

KaLM-Embedding-V2: Superior Training Techniques and Data Inspire A Versatile Embedding Model

Xinping Zhao, Xinshuo Hu, Zifei Shan, Shouzheng Huang, Yao Zhou, Xin Zhang, Zetian Sun, Zhenyu Liu, Dongfang Li, Xinyuan Wei, Youcheng Pan, Yang Xiang, Meishan Zhang, Haofen Wang, Jun Yu, Baotian Hu, Min Zhang

¹Shenzhen Loop Area Institute (SLAI), ²Tencent

深圳河套学院
Shenzhen Loop Area Institute

Tencent 腾讯

Superior Training Techniques and High-quality Data

Focal-style Reweighting Mechanism

- It emphasizes difficult samples, where the more difficult the sample, the larger the weight.

$$w_i = (1 - \frac{e^{s(\mathbf{q}_i, \mathbf{p}_i^+) / \tau}}{Z_i})^\gamma, \quad \mathcal{L} = \mathbb{E}_{i \in N} \left[-w_i \log \frac{e^{s(\mathbf{q}_i, \mathbf{p}_i^+) / \tau}}{Z_i} \right]$$

Online Hard Negative Mixing Strategy

- It synthesizes new informative hard negatives via pair-wise or list-wise mixing.

$$\tilde{\mathbf{h}}_i^- = \frac{\mathbf{h}_i^-}{\|\mathbf{h}_i^-\|_2}, \quad \tilde{\mathbf{h}}_i^- = \lambda \mathbf{p}_{i,j}^- + (1 - \lambda) \mathbf{p}_{i,k}^-, \quad j, k \in [1, M]$$

$$\tilde{\mathbf{s}}_i^- = \frac{\mathbf{s}_i^-}{\|\mathbf{s}_i^-\|_2}, \quad \tilde{\mathbf{s}}_i^- = \sum_{m=1}^M \lambda_m \mathbf{p}_{i,m}^-, \quad \text{s.t.} \quad \sum_{m=1}^M \lambda_m = 1,$$

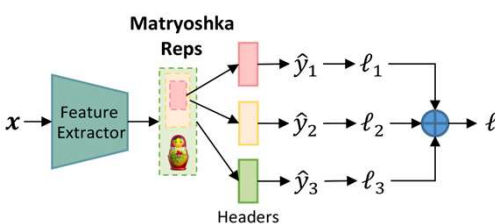
Contrastive Distillation

- It distills fine-grained soft signals from the teacher model that capture nuanced differences.

$$\mathcal{L}_{KL} = D_{KL}(P_i \| P_s) = \sum_i P_i(i) \log \frac{P_i(i)}{P_s(i)}, \quad P_i(i) = \frac{e^{s_{i,i} / \tau}}{\sum_j e^{s_{i,j} / \tau}}, \quad P_s(i) = \frac{e^{s_{s,i} / \tau}}{\sum_j e^{s_{s,j} / \tau}}$$

Matryoshka Representation Learning

- It enables flexible-dimensional embeddings.



Comprehensive Data Curation Recipe

- Around 20 categories of weakly supervised data as well as 100 categories of supervised data
- Hard Negative Mining
- Task-specific instructions
- Persona-based Data Synthesis
- Example-based multi-class labeling

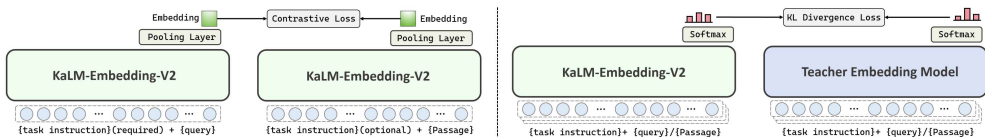
Lychee-KaLM-Embedding-Data

- KaLM-Embedding/KaLM-embedding-finetuning-data
- Viewer · Updated Nov 27, 2025 · 6.34M · ± 852 · ♥ 25
- HIT-TMG/KaLM-embedding-pretrain-data
- Viewer · Updated Nov 27, 2025 · 23.7M · ± 1.82k · ♥ 21

Training Framework and Recipe

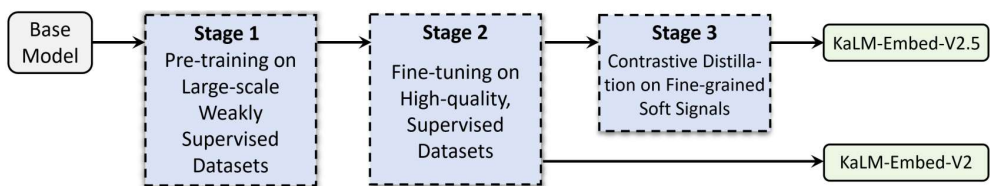
Overall Training Framework

- The left illustrates the workflow of contrastive learning, while the right shows that of contrastive distillation.



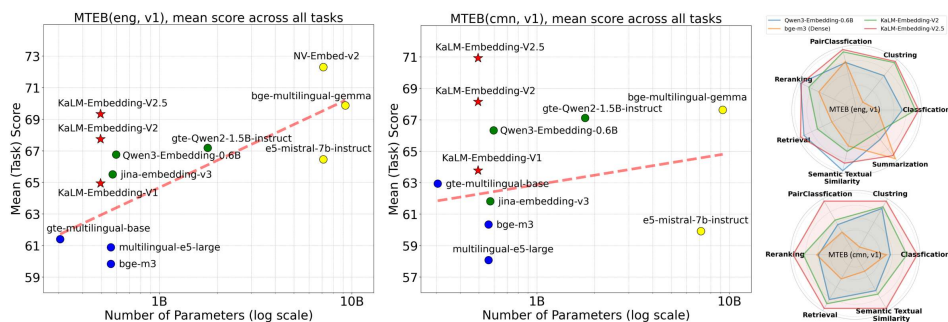
Progressive Multi-stage Training Recipe

- Pre-training on large-scale, weakly supervised datasets over 20 categories to endow the model with strong generalization.
- Fine-tuning on over 100 categories of high-quality supervised datasets to further improve the overall model performance.
- Contrastive Distillation with fine-grained soft knowledge from a stronger teacher model to capture nuanced semantic differences.



Experimental Results

- Our KaLM-Embedding-V2 series significantly outperforms models of comparable size, rivaling models 3–26x larger.



- KaLM-Embed-Gemma3 further advances SOTA on MMEB. KaLM-Embed-V2.5 shows superior intra-class compactness and inter-class separability.

Embedding Leaderboard

This leaderboard compares 180+ text and image embedding models across 100+ languages. We refer to the publication for details on metrics, languages, tasks, and task types. Anyone is welcome to add a model. See [help](#) or [update your own association](#) or [update other entries](#) in the leaderboard.

MTEB(Multilingual, v2)

A large-scale multilingual expansion of MTEB, driven mainly by highly-curated community contributions covering 250+ languages. Cite and share this benchmark.

Number of Languages: 1838
Number of Tasks: 121
Number of Task Types: 9
Number of Datasets: 20

Click for More Info

Rank (St.)	Model	Zero-shot	Memory Bn.	Number of P.	Embedding B.	Max Tokens	Mean (T.)	Max (Task)	Bitext	Classification	Clustering	Instruction R.	Multilabel Class.	Pair Classification	Ret.
1	KaLM-Embedding-Gemma3-1.7B-Inst	95%	4484	11.8	2060	32760	82.28	82.50	88.76	77.69	55.77	5.49	21.03	144.73	12
2	Time-vect-embedding-1B	95%	2829	7.5	4896	32760	81.46	81.89	85.72	73.21	54.25	58.82	29.86	83.97	17
3	OpenAI-Embedding-3	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15
4	OpenAI-Embedding-3.5	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15
5	OpenAI-Embedding-3.5	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15
6	OpenAI-Embedding-3.5	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15
7	OpenAI-Embedding-3.5	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15
8	OpenAI-Embedding-3.5	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15
9	OpenAI-Embedding-3.5	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15
10	OpenAI-Embedding-3.5	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15
11	OpenAI-Embedding-3.5	95%	14413	7.6	4896	32760	80.28	81.09	88.89	74.88	57.85	58.86	28.66	86.46	15

