Negative Matters: Multi-Granularity Hard-Negative Synthesis for Enhanced Text Embeddings

Anonymous ACL submission

Abstract

Text embedding models are essential for various natural language processing tasks, enabling the effective encoding of semantic information into dense vector representations. These models are typically optimized using triplets of (query, positive, negative) data pairs for contrastive learning, where the negative samples play a critical role in enhancing the model's ability to discern subtle semantic distinctions. In this work, we introduce a Multi-Granularity Hard-negative (MGH) synthesis framework that leverages large language models (LLMs) to generate diverse negative samples with varying levels of similarity with the query. This approach facilitates a coarse-to-fine curriculum learning strategy during supervised training, allowing the embedding model to progressively learn more nuanced semantic representations. Meanwhile, we propose an Anchor Token Aware (ATA) pooling method that assigns higher weights to anchor tokens based on aggregation patterns observed in LLMs, improving text embedding accuracy without increasing model complexity. Comprehensive experiments on the MTEB benchmark demonstrate that our methods achieve state-of-theart performance, surpassing existing synthesis strategies both with synthetic data and when combined with public retrieval datasets.

1 Introduction

004

011

012

014

018

040

043

Text embedding models are designed to encode the semantic meaning of a given sequence of natural language words, sentences, or larger text spans into dense vector representations. These vector representations capture not only the lexical content of the text but also its syntactic and semantic nuances, facilitating a wide range of downstream natural language processing (NLP) tasks such as sentiment analysis, text clustering, and content-based information retrieval.

Previous studies havehave explored the potential of leveraging large language models (LLMs) to



Figure 1: Illustration of our proposed multi-granularity hard-negative sample generation and coarse-to-fine learning paradigm.

generate synthetic data for text embedding training tasks (Bonifacio et al., 2022; Wang et al., 2024; Lee et al., 2024b). The substantial volume of synthetic data contributes to increased diversity, thereby improving the model's robustness across various encoding tasks. However, generating high-quality hard negative samples required by contrastive learning is a challenging task for synthetic models, as an effective hard negative must maintain the right balance in its distinction from the positive examples. If the hard negative is overly similar to the query, it may confuse the model with positive samples, whereas if there is a significant difference in the content, the model may struggle to extract useful information from hard negative samples.

In this light, we propose a data synthesis framework to fully leverage the LLMs' ability to identify

the partial order of text similarity, facilitating the 061 generation of multi-granularity synthetic data. By 062 simultaneously generating hard negative samples 063 at multiple similarity levels, the synthesized data not only improves overall quality but also allows for controlled difficulty of the negatives. In the subsequent supervised training, we use the generated 067 difficulty-controllable data to implement coarse-tofine curriculum learning. By progressively moving from simple to hard training samples, the embedding model learns increasingly complex representations, resulting in improved stability and effectiveness.

074

101

102

104

105

106

108

109

110

111

112

Moreover, as research increasingly focuses on effectively adapting large models into text embedding models, studies (Lee et al., 2024a; BehnamGhader et al., 2024) have explored converting causal attention to bidirectional attention to enhance embedding performance. In this context, two common methods to obtain embeddings from the hidden states of a sequence of tokens are mean pooling, which averages the final hidden states, and last token pooling, which uses the last hidden state of the <EOS> token as the sentence representation vector. However, mean pooling tends to dilute the critical information tokens when averaging across all tokens, which leads to a loss of significant features (Lee et al., 2024a). In contrast, the last token pooling method is sensitive to noisy information within the sentence, resulting in instability in the encoding (Springer et al., 2024). Recently, NV-Embed (Lee et al., 2024a) addressed the insufficient pooling issue by adding a crossattention layer over the tokens' final hidden states in a dictionary learning method.

However, we suggest that a simple yet effective pooling method can be utilized without introducing additional parameters. Recent studies have revealed the aggregation pattern of large language models (Wang et al., 2023a; Huang et al., 2024), indicating that decoder-only models tend to aggregate textual information into anchor tokens at shallow layers and use these tokens to generate the next token in deeper layers.

In this paper, we observe that the aggregation pattern still holds in models transformed from causal to bidirectional attention. By simply assigning greater weight to anchor tokens, we achieve improved accuracy in text embedding tasks compared to conventional pooling methods.

Our contributions are summarized as follows:

1. We propose a Multi-Granularity Hard-

negative (MGH) synthesis framework that effectively generates diverse negative samples of varying difficulty levels, fully leveraging the large model's capability to discern the partial order of text similarity. The framework allows for controlled progression in the difficulty of the generated negative examples, enabling subsequent text embedding models to learn a more accurate embedding representation through a coarse-to-fine manner. 113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

2. We propose an Anchor Token Aware (ATA) pooling method that effectively leverages the aggregation pattern of LLMs to acquire a more accurate sentence representation. The model trained with ATA pooling outperformed previous pooling methods and, when trained solely on publicly available retrieval data, ranked among the top models on the MTEB leaderboard.¹

2 Method

2.1 Multi-granularity Synthetic Data Generation

The overall framework of our proposed data synthesis method is illustrated in Figure 1, which consists of two primary stages. In accordance with the setup of Wang et al. (2024), we begin by querying LLMs to generate a list of potential tasks, categorized into two types: asymmetric tasks (e.g., shortlong match, long-short match, long-long match and short-short match), and symmetric tasks (e.g., semantic textual similarity, STS). The task brainstorming process in stage 1 generates a wide range of embedding tasks to enrich the diversity of synthetic data, enriching the diversity of synthetic data.

In stage 2, using a diverse set of tasks as seeds, we query the LLM to generate (query, positive, negative) samples for subsequent contrastive supervised training, which are then used in contrastive supervised training with a standard infoNCE objective to minimize the distance between query and positive samples, while maximizing the separation from negative samples:

$$\mathcal{L} = -\log \frac{\phi(q, d^+)}{\phi(q, d^+) + \sum_{d^- \in N} (\phi(q, d^-))} \quad (1)$$

¹The concurrent work by Li et al. (2024) achieved the state-of-the-art result on the MTEB benchmark using publicly available retrieval dataset. However, they adopt an in-context-learning prompting method, which we believe is largely or-thogonal to our approach.

where ϕ denotes the similarity function, q, d^+ and d^- represent the query, positive, and negative samples respectively. Recognizing that the quality of negative samples significantly impacts the supervised training of text-embedding models, we aim to fully leverage the large model's ability to distinguish between different granularities of negative samples during the generation process. To be specific, we formulate the synthetic target as follows:

156

157

158

159

161

162

163

164

165

167

168

169

170

171

172

173

174

175

176

177

180

181

188

189

193

194

195

197

198

205

$$\mathcal{S}^S = \{ (\mathcal{Q}^S, \mathcal{P}^S, \{\mathcal{N}_k^S\}) \}$$
(2)

where Q^S represents an example query, \mathcal{P}^S denotes its corresponding positive sample, and $\{\mathcal{N}_k^S\}$ refers to a set of hard negative samples with varying levels of granularity indexed by k, which is set to 4 in following experiments.

We use the template illustrated in Figure 2 to constrain the synthesizing format. Differing from Wang et al. (2024), our approach prompts LLM to simultaneously generate multiple hard negative samples $\{\mathcal{N}_k^S\}$ for a single query \mathcal{Q}^S . These negative samples are ranked based on their similarity to the query, arranged in descending order from highest to lowest. This approach allows large models to enforce similarity constraints when generating hard negative samples, effectively mitigating the uncontrolled variation in the hard negative sample similarity during the synthetis process across different query samples. The effectiveness of this approach is demonstrated in Section 5.1 through a detailed case study.

In addition to synthetic data, we also incorporate public retrieval datasets when training the text embedding model. To integrate coarse-tofine hard negative samples into the subsequent training process, we regenerate the negative samples of the retrieval dataset, denoted as $\{\mathcal{N}_k^R\}$, by querying LLM using the same multi-granularity approach. The refined retrieval dataset, $S^R =$ $\{(\mathcal{Q}^R, \mathcal{P}^R, \{\mathcal{N}_k^R\})\}$, is then combined with the synthetic dataset S^S to form the complete training dataset S.

With the multi-granularity dataset S, we adopt a curriculum learning paradigm for the supervised training of the text embedding model. By adjusting the difficulty level k of the hard negatives, we can progressively train a more stable and effective model using a coarse-to-fine approach. Specifically, during the supervised learning process, we gradually transition the negative samples fed into

Omitted for space limitations The output should be formatted as a JSON object with a field indicating the relative similarity level. Use the following format as a guide: "query": "QUERY_TEXT", "positive_example": "POSITIVE_EXAMPLE_TEXT", "hard_negative_examples": ["similarity_level": "high", "text": "HIGH_SIMILARITY_NEGATIVE_EXAMPLE_TEXT" },{ 'similarity_level": "medium", "text": "MEDIUM_HIGH_SIMILARITY_NEGATIVE_EXAMPLE_TEXT" },{ . "similarity_level": "medium", "text": "MEDIUM_LOW_SIMILARITY_NEGATIVE_EXAMPLE_TEXT" },{ similarity_level": "low", "text": "LOW_SIMILARITY_NEGATIVE_EXAMPLE_TEXT"} 1 } Omitted for space limitations

Please generate four negative examples for contrastive learning based on the generated query and positive example. These examples should be arranged in order of decreasing similarity to the query, ranging from highly similar to dissimilar. Ensure the similarity spans a broad spectrum, and every negative example should be different, without repeating words from previous examples. Be creative! Omitted for space limitations

Figure 2: The core template used for prompting LLMs to generate multi-granularity hard negatives. Due to space constraints, the full prompts are presented in Appendix A.1.

the model from N_4 to N_1 , with each difficulty level occupying a quarter of the training data. Further analysis in Section 4.1 demonstrates the effectiveness of the proposed scheduling method.

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

224

225

226

227

228

230

2.2 Anchor Token Aware Pooling

After obtaining the hidden states of multiple tokens from the model's final layer, an appropriate approach is required to derive the sentence representation vector v. In the widely adopted mean pooling method, the last hidden states of all tokens are averaged to form a sentence vector representation. This approach can result in the dilution of key information in the text, as non-critical tokens are averaged along with the more significant ones.

The proposed ATA pooling method aims to assign greater weight to anchor tokens (Wang et al., 2023a), which aggregate more semantic information compared to other tokens. This approach allows the model to adaptively allocate weights to the parts that are most pertinent to the task, resulting in a more effective pooling operation.

Motivated by Huang et al. (2024), we calculate the attention weight along the Key dimension of the attention matrix to identify anchor tokens with stronger representational capabilities. Compared to traditional mean pooling, this allows us to assign higher weights to anchor tokens and more effectively filter out trivial tokens that do not contribute to the semantic information.

231

240

241

244

245

247

248

249

250

251

257

259

261

262

265

Specifically, let $A_{\mathcal{L}}^{h}$ denote the attention matrix for the attention head h in the model's final layer, and let a_{ij}^{h} represent the corresponding value of $A_{\mathcal{L}}^{h}[i][j]$, which indicates the attention score between query i and key j. We define the anchor weight of each query as follows:

$$N_i = \sum_{h=1}^{H} \sum_{j=1}^{S} \log(a_{ij}^h \cdot S + 1)$$
(3)

where S represents the length of the token sequence, and H denotes the number of attention heads. We multiply a_{ij}^h by S because, while the sum of attention weights in the query dimension remains constant (i.e. $\sum_{i=0}^{S} a_{ij}^h = 1$),the expected value of each individual element decreases to $\frac{1}{S}$ as the sentence length increases. Multiplying by the sentence length helps to ensure stability across different sentence lengths by maintaining consistent scaling in subsequent computations.

After obtaining the anchor weights, we normalize them by applying a linear weight adjustment along the Query dimension to compute the weight corresponding to each token, denoted as w_i , such that the sum of all w_i equals 1:

$$w_i = \frac{N_i}{\sum_{i=1}^S N_i} \tag{4}$$

We then apply the normalized weights to reweight the hidden states and obtain the final sentence embedding v:

$$v = \sum_{i=1}^{S} w_i d_i \tag{5}$$

where d_i denotes the model's last layer hidden state of token *i*.

3 Experiments

3.1 Data Synthetic Details

266To ensure a fair comparison with Wang et al. (2024),267we generated an equivalent volume of synthetic268data, maintaining a consistent total token consump-269tion of 180M. In order to minimize the costs asso-270ciated with data generation, we utilized the APIs271of GPT-40 and DeepSeek. We observed that com-272pared to GPT-40, DeepSeek is relatively less cre-273ative when generating the potential tasks in stage



Figure 3: Distribution of the task categories in the synthetic data

274

275

276

277

278

279

280

281

282

283

285

290

291

292

293

294

295

296

297

299

300

301

302

304

305

307

308

1 and tends to produce repetitive negative samples during stage 2. Consequently, we relied on GPT-40 to complete all stage 1 generation processes. In stage 2, we initially used GPT-40 to generate a sufficient amount of data, which was then input as seeds into DeepSeek. Leveraging the data cache provided by the DeepSeek API, introducing additional seeds as input did not incur significant additional costs. The distribution of synthetic data across different task types is illustrated in Figure 3.

The retrieval dataset used in supervised learning was curated by Springer et al. (2024), which consist of approximately 1.5 million samples, covering a variety of languages and retrieval scenarios. In line with LLM2Vec (BehnamGhader et al., 2024), we only used about one-third of the curated retrieval dataset. For more details on the dataset composition, please refer to Appendix B.2.

3.2 Experimental Details

To validate the effectiveness of our proposed finegrained data synthesis framework and the ATA pooling method, we perform experiments following the open source LLM2VEC (BehnamGhader et al., 2024) model, while replacing the supervised training stage with our method. Specifically, we adopt Mistral-7B-Instruct-v0.2 (Jiang et al., 2023) as the base model. We then transform the model's attention pattern from causal to bidirectional and integrate the LoRA weights trained on the masked nexttoken prediction task introduced in LLM2VEC, enabling the model to better adapt to bidirectional attention patterns. This setup serves as the starting point for our model.

For supervised training, we adopt the standard InfoNCE loss, with both in-batch negatives and hard

	FT. Data	Class.	Clust.	Pair.	Rerank.	Retr.	STS	Summ.	Avg.
	Sample Num	Acc.	V-Meas.	AP	MAP	nDCG@10	Spear.	Spear.	
Models trained with synthetic data only									
Mistral _{gpt-40} (Chen et al., 2024)	230K	77.7	47.7	83.9	58.7	46.7	80.9	30.7	62.2
Gecko _{01b-768} (Lee et al., 2024b)	6.6M	70.3	46.8	86.2	57.6	53.2	83.1	32.2	62.6
E5 _{mistral-7b} (Wang et al., 2024)	500K	78.2	50.5	86.0	59.0	46.9	81.2	<u>31.9</u>	63.1
SPEED (Chen et al., 2024)	920K	<u>78.3</u>	48.6	86.3	<u>59.8</u>	48.1	<u>82.6</u>	31.7	<u>63.4</u>
MGH(Ours)	310K	78.6	<u>49.7</u>	86.1	60.1	<u>51.2</u>	82.3	31.6	64.5
Models trained with synthetic data & public available retrieval data									
GTR _{xxl} (Ni et al., 2022)	662K	67.4	42.4	86.1	56.7	48.5	78.4	30.6	59.0
text-embedding-3 _{large} ²	-	75.5	49.0	85.7	59.2	55.4	81.7	29.9	64.6
jina-embeddings-v3 (Sturua et al., 2024)	-	82.6	45.3	84.0	58.1	53.9	85.8	29.7	65.5
Gecko _{01b-768} (Lee et al., 2024b)	>6.6M	<u>81.2</u>	47.5	87.6	58.9	55.7	85.1	32.6	66.3
E5 _{mistral-7b} (Wang et al., 2024)	1.8M	78.5	50.3	88.3	<u>60.2</u>	<u>56.9</u>	84.6	<u>31.4</u>	<u>66.6</u>
SPEED (Chen et al., 2024)	2.2M	78.4	49.3	88.2	60.8	56.5	85.5	31.1	66.5
MGH(Ours)	820K	78.8	<u>50.1</u>	87.9	59.8	57.5	85.6	31.3	67.0

Table 1: MTEB performance comparison of different synthesis models, with training conducted on synthetic data only and on both synthetic and public retrieval dataset. The highest performances are highlighted in bold, while the second-highest are indicated with underlines. *FT. Data Sample Num.* refers to the number of (query, positive, negative) sample pairs used for training.

negatives utilized for training. To ensure a fair comparison, the prompt template follows Wang et al. (2024), as illustrated in Appendix B.2. Supervised training on the full dataset is conducted for 1600 steps, while training on public retrieval dataset is performed for 1000 steps, with a batch size of 64 and gradient accumulation of 8 in both cases. All training was performed on a single 80GB H100 GPU, taking approximately 32 hours to complete 1600 training steps. Further training hyperparameters are represented in Appendix B.1.

For evaluation, we assess our model on the widely used MTEB benchmark (Muennighoff et al., 2023), which encompasses a wide variety of text embedding tasks across different scenarios and domains. The benchmark includes 56 English embedding tasks organized into 7 categories: classification (12), clustering (11), pair classification (3), reranking (4), retrieval (15), semantic textual similarity (10), and summarization (1).

3.3 Main Results

311 312

313

314

315

318

319

320

321

326

327

332

333

335

336

338

339

Table 1 presents a comparison of the performance of our proposed method against existing data synthesis approaches on the MTEB dataset. The results underscore the effectiveness of our data generation framework, demonstrating superior performance in both the synthetic-only and full-data settings. In particular, on the more challenging retrieval tasks, the full-data results achieved the best performance, underscoring the efficacy of our synthesis method.

Model	MTEB Score
SGPT (Muennighoff, 2022)	58.93
UDEVER-bloom-7b (Zhang et al., 2023a)	60.63
ECHO (Springer et al., 2024)	64.68
LLM2Vec (BehnamGhader et al., 2024)	64.80
NV-Embed (Lee et al., 2024a)	64.18
bge-large-en-v1.5 (Xiao et al., 2024)	64.23
bge-en-icl w/o icl (Li et al., 2024)	64.83
bge-en-icl w/ icl (Li et al., 2024)	66.08
MGH(Ours)	<u>65.87</u>

Table 2: Performance comparison of MTEB scores for different models trained on publicly available retrieval corpora. The highest performances are highlighted in bold, while the second-highest are indicated with underlines.

To evaluate the effectiveness of the proposed ATA pooling method, we compare the ATA pooling enhanced model over previous models that leverage LLM for text embedding. However, prior research has indicated that incorporating the training split of the MTEB dataset during the supervised training process introduces a significant amount of MTEB-related data, thereby increasing the risk of overfitting (Li et al., 2024). Therefore, we opted not to include the second training phase proposed by NV-Embed, which utilizes the MTEB training split for a continuous training.

Table 2 compares the performance of our approach to existing models trained exclusively on publicly available retrieval data. The results under the MTEB training split free setting indicate that our model achieved performance slightly below the state-of-the-art model, exhibiting the effectiveness of the ATA pooling method. Notably, the

358

340

341

342

²https://platform.openai.com/docs/guides/embeddings



Figure 4: Trend of MTEB subset scores during supervised training across four experimental settings.

method proposed by Li et al. (2024) employed an in-context learning approach to enhance the text embedding quality, which operates orthogonally to our proposed method.

4 Ablation Study

364

367

372

376

386

4.1 Data Synthesis and Training Strategies

In this section, we evaluate the effectiveness of the proposed MGH synthesis framework through an ablation study on a subset of the MTEB benchmark, as used in Springer et al. (2024)³. We evaluate the impact of data training order on supervised learning through the following experimental settings, analyzing how the results evolve throughout the training process over 1600 steps:

- 1. **Curriculum Learning:** Progressing from easier to more challenging hard negatives, as adopted in our main approach.
- 2. **Reverse Curriculum Learning:** Progressing hard negative samples from harder to easier.
- Random Ordering: Randomly inputting hard negative examples of varying difficulty.
- 4. **Fixed Difficulty Level:** Consistently using hard negatives of a fixed difficulty level.

The results in Figure 4 demonstrate that the curriculum learning strategy adopted by MGH not only achieves the best performance but also maintains greater stability during training. In contrast, the random ordering strategy predictably exhibits fluctuations in the results, while the reverse curriculum learning method also fails to yield satisfactory training outcomes, as the model struggles to adapt to the reversed difficulty progression. 387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

Additionally, none of the four fixed difficulty settings outperformed the curriculum learning approach. Among them, maintaining a difficulty level of 2 or 3 yielded relatively better results, while using excessively low or high difficulty levels failed to converge to optimal outcomes. This suggests that overly simple negative examples may prevent the model from learning useful knowledge, and excessively difficult synthetic negatives may hinder the model's ability to distinguish them from positive examples at early training stages.

4.2 Pooling Methods

This section presents a detailed comparison of four pooling methods employed in embedding tasks, namely mean pooling, last token pooling, NV-Embed pooling (Lee et al., 2024a), and the proposed ATA pooling. The ablation experiments leveraged publicly available retrieval datasets only, adhering to the hyperparameter settings outlined in Section 3.2. The evaluation is conducted using the full MTEB benchmark, with Table 3 summarizing their respective performances.

Pooling Method	MTEB Score
mean pooling	65.41
last pooling	64.97
NV-Embed pooling	<u>65.80</u>
ATA pooling (Ours)	65.87

Table 3: Performance comparison of MTEB scores for different pooling methods trained on publicly available retrieval corpora. The highest performances are highlighted in bold, while the second-highest are indicated with underlines.

The results suggest that conventional mean pooling and last token pooling yield subpar performance in text embedding tasks using bidirectional models. On the other hand, both the NV-Embed pooling and ATA pooling methods demonstrate favorable results, validating the necessity of adaptive reweighting for the last hidden states.

5 Further Analysis

5.1 Case Study

How does MGH enhance hard negative sample quality? We illustrate this through the example presented in Figure 5, which demonstrates

 $^{^{3}}$ Conducting a full evaluation on the MTEB dataset is computationally expensive, requiring over 200 hours on a single H100 GPU. Therefore, a subset of the dataset was selected for this study. Details of the subset composition can be found in Appendix C.2.

Query: What is the Latin name for the element Sodium, giving rise to its chemical formula Na?

Positive: The Latin name for Sodium is *Natrium*, which is why its chemical symbol is Na. Sodium is a highly reactive alkali metal with atomic number 11, commonly found in nature as part of compounds like table salt (sodium chloride, NaCl). The name "Natrium" originates from the Arabic word *natron*, which refers to natural mineral deposits rich in sodium carbonate.

Hard Negative 1: The Latin name for Sodium is *Kalium*, giving rise to its chemical symbol K. Sodium is an alkali metal with the atomic number 19 and is essential in the human body for regulating blood pressure and nerve function.

Hard Negative 2: Sodium's chemical formula is *Sod*, derived from its Latin name *Sodium*. As an alkali metal, it is highly reactive and never found in its elemental state in nature.

Hard Negative 3: The Latin name for Hydrogen is *Hydrogenium*, which corresponds to its symbol H. Hydrogen is the lightest element, and while not related to Sodium, it shares its essential role in chemical reactions.

Hard Negative 4: Sodium is often associated with salt in cooking and baking, as it is a crucial element for flavor and preservation. Its role in chemistry includes forming compounds such as NaCl, but it is better known for its culinary uses.

Figure 5: An example of a multi-granularity hard negative synthesis, randomly selected for illustration.

how MGH effectively improves the quality of hard negative samples by leveraging multi-granularity similarity constraints.

425

426

497

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

In this example, the data synthetic model is tasked with generate negative samples for the Latin name of the element *Sodium*. The GPT-40 model used in this case selects *Potassium*, which shares similar chemical properties with *Sodium*, as the most challenging hard negative example. Subsequently, the model generates a fake Latin name for *Sodium* as a moderately confusable negative sample, followed by answering element *Hydrogen* as a more distinguishable example. While the first three negative samples involve answering the Latin names of chemical elements, the last simplest negative sample generated by the model focuses on *Sodium* but lacks any reference to its Latin name.

As demonstrated in the example above, the MGH approach effectively distills world knowledge from LLMs, enabling the generation of multiple negative samples with varying granularities. As the examples progress from challenging to simple, the synthetic model's outputs range from showing subtle differences in detail to being more easily distinguishable. In this process, the subsequent negative samples are adjusted based on previously generated ones, enabling a dynamic progression of negative sample difficulty, further enhancing the quality of negative sample generation.

How does ATA reweight using aggregation pat-

tern? As shown in the example from Figure 6,
the aggregation pattern still remains the base model
is transformed from causal to bidirectional atten-

tion. The figure illustrates three prominent anchor tokens: the initial token, the punctuation between the two sentences, and the [INST] template appended to the end of the sentence. Accordingly, these anchor tokens receive higher weight values in the ATA weight calculation, contributing more significantly to the subsequent computation of text embeddings. 458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

Through observations of numerous examples, we found that most samples allocate a greater proportion of the ATA weight to the three anchor patterns mentioned above, with particular emphasis on the [INST] token at the sentence's end. Therefore, the ATA pooling method captures the important last token while also dynamically identifying key positions within the preceding text. This approach not only mitigates the stability issues associated with relying solely on the last token but also assigns greater weight to tokens that are essential for capturing the entire semantic meaning of the input sequence, thereby facilitating more effective embedding learning.

5.2 Cost of Synthetic Data

Although our model entails additional tokens to generate multiple hard negatives per synthetic sample, the cost is offset by maintaining the same token consumption (180M) as Wang et al. (2024), which results in fewer synthetic samples being generated. Under this fair comparison, the embedding model trained with our MGH strategy outperforms previous state-of-the-art results, demonstrating that our method is more effective under the same token consumption.



Figure 6: The upper part illustrates the summed results of the model's final layer attention weights across 32 attention heads⁴, while the lower part shows the corresponding ATA weights for each token. Example is randomly selected in STS13 evaluation split.

6 Related Work

491

492

493

494

495

496

497

499

500

501

508

511

512

513

514

515

516

517

LLM Based Text Embedding Models In recent years, as decoder-only models have scaled up in terms of parameters and training data, researchers have explored the possibility of transforming nexttoken prediction models into effective text embedding models through continued training. Neelakantan et al. (2022) was the first to apply the GPT-3 model to text embedding tasks, leveraging the <EOS> token as the representation vector. Subsequent work by Ma et al. (2024) employed a similar last-token pooling method, fine-tuning the LLaMA-2 model.

However, the autoregressive training objective imposes an inherent limitation on the model's performance, as the causal attention mask prevents earlier tokens from accessing subsequent tokens. SGPT (Muennighoff, 2022) addressed this limitation by linearly assigning more weight to tokens at later positions, a strategy subsequently adopted by E5_{mistral-7b} (Wang et al., 2024). LLM2Vec (BehnamGhader et al., 2024) transformed the model from causal to bidirectional attention by employing a masked next-token prediction approach, followed by mean pooling for supervised learning. Recent work by NV-Embed (Lee et al., 2024a) introduced an additional cross-attention layer for hidden state pooling, simultaneously removing the causal mask. Additionally, Echo embeddings (Springer et al., 2024) repeated a text twice and used the second instance to compute the representation vector. In this work, we attempt to address the insufficient pooling problem in LLM based embedding models by an adaptive weighting strategy using the model's aggregation pattern.

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

559

560

561

562

563

564

565

567

Data Synthesis for Embedding Models Highquality data is crucial for training effective text embedding models. Previous studies (Nogueira et al., 2019; Wang et al., 2023b) have explored expansion-based approaches to augment document and query data. With the advancements in large language model (LLM) capabilities, recent research has focused on leveraging LLMs to generate large amounts of high-quality supervised training data (Wang et al., 2024; Jeronymo et al., 2023; Sturua et al., 2024). In domain-specific retrieval, studies (Dai et al., 2022; Khramtsova et al., 2024) have shown that LLM-generated query-document pairs significantly improve embedding quality in domainspecific retrieval tasks. Additionally, Gecko (Lee et al., 2024b) curates web data to enable LLMs to produce high-quality synthetic samples. Our work focuses on how to enable large models to generate multi-granularity synthetic negative examples, and achieves more efficient and stable text embedding model training by controlling the training difficulty.

7 Conclusion

In this work, we evaluate the importance of hard negative granularity when training text embedding models using contrastive learning. Through the proposed MGH synthesis framework, we generate diverse negative samples at varying levels of similarity, enabling the embedding model to learn more nuanced semantic representations by coarseto-fine curriculum learning approach. Experimental results demonstrate that our methods achieve state-of-the-art performance on the MTEB benchmark, outperforming existing synthesis strategies both with synthetic-only and combined datasets. Additionally, our proposed ATA pooling method effectively leverages the aggregation patterns inherent in large language models, improving sentence pooling efficacy without introducing extra parameters. Ablation studies confirm the effectiveness of our MGH framework and ATA pooling method in enhancing text embedding model performance and training stability.

⁴For space limitations, the individual attention maps from all 32 attention heads are provided in Appendix D.2.

568

589

590

593

596

600

603

607

611

612

613

614

615

616

617

618

8 Limitations

Despite the effectiveness of our method, there are several limitations that should be acknowledged: 570 (1) Due to the costs associated with API usage, we 571 limited the synthetic data generation to the same token volume as used in previous studies (Wang 573 et al., 2024). This constraint prevented us from 574 exploring whether a larger synthetic dataset could 575 further improve the performance of the text embedding model. We leave this exploration to future work, where more extensive synthetic data could be 578 generated to assess the scalability and potential per-579 formance gains. (2) To facilitate comparisons with 580 prior work (BehnamGhader et al., 2024; Springer et al., 2024), we used Mistral-7b-v0.2-Instruct as our base text embedding model. Given the con-583 584 tinuous advancements in 7b-level models, we plan to investigate more powerful models as our base embedding model in our future work.

References

- Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. 2016. Ms marco: A human generated machine reading comprehension dataset. arXiv preprint arXiv:1611.09268.
- Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. 2024. Llm2vec: Large language models are secretly powerful text encoders. *arXiv preprint arXiv:2404.05961*.
- Luiz Henrique Bonifacio, Hugo Queiroz Abonizio, Marzieh Fadaee, and Rodrigo Nogueira. 2022. Inpars: Data augmentation for information retrieval using large language models. *ArXiv*, abs/2202.05144.
- Haonan Chen, Liang Wang, Nan Yang, Yutao Zhu, Ziliang Zhao, Furu Wei, and Zhicheng Dou. 2024. Little giants: Synthesizing high-quality embedding data at scale. *Preprint*, arXiv:2410.18634.
- Zhuyun Dai, Vincent Y. Zhao, Ji Ma, Yi Luan, Jianmo Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith B. Hall, and Ming-Wei Chang. 2022. Promptagator: Few-shot dense retrieval from 8 examples. *Preprint*, arXiv:2209.11755.
- DataCanary, hilfialkaff, Lili Jiang, Meg Risdal, Nikhil Dandekar, and tomtung. 2017. Quora question pairs. https://kaggle.com/competitions/ quora-question-pairs. Kaggle.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. ELI5: Long form question answering. In *Proceedings of*

the 57th Annual Meeting of the Association for Computational Linguistics, pages 3558–3567, Florence, Italy. Association for Computational Linguistics. 619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2024. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13418– 13427.
- Vitor Jeronymo, Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, Roberto Lotufo, Jakub Zavrel, and Rodrigo Nogueira. 2023. Inpars-v2: Large language models as efficient dataset generators for information retrieval. *Preprint*, arXiv:2301.01820.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for opendomain question answering. In *Proceedings of the* 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6769–6781, Online. Association for Computational Linguistics.
- Ekaterina Khramtsova, Shengyao Zhuang, Mahsa Baktashmotlagh, and Guido Zuccon. 2024. Leveraging llms for unsupervised dense retriever ranking. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1307–1317.
- Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024a. Nv-embed: Improved techniques for training llms as generalist embedding models. *arXiv preprint arXiv:2405.17428*.
- Jinhyuk Lee, Zhuyun Dai, Xiaoqi Ren, Blair Chen, Daniel Cer, Jeremy R Cole, Kai Hui, Michael Boratko, Rajvi Kapadia, Wen Ding, et al. 2024b. Gecko: Versatile text embeddings distilled from large language models. *arXiv preprint arXiv:2403.20327*.
- Chaofan Li, MingHao Qin, Shitao Xiao, Jianlyu Chen, Kun Luo, Yingxia Shao, Defu Lian, and Zheng Liu.
 2024. Making text embedders few-shot learners. *arXiv preprint arXiv:2409.15700*.

731

732

733

Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2024. Fine-tuning llama for multi-stage text retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2421– 2425.

674

675

686

697

698

703

705

706

709

710

711

712

714

717

718

719

720

721

723

724

725

726

727

729

730

- Niklas Muennighoff. 2022. Sgpt: Gpt sentence embeddings for semantic search. *Preprint*, arXiv:2202.08904.
- Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. Mteb: Massive text embedding benchmark. In *Proceedings of the 17th Conference* of the European Chapter of the Association for Computational Linguistics, pages 2014–2037.
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr, Felipe Petroski Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter Welinder, and Lilian Weng. 2022. Text and code embeddings by contrastive pre-training. *Preprint*, arXiv:2201.10005.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. 2022. Large dual encoders are generalizable retrievers. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 9844–9855, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019. Document expansion by query prediction. *arXiv preprint arXiv:1904.08375*.
- Yifu Qiu, Hongyu Li, Yingqi Qu, Ying Chen, QiaoQiao She, Jing Liu, Hua Wu, and Haifeng Wang. 2022.
 DuReader-retrieval: A large-scale Chinese benchmark for passage retrieval from web search engine. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5326–5338, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jacob Mitchell Springer, Suhas Kotha, Daniel Fried, Graham Neubig, and Aditi Raghunathan. 2024. Repetition improves language model embeddings. *arXiv preprint arXiv*:2402.15449.
- Saba Sturua, Isabelle Mohr, Mohammad Kalim Akram, Michael Günther, Bo Wang, Markus Krimmel, Feng Wang, Georgios Mastrapas, Andreas Koukounas, Nan Wang, and Han Xiao. 2024. jina-embeddingsv3: Multilingual embeddings with task lora. *Preprint*, arXiv:2409.10173.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In *Proceedings of the 2018*

Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.

- Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. 2023a. Label words are anchors: An information flow perspective for understanding in-context learning. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. Improving text embeddings with large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11897–11916, Bangkok, Thailand. Association for Computational Linguistics.
- Liang Wang, Nan Yang, and Furu Wei. 2023b. Query2doc: Query expansion with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9414–9423, Singapore. Association for Computational Linguistics.
- Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. 2024. C-pack: Packed resources for general chinese embeddings. In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24, page 641–649, New York, NY, USA. Association for Computing Machinery.
- Xiaohui Xie, Qian Dong, Bingning Wang, Feiyang Lv, Ting Yao, Weinan Gan, Zhijing Wu, Xiangsheng Li, Haitao Li, Yiqun Liu, and Jin Ma. 2023. T2ranking: A large-scale chinese benchmark for passage ranking. *Preprint*, arXiv:2304.03679.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.
- Xin Zhang, Zehan Li, Yanzhao Zhang, Dingkun Long, Pengjun Xie, Meishan Zhang, and Min Zhang. 2023a. Language models are universal embedders. *Preprint*, arXiv:2310.08232.
- Xinyu Zhang, Xueguang Ma, Peng Shi, and Jimmy Lin. 2021. Mr. TyDi: A multi-lingual benchmark for dense retrieval. In *Proceedings of the 1st Workshop* on Multilingual Representation Learning, pages 127– 137, Punta Cana, Dominican Republic. Association for Computational Linguistics.

787

789

790

793

794

797

799

801

804

810

811

812

817

821

823

824

830

Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. 2023b. MIRACL: A multilingual retrieval dataset covering 18 diverse languages. Transactions of the Association for Computational Linguistics, 11:1114–1131.

Α **Experimental Details for Data** Synthesis

A.1 Prompts for Data Synthesis

To enable large models to generate multiple hard negatives with varying granularities, we extend and refine the prompt template proposed by Wang et al. (2024). Specifically, we modify the original template to guide the model in producing negative samples with different levels of similarity to the query, thereby enhancing the diversity and difficulty of the generated data. Table 5 illustrates the complete prompt template used to generate short-long match tasks as an example.

Experimental Details for Supervised B Training

B.1 Hyperparameters

We present the hyperparameters involved in the supervised training in Table 4. The max sequence *length* specifies that any text sequence exceeding this number of tokens is truncated in our text embedding model.

B.2 Public Retrieval Datasets

We follow the training data setup from previous 815 work (Springer et al., 2024; BehnamGhader et al., 816 2024), adopting the dataset configuration used by Wang et al. (2024), which includes the following 818 datasets: ELI5 (Fan et al., 2019) (sample ratio 0.1) 819 , HotpotQA (Yang et al., 2018), FEVER (Thorne et al., 2018), MIRACL (Zhang et al., 2023b), MS-MARCO (Bajaj et al., 2016) passage ranking (sam-822 ple ratio 0.5) and document ranking (sample ratio 0.2), NQ (Karpukhin et al., 2020), NLI (Gao et al., 2021), SQuAD (Karpukhin et al., 2020), TriviaQA (Karpukhin et al., 2020), Quora Duplicate Ques-826 tions (DataCanary et al., 2017) (sample ratio 0.1), Mr-TyDi (Zhang et al., 2021), DuReader (Qiu et al., 2022), and T2Ranking (Xie et al., 2023) (sample ratio 0.5). The full supervised training data has approximately 1.5M training examples. The instructions applied to each dataset are in line with BehnamGhader et al. (2024), which are listed in Table 6. 834

Hyperparameter	Value
Batch Size	64
Gradient Accumulation Steps	8
Learning Rate	2e-5
Max Sequence Length	512
LoRA rank	16
LoRA α	32
Optimizer	Adam
Training Steps - Synthetic Only	600
Training Steps - Public Only	1000
Training Steps - Synthetic & Public	1600
Warmup Steps - Synthetic Only	200
Warmup Steps - Public Only	300
Warmup Steps - Synthetic & Public	300

Table 4: Hyperparameters used in the experiments

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

С **Experimental Details for Evaluation**

C.1 Prompts for MTEB Evaluation

For a fair comparison with previous work (Wang et al., 2024; BehnamGhader et al., 2024; Lee et al., 2024a) evaluated on MTEB, we adopted the same set of prompt instructions used in their evaluations when assessing our model's performance. The instructions applied to each evaluation dataset are listed in Table 7.

C.2 Subset Used for Ablation Study

To speed up evaluation in the ablation study, we follow Springer et al. (2024) by selecting a representative subset of the MTEB evaluation benchmark, which includes the following datasets: FiQA2018, SCIDOCS, SciFact, NF-Corpus, TwitterSemEval2015, TwitterURLCorpus, ImdbClassification, AmazonReviewsClassification, TweetSentimentExtractionClassification, MTOP-DomainClassification, TwentyNewsgroupsClustering, BiorxivClusteringS2S, MedrxivClusteringS2S, StackOverflowDupQuestions, AskUbuntuDupQuestions, SciDocsRR, BIOSSES, STS12, STS13, STS14, STS15, STS16, STS17, STS22, STSBenchmark, and SICK-R.

D **Additional Results**

D.1 Full MTEB Results

In this section, we present the complete results for all 56 MTEB datasets across the three experimental settings of our main experiment: public retrieval data only, synthetic data only, and full data. The corresponding results are shown in Table 8.

D.2 Full Attention Matrices

866

As shown in Figure 7, after transforming the at-867 tention mask of the base model (i.e. Mistral-7B-868 Instruct-v0.2) from causal to bidirectional, the at-869 tention heads continue to exhibit distinct patterns, 870 with some heads focus on tokens in their origi-871 nal positions, while others show higher attention 872 scores across all query dimensions at the current 873 position (e.g. Head 1, Head 6, Head 21 and Head 874 22). The latter pattern reflects the characteristics 875 of anchor tokens, allowing for more effective ag-876 gregation of information from the entire sentence. 877 Consequently, in our ATA pooling method, these 878 879 attention heads are assigned with greater pooling 880 weight.

Brainstorm a list of potentially useful text retrieval tasks.

Here are a few examples for your reference:

- Retrieve relevant documents for a short keyword web search query that asks for weather information.
- Search for documents that answers a FAQ-style query on children's nutrition.

Please adhere to the following guidelines:

- Specify what the query is, and what the desired documents are.
- Each retrieval task should cover a wide range of queries, and should not be too specific.

Your output must always be a python list of strings only, with about 20 elements, and each element corresponds to a distinct retrieval task in one sentence. Do not explain yourself or output anything else. Be creative! You have been assigned a retrieval task: {task}

Your mission is to write one text retrieval example for this task in the following JSON format. The JSON object must contain the following keys:

- "user_query": a string, a random user search query specified by the retrieval task.

- "positive_document": a string, a relevant document for the user query.

- "hard_negative_document": a list of strings, hard negative documents that only appears relevant to the query.

The output should be formatted as a JSON object with a field indicating the relative similarity level of hard negative examples. Use the following format as a guide:

```
{
  "user_query": "QUERY_TEXT".
  "positive document": "POSITIVE EXAMPLE TEXT",
  "hard_negative_document": [
    {
      "similarity_level": "high",
       "text": "HIGH_SIMILARITY_NEGATIVE_EXAMPLE_TEXT"
    },{
    "similarity_level": "medium",
    "STATE HIGH SITE
      "text": "MEDIUM_HIGH_SIMILARITY_NEGATIVE_EXAMPLE_TEXT"
    },{
    "similarity level": "medium".
    "text": "MEDIUM_LOW_SIMILARITY_NEGATIVE_EXAMPLE_TEXT"
    "text": "LOW_SIMILARITY_NEGATIVE_EXAMPLE_TEXT"
    }
  ]
}
```

Please adhere to the following guidelines:

- The "user_query" should be {query_type}, {query_length}, {clarity}, and diverse in topic.

- All documents must be created independent of the query. Avoid copying the query verbatim. It's acceptable if some parts of the "positive_document" are not topically related to the query.

- All documents should be at least {num_words} words long.

- The "hard_negative_document" contains some useful information, but it should be less useful or comprehensive compared to the "positive_document". Please generate four hard negative documents for contrastive learning based on the generated query and positive example. These examples should be arranged in order of decreasing similarity to the query, ranging from highly similar to dissimilar. Ensure the similarity spans a broad spectrum, and every negative example should be different, without repeating words from previous examples.

- Both the query and documents should be in {language}.

- Do not provide any explanation in any document on why it is relevant or not relevant to the query.

- Both the query and documents require {difficulty} level education to understand.

Your output must always be a JSON object only, do not explain yourself or output anything else. Be creative!

Table 5: Prompt template for the short-long matching task. For placeholders, "*{query_type}*" \in {extremely long-tail, long-tail, common}, "*{query_length}*" \in {less than 5 words, 5 to 15 words, at least 10 words}, "*{difficulty}*" \in {high school, college, PhD}, "*{clarity}*" \in {clear, understandable with some effort, ambiguous}, "*{num_words}*" \in {50, 100, 200, 300, 400, 500}.

Task Name	Instruction
NLI	Given a premise, retrieve a hypothesis that is entailed by the premise
	Retrieve semantically similar text
DuReader	Given a Chinese search query, retrieve web passages that answer the question
ELI5	Provided a user question, retrieve the highest voted answers on Reddit ELI5 forum
FEVER	Given a claim, retrieve documents that support or refute the claim
HotpotQA	Given a multi-hop question, retrieve documents that can help answer the question
MIRACL	Given a question, retrieve Wikipedia passages that answer the question
MrTyDi	Given a question, retrieve Wikipedia passages that answer the question
MSMARCO Passage	Given a web search query, retrieve relevant passages that answer the query
MSMARCO Document	Given a web search query, retrieve relevant documents that answer the query
NQ	Given a question, retrieve Wikipedia passages that answer the question
QuoraDuplicates	Given a question, retrieve questions that are semantically equivalent to the given
	question
	Find questions that have the same meaning as the input question
SQuAD	Retrieve Wikipedia passages that answer the question
T2Ranking	Given a Chinese search query, retrieve web passages that answer the question
TriviaQA	Retrieve Wikipedia passages that answer the question

Table 6: The prompt instructions used for public retrieval datasets, following BehnamGhader et al. (2024)

Task Name	Instruction
AmazonCounterfactualClassification	Classify a given Amazon customer review text as either counterfactual or not-
	counterfactual
AmazonPolarityClassification	Classify Amazon reviews into positive or negative sentiment
AmazonReviewsClassification	Classify the given Amazon review into its appropriate rating category
Banking77Classification	Given a online banking query, find the corresponding intents
EmotionClassification	Classify the emotion expressed in the given Twitter message into one of the six
	emotions: anger, fear, joy, love, sadness, and surprise
ImdbClassification	Classify the sentiment expressed in the given movie review text from the IMDB
	dataset
MassiveIntentClassification	Given a user utterance as query, find the user intents
MassiveScenarioClassification	Given a user utterance as query, find the user scenarios
MTOPDomainClassification	Classify the intent domain of the given utterance in task-oriented conversation
MTOPIntentClassification	Classify the intent of the given utterance in task-oriented conversation
ToxicConversationsClassif.	Classify the given comments as either toxic or not toxic
TweetSentimentClassification	Classify the sentiment of a given tweet as either positive, negative, or neutral
ArxivClusteringP2P	Identify the main and secondary category of Arxiv papers based on the titles and abstracts
ArxivClusteringS2S	Identify the main and secondary category of Arxiv papers based on the titles
BiorxivClusteringP2P	Identify the main category of Biorxiv papers based on the titles and abstracts
BiorxivClusteringS2S	Identify the main category of Biorxiv papers based on the titles
MedrxivClusteringP2P	Identify the main category of Medrxiv papers based on the titles and abstracts
MedrxivClusteringS2S	Identify the main category of Medrxiv papers based on the titles
RedditClustering	Identify the topic or theme of Reddit posts based on the titles
RedditClusteringP2P	Identify the topic or theme of Reddit posts based on the titles and posts
StackExchangeClustering	Identify the topic or theme of StackExchange posts based on the titles
StackExchangeClusteringP2P	Identify the topic or theme of StackExchange posts based on the given paragraphs
TwentyNewsgroupsClustering	Identify the topic or theme of the given news articles
SprintDuplicateQuestions	Retrieve duplicate questions from Sprint forum
TwitterSemEval2015	Retrieve tweets that are semantically similar to the given tweet
TwitterURLCorpus	Retrieve tweets that are semantically similar to the given tweet
AskUbuntuDupQuestions	Retrieve duplicate questions from AskUbuntu forum
MindSmallReranking	Retrieve relevant news articles based on user browsing history
SciDocsRR	Given a title of a scientific paper, retrieve the titles of other relevant papers
StackOverflowDupQuestions	Retrieve duplicate questions from StackOverflow forum
ArguAna	Given a claim, find documents that refute the claim
ClimateFEVER	Given a claim about climate change, retrieve documents that support or refute the claim
COADupstackRetrieval	Given a question, retrieve detailed question descriptions from Stackexchange that
	are duplicates to the given question
DBPedia	Given a query, retrieve relevant entity descriptions from DBPedia
FEVER	Given a claim, retrieve documents that support or refute the claim
FiQA2018	Given a financial question, retrieve user replies that best answer the question
HotpotQA	Given a multi-hop question, retrieve documents that can help answer the question
MSMARCO	Given a web search query, retrieve relevant passages that answer the query
NFCorpus	Given a question, retrieve relevant documents that best answer the question
NQ	Given a question, retrieve Wikipedia passages that answer the question
QuoraRetrieval	Given a question, retrieve questions that are semantically equivalent to the given
SCIDOCS	question
SCIDUCS SciEnat	Given a scientific paper title, retrieve paper abstracts that are cited by the given paper
SCIFACI Touche 2020	Given a scientific claim, refrieve documents that support or refute the claim
	Given a question, retrieve detailed and persuasive arguments that answer the question Given a query on COVID 10, retrieve documents that answer the query
TRECCUTID STS*	Diven a query on COVID-19, retrieve documents that answer the query Detrieve sementically similar text
SummEval	Reuleve semantically similar text.
Summervar	Given a news summary, reureve other semanticarry similar summaries

Table 7: The prompt instructions used for MTEB benchmark evaluation, following Wang et al. (2024). The "STS*" instruction applies to all STS tasks.

Dataset	Public Retrieval Data Only	Synthetic Data Only	Full Dataset
AmazonCounterfactualClassification	80.1	78.7	79.2
AmazonPolarityClassification	94.0	94.4	95.9
AmazonReviewsClassification	51.8	54.1	55.8
ArguAna	60.2	51.4	61.3
ArxivClusteringP2P	48.1	50.7	50.3
ArxivClusteringS2S	46.0	47.2	46.9
AskUbuntuDupQuestions	64.2	66.3	66.1
BIOSSES	85.6	85.1	87.5
Banking77Classification	88.5	88.3	89.2
BiorxivClusteringP2P	37.7	43.8	42.7
BiorxivClusteringS2S	36.9	41.4	41.2
CQADupstackRetrieval	48.8	44.3	47.1
ClimateFEVER	35.4	26.0	37.8
DBPedia	51.5	45.8	52.3
EmotionClassification	51.2	53.4	51.9
FEVER	91.2	79.1	89.4
FiQA2018	54.1	45.8	55.8
HotpotQA	77.6	57.9	75.9
ImdbClassification	90.3	93.4	94.2
MSMARCO	43.4	29.3	42.4
MTOPDomainClassification	96.3	95.7	96.6
MTOPIntentClassification	86.5	87.4	87.0
MassiveIntentClassification	80.1	80.6	80.3
MassiveScenarioClassification	82.1	81.8	82.4
MedrxivClusteringP2P	32.2	34.8	33.6
MedrxivClusteringS2S	32.5	35.4	34.8
MindSmallReranking	32.5	33.8	33.3
NFCorpus	39.4	37.9	38.5
NQ	65.9	57.7	66.9
QuoraRetrieval	89.5	86.0	89.1
RedditClustering	63.9	61.7	64.8
RedditClusteringP2P	66.8	64.1	67.3
SCIDOCS	22.0	23.7	22.7
SICK-R	83.5	80.4	83.8
STS12	76.6	75.4	79.8
STS13	86.8	86.6	88.3
STS14	83.1	82.4	85.6
STS15	88.5	88.6	91.3
STS16	85.9	86.6	88.1
STS17	91.7	87.0	91.9
STS22	67.9	66.5	69.7
STSBenchmark	87.9	84.4	89.7
SciDocsRR	84.4	85.7	84.7
SciFact	78.6	74.1	76.4
SprintDuplicateQuestions	95.3	94.7	95.3
StackExchangeClustering	72.9	71.3	72.6
StackExchangeClusteringP2P	37.1	43.5	42.9
StackOverflowDupQuestions	54.1	54.7	55.0
SummEval	31.1	31.6	31.3
TRECCOVID	81.4	81.3	82.2
Touche2020	23.3	26.6	24.6
ToxicConversationsClassification	66.9	69.8	68.8
TweetSentimentExtractionClassification	63.7	65.3	64.8
TwentyNewsgroupsClustering	53.4	53.2	53.8
TwitterSemEval2015	81.1	77.5	81.5
TwitterURLCorpus	87.1	86.3	87.0
Average	65.9	64.5	67.0

Table 8: Complete MTEB evaluation results for each dataset. Detailed evaluation metrics and dataset information are available in Muennighoff et al. (2023).



Figure 7: The individual attention matrices from all 32 attention heads in the last layer of the bidirectional model, obtained from a randomly selected example in the STS13 dataset.