# **Enhancing Lexicon-Based Text Embeddings with Large Language Models**

**Anonymous ACL submission** 

## Abstract

Recent large language models (LLMs) have demonstrated exceptional performance on general-purpose text embedding tasks. While dense embeddings have dominated related research, we introduce the first Lexicon-based EmbeddiNgS (LENS) leveraging LLMs that achieve competitive performance on these tasks. Regarding the inherent tokenization redundancy issue and unidirectional attention limitations in traditional causal LLMs, LENS consolidates the vocabulary space through token embedding clustering, and investigates bidirectional attention and various pooling strategies. Specifically, LENS simplifies lexicon matching by assigning each dimension to a specific token cluster, where semantically similar tokens are 016 grouped together, and unlocking the full po-017 tential of LLMs through bidirectional attention. Extensive experiments demonstrate that LENS outperforms dense embeddings on the Massive Text Embedding Benchmark (MTEB), delivering compact feature representations that match 022 the sizes of dense counterparts. Notably, combining LENS with dense embeddings achieves state-of-the-art performance on the retrieval subset of MTEB (i.e. BEIR).<sup>1</sup>

### 1 Introduction

034

038

Text embeddings are vector representations of text that power a wide range of applications, including retrieval, question answering, semantic textual similarity, and clustering. Recent advances in LLMs have shown that a single model can generate embeddings excelling across diverse tasks, highlighting their versatility (Li et al., 2024; BehnamGhader et al., 2024; Wang et al., 2023; Muennighoff et al., 2024; Meng et al., 2024; Lee et al., 2024a).

While dense embeddings that encode texts into low-dimensional, real-valued latent semantic spaces dominate recent research, lexicon-based



Figure 1: The redundancy and noise in LLM tokenizers, as well as the absence of bidirectional dependencies in causal LLMs motivate LENS.

040

041

042

044

045

046

048

051

052

054

056

057

060

061

062

063

embeddings (Formal et al., 2021b,a; Shen et al., 2023a; Lassance et al., 2024) offer distinctive advantages. These high-dimensional representations, where each dimension corresponds to a specific token of the vocabulary, align more closely with the pre-training objectives of language models due to their shared use of the vocabulary space and the language modeling head (Shen et al., 2023a). Recent studies have demonstrated that lexicon-based embeddings can surpass their dense counterparts, utilizing masked language models under specific control (Déjean et al., 2023). Additionally, lexiconbased embeddings can offer better transparency, providing clearer insights into the model's decisions via the weight of each token. Moreover, the combination of dense and lexicon-based embeddings has also been proven to be promising in prior studies, as they effectively complement each other (Lin, 2021; Shen et al., 2023b).

Despite these benefits, lexicon-based embeddings remain underexplored beyond retrieval tasks. To unlock their full potential in more scenarios, it is essential to address the challenges posed by LLMs, as shown in Fig. 1. The first one is the inherent

<sup>&</sup>lt;sup>1</sup>Our anonymous code is available at https://anonymous. 40pen.science/r/lens.

redundancy of LLM vocabularies. Since most modern tokenizers rely on subword tokenization (e.g., 065 "education" is split into "edu" and "cation"), 066 it fragments the entire vocabulary space (Soler et al., 2024). And semantically equivalent tokens can appear in multiple forms in the tokenizer (e.g., "what", "What", " what" and "review", "reviews"), introducing inconsistencies and difficulties in lexicon matching. Consequently, recent studies indicate that replacing the original tokenization of BM25 (Robertson et al., 1995) with the XLM-R tokenizer (Conneau et al., 2020) can lead to a significant performance drop due to the noisier vocabulary (Chen et al., 2024). The second 077 challenge is that LLMs typically employ unidirectional attention during pre-training, where tokens can only attend to preceding tokens. This limitation prevents each token from fully leveraging the surrounding context, which is crucial as lexicon-based embeddings are always derived from the outputs of all tokens.

> To address these challenges, we first explore the potential of LLMs generating embeddings where each dimension corresponds to a token cluster instead of the traditional single token, with each cluster grouping tokens that share similar meanings or stem from the same lexeme. To achieve this, we utilize a simple yet effective approach that directly clusters the token embeddings and leverages the centroids of these clusters as the new token embeddings for the language modeling head. As shown in Table 4, the resulting clusters naturally group tokens with similar meanings, forming more coherent and compact embeddings. At the meanwhile, these cluster-based embeddings can achieve the equivalent feature size as dense embeddings (e.g., 4,000d), which is much smaller than previous lexicon-based embeddings. Such a property not only i) facilitates the integration of LENS into existing dense frameworks like FAISS, freeing us from the sparsity constraints that, while essential for efficient retrieval, can limit expressiveness and effectiveness of models (Formal et al., 2024), but also ii) eliminates computational overhead in tasks such as clustering and classification, where inverted indices cannot be used.

094

100

102

104

105

107

108

109

110 111

112

114

115

Furthermore, to address the interior LLM architecture drawbacks, we also conduct extensive investigations into modifying the model frameworks. Given the recent studies highlight the significant 113 impact of attention mechanisms and pooling strategies on dense embeddings (Li et al., 2024; Muennighoff et al., 2024; BehnamGhader et al., 2024; Lee et al., 2024a), we incorporate variants of these two factors in our framework to examine how they affect lexicon-based embeddings. Contrary to prior findings (Li et al., 2024), which suggest that preserving the original architecture of LLMs typically vields optimal performance for dense embeddings, our results indicate that bidirectional attention is critical for achieving superior performance with lexicon-based embeddings.

116

117

118

119

120

121

122

123

124

125

126

127

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

156

157

158

159

160

161

162

163

164

165

Built on these techniques, we introduce LENS, a framework designed to generate low-dimensional lexicon-based embeddings that achieve impressive results across a variety of tasks. Specifically, our experiments demonstrate that LENS outperforms dense embeddings on the Massive Text Embedding Benchmark (MTEB) (Muennighoff et al., 2023), achieving state-of-the-art (SOTA) zero-shot performance among models trained exclusively on public data, as of December 1, 2024. Qualitative examples also illustrate that LENS produces grounded and meaningful representations. Further analysis demonstrates that LENS, even when using 2000 clusters, still outperforms embeddings that leverage the original vocabulary space. Moreover, combining LENS with dense embeddings achieves SOTA performance on the retrieval subset of MTEB (specifically, BEIR).

#### 2 **Related Work**

Lexicon-Based Embeddings. Lexicon-based embeddings assign each dimension of the embedding vector to a specific token in the vocabulary. With the advancements in masked language models, recent studies have demonstrated that lexicon-based embeddings (Mallia et al., 2021; Lin and Ma, 2021; Zhuang and Zuccon, 2021; Formal et al., 2021b,a; Shen et al., 2023a; Nguyen et al., 2023; Lassance et al., 2024) can deliver superior performance. Among these approaches, SPLADE (Formal et al., 2021b,a; Lassance et al., 2024) stands out as one of the most effective methods, often outperforming dense embeddings (Déjean et al., 2023). Moreover, lexicon-based embeddings have been shown to complement dense embeddings, with their combination yielding substantial performance improvements (Chen et al., 2024; Shen et al., 2023b; Lin, 2021). Despite these advances, research on lexicon-based embeddings has largely focused on retrieval tasks, leaving other applications relatively

193

194

195

196

197

198

199

201

203

204

210

211

212

214

such as clustering and classification underexplored.

LLM-Based Embeedings. As decoder-only LLMs continue to advance, recent work has in-168 vestigated their potential for generating dense text 169 embeddings capable of performing well across dif-170 ferent tasks. To align LLMs with text embedding 171 tasks, LLM2Vec (BehnamGhader et al., 2024) em-172 ploys masked next-token prediction training and un-173 supervised contrastive learning, while LLaRA (Li 174 et al., 2023a) leverages an auto-encoding objective 175 to enhance embedding quality. Recent efforts, such 176 as E5-Mistral (Wang et al., 2023) and Gecko (Lee 177 et al., 2024b), focus on improving embedding mod-178 els by using LLMs to generate diverse training data. 179 Additionally, GRIT (Muennighoff et al., 2024) explores the combination of contrastive learning and 181 182 language modeling objectives to train a single LLM that performs well on both embedding and generation tasks. Meanwhile, studies (Muennighoff et al., 184 2024; Lee et al., 2024a; BehnamGhader et al., 2024; Li et al., 2024) highlight the significant influence 186 of architectural choices on embedding model performance, with findings (Li et al., 2024) indicating that retaining the original unidirectional attention 189 often yields the best results. 190

Research on leveraging LLMs for lexicon-based embeddings remains limited. PromptReps (Zhuang et al., 2024) and Mistral-SPLADE (Doshi et al., 2024) use prompt engineering to generate lexiconbased embeddings from LLMs. However, these methods often perform worse than their dense counterparts, introduce additional computational overhead, and are limited to exploring only retrieval tasks.

# 3 Methodology

In this section, we first introduce preliminaries for a better understanding of the design of our framework, then formally describe the details of LENS.

# 3.1 Preliminaries

# 3.1.1 Lexicon-Based Embeddings Using Masked Language Models

SPLADE (Formal et al., 2021b,a; Lassance et al., 2024) is a representative method that utilizes Masked Language Models (MLMs) and regards the logits from the masked language modeling head as lexicon-based embeddings, leveraging the bidirectional attention. The MLM produces a sequence of logits  $L = (l_1, l_2, ..., l_n), l_i \in \mathbb{R}^{|V|}$  given the input sequence, where |V| is the vocabulary size. Each logit value  $l_{ij}$  represents the likelihood of the vocabulary token j being relevant to the position i. Specifically, these scores are produced by the language modeling head, which maps the output hidden states to the vocabulary space using the token embedding matrix.

To obtain the lexicon-based embeddings, SPLADE first applies a log-saturation transformation to the logits to scale the weight and enforce it as non-negative,

$$w_{ij} = \log\left(1 + ReLU(l_{ij})\right). \tag{1}$$

215

216

217

218

219

220

221

223

224

225

226

227

228

229

230

231

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

257

258

259

261

262

Then it performs max-pooling across logits of all tokens to derive the final weight for each vocabulary token,

$$w_j = \max_{i \in n} w_{ij}.$$
 (2)

Despite its proven effectiveness, former research on lexicon-based embeddings using MLMs primarily focused on small-scaled models, leaving the performance of larger models mostly unexplored.

# 3.1.2 Lexicon-Based Embeddings Using Causal Language Models

Motivated by the growing capability of largerscaled models, recent works have begun to use causal language models with significantly more parameters, such as LLaMA (Touvron et al., 2023) and Mistral (Jiang et al., 2023), to derive lexiconbased embeddings. Two notable methods are **PromptReps** (Zhuang et al., 2024) and **Mistral-SPLADE** (Doshi et al., 2024), which employ prompts to alleviate the limitations brought by the unidirectional attention.

**PromptReps** enables LLMs to generate both dense and lexicon-based embeddings through carefully designed prompts such as *"This sentence [IN-PUT] means in one word:"*. Dense embeddings are derived from the hidden states of the final token *"*, and lexicon-based ones are the logits for the next token prediction. Nevertheless, such a method relying solely on prompt causes a substantial performance drop of lexicon-based embeddings compared to their dense counterparts, e.g., MRR@10 of 34.15 vs. 41.86 on the MS MARCO dataset.

**Mistral-SPLADE** adapts SPLADE to large causal models like Mistral by using echo prompting (Springer et al., 2024). It enables full-context visibility of each token by duplicating the input sequence and regards the representations of the second occurrence as the output. Despite getting advancements on BEIR benchmark, it still lags behind dense embeddings like E5-Mistral (Wang et al., 2023) and LLM2Vec (BehnamGhader et al., 2024), demonstrating that lexicon-based embeddings using large models cannot solely rely on prompting.

Hence, LENS systematically investigate the architecture of LLMs, including attention mechanisms and pooling methods, rather than exterior prompting. We try to unlock the full potential of LLMs for lexicon-based embeddings, not only on retrieval tasks that have been widely examined before, but also on clustering and classification tasks which remains unexplored.

### 3.1.3 Tokenization in LLMs

263

264

265

269

270

271

272

273

275 276

277

279

284

286

289

290

291

293

296

297

301

303

LLM tokenizers, though designed to cover all possible text forms for the language modeling objective, may hinder the effectiveness of lexicon-based embedding. i) Extra redundancy can be introduced under the same lexeme and further affect the token matching. E.g., "What", "what", and " what" can be regarded as distinct tokens due to differences in case or whitespace, even though they represent the same word. ii) Subword fragmentation (Soler et al., 2024) split a common word into pieces like "education" into "edu" and "cation", posing additional matching complexity. iii) Tokenizers trained on large corpora often include rare tokens, which inflate the vocabulary size and make the embedding larger and slower to match.

Therefore, instead of directly using the original language modeling head, we simply cluster original tokens to form clusters and use their centroid embeddings to replace the original token embeddings of the language modeling head. This approach reduces the redundancy by merging related tokens and decreases the size of embeddings by using a smaller clustered vocabulary.

3.2 Framework of LENS

After discussing the background, we introduce the framework of our method, as shown in Fig. 2.

# 3.2.1 Architecture Design

Language Modeling Head. Motivated by the redundancy and noise in LLM tokenizer mentioned above, LENS assigns weights to groups of tokens with similar meanings, whose effectiveness has been verified in Zhang et al. (2024). Specifically, we apply KMeans clustering (Hartigan and Wong, 1979) to the token embeddings from the language



Figure 2: The model framework of LENS.

modeling head, where k is our desired lexiconbased embedding size. Then the original token embeddings in the LM head are replaced by the cluster centroids, while the input token embeddings remain unchanged. Such a substitution reduces the dimensionality of the lexicon-based embeddings, as the logits now represent scores over fewer clusters rather than the original huge vocabulary. Check Table 4 and Appendix A.1 for detailed cluster results.

Attention Mechanism. Given the former illustration on the limitations of unidirectional attention in typical causal LLMs, we emphasize it restricts the visibility of each token to the entire context. Hence, unlike previous works that rely on nonfundamental solutions like prompt engineering, we address this issue by directly modifying attention to be bidirectional during fine-tuning, which makes prompt design easier and inference more efficient.

### 3.2.2 Representation Generation

Following Wang et al. (2023) and Li et al. (2024), given a raw query-passage pair (q, p) for a specific embedding task, we first construct the instructed query input text as

$$q_{\text{ins}} = \langle \text{Instruct} \rangle \{ \text{task\_definition} \} \langle \text{query} \rangle \{ q \}.$$
(3)

Here *task\_definition* refers to the definition of the specific embedding task, guiding the model to adapt towards that task. On the other hand, the input of the passage part is solely the original text. Following Wang et al. (2023) and Li et al. (2024), a [EOS] token is also appended to the end of the sequence.

We then feed such an input into the modified LLM, and derive a series of logits vectors  $L = (l_1, l_2, ..., l_n), l_i \in \mathbb{R}^k$ , where n is the sequence

346

312

313

314

315

316

317

318

length and k is our clustering size. To obtain the final embeddings, following Formal et al. (2021a), log-saturation and max-pooling will be applied to L along the sequence dimension which is similar to Eq. 1, and 2.

It is also worth noting that we only employ tokens corresponding to the original query qto derive the output of the query, avoiding the noise brought by *task\_definition* tokens, inspired by BehnamGhader et al. (2024). Moreover, considering the autoregressive nature of LLMs that each logit is used for the prediction of the subsequent position, we shift the logits during pooling. In other words, we regard the logit corresponding to the neighboring on the left of each token as its feature during computation.

### 3.2.3 Training

347

348

349

356

360

361

363

364

366

370

374

376

377

384

Recent research has explored various complex methods for training embedding models. For example, NV-Embed-v2 (Lee et al., 2024a) employs a two-stage training pipeline while also incorporating positive-aware hard-negative mining and synthetic data generation. In contrast, for simplicity and fair comparison, the training of LENS strictly adheres to the training procedure of BGE-en-ICL (Li et al., 2024), an SOTA LLM-based dense embedding model. It uses a single-stage training process and relies exclusively on publicly available data.

Given a processed input pair  $(q_{ins}, p)$ , we utilize the InfoNCE loss as our objective,

$$\mathcal{L} = -\log \frac{\exp(\operatorname{sim}(q_{\operatorname{ins}}, p) / \tau)}{\exp(\frac{\operatorname{sim}(q_{\operatorname{ins}}, p)}{\tau}) + \sum_{j=1}^{N} \exp(\frac{\operatorname{sim}(q_{\operatorname{ins}}, p_j^-)}{\tau})}{(4)}$$

Here  $p_j^-$  and N denote the negative passage and number of negative passages, respectively. sim() is the cosine similarity function, defined as  $sim() = cos(h_{q_{ins}}, h_p)$ , where  $h_{q_{ins}} \in \mathbb{R}^k$  and  $h_p \in \mathbb{R}^k$  are the lexicon-based embeddings from the LLM for the instructed query and passage. The temperature  $\tau$  is set to 0.02 in our experiments.

### 4 Experiments

### 4.1 Setups

To ensure a fair comparison between dense embeddings and LENS, we strictly adhere to the training recipe of the SOTA dense model, BGE-en-ICL.

Model Setup. The Mistral-7B-v0.1 (Jiang et al.,
2023) model is used as the backbone in LENS,

in line with recent works such as BGE-en-392 ICL (Li et al., 2024), E5-Mistral (Wang et al., 393 2023), NV-Embed-v2 (Lee et al., 2024a), and 394 LLM2Vec (BehnamGhader et al., 2024). To in-395 vestigate the effect of different clustering sizes to 396 consolidate the output token embeddings, we set 397 k in KMeans clustering to 4,000 and 8,000 clus-398 ters, referred to as LENS-4000 and LENS-8000, 399 respectively. LENS-4000 can output 4000-d em-400 beddings, which is comparable to the 4096-d dense 401 embeddings produced by the same backbone LLM. 402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

**Training Data.** We directly utilize the publicly available training data provided by BGE-en-ICL. This dataset is a mixture of retrieval, reranking, clustering, classification, and semantic textual similarity (STS) tasks. Details about the training data can be found in Appendix A.2. We use the same set of task instructions as BGE-en-ICL, refer to Appendix A.3 for details.

Training Configurations. Following BGE-en-ICL, our model is trained for one epoch using LoRA (Hu et al., 2021), where the LoRA rank is 32 and the alpha is 64, and the learning rate is set to 1e-4. Each training sample is composed of 1 positive and 7 hard negatives. For retrieval tasks, we use a batch size of 512, whereas a batch size of 256 is used for the rest tasks. All data are drawn from the same dataset within the same batch. In retrieval tasks, we employ in-batch negatives and apply a KL-divergence loss to distill ranking scores from the BGE-reranker model<sup>2</sup>. The maximum length for both the query and passage is set to 512. It should be noted that we deviate from BGE-en-ICL by omitting in-context learning samples during training and concentrate on zero-shot scenarios solely. It enables us to exclusively evaluate LENS performance, free from extraneous signals.

**Evaluations.** We evaluate the performance of various embedding models using MTEB (Muennighoff et al., 2023) and AIR-Bench (Zeng et al., 2024). MTEB is a comprehensive text embedding benchmark encompassing seven task types across a total of 56 datasets. AIR-Bench, on the other hand, spans diverse domains for retrieval tasks, including law, healthcare, and books, having no overlap with MTEB. Notably, the ground truth for the test set in AIR-Bench is hidden, and we use the 24.04 version to assess the model's out-of-domain capabilities.

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/BAAI/bge-reranker-large

Task	#Dims	Retr.	Rerank.	Clust.	PairClass.	Class.	STS	Summ.	Avg.
# of datasets $\rightarrow$		15	4	11	3	12	10	1	56
Non-Fully Public Training Data									
E5-mistral-7b-instruct	4096	56.90	60.21	50.26	88.34	78.47	84.66	31.40	66.63
Linq-Embed-Mistral	4096	60.19	60.29	51.42	88.35	80.20	84.97	30.98	68.17
voyage-large-2-instruct	1024	58.28	60.09	53.35	89.24	81.49	84.31	30.84	68.23
stella_en_400M_v5	8192	58.97	60.16	56.70	87.74	86.67	84.22	31.66	70.11
gte-Qwen2-7B-instruct	3584	60.25	61.42	56.92	85.79	86.58	83.04	31.35	70.24
SFR-Embedding-2_R	4096	60.18	60.14	56.17	88.07	89.05	81.26	30.71	70.31
stella_en_1.5B_v5	8192	61.01	61.21	57.69	88.07	87.63	84.51	31.49	71.19
NV-Embed-v2	4096	62.65	60.65	58.46	88.67	90.37	84.31	30.70	72.31
		Fully	y Public Tr	aining D	ata				
LLM2Vec-Mistral-supervised	4096	55.99	58.42	45.54	87.99	76.63	84.09	29.96	64.80
GritLM-7B	4096	57.41	60.49	50.61	87.16	79.46	83.35	30.37	66.76
NV-Embed-v1	4096	59.36	60.59	52.80	86.91	87.35	82.84	31.20	69.32
bge-multilingual-gemma2	3584	59.24	59.72	54.65	85.84	88.08	83.88	31.20	69.88
BGE-en-ICL (zero-shot)	4096	<u>61.67</u>	59.66	57.51	86.93	88.62	83.74	30.75	71.24
LENS-4000 (Ours)	4000	60.76	<u>60.86</u>	<u>57.92</u>	87.93	88.13	<u>84.35</u>	31.56	71.22
LENS-8000 (Ours)	8000	61.86	60.91	58.02	87.98	88.43	84.67	29.54	71.63

Table 1: Top-performing models on the MTEB leaderboard as of December 1, 2024 compared to LENS. #Dims refers to the embedding dimensions. Abbreviations: Retr. = Retrieval; Rerank. = Reranking; Clust. = Clustering; PairClass. = Pair Classification; Class. = Classification; STS = Semantic Textual Similarity; Summ. = Summarization. The best and the second best results using public data are in **bold** and <u>underlined</u> font respectively.

Domain	#Dims	wiki	web	news	healthcare	law	finance	arxiv	msmarco	Avg.
# of datasets $\rightarrow$		1	1	1	1	1	1	1	1	8
E5-mistral-7b-instruct	4096	61.67	44.41	48.18	56.32	19.32	54.79	44.78	59.03	48.56
Linq-Embed-Mistral	4096	61.04	48.41	49.44	60.18	20.34	50.04	47.56	60.50	49.69
NV-Embed-v1	4096	62.84	50.42	51.46	58.53	20.65	49.89	46.10	60.27	50.02
gte-Qwen2-7B-instruct	3584	63.46	51.20	54.07	54.20	22.31	58.20	40.27	58.39	50.26
stella_en_1.5B_v5	8192	61.99	50.88	53.87	58.81	23.22	<u>57.26</u>	44.81	61.38	51.53
SFR-Embedding-Mistral	4096	63.46	51.27	52.21	58.76	23.27	56.94	47.75	58.99	51.58
NV-Embed-v2	4096	65.19	52.58	53.13	59.56	25.00	53.04	48.94	60.80	52.28
BGE-en-ICL (zero-shot)	4096	64.61	54.40	55.11	57.25	25.10	54.81	48.46	63.71	52.93
LENS-4000 (Ours)	4000	62.60	52.06	52.49	57.23	24.08	48.87	43.78	61.17	50.28
LENS-8000 (Ours)	8000	65.50	54.52	55.16	58.20	25.62	54.57	45.45	<u>63.00</u>	<u>52.75</u>

Table 2: QA performance on AIR-Bench 24.04 (English) across different models, where nDCG@10 is used as the metric. #Dims refers to the embedding dimensions. The best and the second best results across all models are in **bold** and <u>underlined</u> font respectively.

We compare LENS to numerous baselines, 440 including E5-mistral-7b-instruct (Wang et al., 441 2023), NV-Embed-v1/v2 (Lee et al., 2024a), 442 gte-Qwen2-7B-instruct (Li et al., 443 2023b), LLM2Vec (BehnamGhader al., 2024), et 444 SFR-Embedding-2\_R (Meng et al., 2024), 445 GritLM-7B (Muennighoff et al., 2024), and 446 BGE-en-ICL (Li et al., 2024). The results of 447 PromptReps and Mistral-Splade are excluded, as 448 they are designed specifically for retrieval tasks 449 and their performance falls below of the weakest 450 baseline, namely LLM2Vec-Mistral-supervised. 451 Some of the baselines use private data during 452 training or involve in-context learning. To ensure a 453 fair comparison, we focus on zero-shot scenarios 454 where no few-shot sample is included in the 455 prompt, e.g., BGE-en-ICL. 456

#### 4.2 Main results

MTEB. Table 1 demonstrates the results of a va-458 riety of models on MTEB. LENS-8000 achieves 459 the highest average performance among all mod-460 els trained on fully public data as of December 461 1, 2024. Notably, LENS-8000 outperforms BGE-462 en-ICL, its dense embedding counterpart trained 463 with the same data and hyperparameters. 6 among 464 7 categories of tasks also demonstrate consistent 465 superiority. Besides, LENS-4000 yields compara-466 ble performance as BGE-en-ICL, both share equiv-467 alent feature dimensions, but our lexicon-based 468 method can deliver better transparency. Further-469 more, LENS-8000 ranks second among all models 470 in overall average performance. The leading model, 471 NV-Embed-v2, attains its superiority through a sig-472 nificantly more complex training pipeline, which in-473

Text	Top-weighted clusters
most dependable affordable car	s(cars, Cars), (cheap, affordable), (reliable, reli), (depend, depends), (aff, afford)
fastest growing bonsai trees	(faster, fastest), (grow, growing), (fast, Fast), (tree, trees), (quickly, rapid)
causes of hypoxia in adults	(adult, adults), (oxygen, oxy), (cause, caused), (hyp, yp), (ox, 0x)
weather in lisbon april	(Portug, Portuguese), (bon, Bon), (weather, rather), (Spring, spring), (AP, #AP)
other hot flashes causes	(hot, Hot), (cause, causes), (flash, Flash), (flush, #flush), (heat, Heat)

Table 3: Qualitive examples of LENS-8000. For each example, the top-5 clusters with the largest weights in the embeddings are shown, with two tokens from each cluster included.

# Clusters

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

505

quickly, rapid, rapidly, swift
cannot, impossible, Unable, Cannot, Unable
shows, shown, showed, showing
review, Review, reviews, reviewed, Reviews
educ, education, Educ, Education, educational, Edu

Table 4: Cluster examples of LENS-8000. Each row presents tokens belonging to a single cluster. More cluster examples are provided in Appendix A.1.

cludes a two-stage training pipeline, positive-aware hard-negative mining, and synthetic data generation. By contrast, LENS uses fully public data and adopts a simpler training procedure.

**AIR-Bench.** We also evaluate LENS along with baselines on the QA tasks of AIR-Bench. As shown in Table 2, LENS-8000 outperforms the top-performing model on MTEB, NV-Embed-v2, demonstrating its promising generalization capabilities. Despite slightly lagging behind its dense counterpart BGE-en-ICL, LENS still remains competitive in several sub-tasks. However, LENS-4000 performs less competitively, potentially because a smaller number of clusters may result in overgeneralized clusters and information loss.

### 4.3 Qualitative Examples

We present some clustering results in Table 4. It can be found that: i) LENS groups semantically equivalent tokens (e.g., rapid and quickly, cannot and impossible); ii) it groups morphologically similar tokens (e.g., shows and showed); and iii) group uppercase/lowercase variants and wholeword/subword forms (e.g., review and Review, Edu and education). Such an observation proves the effect of clustering to eliminate the redundancy and noise of the tokenizer in some ways as we expected.

In addition, qualitative examples from MS MARCO of LENS-8000 embeddings are given in Table 3. The top-5 clusters with the largest weights in the embeddings are presented for each sample, where two tokens from each cluster are included. Obviously, these clusters are highly semantically relevant to the input texts, which can be regarded as some keywords. There are also interesting findings that the embeddings show some deep understanding of the text such as oxygen in response to the input "causes of hypoxia in adults", and some knowledge expansion capabilities like Portuguese and spring for the input "weather in lisbon april". These qualitative samples demonstrate that lexicon-based embeddings from LLMs captures more contextual features rather than some shallow token meanings. 506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534



Figure 3: Influence of the number of clusters. The configuration with 32,000 clusters retains the original token embeddings without clustering.

# 5 Analysis

In this section, we conduct a detailed investigation of LENS. For the sake of computational resources, we reduce the training data of each dataset to 10% of its original size in this part. Besides, we use the same MTEB subset as Jiang et al. (2024) for faster evaluation, as it correlates well with the overall performance of MTEB (details in Appendix A.4).

### 5.1 Influence of the Number of Clusters

We investigate how the number of clusters k affects the performance, as shown in Figure 3. The configuration with 32,000 clusters retains the original token embeddings without applying clustering. Our analysis reveals that decreasing the number of output entries from 32,000 tokens consistently improves performance, even when the number of clusters is reduced to as low as 2,000. However, it

Task # of datasets $\rightarrow$	Retr. 1	Rerank. 1	Clust. 1	PairClass. 1	Class. 1	STS 1	Summ. 1	Avg. 7			
Unidirectional Attention											
Last-token pooling	73.84	65.19	60.46	96.69	58.66	89.26	30.05	67.73			
Sum-pooling	72.46	59.57	50.55	89.90	54.64	80.55	29.70	62.48			
Max-pooling	75.18	59.68	50.93	92.06	57.58	82.74	30.89	64.15			
Bidirectional Attention											
Last-token pooling	76.89	64.21	61.57	96.62	58.33	88.72	30.72	68.15			
Sum-pooling	75.65	63.64	61.77	96.97	60.05	89.58	30.98	68.38			
Max-pooling	76.19	64.53	63.05	97.03	62.30	88.92	31.49	69.07			

Table 5: Influence of attention mechanisms and pooling methods.

Dataset	ARG	CLI	CQA	DBP	FEV	FIQ	HOT	MSM	NFC	NQ	QUO	SCD	SCF	TOU	COV	Avg.
BGE-en-ICL	82.76	45.35	47.23	50.42	91.96	58.77	84.98	46.72	40.69	73.85	91.02	25.25	78.33	29.67	78.11	61.67
LENS-8000 (Ours)	76.02	45.77	48.67	49.75	92.32	61.57	85.71	47.24	40.61	74.64	90.79	28.54	79.75	29.34	77.18	61.86
NV-Embed-v2	70.07	45.39	50.24	53.50	93.75	65.73	85.48	45.63	45.17	73.57	89.04	21.90	80.13	31.78	88.44	62.65
LENS (Ours) + BGE	81.37	47.14	48.57	51.79	93.12	62.00	87.12	47.66	41.55	75.81	91.07	28.41	80.19	30.51	78.72	63.00

Table 6: Results in terms of nDCG@10 on the retrieval subset (i.e. BEIR) of MTEB. We use the first three letters of each dataset's name as its abbreviation, except SCIDOCS (abbreviated as SCD) and SciFact (abbreviated as SCF). **Bolded values** indicate datasets where the combinations outperform both LENS-8000 and BGE-en-ICL individually. On 12 out of 15 datasets, combining LENS-8000 and BGE-en-ICL results in improved performance.

is crucial to maintain an adequate number of clusters to prevent information loss that can arise from overgeneralization. A configuration of 8,000 clusters strikes a good balance between effectiveness and efficiency (dimensionality). Consequently, we employ k = 8,000 in the subsequent experiments.

### 5.2 Influence of Model Architcure

535

536

539

540

541

543

544

545

546

547

549

550

551

553

555

We extensively investigate the effects of the attention mechanism and pooling methods, as illustrated in Table 5. For attention, we examine both unidirectional and bidirectional attention. Regarding pooling strategies, we assess max-pooling, sumpooling, and last-token pooling. The results highlight the critical role of bidirectional attention in achieving strong performance with lexicon-based embeddings, as evidenced by its superiority across all pooling methods. Among these pooling methods, max-pooling emerges as the most effective strategy. This finding partially explains the poor performance of lexicon-based embeddings from PromptReps, which relies on last-token pooling with unidirectional attention.

### 557 5.3 Hybrid Lexicon-Dense Embeddings

Previous studies have demonstrated that lexiconbased embeddings and dense embeddings are complementary, and combining them can lead to significant performance improvements. In this section,
we explore the effectiveness of combining LENS
with BGE-en-ICL, both trained on the same data
but representing different types of embeddings. To

evaluate general-use cases, we concatenate the two embeddings into a single embedding, without applying any additional operations. We hypothesize that enhanced performance could be achieved by tuning the combination weights of the two embeddings. 565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

The results are presented in Table 6. Combining LENS-8000 with BGE-en-ICL yields a substantial performance improvement, increasing from 61.67/61.86 to 63.00, which surpasses NV-Embedv2 and achieves SOTA results on the retrieval subset of MTEB as of December 1, 2024. Furthermore, such an improvement is consistent, as evidenced by performance gains on 12 out of 15 datasets.

# 6 Conclusion

In this work, we introduce LENS, a simple yet effective framework for generating lexicon-based text embeddings using LLMs. Our approach leverages token embedding clustering to address the redundancy challenges inherent in LLM tokenizers, while also enabling bidirectional attention to fully unlock the potential of LLMs. Extensive experiments demonstrate the promising effectiveness and generalization capabilities of LENS compared to SOTA dense embeddings. Qualitative examples reveal that LENS produces embeddings that are grounded and demonstrate a deep understanding of the input. Further analyses show the superiority of fusing lexicon-based LENS and dense embeddings, which surpasses each individual model on the retrieval subset of MTEB (i.e., BEIR).

# Limitations

596

610

611

612

614

615

616

617

619

621

623

626

633

634

635

637

638

639

641

642

647

We acknowledge the following limitations of our work. First, our training and evaluation are limited to English, leaving multilingual datasets, such as Miracl (Zhang et al., 2023), unexplored. This restricts the generalizability of our findings to non-English contexts. Second, we applied LENS exclusively to the widely used Mistral-7B model, leaving other models unexplored. Additionally, compared to previous lexicon-based models like SPLADE, utilizing LLMs as the backbone significantly increases computational costs.

# References

- C.J. Adams, Daniel Borkan, Jeffrey Sorensen, Lucas Dixon, Lucy Vasserman, and Nithum Thain. 2019. Jigsaw unintended bias in toxicity classification.
- Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. 2012. Semeval-2012 task 6: A pilot on semantic textual similarity. in\* sem 2012: The first joint conference on lexical and computational semantics–volume 1: Proceedings of the main conference and the shared task, and volume 2: Proceedings of the sixth international workshop on semantic evaluation (semeval 2012). Association for Computational Linguistics.
  - Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2018. 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*.
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. Efficient intent detection with dual sentence encoders. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 38–45.
- Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. *arXiv preprint arXiv:1708.00055*.
- Jianlyu Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. M3embedding: Multi-linguality, multi-functionality, multi-granularity text embeddings through selfknowledge distillation. In *Findings of the Association for Computational Linguistics: ACL 2024*.

Xi Chen, Ali Zeynali, Chico Camargo, Fabian Flöck, Devin Gaffney, Przemyslaw Grabowicz, Scott Hale, David Jurgens, and Mattia Samory. 2022. Semeval-2022 task 8: Multilingual news article similarity. In Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022), pages 1094– 1106. 649

650

651

652

653

654

655

656

657

658

659

660

661

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

- Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S Weld. 2020. Specter: Document-level representation learning using citation-informed transformers. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 2270–2282.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics.
- DataCanary, hilfialkaff, Lili Jiang, Meg Risdal, Nikhil Dandekar, and tomtung. 2017. Quora question pairs.
- Hervé Déjean, Stephane Clinchant, Carlos Lassance, Simon Lupart, and Thibault Formal. 2023. Benchmarking middle-trained language models for neural search. In Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '23.
- Meet Doshi, Vishwajeet Kumar, Rudra Murthy, Vignesh P, and Jaydeep Sen. 2024. Mistral-splade: Llms for better learned sparse retrieval. *arXiv preprint arXiv:2110.01529*.
- 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.
- Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2024. Towards effective and efficient sparse neural information retrieval. ACM Trans. Inf. Syst.
- Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2021a. Splade v2: Sparse lexical and expansion model for information retrieval. *arXiv preprint arXiv:2109.10086*.
- Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021b. Splade: Sparse lexical and expansion model for first stage ranking. In *Proceedings* of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '21.
- John A Hartigan and Manchek A Wong. 1979. Algorithm as 136: A k-means clustering algorithm. *Journal of the royal statistical society. series c (applied statistics)*, 28(1):100–108.

811

812

813

814

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

704

705

708

710

711

712

714

715

716

717

718

720

721

722

725

726

727

730

731 732

733

737

738

739

740

741

742

743

744

745

746

747

748

749

751

752

753 754

755

- 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.
- Albert Q. Jiang, Alicja Ziarko, Bartosz Piotrowski, Wenda Li, Mateja Jamnik, and Piotr Miłoś. 2024.
   Repurposing language models into embedding models: Finding the compute-optimal recipe. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems.*
  - Ken Lang. 1995. Newsweeder: Learning to filter netnews. In *Machine learning proceedings 1995*, pages 331–339. Elsevier.
  - Carlos Lassance, Hervé Déjean, Thibault Formal, and Stéphane Clinchant. 2024. Splade-v3: New baselines for splade. *arXiv preprint arXiv:2403.06789*.
  - 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, Yi Luan, Sai Meher Karthik Duddu, Gustavo Hernandez Abrego, Weiqiang Shi, Nithi Gupta, Aditya Kusupati, Prateek Jain, Siddhartha Reddy Jonnalagadda, Ming-Wei Chang, and Iftekhar Naim. 2024b. Gecko: Versatile text embeddings distilled from large language models. arXiv preprint arXiv:2403.20327.
- Chaofan Li, Zheng Liu, Shitao Xiao, and Yingxia Shao. 2023a. Making large language models a better foundation for dense retrieval. *arXiv preprint arXiv:2312.15503*.
- 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.
- Haoran Li, Abhinav Arora, Shuohui Chen, Anchit Gupta, Sonal Gupta, and Yashar Mehdad. 2021. Mtop: A comprehensive multilingual task-oriented semantic parsing benchmark. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 2950–2962.
- Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023b. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*.

- Jimmy Lin. 2021. A proposed conceptual framework for a representational approach to information retrieval. *arXiv preprint arXiv:2110.01529*.
- Jimmy Lin and Xueguang Ma. 2021. A few brief notes on deepimpact, coil, and a conceptual framework for information retrieval techniques. *arXiv preprint arXiv:2106.14807*.
- Xueqing Liu, Chi Wang, Yue Leng, and ChengXiang Zhai. 2018. Linkso: a dataset for learning to retrieve similar question answer pairs on software development forums. In *Proceedings of the 4th ACM SIG-SOFT International Workshop on NLP for Software Engineering*, pages 2–5.
- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies, pages 142–150.
- Wei Chen Maggie, Phil Culliton. 2020. Tweet sentiment extraction.
- Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. Www'18 open challenge: financial opinion mining and question answering. In *Companion proceedings of the the web conference* 2018, pages 1941–1942.
- Antonio Mallia, Omar Khattab, Torsten Suel, and Nicola Tonellotto. 2021. Learning passage impacts for inverted indexes. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 1723–1727.
- Julian McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: understanding rating dimensions with review text. In *Proceedings of the 7th ACM conference on Recommender systems*, pages 165–172.
- Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. 2024. Sfr-embedding-2: Advanced text embedding with multi-stage training.
- Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. 2024. Generative representational instruction tuning.
- 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.
- Thong Nguyen, Sean MacAvaney, and Andrew Yates. 2023. A unified framework for learned sparse retrieval. In Advances in Information Retrieval: 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2–6, 2023, Proceedings, Part III, pages 101–116. Springer.

900

901

902

903

904

905

906

907

908

870

871

872

- 815 816 817
- 818 819
- 82
- 82
- 82
- 8
- 8 8
- 8
- 831 832
- 8
- 834 835
- 8
- 838 839
- 841 842
- 843 844
- 84
- 846 847
- 8
- 850 851

852 853

854

- 856
- 857
- 8
- 8

8

864 865

1

- James O'Neill, Polina Rozenshtein, Ryuichi Kiryo, Motoko Kubota, and Danushka Bollegala. 2021. I wish i would have loved this one, but i didn't-a multilingual dataset for counterfactual detection in product review. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7092–7108.
- Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. 1995. Okapi at trec-3. *Nist Special Publication Sp.*
- Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. Carer: Contextualized affect representations for emotion recognition. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 3687–3697.
  - Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Xiaolong Huang, Binxing Jiao, Linjun Yang, and Daxin Jiang. 2023a. LexMAE: Lexicon-bottlenecked pretraining for large-scale retrieval. In *The Eleventh International Conference on Learning Representations*.
- Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Guodong Long, Kai Zhang, and Daxin Jiang. 2023b. Unifier: A unified retriever for large-scale retrieval. KDD '23.
- Aina Garí Soler, Matthieu Labeau, and Chloé Clavel. 2024. The impact of word splitting on the semantic content of contextualized word representations. *Transactions of the 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*.
- 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.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Henning Wachsmuth, Shahbaz Syed, and Benno Stein. 2018. Retrieval of the best counterargument without prior topic knowledge. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 241–251.

- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368*.
- 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.
- Yi Zeng, Yu Yang, Andy Zhou, Jeffrey Ziwei Tan, Yuheng Tu, Yifan Mai, Kevin Klyman, Minzhou Pan, Ruoxi Jia, Dawn Song, et al. 2024. Airbench 2024: A safety benchmark based on risk categories from regulations and policies. *arXiv preprint arXiv:2407.17436*.
- Xinyu Zhang, Jing Lu, Vinh Q. Tran, Tal Schuster, Donald Metzler, and Jimmy Lin. 2024. Tomato, tomahto, tomate: Measuring the role of shared semantics among subwords in multilingual language models. *arXiv preprint arXiv:2411.04530*.
- Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. 2023. Miracl: A multilingual retrieval dataset covering 18 diverse languages. *Transactions of the Association for Computational Linguistics*, 11:1114–1131.
- Shengyao Zhuang, Xueguang Ma, Bevan Koopman, Jimmy Lin, and Guido Zuccon. 2024. PromptReps: Prompting large language models to generate dense and sparse representations for zero-shot document retrieval. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*.
- Shengyao Zhuang and Guido Zuccon. 2021. Fast passage re-ranking with contextualized exact term matching and efficient passage expansion. *arXiv preprint arXiv:2108.08513*.

911

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

933

935

937

# A Appendix

# A.1 Clustering results

# Clusters

impact, Impact, impacts
Entity, entity, #Entity, Entities, entities
TV, television, tv, Television, televis
comfort, comfortable, comfort
beautiful, lovely, gorgeous, handsome, beautifully
guy, guys, Guy, dude
fit, FIT, fits, fitting, fitted
recomm, recommend, recommended, recommendation
star, stars, Stars
reach, reached, reaching, reaches, reach

Table 7: Cluster examples of LENS. Each row presents tokens belonging to a single cluster.

# A.2 Training data details

912We leverage the public training data provided by913BGE-en-ICL (Li et al., 2024). Specifically, the914training data is a mixture of retrieval, reranking,915classification, clustering, and STS data.

- **Retrieval**: ELI5 (Fan et al., 2019), HotpotQA (Yang et al., 2018), FEVER (Thorne et al., 2018), MSMARCO passage and document ranking (Bajaj et al., 2018), NQ, NLI, SQuAD, TriviaQA, Quora Duplicate Questions (DataCanary et al., 2017), Arguana (Wachsmuth et al., 2018), FiQA (Maia et al., 2018).
- **Reranking**: SciDocsRR (Cohan et al., 2020), StackOverFlowDupQuestions (Liu et al., 2018).
- AmazonReviews-• Classification: Classification (McAuley and Leskovec, 2013), AmazonCounterfactual-Classification (O'Neill al.. 2021). Banking77et Classification (Casanueva et al., 2020), Emotion-Classification (Saravia et al., 2018), TweetSentimentExtraction-Classification (Maggie, 2020), MTOPIntent-Classification (Li et al., 2021), IMDB-Classification (Maas et al., 2011), ToxicConversations-Classification (Adams et al., 2019).
- 938• Clustering:TwentyNewsgroups-939Clustering(Lang, 1995),940{Arxiv/Biorxiv/Medrxiv/Reddit/StackExchange}-941Clustering-{S2S/P2P}

• STS: STS12 (Agirre et al., 2012), STS22	942
(Chen et al., 2022), STS-Benchmark (Cer	943
et al., 2017).	944

945

946 947

948

949

950

951

952

953

954

# A.3 Task Instructions

We present the task instructions we used in Table 8.

# A.4 MTEB Subset Details

Following Jiang et al. (2024), for each task category, we select one dataset for evaluation. The chosen dataset is determined based on the model results presented in the original MTEB paper, focusing on the dataset with the highest correlation to the category's average performance.

• Classification: EmotionClassification 955 • Clustering: TwentyNewsgroupsClustering 956 • Pair classification: SprintDuplicateQuestions 957 • Reranking: AskUbuntuDupQuestions 958 • Retrieval: SciFact 959 • Semantic text similarity: STS15 960 • Summarization: SummEval 961 A.5 Detailed MTEB Results 962

We present the detailed MTEB results in Table 9. 963

Task Name	Instruction
ArguAna	Given a claim, find documents that refute the claim.
ClimateFEVER	Given a claim about climate change, retrieve documents that support or refute
CQADupStack	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.
Natural Question	Given a question, retrieve Wikipedia passages that answer the question.
QuoraRetrieval	Given a question, retrieve questions that are semantically equivalent to the given
SCIDOCS	Given a scientific paper title, retrieve paper abstracts that are cited by the given paper
SciFact	Given a scientific claim retrieve documents that support or refute the claim
Touche2020	Given a question, retrieve detailed and persuasive arguments that answer the question.
TREC-COVID	Given a query, retrieve documents that answer the query.
STS*	Retrieve semantically similar text.
SummEval	Given a news summary, retrieve other semantically similar summaries.
AmazonCounterfactualClassification	Classify a given Amazon customer review text as either 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.
	Classify the emotion expressed in the given Twitter message into one of the six
EmotionClassification	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.
ToxicConversationsClassification	Classify the given comments as either toxic or not toxic.
TweetSentimentExtractionClassification	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.
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.
SprintDuplicateQuestions	Retrieve duplicate questions from Sprint forum.
TwitterSemEval2015	Retrieve tweets that are semantically similar to the given tweet.
IwitterUKLCorpus	Retrieve tweets that are semantically similar to the given tweet.
AIR-Bench	Given a question, retrieve passages that answer the question.

Table 8: Task instructions for MTEB and AIR-Bench benchmarks.

Dataset $7B$ -instruct $dding-2_R$ $1.5B_v5^-$ (zero-shot) $bed-v2$ $-4000$ $-8000$ ArguAna $64.27$ $62.34$ $65.27$ $82.76$ $70.07$ $77.32$ $76.02$ ClimateFEVER $45.88$ $34.43$ $46.11$ $45.35$ $45.39$ $44.62$ $45.77$ CQADupStack $46.43$ $46.11$ $47.75$ $47.23$ $50.24$ $47.39$ $48.67$ DBPEDIA $52.42$ $51.21$ $52.28$ $50.42$ $53.50$ $50.10$ $49.75$ FEVER $95.11$ $92.16$ $94.83$ $91.96$ $93.75$ $92.37$ $92.32$ FiQA2018 $62.03$ $61.77$ $60.48$ $58.77$ $65.73$ $60.43$ $61.57$ HotpotQA $73.08$ $81.36$ $76.67$ $84.98$ $85.48$ $85.07$ $85.71$ MSMARCO $45.98$ $42.18$ $45.22$ $46.72$ $45.63$ $46.95$ $47.24$ NFCorpus $40.60$ $41.34$ $42.00$ $40.69$ $45.17$ $41.64$ $40.61$ Natural Question $67.00$ $73.96$ $71.80$ $73.85$ $73.57$ $73.13$ $74.64$ QuoraRetrieval $90.09$ $89.58$ $90.03$ $91.02$ $89.04$ $90.84$ $90.79$ SCIDOCS $28.91$ $24.87$ $26.64$ $25.25$ $21.90$ $27.51$ $28.54$ SciFact $79.06$ $85.91$ $80.09$ $78.33$ $80.13$ $77.51$ $28.54$ BIOSSES $81.37$ $87.60$ $83.11$ $8$	-	gte-Owen2-	SFR-Embe	stella en	BGE-en-ICL	NV-Em	LENS	LENS
ArguAna64.2762.3465.2782.7670.0777.3276.02ClimateFEVER45.8834.4346.1145.3545.3944.6245.77CQADupStack46.4346.1147.7547.2350.2447.3948.67DBPEDIA52.4251.2152.2850.4253.5050.1049.75FEVER95.1192.1694.8391.9693.7592.3792.32FiQA201862.0361.7760.4858.7765.7360.4361.57HotpotQA73.0881.3676.6784.9885.4885.0785.71MSMARCO45.9842.1845.2246.7245.6346.9547.24NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.287	Dataset	7B-instruct	dding-2_R	1.5B_v5	(zero-shot)	bed-v2	-4000	-8000
ClimateFEVER45.8834.4346.1145.3545.3944.6245.77CQADupStack46.4346.1147.7547.2350.2447.3948.67DBPEDIA52.4251.2152.2850.4253.5050.1049.75FEVER95.1192.1694.8391.9693.7592.3792.32FiQA201862.0361.7760.4858.7765.7360.4361.57HotpotQA73.0881.3676.6784.9885.4885.0785.71MSMARCO45.9842.1845.2246.7245.6346.9547.24NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.	ArguAna	64.27	62.34	65.27	82.76	70.07	77.32	76.02
CQADupStack46.4346.1147.7547.2350.2447.3948.67DBPEDIA52.4251.2152.2850.4253.5050.1049.75FEVER95.1192.1694.8391.9693.7592.3792.32FiQA201862.0361.7760.4858.7765.7360.4361.57HotpotQA73.0881.3676.6784.9885.4885.0785.71MSMARCO45.9842.1845.2246.7245.6346.9547.24NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1483.8779.9385.07	ClimateFEVER	45.88	34.43	46.11	45.35	45.39	44.62	45.77
DBPEDIA52.4251.2152.2850.4253.5050.1049.75FEVER95.1192.1694.8391.9693.7592.3792.32FiQA201862.0361.7760.4858.7765.7360.4361.57HotpotQA73.0881.3676.6784.9885.4885.0785.71MSMARCO45.9842.1845.2246.7245.6346.9547.24NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289	CQADupStack	46.43	46.11	47.75	47.23	50.24	47.39	48.67
FEVER95.1192.1694.8391.9693.7592.3792.32FiQA201862.0361.7760.4858.7765.7360.4361.57HotpotQA73.0881.3676.6784.9885.4885.0785.71MSMARCO45.9842.1845.2246.7245.6346.9547.24NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1388.8382.4089.6885.9888.3086.5487.34STS1588.5485.8289.3987.3589.0484.3284.39	DBPEDIA	52.42	51.21	52.28	50.42	53.50	50.10	49.75
FiQA201862.0361.7760.4858.7765.7360.4361.57HotpotQA73.0881.3676.6784.9885.4885.0785.71MSMARCO45.9842.1845.2246.7245.6346.9547.24NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	FEVER	95.11	92.16	94.83	91.96	93.75	92.37	92.32
HotpotQA73.0881.3676.6784.9885.4885.0785.71MSMARCO45.9842.1845.2246.7245.6346.9547.24NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1588.5485.5485.8289.3987.3589.0489.6989.75	FiQA2018	62.03	61.77	60.48	58.77	65.73	60.43	61.57
MSMARCO45.9842.1845.2246.7245.6346.9547.24NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	HotpotQA	73.08	81.36	76.67	84.98	85.48	85.07	85.71
NFCorpus40.6041.3442.0040.6945.1741.6440.61Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	MSMARCO	45.98	42.18	45.22	46.72	45.63	46.95	47.24
Natural Question67.0073.9671.8073.8573.5773.1374.64QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	NFCorpus	40.60	41.34	42.00	40.69	45.17	41.64	40.61
QuoraRetrieval90.0989.5890.0391.0289.0490.8490.79SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	Natural Question	67.00	73.96	71.80	73.85	73.57	73.13	74.64
SCIDOCS28.9124.8726.6425.2521.9027.5128.54SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	QuoraRetrieval	90.09	89.58	90.03	91.02	89.04	90.84	90.79
SciFact79.0685.9180.0978.3380.1378.3979.75Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	SCIDOCS	28.91	24.87	26.64	25.25	21.90	27.51	28.54
Touche202030.5728.1829.9429.6731.7825.8629.34TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	SciFact	79.06	85.91	80.09	78.33	80.13	78.39	79.75
TREC-COVID82.2687.2885.9878.1188.4469.7377.18BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	Touche2020	30.57	28.18	29.94	29.67	31.78	25.86	29.34
BIOSSES81.3787.6083.1186.3587.4284.4785.83SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	TREC-COVID	82.26	87.28	85.98	78.11	88.44	69.73	77.18
SICK-R79.2877.0182.8983.8782.1583.8183.30STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	BIOSSES	81.37	87.60	83.11	86.35	87.42	84.47	85.83
STS1279.5575.6780.0977.7377.8979.0780.99STS1388.8382.4089.6885.9888.3086.5487.34STS1483.8779.9385.0782.3484.3084.3284.39STS1588.5485.8289.3987.3589.0489.6989.75	SICK-R	79.28	77.01	82.89	83.87	82.15	83.81	83.30
STS13       88.83       82.40       89.68       85.98       88.30       86.54       87.34         STS14       83.87       79.93       85.07       82.34       84.30       84.32       84.39         STS15       88.54       85.82       89.39       87.35       89.04       89.69       89.75	STS12	79.55	75.67	80.09	77.73	77.89	79.07	80.99
STS14         83.87         79.93         85.07         82.34         84.30         84.32         84.39           STS15         88.54         85.82         89.39         87.35         89.04         89.69         89.75	STS13	88.83	82.40	89.68	85.98	88.30	86.54	87.34
STS15 88 54 85 82 89 39 87 35 89 04 89 69 89 75	STS14	83.87	79.93	85.07	82.34	84.30	84.32	84.39
	STS15	88 54	85.82	89 39	87.35	89.04	89.69	89.75
STS16 8649 8450 8715 8654 8677 8723 8763	STS16	86.49	84 50	87.15	86 54	86 77	87.23	87.63
STS17 88.73 88.93 91.35 91.25 90.67 91.55 90.87	STS17	88 73	88.93	91.35	91.25	90.67	91.55	90.87
STS22 66.8 67.10 68.10 68.08 68.12 68.60 68.00	STS17 STS22	66.88	67.10	68 10	68.08	68.12	68.60	68.00
STSRanchmark 86.5 83.60 88.23 87.02 88.47	STSBenchmark	86.85	83.60	88 23	87.02	88 /1	88 22	88 47
SummEval 31.35 30.71 31.40 30.75 30.70 31.55 20.54	SummEval	31 35	30.71	31 40	30.75	30.70	31.55	20.57
SummEval 51.55 50.71 51.47 50.75 50.70 51.55 27.54	Summe var	02.82	07.62	06.04	JU.75 05.06	07.02	06.09	29.34
Spinitzue are Evolutions 92.62 97.02 90.04 95.00 97.02 90.96 97.00 10.56 97.00 97.02 90.96 97.00 97.00 90.96 97.00 90.96 97.00 90.96 97.00 90.96 97.00 90.96 97.00 90.96 97.00 90.96 97.00	TwitterSemEval2015	92.02	97.02	90.04	95.00	97.02	90.90 70.21	97.00 70.56
Iwitter Semileval2013         //.90         /6.3/         80.36         /6.34         81.11         /9.31         /9.30           Twitter Semileval2013         //.90         /6.3/         80.36         /6.34         81.11         /9.31         /9.30           Twitter Semileval2013         //.90         /8.02         97.50         97.70         97.50         97.77         97.50         97.77         97.50         97.77         97.50         97.77         97.50         97.77         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.50         97.70         97.70         97.50         97.70	TwitterUDL Comus	77.90	/8.3/	00.30 07.50	/8.34	01.11 07.07	19.51	19.30
Iwitter/ORLCorpus         80.59         88.05         87.19         87.87         87.30         87.30           Among Country C	A manual Counterfeature	80.39	88.03	87.58	87.19	8/.8/	87.50	87.37
Amazon Ounterraciual 91.51 92.72 92.87 92.88 94.28 95.01 95.09	AmazonCounterlactual	91.51	92.72	92.87	92.88	94.28	93.01	93.09
AmazonPolarity 97.50 97.31 97.16 96.86 97.74 97.05 97.07	AmazonPolarity	97.50	97.31	97.16	96.86	97.74	97.05	97.07
Amazonkeviews         62.56         61.04         59.56         61.28         63.96         62.83         63.01	AmazonKeviews	62.56	61.04	59.36	61.28	63.96	62.83	03.01
Banking // 87.57 90.02 89.79 91.42 92.42 90.43 90.19	Banking //	87.57	90.02	89.79	91.42	92.42	90.43	90.19
Emotion 79.45 93.37 84.29 93.31 93.38 92.33 91.87	Emotion	79.45	93.37	84.29	93.31	93.38	92.33	91.87
Imdb 96.75 96.80 96.66 96.91 97.14 97.12 97.00	Imdb	96.75	96.80	96.66	96.91	97.14	97.12	97.00
Massivelntent 85.41 85.97 85.83 82.26 86.10 79.65 81.14	MassiveIntent	85.41	85.97	85.83	82.26	86.10	79.65	81.14
MassiveScenario 89.77 90.61 90.20 83.92 92.17 81.97 83.53	MassiveScenario	89.77	90.61	90.20	83.92	92.17	81.97	83.53
MTOPDomain 99.04 98.58 99.01 97.99 99.25 97.49 97.44	MTOPDomain	99.04	98.58	99.01	97.99	99.25	97.49	97.44
MTOPIntent 91.88 91.30 92.78 93.56 94.37 92.59 92.81	MTOPIntent	91.88	91.30	92.78	93.56	94.37	92.59	92.81
ToxicConversations         85.12         91.14         88.76         93.16         92.74         92.29         92.37	ToxicConversations	85.12	91.14	88.76	93.16	92.74	92.29	92.37
TweetSentimentExtraction         72.58         79.70         74.84         79.90         80.87         80.17         80.42	TweetSentimentExtraction	72.58	79.70	74.84	79.90	80.87	80.17	80.42
Arxiv-P2P 54.46 54.02 55.44 54.42 55.80 54.87 54.81	Arxiv-P2P	54.46	54.02	55.44	54.42	55.80	54.87	54.81
Arxiv-S2S         51.74         48.82         50.66         49.17         51.26         50.25         50.14	Arxiv-S2S	51.74	48.82	50.66	49.17	51.26	50.25	50.14
Biorxiv-P2P 50.09 50.76 50.68 52.32 54.09 52.39 52.48	Biorxiv-P2P	50.09	50.76	50.68	52.32	54.09	52.39	52.48
Biorxiv-S2S 46.65 46.57 46.87 48.38 49.60 48.35 48.52	Biorxiv-S2S	46.65	46.57	46.87	48.38	49.60	48.35	48.52
Medrxiv-P2P 46.23 46.66 46.87 46.13 46.09 46.35 46.38	Medrxiv-P2P	46.23	46.66	46.87	46.13	46.09	46.35	46.38
Medrxiv-S2S 44.13 44.18 44.65 44.20 44.86 44.54 44.89	Medrxiv-S2S	44.13	44.18	44.65	44.20	44.86	44.54	44.89
Reddit 73.55 62.92 72.86 71.20 71.10 72.32 72.37	Reddit	73.55	62.92	72.86	71.20	71.10	72.32	72.37
Reddit-P2P 74.13 72.74 75.27 72.17 74.94 73.20 73.89	Reddit-P2P	74.13	72.74	75.27	72.17	74.94	73.20	73.89
StackExchange79.8676.4880.2981.2982.1081.7081.60	StackExchange	79.86	76.48	80.29	81.29	82.10	81.70	81.60
StackExchange-P2P 49.41 48.29 49.57 45.53 48.36 43.73 44.41	StackExchange-P2P	49.41	48.29	49.57	45.53	48.36	43.73	44.41
TwentyNewsgroups 53.91 66.42 61.43 68.51 64.82 69.44 68.78	TwentyNewsgroups	53.91	66.42	61.43	68.51	64.82	69.44	68.78
AskUbuntuDupQuestions 67.58 66.71 67.33 64.80 67.46 65.45 65.74	AskUbuntuDupQuestions	67.58	66.71	67.33	64.80	67.46	65.45	65.74
MindSmallRerank 33.36 31.26 33.05 30.60 31.76 31.92 31.46	MindSmallRerank	33.36	31.26	33.05	30.60	31.76	31.92	31.46
SciDocsRR 89.09 87.29 89.20 86.90 87.59 87.92 87.63	SciDocsRR	89.09	87.29	89.20	86.90	87.59	87.92	87.63
StackOverflowDupOuestions 55.66 55.32 55.25 56.32 55.79 58.15 58.79	StackOverflowDunOuestions	55.66	55.32	55.25	56.32	55.79	58.15	58.79
MTEB Average (56) 70.24 70.31 71.19 71.24 72.31 71.22 71.63	MTEB Average (56)	70.24	70.31	71.19	71.24	72.31	71.22	71.63

Table 9: Detailed MTEB results.