

# GASE: Generatively Augmented Sentence Encoding

Anonymous ACL submission

## Abstract

We propose a training-free approach to improve sentence embeddings leveraging test-time compute by applying generative text models for data augmentation at inference time. Unlike conventional data augmentation that utilises synthetic training data, our approach does not require access to model parameters or the computational resources typically required for fine-tuning state-of-the-art models. Generatively Augmented Sentence Encoding variates the input text by paraphrasing, summarising, or extracting keywords, followed by pooling the original and synthetic embeddings. Experimental results on the Massive Text Embedding Benchmark for Semantic Textual Similarity (STS) demonstrate performance improvements across a range of embedding models using different generative models for augmentation. We find that generative augmentation leads to larger performance improvements for embedding models with lower baseline performance. These findings suggest that integrating generative augmentation at inference time adds semantic diversity and can enhance the robustness and generalisability of sentence embeddings for embedding models. Our results show that performance gains depend on the embedding model and the dataset.

## 1 Introduction

Representation learning has emerged as a fundamental technique in natural language processing (NLP). However, the quality and robustness of the embeddings are highly dependent on the richness and diversity of the training data. Recent advancements of generative Large Language Models (LLMs) led to remarkable capabilities in generating human-like text. LLMs were also used for data augmentation by Dai et al. (2023), who applied paraphrasing techniques to train a BERT model (Devlin et al., 2019). Wahle et al. (2022) showed that LLM-generated paraphrases are harder

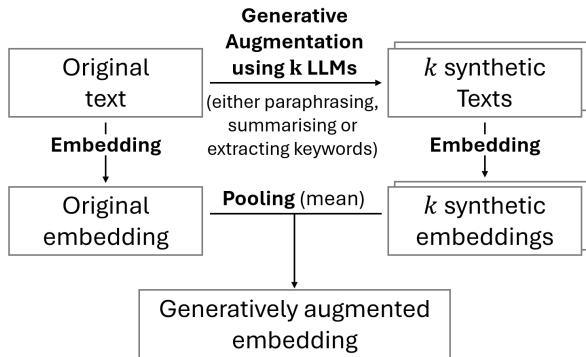


Figure 1: Approach for Generatively Augmented Sentence Encoding.

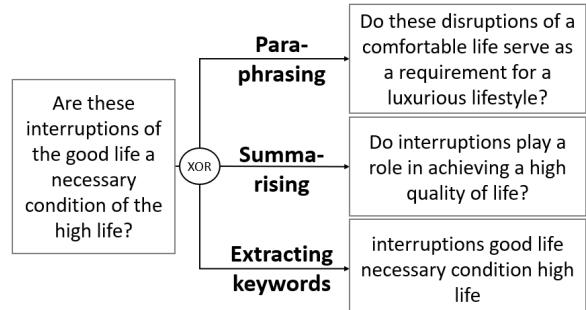


Figure 2: Augmentation examples for paraphrasing, summarising, and extracting keywords.

to detect for humans than paraphrases generated with simple techniques like synonym replacement. The integration of generative models with representation learning has been explored by augmenting generative models, e.g., in Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Guu et al., 2020).

In contrast, we introduce an approach to augment embedding models that applies generative models using test-time compute: Generatively Augmented Sentence Encoding (GASE). Instead of generating synthetic training data, GASE creates textual variants of an input text through *paraphrasing*, *summarising*, or *extracting keywords* at inference time. A joint embedding is derived by pooling the embed-

057 dings of the original text and the generated transformation(s). The underlying hypothesis of our work  
058 is that adding textual diversity using generative  
059 models increases the ability of current embedding  
060 models to model semantics and hence benefits the  
061 performance of downstream STS tasks.  
062

063 Also aiming to propose a training-free method,  
064 Lei et al. (2024) introduced Meta-Task Prompting,  
065 which employs use-case-specific prompts to generate  
066 multiple embeddings via a generative model, after  
067 which the embeddings are pooled. For retrieval  
068 systems, Gao et al. (2023) and Wang et al. (2023)  
069 described query expansion approaches using  
070 generative models that generate synthetic documents  
071 as a response to a retrieval request. Extending this  
072 work, GASE combines generative and embedding  
073 models without any assumptions on the type of  
074 models (unlike Meta-Task Prompting), nor is it limited  
075 to augmentation of retrieval queries such as  
076 Gao et al. (2023) and Wang et al. (2023). More  
077 recently, GenEOL (Thirukovalluru and Dhingra,  
078 2025), a related approach, was introduced (see Ta-  
079 ble 7 for a comparison with GASE).

## 080 2 Method

081 GASE performs the following steps (see Figure 1):

082 **Generative Augmentation.** We apply  $k$  genera-  
083 tive models to produce  $k$  variations of the original  
084 input sequence using exactly one of the following  
085 transformations: *Paraphrasing*, which provides a  
086 semantically equivalent but lexically or syntac-  
087 cally different variation of the input text. *Sum-  
088 marising*, which produces a shorter output text  
089 that captures the most important information of  
090 the input text. *Extracting Keywords*, which lists  
091 the most relevant words from a given text. Figure  
092 2 provides an example for each. We evaluated  
093  $k \in \{0, 1, 2, 3\}$ <sup>1</sup> using GPT-3.5 Turbo (OpenAI,  
094 2024b), Reka-Flash (Reka, 2024), and GPT-4o  
095 mini (OpenAI, 2024a)<sup>2</sup>.

096 **2. Sentence Embedding.** Generating embed-  
097 dings for the original and  $k$  synthetic texts with one  
098 of 12 encoders.<sup>3</sup>

099 **3. Pooling.** The embeddings of the original text  
100 and the  $k$  synthetic texts are pooled by computing  
101 their arithmetic mean.<sup>4</sup>

<sup>1</sup>Using a single generative model for  $k = 2$  did not yield  
textually diverse and meaningful variations.

<sup>2</sup>For model versions and hyperparameters see Appendix F.

<sup>3</sup>See Appendix D for a list incl. references.

<sup>4</sup>Mean pooling dominates max pooling (see Table 8).

We evaluate our approach on the English sub-tasks of the MTEB STS task<sup>5</sup> using cosine similarity and Spearman’s rank correlation<sup>6</sup>.

## 105 3 Results

106 Table 1 shows the results for different augmentation  
107 strategies vs. the respective baseline (see Table 9  
108 in the Appendix for the results per dataset). All em-  
109 bedding models improved their avg. scores through  
110 augmentation with the greatest improvements for  
111 Llama-3.1 (+6.49pp), BERT (+6.11pp), and GloVe  
112 (+2.96pp). Paraphrase augmentation worked best  
113 for most models; however, GloVe improved most  
114 with keyword extraction, while BERT, MPNET,  
115 and Llama-3.1 excelled with summarisation.<sup>7</sup>

116 Summarising and extracting keywords signif-  
117 icantly reduced the average word count across  
118 datasets (see Table 2). However, GPT-3.5 Turbo  
119 could not maintain STS22’s text length when para-  
120 phrasing and at the same time only effectively sum-  
121 marised STS22.<sup>8</sup>

122 Evaluating the different generative models for  
123 paraphrase augmentation we find that GPT-3.5  
124 Turbo performed best for all embedding models,  
125 except GloVe and Voyage-Large-2 (see Table 3 and  
126 for results per dataset Table 13 and Table 14 in the  
127 Appendix).

128 Using weighted averages between paraphrases  
129 and originals peaked at 25%/75% (excl. BERT)  
130 and was worst with paraphrases only (Figure 3).

131  $k > 1$  generative models yielded higher  
132 scores than a single model (see Table 3). BERT  
133 (+7.23pp), Llama-3.1 (+6.49pp) and GloVe  
134 (+4.11pp) showed the largest gains. Using more  
135 than one generative models yielded higher scores  
136 than a single model (see Table 3). BERT  
137 (+7.23pp), Llama-3.1 (+6.49pp) and GloVe  
138 (+4.11pp) showed the largest gains. Moreover, the  
139 overall performance of embedding models with  
140 lower non-augmented performance was associated  
141 with larger performance improvements through  
142 generative augmentation.

<sup>5</sup>See Appendix E for details on the datasets.

<sup>6</sup>See Reimers et al. (2016) for why Spearman’s rank cor-  
relation is preferable over Pearson correlation for STS tasks.

<sup>7</sup>Additionally, we compared LLM-based keyword extrac-  
tion with random keyword extraction (RKE) (see Appendix B  
for implementation details) and stopword removal showing  
that it generates semantically richer extractions (see Table 10).

<sup>8</sup>A word overlap analysis using the Jaccard Similarity (JS)  
can be found in Table 12 in the Appendix.

Embedding model	No augmentation		Paraphrasing		Summarising		Extracting keywords	
	Avg.	STS22	Avg.	STS22	Avg.	STS22	Avg.	STS22
glove.840B.300d	57.86	54.08	59.98	57.78	60.64	<b>61.21</b>	<b>60.82</b>	56.94
bert-large-cased	60.78	55.48	66.29	58.57	<b>66.89</b>	<b>63.47</b>	62.62	54.64
all-MiniLM-L6-v2	78.91	67.26	<b>80.32</b>	68.47	79.90	<b>69.27</b>	77.96	64.57
all-mpnet-base-v2	80.28	68.00	81.39	69.01	<b>81.48</b>	<b>70.12</b>	79.95	66.19
embed-english-light-v3.0	78.75	67.88	<b>80.30</b>	68.42	80.21	<b>68.97</b>	78.53	67.27
embed-english-v3.0	81.25	68.22	<b>82.10</b>	68.63	82.08	<b>70.29</b>	81.26	66.91
voyage-2	82.50	65.26	<b>83.38</b>	66.80	82.74	<b>68.50</b>	81.21	65.29
voyage-lite-02-instruct	85.95	78.62	<b>85.99</b>	79.08	85.43	<b>79.87</b>	83.91	76.60
voyage-large-2	83.60	63.97	<b>83.92</b>	64.69	83.66	<b>65.90</b>	81.98	63.68
voyage-large-2-instruct	84.61	66.59	<b>84.80</b>	67.72	84.54	<b>68.66</b>	83.40	65.73
mxbai-embed-large-v1	84.63	68.73	<b>85.01</b>	70.03	84.45	<b>70.36</b>	83.42	67.83
llama-3.1-8b	71.36	33.39	75.28	50.55	<b>76.74</b>	<b>59.44</b>	74.95	51.79

Table 1: Average scores across MTEB STS datasets and STS22 individually for all embedding models with and without generative augmentation using GPT-3.5 Turbo in %, bold: highest average scores for average and STS22 respectively).

Dataset	Original	Para-phrases	Summa-ries	Extracted keywords
STSB	10.1	11.5	10.7	5.7
STS12	11.1	11.6	11.1	6.1
STS13	9.0	10.5	11.1	5.4
STS14	9.3	11.0	11.4	5.7
STS15	10.6	12.0	11.5	5.7
STS16	11.6	12.6	12.8	5.7
STS17	8.7	9.6	9.0	4.5
STS22	477.2	216.5	58.4	103.8
SICK-R	9.6	10.0	9.3	4.7
BIOSSES	24.5	25.2	19.2	13.9
Average	58.2	33.1	16.4	16.1

Table 2: Mean word count (using GPT-3.5 Turbo for paraphrasing, summarising and extracting keywords).

### 3.1 Extensions

To evaluate the generalisability of GASE we extended our work as follows:

**Multilingual.** Consistent with prior findings, GASE applied to 2 multilingual datasets with 3 multilingual models yields greater improvements for weaker encoders and summarisation remains most effective for BERT and STS22. (Table 4).<sup>9</sup>

**Generative model ablation.** Ablating the generative model size shows that GASE’s performance correlates with generative model size (Table 5).<sup>10</sup>

**Additional task.** We evaluated GASE on the Pair Classification (PC) task (measured in average precision). In line with STS, improvements are negatively correlated with baseline performance and lower for STS22 (Table 6)<sup>11</sup>.

<sup>9</sup>See Table 15 and Table 16 for results per dataset and Appendix D for model properties.

<sup>10</sup>See Appendix F for implementation details.

<sup>11</sup>See Table 17 for scores per dataset and (Muennighoff

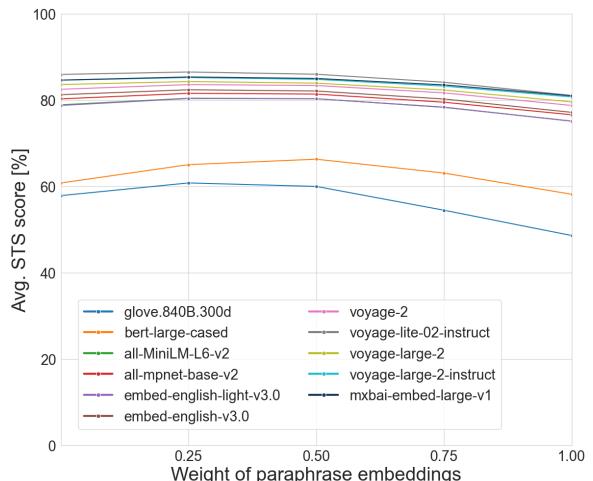


Figure 3: Scores for different weights for a weighted average between embeddings of paraphrases generated with GPT-3.5-Turbo and original texts ( $k = 1$ ).

## 4 Discussion

We believe that embedding models with lower baseline performance benefited more from the semantic diversity induced by augmentation due to their inherent lower capability to model diverse semantics. Similarly, they gained more from the ensemble effect using with  $k \geq 2$  generative models.

The effectiveness of keyword extraction for GloVe, BERT, and Llama-3.1 may be explained by the model’s limited capability to handle complex semantics. Therefore, the reduction to keywords might help the model to reduce noise.

Longer texts likely caused summary augmentation to outperform paraphrase augmentation on

et al., 2023) for details for score calculations.

Embedding model	No Augmentation	gpt-3.5-turbo	reka-flash	gpt-4o-mini	gpt-3.5-turbo + reka-flash	gpt-3.5-turbo + reka-flash + gpt-4o-mini
glove.840B.300d	57.86	59.98	<u>60.73</u>	60.04	61.97	<b>62.18</b>
bert-large-cased	60.78	<b>66.29</b>	65.76	65.01	68.01	<b>68.41</b>
all-MiniLM-L6-v2	78.91	<u>80.32</u>	79.74	79.81	80.64	<b>80.76</b>
all-mpnet-base-v2	80.28	<u>81.39</u>	81.01	81.07	81.72	<b>81.85</b>
embed-english-light-v3.0	78.75	<u>80.30</u>	79.72	79.90	80.54	<b>80.66</b>
embed-english-v3.0	81.25	<b>82.10</b>	82.00	81.82	<b>82.53</b>	82.08
voyage-2	82.50	<u>83.38</u>	83.37	83.07	<b>83.80</b>	83.72
voyage-lite-02-instruct	85.95	<u>85.99</u>	85.92	85.97	<b>86.05</b>	86.01
voyage-large-2	83.60	<u>83.92</u>	<u>83.98</u>	83.95	84.12	<b>84.17</b>
voyage-large-2-instruct	84.61	<u>84.80</u>	84.68	84.78	84.86	<b>84.91</b>
mxbai-embed-large-v1	84.63	<u>85.01</u>	84.80	84.90	<b>85.17</b>	85.16
llama-3.1-8b	71.36	75.28	<u>77.02</u>	74.71	<b>77.85</b>	77.52

Table 3: Average STS scores without and with paraphrase augmentation using different generative models in % (bold: highest score per emb. model, underlined: highest score for  $k \leq 1$  per emb. model).

Embedding model	Augment.	STS17	STS22
bert-base-multilingual-cased	None	37.16	29.50
	Paraphrase	43.93	30.36
	Summarise	<b>47.13</b>	<b>42.03</b>
paraphrase-multilingual-mpnet-base-v2	None	82.51	61.56
	Paraphrase	<b>82.95</b>	62.72
	Summarise	82.88	<b>66.70</b>
multilingual-e5-large-instruct	None	83.39	69.64
	Paraphrase	<b>83.95</b>	66.48
	Summarise	83.60	<b>70.26</b>
voyage-multilingual-2	None	85.14	69.88
	Paraphrase	<b>85.73</b>	70.21
	Summarise	85.41	<b>71.92</b>

Table 4: Avg. STS scores for multi-lingual datasets with and without augmentation with GPT-4o mini in % (bold: highest scores per embedding model and dataset).

Generative Model Size (Qwen2.5)					
None	0.5b	1.5b	3b	7b	14b
80.28	79.2	79.79	80.26	80.33	80.79

Table 5: Avg. STS scores of all-mpnet-base-v2 without and with paraphrase augmentation using Qwen2.5 models with different number of parameters (in %).

173 STS22, as summarising significantly reduced the  
 174 length. On other datasets, however, summaries  
 175 were only slightly shorter, diminishing its augmen-  
 176 tation value. Our results show pooling original and  
 177 synthetic embeddings outperforms individual use  
 178 ([Figure 3](#)). Across datasets, tasks, and languages,  
 179 GASE particularly benefits embedding models with  
 180 lower baseline performance and shorter texts, while  
 181 summary augmentation is more effective for longer  
 182 texts. Moreover, GASE requires a certain size of  
 183 the generative model to be effective ([Table 5](#)).

184 Since GenEOL used different, mostly weaker,  
 185 embedding models, our results cannot be directly

compared with their results.

Embedding model	Augmentation	PC Average
glove.840B.300d	None	58.41
	Paraphrase	<b>64.64</b>
bert-large-cased	None	55.07
	Paraphrase	<b>62.36</b>
all-MiniLM-L6-v2	None	82.34
	Paraphrase	<b>83.80</b>
all-mpnet-base-v2	None	83.00
	Paraphrase	<b>84.61</b>
embed-english-light-v3.0	None	80.18
	Paraphrase	<b>82.88</b>
embed-english-v3.0	None	83.54
	Paraphrase	<b>85.12</b>
voyage-2	None	83.81
	Paraphrase	<b>85.24</b>
voyage-lite-02-instruct	None	<b>93.08</b>
	Paraphrase	92.22
voyage-large-2	None	<b>92.96</b>
	Paraphrase	92.36
voyage-large-2-instruct	None	85.21
	Paraphrase	<b>86.22</b>
mxbai-embed-large-v1	None	87.16
	Paraphrase	<b>87.38</b>
llama-3.1-8B	None	70.10
	Paraphrase	<b>72.78</b>

Table 6: Average Pair Classification (PC) scores with and without paraphrase augmentation in % (bold: highest scores per embedding model).

## 5 Conclusion

We propose GASE, which augments sentence encoders at inference with generative models for paraphrasing, summarising, or keyword extraction. On MTEB STS, GASE significantly improves lower-performing embedding models and incrementally enhances SOTA models. It applies broadly, needing only embedding and generative model outputs.

## 6 Limitations

This study presents several limitations:

1. Our approach significantly increases the computational cost for sentence encoding since each generative augmentation requires an additional inference by the embedding model plus inference using a generative LLM. For a basic estimate of the runtime increase based on STS17 see [Table 18](#). The results highlight GASE’s particular suitability for smaller embedding models, where it yields significant performance gains with minimal encoder runtime.
2. While deterministic embedding models such as GloVe yield consistent results, other models like GPT-3.5 Turbo examined in this research exhibit stochastic behaviour. Consequently, additional experiments are necessary to corroborate our findings, as LLMs can produce diverse outcomes, thereby limiting the conclusiveness of a single experimental run ([Reimers and Gurevych, 2018](#)). This caveat is particularly pertinent to observations based on marginal performance differentials.
3. Our experiments demonstrate that the effectiveness of GASE is limited to embedding models with lower baseline performance while encoders models with higher baseline performance exhibit minimal to no improvement. Further investigation is needed to understand the underlying reasons for this observation.
4. Our investigation was confined to evaluating GPT-3.5 Turbo, Reka-Flash, and GPT-4o mini for generative augmentation. Alternative LLMs, including the Claude and Gemini model families, may lead to different results.
5. Furthermore, our experimental design was restricted to  $k \in \{0, 1, 2, 3\}$ . While values of  $k > 3$  were not explored, we posit that such configurations may be impractical for real-world applications due to prohibitive computational costs.
6. Our experiments for generative model ablation are limited to models with a parameter count  $\in \{0.5, 1.5, 3, 7, 14\}$ . Extending the analysis to larger generative models would

yield insights into whether the observed pattern holds for models with more than 14b parameters.

7. Summaries typically comprise 50% or fewer words relative to its source text ([Radev et al., 2002](#)). However, due to the concise nature of the texts within the examined STS datasets the summaries generated in this work approximately maintain the original word count. (Except for STS22 which contains longer text sequences.) Hence, the efficacy of augmentation through summarisation may be more pronounced when applied to datasets comprising longer textual inputs (e.g., paragraphs) compared to the predominantly short text sequences examined within the scope of this study.

## Acknowledgments

Code development for this work has been assisted by GPT-3.5 Turbo, GPT-4 Omni, Claude 3.5 Sonnet, Gemini Pro 1.5, and Github Copilot.

## References

- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Iñigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, German Rigau, Larraitz Uria, and Janyce Wiebe. 2015. [SemEval-2015 Task 2: Semantic Textual Similarity, English, Spanish and Pilot on Interpretability](#). In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 252–263, Denver, Colorado. Association for Computational Linguistics.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. 2014. [SemEval-2014 Task 10: Multilingual Semantic Textual Similarity](#). In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 81–91, Dublin, Ireland. Association for Computational Linguistics.
- Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Rada Mihalcea, German Rigau, and Janyce Wiebe. 2016. [SemEval-2016 Task 1: Semantic Textual Similarity, Monolingual and Cross-Lingual Evaluation](#). In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 497–511, San Diego, California. Association for Computational Linguistics.
- Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. 2012. SemEval-2012 Task 6: A pilot on semantic textual similarity. In *\*SEM 2012 -*

293	<i>1st Joint Conference on Lexical and Computational Semantics</i> , volume 2, pages 385–393.	349
294		350
295	Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. 2013. *SEM 2013 shared task: Semantic Textual Similarity. In <i>Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity</i> , pages 32–43, Atlanta, Georgia, USA. Association for Computational Linguistics.	351
296		352
297		353
298		354
299		
300		
301		
302		
303	Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017a. <i>SemEval-2017 Task 1: Semantic Textual Similarity Multilingual and Crosslingual Focused Evaluation</i> . In <i>Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)</i> , pages 1–14, Vancouver, Canada. Association for Computational Linguistics.	355
304		356
305		357
306		358
307		
308		
309		
310		
311		
312		
313		
314		
315		
316		
317	Xi Chen, Ali Zeynali, Chico Camargo, Fabian Flöck, Devin Gaffney, Przemysław Grabowicz, Scott Hale, David Jurgens, and Mattia Samorøy. 2022. <i>SemEval-2022 Task 8: Multilingual news article similarity</i> . In <i>Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)</i> , pages 1094–1106, Seattle, United States. Association for Computational Linguistics.	359
318		360
319		361
320		362
321		363
322		364
323		365
324		366
325	Cohere. 2024. Documentation: Model Overview. <a href="https://docs.cohere.com/docs/cohere-embed">https://docs.cohere.com/docs/cohere-embed</a> (accessed 2024-08-08).	367
326		368
327		369
328	Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke Huang, Zihao Wu, Lin Zhao, Wei Liu, Ninghao Liu, Sheng Li, Daqiang Zhu, Hongmin Cai, Quanzheng Li, Dinggang Shen, Tianming Liu, and Xiang Li. 2023. <i>ChatAug: Leveraging ChatGPT for Text Data Augmentation</i> .	370
329		371
330		372
331		373
332		374
333		375
334	Jacob Devlin, Ming Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. <i>BERT: Pre-training of deep bidirectional transformers for language understanding</i> . <i>NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference</i> , 1(Mlm):4171–4186.	376
335		377
336		378
337		379
338		380
339		
340		
341		
342	Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023. <i>Precise Zero-Shot Dense Retrieval without Relevance Labels</i> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1762–1777, Toronto, Canada. Association for Computational Linguistics.	381
343		382
344		383
345		
346		
347		
348		
510	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. <i>REALM: Retrieval-augmented language model pre-training</i> . In <i>Proceedings of the 37th International Conference on Machine Learning</i> , volume 119 of <i>ICML’20</i> , pages 3929–3938. JMLR.org.	384
511		385
512		386
513		387
514		
515	{Hugging Face}. 2024. MTEB Leaderboard - a Hugging Face Space by mteb. <a href="https://huggingface.co/spaces/mteb/leaderboard">url{https://huggingface.co/spaces/mteb/leaderboard}</a> (accessed 2024-09-06).	388
516		389
517		390
518	Yibin Lei, Di Wu, Tianyi Zhou, Tao Shen, Yu Cao, Chongyang Tao, and Andrew Yates. 2024. <i>Meta-Task Prompting Elicits Embeddings from Large Language Models</i> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10141–10157, Bangkok, Thailand. Association for Computational Linguistics.	391
519		392
520		393
521		394
522		
523	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kütller, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In <i>Proceedings of the 34th International Conference on Neural Information Processing Systems</i> , NIPS ’20, pages 9459–9474, Red Hook, NY, USA. Curran Associates Inc.	395
524		396
525		397
526		398
527		399
528	M. Marelli, Stefano Menini, Marco Baroni, L. Bentivogli, R. Bernardi, and Roberto Zamparelli. 2014. A SICK cure for the evaluation of compositional distributional semantic models. In <i>International Conference on Language Resources and Evaluation</i> .	400
529		401
530		402
531		403
532		
533	Meta. The Llama 3 Herd of Models   Research - AI at Meta. <a href="https://ai.meta.com/research/publications/the-llama-3-herd-of-models/">https://ai.meta.com/research/publications/the-llama-3-herd-of-models/</a> .	404
534		405
535		406
536		407
537		
538	Mixedbread. 2024. Documentation: Mxbai-embed-large-v1. <a href="https://www.mixedbread.ai/embeddings/mxbai-embed-large-v1">url{https://www.mixedbread.ai/embeddings/mxbai-embed-large-v1}</a> (accessed 2024-08-08).	408
539		409
540		410
541		411
542	Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2023. <i>MTEB: Massive Text Embedding Benchmark</i> . <i>Preprint</i> , arXiv:2210.07316.	412
543		413
544		414
545		415
546		
547	OpenAI. 2024a. GPT-4o mini: Advancing cost-efficient intelligence. <a href="https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/">url{https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/}</a> (accessed 2024-07-24).	416
548		417
549		418
550		419
551	OpenAI. 2024b. OpenAI - GPT-3.5 Turbo. <a href="https://platform.openai.com/docs/models/gpt-3-5-turbo">https://platform.openai.com/docs/models/gpt-3-5-turbo</a> (accessed 2024-10-08).	420
552		421
553		422
554		423
555		424
556		425
557		426
558		427
559		428
560		429
561		430
562		431
563		432
564		433
565		434
566		435
567		436
568		437
569		438
570		439
571		440
572		441
573		442
574		443
575		444
576		445
577		446
578		447
579		448
580		449
581		450
582		451
583		452
584		453
585		454
586		455
587		456
588		457
589		458
590		459
591		460
592		461
593		462
594		463
595		464
596		465
597		466
598		467
599		468
600		469
601		470
602		471
603		472
604		473
605		474
606		475
607		476
608		477
609		478
610		479
611		480
612		481
613		482
614		483
615		484
616		485
617		486
618		487
619		488
620		489
621		490
622		491
623		492
624		493
625		494
626		495
627		496
628		497
629		498
630		499
631		500
632		501
633		502
634		503
635		504
636		505
637		506
638		507
639		508
640		509
641		510
642		511
643		512
644		513
645		514
646		515
647		516
648		517
649		518
650		519
651		520
652		521
653		522
654		523
655		524
656		525
657		526
658		527
659		528
660		529
661		530
662		531
663		532
664		533
665		534
666		535
667		536
668		537
669		538
670		539
671		540
672		541
673		542
674		543
675		544
676		545
677		546
678		547
679		548
680		549
681		550
682		551
683		552
684		553
685		554
686		555
687		556
688		557
689		558
690		559
691		560
692		561
693		562
694		563
695		564
696		565
697		566
698		567
699		568
700		569
701		570
702		571
703		572
704		573
705		574
706		575
707		576
708		577
709		578
710		579
711		580
712		581
713		582
714		583
715		584
716		585
717		586
718		587
719		588
720		589
721		590
722		591
723		592
724		593
725		594
726		595
727		596
728		597
729		598
730		599
731		600
732		601
733		602
734		603
735		604
736		605
737		606
738		607
739		608
740		609
741		610
742		611
743		612
744		613
745		614
746		615
747		616
748		617
749		618
750		619
751		620
752		621
753		622
754		623
755		624
756		625
757		626
758		627
759		628
760		629
761		630
762		631
763		632
764		633
765		634
766		635
767		636
768		637
769		638
770		639
771		640
772		641
773		642
774		643
775		644
776		645
777		646
778		647
779		648
780		649
781		650
782		651
783		652
784		653
785		654
786		655
787		656
788		657
789		658
790		659
791		660
792		661
793		662
794		663
795		664
796		665
797		666
798		667
799		668
800		669
801		670
802		671
803		672
804		673
805		674
806		675
807		676
808		677
809		678
810		679
811		680
812		681
813		682
814		683
815		684
816		685
817		686
818		687
819		688
820		689
821		690
822		691
823		692
824		693
825		694
826		695
827		696
828		697
829		698
830		699
831		700
832		701
833		702
834		703
835		704
836		705
837		706
838		707
839		708

404 Dragomir R. Radev, Eduard Hovy, and Kathleen McKeown. 2002. *Introduction to the Special Issue on Summarization*. *Computational Linguistics*, 28(4):399–408.

405

406 Nils Reimers, Philip Beyer, and Iryna Gurevych. 2016. Task-Oriented Intrinsic Evaluation of Semantic Textual Similarity. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 87–96, Osaka, Japan. The COLING 2016 Organizing Committee.

407

408 Nils Reimers and Iryna Gurevych. 2018. Why Comparing Single Performance Scores Does Not Allow to Draw Conclusions About Machine Learning Approaches. *preprint, arXiv:1803.09578*.

409

410 Nils Reimers and Iryna Gurevych. 2020. Sentence-BERT: Sentence embeddings using siamese BERT-networks. In *EMNLP-IJCNLP 2019 - 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, Proceedings of the Conference*, pages 3982–3992.

411

412 Reka. 2024. Model Catalogue. <https://www.reka.ai/ourmodels> (accessed 2024-06-14).

413

414 Sentence Transformers. 2024. Pretrained Models — Sentence Transformers documentation.

415

416 Gizem Soğancıoğlu, Hakime Öztürk, and Arzucan Özgür. 2017. BIOSSES: A semantic sentence similarity estimation system for the biomedical domain. *Bioinformatics (Oxford, England)*, 33(14):i49–i58.

417

418 Raghveer Thirukovalluru and Bhuwan Dhingra. 2025. GenEOL: Harnessing the generative power of LLMs for training-free sentence embeddings. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 2295–2308, Albuquerque, New Mexico. Association for Computational Linguistics.

419

420 Voyage AI. 2024. Voyage AI Embedding Models. <https://docs.voyageai.com/docs/introduction> (accessed 2024-07-20).

421

422 Jan Philip Wahle, Terry Ruas, Frederic Kirstein, and Bela Gipp. 2022. How Large Language Models are Transforming Machine-Paraphrased Plagiarism.

423

424 Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. Multilingual e5 text embeddings: A technical report. *arXiv preprint arXiv:2402.05672*.

425

426 Liang Wang, Nan Yang, and Furu Wei. 2023. 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.

427

428 Bowen Zhang, Kehua Chang, and Chunping Li. 2024. Simple Techniques for Enhancing Sentence Embeddings in Generative Language Models. *Preprint, arXiv:2404.03921*.

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

455

456

457

458

459

A Comparison of GASE and GenEOL

A comparison of our approach, GASE, and GenEOL can be found in Table 7. Based on 7 STS datasets, Thirukovalluru and Dhingra (2025) report a Spearman’s rank correlation for the weakest baseline model of 47.06 and 80.37 for their best model using 32 generative augmentations. Using the same 7 datasets, we find that without any augmentation their best performing model is outperformed by 7 of the 12 embedding models evaluated in this work (e.g. by all-mpnet-base-v2; see Table 13 and Table 14).

460

461

462

463

464

465

466

467

468

469

470

	<b>GASE (ours)</b>	<b>GenEOL</b>
<b>Transformations</b>	Paraphrasing Summarising Extracting keywords	Paraphrasing Concise paraphrasing Changing the sentence structure Entailment
		Each one followed by summarisation.
<b>Number of transformations</b>	$k \in \{0, 1, 2, 3\}$	$m \in \{0, 2, 4, 8, 16, 24, 32\}$
<b>Pooling</b>	Mean pooling over variations generated <i>using k different generative LLMs</i>	Mean pooling over variations generated <i>using m different transformations</i>
<b>Embedding models</b>	English-language models: <code>glove.840B.300d</code> <code>bert-large-cased</code> <code>all-MiniLM-L6-v2</code> <code>all-mpnet-base-v2</code> <code>embed-english-light-v3.0</code> <code>embed-english-v3.0</code> <code>voyage-2</code> <code>voyage-lite-02-instruct</code> <code>voyage-large-2</code> <code>voyage-large-2-instruct</code> <code>mxbai-embed-large-v1</code> <code>llama-3.1-8B</code>  Multilingual models: <code>paraphrase-multilingual-mpnet-base-v2</code> <code>multilingual-e5-large-instruct</code> <code>voyage-multilingual-2</code>	<code>mistral0.1-7B</code> <code>llama-2-7B</code> <code>llama-3-8B</code> <code>bert-large</code>
<b>Generative models</b>	<code>gpt-3.5-turbo-0125</code> <code>reka-flash-20240226</code> <code>gpt-4o-mini-2024-07-18</code> <code>qwen2.5-0.5b</code> <code>qwen2.5-1.5b</code> <code>qwen2.5-3b</code> <code>qwen2.5-7b</code> <code>qwen2.5-14b</code> <code>qwen2.5-32b</code>	<code>gpt-3.5-turbo-0125</code> <code>mistral0.1-1-7B</code>
<b>STS datasets</b>	Complete English MTEB for STS: <code>STSB</code> <code>STS12</code> <code>STS13</code> <code>STS14</code> <code>STS15</code> <code>STS16</code> <code>STS17</code> <code>STS22</code> <code>Sick-R</code> <code>BIOSSES</code>  Multi-lingual datasets: <code>STS17</code> <code>STS22</code>	Subset of English MTEB for STS: <code>STSB</code> <code>STS12</code> <code>STS13</code> <code>STS14</code> <code>STS15</code> <code>STS16</code> <code>Sick-R</code>
<b>Other tasks</b>	TextPairClassification Summarisation	TextPairClassification Classification Clustering Reranking

Table 7: Comparison of GASE and GenEOL.

## 472 B Random keyword extraction

473 Random keyword extraction randomly selects  
474 16.1/58.2 of the words from a text sequence.  
475 This share is equal to the average share of LLM-  
476 extracted keywords and the total word count of a  
477 sentence (see Table 2). As this results in a low  
478 number of keywords for short sentences, the min-  
479 imal number of randomly selected keywords was  
480 set to 3. All punctuation has been removed before  
481 extracting the keywords.

## 482 C Prompts

483 Paraphrasing with GPT-3.5 Turbo and GPT-4o  
484 mini :

485 *"Rephrase the following text while maintaining its  
486 original meaning. If the text contains only a single  
487 word, provide a definition or a synonym. When  
488 done, check and make sure that the length of the  
489 original is approximately maintained. Text:"*

490 Paraphrasing with Reka-Flash:

491 *"Rephrase the following text while maintaining  
492 its original meaning. Do not provide multiple  
493 alternatives. Before you reply, remove any  
494 explanations. Do only reply with the paraphrased  
495 text. If the text contains only a single word, provide  
496 a definition or a synonym. Text:"*

497 Summarising with GPT-3.5 Turbo:

498 *"summarise the following text. Do not include any  
499 meta text and only output the summary. If the text  
500 is too short to summarise, paraphrase it instead.  
501 Text:"*

502 Extracting keywords with GPT-3.5 Turbo:

503 *"Extract the keywords from the following sentence.  
504 Do NOT include any commas or fullstops in  
505 your response and do not start your answer with  
506 "keywords". Text: "*

507 Model-specific post-processing is performed to  
508 remove meta-text not used for augmentation.

## 513 D Embeddings models

- 514 • glove.840B.300d (Pennington et al., 2014)
- 515 • bert-large-cased (Devlin et al., 2019)
- 516 • all-MiniLM-L6-v2 (Sentence Transformers,  
517 2024)

- 518 • all-mpnet-base-v2 (Sentence Transform-  
519 ers, 2024) 519
  - 520 • embed-english-light-v3.0 (Cohere,  
521 2024) 521
  - 522 • embed-english-v3.0 (Cohere, 2024) 522
  - 523 • voyage-2 (Voyage AI, 2024) 523
  - 524 • voyage-lite-02-instruct (Voyage AI,  
525 2024) 524
  - 526 • voyage-large-2 (Voyage AI, 2024) 526
  - 527 • voyage-large-2-instruct (Voyage AI,  
528 2024) 528
  - 529 • mxbai-embed-large-v1 (Mixedbread, 2024) 529
  - 530 • llama-3.1-8b (Meta) 530
  - 531 • paraphrase-multilingual-mpnet-  
532 base-v2 (Reimers and Gurevych, 2020) 531
  - 533 • intfloat/multilingual-e5-  
534 large-instruct (Wang et al., 2024) 534
  - 535 • voyage-multilingual-2 (Voyage AI, 2024) 535
- 536 For llama-3.1-8b we used the following  
537 KEEOL (Zhang et al., 2024) prompt:
- 538 The essence of a sentence is often cap-  
539 tured by its main subjects and actions,  
540 while descriptive terms provide addi-  
541 tional but less central details. With this  
542 in mind, this sentence: "input\_text"  
543 means in one word:
- 544 Following Zhang et al. (2024), we used the penulti-  
545 mate layer of llama-3.1-8b to extract the embed-  
546 dings.
- ## 547 E Datasets
- 548 The following datasets from the MTEB STS leader-  
549 board ({Hugging Face}, 2024) were used:
- 550 • STSBenchmark (Cer et al., 2017a) 550
  - 551 • STS12 (Agirre et al., 2012) 551
  - 552 • STS13 (Agirre et al., 2013) 552
  - 553 • STS14 (Agirre et al., 2014) 553
  - 554 • STS15 (Agirre et al., 2015) 554
  - 555 • STS16 (Agirre et al., 2016) 555

- 556
- STS17 (Cer et al., 2017b)

557

  - STS22 (Chen et al., 2022)

558

  - SICK-R (Marelli et al., 2014)

559

  - BIOSSES (Soğancioğlu et al., 2017)

560 The original MTEB paper (Muennighoff et al.,  
 561 2023) also included the STS11 dataset in Figure 1  
 562 of their paper where the authors specify the MTEB  
 563 datasets, but excluded it in the evaluations. Sim-  
 564ilarly, the MTEB Leaderboard ({Hugging Face},  
 565 2024) does not include STS11 either. To make  
 566 our results comparable to other models evaluated  
 567 on MTEB we therefore excluded STS11 from our  
 568 evaluations.

569 **F Generative Model Versions and**  
 570 **Hyperparameters**

571 The following versions have been used for genera-  
 572 tive augmentation:

- 573
- GPT-3.5 Turbo: gpt-3.5-turbo-0125

574

  - Reka-Flash: reka-flash-20240226

575

  - GPT-4o mini: gpt-4o-mini-2024-07-18

576

  - Qwen2.5 (all models running via Ollama with  
 577 quantization Q4\_K\_M):
    - qwen2.5:0.5b
    - qwen2.5:1.5b
    - qwen2.5:3b
    - qwen2.5:7b
    - qwen2.5:14b

583 GPT-3.5 Turbo, Reka-Flash, and GPT-4o mini  
 584 have been accessed through the respective APIs  
 585 with temperature set to 0 and top\_p to 1 to make  
 586 the experiments as reproducible as possible. Where  
 587 applicable, a random seed of 1337 was used. Other  
 588 hyperparameters were used with default values.  
 589 The Qwen2.5 model were run with Ollama on a  
 590 local computer.

591 **G Additional Experimental Results**

Embedding model	Pooling	Avg. Score
glove.840B.300d	mean	59.98
	max	59.02
bert-large-cased	mean	66.29
	max	65.16
all-MiniLM-L6-v2	mean	80.32
	max	79.88
all-mpnet-base-v2	mean	81.39
	max	81.03
embed-english-light-v3.0	mean	80.30
	max	79.74
embed-english-v3.0	mean	82.10
	max	81.62
voyage-2	mean	83.38
	max	82.69
voyage-lite-02-instruct	mean	85.99
	max	85.53
voyage-large-2	mean	83.92
	max	83.48
voyage-large-2-instruct	mean	84.80
	max	84.41
mxbai-embed-large-v1	mean	85.01
	max	84.69
llama-3.1-8b	mean	75.28
	max	75.07

Table 8: Average STS scores across tasks for different pooling methods.

Embedding model	Augmentation	Average	STSB	STS12	STS13	STS14	STS15	STS16	STS17	STS22	SICK-R	BIOSSES
glove.840B.300d	none	57.86	50.73	57.50	70.98	60.69	70.85	63.85	62.05	54.08	55.42	32.45
	paraphrasing	59.98	56.78	57.33	<b>71.17</b>	61.16	71.52	<b>67.66</b>	62.28	57.78	58.92	<b>35.21</b>
	summarising	60.64	57.67	55.00	70.24	<b>61.18</b>	71.64	66.81	67.08	<b>61.21</b>	<b>60.92</b>	34.65
	extracting keywords	<b>60.82</b>	<b>59.26</b>	<b>58.09</b>	71.12	60.34	<b>72.44</b>	65.31	<b>75.92</b>	56.94	55.34	33.50
bert-large-cased	none	60.78	59.18	45.71	63.94	54.42	68.98	65.46	71.30	55.48	64.16	59.14
	paraphrasing	66.29	67.97	<b>53.62</b>	68.56	59.73	74.01	<b>71.97</b>	74.14	58.57	69.79	<b>64.56</b>
	summarising	<b>66.89</b>	<b>68.23</b>	51.39	<b>70.58</b>	<b>61.98</b>	<b>74.27</b>	70.60	<b>74.48</b>	<b>63.47</b>	<b>70.74</b>	63.23
	extracting keywords	62.62	63.03	53.50	63.32	58.51	72.62	68.13	71.99	54.64	58.24	62.22
all-MiniLM-L6-v2	none	78.91	82.03	72.37	80.60	75.59	85.39	78.99	87.59	67.26	77.58	81.64
	paraphrasing	<b>80.32</b>	<b>82.33</b>	<b>73.33</b>	<b>83.33</b>	<b>77.85</b>	86.23	<b>79.88</b>	87.02	68.47	79.63	<b>85.07</b>
	summarising	79.90	82.17	72.14	82.89	77.24	<b>86.34</b>	78.99	<b>87.70</b>	<b>69.27</b>	<b>79.86</b>	82.41
	extracting keywords	77.96	81.99	73.29	81.14	75.77	85.87	78.70	87.45	64.57	70.92	79.88
all-mpnet-base-v2	none	80.28	83.42	72.63	83.48	78.00	85.66	80.03	<b>90.60</b>	68.00	80.59	80.43
	paraphrasing	81.39	<b>84.10</b>	74.04	85.52	<b>80.52</b>	<b>86.82</b>	<b>81.15</b>	89.08	69.01	81.25	82.39
	summarising	<b>81.48</b>	83.94	73.76	<b>85.54</b>	79.78	86.73	80.20	90.41	<b>70.12</b>	<b>81.32</b>	<b>82.96</b>
	extracting keywords	79.95	82.56	<b>74.82</b>	84.82	78.74	86.31	79.01	90.23	66.19	77.07	79.74
embed-english-light-v3.0	none	78.75	83.52	72.81	77.69	76.91	83.51	78.49	<b>88.49</b>	67.88	77.98	80.20
	paraphrasing	<b>80.30</b>	<b>84.65</b>	74.65	80.73	<b>79.41</b>	85.45	81.14	87.45	68.42	78.75	<b>82.36</b>
	summarising	80.21	84.08	73.45	<b>81.36</b>	78.41	<b>86.26</b>	<b>81.52</b>	88.19	<b>68.97</b>	<b>79.08</b>	80.79
	extracting keywords	78.53	82.81	<b>75.03</b>	80.17	78.29	84.14	76.38	88.14	67.27	74.45	78.64
embed-english-v3.0	none	81.25	86.54	74.76	81.68	78.81	87.18	83.01	<b>89.70</b>	68.22	77.52	85.05
	paraphrasing	<b>82.10</b>	<b>86.69</b>	<b>76.37</b>	83.71	<b>81.29</b>	87.96	<b>83.95</b>	88.47	68.63	<b>78.21</b>	<b>85.75</b>
	summarising	82.08	86.47	74.40	84.41	80.76	<b>88.49</b>	83.55	89.08	<b>70.29</b>	78.06	85.25
	extracting keywords	81.26	85.50	74.77	<b>84.66</b>	79.98	87.72	81.09	88.75	66.91	77.53	85.68
voyage-2	none	82.50	87.06	77.68	86.25	80.11	88.38	85.72	<b>89.78</b>	65.26	78.95	85.83
	paraphrasing	<b>83.38</b>	<b>87.17</b>	<b>79.23</b>	<b>87.13</b>	<b>81.67</b>	<b>88.57</b>	<b>86.27</b>	89.07	66.80	80.33	<b>87.56</b>
	summarising	82.74	86.51	76.81	86.25	80.89	88.46	84.88	89.10	<b>68.50</b>	<b>80.48</b>	85.55
	extracting keywords	81.21	85.71	77.76	86.17	79.63	87.49	84.37	88.66	65.29	72.06	85.03
voyage-lite-02-instruct	none	85.95	88.74	<b>86.19</b>	88.84	<b>86.84</b>	<b>89.84</b>	86.07	<b>87.19</b>	78.62	77.54	89.60
	paraphrasing	<b>85.99</b>	<b>88.81</b>	84.47	<b>89.05</b>	86.68	89.52	<b>86.53</b>	86.79	79.08	79.09	89.89
	summarising	85.43	87.93	82.77	88.06	85.05	88.98	85.49	86.70	<b>79.87</b>	<b>79.16</b>	<b>90.32</b>
	extracting keywords	83.91	86.75	83.82	87.97	84.47	88.98	84.55	86.28	76.60	70.77	88.88
voyage-large-2	none	83.60	<b>87.84</b>	78.66	86.98	82.25	<b>89.94</b>	85.71	<b>91.29</b>	63.97	79.83	<b>90.49</b>
	paraphrasing	<b>83.92</b>	87.62	<b>79.18</b>	<b>87.66</b>	<b>83.60</b>	88.91	<b>86.00</b>	90.22	64.69	81.00	90.29
	summarising	83.66	87.56	77.83	86.93	82.88	88.93	84.88	90.49	<b>65.90</b>	<b>81.09</b>	90.15
	extracting keywords	81.98	85.92	78.00	86.45	81.33	88.28	84.44	89.90	63.68	72.50	89.33
voyage-large-2-instruct	none	84.61	89.21	76.15	88.49	86.50	91.13	85.68	<b>90.05</b>	66.59	83.17	89.09
	paraphrasing	<b>84.80</b>	<b>89.31</b>	<b>77.61</b>	88.40	<b>86.62</b>	<b>91.17</b>	<b>86.20</b>	89.58	67.72	82.83	88.58
	summarising	84.54	88.84	76.27	88.29	85.56	91.06	85.24	89.40	<b>68.66</b>	<b>83.28</b>	88.78
	extracting keywords	83.40	87.38	76.63	<b>88.58</b>	85.55	90.51	84.03	88.57	65.73	76.58	<b>90.46</b>
mxbai-embed-large-v1	none	84.63	<b>89.29</b>	79.07	<b>89.80</b>	85.22	89.34	86.77	<b>89.21</b>	68.73	82.78	86.13
	paraphrasing	<b>85.01</b>	89.09	<b>80.83</b>	89.43	<b>85.88</b>	<b>89.46</b>	<b>86.94</b>	88.75	70.03	82.67	87.02
	summarising	84.45	88.60	78.51	89.07	84.82	89.09	85.96	88.43	<b>70.36</b>	<b>82.82</b>	86.82
	extracting keywords	83.42	87.43	78.64	88.75	83.82	88.44	84.84	87.40	67.83	78.17	<b>88.90</b>
llama-3.1-8b	none	71.36	79.79	62.81	79.08	69.72	78.49	80.88	84.12	33.39	76.27	69.04
	paraphrasing	75.28	<b>81.37</b>	67.03	81.29	72.88	80.91	<b>82.60</b>	83.04	50.55	78.25	74.84
	summarising	<b>76.74</b>	80.39	65.02	81.18	74.51	<b>82.07</b>	81.00	84.48	<b>59.44</b>	<b>79.36</b>	<b>79.92</b>
	extracting keywords	74.95	78.57	<b>67.39</b>	<b>82.68</b>	<b>74.74</b>	81.67	80.31	<b>85.71</b>	51.79	71.05	75.60

Table 9: Scores per dataset with and without generative augmentation using GPT-3.5 Turbo (Spearman’s rank correlation in %, bold: highest scores).

Dataset	No Augmentation	Extracting Keywords (GPT-3.5 Turbo)	Random Keyword Extraction
STSB	50.73	<b>59.26</b>	23.95
STS12	57.50	<b>58.09</b>	34.78
STS13	70.98	<b>71.12</b>	22.25
STS14	<b>60.69</b>	60.34	28.48
STS15	70.85	<b>72.44</b>	33.86
STS16	63.85	<b>65.31</b>	24.99
STS17	62.05	<b>75.92</b>	36.47
STS22	54.08	<b>56.94</b>	38.44
SICK-R	<b>55.42</b>	55.34	40.04
BIOSSES	32.45	<b>33.50</b>	29.02
<b>Average</b>	<b>57.86</b>	<b>60.82</b>	31.23

Table 10: Scores per dataset for GloVe embeddings without generative augmentation vs. generative augmentation using LLM-extracted keywords with GPT-3.5 Turbo or random keyword extraction (Spearman’s rank correlation in %, bold: highest scores).

Embedding model	Augment- ation	Stop- word removal	Average	STSB	STS12	STS13	STS14	STS15	STS16	STS17	STS22	SICK-R	BIOSSES
glove.840B. 300d	none	no	57.86	50.73	57.50	70.98	60.69	70.85	63.85	62.05	54.08	55.42	32.45
		yes	60.07	57.13	57.43	70.91	60.68	70.86	63.84	75.83	54.66	54.97	34.39
	para- phrasing	no	59.98	56.78	57.33	71.17	61.16	71.52	67.66	62.28	57.78	58.92	35.21
		yes	65.43	67.05	59.98	72.72	63.75	73.70	69.98	75.01	61.01	61.94	49.17
bert-large-cased	none	no	60.78	59.18	45.71	63.94	54.42	68.98	65.46	71.30	55.48	64.16	59.14
		yes	53.08	50.11	41.37	53.10	50.11	63.68	59.66	63.73	49.69	48.75	50.54
	para- phrasing	no	66.29	67.97	53.62	68.56	59.73	74.01	71.97	74.14	58.57	69.79	64.56
		yes	60.25	61.42	50.45	58.43	55.20	69.75	67.19	67.98	57.01	58.25	56.82
all-MiniLM- L6-v2	none	no	78.91	82.03	72.37	80.60	75.59	85.39	78.99	87.59	67.26	77.58	81.64
		yes	75.36	78.80	68.29	77.62	73.05	84.35	74.92	87.05	65.74	65.58	78.24
	para- phrasing	no	80.32	82.33	73.33	83.33	77.85	86.23	79.88	87.02	68.47	79.63	85.07
		yes	77.78	80.69	71.61	81.15	75.92	84.86	77.51	85.89	66.36	70.61	83.21
all-mpnet-base- v2	none	no	80.28	83.42	72.63	83.48	78.00	85.66	80.03	90.60	68.00	80.59	80.43
		yes	75.51	78.32	65.71	79.49	73.99	83.98	73.28	88.98	65.82	67.01	78.47
	para- phrasing	no	81.39	84.10	74.04	85.52	80.52	86.82	81.15	89.08	69.01	81.25	82.39
		yes	78.20	80.87	70.76	83.09	77.53	85.00	76.72	87.05	66.28	71.99	82.70
embed-english- light-v3.0	none	no	78.75	83.52	72.81	77.69	76.91	83.51	78.49	88.49	67.88	77.98	80.20
		yes	74.36	78.16	68.36	76.23	73.84	81.67	71.69	88.18	66.87	64.32	74.34
	para- phrasing	no	80.30	84.65	74.65	80.73	79.41	85.45	81.14	87.45	68.42	78.75	82.36
		yes	77.58	81.43	73.12	79.15	76.81	83.77	76.89	87.55	67.81	70.40	78.87
embed-english- v3.0	none	no	81.25	86.54	74.76	81.68	78.81	87.18	83.01	89.70	68.22	77.52	85.05
		yes	76.70	79.52	67.15	80.84	75.37	83.76	74.40	88.51	66.99	66.56	83.95
	para- phrasing	no	82.10	86.69	76.37	83.71	81.29	87.96	83.95	88.47	68.63	78.21	85.75
		yes	79.26	83.00	72.79	83.09	78.72	85.21	79.03	87.36	68.03	71.45	83.96
voyage-2	none	no	82.50	87.06	77.68	86.25	80.11	88.38	85.72	89.78	65.26	78.95	85.83
		yes	78.11	81.52	71.71	82.73	77.23	85.59	79.12	87.91	66.40	67.22	81.64
	para- phrasing	no	83.38	87.17	79.23	87.13	81.67	88.57	86.27	89.07	66.80	80.33	87.56
		yes	80.42	84.16	76.32	84.89	79.29	86.33	82.93	87.58	67.55	71.42	83.69
voyage-lite-02- instruct	none	no	85.95	88.74	86.19	88.84	86.84	89.84	86.07	87.19	78.62	77.54	89.60
		yes	81.35	83.24	79.33	86.13	82.37	87.95	80.01	85.84	74.72	66.73	87.14
	para- phrasing	no	85.99	88.81	84.47	89.05	86.68	89.52	86.53	86.79	79.08	79.09	89.89
		yes	82.88	85.57	81.64	87.14	83.29	87.58	83.46	85.18	76.02	70.83	88.11
voyage-large-2- instruct	none	no	83.60	87.84	78.66	86.98	82.25	88.94	85.71	91.29	63.97	79.83	90.49
		yes	79.39	82.25	72.85	82.94	78.99	86.54	79.68	89.63	65.38	67.42	88.25
	para- phrasing	no	83.92	87.62	79.18	87.66	83.60	88.91	86.00	90.22	64.69	81.00	90.29
		yes	81.19	84.57	76.72	84.73	80.84	87.17	82.95	88.78	65.43	71.64	89.04
voyage-large-2- instruct	none	no	84.61	89.21	76.15	88.49	86.50	91.13	85.68	90.05	66.59	83.17	89.09
		yes	81.35	83.24	79.33	86.13	82.37	87.95	80.01	85.84	74.72	66.73	87.14
	para- phrasing	no	84.80	89.31	77.61	88.40	86.62	91.17	86.20	89.58	67.72	82.83	88.58
		yes	81.88	85.70	74.85	87.76	83.91	88.69	82.40	88.41	67.66	71.72	87.73
mxbai-embed- large-v1	none	no	84.63	89.29	79.07	89.80	85.22	89.34	86.77	89.21	68.73	82.78	86.13
		yes	79.54	83.54	72.88	85.73	81.33	86.72	78.69	86.36	67.41	68.76	83.99
	para- phrasing	no	85.01	89.09	80.83	89.43	85.88	89.46	86.94	88.75	70.03	82.67	87.02
		yes	81.36	86.15	78.11	87.20	83.03	86.93	78.69	86.40	68.47	72.70	85.97
llama-3.1-8b	none	no	71.36	79.79	62.81	79.08	69.72	78.49	80.88	84.12	33.39	76.27	69.04
		yes	69.50	73.48	59.66	77.37	69.66	77.17	72.94	83.00	50.94	62.15	68.58
	para- phrasing	no	75.28	81.37	67.03	81.29	72.88	80.91	82.60	83.04	50.55	78.25	74.84
		yes	74.57	78.29	66.98	80.90	73.61	79.63	78.76	83.86	61.88	68.92	72.90

Table 11: Scores per dataset with and without stopword removal using either no generative augmentation or paraphrase augmentation with GPT-3.5 Turbo (Spearman’s rank correlation in %).

Dataset	Paraphrases	Summaries
STSBenchmark	0.39	0.52
STS12	0.39	0.49
STS13	0.33	0.41
STS14	0.37	0.47
STS15	0.42	0.53
STS16	0.37	0.40
STS17	0.41	0.59
STS22	0.39	0.21
SICK-R	0.46	0.67
BIOSSES	0.46	0.52
<b>Average</b>	<b>0.40</b>	<b>0.48</b>

Table 12: Jaccard Similarity between original texts and paraphrases and summaries generated with GPT-3.5 Turbo respectively.

Embedding model	Gen model	Average	STSB	STS12	STS13	STS14	STS15	STS16	STS17	STS22	SICK-R	BIOSES
glove.840B.300d	none	57.86	50.73	57.50	70.98	60.69	70.85	63.85	62.05	54.08	55.42	32.45
	gpt-3.5-turbo	59.98	56.78	57.33	71.17	61.16	71.52	<b>67.66</b>	62.28	57.78	58.92	35.21
	reka-flash	60.73	55.02	59.53	71.09	61.22	72.81	65.20	64.27	60.11	61.65	36.43
	gpt-4o-mini	60.04	55.09	60.12	72.65	62.45	71.60	66.49	63.39	55.48	60.36	32.76
	gpt-3.5-turbo + reka-flash	61.97	59.46	60.06	72.01	62.68	72.46	67.62	<b>65.43</b>	61.46	61.51	<b>36.95</b>
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	<b>62.18</b>	<b>59.99</b>	<b>61.79</b>	<b>72.88</b>	<b>64.04</b>	<b>72.03</b>	67.38	64.75	<b>60.76</b>	<b>61.79</b>	36.42
	none	60.78	59.18	45.71	63.94	54.42	68.98	65.46	71.30	55.48	64.16	59.14
bert-large-cased	gpt-3.5-turbo	66.29	67.97	53.62	68.56	59.73	74.01	71.97	74.14	58.57	69.80	64.56
	reka-flash	65.76	65.20	52.92	68.55	59.77	72.67	71.80	71.46	58.93	70.09	66.19
	gpt-4o-mini	65.01	65.37	54.27	68.60	58.18	72.26	71.12	73.54	55.72	70.50	60.54
	gpt-3.5-turbo + reka-flash	68.01	69.39	56.86	70.04	62.30	<b>74.61</b>	73.86	74.52	60.49	71.32	66.72
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	<b>68.41</b>	<b>69.40</b>	<b>59.01</b>	<b>71.25</b>	<b>62.60</b>	74.13	<b>73.86</b>	<b>74.99</b>	<b>60.35</b>	<b>71.74</b>	<b>66.73</b>
	none	78.91	82.03	72.37	80.60	75.59	85.39	78.99	87.59	67.26	77.58	81.64
	gpt-3.5-turbo	80.32	82.33	73.33	83.33	77.85	86.23	79.88	87.02	68.47	79.63	85.07
all-MiniLM-L6-v2	reka-flash	79.74	81.30	73.52	82.42	76.83	85.72	79.02	87.01	68.83	79.66	83.11
	gpt-4o-mini	79.81	81.97	74.03	83.31	76.70	86.03	79.19	<b>87.69</b>	67.01	80.08	82.11
	gpt-3.5-turbo + reka-flash	80.64	82.30	74.59	83.72	78.11	86.20	<b>79.95</b>	87.23	<b>69.23</b>	80.12	84.97
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	<b>80.76</b>	<b>82.31</b>	<b>75.16</b>	<b>84.29</b>	<b>78.22</b>	<b>86.18</b>	79.98	87.31	68.84	<b>80.34</b>	<b>85.02</b>
	none	80.28	83.42	72.63	83.48	78.00	85.66	80.03	90.60	68.00	80.59	80.43
	gpt-3.5-turbo	81.39	84.10	74.04	85.52	80.52	86.82	81.15	89.08	69.01	81.25	82.39
	reka-flash	81.01	83.72	74.27	85.12	79.33	86.41	80.41	90.17	68.42	81.37	80.87
all-mpnet-base-v2	gpt-4o-mini	81.07	83.75	74.19	85.19	79.28	86.45	80.47	<b>90.89</b>	67.81	81.23	81.43
	gpt-3.5-turbo + reka-flash	81.72	84.43	75.46	86.12	80.63	87.00	81.21	89.37	<b>69.33</b>	<b>81.55</b>	82.13
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	<b>81.85</b>	<b>84.42</b>	<b>75.96</b>	<b>86.37</b>	<b>80.72</b>	<b>86.92</b>	<b>81.33</b>	89.68	69.05	81.50	<b>82.51</b>
	none	78.75	83.52	72.81	77.69	76.91	83.51	78.49	88.49	67.88	77.98	80.20
	gpt-3.5-turbo	80.30	<b>84.65</b>	74.65	80.73	79.41	85.45	81.14	87.45	68.42	78.75	<b>82.36</b>
	reka-flash	79.72	83.51	75.74	78.60	78.38	84.63	80.09	87.98	69.06	78.76	80.42
	gpt-4o-mini	79.90	84.14	75.76	79.65	77.89	85.12	79.83	87.90	67.91	78.68	82.10
embed-english-light-v3.0	gpt-3.5-turbo + reka-flash	80.54	84.48	76.45	80.61	<b>79.65</b>	85.54	<b>81.21</b>	87.60	<b>69.23</b>	<b>79.06</b>	81.55
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	<b>80.66</b>	84.59	<b>77.10</b>	<b>81.00</b>	79.75	<b>85.68</b>	81.16	<b>87.70</b>	68.90	79.09	81.61
	none	81.25	86.54	74.76	81.68	78.81	87.18	83.01	89.70	68.22	77.52	85.05
	gpt-3.5-turbo	82.10	86.69	76.37	83.71	81.29	87.96	83.95	88.47	68.63	78.21	85.75
	reka-flash	82.00	86.31	76.91	82.22	80.48	87.82	83.69	88.93	69.70	78.06	85.90
	gpt-4o-mini	81.82	86.52	76.80	82.57	80.11	87.77	83.72	<b>89.42</b>	68.36	77.90	85.05
	gpt-3.5-turbo + reka-flash	<b>82.53</b>	86.68	77.88	83.56	81.64	88.05	84.13	88.59	<b>69.87</b>	<b>78.44</b>	<b>86.48</b>
embed-english-v3.0	gpt-3.5-turbo + reka-flash + gpt-4o-mini	82.08	<b>86.69</b>	<b>78.49</b>	<b>83.60</b>	<b>81.65</b>	<b>87.97</b>	<b>84.10</b>	88.76	69.67	78.35	86.36

Table 13: Scores per dataset with and without paraphrase augmentation using GPT-3.5 Turbo, Reka-Flash, and GPT-4o mini (Spearman’s rank correlation in %, bold: highest scores) (part 1).

Embedding model	Gen model	Average	STSB	STS12	STS13	STS14	STS15	STS16	STS17	STS22	SICK-R	BIOSES
voyage-2	none	82.50	87.06	77.68	86.25	80.11	88.38	85.72	<b>89.78</b>	65.26	78.95	85.83
	gpt-3.5-turbo	83.38	87.17	79.23	87.13	81.67	88.57	86.27	89.07	66.80	80.33	<b>87.56</b>
	reka-flash	83.37	86.93	79.65	87.83	81.79	88.48	86.18	88.98	67.42	80.67	85.72
	gpt-4o-mini	83.07	86.89	79.59	87.27	81.23	88.46	85.81	89.37	65.25	80.85	86.02
	gpt-3.5-turbo + reka-flash	<b>83.80</b>	<b>87.27</b>	80.62	87.99	82.40	88.51	<b>86.40</b>	88.77	<b>67.94</b>	80.87	87.28
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	83.72	87.14	<b>81.06</b>	<b>88.14</b>	<b>82.50</b>	<b>88.40</b>	86.14	88.81	66.91	<b>81.00</b>	87.08
	none	85.95	88.74	86.19	88.84	86.84	<b>89.84</b>	86.07	<b>87.19</b>	78.62	77.54	89.60
voyage-lite-02-instruct	gpt-3.5-turbo	85.99	<b>88.81</b>	84.47	89.05	86.68	89.52	86.53	86.79	<b>79.08</b>	79.09	89.89
	reka-flash	85.92	88.24	84.88	88.79	86.58	89.46	86.28	86.74	78.76	78.97	90.53
	gpt-4o-mini	85.97	88.25	85.24	89.06	86.62	89.57	85.99	86.90	78.59	<b>79.45</b>	90.08
	gpt-3.5-turbo + reka-flash	<b>86.05</b>	88.74	84.83	89.14	86.69	89.36	<b>86.57</b>	86.72	78.81	79.34	90.34
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	86.01	88.48	<b>84.95</b>	<b>89.26</b>	<b>86.74</b>	89.25	86.32	86.75	78.31	79.50	<b>90.55</b>
voyage-large-2	none	83.60	<b>87.84</b>	78.66	86.98	82.25	88.94	85.71	<b>91.29</b>	63.97	79.83	90.49
	gpt-3.5-turbo	83.92	87.62	79.18	87.66	83.60	88.91	86.00	90.22	64.69	81.00	90.29
	reka-flash	83.98	87.32	79.52	88.25	83.49	88.77	85.90	90.40	64.46	81.35	90.37
	gpt-4o-mini	83.95	87.54	80.02	87.91	83.22	<b>89.11</b>	85.68	90.47	63.99	81.54	90.08
	gpt-3.5-turbo + reka-flash	84.12	87.56	80.05	88.31	83.98	88.76	<b>85.97</b>	90.04	<b>64.84</b>	81.55	90.14
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	<b>84.17</b>	87.49	<b>80.53</b>	<b>88.42</b>	<b>84.00</b>	88.73	85.77	89.82	64.69	<b>81.75</b>	<b>90.49</b>
	none	84.61	89.21	76.15	88.49	86.50	91.13	85.68	<b>90.05</b>	66.59	83.17	<b>89.09</b>
voyage-large-2-instruct	gpt-3.5-turbo	84.80	<b>89.31</b>	77.61	88.40	86.62	91.17	<b>86.20</b>	89.58	<b>67.72</b>	82.83	88.58
	reka-flash	84.68	88.81	78.41	88.30	86.40	90.95	85.42	89.71	67.04	<b>83.22</b>	88.52
	gpt-4o-mini	84.78	88.83	78.71	88.43	86.45	<b>91.17</b>	85.89	89.14	67.09	83.19	88.87
	gpt-3.5-turbo + reka-flash	84.86	89.16	79.11	88.56	<b>86.67</b>	90.97	85.86	89.67	67.56	83.00	88.08
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	<b>84.91</b>	88.98	<b>79.79</b>	<b>88.59</b>	86.69	90.89	85.83	89.23	67.62	82.97	88.51
mxbai-embed-large-v1	none	84.63	<b>89.29</b>	79.07	89.80	85.22	89.34	86.77	<b>89.21</b>	68.73	82.78	86.13
	gpt-3.5-turbo	85.01	89.09	80.83	89.43	85.88	<b>89.46</b>	<b>86.94</b>	88.75	<b>70.03</b>	82.67	87.02
	reka-flash	84.80	88.93	81.04	89.36	85.51	89.27	86.48	89.04	68.67	82.66	87.07
	gpt-4o-mini	84.90	88.95	81.08	<b>89.76</b>	85.37	89.43	86.58	89.14	69.10	<b>82.73</b>	86.89
	gpt-3.5-turbo + reka-flash	85.17	89.16	81.96	89.37	85.96	89.39	86.79	88.97	69.51	82.72	87.81
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	<b>85.16</b>	89.13	<b>82.24</b>	89.54	<b>85.98</b>	89.36	86.70	88.92	69.44	82.65	<b>87.63</b>
	none	71.36	79.79	62.81	79.08	69.72	78.49	80.88	84.12	33.39	76.27	69.04
llama-3.1-8b	gpt-3.5-turbo	75.28	81.37	67.03	81.29	72.88	80.91	82.60	83.04	50.55	78.25	74.84
	reka-flash	77.02	80.44	68.26	82.29	74.75	81.65	82.53	<b>85.77</b>	58.68	78.43	77.38
	gpt-4o-mini	74.71	80.38	68.15	82.09	73.39	81.20	81.87	84.60	41.34	<b>78.82</b>	75.23
	gpt-3.5-turbo + reka-flash	<b>77.85</b>	81.68	69.82	82.99	75.64	82.10	<b>83.27</b>	84.46	59.90	78.91	<b>79.68</b>
	gpt-3.5-turbo + reka-flash + gpt-4o-mini	77.52	<b>81.41</b>	<b>70.96</b>	<b>83.47</b>	<b>75.76</b>	<b>81.99</b>	82.84	84.14	<b>56.86</b>	78.99	78.77

Table 14: Scores per dataset with and without paraphrase augmentation using GPT-3.5 Turbo, Reka-Flash, and GPT-4o mini (Spearman’s rank correlation in %, bold: highest scores) (part 2).

Dataset	paraphrase-multilingual-mpnet-base-v2		intfloat/multilingual-e5-large-instruct		voyage-multilingual-2	
	None	Paraph.	None	Paraph.	None	Paraph.
STS17_en-en	<b>86.99</b>	86.01	<b>87.64</b>	87.32	<b>90.81</b>	90.20
STS17_ar-ar	79.09	<b>81.04</b>	81.61	<b>83.36</b>	80.53	<b>81.58</b>
STS17_en-ar	80.85	<b>81.76</b>	78.66	<b>79.50</b>	81.59	<b>83.48</b>
STS17_en-de	83.28	<b>83.81</b>	<b>85.09</b>	85.05	88.09	<b>88.11</b>
STS17_en-tr	74.90	<b>75.44</b>	77.04	<b>78.88</b>	74.12	<b>77.09</b>
STS17_es-en	86.11	<b>86.25</b>	84.72	<b>85.35</b>	<b>88.24</b>	87.74
STS17_es-es	<b>85.14</b>	85.00	<b>88.09</b>	87.87	<b>88.71</b>	88.63
STS17_fr-en	81.17	<b>82.47</b>	83.26	<b>84.09</b>	87.10	<b>87.65</b>
STS17_it-en	<b>84.24</b>	84.20	84.94	<b>85.00</b>	88.98	<b>89.43</b>
STS17_ko-ko	83.29	<b>83.83</b>	81.58	<b>82.47</b>	80.17	<b>80.69</b>
STS17_nl-en	82.51	<b>82.65</b>	<b>84.65</b>	84.59	88.22	<b>88.38</b>

Table 15: Multilingual STS17 Scores with and without paraphrase augmentation using GPT-4o mini in % (bold: highest score per embedding model and dataset).

Dataset	paraphrase-multilingual-mpnet-base-v2		multilingual-e5-large-instruct		voyage-multilingual-2	
	None	Paraph	None	Paraph.	None	Paraph.
STS_ar-ar	<b>64.77</b>	56.36	<b>81.61</b>	65.02	59.65	<b>60.19</b>
STS_de-de	<b>62.30</b>	48.16	<b>78.66</b>	62.00	63.01	<b>63.08</b>
STS_de-en	<b>62.85</b>	57.03	<b>85.09</b>	61.72	<b>62.87</b>	62.56
STS_de-fr	<b>64.19</b>	62.68	<b>77.04</b>	65.52	<b>64.88</b>	64.12
STS_de-pl	<b>56.01</b>	41.26	<b>84.72</b>	56.97	54.14	<b>56.35</b>
STS_es-es	<b>73.27</b>	65.03	<b>88.09</b>	73.07	75.32	<b>75.87</b>
STS_es-en	<b>80.45</b>	74.22	<b>83.26</b>	79.99	82.43	<b>83.79</b>
STS_es-it	<b>78.71</b>	56.99	<b>84.94</b>	78.76	79.51	<b>79.78</b>
STS_fr-fr	<b>82.34</b>	76.56	81.58	<b>82.35</b>	<b>81.36</b>	81.35
STS_fr-pl	73.25	<b>84.52</b>	<b>84.65</b>	28.17	<b>73.25</b>	73.25
STS_it-it	<b>80.69</b>	62.82	64.77	<b>80.23</b>	79.92	<b>80.15</b>
STS_pl-pl	<b>46.83</b>	39.36	<b>62.30</b>	45.29	<b>46.30</b>	46.00
STS_pl-en	76.42	<b>76.90</b>	62.85	<b>75.61</b>	80.01	<b>80.81</b>
STS_ru-ru	<b>70.26</b>	63.33	64.19	<b>70.74</b>	68.35	<b>68.91</b>
STS_tr-tr	<b>69.56</b>	63.49	56.01	<b>70.16</b>	67.98	<b>68.34</b>
STS_zh-zh	<b>66.92</b>	66.89	<b>73.27</b>	68.53	72.22	<b>72.49</b>
STS_zh-en	69.56	<b>70.68</b>	<b>80.45</b>	66.06	<b>76.81</b>	76.53

Table 16: Multilingual STS22 Scores with and without paraphrase augmentation using GPT-4o mini in % (bold: highest score per embedding model and dataset).

Embedding model	Augmen-tation	Pair Classification		
		Twitter Sem-Eval 2015	Twitter URL-Corpus	Sprint duplicate questions
glove.840B.300d	none	43.81	63.09	68.33
	paraphrase	<b>50.04</b>	<b>70.49</b>	<b>73.39</b>
bert-large-cased	none	47.27	69.81	48.15
	paraphrase	<b>57.21</b>	<b>76.28</b>	<b>53.61</b>
all-MiniLM-L6-v2	none	67.86	84.61	94.55
	paraphrase	<b>71.50</b>	<b>84.70</b>	<b>95.20</b>
all-mpnet-base-v2	none	73.85	84.99	90.15
	paraphrase	<b>77.34</b>	<b>85.08</b>	<b>91.42</b>
embed-english-light-v3.0	none	68.26	83.75	88.54
	paraphrase	<b>72.76</b>	<b>84.49</b>	<b>91.40</b>
embed-english-v3.0	none	73.43	84.98	92.21
	paraphrase	<b>76.58</b>	<b>85.45</b>	<b>93.33</b>
voyage-2	none	71.29	<b>86.25</b>	93.88
	paraphrase	<b>74.83</b>	86.19	<b>94.70</b>
voyage-lite-02-instruct	none	<b>91.38</b>	<b>89.14</b>	98.35
	paraphrase	90.31	88.37	<b>98.39</b>
voyage-large-2	none	75.10	86.57	93.96
	paraphrase	<b>76.82</b>	<b>86.58</b>	<b>95.25</b>
voyage-large-2-instruct	none	<b>92.87</b>	<b>90.10</b>	96.28
	paraphrase	90.78	89.09	<b>96.78</b>
mxbai-embed-large-v1	none	78.55	<b>86.12</b>	96.82
	paraphrase	<b>79.17</b>	86.03	<b>96.95</b>
llama-3.1-8b	none	69.70	82.73	57.87
	paraphrase	<b>73.00</b>	<b>84.26</b>	<b>61.09</b>

Table 17: Performance on Pair Classification with and without paraphrase generation using GPT-3.5-Turbo (in %).

## H Runtime analysis

To provide a basic estimate of the additional runtime introduced by GASE over a non-augmented baseline, we performed five runs using GPT-3.5-Turbo for generation and all-mpnet-base-v2 for embeddings and report the mean and standard deviation. Because the difference in runtimes between the augmented and non-augmented setups was substantial and the standard deviation across runs was small, we did not carry out statistical significance testing (see [Table 18](#)).

All runtime experiments were conducted on a computer with the following specification: GeForce RTX 3090 MSI Gaming X Trio 24GB GPU, AMD Ryzen 9 9950X CPU, DDR5-6400 64GB RAM.

<b>Augmentation</b>	<b>Average runtime [sec]</b>	<b>Stand. deviation of runtime [sec]</b>
None	2.30	0.40
paraphrase	132.57	15.87
summarise	140.48	14.15
keywords	127.17	18.16

Table 18: Average runtime and corresponding standard deviation without and with generative augmentation using GPT-3.5-Turbo and, as an encoder, all-mpnet-base-v2 on STS17 (based on runs, in seconds).