]:

- Single task retriever → instruction-tuned retriever
- Condition on a task-relevant instruction:  $s_\phi : (x, l) \rightarrow z_e$

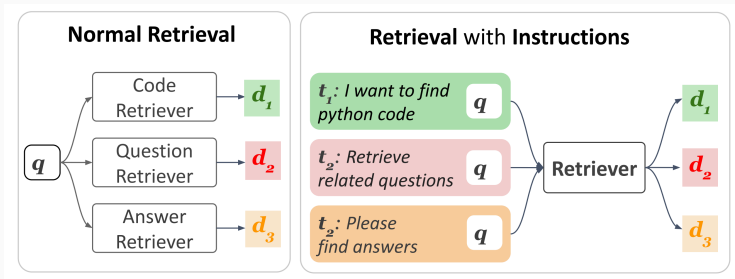
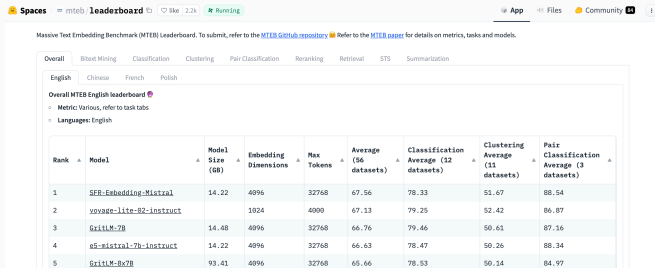


Figure 19: Task Aware Retrieval with Instructions, Asai et al. ACL 2023

# Retrieval-augmented generation | retrievers

LLM-based retrieval with instructions:

- Adapt LLM (e.g. Mistral 7b) to produce embeddings  $z_e \in \mathbb{R}^d$
- Condition on a task-relevant instruction  $s_\phi : (x, I) \rightarrow z_e$



Massive Text Embedding Benchmark (MTEB) Leaderboard. To submit, refer to the [MTEB GitHub repository](#). Refer to the [MTEB paper](#) for details on metrics, tasks and models.

Overall | [Bbox Mining](#) | [Classification](#) | [Clustering](#) | [Pair Classification](#) | [Reranking](#) | [Retrieval](#) | [STS](#) | [Summarization](#)

English | Chinese | French | Polish

Overall MTEB English leaderboard

- Metric: Various, refer to task tabs
- Languages: English

Rank	Model	Model Size (GB)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)	Pair Classification Average (3 datasets)
1	SFR-Embedding-Mistral	14.22	4096	32768	67.56	78.33	51.67	88.54
2	voyage-lite-82-instruct		1024	4000	67.13	79.25	52.42	86.87
3	GritLM-7B	14.48	4096	32768	66.76	79.46	50.61	87.16
4	e5-mistral-7b-instruct	14.22	4096	32768	66.63	78.47	50.26	88.34
5	GritLM-8x7B	93.41	4096	32768	65.66	78.53	50.14	84.97



# Retrieval-augmented generation | retrievers

Example: Retrieval with embeddings from SFR-Embedding-Mistral

```
from datasets import load_dataset
from sentence_transformers import SentenceTransformer, util

model = SentenceTransformer("Salesforce/SFR-Embedding-Mistral", device="cuda")

ds = load_dataset("openai_humaneval")
datastore = [x['canonical_solution'] for x in ds['test']]
prompts = [x['prompt'] for x in ds['test']]

query = """Instruct: Given a search query, return relevant source code.
Query: code involving Fibonacci numbers"""

embeddings = model.encode([query] + datastore)
scores = util.cos_sim(embeddings[:1], embeddings[1:]) * 100

k = 3

for query_i, scores_i in zip(queries, scores):
    retrieved = [prompts[i] + datastore[i] for i in scores_i.argsort(descending=True)[:k]]
    for item in retrieved:
        print('===== QUERY: ' + query_i)
        print(item)
        print('-----\n')
```

# Retrieval-augmented generation | retrievers

## Example: Retrieval with embeddings from SFR-Embedding-Mistral

```
===== QUERY: Instruct: Given a search query, return relevant source code.  
Query: code involving Fibonacci numbers  
  
def fib(n: int):  
    """Return n-th Fibonacci number.  
    >>> fib(10)  
    55  
    >>> fib(1)  
    1  
    >>> fib(8)  
    21  
    """  
    if n == 0:  
        return 0  
    if n == 1:  
        return 1  
    return fib(n - 1) + fib(n - 2)  
  
-----  
  
===== QUERY: Instruct: Given a search query, return relevant source code.  
Query: code involving Fibonacci numbers  
  
def prime_fib(n: int):  
    """  
    prime_fib returns n-th number that is a Fibonacci number and it's also prime.  
    >>> prime_fib(1)  
    2  
    >>> prime_fib(2)  
    3  
    >>> prime_fib(3)  
    5  
    >>> prime_fib(4)  
    13  
    >>> prime_fib(5)  
    89  
    """  
    import math  
  
    def is_prime(p):  
        if p < 2:  
            return False  
        for k in range(2, min(int(math.sqrt(p)) + 1, p - 1)):  
            if p % k == 0:  
                return False  
        return True  
    f = [0, 1]  
    while True:
```

# Retrieval-augmented generation for code | DocPrompting

## Case study: retrieving *documentation* for code generation

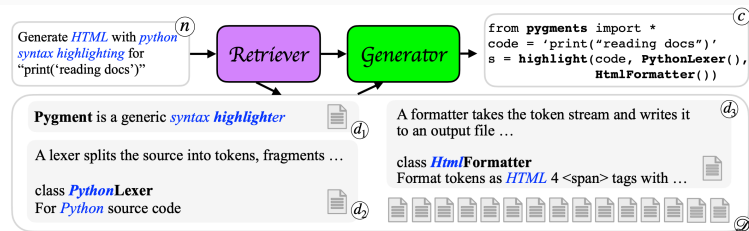


Figure 1: DocPrompting: given an NL intent  $\mathcal{N}$ , the retriever retrieves a set of relevant documentation  $\{d_1, d_2, d_3\}$  from a documentation pool  $\mathcal{D}$ . Then, the generator generates the code  $\mathcal{C}$  based on the NL and retrieved docs. DocPrompting allows the model to generalize to previously unseen usages by reading those docs. *Italic blue* highlights the shared tokens between NL and docs; **Bold** shows shared tokens between docs and the code snippet.

Figure 20: DocPrompting: Generating Code by Retrieving the Docs, Zhou et al. ICLR 2023

- Retriever:
  - Dense (CodeT5 encoder, InfoNCE loss + in-batch negatives)
- Generator:
  - Finetune CodeT5 with fusion-in-decoder [5]
    - Attends over  $k$  encoded retrieved docs

- Retriever:
  - Dense (CodeT5 encoder, InfoNCE loss + in-batch negatives)
- Generator:
  - Finetune CodeT5 with fusion-in-decoder [5]
    - Attends over  $k$  encoded retrieved docs
- Task:
  - StackOverflow questions → Python (CoNaLa dataset)
  - Split data so that test problems have **at least 1 unseen function**
- Datastore:
  - Docs for all Python functions on *devdocs.io*

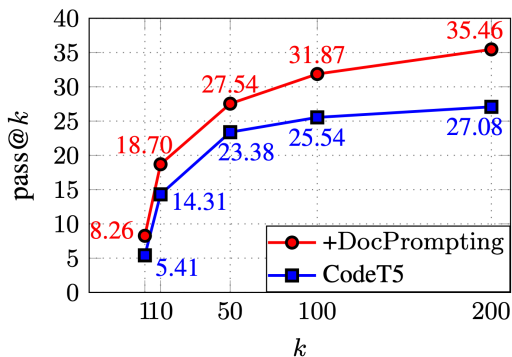


Figure 3: Pass@ $k$  of CodeT5 with and without DocPrompting on 100 CoNaLa examples.

Figure 21: Pass@ $k$  on CoNaLa

Table 5: Examples of predictions from CoNaLa, of the base CodeT5 compared to CodeT5+DocPrompting. Unseen functions are underscored.

<b>NL Intent: Open image "picture.jpg"</b>	
Ground truth:	<code>img = <u>Image.open</u>('picture.jpg') \n Img.show</code>
CodeT5:	<code>os.open('picture.jpg', 'r')</code>
CodeT5+DocPrompting:	<code>image = <u>Image.open</u>('picture.jpg', 'rb')</code>
<b>NL Intent: Exclude column names when writing dataframe 'df' to a csv file 'filename.csv'</b>	
Ground truth:	<code><u>df.to_csv</u>('filename.csv', header=False)</code>
CodeT5:	<code>df.drop(['col1', 'col2'], axis=1, inplace=True)</code>
CodeT5+DocPrompting:	<code><u>df.to_csv</u>('filename.csv', skiprows=1)</code>

Figure 22: Qualitative example

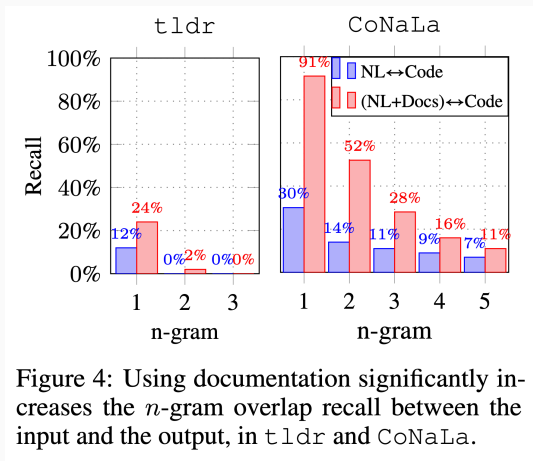


Figure 23: Analysis: why does reading docs help generate code (with these models)?



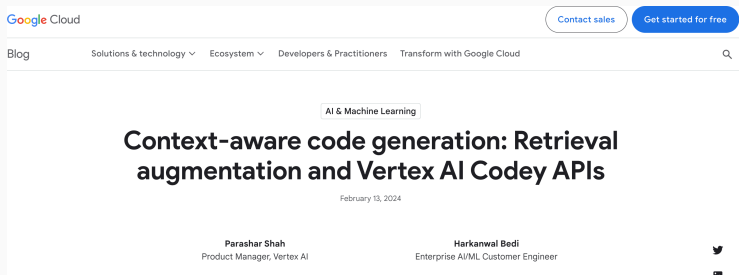


Figure 24: RAG with Blackbox API models (Feb 2024) <sup>a</sup>

<sup>a</sup><https://cloud.google.com/blog/products/ai-machine-learning/context-aware-code-generation-rag-and-vertex-ai-codey-apis>

**Prompt:** "Create python function that takes a prompt and predicts using langchain.llms interface for VertexAI text-bison model"

**Output without RAG:**

```
[86] response = code_llm.predict(text=user_question, max_output_tokens=2048, temperature=0.1)
      print(response)

```python
import langchain.llms as llms

def predict(prompt):
    # Load the text-bison model
    model = llms.TextBisonModel()

    # Generate predictions
    predictions = model.predict(prompt, num_return_sequences=3)

    # Return the predictions
    return predictions
```
```

*Figure 2: Output from the model without any external context*

**Figure 25: RAG with Blackbox API models (Feb 2024)<sup>a</sup>**

<sup>a</sup><https://cloud.google.com/blog/products/ai-machine-learning/context-aware-code-generation-rag-and-vertex-ai-codey-apis>

Output with RAG:

```
[87] results = qa_chain({"query": user_question})
      print(results["result"])

```python
def predict_with_langchain_llms(prompt):
    """Predicts using langchain.llms interface for VertexAI text-bison model.

    Args:
        prompt: The prompt to predict.

    Returns:
        The prediction.
    """

    # Initialize the VertexAI LLMs client.
    llm = VertexAI(
        model_name="text-bison-32k",
        max_output_tokens=256,
        temperature=0.1,
        top_p=0.8,
        top_k=40,
        verbose=True,
    )

    # Predict the prompt.
    prediction = llm.predict(prompt)

    return prediction
```
```

Figure 26: RAG with Blackbox API models (Feb 2024) <sup>a</sup>

$$p(y|x, C)$$

- Approach 1: place a long context  $C$  in the input of a LM
  - Use techniques to expand the context size of transformers
- Approach 2: retrieve relevant contexts, place in the input of a LM
  - Define a datastore and retriever



A. Asai, T. Schick, P. Lewis, X. Chen, G. Izacard, S. Riedel, H. Hajishirzi, and W.-t. Yih.

**Task-aware retrieval with instructions.**




In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3650–3675, Toronto, Canada, July 2023. Association for Computational Linguistics.



T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Re.

**Flashattention: Fast and memory-efficient exact attention with IO-awareness.**

In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022.

-  D. Gillick, S. Kulkarni, L. Lansing, A. Presta, J. Baldridge, E. Ie, and D. Garcia-Olano.  
**Learning dense representations for entity retrieval, 2019.**
-  M. Henderson, R. Al-Rfou, B. Strope, Y.-H. Sung, L. Lukács, R. Guo, S. Kumar, B. Miklos, and R. Kurzweil.  
**Efficient natural language response suggestion for smart reply.**  
*ArXiv, abs/1705.00652, 2017.*
-  G. Izacard and E. Grave.  
**Leveraging passage retrieval with generative models for open domain question answering, 2020.**



V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W.-t. Yih.

**Dense passage retrieval for open-domain question answering.**

In B. Webber, T. Cohn, Y. He, and Y. Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online, Nov. 2020.

Association for Computational Linguistics.



H. Liu and P. Abbeel.

**Blockwise parallel transformer for large context models, 2023.**



H. Liu, M. Zaharia, and P. Abbeel.

**Ring attention with blockwise transformers for near-infinite context, 2023.**



M. Milakov and N. Gimelshein.

**Online normalizer calculation for softmax.**

*ArXiv*, abs/1805.02867, 2018.



M. N. Rabe and C. Staats.

**Self-attention does not need  $\mathcal{O}(n^2)$  memory.**

2021.



A. van den Oord, Y. Li, and O. Vinyals.

**Representation learning with contrastive predictive coding.**

*ArXiv*, abs/1807.03748, 2018.