# Generative Representational Instruction Tuning
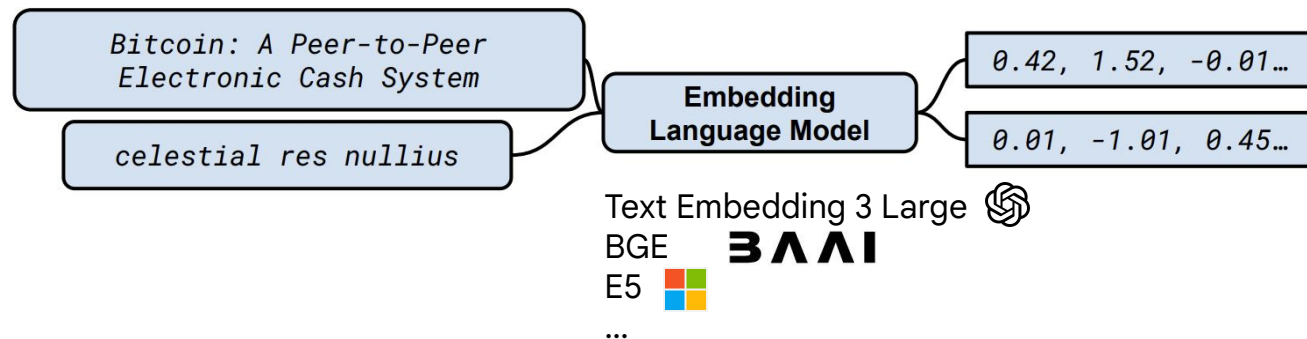
**Niklas Muennighoff (**Twitter: @Muennighoff**)**

Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, Douwe Kiela

This talk is being recorded.

# Two types of language models

*Bitcoin: A Peer-to-Peer Electronic Cash System*

*celestial res nullius*

**Embedding Language Model**

`0.42, 1.52, -0.01…`

`0.01, -1.01, 0.45…`

Clustering, Semantic Search, Text Mining…

Text Embedding 3 Large

BGE **ƎΛΛI**

E5

…

Please write me a blog post about my recent hike of Mt. Fuji at midnight.

You have two ropes, each takes exactly 1 hour to burn. How would you use them to time exactly 15 minutes? The ropes are of uneven densities, so half the rope does not necessarily take half the time.

**Generative Language Model**

GPT-4

Gemini

Llama

…

*Sure, here is the blog post.*
*It was August the 10th when I arrived at Lake Kawaguchi from…*

*You start by…*

Story Generation, Question Answering, Chat…

# Advantages of combining them

**Performance:**
Get better on both?

**Efficiency:**
Speed-up joint use cases

**Simplicity:**
Unify endpoints

**Embedding benchmarks: MTEB..**

**Generation benchmarks: AlpacaEval..**

**Traditional RAG**

To slow down your speed of aging, you can…

**Generative Model**

How to prevent aging?

Technological and lifestyle factors may influence an individual's longevity. Cellular reprogramming…

**Index**
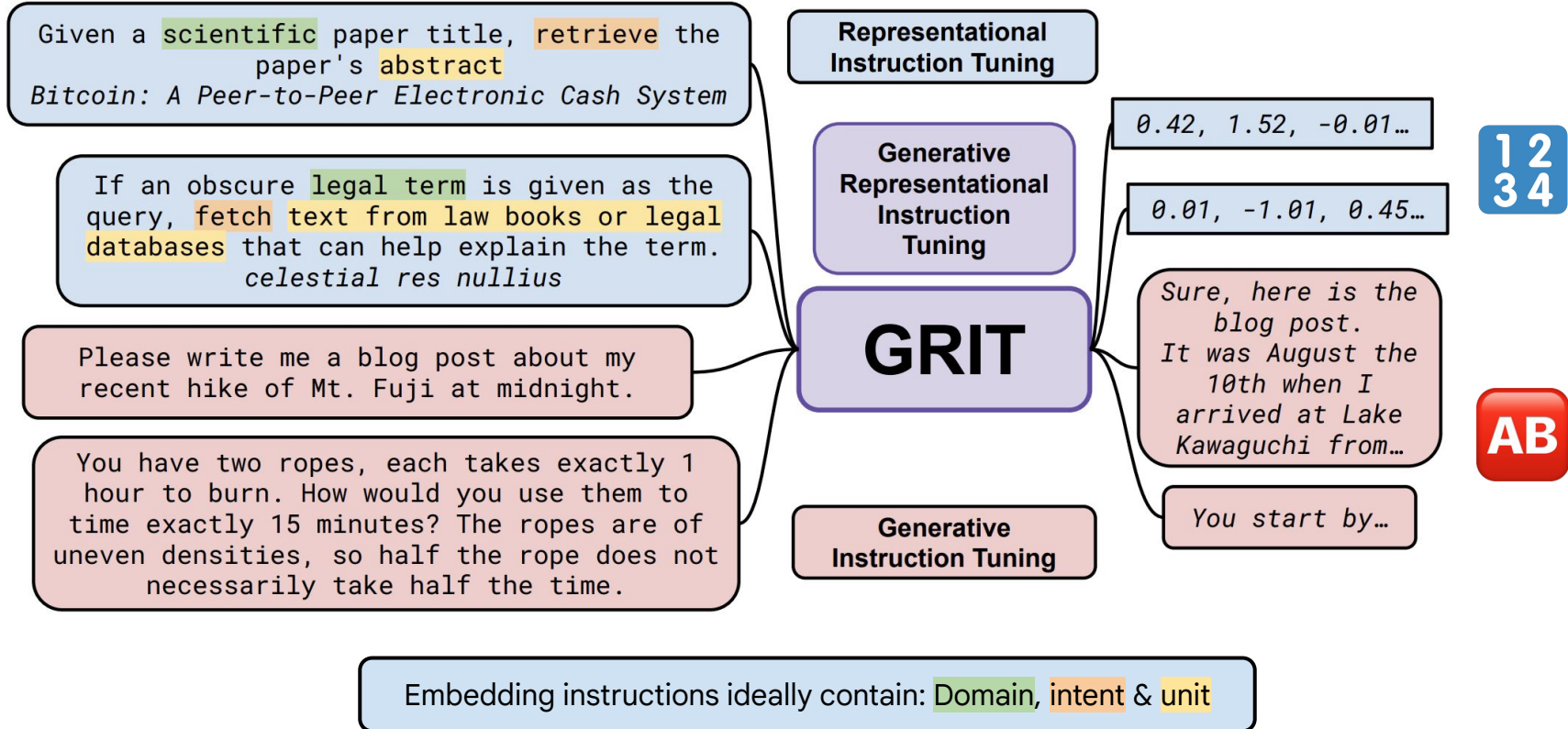
**Embedding Model**

How to prevent aging?

**Embedding endpoint**

```
1  curl https://api.openai.com/v1/embeddings \
2    -H "Content-Type: application/json" \
3    -H "Authorization: Bearer $OPENAI_API_KEY" \
4    -d '{
5      "input": "Your text string goes here",
6      "model": "text-embedding-3-small"
7    }'
```

**Generation endpoint**

```
1  curl https://api.openai.com/v1/chat/completions \
2    -H "Content-Type: application/json" \
3    -H "Authorization: Bearer $OPENAI_API_KEY" \
4    -d '{
5      "model": "gpt-3.5-turbo",
6      "messages": [
7        {
8          "role": "system",
9          "content": "You are a helpful assistant."
10        },
11        {
```

MTEB
**Massive Text Embedding Benchmark**

AlpacaFarm

# Unifying representation & generation

Given a scientific paper title, retrieve the paper's abstract
*Bitcoin: A Peer-to-Peer Electronic Cash System*

If an obscure legal term is given as the query, fetch text from law books or legal databases that can help explain the term.
*celestial res nullius*

Please write me a blog post about my recent hike of Mt. Fuji at midnight.

You have two ropes, each takes exactly 1 hour to burn. How would you use them to time exactly 15 minutes? The ropes are of uneven densities, so half the rope does not necessarily take half the time.

**Representational Instruction Tuning**

**Generative Representational Instruction Tuning**

**GRIT**

**Generative Instruction Tuning**

0.42, 1.52, -0.01…

0.01, -1.01, 0.45…

*Sure, here is the blog post.*
*It was August the 10th when I arrived at Lake Kawaguchi from…*

*You start by…*

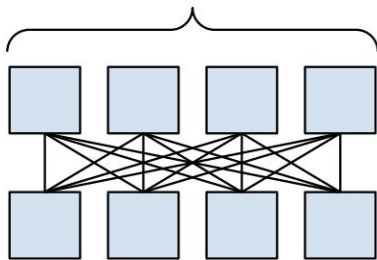Embedding instructions ideally contain: Domain, intent & unit

# Differing instructions, format & attention

**1234** **Representation**

Mean Pooling
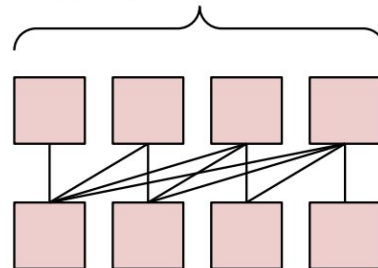
```
<s><|user|>
{instruction}
<|embed|>
{sample to represent}
```

**AB** **Generation**

Language Modeling Head

```
<s><|user|>
{instruction}
<|assistant|>
{response}</s>
<|user|>
…
```

# Combining losses

**Contrastive embedding loss with hard negatives**

$$\mathcal{L}_{\text{Rep}} = -\frac{1}{M} \sum_{i=1}^{M} \log \frac{\exp(\tau \cdot \sigma(f_\theta(q^{(i)}), f_\theta(d^{(i)})))}{\sum_{j=1}^{M} \exp(\tau \cdot \sigma(f_\theta(q^{(i)}), f_\theta(d^{(j)})))}$$

**Next token prediction loss for generation**

$$\mathcal{L}_{\text{Gen}} = -\frac{1}{N} \sum_{i=1}^{N} \log P(f_{\theta,\eta}(x^{(i)}) | f_{\theta,\eta}(x^{(<i)}))$$

**Combining both with adjustable weights**

$$\mathcal{L}_{\text{GRIT}} = \lambda_{\text{Rep}} \mathcal{L}_{\text{Rep}} + \lambda_{\text{Gen}} \mathcal{L}_{\text{Gen}}$$

# GritLM-7B & GritLM-8x7B

**Embedding Data**

**E5S:** ELI5, GPT4 Synthetic, MSMARCO, S2ORC, SQuAD...

**Generative Data**

**Tülu-2:** FLAN, Open Assistant, ShareGPT, LIMA, Open-Orca…

**Finetune Mistral 7B with GRIT**

**GritLM-7B**

**… Mixtral 8x7B …**

**GritLM-8x7B**

**… GRIT should work with any LM …**

# GritLM-7B & GritLM-8x7B

**Embedding Performance** 🔢

**Generative Performance** 🆎

| Task (→) | CLF Acc. | Clust. V-Meas. | PairCLF AP | Rerank MAP | Retrieval nDCG | STS Spear. | Summ. Spear. | Avg. |
|---|---|---|---|---|---|---|---|---|
| Metric (→) Dataset # (→) | 12 | 11 | 3 | 4 | 15 | 10 | 1 | 56 |
| *Proprietary models* ♥ | | | | | | | | |
| OpenAI v3 | 75.5 | 49.0 | 85.7 | 59.2 | 55.4 | 81.7 | 29.9 | 64.6 |
| *Other Open Models* ♥ | | | | | | | | |
| Llama 2 70B | 60.4 | 29.0 | 47.1 | 38.5 | 9.0 | 49.1 | 26.1 | 35.6 |
| Mistral 7B | 63.5 | 34.6 | 53.5 | 43.2 | 13.2 | 57.4 | 19.7 | 40.5 |
| Mistral 7B Instruct | 67.1 | 34.6 | 59.6 | 44.8 | 16.3 | 63.4 | 25.9 | 43.7 |
| GPT-J 6B | 66.2 | 39.0 | 60.6 | 48.9 | 19.8 | 60.9 | 26.3 | 45.2 |
| SGPT BE 5.8B | 68.1 | 40.3 | 82.0 | 56.6 | 50.3 | 78.1 | 31.5 | 58.9 |
| Instructor XL 1.5B | 73.1 | 44.7 | 86.6 | 57.3 | 49.3 | 83.1 | **32.3** | 61.8 |
| BGE Large 0.34B | 76.0 | 46.1 | 87.1 | 60.0 | 54.3 | 83.1 | <u>31.6</u> | 64.2 |
| E5 Mistral 7B | 78.5 | 50.3 | **88.3** | 60.2 | 56.9 | **84.6** | 31.4 | 66.6 |
| **GRITLM** | | | | | | | | |
| Gen.-only 7B | 65.4 | 32.7 | 54.2 | 43.0 | 13.7 | 60.2 | 21.1 | 41.2 |
| Emb.-only 7B | <u>78.8</u> | **51.1** | 87.1 | **60.7** | **57.5** | <u>83.8</u> | 30.2 | **66.8** |
| GRITLM 7B | **79.5** | <u>50.6</u> | <u>87.2</u> | <u>60.5</u> | <u>57.4</u> | 83.4 | 30.4 | **66.8** |
| GRITLM 8x7B | 78.5 | 50.1 | 85.0 | 59.8 | 55.1 | 83.3 | 29.8 | 65.7 |

| Dataset (→) | MMLU 0 FS | GSM8K 8 FS, CoT | BBH 3 FS, CoT | TyDi QA 1 FS, GP | HumanEval 0 FS | Alpaca 0 FS, 1.0 | Avg. |
|---|---|---|---|---|---|---|---|
| Setup (→) Metric (→) | EM | EM | EM | F1 | pass@1 | % Win | |
| *Proprietary models* ♥ | | | | | | | |
| GPT-4-0613 | 81.4 | 95.0 | 89.1 | 65.2 | 86.6[†] | 91.2 | 84.8 |
| *Other Open Models* ♥ | | | | | | | |
| Zephyr 7B β | 58.6 | 28.0 | 44.9 | 23.7 | 28.5 | 85.8 | 44.9 |
| Llama 2 70B | 64.5 | 55.5 | 66.0 | **62.6** | 29.9♦ | 0.0 | 46.4 |
| Llama 2 Chat 13B | 53.2 | 9.0 | 40.3 | 32.1 | 19.6[†] | 91.4 | 40.9 |
| Llama 2 Chat 70B | 60.9 | 59.0 | 49.0 | 44.4 | 34.3[†] | <u>94.5</u> | 57.0 |
| Tülu 2 13B | 55.4 | 46.0 | 49.5 | 53.2 | 31.4 | 78.9 | 52.4 |
| Tülu 2 70B | <u>67.3</u> | **73.0** | <u>68.4</u> | 53.6 | 41.6 | 86.6 | <u>65.1</u> |
| Mistral 7B Inst. | 53.0 | 36.0 | 38.5 | 27.8 | 34.0 | 75.3 | 44.1 |
| Mixtral 8x7B Inst. | **68.4** | <u>65.0</u> | 55.9 | 24.3 | **53.5** | **94.8** | 60.3 |
| **GRITLM** | | | | | | | |
| Emb.-only 7B | 23.5 | 1.0 | 0.0 | 21.0 | 0.0 | 0.0 | 7.6 |
| Gen.-only 7B | 57.5 | 52.0 | 55.4 | 56.6 | 34.5 | 75.4 | 55.2 |
| GRITLM 7B | 57.6 | 57.5 | 54.8 | 55.4 | 32.8 | 74.8 | 55.5 |
| GRITLM 8x7B | 66.7 | 61.5 | **70.2** | <u>58.2</u> | <u>53.4</u> | 84.0 | 65.7 |

# GritLM-7B & GritLM-8x7B



**Embedding Performance** 🔢

**Generative Performance** 🅰🅱

| Task (→) | CLF | Clust. |
|---|---|---|
| Metric (→) | Acc. | V-Meas. |
| Dataset # (→) | 12 | 11 |
| | | |
| OpenAI v3 | 75.5 | 49.0 |
| | | |
| Llama 2 70B | 60.4 | 29.0 |
| Mistral 7B | 63.5 | 34.6 |
| Mistral 7B Instruct | 67.1 | 34.6 |
| GPT-J 6B | 66.2 | 39.0 |
| SGPT BE 5.8B | 68.1 | 40.3 |
| Instructor XL 1.5B | 73.1 | 44.7 |
| BGE Large 0.34B | 76.0 | 46.1 |
| E5 Mistral 7B | 78.5 | 50.3 |
| | | |
| Gen.-only 7B | 65.4 | 32.7 |
| Emb.-only 7B | 78.8 | 51.1 |
| GRITLM 7B | 79.5 | 50.6 |
| GRITLM 8x7B | 78.5 | 50.1 |

| QA | HumanEval | Alpaca | Avg. |
|---|---|---|---|
| GP | 0 FS | 0 FS, 1.0 | |
| | pass@1 | % Win | |
| | | | |
| 86.6† | 91.2 | 84.8 |
| | | | |
| 28.5 | 85.8 | 44.9 |
| 29.9♦ | 0.0 | 46.4 |
| 19.6† | 91.4 | 40.9 |
| 34.3† | 94.5 | 57.0 |
| 31.4 | 78.9 | 52.4 |
| 41.6 | 86.6 | 65.1 |
| 34.0 | 75.3 | 44.1 |
| 53.5 | 94.8 | 60.3 |
| | | | |
| 0.0 | 0.0 | 7.6 |
| 34.5 | 75.4 | 55.2 |
| 32.8 | 74.8 | 55.5 |
| 53.4 | 84.0 | 65.7 |

# Questions thus far?

Next: RAG with GRIT

# RAG with GRIT



**Traditional RAG**

To slow down your speed of aging, you can…

**Generative Model**

How to prevent aging?

Technological and lifestyle factors may influence an individual's longevity. Cellular reprogramming…

**Index**

**Embedding Model**

How to prevent aging?

# RAG with GRIT

# RAG with GRIT



**Traditional RAG**

To slow down your speed of aging, you can…

**Generative Model**

How to prevent aging?

Technological and lifestyle factors may influence an individual's longevity. Cellular reprogramming…

**Index**

**Embedding Model**

How to prevent aging?

**Query Caching**

**1st Cache:** Reuse query representation for retrieval

**Index**

**GritLM**

How to prevent aging?

To slow down your speed of aging, you can…

Technological and lifestyle factors may influence an individual's longevity. Cellular reprogramming…

**Query-Doc Caching**

**1st Cache:** Reuse query representation for retrieval

**Index**

**GritLM**

How to prevent aging?

To slow down your speed of aging, you can…

**2nd Cache:** Reuse document key-value states for generation

# Attention mismatch problem

## 1) Combining bidirectional & causal attention

**Embed query/doc bidirectionally & cache**



**Reuse bidirectionally attended cache for causal generation**

## 2) Combining separately attended texts (only if caching both, query-doc/doc-query)

**Embed query & cache**



**Reuse separately attended caches for causal generation**

**Embed doc & cache**

# Caching Performance

|  | Match (0-shot, ↑) |
|---|---|
| No RAG | 21.00 |

**Query Caching Generation Prompt**

```
<|user|>
GRIT is…
Optionally using the prior
context answer the query
prior to it
<|assistant|>….
```

|  | Match (0-shot, ↑) |
|---|---|
| RAG | <u>30.50</u> |
| Query Caching | 25.46 |
| Query-Doc Caching | 21.63 |

**Doc Caching Generation Prompt**

```
<|user|>
What is "GRIT"?
Answer the prior query while
optionally using the context
prior to it
<|assistant|>….
```
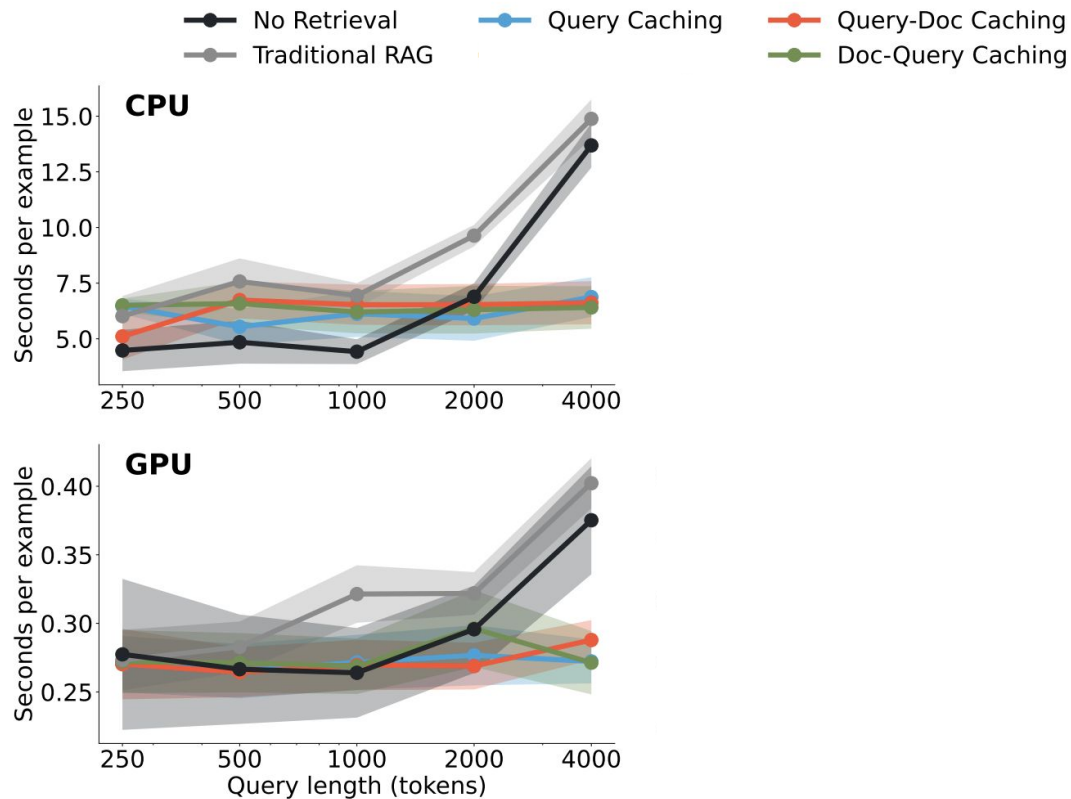
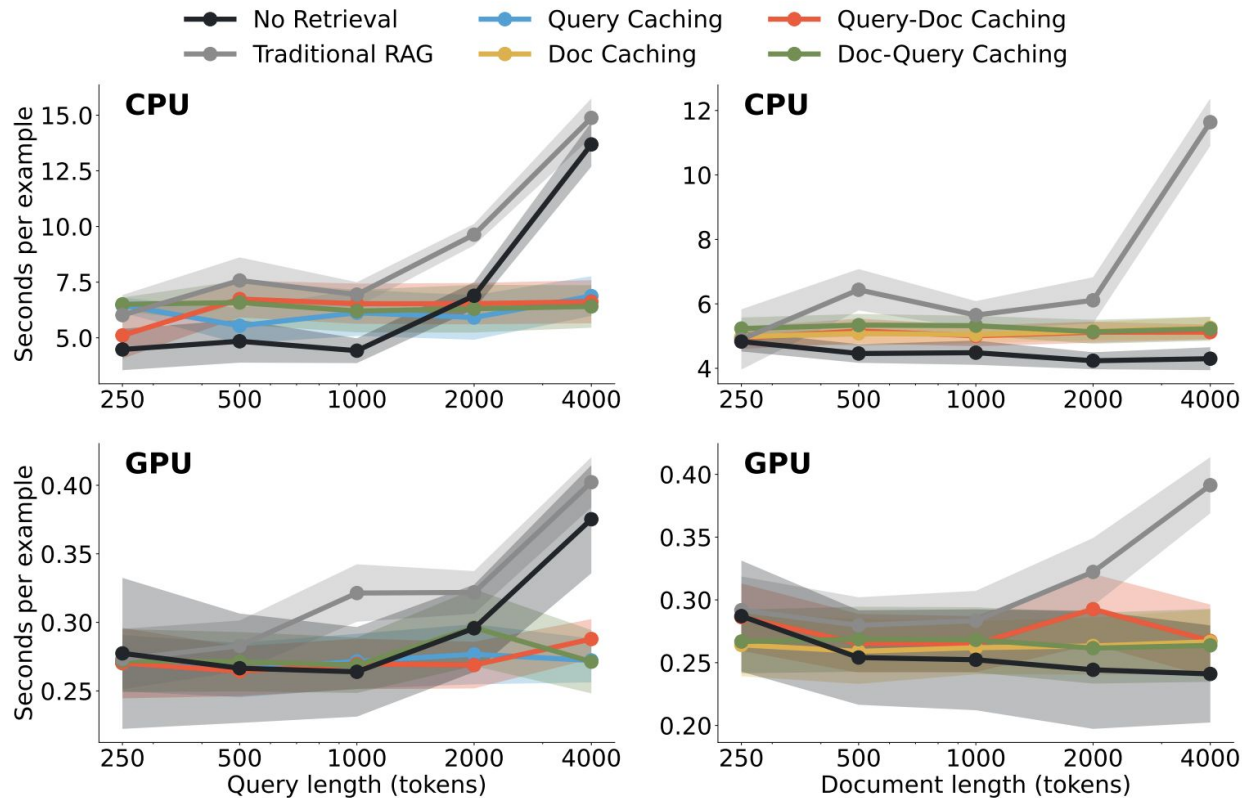|  | Match (0-shot, ↑) |
|---|---|
| RAG | 30.47 |
| Doc Caching | **33.38** |
| Doc-Query Caching | 18.39 |

# Scaling Query Length

# Scaling Document Length

# Questions thus far?

Next: GRIT Ablations

# Ablations: Attention & Pooling

| Attention Emb | | Attention Gen | | Pooling | Emb | Gen |
|---|---|---|---|---|---|---|
| Instruction | Sample | Instruction | Sample | | | |
| *Embedding Only* | | | | | | |
| Causal | | | | Wmean | 60.0 | - |
| Causal | Bidirectional | | | Mean | 61.0 | - |
| Bidirectional | | | | Mean | 61.8 | - |

# Ablations: Attention & Pooling

| Attention Emb | | Attention Gen | | Pooling | Emb | Gen |
|---|---|---|---|---|---|---|
| Instruction | Sample | Instruction | Sample | | | |
| *Embedding Only* | | | | | | |
| Causal | | | | Wmean | 60.0 | - |
| Causal | Bidirectional | | | Mean | 61.0 | - |
| Bidirectional | | | | Mean | 61.8 | - |
| *Generative Only* | | | | | | |
| | | Causal | | | - | **55.2** |
| | | Bidirectional | Causal | | - | 50.7 |

# Ablations: Attention & Pooling

| | Attention Emb | | Attention Gen | | Pooling | Emb | Gen |
| Instruction | Sample | Instruction | Sample | | | |
|---|---|---|---|---|---|---|---|
| | | | *Embedding Only* | | | | |
| | Causal | | | | Wmean | 60.0 | - |
| Causal | Bidirectional | | | | Mean | 61.0 | - |
| | Bidirectional | | | | Mean | 61.8 | - |
| | | | *Generative Only* | | | | |
| | | | Causal | | | - | **55.2** |
| | | Bidirectional | Causal | | | - | 50.7 |
| | | | *Unified* | | | | |
| | Causal | | Causal | | Last token | 61.2 | 53.0 |
| | Causal | | Causal | | Wmean | 62.8 | 52.8 |
| | **Bidirectional** | | **Causal** | | **Mean** | **64.0** | 52.9 |

# Ablations: Base Model

| Variant | Emb | Gen |
|---|---|---|
| **Mistral 7B** | **54.6** | **22.4** |
| Llama 2 7B | 48.2 | 20.8 |
| GPT-J 6B | 51.9 | 14.0 |

**Finetuned with GRIT**

# Embedding Performance after GRIT ≠ Raw performance

| Task (→) | CLF | Clust. | | Avg. |
|---|---|---|---|---|
| Metric (→) | Acc. | V-Meas. | | |
| Dataset # (→) | 12 | 11 | | 56 |
| Llama 2 70B | 60.4 | 29.0 | | 35.6 |
| Mistral 7B | 63.5 | 34.6 | | 40.5 |
| Mistral 7B Instruct | 67.1 | 34.6 | | 43.7 |
| GPT-J 6B | 66.2 | 39.0 | | 45.2 |

**Only pretrained**

# Ablations: Base Dataset

| Dataset | Emb |
| --- | --- |
| MEDI | 64.0 |
| MEDI2 | 64.7 |
| **E5** | **66.0** |

| Dataset | Gen |
| --- | --- |
| **Tülu 2** | **55.2** |
| OASST | 37.7 |
| UltraChat | 47.4 |

**Includes GPT-4 generated samples**

**Collection of many instruction datasets**

# Ablations: Embedding Head

🔢 🔤

| Variant | Emb | Gen |
|---------|-----|-----|
| **No head** | **62.7** | **49.2** |
| -> 1024 | 62.1 | 48.0 |

**GritLM** → `0.01, -1.01, 0.45…` → `0.42, 1.52…`

Linear downprojection

4096 dimensions

1024 dimensions

# Ablations: Batch Size (BS)

| BS Emb:Gen | Emb | Gen |
|---|---|---|
| 256:256 | 63.2 | **53.4** |
| 4096:256 | **64.2** | 53.3 |

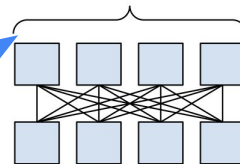**Increases M**
➡️ **more negatives**

$$\mathcal{L}_{\text{Rep}} = -\frac{1}{M} \sum_{i=1}^{M} \log \frac{\exp(\tau \cdot \sigma(f_\theta(q^{(i)}), f_\theta(d^{(i)})))}{\sum_{j=1}^{M} \exp(\tau \cdot \sigma(f_\theta(q^{(i)}), f_\theta(d^{(j)})))}$$

# Ablations: Precision

| Precision | 🔢 Emb | 🆎 Gen |
|-----------|--------|--------|
| FP32      | 66.3   | 52.4   |
| **BF16*** | **66.5** | **55.0** |

**\*Pooling & cosine similarity still in FP32 i.e. cast BF16->FP32 here**

Mean Pooling

```
<s><|user|>
{instruction}
<|embed|>
{sample to represent}
```

# Ablations: In-Batch-Negatives (IBN)

| IBN origin | Emb | Gen |
|---|---|---|
| Any dataset | **66.0** | 50.9 |
| **Same dataset** | **66.0** | **51.1** |

**Massive boost on Retrieval though:**

| Retrieval nDCG |
|---|
| 54.9 |
| 56.2 |

# Ablations: Format

| Format | Gen |
|---|---|
| **Tülu 2** | **55.2** |
| Zephyr $\beta$ | 49.0 |

**Additional end-of-sequence token after user utterance**

# Ablations: Loss

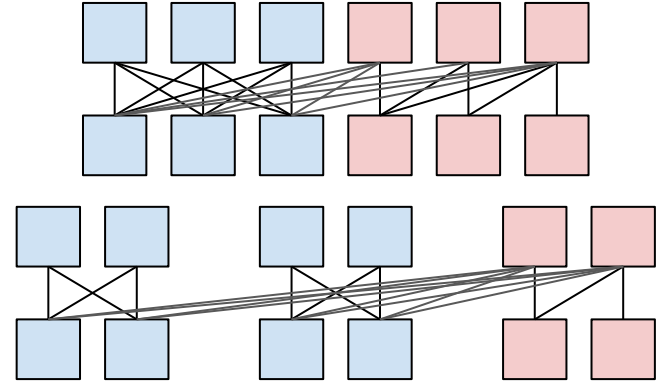| Gen loss type | $\mathcal{L}_{\text{Rep}}/\mathcal{L}_{\text{Gen}}$ | Emb | Gen |
|---|---|---|---|
| Token | 2.4 | 66.1 | 54.4 |
| Token | 6.0 | 66.5 | 55.0 |
| **Mix (32 -> 8)** | 4.1 | **66.7** | **55.4** |

**Ratio adjusted via weights:**

$$\mathcal{L}_{\text{GRIT}} = \lambda_{\text{Rep}}\mathcal{L}_{\text{Rep}} + \lambda_{\text{Gen}}\mathcal{L}_{\text{Gen}}$$

# Future Directions



## 1) Solving the attention mismatch issues

Simple finetuning may suffice?

## 2) GritLM Agents

Teaching the model to invoke its own embedding capabilities, maybe via finetuning

## 3) GRIT Doc Caching with multiple documents

Likely solved once the attention mismatch issue is solved

# Thanks!

**More about GRIT:**



**Paper**

**Open-source code:**
**ContextualAI/gritlm**

**Open-source models:**
**hf.co/GritLM** 🤗

**Niklas Muennighoff (**Twitter: @Muennighoff**)**

Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, Douwe Kiela

# Reranking with GritLM further boosts performance

**Table 3: Reranking (Rerank) using GRITLM as both Bi- and Cross-Encoder.**

| MTEB DS (↓) | No Rerank | Rerank top 10 |
|---|---|---|
| ArguAna | 63.24 | **64.39** |
| ClimateFEVER | 30.91 | **31.85** |
| CQADupstack | 49.42 | **50.05** |
| DBPedia | 46.60 | **47.82** |
| FiQA2018 | 59.95 | **60.39** |
| FEVER | 82.74 | **82.85** |
| HotpotQA | 79.40 | **80.46** |
| NFCorpus | 40.89 | **41.23** |
| NQ | 70.30 | **71.49** |
| MSMARCO | 41.96 | **42.47** |
| QuoraRetrieval | **89.47** | 88.67 |
| SCIDOCS | 24.41 | **24.54** |
| SciFact | 79.17 | **79.28** |
| TRECCOVID | 74.80 | **75.24** |
| Touche2020 | 27.93 | **28.41** |
| Average | 57.4 | **57.9** |

# Few-shot embedding does not work

Table 4: **Few-shot embedding.** The 12 MTEB datasets ("DS") are grouped by the 7 main MTEB tasks in the same order as in Table 1.

| Train DS (→) | E5S | | MEDI2 | |
| --- | --- | --- | --- | --- |
| MTEB DS (↓) | 0 FS | 1 FS | 0 FS | 1 FS |
| Banking77 | **88.5** | 88.3 | **88.1** | 87.9 |
| Emotion | **52.8** | 51.0 | **52.5** | 51.9 |
| IMDB | **95.0** | 93.9 | **94.3** | 92.2 |
| BiorxivS2S | **39.8** | 39.4 | **37.6** | 37.4 |
| SprintDup. | 93.0 | **94.9** | 95.2 | **95.7** |
| TwitterSem | **81.1** | 77.9 | **76.8** | 73.9 |
| TwitterURL | **87.4** | 87.1 | 85.9 | **86.1** |
| ArguAna | **63.2** | 51.7 | **53.5** | 53.2 |
| SCIDOCS | **24.4** | 19.7 | **25.5** | 25.5 |
| AskUbuntu | **67.3** | 64.7 | **66.6** | 66.0 |
| STS12 | 77.3 | **78.0** | **76.6** | 73.5 |
| SummEval | **30.4** | 29.5 | 29.1 | **31.5** |

# GRIT + KTO

Table 5: **Aligning GRITLM with KTO after GRIT.** The upper table depicts embedding performance while the lower depicts generative performance.

| Task (→)<br>Metric (→)<br>Dataset # (→) | CLF<br>Acc.<br>12 | Clust.<br>V-Meas.<br>11 | PairCLF<br>AP<br>3 | Rerank<br>MAP<br>4 | Retrieval<br>nDCG<br>15 | STS<br>Spear.<br>10 | Summ.<br>Spear.<br>1 | Avg.<br><br>56 |
|---|---|---|---|---|---|---|---|---|
| GritLM-7B | 79.5 | 50.6 | 87.2 | 60.5 | 57.4 | 83.4 | 30.4 | 66.8 |
| GritLM-7B-KTO | 79.6 | 50.1 | 87.1 | 60.5 | 57.1 | 83.5 | 30.5 | 66.7 |
| GritLM-8x7B | 78.5 | 50.1 | 85.0 | 59.8 | 55.1 | 83.3 | 29.8 | 65.7 |
| GritLM-8x7B-KTO | 78.7 | 50.0 | 84.4 | 59.4 | 54.1 | 82.5 | 30.8 | 65.2 |

| Dataset (→)<br>Setup (→)<br>Metric (→) | MMLU<br>0 FS<br>EM | GSM8K<br>8 FS, CoT<br>EM | BBH<br>3 FS, CoT<br>EM | TyDi QA<br>1 FS, GP<br>F1 | HumanEval<br>0 FS<br>pass@1 | Alpaca<br>0 FS, 1.0<br>% Win | Avg.<br><br> |
|---|---|---|---|---|---|---|---|
| GritLM-7B | 57.6 | 57.5 | 54.8 | 55.4 | 32.8 | 74.8 | 55.5 |
| GritLM-7B-KTO | 57.6 | 57.5 | 55.4 | 55.8 | 31.5 | 86.7 | 57.4 |
| GritLM-8x7B | 66.7 | 61.5 | 70.2 | 58.2 | 53.4 | 84.0 | 65.7 |
| GritLM-8x7B-KTO | 66.8 | 79.5 | 67.1 | 31.4 | 56.8 | 95.3 | 66.2 |