# LLMs are Also Effective Embedding Models: An In-depth Overview

CHONGYANG TAO, SKLSDE Lab, Beihang University, China

TAO SHEN, University of Technology Sydney, Australia

SHEN GAO, University of Electronic Science and Technology of China, China

JUNSHUO ZHANG, University of Electronic Science and Technology of China, China

ZHEN LI, Peking University, China

KAI HUA, Peking University, China

WENPEN HU, Peking University, China

ZHANGWEI TAO, Peking University, China

SHUAI MA, SKLSDE Lab, Beihang University, China

Large language models (LLMs) have revolutionized natural language processing by achieving state-of-the-art performance across various tasks. Recently, their effectiveness as embedding models has gained attention, marking a paradigm shift from traditional encoder-only models like ELMo and BERT to decoder-only, large-scale LLMs such as GPT, LLaMA, and Mistral. This survey provides an in-depth overview of this transition, beginning with foundational techniques before the LLM era, followed by LLM-based embedding models through two main strategies to derive embeddings from LLMs. 1) Direct prompting: We mainly discuss the prompt designs and the underlying rationale for deriving competitive embeddings. 2) Data-centric tuning: We cover extensive aspects that affect tuning an embedding model, including model architecture, training objectives, data constructions, etc. Upon the above, we also cover advanced methods for producing embeddings from longer texts, multilingual, code, cross-modal data, as well as reasoning-aware and other domain-specific scenarios. Furthermore, we discuss factors affecting choices of embedding models, such as performance/efficiency comparisons, dense vs sparse embeddings, pooling strategies, and scaling law. Lastly, the survey highlights the limitations and challenges in adapting LLMs for embeddings, including cross-task embedding quality, trade-offs between efficiency and accuracy, low-resource, long-context, data bias, robustness, etc. This survey serves as a valuable resource for researchers and practitioners by synthesizing current advancements, highlighting key challenges, and offering a comprehensive framework for future work aimed at enhancing the effectiveness and efficiency of LLMs as embedding models.

CCS Concepts: • **Computing methodologies** → **Semantic networks**; • **Information systems** → **Document representation**; *Language models*; **Retrieval tasks and goals**.

## 1   Introduction

Representation learning is a key concept in deep learning, where models learn to capture meaningful features or patterns from raw data in a compressed, low-dimensional form, known as embeddings [11, 15]. In the context of information retrieval (IR), natural language processing (NLP) and computer vision (CV), representation learning is used to encode a piece of text or images into embedding vectors that capture the semantic meaning and syntactic structure of the input, enabling various downstream tasks such as classification [109], retrieval [48, 109], clustering [122, 123], anomaly detection [17, 39], reward model [149, 150], recommendation [90, 119, 195], and retrieval-augmented generation (RAG) [18, 43, 72]. Embeddings are crucial in modern deep learning literature because they enable efficient representation of high-dimensional data in a compact, dense format [59]. This compression not only reduces storage requirements but also allows for offline computation, which can then be easily used in real-time applications, including retrieval and recommendation systems, with online lightweight operations, e.g., dot-product similarity between embedded vectors. That is, embeddings preserve essential semantic and syntactic information, making it possible to perform complex operations like similarity comparison or clustering with significantly lower computational overhead.

Attributed to the parallelizable Transformer architecture [164] and the availability of high-performance computational resources, representation learning has been moved from shallow-contextualization word2vec [104] to a large-scale pre-training [28] era in the past years, where models trained on large-scale general corpora could generate generic embeddings that better capture both word- and sequence-level semantics meanings, and could be further enhanced through domain-specific or task-specific fine-tuning. Essentially, pre-trained Transformer encoders, including BERT [28], RoBERTa [94] and T5 enocder [131], outperform their RNN-based pioneering works, e.g., CoVe [101] and ELMo [121], in contextual representation learning, significantly improving performance across various NLP tasks like classification and semantic relatedness.

Nonetheless, pre-training Transformer encoders mainly depend on masked language modeling (MLM), where only a small proportion (e.g., 15% in BERT) of words or tokens are masked as learning objectives [28]. This deterministic learning process cannot fully perceive the rich contextual dependencies present in unmasked tokens, making these pre-training approaches less efficient and thus quickly reaching performance saturation even scaling model and data size. On the other hand, causal language modeling (CLM), learning to predict every next token given its preceding context, is more effective in utilizing both model parameters and training corpora to pre-train generative models [129], which is proven to have scaling law w.r.t. model performance. The resulting LLMs [14], mainly built upon the much deeper Transformer decoder architecture and pre-trained over trillions of tokens, have been proven to have emergent capabilities [173] in understanding and reasoning. The models, including GPT [2, 14], LLaMA [33, 159], Mistral [55] and Qwen [23, 194], achieve remarkable milestones in a broad spectrum of benchmarking tasks with excellent zero-/few-shot capability and state-of-the-art performance, including question-answering, coding, math, reasoning, dialogue generation, etc.

The success of LLMs has been extended to representation learning, which tunes the LLMs to generate expressively powerful embeddings [67, 98, 113, 168, 203]. Intuitively, the LLMs, with

significantly more parameters, larger pretraining corpora, and extended training durations, are expected to outperform the previous BERT-family models in capturing richer semantic representations and contextual nuances. This enhanced capacity enables LLMs to generate more accurate and generalizable embeddings, which is expected to improve performance across a wide range of downstream tasks compared to traditional encoder-only models.

In the context of how to leverage LLMs for representation learning, there are two distinct yet complementary perspectives. On the one hand, as LLMs can be viewed as extensions of previous encoder-based models like BERT, several established techniques from the BERT era are still applicable to LLMs. For instance, methods such as corpus-aware pre-training [75], multi-task learning [108], hard negative construction [168], and distillation from cross-encoder models [69] are still natural and practical to LLM-based approaches. Some of them have been confirmed in recent studies in terms of effectiveness, which demonstrates the utility of these techniques in fine-tuning LLMs for more robust embeddings. However, unique challenges arise when using LLMs for representation learning, particularly in determining how to extract efficient representations from a CLM [96]. Unlike BERT's natural use of the [CLS] token to generate embeddings [28], extracting useful representations from an LLM trained on CLM objectives [129] is less straightforward.

On the other hand, the powerful expressiveness of LLMs has opened up entirely new paradigms for representation learning, one of which is termed "Direct Prompt for Embedding" [56, 203]. Thanks to instruction-tuning, many LLMs are endowed with the ability to follow instructions [117, 155, 182]. This allows practitioners [56] to prompt the models to generate a specific category or a topic word, and then utilize its contextualized representations as the final embedding. This prompt-based approach presents exciting opportunities but also introduces several open research challenges. For example, the effectiveness and generalization of embeddings generated through direct prompting are still under exploration, with areas like in-context learning [76] and meta-learning [70] offering potential pathways for improving these embeddings in diverse tasks and domains.

With the growing reliance on LLMs in natural language understanding and information retrieval, this survey aims to provide timely and comprehensive insights into the paradigm shift in embedding methods, offering a deeper understanding of how embedding techniques with LLMs can enhance performance across a wide range of tasks, while also highlighting critical open problems in the field. The paper begins with a brief introduction of the foundational techniques that shaped the field before the LLM era. It then introduces two primary approaches: tuning-free embedding methods, which extract meaningful text embeddings from the hidden states of LLMs through direct prompting, without the need for explicit training on embedding-specific tasks; and tuning-based embedding methods, which involve continued supervised fine-tuning, focusing on optimizing model architecture, improving training objectives, and refining training data. Next, we summarize advanced techniques developed to handle longer texts, multiple languages, cross-modal data, and codes. We then compare the performance of various LLM-based embedding methods, considering key aspects in adapting LLMs as embedding models, such as contrasting dense versus sparse embeddings, assessing different pooling strategies, and exploring the implications of scaling laws as LLMs increase in size. Finally, the survey highlights the limitations and emerging challenges involved in adapting LLMs to be more effective as embedding models.

## 2 Background

Over the past decade, the paradigms for representation learning have shifted multiple times because of the rapid advancements in neural architectures, and the availability of large-scale datasets and computational resources: the first being shallow contextualization (e.g., word2vec and specific tuning), the second marked by BERT's pre-training methods, and now, the third transition towards

LLMs as embedding models. This section provides a brief overview of foundational knowledge and techniques that shaped the field before the LLM era.

## 2.1    Shallow Contextualization

Initially, word-level representations were learned using shallow models like word2vec [103, 104], GloVe [120], and FastText [13], which employed techniques such as skip-gram or continuous bag-of-words (CBoW) to capture context, with sequence-level representations typically obtained through a weighted sum of word vectors, where the weights were calculated using heuristic methods like TF-IDF [27]. However, these methods' shallow architecture limits their ability to fully leverage contextual information and struggle with capturing polysemy and ambiguity, leading to unsatisfactory performance in a wide range of NLP and IR tasks. On the other hand, some follow-up works constructed unsupervised training objectives at sequence (e.g., sentence) level, such as Skip-Thought [63]. Later, contextualized models like CoVe [101] and ELMo [121] addressed these limitations by using bi-directional LSTMs to capture richer contextual information and generate word representations that adapt to the surrounding context. ELMo, in particular, significantly improved tasks requiring nuanced semantic understanding due to its contextual flexibility.

## 2.2    Prominent Techniques in BERT Era

The advent of BERT [28] and its successors, such as RoBERTa [94] and T5-encoder [131], marked a significant leap in representation learning, leveraging large-scale pre-training to capture deep contextualized embeddings that could be adapted to different domains and tasks. In the following, we will introduce several prominent techniques to improve embedding quality, which have the potential to benefit LLMs as embedding models.

**Corpus-aware pre-training.** These pre-trained models are all trained on general text corpora, but their performance on specific target domains is often suboptimal. A straightforward approach to improve domain-specific performance is to do continual pre-training over the corpus in the target domain [40, 51]. The training objectives can be quite flexible, ranging from general tasks like masked language modeling or next sentence prediction to more specialized methods that generate pseudo-contrastive learning examples through heuristics. Furthermore, research has shown that a novel bottleneck-enforced pre-training approach can be effective [41, 139, 166, 178], where the original encoder structure is retained but a weaker decoder is introduced. This decoder reconstructs the original input based on a bottleneck, i.e., a single-vector embedding from the encoder, which forces semantic knowledge to be retained in the embedding.

**Hard negative mining.** Negative samples are essential in contrastive learning, as they are examples that, compared to the positive sample, are far from the anchor according to a specific metric. In contrast to random negative samples, learning with hard negatives has proven to be highly effective in representation learning for embeddings [140, 154, 166, 181]. Intuitively, hard negatives, which are closer to the anchor but belong to different classes, force the model to make more refined distinctions, thus leading to more robust and accurate embeddings. Therefore, how to construct or sample hard negatives has remained a popular research topic, with methods like using a BM25 retriever or the strongest available retriever to select challenging negatives from large-scale collections.

**Supervision from re-ranker.** In model training, knowledge distillation from a teacher (stronger or larger) model to a student model (weaker or smaller) has been proven effective in various scenarios and tasks of deep learning [183]. In the field of embedding, one unique opportunity is to distill knowledge from a cross-encoder-based re-ranker to a bi-encoder-based encoding/embedding model by applying a KL divergence loss between the score distribution over several candidates

against their anchor [127]. Here, a cross-encoder [200] directly computes interaction between every pair of input sequences, resulting in richer contextual embeddings but with higher computational cost, while a bi-encoder [59] independently encodes two sequences and computes their similarity in the embedding space, offering a more efficient yet slightly less expressive representation. The supervision from the cross-encoder allows the bi-encoder to learn more nuanced distinctions and improve its embedding quality. Furthermore, studies have shown that mutual supervision or distillation between cross- and bi-encoders can simultaneously lead to boosted re-ranking performance for the cross-encoder and enhanced encoding for the bi-encoder, creating a mutual learning framework that benefits both models [16, 37, 133].

**Multi-task learning.** Multi-task learning is a machine learning approach where a model is trained to perform multiple tasks simultaneously, leveraging shared information across tasks to improve its overall learning efficiency and generalization [91, 131, 136]. This approach is effective because it allows the model to capture commonalities between tasks while also learning task-specific nuances, leading to more robust performance across diverse scenarios. In the field of embedding models, multi-task learning has shown great potential for improving generalization and robustness by enabling models to learn representations that are not only effective for a single task but adaptable across various tasks [16, 100]. This adaptability helps the model capture richer and more diverse semantic features, making the embeddings more versatile for downstream applications like classification, retrieval, and clustering. Going beyond and following the inspiration of prompt-tuning [71, 92], some pioneering works [146, 179, 192] propose to unify a broad spectrum of retrieval tasks by augmenting a piece of text with an instruction/explanation of the corresponding task or domain, which is still the most prevalent paradigm up to now.

In the remaining sections, starting with a task formulation, we will explore extensions and variants of prior techniques in the LLM era, as well as advancements and challenges in utilizing LLMs as effective embedding models. In general, we will delve into both direct prompting and fine-tuning approaches, and provide insights into their performance across various downstream tasks.

## 3 Problem Formalization and Survey Overview

Let $\mathbf{x} = \{x_1, x_2, \ldots, x_n\}$ be an input sequence of tokens where $x_i$ represents the $i$-th token in the sequence. The decoder-only LLM $\mathcal{M}$ processes this sequence to produce a contextualized representation for each token, denoted as

$$\{\mathbf{h}_1, \mathbf{h}_2, \ldots, \mathbf{h}_n\} = \mathcal{M}(\mathbf{x}, \theta),$$
$$\mathbf{v} = f_{\text{agg}}(\{\mathbf{h}_1, \ldots, \mathbf{h}_n\}),$$

where $\mathbf{h}_i$ is the output hidden state or logit distribution for token $x_i$, a fixed-length embedding vector $\mathbf{v}$ for the entire sequence can be obtained by the aggregation function $f_{\text{agg}}$, such as special-token selection (e.g., [EOS]), mean pooling, or attention-based weighting.

This formulation provides a general framework for using decoder-only LLMs as embedding models. However, practical instantiations of this process vary widely depending on how the model is used, adapted, or trained. As illustrated in Figure 1, recent research has proposed three major paradigms for generating embeddings from LLMs: direct prompting, fine-tuning, and specialization. First, we explore how LLMs can be directly prompted to produce embeddings, with a focus on prompt design and embedding derivation strategies that require no parameter updates. Second, we examine tuning-based methods that convert LLMs into dedicated embedding models through architectural adaptations (e.g., bi-attention, low-rank adaptation), refined embedding extraction mechanisms, and training objectives such as multi-task learning, multi-stage optimization, and knowledge distillation.
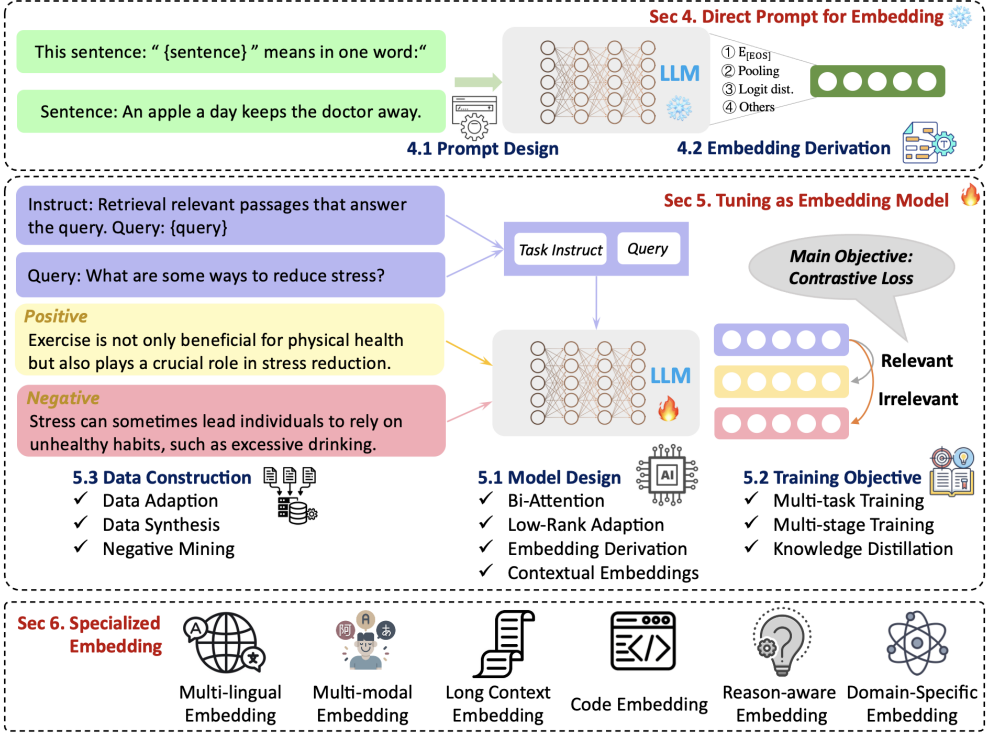
Fig. 1. An overview of our survey on deploying LLMs as embedding models, covering both ① direct prompting as embedding models, ② fine-tuning as embedding models and ③ specialized embedding mdels.

These approaches are supported by carefully constructed datasets, leveraging both adaptation of existing data and synthesis via LLMs, often with hard negative mining to improve contrastive learning. Finally, we discuss specialized embeddings developed for multilingual, multimodal, long-context, code, reason-aware, and domain-specific scenarios, which address real-world needs and highlight the flexibility of embedding models in capturing complex semantics beyond general-purpose tasks.

In addition to these paradigms and techniques, this survey examines critical considerations in adapting LLMs as embedding models, including the comparison of dense and sparse embeddings, evaluation of various aggregation strategies, and analysis of scaling law implications. Furthermore, we also highlight open challenges and future directions, shedding light on key limitations and potential research opportunities in LLM-based embedding techniques.

## 4  Tuning-Free Embedding Models

Previous masked language models, like BERT and RoBERTa, employ a mask prediction task to capture contextual information for specific tokens. Building on this factual, PromptBERT [57] frames the text embedding extraction as a similar task, utilizing the following template for prompting: `This sentence : "[X]" means [MASK] .` In this template, [X] and [MASK] represent placeholders for the input sequence and the mask token, respectively. The hidden vector of the [MASK] token from the final layer is directly used as the representation of the text sequence.

<div style="border:1px solid #000; padding:10px;">

**PromptEOL** [56]
```
This sentence : "[X]" means in one word:"
```

**PromptSTH** [190]
```
This sentence : "[X]" means something
```

**PromptSUM** [190]
```
This sentence : "[X]" can be summarized as
```

---

**Pretended Chain of Thought** [190]
```
After thinking step by step , this sentence : "[X]" means in one word:"
```

**Knowledge Enhancement** [190]
```
The essence of a sentence is often captured by its main subjects and actions, while
descriptive terms provide additional but less central details. With this in mind ,
this sentence : "[X]" means in one word:"
```

</div>

Fig. 2. Prompts used in various tuning-free methods.

Recently, researchers have started to explore the potential of extracting meaningful text embeddings directly from the hidden states of LLMs (e.g., OPT or LLAMA) without requiring explicit training on embedding-specific tasks. These embeddings have been applied in various tasks, including clustering [122], recommendation [83, 90, 119], and retrieval [70, 203]. These methods typically involve using fill-in-the-blank prompts. Inspired by PromptBERT, [56] introduced PromptEOL, which enhances the prompt-based method in BERT by incorporating an "explicit one word limitation" (EOL) to extract text representations for LLMs. As shown in the first line of Figure 2, PromptEOL incorporates ": "" at the end of the template to prevent the model from generating punctuation in the next token. Additionally, it considers an in-context learning framework to automatically create and search for demonstration sets to improve embeddings in LLMs. This method improves performance across all OPT models, allowing them to match or even outperform BERT in prompt-based embedding tasks.

More recently, [190] reveals that EOL does not consistently yield optimal performance when fine-tuning generative models. They introduced two templates, PromptSTH and PromptSUM, which intentionally omit the "in one word" constraint, along with two powerful prompt engineering methods, Pretended CoT and Knowledge Enhancement, to enable the model to analyze diverse semantic aspects, as shown in Figure 2. To effectively capture multiple representations of sentences from distinct perspectives, [70] introduced MetaEOL, which utilizes a diverse set of meta-task prompts, including text classification, sentiment analysis, paraphrase identification, and information extraction, to generate embeddings. [157] further present GenEOL, which leverages LLMs to generate diverse transformations of a sentence that retain its original meaning, and then aggregates the embeddings of these transformations to enhance the final sentence embedding. [203] directly prompt LLMs to generate both dense embedding representations and sparse bag-of-words representations for document retrieval. [143] propose echo embeddings that repeat the input twice in context and extract embeddings from the second occurrence. [87] introduced MOE Embedding (MOEE), which combines the routing weights and hidden states of Mixture-of-Experts LLMs to form a powerful zero-shot embedding model that surpasses traditional approaches without requiring any additional fine-tuning.

## 5   Tuning-based Embedding Models

While direct prompting of LLMs can yield useful embeddings, the true potential of these models is unlocked through a more refined approach: tuning them with existing or synthetic paired text data. This process typically employs contrastive learning to enhance the models' ability to distinguish between semantically similar and dissimilar text pairs, resulting in more accurate and meaningful embeddings. For example, in the early stage of LLMs outbreak, [107] initialized the embedding models with pre-trained GPT-3 models and applied continued contrastive training. The hidden state from the last layer corresponding to the special token [EOS] at the end of the sequence is taken as the embedding of the input sequence. Instructor [146] is the first work to explore a unified approach for generating text embeddings using task instructions based on GTR [112]. It annotates instructions for 330 diverse tasks and trains the model on this multitask mixture with a contrastive loss.

With the continuous progress of open-source LLMs, many researchers have attempted to use these models to build better text representation models. These advances primarily focus on optimizing model architecture [69, 108], improving training objectives [10, 67, 69, 86], and further enhancing training data [52, 62, 168]. In this section, we will delve into these three foundational techniques that underpin the fine-tuning of LLMs as powerful representation models. Table 1 summarize summarizes the configurations of representative models in terms of these three aspects.

### 5.1   Model Architecture



Fig. 3. An illustration of the evolution of text embeddings from ELMo and BERT to GPT-style models, from [188].

The development of language models has seen significant changes in neural architecture design. As illustrated in Figure 3, early models like ELMo [121] and BERT [28] primarily utilized encoder-only architectures, which focus on bidirectional context, enabling the model to consider both left and right context during training. In contrast, recent LLMs have adopted decoder-only architectures [55, 116, 160], which are driven by the need for models to excel at generative tasks, such as text completion and creative writing. While encoder-only models provide strong contextual embeddings, decoder-only models have demonstrated superior performance in generative tasks due to their autoregressive nature.

**Embedding with Bi-directional Contextualization.** Current generative LLMs predominantly use mono-directional attention, focusing on a unidirectional flow of information. This approach simplifies the model architecture and aligns well with autoregressive tasks where future tokens are predicted based on past context. However, the lack of bidirectional attention can limit the model's ability to fully capture dependencies within the entire sequence. To address this, some models like Gecko [69] and LLM2vec [10] propose incorporating bidirectional attention mechanisms within

Table 1. Detail of representative models that tune LLMs as embedding models. "sup.CL" and "unsup.CL" means supervised contrastive loss and unsupervised contrastive loss respectively. "Bi-Att" indicates whether the model enables bi-directional attention in LLMs, while "LoRA" specifies whether the model was trained using the LoRA technique. For the manner of embedding derivation, "MP" is mean pooling operation and "P-MP" means Position-weighted mean pooling, "LAT" means Latent Attention Layer.

| Model | Dim. | # Params. | Base | Fusion | Bi-Att | LoRA | Training Data | Training Obj. (Neg) |
|---|---|---|---|---|---|---|---|---|
| SGPT [107] | 4,096 | 5.8B | GPT | P-MP | ✗ | ✗ | SNLI/MNLI | sup.CL |
| Instructor-XL [146] | 768 | 1.5B | GTR | MP | ✗ | ✗ | MEDI | sup.CL |
| GTE-Qwen2-7B [86] | 3,584 | 7B | Qwen2 | MP | ✓ | ✗ | Unsup/Sup Pair | unsup.CL & sup.CL |
| E5-mistral-7b [168] | 4,096 | 7B | Mistral | [EOS] | ✗ | ✓ | E5(Public&Synthetic) | sup.CL |
| GritLM-7B [108] | 4,096 | 7B | Mistral | MP | ✓ | ✗ | E5S(E5&S2ORC)/Tülu 2 | sup.CL & NTP |
| Echo-mistral-7b [143] | 4,096 | 7B | Mistral | [EOS] | ✗ | ✓ | E5 | sup.CL |
| LLM2Vec [10] | 4,096 | 8B | Llama3 | MP | ✓ | ✓ | Wikipedia | NTP & unsup.CL |
| SFR-Embedding | 4,096 | 7B | E5 | [EOS] | ✗ | ✗ | Specially Curated Dataset | sup.CL |
| NV-Embed [67] | 4,096 | 7B | Mistral | LAT | ✓ | ✓ | Public Sup Datasets | Two-stage sup.CL |
| Linq-Embed-Mistral [62] | 4,096 | 7B | E5 | [EOS] | ✗ | ✓ | E5S&Refined Synthetic | sup.CL |
| Gecko [69] | 256 | 1.2B | Transformer | MP | ✓ | ✗ | FRet&Public Datasets | sup.CL |
| NV-Retriever [105] | 4096 | 7B | Mistral | MP | ✓ | ✓ | Public Sup Datasets | Sup CL |
| BGE-ICL [76] | 4,096 | 7B | Mistral | [EOS] | ✗ | ✓ | Public Sup Datasets | sup.CL & KD |
| Gemini-Embedding [68] | 3072 | UNK | Gemini | MP | ✗ | ✗ | Pub & Syn | sup.CL |
| Qwen3-Embedding [194] | 4096 | 4B/8B | Qwen3 | [EOS] | ✗ | ✗ | Pub & Syn | sup.CL |

existing LLMs, enabling the model to consider both past and future tokens simultaneously. This enhancement aims to improve the quality of sequence embeddings by leveraging a more comprehensive understanding of the input text. In addition, GritLM [108] unifies embedding tasks and generative tasks into a single model with bidirectional attention through generative representational instruction tuning. Recently, however, BGE-ICL [76] argue that enabling bidirectional attention during embedding fine-tuning misalign with the model's original pre-training setup, which could compromise its effectiveness in generation tasks.

**Low-Rank Adaption.** Traditional fine-tuning methods often require extensive computational resources and large amounts of labeled data. As a result, there is a growing interest in parameter-efficient tuning approaches (e.g., LoRA) that enable the adaptation of LLMs for embedding tasks while minimizing resource consumption [62, 76, 168]. These parameter-efficient tuning methods offer promising solutions by reducing the number of trainable parameters while maintaining comparable model performance or even better generalization performance [168].

**Embedding Derivation.** Obtaining sequence embeddings from these architectures involves different techniques. In BERT, the output states of [CLS] token is commonly used as a representation of the entire input sequence, providing a summary embedding. For LLMs, the [EOS] token, which signifies the end of a sequence, serves a similar purpose. Additionally, mean pooling of the last hidden layer can be employed to aggregate the contextual information across all tokens in the sequence [10, 69, 143, 146].Recent studies have explored more advanced methods to obtain sequence representations, such as using latent attention layer [67] that adaptively combine token embeddings to enhance the final sequence representation, or sparse representations based on lexicon-importance distribution [19, 29, 97, 203].

**Contextual Expansion.** Compared with previous text embedding models which often only consider the text itself, some recent efforts expand the context of query or document to obtain enhanced text semantic information embedding. [76] incorporated few-shot examples on the query side to enhance query embeddings during supervised contrastive training. [106] designed
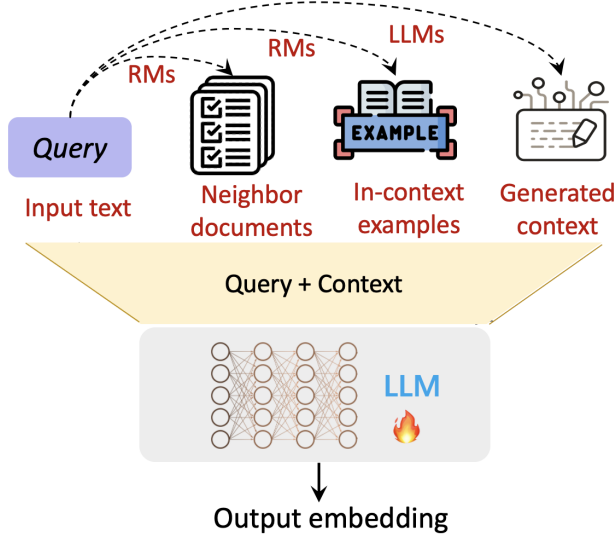
Fig. 4. An overview of contextual expansion for embedding models. RMs means retrieval models.

the alternative contrastive learning objective and the new contextual architecture to incorporate neighbor document information into document embeddings. [185] generated thoughts towards the input query, then obtained the embedding of the query and generated thought separately, and finally incorporated both embeddings to produce the embedding that better reflects the query semantics. [54] introduced the Chain-of-Deliberation mechanism, which encouraged LLM to encode documents with conducting reasoning process, enabling to generate more fine-grained document embeddings. [110] generated potential user queries based on each document, and through incorporating the embeddings of generated queries with the document embeddings, it's more effective to capture diverse user queries that reference the same document in different ways. [156] proposed a simple yet effective method that enhances retriever models by training them with semantically similar in-context examples, enabling better utilization of contextual signals for improved retrieval performance. [74] proposed a Reinforced Information Retrieval framework in which the query-expansion model and embedding model mutually enhance each other. Expanded queries generate more informative embeddings, while refined embeddings guide more effective query expansions, creating a synergistic loop that improves retrieval accuracy. An overview of these contextual expansion techniques and their integration strategies within embedding-based retrieval models is illustrated in Figure 4.

**Matryoshka Embedding (ME).** Inspired by Russian nesting dolls, where smaller parts nest within larger ones. Matryoshka Embeddings [65] are trained to store the most important information in the early dimensions of an embedding. This allows the embeddings to be truncated to smaller sizes (e.g., 64, 128, 256 dimensions) while still retaining useful information. As shown in Figure 5, by applying loss functions to multiple truncated versions during training, ME enables flexible trade-offs among performance, speed, and memory usage. It's especially useful for tasks like retrieval or classification, where quick filtering can be done with smaller embeddings, followed by re-ranking with full-size ones. More recently, The Qwen3 Embedding model series [194] employs the Matryoshka Embedding training approach and offers flexible, user-defined output dimensions
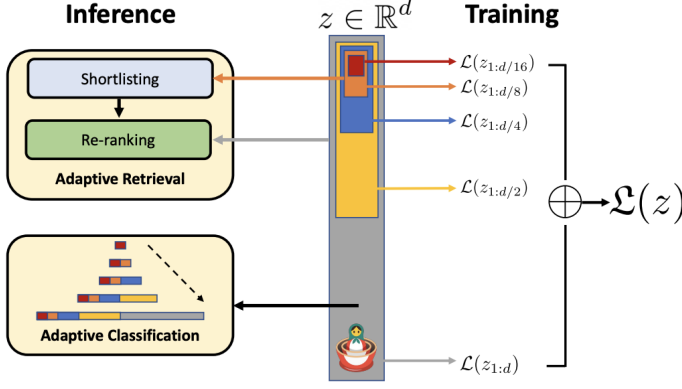
Fig. 5. An illustration of Matryoshka Embedding, from [65].

ranging from 32 to 4096. This enables the model to better balance performance and computational efficiency across different applications.

## 5.2 Training Objective

Most existing methods for fine-tuning LLMs as embedding models adopt contrastive learning [20]. The core idea behind contrastive learning is to bring representations of similar sequences closer in the embedding space while pushing dissimilar sequences apart, thereby learning robust and distinctive sequence embeddings. Formally, given the query $q_i$ and each of the candidate passage $p_j$, we can utilize a LLM to produces contextualized token embeddings for a sequence of tokens, then we can obtain the vector representations of $q_i$ and $p_i$ by taking the embedding of [EOS] token or averaging the token embeddings. To incorporate task-specific context, a task instruct $t$ is prepended to each query. The contrastive training process for the query $q_i$ is formulated as the minimization of the following InfoNCE loss:

$$\mathcal{L} = -\log \frac{\exp^{\mathcal{S}(q_i, p_i^+)/\tau}}{\exp^{\mathcal{S}(q_i, p_i^+)/\tau} + \sum_{p_{i,j}^- \in \mathbb{P}_i^-} \exp^{\mathcal{S}(q_i, p_{i,j}^-)/\tau}}.$$

where $p_i^+$ is the labeled positive document paired with $q_i$ and $\mathbb{P}_i^-$ denotes the set of candidate documents for $q_i$ which is typically constructed during training by random negative sampling or hard negative mining methods. $\mathcal{S}(q_i, p_i)$ denotes the similarity measure (e.g., cosine similarity) between the query and the positive. $\tau$ is a temperature parameter that influences the sharpness of the similarity distribution. By minimizing this contrastive loss, the model learns to bring embeddings of positive pairs closer together while pushing embeddings of negative pairs further apart. This process aids in creating more effective and discriminative text embeddings, enhancing the model's performance on a range of downstream tasks.

**Multi-task/stage Training.** Besides the standard contrastive objective, much effort has been made to improve the training procedures to enhance LLMs as versatile embedding models. For example, GTE [86] and Conan-embedding [85] introduce a multi-stage training approach in which the model is first pretrained with an InfoNCE loss using in-batch negatives on weakly supervised text relevance data, followed by fine-tuning with a CoSENT loss on the STS task. LLM2Vec [10] transforms a pre-trained decoder-only LLM into a universal text encoder through masked next

token prediction and unsupervised contrastive learning. This method uses only publicly available data and applies unsupervised contrastive training similarly to SimCSE [42]. NV-Embed [67] also focuses on training procedures to enhance LLMs. Specifically, the model first performs contrastive training with instructions on retrieval datasets and then integrates carefully curated non-retrieval datasets into the stage-one training data. GritLM [108] proposes a unified model for both embedding tasks and generative tasks, which is jointly optimized with both NLL objective and contrastive loss. [194] also employs a multi-stage training pipeline, which includes weakly supervised training using large-scale synthetic paired data (over 150 million pairs), followed by supervised fine-tuning with high-quality synthetic and labeled data.

**Knowledge Distillation.** Additionally, several studies explore using knowledge distillation to enhance embedding performance by leveraging larger embedding models or cross-attention ranker models. Gecko [69], for instance, builds a smaller bidirectional embedding model (with 1.2B parameters) by distilling knowledge from a decoder-only LLM by generating synthetic paired data. BGE-ICL [76] introduce few-shot examples into the query side to enhance the query embeddings and also consider distilling the relevance score from the reranker for retrieval task during training. Zhang et al. [192] introduce a unified embedding model designed to support the diverse retrieval augmentation needs of LLMs, including document knowledge, tools, in-context learning examples, and memory knowledge, through multi-task learning. The model incorporates a reward formulation based on LLM feedback and stabilizes knowledge distillation by integrating both soft reward-based labels and hard ranking-based labels during contrastive training. Lupart et al. [97] propose distilling the scores of rewritten queries paired with documents from a teacher model and utilizing knowledge distillation from multiple teachers to enhance conversational representations in conversational search tasks.



Fig. 6. An illustration of various data construction methods for training LLM-based embedding models. Here, $d^+$ denotes the positive document, and $\{d^-\}$ represents the set of hard negative documents. The variable $m$ refers to diverse meta-information used to guide synthetic data generation. $I$ is the instruction provided to the LLM for generating either a query $q$ or a document $d$, while $c$ optionally includes few-shot demonstrations. Finally, $t$ denotes the task-specific instruction.

## 5.3   Data Construction

Training data plays a crucial role in the success of fine-tuning LLMs as embedding models. The methods typically employ two types of data: (1) adapting existing datasets, and (2) synthesizing

Table 2. Comparison of datasets used for supervised fine-tuning of several representative embedding models. QuestionRR means StackOverFlowDupQuestionRR. "−" indicates that the model uses this type of data, but specific datasets are not specified.

| TASK | Dataset | E5-Mistral | GritLM | Linq | SFR | Gecko | LLM2Vec | BGE-en-ICL | NV-Embed | NV-Retriever | MTEB |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Retrieval | DuReader | ✓ | ✓ | ✓ | × | × | ✓ | × | × | × | × |
| | ELI5 | ✓ | ✓ | ✓ | × | × | ✓ | ✓ | × | × | × |
| | FEVER | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | HotpotQA | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | NLI | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | × |
| | MIRACL | ✓ | ✓ | ✓ | × | ✓ | ✓ | × | × | × | × |
| | MSMARCO | ✓ | ✓ | ✓ | ✓ | × | ✓ | ✓ | ✓ | ✓ | × |
| | Mr.TyDi | ✓ | ✓ | ✓ | × | × | ✓ | × | × | × | × |
| | NQ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | S2ORC | × | ✓ | ✓ | × | × | × | × | × | × | × |
| | SQuAD | ✓ | ✓ | ✓ | × | × | ✓ | ✓ | ✓ | ✓ | × |
| | T2Ranking | ✓ | ✓ | ✓ | × | × | ✓ | × | × | × | × |
| | TriviaQA | ✓ | ✓ | ✓ | × | × | ✓ | ✓ | × | ✓ | × |
| | Quora | ✓ | ✓ | ✓ | × | × | ✓ | ✓ | × | × | ✓ |
| | FiQA | × | × | × | ✓ | × | × | ✓ | ✓ | ✓ | ✓ |
| | SciFact | × | × | × | ✓ | × | × | × | × | ✓ | ✓ |
| | NFCorpus | × | × | × | ✓ | × | × | × | × | ✓ | ✓ |
| | DBPedia | × | × | × | ✓ | × | × | × | × | × | ✓ |
| | MedMCQA | × | × | × | × | ✓ | × | × | × | × | × |
| | PAQ | × | × | × | × | × | × | × | ✓ | ✓ | × |
| | Stackexchange | × | × | × | × | × | × | × | ✓ | ✓ | ✓ |
| | ArguAna | × | × | × | × | × | × | ✓ | ✓ | ✓ | ✓ |
| | GOOAQ | × | × | × | × | × | × | × | × | ✓ | × |
| | BioASQ | × | × | × | × | × | × | × | ✓ | ✓ | × |
| Clustering | arXiv | × | × | × | ✓ | × | × | ✓ | × | ✓ | ✓ |
| | Reddit | × | × | × | × | × | × | ✓ | × | × | ✓ |
| | StackExchange | × | × | × | × | × | × | ✓ | × | × | ✓ |
| | bioRxiv | × | × | × | ✓ | × | × | ✓ | × | ✓ | ✓ |
| | medRxiv | × | × | × | ✓ | × | × | ✓ | × | ✓ | ✓ |
| Classification | AmazonReview | × | × | × | ✓ | - | × | ✓ | ✓ | ✓ | ✓ |
| | Emotion | × | × | × | ✓ | - | × | ✓ | ✓ | ✓ | ✓ |
| | MTOPIntent | × | × | × | ✓ | - | × | ✓ | ✓ | ✓ | ✓ |
| | ToxicConversation | × | × | × | ✓ | - | × | ✓ | ✓ | ✓ | ✓ |
| | TweetSentiment | × | × | × | ✓ | - | × | ✓ | ✓ | ✓ | ✓ |
| | AmazonCounterfactual | × | × | × | ✓ | - | × | ✓ | ✓ | ✓ | ✓ |
| | Banking77 | × | × | × | × | - | × | ✓ | ✓ | ✓ | ✓ |
| | IMDB | × | × | × | × | - | × | ✓ | ✓ | ✓ | ✓ |
| STS | STS12 | × | × | × | ✓ | × | × | ✓ | × | ✓ | ✓ |
| | STS22 | × | × | × | ✓ | × | × | ✓ | × | ✓ | ✓ |
| | STSBenchmark | × | × | × | ✓ | × | × | ✓ | × | ✓ | ✓ |
| Reranking | SciDocs | × | × | × | ✓ | × | × | ✓ | × | × | ✓ |
| | QuestionRR | × | × | × | ✓ | × | × | ✓ | × | × | ✓ |
| Synthesic | - | | ✓ | ✓ | ✓ | × | ✓ | × | × | × | × |

datasets via LLMs. The latter often involves both data curation and automated labeling to construct training pairs. Figure 6 illustrates representative methods for constructing contrastive training datasets under these paradigms.

**Adapting Existing Datasets.** Table 2 presents the most commonly used public datasets, which include retrieval, clustering, classification, STS, and reranking tasks that align with the evaluation criteria in the Massive Text Embedding Benchmark (MTEB)[1] [109]. These datasets provide diverse

---

[1]https://huggingface.co/spaces/mteb/leaderboard

labeled examples that can be used to fine-tune models in a supervised manner, enhancing their ability to capture specific semantic relationships and improve performance on downstream embedding tasks. Specifically, these tasks are typically divided into symmetric and asymmetric tasks [168]. Symmetric tasks consist of queries and documents that convey similar meanings but differ in their wording, while asymmetric tasks include tasks where the query and document share a semantic relationship without being direct paraphrases. Besides paired data, some works leverage self-supervised learning methods like SimCSE, using natural language text from sources such as Wikipedia [10]. These tasks often involve generating different views of the same input sequence through augmentations, such as dropout or token masking, to create positive pairs for contrastive learning. Moreover, SFR-Embedding-Mistral [102] fine-tunes the E5 model across a diverse array of tasks, and demonstrates that the effectiveness of embedding models can be enhanced through knowledge transfer from multiple tasks.

**Synthesizing Datasets via LLMs** Considering the limited quantity of supervised data and the uneven distribution of corresponding pair relevance, recent research has addressed this challenge by leveraging LLMs to automatically synthesize large-scale, diverse, and contextually rich datasets [168]. This approach draws inspiration from LLM-based data augmentation in the supervised fine-tuning process [183]. Synthetic datasets for training embedding models typically involve two primary methods: data labeling and data curation. In the labeling approach, a query is provided for the LLM to generate a relevant document as a positive example[168], optionally producing hard negative documents[168], or to generate a query[69, 194] for a sampled document from the corpus. In the curation method, the LLM is prompted to generate both the query and the relevant document, utilizing few-shot query-document examples as optional information. By using model-generated data, researchers can produce large-scale, high-quality training datasets tailored to specific embedding tasks, improving the robustness and generalizability of the fine-tuned models.

In the early stage of LLM area, [52, 53] utilize GPT-J to generate one synthetic query per given document, prompted with three examples from MS MARCO, and then leverage a powerful reranker to select the synthetic query-document pairs for training. E5-Mistral [168] introduce a two-step prompting method in which proprietary LLMs (GPT-4) are initially prompted to create a pool of candidate tasks. From this pool, a specific task is chosen, and the LLMs are prompted again to generate triples consisting of a query, a positive document, and a hard-negative document. This method ultimately supports a broad range of text embedding tasks across 93 languages, covering hundreds of thousands of embedding tasks. Gecko [69] prompted LLMs to read a sampled passage from a web corpus to generate both a task description and a relevant query, after which the top-K relevant passages were retrieved for these synthetic queries, and LLMs were employed to relabel more relevant positive passages along with a better hard negative. Linq-Embed-Mistral [62] utilized GPT-4-turbo to generate the query-positive-negative triplet as training samples across six tasks, where the synthetic strategy was closely adhering to the [168]. Additionally, [62] conducted a data refinement on synthetic data to improve overall data quality through employing various types of few-shot prompt engineering, filtering, and negative mining. [61] explored the zero-shot capabilities of LLMs for generating queries, relevance judgments, and reference lists. [68] leveraged multi-stage prompting strategies to generate training samples across retrieval and classification tasks. For retrieval, Gemini was prompted to generate queries for web passages followed by a Gemini auto-rater to filter lower-quality examples and for classification, counterfactual, sentiment, and review classification training samples were generated. [194] prompted Qwen3-32B to generate queries for the documents by employing diverse prompting strategies, where specific roles were assigned to each document with incorporating various generative dimensions including query type, length, difficulty, and language.

**Negative Mining.** Negative mining is crucial for constructing high-quality training data for embedding models, as appropriate negative samples enhance the representation capacity of embedding models trained with contrastive loss. Consistent with previous representation learning methods [59, 181], current LLM-based embedding models are exploring various methods to mine negative to improve the training process. A naive negative mining technique is using positive document examples from other queries within the same training batch as negatives [20]. While efficient, this method often yields overly simple and uninformative negatives, which may hinder effective contrastive learning. To address this limitation, recent works focus on constructing hard negatives using two primary methods: selection and generation.

The selection method involves identifying hard negatives from the candidate corpus. [146] employed Sentence-T5 embeddings [111] to encode pairs of texts, calculating pairwise cosine similarity scores to determine suitable hard negatives. [69] retrieved the top-K relevant passages for synthetic queries and utilized LLMs to relabel more relevant positive passages and identify better hard negatives. [85] introduce dynamic hard negative mining, which selects more challenging negative examples during contrastive training. During the selection process, some studies have noted the issue of false negatives, where certain hard negatives should actually be classified as positive examples [127]. To address this, [105] examined various filtering methods for hard negatives, proposing the use of similarity scores from positive examples as thresholds with specific margins or percentages.

In contrast, the generation method leverages LLMs to directly create hard negatives for queries. [168] utilized LLMs to generate query, positive, hard-negative triples based on a selected task from pre-generted candidate tasks pool. This synthetic strategy has been applied to various embedding models [62, 108, 143], demonstrating its versatility. Building on the approach of [168], [62] introduced a data refinement process, such as such as few-shot prompt engineering, filtering, and improved negative mining, to enhance the quality of synthetic data. This integration of hard negative mining techniques directly exploits the powerful generative capabilities of LLMs, ultimately facilitates more robust representational learning, enhancing the capabilities of embedding models.

## 6 Specialized Embeddings

Beyond regular text embeddings, advanced techniques have been developed to handle longer texts, multiple languages, cross-modal data, and programming code. This section delves into these specialized embedding techniques, highlighting their methodologies and applications.

### 6.1 Multi-lingual Embedding

Multi-lingual embeddings are designed to represent text from different languages in a shared vector space. These embeddings are useful for various tasks such as cross-lingual information retrieval, machine translation, and multi-lingual sentiment analysis. The development of monolingual language models has substantially progressed in learning embeddings enriched with contextual information across a range of specific languages.

Early pretrained LMs like BERT [28] and GPT [129] learn fine-grained contextual monolingual representations through masked language modeling (MLM) and next-token prediction (NTP). Built on the foundation of monolingual LMs, a series of advanced works have been proposed to learn universal multi-lingual embeddings. We categorize these approaches based on their alignment techniques into two types: implicit alignment and explicit alignment.

*6.1.1 Implicit alignment.* Instead of direct translations or aligned sentence pairs, implicit alignment leverages shared representations and latent relationships within the training data. By utilizing techniques such as shared vocabularies, cross-lingual embeddings, and self-supervised learning,

multilingual LLMs can well capture semantic similarities and syntactic structures across diverse languages and implicitly align different languages into a unified semantic embedding space. This approach enhances the model's ability to perform tasks in low-resource languages, where paired data may be scarce, ultimately leading to more robust and versatile multilingual capabilities.

Multi-lingual Language Models such as mBERT [124], XLM-R [26], mT5 [184] have explored the capabilities of language models on multilingual NLP tasks. Recently, with model parameters and training data scaled up, numerous multilingual large language models (MLLMs), such as GPT-3 [14], Gopher [130], LaMDA [158], InstructGPT [117], PaLM [22], BLOOM [176], LLaMA [160], PaLM 2 [6], LLaMA 2 [161], GLM-130B [189], have achieved impressive multilingual performance.

Multilingual capabilities learned through implicit alignment in MLMs can be evidenced by some works that use extracted multilingual embeddings to complete downstream multilingual tasks. For instance, M3-Embedding [19], based on XLM-RoBERTa [26] and trained on a large-scale, diverse multilingual dataset, employs self-distillation to address multilingual hybrid retrieval tasks, including dense retrieval, lexical (sparse) retrieval, and multi-vector retrieval. Additionally, to overcome the limitations of existing embedding models in Malay retrieval tasks, the Multi-Lingual Malaysian Embedding [205] pre-trained and fine-tuned Malaysian Llama2 on collected synthetic Malaysian RAG and hard mining datasets. Features from different hidden layers of the Malaysian Llama2 were extracted to cater to the needs of various application scenarios.

[144] introduced jina-embeddings-v3, a multilingual embedding model employs task-specific Low-Rank Adaptation (LoRA) adapters to enhance performance across various downstream tasks. The model supports long-context retrieval with sequences up to 8192 tokens by integrating RoPE [148]. Additionally, jina-embeddings-v3 incorporates Matryoshka Representation Learning [66] to allow flexible reduction of embedding dimensions without compromising performance, enabling embeddings to be truncated to as low as 32 dimensions. By addressing common retrieval failures using synthetic data and leveraging instruction tuning [145, 172], the model achieves state-of-the-art results on the MTEB benchmark [109], outperforming proprietary embeddings from OpenAI and Cohere on English tasks and surpassing multilingual-e5-large-instruct [167] across all multilingual tasks.

*6.1.2   Explicit Alignment.* Unfriendly to low-resource languages, implicit alignment often requires training a multilingual model from scratch, which typically demands extensive computational costs and multilingual data. Specifically, both the large corpus used for pretraining and the fine-grained data for supervised fine-tuning (SFT) or reinforcement Learning with Human Feedback (RLHF) are hard to obtain. MLLMs tend to perform poorly on low-resource languages due to the imbalance in the corpus. Therefore, recent work has focused on bootstrapping multilingual embedding learning with enrich multilingual representations from pre-trained models, such as monolingual expert language models or machine translation models.

MT-LLM [137] aligns the multilingual representations, which are extracted by machine translation encoder, into the semantic embedding space of LLMs via self-distillation. This integration enables the LLMs to perform zero-shot ability to any language supported by the machine translation encoder. [163] employs a methodology that exclusively fine-tunes the query encoder while keeping the text encoder frozen on an english-only dataset, finding that this approach not only preserves but significantly enhances the multilingual embedding capabilities of the model. Based on the assumption that monolingual embeddings are well-structured, [118] proposes to distill alignment information from the monolingual similarity matrix into cross-modal embeddings to guide the cross-lingual alignment process. This method resolves issues in prior contrastive learning approaches that treated non-exact translation pairs as negative samples, which disrupted the monolingual embedding space. In addition, to further assess the capabilities of multilingual language models, [186]

propose MINERS, a multilingual, multi-task, tuning-free benchmark specifically designed to evaluate semantic retrieval tasks including bitext mining and classification via retrieval-augmented contexts.

Table 3. Recent Benchmarks for Evaluating Multilingual and Language-Specific.

| Benchmark | Language Coverage | Year |
|---|---|---|
| MMTEB [34] | 250+ languages | 2025 |
| mFollowIR [175] | Persian, Chinese, Russian | 2025 |
| FaMTEB [204] | Persian | 2025 |
| C-MTEB [180] | Chinese | 2024 |
| PL-MTEB [126] | Polish | 2024 |
| Scandinavian Embedding Benchmark (SEB) [35] | Nordic languages | 2024 |
| ArabicMTEB [12] | Arabic dialects, cross-lingual | 2024 |
| RusBEIR [64] | Russian | 2024 |
| JMTEB [84] | Japanese | 2024 |
| Mteb-french [25] | French | 2024 |
| Amharic Passage Retrieval Benchmark [5] | Amharic | 2025 |

*6.1.3 Multilingual Text Embedding Benchmark.* Recent efforts have also focused on benchmarking multilingual embeddings to evaluate and advance their effectiveness. Among these, the Massive Multilingual Text Embedding Benchmark (MMTEB) [34], an extension of the well-known MTEB [109], provides a comprehensive evaluation across over 250 languages and 500 tasks spanning 10 categories. In addition to MMTEB, specialized multilingual benchmarks have been proposed for specific languages or language groups, such as mFollowIR [175] for instruction-following retrieval, FaMTEB [204] for Persian, C-MTEB [180] for Chinese, PL-MTEB [126] for Polish, the Scandinavian Embedding Benchmarks [35] covering Nordic languages, Swan and ArabicMTEB [12] focusing on Arabic dialects and cross-cultural embeddings, as well as JMTEB [84] for Japanese and various benchmarks targeting German [171], targeting French [25], Russian [64], and Amharic [5]. These benchmarks collectively address the linguistic diversity, dialectal variations, and cultural contexts inherent in multilingual NLP, providing crucial standardized frameworks for evaluating embedding quality and facilitating progress in universal and language-specific embedding models.

## 6.2 Code Embedding

LLM-based code embeddings transform programming code into continuous vector representations that capture semantic, syntactic, functional properties and logic rules. These embeddings enable retrieval systems to encode code similarly to natural language, powering various types of code retrieval tasks. Traditional text embedding approaches struggle with code retrieval tasks which is largely due to a mismatch between the semantic structures of code and natural language. Therefore, it's important to align features within and between code and natural languages. Early code embedding works like code2vec [4] and Codebert [38], adapted NLP techniques to code by treating code as text. However, these methods struggled with structural nuances (e.g., control flow, data dependencies) and multi-language support. GraphCodeBERT [47] improved this by incorporating abstract syntax trees (ASTs). Further, UniXcoder [46] unifies code representations (tokens, ASTs, comments) into a single multimodal transformer, enabling diverse code understanding and generation tasks like code search and summarization.

With the development of LLMs, some LLM-based code embedding models were proposed to alleviate the limitation. [82] released StarCoderBase and StarCoder models, trained on large amounts

of code data across more than 80 programming languages which were sourced from GitHub issues, Git commits, and Jupyter notebooks. They also used multi-query attention for efficient long-context embedding (8K tokens). [93] proposed a generalizable training framework which converted diverse code-related tasks into retrieval tasks, and constructed training samples from code-to-text, text-to-code, code-to-code, text-to-text and hybrid text and code tasks across multiple programming languages. [73] introduced CodeR trained on a large-scale synthetic training dataset, which was created via a novel pipeline using LLMs for task design and sample generation. Meanwhile they employed a three-stage curriculum learning strategy to optimize training process.

## 6.3    Long Context Embedding

The ability to effectively capture and utilize long context information has become increasingly vital. Traditional embedding techniques, which typically focus on limited length of text, often struggle to maintain coherence and capture the full semantic meaning when applied to texts with extended contexts. Long context embeddings are designed to address these challenges, enabling models to handle extended sequences of text while preserving contextual integrity and meaning. A popular method for extending the length of context is training the embedding models with the long-context backbone, which can be obtained either by using the existing model or by pre-training with long inputs from scratch [19, 45, 114, 168]. Typically, [193] introduced a text encoder enhanced with RoPE [148] and unpadding [125], which is pre-trained by masked language model based on much longer (e.g. 8K tokens) and multilingual context with a two-stage curriculum. Based on the text encoder, they construct the embedding model through contrastive pre-training and finetuning utilizing InfoNCE as the loss function.

However, fine-tuning LLMs may result in high computational costs, where there are more and more efforts are proposed to extend the contexts of LLM embedding models through the plug-and-play methods [58, 132, 169]. [95] proposes a chunking-free architecture to process long context, which adds a special token <LMK> at the end of each sentence in the long context to capture the coherent semantic. During training, this work introduces the positional weight and utilizes multi-stage training to achieve a superior performance. [138] proposes a plug-in module extensible embedder to process the long context which is partitioned into multiple chunks. Each chunk is embedded and down-scaled by extensible embedder as the compact representation. This work train extensible embedder through two-stream auto-regression with fixed downstream LLM's parameters, which does not affect the LLM's original capabilities. Furthermore, [202] explores the extensive plug-and-play strategies to extend existing embedding models to long context, which includes parallel context windows, grouped positions, recurrent positions, linear position interpolation and so on. This work also proposes a newly constructed LONGEMBED benchmark for long context retrieval evaluation, which includes two synthetic tasks with flexible document length and four real tasks with featuring dispersed target information.

## 6.4    Cross-modal Embedding

Cross-modal embeddings aim to create a unified representation for data from different modalities, such as text, images, and audio. These embeddings enable the integration of multimodal information, which is crucial for applications like image captioning, visual QA, and multimodal search.

The potential of Transformers for efficiently learning cross-modal embeddings has been demonstrated by works like VATT [3], ViT [30]. Following the introduction of CLIP [128], which uses language as supervision for visual models to bridge the gap between vision and language, a surge of vision-language cross-modal models such as ALBEF [79], VLMO [9], CoCa [187], BLIP [78], and BEiT [8] emerged. However, these models have significant limitations in representing purely textual or purely visual data. The VISTA [198] proposes a new embedding approach for universal

multimodal retrieval, allowing the pretrained LMs to recognize image tokens by utilizing ViT as an image tokenizer. It undergoes two-stage training on datasets with weak labels and high-quality synthetically composed image-text datasets, achieving superior performance across a variety of multimodal retrieval tasks.

With the development of LLMs, many efforts have extended LLMs to various modalities, such as image (BLIP-2 [77], MiniGPT4 [201], Llava [89]), audio (AudioPaLM [134], Qwen-Audio [23]), and video (Video-Llama [191], Video-ChatGPT [99], Video-Llava [88]). However, these efforts primarily focus on generating text related to these modalities. PaLM 2 DE [44] incorporates LLMs into dual-encoder architecture, utilizing PaLM to initialize and augment multilingual text comprehension for cross-modal retrieval tasks. This integration results in enhanced performance across an impressive range of 102 languages, even though it was trained on a limited dataset of 21 languages. Furthermore, by leveraging machine translation data, PaLM 2 DE significantly boosts its cross-lingual capabilities.

Table 4. Summary of recent reasoning-based retrieval benchmarks.

| Dataset | Year | Task Type | Reasoning Required | Domain |
|---|---|---|---|---|
| **RAR-b** [177] | 2024 | Reasoning-Intensive Retrieval | temporal, numerical, spatial, symbolic | General |
| **ImpliRet** [152] | 2025 | Reasoning-Intensive Retrieval | arithmetic, temporal, commonsense | General |
| **BRIGHT** [147] | 2025 | Reasoning-Intensive Retrieval | Symbolic, math, code | General |
| **InstructIR** [115] | 2024 | Instructional IR | Instruction following | General |
| **InfoSearch** [199] | 2024 | Instructional IR | Customized instruction following | General |
| **FollowIR** [174] | 2024 | Instructional IR | Long-Form instruction comprehension | General |
| **MAIR** [151] | 2024 | Instructional IR | Diverse instruction following | General |
| **BIRCO** [170] | 2024 | Instructional IR | Compositional, Complex Objectives | General |
| **mFollowIR** [175] | 2025 | Instructional IR | Multi-lingual instruction following | General |
| **IFIR** [142] | 2025 | Instructional IR | Domain-specific instruction following | Biomedical, Legal, Technical |
| **CLERC** [49] | 2024 | Case/Legal Retrieval | Precedent Reasoning, Legal Logic | Legal |
| **Bar Exam/ Housing Statute** [196] | 2025 | Case/Legal Retrieval | Symbolic, Legal Reasoning | Legal |

## 6.5 Reasoning-aware Embedding

The landscape of IR is rapidly evolving with the emergence of language models capable of complex reasoning and instruction following. Traditional embedding methods, which primarily encode surface-level semantic similarity, often fall short in scenarios requiring logical inference, multi-hop reasoning, or constraint satisfaction. This has motivated the development of reasoning-aware embedding techniques—embedding models explicitly designed to encode latent reasoning paths, logical structures, and task-specific constraints into vector representations.

Two complementary research directions highlight the demand for such embeddings: Reasoning-Intensive Retrieval and Instruction-Following IR. Both emphasize retrieval tasks where shallow matching is insufficient, and embeddings must internalize complex relationships.

**Reasoning-Intensive Retrieval** involves tasks requiring deep inference over queries and documents. Benchmarks such as RAR-b [177], BRIGHT [147], and ImpliRet [152] evaluate systems

on arithmetic reasoning, implicit fact chaining, and temporal logic. In high-stakes domains like law, datasets like CLERC [49] and the Legal Reasoning Benchmark [196] demand embeddings that capture legal concepts, precedents, and inferable relationships beyond what surface-level semantics can provide. These challenges necessitate reasoning-aware embedding spaces that can support multi-step deductive or abductive retrieval.

**Instruction-Following Retrieval**, on the other hand, emphasizes alignment with user intent expressed through natural language instructions. Systems like InstructIR [115], MAIR [151], and InfoSearch [199] show that instructions often introduce constraints or goals requiring interpretation and reasoning. Benchmarks such as FollowIR, IFIR [142], and mFollowIR [175] evaluate how well models can incorporate instruction semantics into retrieval behavior. Embedding-based approaches here must go beyond vanilla query representation—they must fuse instruction semantics and reasoning signals into a unified vector form, again underscoring the need for reasoning-aware embeddings.

These developments collectively reveal a growing shift from simple semantic encoding toward embedding architectures that encode inferential and procedural knowledge, thereby enhancing retrieval in complex, instruction-rich, or high-reasoning domains.

## 6.6 Other Domain-Specific Embeddings

Beyond general-purpose language and multi-modal embeddings, there has been growing interest in developing domain-specific embedding models tailored for specialized fields such as finance and chemistry [2]. These domains pose unique challenges due to their specialized terminology, complex semantics, and domain-specific knowledge requirements. The Finance Massive Text Embedding Benchmark (FinMTEB) [153] targets the financial domain, covering 64 datasets across seven tasks involving diverse financial text types such as news articles, annual reports, ESG disclosures, and earnings calls in both Chinese and English. FinMTEB evaluations reveal that general-purpose embeddings perform poorly on financial tasks, while domain-adapted models like Fin-E5, trained via persona-based synthetic data, achieve significantly better results. Interestingly, simple Bag-of-Words methods sometimes outperform complex dense embeddings on financial semantic similarity tasks, indicating limitations in current embedding techniques for finance. Similarly, the Chemical Text Embedding Benchmark (ChemTEB) [60] addresses the chemical sciences, a domain with specialized linguistic and semantic challenges. ChemTEB evaluates 34 models on chemical literature and data, providing insights into their ability to handle chemical terminology, molecular descriptions, and scientific language. This benchmark highlights the need for domain-aware embedding strategies in chemistry, as generic models often fall short in capturing domain-specific knowledge. Both benchmarks contribute standardized evaluation frameworks that drive the advancement of more accurate and efficient embeddings tailored for expert domains.

## 7 Discussions

In this section, we will discuss the performance of representative embedding models on MTEB benchmarks, compare dense and sparse embeddings, and different pooling strategies, and review the implications of scaling laws for larger LLMs.

**Performance and Efficiency Comparison.** Table 5 summarizes the performance of both traditional pretrained models and recent LLM-based models on the MTEB benchmarks. It is evident that supervised methods generally outperform traditional self-supervised methods. Moreover, more

---

[2]Code embeddings focus on programming languages and code semantics across tasks, emphasizing data modality rather than a specific application domain, unlike domain-specific embeddings tailored to fields like finance or chemistry.

Table 5. Results on MTEB leaderboard as of July 2025. The models are categorized into two groups based on their embedding dimension (Dim.) and the number of parameters (# Params.).

| | Dim. | # Params. | Class. | Cluster. | Pair. | Rerank. | Retrieval | STS | Summary | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| Self-Supervised methods | | | | | | | | | | |
| Glove | 200 | - | 57.29 | 27.73 | 70.92 | 43.29 | 21.62 | 61.85 | 28.87 | 41.97 |
| BERT | 768 | 110M | 61.66 | 30.12 | 56.33 | 43.44 | 10.59 | 54.36 | 29.82 | 38.33 |
| SimCSE-BERT-unsup | 768 | 110M | 62.50 | 29.04 | 70.33 | 46.47 | 20.29 | 74.33 | 31.15 | 45.45 |
| Supervised methods | | | | | | | | | | |
| SimCSE-BERT-sup | 768 | 110M | 67.32 | 33.43 | 73.68 | 47.54 | 21.82 | 79.12 | 23.31 | 48.72 |
| Contriever | 768 | 110M | 66.68 | 41.10 | 82.53 | 53.14 | 41.88 | 76.51 | 30.36 | 56.00 |
| SGPT-1.3B | | 1.3B | 66.52 | 39.92 | 79.58 | 54.00 | 44.49 | 75.74 | 25.44 | 56.11 |
| GTR-T5-XXL | 768 | 5B | 67.41 | 42.42 | 86.12 | 56.66 | 48.48 | 78.38 | 30.64 | 58.97 |
| GTR-T5-XL | 768 | 1.2B | 67.11 | 41.51 | 86.13 | 55.97 | 47.96 | 77.80 | 30.21 | 58.42 |
| Instructor-XL [146] | 768 | 1.5B | 73.12 | 44.74 | 86.62 | 57.29 | 49.26 | 83.06 | 32.32 | 61.79 |
| Text-embedding-3-large | 3,072 | n/a | 75.45 | 49.01 | 85.72 | 59.16 | 55.44 | 81.73 | 29.92 | 64.59 |
| E5-mistral-7b-instruct [168] | 4,096 | 7B | 78.47 | 50.26 | 88.34 | 60.21 | 56.89 | 84.63 | 31.40 | 66.63 |
| GritLM-7B [108] | 4,096 | 7B | 79.46 | 50.61 | 87.16 | 60.49 | 57.41 | 83.35 | 30.37 | 66.76 |
| Echo-mistral-7b-instruct [143] | 4,096 | 7B | 77.43 | 46.32 | 87.34 | 58.14 | 55.52 | 82.56 | 30.73 | 64.69 |
| Gecko [69] | 256 | 1.2B | 78.99 | 45.07 | 87.25 | 57.78 | 52.44 | 84.93 | 32.36 | 64.37 |
| LLM2Vec-Llama-3 [10] | 4,096 | 8B | 75.92 | 46.45 | 87.79 | 59.69 | 56.63 | 83.58 | 30.94 | 65.01 |
| GTE-Qwen1.5-7B-instruct [86] | 4,096 | 7B | 79.60 | 55.83 | 87.38 | 60.13 | 56.24 | 82.42 | 31.46 | 67.34 |
| SFR-Embedding | 4,096 | 7B | 78.33 | 51.67 | 88.54 | 60.64 | 59.00 | 85.05 | 31.16 | 67.56 |
| Voyage-lite-02-instruct [165] | 1,024 | - | 79.25 | 52.42 | 86.87 | 58.24 | 56.6 | 85.79 | 31.01 | 67.13 |
| Voyage-large-2-instruct [165] | 1,024 | - | 81.49 | 53.35 | 89.24 | 60.09 | 58.28 | 84.58 | 30.84 | 68.28 |
| Linq-Embed-Mistral [62] | 4096 | 7B | 80.2 | 51.4 | 88.4 | 60.3 | 60.2 | 85.0 | 31.0 | 68.2 |
| NV-Embed [67] | 4,096 | 7B | 87.35 | 52.80 | 86.91 | 60.59 | 59.36 | 82.84 | 31.20 | 69.32 |
| BGE-en-icl [76] | 4,096 | 7B | 88.95 | 57.89 | 88.14 | 59.86 | 62.16 | 84.24 | 30.77 | 71.67 |
| NV-Retriever [105] | 4,096 | 7B | 90.37 | 58.46 | 88.67 | 60.65 | 62.65 | 84.31 | 30.70 | 72.31 |
| Gemini-embedding [68] | 3072 | - | 90.05 | 59.39 | 87.7 | 48.59 | 64.35 | 85.29 | 38.28 | 73.3 |
| Qwen-embedding-8B [194] | 1024 | 8B | 90.43 | 58.57 | 87.52 | 51.56 | 69.44 | 88.58 | 34.83 | 75.22 |
| Qwen-embedding-4B [194] | 2560 | 4B | 89.84 | 57.51 | 87.01 | 50.76 | 68.46 | 88.72 | 34.39 | 74.60 |
| Qwen-embedding-0.6B [194] | 4096 | 0.6B | 85.76 | 54.05 | 84.37 | 48.18 | 61.83 | 86.57 | 33.43 | 70.70 |

recent LLM-based embedding models, such as those built on LLama and Mistral, demonstrate significantly better performance compared to earlier smaller-scale pretrained models like BERT, SGPT, and GTR. This performance boost can be attributed to the increased model size and higher embedding dimensionality of the newer models, which enable them to capture more nuanced semantic relationships and represent richer contextual information within the embeddings. Additionally, through distillation from the output signals of a cross-encoder, a relatively small-scale model like Gecko (1.5B) can achieve performance comparable to that of larger models (7B). For specific task types, recent LLM-based embedding models demonstrate close performance in pair, reranking, STS, and summary tasks. However, these models (e.g., BGE-en-icl, NV-Embed, NV-Retriever) show significant improvements in clustering, classification, and retrieval tasks, likely due to training datasets that align closely with the MTEB data sources (as seen in Table 2).

**Dense vs. Sparse Embedding.** While constructing dense embeddings based on LLMs has gained significant attention, some researchers are also investigating the acquisition of sparse embeddings by leveraging lexicon-importance distributions derived from LLMs [19, 29, 113, 203]. [203] evaluate various embedding types by directly prompting or fine-tuning LLMs. Their findings

reveal that in zero-shot retrieval settings, dense and sparse embeddings perform differently across different benchmarks: for instance, dense embeddings outperform sparse ones on BEIR and TREC-DL 2019, while sparse embeddings show better results on TREC-DL 2020 and MSMARCO. However, in supervised fine-tuning settings, dense retrieval significantly outperforms sparse retrieval on datasets like MSMARCO and TREC-DL 2019/2020 when using models such as Llama3-8B and Llama3-70B. Notably, combining both representation types further enhances performance. [19] fine-tune multilingual, multi-granularity text embeddings, integrating both dense and sparse retrieval models based on XLM-RoBERTa. Their results show that dense retrieval notably outperforms sparse retrieval on cross-lingual tasks such as MKQA and MIRACL, and that combining both methods also yields superior results, consistent with the findings reported in [203].

**Last token vs. Mean pooling.** Two common methods for obtaining embeddings from a sequence of tokens are: i) mean pooling and ii) the embedding of the last [EOS] token. The former calculates the average of token embeddings, which can potentially obscure key information from phrases. In contrast, [EOS] token embedding is susceptible to recency bias, as it heavily relies on the embedding of the final token in the sequence. [67] conducted experiments comparing several methods for obtaining embeddings from the last layer of LLMs, including the [EOS] token embedding, mean pooling, and their proposed latent attention method. The results demonstrate that mean pooling can always better than [EOS] token embedding based on average MTEB scores, regardless of whether causal attention or bidirectional attention settings are used in the LLMs. Furthermore, their latent attention method enhances embedding capabilities, yielding better results than both the [EOS] token embedding and mean pooling. This research highlights the potential for exploring more advanced designs to improve last-layer embeddings for enhanced performance.

**In-context Learning (ICL) for Embedding Model.** ICL has emerged as a powerful paradigm for enhancing the performance of LLMs in various generative tasks. It enables LLMs to adapt dynamically to the input context by leveraging real-time examples and instructions during inference. This capability to learn from the surrounding context potentially allows LLMs to produce more accurate and context-sensitive embeddings. Recently, researchers have also explored the effectiveness of ICL for LLM-based representation models. For instance, [56] first investigated incorporating ICL examples into prompt-based methods and demonstrated their effectiveness for STS tasks. Their approach involved constructing ICL examples by generating word-sentence pairs using ChatGPT to align words with sentence semantics from the STS-B training set, as well as utilizing word-definition pairs from the Oxford dictionary, inspired by DefSent [162]. More recently, [76] incorporated few-shot examples on the query side to enhance query embeddings during supervised contrastive training. Despite these few-shot examples being randomly sampled, they still improved the model's performance. Besides, [143] proposed repeating the input twice within the context to obtain sequence embeddings, a technique that can be viewed as a specialized form of ICL, and validated its effectiveness in both training-free and supervised contrastive training settings. These studies collectively indicate that ICL is a promising strategy for enhancing representation learning with LLMs. Nevertheless, exploring more diverse and high-quality contextual information holds potential for further improving the effectiveness of this approach.

**Scaling Law of Embedding Models.** The question of whether scaling laws hold for embedding models based on LLMS is both significant and intriguing. The effectiveness of embedding models is highly sensitive to several training factors, including model size, data size, and data quality (e.g., diversity and quality of negatives), making it difficult to isolate the effect of each factor independently. [36] evaluate the quality of dense retrieval models using a contrastive entropy metric and observe a power law relationship between model performance, model size, and data size across various annotation methods and datasets. However, their experiments were conducted
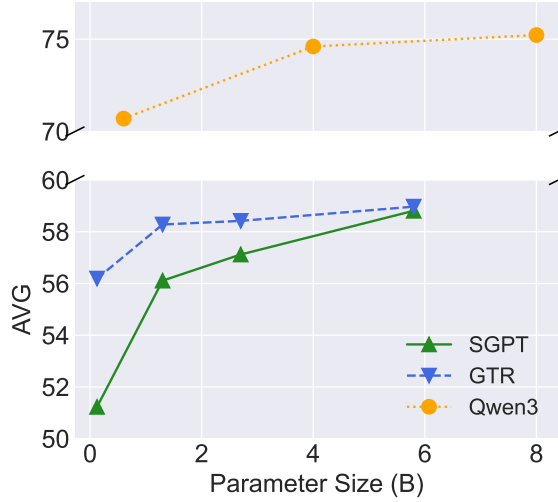
Fig. 7. The performance trend of models with varying sizes on MTEB benchmark.

on relatively small-scale models, such as BERT models ranging from 0.5M to 110M parameters. Additionally, [109] report the performance of models with varying sizes, such as SGPT (125M, 1.3B, 2.7B, 5.8B) and GTR (Base, Large, XL, XXL), on the MTEB benchmark shown in Figure 7, showing consistent improvements as model size increases, even up to 5.8B parameters. [203] further test performance by directly prompting LLMs for embeddings, finding that LLaMA3-70B-Instruct outperforms LLaMA3-8B-Instruct across different embedding types. Similarly, the Qwen3-Embedding [194] models demonstrate a positive correlation between scale and performance, improving steadily from 0.6B to 4B and 8B; however, the marginal gain from 4B to 8B is notably smaller, suggesting diminishing returns at higher scales. Taken together, this evidence suggests that the scaling law may also apply to embedding models built on large LMs, although more systematic experiments are needed to rigorously validate this trend across broader configurations.

## 8 Open Problems

Embedding techniques based on LLMs have emerged as a crucial area of research. Despite their impressive performance across various tasks, the process of generating more effective embeddings still presents a number of challenges.

**Embedding Quality Across Different Tasks.** While LLMs have shown impressive performance in generating embeddings for specific tasks, such as semantic similarity or text classification, their effectiveness across a broader range of tasks remains uncertain. Tasks like clustering, reranking, and summary require embeddings to capture different kinds of relationships between data points, and it is not always clear if embeddings generated by LLMs can consistently perform well in these varied contexts. Moreover, different tasks may prioritize different aspects of the embedding space, such as local versus global structure or task-specific nuances [16], which might not be fully captured by general-purpose LLM embeddings. This raises the open question of how LLM embeddings can be optimized or adapted to maintain high-quality representations across diverse tasks [81], and whether task-specific tuning or hybrid approaches might be necessary to improve overall effectiveness in more complex and specialized applications.

**Efficiency vs. Accuracy Trade-offs.** While LLM-generated embeddings can achieve better accuracy in capturing semantic relationships, they often come at a significant computational cost due to the large model sizes and high-dimensional representations (e.g., 4096 for most recent LLMs). This creates a critical challenge in balancing efficiency and accuracy. For tasks that require real-time processing or deployment in resource-constrained environments, the time and memory demands of using LLMs can be prohibitive. Reducing the computational burden without sacrificing the quality of embeddings is an open problem. Techniques such as model distillation [69], pruning, or dimensionality reduction have been proposed, but each comes with trade-offs in terms of performance loss or reduced representational power. Moreover, it remains unclear to what extent the high-dimensional nature of LLM embeddings contributes to their effectiveness, and whether lower-dimensional embeddings could perform comparably in specific tasks. Thus, research is needed to explore strategies that make LLM embeddings more efficient while preserving their robustness, especially in applications where both accuracy and real-time performance are critical. Additionally, exploring the embedding performance of relatively smaller models, such as Phi3 [1] and MiniCPM [50], could offer valuable insights into achieving a balance between computational efficiency and embedding quality.

**Long Context Embedding.** LLMs are known for their ability to handle long-context dependencies in text, but how well these long contexts are represented in the embeddings they generate is still an open question. Tasks such as document retrieval and multi-hop question answering require embeddings that capture the relationships and themes spanning lengthy passages, yet traditional embedding approaches may struggle to encode this information effectively. One of the main challenges is that long-context embeddings need to balance local context (word-level relationships) with global coherence (document-wide meaning), which is difficult to achieve in a single embedding space. Additionally, handling long contexts often increases the computational complexity and dimensionality of the embeddings, making them harder to use in resource-constrained applications. While techniques such as efficient attention mechanisms [21] and chunk-based processing [138] have been proposed to address these issues, they often introduce trade-offs by either sacrificing fine-grained contextual information or losing coherence across chunks, limiting their ability to fully capture the intricacies of long-context dependencies This raises the question of how to design more efficient embedding methods that can faithfully represent long-context information while remaining computationally feasible.

**Reasoning Gaps in Embedding-Based Retrieval.** Despite progress in retrieval modeling, current systems still struggle to perform robustly on reasoning-intensive tasks. A fundamental issue is the lack of reasoning capability in the first-stage retrievers, which are responsible for high-recall candidate generation. These retrievers, typically based on dense or hybrid embeddings, often fail to capture implicit connections, logical dependencies, or multi-step inference chains required for identifying relevant documents. As a result, important evidence is missed early on, limiting the performance ceiling of downstream rerankers or generators. This challenge is especially prominent in benchmarks like BRIGHT and ImpliRet, where retrieval requires understanding implicit facts or combining dispersed clues. The absence of fine-grained supervision for reasoning types during training further compounds the issue, as retrievers are usually optimized with weak relevance signals that do not reflect the nature of the reasoning required. Moreover, standard IR metrics such as recall or nDCG are not designed to evaluate reasoning quality, making progress difficult to measure. Addressing these challenges will require retrieval models that are not only semantically aware but also reasoning-aware—capable of aligning with the latent inference patterns embedded in user queries.

**Impact of Training Data Bias on Embeddings.** LLMs are trained on large text corpora, which often contain subtle biases related to factors like gender, race, and culture [31, 32, 135]. These biases can inadvertently be reflected in the embeddings generated by LLMs, potentially leading to unfair or unbalanced results in downstream tasks such as sentiment analysis, search, recommendation systems, or hiring platforms. For example, if the training data is skewed towards certain demographics, the embeddings may fail to represent minority groups accurately, resulting in biased outputs. Addressing this issue is challenging, as efforts to reduce bias, such as using de-biasing techniques or filtering the training data, can sometimes come at the cost of model performance or generalization. Furthermore, defining what constitutes "fair" or "unbiased" embeddings is complex and may vary across different tasks and contexts. As a result, balancing fairness with performance remains an open problem. Research into better methods for detecting and mitigating bias in LLM-generated embeddings, while preserving their effectiveness in real-world applications, is essential for ensuring that LLMs are more equitable and reliable in diverse use cases.

**Adapting LLM Embeddings for Low-Resource Domains.** One of the major challenges in using LLM-based embeddings is their adaptability to low-resource domains, where training data is sparse, or under-representative of the target domain. While LLMs have been pre-trained on large, diverse datasets, these models often struggle when applied to specialized domains, such as medical, legal, or technical fields where pair data is limited [80, 141]. The embeddings generated in these settings may fail to capture the necessary domain-specific nuances, leading to limited performance on downstream tasks like question answering or retrieval. Fine-tuning LLMs for such low-resource domains requires domain-specific data that may not always be available. Additionally, overfitting to small datasets is a concern, which raises the question of how to adapt LLM embeddings to low-resource settings without sacrificing generalization.

**Robustness to Adversarial Attacks.** LLM-generated embeddings, while powerful, are often vulnerable to adversarial attacks, where small, carefully crafted perturbations in the input can lead to significant changes in the embeddings and downstream task performance [7, 24]. This raises concerns about the robustness of LLM embeddings in real-world applications, where inputs may be noisy or manipulated. These vulnerabilities pose significant risks, especially in critical applications like healthcare or security information retrieval or retrieval-augmented generation. While there have been some efforts to improve the robustness of embeddings through techniques such as adversarial training, regularization, or perturbation-based testing, these methods are still far from fully addressing the problem. Moreover, ensuring that embeddings remain robust without sacrificing their generalization ability or making them too rigid to handle legitimate variations in input is a complex challenge. The development of more resilient embedding methods that can withstand adversarial manipulations while maintaining high performance across a variety of tasks is an open problem and a vital area for future research.

**How to effective training for embedding models?** LLMs are typically pre-trained on large text corpora for general language modeling tasks, but adapting them specifically for generating high-quality embeddings often requires resource-intensive contrastive training. Recent models increasingly rely on large-scale paired data for contrastive fine-tuning; for instance, Qwen3-Emb is trained on 150 million paired examples [194]. This poses a significant challenge, particularly in scenarios where only limited data is available. An intriguing open problem is how to effectively tune LLMs for embedding tasks using small-scale datasets. For example, LIMA demonstrates that only about 1,000 high-quality, human-curated samples are sufficient for aligning LLMs in general supervised fine-tuning [197]. This suggests that efficient fine-tuning with a smaller amount of data might be feasible for embedding tasks as well. Additionally, exploring whether there are more

effective training methods or strategies tailored specifically for tuning LLMs as embedding models could provide valuable insights.

**Is Continued Training Beneficial for Embedding Models?** Current LLMs are pretrained primarily on next-token prediction (NTP) tasks, which focus on sequence completion rather than embedding generation. This misalignment raises the question of whether continued training with alternative tasks, such as next sentence prediction or sentence-level reconstruction, could enhance their performance as embedding models. Future research could explore how and what domain-specific or task-specific continued training objectives can be designed to preserve generalization while improving domain relevance, balancing computational efficiency, embedding quality, and downstream task adaptability.

## 9    Conclusion

In conclusion, this survey addresses the paradigm shift toward leveraging LLMs as embedding models, emphasizing their impact on representation learning across diverse tasks in natural language understanding, information retrieval and recommendation. The paper provides an in-depth analysis of both tuning-free and tuning-based approaches for generating effective embeddings, highlighting methods that maximize LLM capabilities without extensive fine-tuning, as well as techniques that further finetune these models for task-specific contexts. By comparing various LLM-based embedding strategies, including dense and sparse embeddings, pooling techniques, and scalability considerations, we provide practical insights for researchers and practitioners. Moreover, this survey identifies key open challenges. As LLMs continue to evolve, this work aims to serve as a valuable guide for understanding and advancing embedding methodologies in the LLM era, ultimately supporting improved performance across a wide range of NLP and IR applications.

## Acknowledgments

## References

[1] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Hassan Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Singh Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio C'esar Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allison Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Young Jin Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Xianmin Song, Olatunji Ruwase, Praneetha Vaddamanu, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Cheng-Yuan Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. ArXiv abs/2404.14219 (2024). https://api.semanticscholar.org/CorpusID:269293048

[2] OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim ing Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Made laine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Benjamin Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Sim'on Posada Fishman, Juston Forte, Is abella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh,

Raphael Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Jo hannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Lukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Ryan Kiros, Matthew Knight, Daniel Kokotajlo, Lukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Ma teusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel P. Mossing, Tong Mu, Mira Murati, Oleg Murk, David M'ely, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Ouyang Long, Cullen O'Keefe, Jakub W. Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alexandre Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Pondé de Oliveira Pinto, Michael Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack W. Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario D. Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin D. Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas A. Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cer'on Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll L. Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. GPT-4 Technical Report. https://api.semanticscholar.org/CorpusID:257532815

[3] Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. 2021. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text. Advances in Neural Information Processing Systems 34 (2021), 24206–24221.

[4] Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. 2019. code2vec: Learning distributed representations of code. Proceedings of the ACM on Programming Languages 3, POPL (2019), 1–29.

[5] Kidist Amde Mekonnen, Yosef Worku Alemneh, and Maarten de Rijke. 2025. Optimized Text Embedding Models and Benchmarks for Amharic Passage Retrieval. arXiv e-prints (2025), arXiv–2505.

[6] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. arXiv preprint arXiv:2305.10403 (2023).

[7] Suriya Ganesh Ayyamperumal and Limin Ge. 2024. Current state of LLM Risks and AI Guardrails. arXiv preprint arXiv:2406.12934 (2024).

[8] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. 2021. Beit: Bert pre-training of image transformers. arXiv preprint arXiv:2106.08254 (2021).

[9] Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, Subhojit Som, and Furu Wei. 2021. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. arXiv preprint arXiv:2111.02358 (2021).

[10] 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 (2024).

[11] Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2013. Representation learning: A review and new perspectives. IEEE transactions on pattern analysis and machine intelligence 35, 8 (2013), 1798–1828.

[12] Gagan Bhatia, El Moatez Billah Nagoudi, Abdellah El Mekki, Fakhraddin Alwajih, and Muhammad Abdul-Mageed. 2024. Swan and ArabicMTEB: Dialect-aware, Arabic-centric, cross-lingual, and cross-cultural embedding models and benchmarks. arXiv preprint arXiv:2411.01192 (2024).

[13] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. Transactions of the association for computational linguistics 5 (2017), 135–146.

[14] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural

information processing systems 33 (2020), 1877–1901.

[15] Nithin Buduma, Nikhil Buduma, and Joe Papa. 2022. Fundamentals of deep learning. " O'Reilly Media, Inc.".

[16] ZeFeng Cai, Chongyang Tao, Tao Shen, Can Xu, Xiubo Geng, Xin Alex Lin, Liang He, and Daxin Jiang. 2022. HypeR: Multitask Hyper-Prompted Training Enables Large-Scale Retrieval Generalization. In The Eleventh International Conference on Learning Representations.

[17] Yang Cao, Sikun Yang, Chen Li, Haolong Xiang, Lianyong Qi, Bo Liu, Rongsheng Li, and Ming Liu. 2025. TAD-Bench: A Comprehensive Benchmark for Embedding-Based Text Anomaly Detection. arXiv preprint arXiv:2501.11960 (2025).

[18] Laura Caspari, Kanishka Ghosh Dastidar, Saber Zerhoudi, Jelena Mitrovic, and Michael Granitzer. 2024. Beyond Benchmarks: Evaluating Embedding Model Similarity for Retrieval Augmented Generation Systems. arXiv preprint arXiv:2407.08275 (2024).

[19] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. Bge m3-embedding: Multilingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. arXiv preprint arXiv:2402.03216 (2024).

[20] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In International conference on machine learning. PMLR, 1597–1607.

[21] Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2023. Longlora: Efficient fine-tuning of long-context large language models. arXiv preprint arXiv:2309.12307 (2023).

[22] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. Journal of Machine Learning Research 24, 240 (2023), 1–113.

[23] Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. 2023. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. arXiv preprint arXiv:2311.07919 (2023).

[24] Jaymari Chua, Yun Li, Shiyi Yang, Chen Wang, and Lina Yao. 2024. Ai safety in generative ai large language models: A survey. arXiv preprint arXiv:2407.18369 (2024).

[25] Mathieu Ciancone, Imene Kerboua, Marion Schaeffer, and Wissam Siblini. 2024. Mteb-french: Resources for french sentence embedding evaluation and analysis. arXiv preprint arXiv:2405.20468 (2024).

[26] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. arXiv preprint arXiv:1911.02116 (2019).

[27] Bijoyan Das and Sarit Chakraborty. 2018. An improved text sentiment classification model using TF-IDF and next word negation. arXiv preprint arXiv:1806.06407 (2018).

[28] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018).

[29] Meet Doshi, Vishwajeet Kumar, Rudra Murthy, Jaydeep Sen, et al. 2024. Mistral-SPLADE: LLMs for for better Learned Sparse Retrieval. arXiv preprint arXiv:2408.11119 (2024).

[30] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020).

[31] Yucong Duan. 2024. The Large Language Model (LLM) Bias Evaluation (Age Bias). DIKWP Research Group International Standard Evaluation. DOI 10 (2024).

[32] Yucong Duan, Fuliang Tang, Kunguang Wu, Zhendong Guo, Shuaishuai Huang, Yingtian Mei, Yuxing Wang, Zeyu Yang, and Shiming Gong. 2023. Ranking of large language model (llm) regional bias. (2023).

[33] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony S. Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Cantón Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab A. AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriele Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guanglong Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Laurens Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Ju-Qing Jia, Kalyan Vasuden Alwala, K. Upasani, Kate Plawiak, Keqian Li, Ken-591 neth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,

Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline C. Muzzi, Mahesh Babu Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melissa Hall Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri S. Chatterji, Olivier Duchenne, Onur cCelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasić, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Chandra Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yiqian Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zhengxu Yan, Zhengxing Chen, Zoe Papakipos, Aaditya K. Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adi Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Ben Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Shang-Wen Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzm'an, Frank J. Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory G. Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Han Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kaixing(Kai) Wu, U KamHou, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, A Lavender, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollár, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sung-Bae Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Andrei Poenaru, Vlad T. Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xia Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary

DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The Llama 3 Herd of Models. ArXiv abs/2407.21783 (2024). https://api.semanticscholar.org/CorpusID:271571434

[34] Kenneth Enevoldsen, Isaac Chung, Imene Kerboua, Márton Kardos, Ashwin Mathur, David Stap, Jay Gala, Wissam Siblini, Dominik Krzemiński, Genta Indra Winata, Saba Sturua, Saiteja Utpala, Mathieu Ciancone, Marion Schaeffer, Gabriel Sequeira, Diganta Misra, Shreeya Dhakal, Jonathan Rystrøm, Roman Solomatin, Ömer Çağatan, Akash Kundu, Martin Bernstorff, Shitao Xiao, Akshita Sukhlecha, Bhavish Pahwa, Rafał Poświata, Kranthi Kiran GV, Shawon Ashraf, Daniel Auras, Björn Plüster, Jan Philipp Harries, Loïc Magne, Isabelle Mohr, Mariya Hendriksen, Dawei Zhu, Hippolyte Gisserot-Boukhlef, Tom Aarsen, Jan Kostkan, Konrad Wojtasik, Taemin Lee, Marek Šuppa, Crystina Zhang, Roberta Rocca, Mohammed Hamdy, Andrianos Michail, John Yang, Manuel Faysse, Aleksei Vatolin, Nandan Thakur, Manan Dey, Dipam Vasani, Pranjal Chitale, Simone Tedeschi, Nguyen Tai, Artem Snegirev, Michael Günther, Mengzhou Xia, Weijia Shi, Xing Han Lù, Jordan Clive, Gayatri Krishnakumar, Anna Maksimova, Silvan Wehrli, Maria Tikhonova, Henil Panchal, Aleksandr Abramov, Malte Ostendorff, Zheng Liu, Simon Clematide, Lester James Miranda, Alena Fenogenova, Guangyu Song, Ruqiya Bin Safi, Wen-Ding Li, Alessia Borghini, Federico Cassano, Hongjin Su, Jimmy Lin, Howard Yen, Lasse Hansen, Sara Hooker, Chenghao Xiao, Vaibhav Adlakha, Orion Weller, Siva Reddy, and Niklas Muennighoff. [n. d.]. MMTEB: Massive Multilingual Text Embedding Benchmark. ([n. d.]). https://doi.org/10.48550/arXiv.2502.13595 arXiv:2502.13595 [cs]

[35] Kenneth Enevoldsen, Márton Kardos, Niklas Muennighoff, and Kristoffer L Nielbo. 2024. The scandinavian embedding benchmarks: Comprehensive assessment of multilingual and monolingual text embedding. Advances in Neural Information Processing Systems 37 (2024), 40336–40358.

[36] Yan Fang, Jingtao Zhan, Qingyao Ai, Jiaxin Mao, Weihang Su, Jia Chen, and Yiqun Liu. 2024. Scaling laws for dense retrieval. In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval. 1339–1349.

[37] Jiazhan Feng, Chongyang Tao, Zhen Li, Chang Liu, Tao Shen, and Dongyan Zhao. 2022. Reciprocal learning of knowledge retriever and response ranker for knowledge-grounded conversations. In Proceedings of the 29th International Conference on Computational Linguistics. 389–399.

[38] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. arXiv preprint arXiv:2002.08155 (2020).

[39] Jianling Gao, Chongyang Tao, Zhenchao Sun, Xiya Jiang, and Shuai Ma. 2025. Semi-Supervised Anomaly Detection through Denoising-Aware Contrastive Distance Learning. In Proceedings of the ACM on Web Conference 2025. 2111–2119.

[40] Luyu Gao and Jamie Callan. 2021. Condenser: a Pre-training Architecture for Dense Retrieval. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 981–993. https://doi.org/10.18653/v1/2021.emnlp-main.75

[41] Luyu Gao and Jamie Callan. 2021. Unsupervised corpus aware language model pre-training for dense passage retrieval. arXiv preprint arXiv:2108.05540 (2021).

[42] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 6894–6910. https://doi.org/10.18653/v1/2021.emnlp-main.552

[43] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-Augmented Generation for Large Language Models: A Survey. arXiv preprint arXiv:2312.10997 (2023).

[44] Frank Palma Gomez, Ramon Sanabria, Yun-hsuan Sung, Daniel Cer, Siddharth Dalmia, and Gustavo Hernandez Abrego. 2024. Transforming LLMs into Cross-modal and Cross-lingual RetrievalSystems. arXiv preprint arXiv:2404.01616 (2024).

[45] Michael Günther, Jackmin Ong, Isabelle Mohr, Alaeddine Abdessalem, Tanguy Abel, Mohammad Kalim Akram, Susana Guzman, Georgios Mastrapas, Saba Sturua, Bo Wang, et al. 2023. Jina embeddings 2: 8192-token general-purpose text embeddings for long documents. arXiv preprint arXiv:2310.19923 (2023).

[46] Daya Guo, Shuai Lu, Nan Duan, Yanlin Wang, Ming Zhou, and Jian Yin. 2022. Unixcoder: Unified cross-modal pre-training for code representation. arXiv preprint arXiv:2203.03850 (2022).

[47] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. 2020. Graphcodebert: Pre-training code representations with data flow. arXiv preprint arXiv:2009.08366 (2020).

[48] Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. 2024. ToolkenGPT: Augmenting Frozen Language Models with Massive Tools via Tool Embeddings. arXiv:2305.11554 [cs.CL]

[49] Abe Bohan Hou, Orion Weller, Guanghui Qin, Eugene Yang, Dawn Lawrie, Nils Holzenberger, Andrew Blair-Stanek, and Benjamin Van Durme. 2024. CLERC: A Dataset for Legal Case Retrieval and Retrieval-Augmented Analysis Generation. arXiv:2406.17186 (2024). https://doi.org/10.48550/arXiv.2406.17186 arXiv:2406.17186

[50] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zhen Leng Thai, Kaihuo Zhang, Chongyi Wang, Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chaochao Jia, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. MiniCPM: Unveiling the Potential of Small Language Models with Scalable Training Strategies. ArXiv abs/2404.06395 (2024). https://api.semanticscholar.org/CorpusID:269009975

[51] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. arXiv preprint arXiv:2112.09118 (2021).

[52] Vitor Jeronymo, Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, Roberto Lotufo, Jakub Zavrel, and Rodrigo Nogueira. 2023. InPars-v2: Large Language Models as Efficient Dataset Generators for Information Retrieval. arXiv preprint arXiv:2301.01820 (2023).

[53] Vitor Jeronymo, Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, Roberto Lotufo, Jakub Zavrel, and Rodrigo Nogueira. 2023. Inpars-v2: Large language models as efficient dataset generators for information retrieval. arXiv preprint arXiv:2301.01820 (2023).

[54] Yifan Ji, Zhipeng Xu, Zhenghao Liu, Yukun Yan, Shi Yu, Yishan Li, Zhiyuan Liu, Yu Gu, Ge Yu, and Maosong Sun. [n. d.]. Learning More Effective Representations for Dense Retrieval through Deliberate Thinking Before Search. ([n. d.]). https://doi.org/10.48550/arXiv.2502.12974 arXiv:2502.12974 [cs]

[55] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. arXiv:2310.06825 [cs.CL].

[56] Ting Jiang, Shaohan Huang, Zhongzhi Luan, Deqing Wang, and Fuzhen Zhuang. 2023. Scaling sentence embeddings with large language models. arXiv preprint arXiv:2307.16645 (2023).

[57] Ting Jiang, Jian Jiao, Shaohan Huang, Zihan Zhang, Deqing Wang, Fuzhen Zhuang, Furu Wei, Haizhen Huang, Denvy Deng, and Qi Zhang. 2022. Promptbert: Improving bert sentence embeddings with prompts. arXiv preprint arXiv:2201.04337 (2022).

[58] Hongye Jin, Xiaotian Han, Jingfeng Yang, Zhimeng Jiang, Zirui Liu, Chia-Yuan Chang, Huiyuan Chen, and Xia Hu. 2024. Llm maybe longlm: Self-extend llm context window without tuning. arXiv preprint arXiv:2401.01325 (2024).

[59] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. arXiv preprint arXiv:2004.04906 (2020).

[60] Ali Shiraee Kasmaee, Mohammad Khodadad, Mohammad Arshi Saloot, Nicholas Sherck, Stephen Dokas, Hamidreza Mahyar, and Soheila Samiee. 2024. ChemTEB: Chemical Text Embedding Benchmark, an Overview of Embedding Models Performance & Efficiency on a Specific Domain. arXiv preprint arXiv:2412.00532 (2024).

[61] Ekaterina Khramtsova, Shengyao Zhuang, Mahsa Baktashmotlagh, and Guido Zuccon. 2024. Leveraging llms for unsupervised dense retriever ranking. In Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval. 1307–1317.

[62] Junseong Kim, Seolhwa Lee, Jihoon Kwon, Sangmo Gu, Yejin Kim, Minkyung Cho, Jy-yong Sohn, and Chanyeol Choi. 2024. Linq-Embed-Mistral Report. (2024).

[63] Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-thought vectors. Advances in neural information processing systems 28 (2015).

[64] Grigory Kovalev, Mikhail Tikhomirov, Evgeny Kozhevnikov, Max Kornilov, and Natalia Loukachevitch. 2025. Building Russian Benchmark for Evaluation of Information Retrieval Models. arXiv preprint arXiv:2504.12879 (2025).

[65] Aditya Kusupati, Gantavya Bhatt, Aniket Rege, Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen, Sham Kakade, Prateek Jain, et al. 2022. Matryoshka representation learning. Advances in Neural Information Processing Systems 35 (2022), 30233–30249.

[66] Aditya Kusupati, Gantavya Bhatt, Aniket Rege, Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen, Sham Kakade, Prateek Jain, and Ali Farhadi. 2024. Matryoshka Representation Learning. arXiv:2205.13147 [cs.LG] https://arxiv.org/abs/2205.13147

[67] Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024. NV-Embed: Improved Techniques for Training LLMs as Generalist Embedding Models. arXiv preprint arXiv:2405.17428 (2024).

[68] Jinhyuk Lee, Feiyang Chen, Sahil Dua, Daniel Cer, Madhuri Shanbhogue, Iftekhar Naim, Gustavo Hernández Ábrego, Zhe Li, Kaifeng Chen, Henrique Schechter Vera, Xiaoqi Ren, Shanfeng Zhang, Daniel Salz, Michael Boratko, Jay Han, Blair Chen, Shuo Huang, Vikram Rao, Paul Suganthan, Feng Han, Andreas Doumanoglou, Nithi Gupta, Fedor

Moiseev, Cathy Yip, Aashi Jain, Simon Baumgartner, Shahrokh Shahi, Frank Palma Gomez, Sandeep Mariserla, Min Choi, Parashar Shah, Sonam Goenka, Ke Chen, Ye Xia, Koert Chen, Sai Meher Karthik Duddu, Yichang Chen, Trevor Walker, Wenlei Zhou, Rakesh Ghiya, Zach Gleicher, Karan Gill, Zhe Dong, Mojtaba Seyedhosseini, Yunhsuan Sung, Raphael Hoffmann, and Tom Duerig. [n. d.]. Gemini Embedding: Generalizable Embeddings from Gemini. ([n. d.]). https://doi.org/10.48550/arXiv.2503.07891 arXiv:2503.07891 [cs]

[69] Jinhyuk Lee, Zhuyun Dai, Xiaoqi Ren, Blair Chen, Daniel Cer, Jeremy R Cole, Kai Hui, Michael Boratko, Rajvi Kapadia, Wen Ding, et al. 2024. Gecko: Versatile text embeddings distilled from large language models. arXiv preprint arXiv:2403.20327 (2024).

[70] Yibin Lei, Di Wu, Tianyi Zhou, Tao Shen, Yu Cao, Chongyang Tao, and Andrew Yates. 2024. Meta-Task Prompting Elicits Embedding from Large Language Models. arXiv preprint arXiv:2402.18458 (2024).

[71] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. arXiv preprint arXiv:2104.08691 (2021).

[72] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. Advances in Neural Information Processing Systems 33 (2020), 9459–9474.

[73] Chaofan Li, Jianlyu Chen, Yingxia Shao, Defu Lian, and Zheng Liu. 2025. Towards A Generalist Code Embedding Model Based On Massive Data Synthesis. arXiv preprint arXiv:2505.12697 (2025).

[74] Chaofan Li, Zheng Liu, Jianlyv Chen, Defu Lian, and Yingxia Shao. 2025. Reinforced Information Retrieval. arXiv:2502.11562 (2025). https://doi.org/10.48550/arXiv.2502.11562 arXiv:2502.11562

[75] Chaofan Li, Zheng Liu, Shitao Xiao, Yingxia Shao, and Defu Lian. 2024. Llama2vec: Unsupervised adaptation of large language models for dense retrieval. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 3490–3500.

[76] 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 (2024).

[77] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In International conference on machine learning. PMLR, 19730–19742.

[78] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In International conference on machine learning. PMLR, 12888–12900.

[79] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. Advances in neural information processing systems 34 (2021), 9694–9705.

[80] Lei Li, Xiangxu Zhang, Xiao Zhou, and Zheng Liu. 2024. AutoMIR: Effective Zero-Shot Medical Information Retrieval without Relevance Labels. arXiv preprint arXiv:2410.20050 (2024).

[81] Mingxin Li, Zhijie Nie, Yanzhao Zhang, Dingkun Long, Richong Zhang, and Pengjun Xie. 2024. Improving General Text Embedding Model: Tackling Task Conflict and Data Imbalance through Model Merging. arXiv:2410.15035 (2024). arXiv:2410.15035 http://arxiv.org/abs/2410.15035

[82] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023. Starcoder: may the source be with you! arXiv preprint arXiv:2305.06161 (2023).

[83] Ruyu Li, Wenhao Deng, Yu Cheng, Zheng Yuan, Jiaqi Zhang, and Fajie Yuan. 2023. Exploring the upper limits of text-based collaborative filtering using large language models: Discoveries and insights. arXiv preprint arXiv:2305.11700 (2023).

[84] Shengzhe Li, Masaya Ohagi, and Ryokan Ri. 2024. JMTEB: Japanese Massive Text Embedding Benchmark. https://huggingface.co/datasets/sbintuitions/JMTEB. (2024).

[85] Shiyu Li, Yang Tang, Shizhe Chen, and Xi Chen. 2024. Conan-Embedding: General Text Embedding with More and Better Negative Samples. arXiv:2408.15710 (2024). arXiv:2408.15710 http://arxiv.org/abs/2408.15710

[86] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards general text embeddings with multi-stage contrastive learning. arXiv preprint arXiv:2308.03281 (2023).

[87] Ziyue Li and Tianyi Zhou. [n. d.]. Your Mixture-of-Experts LLM Is Secretly an Embedding Model For Free. ([n. d.]). arXiv:2410.10814 http://arxiv.org/abs/2410.10814

[88] Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023. Video-llava: Learning united visual representation by alignment before projection. arXiv preprint arXiv:2311.10122 (2023).

[89] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. Advances in neural information processing systems 36 (2024).

[90] Qijiong Liu, Nuo Chen, Tetsuya Sakai, and Xiao-Ming Wu. 2024. Once: Boosting content-based recommendation with both open-and closed-source large language models. In Proceedings of the 17th ACM International Conference on Web Search and Data Mining. 452–461.

[91] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019. Multi-task deep neural networks for natural language understanding. arXiv preprint arXiv:1901.11504 (2019).

[92] Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. arXiv preprint arXiv:2110.07602 (2021).

[93] Ye Liu, Rui Meng, Shafiq Joty, Silvio Savarese, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. 2024. Codexembed: A generalist embedding model family for multiligual and multi-task code retrieval. arXiv preprint arXiv:2411.12644 (2024).

[94] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. CoRR abs/1907.11692 (2019). arXiv:1907.11692 http://arxiv.org/abs/1907.11692

[95] Kun Luo, Zheng Liu, Shitao Xiao, and Kang Liu. 2024. BGE Landmark Embedding: A Chunking-Free Embedding Method For Retrieval Augmented Long-Context Large Language Models. arXiv preprint arXiv:2402.11573 (2024).

[96] Kun Luo, Minghao Qin, Zheng Liu, Shitao Xiao, Jun Zhao, and Kang Liu. 2024. Large Language Models as Foundations for Next-Gen Dense Retrieval: A Comprehensive Empirical Assessment. arXiv preprint arXiv:2408.12194 (2024).

[97] Simon Lupart, Mohammad Aliannejadi, and Evangelos Kanoulas. 2024. DiSCo Meets LLMs: A Unified Approach for Sparse Retrieval and Contextual Distillation in Conversational Search. arXiv:2410.14609 (2024). arXiv:2410.14609 http://arxiv.org/abs/2410.14609

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

[99] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2023. Video-chatgpt: Towards detailed video understanding via large vision and language models. arXiv preprint arXiv:2306.05424 (2023).

[100] Jean Maillard, Vladimir Karpukhin, Fabio Petroni, Wen-tau Yih, Barlas Oğuz, Veselin Stoyanov, and Gargi Ghosh. 2021. Multi-task retrieval for knowledge-intensive tasks. arXiv preprint arXiv:2101.00117 (2021).

[101] Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. 2017. Learned in translation: Contextualized word vectors. Advances in neural information processing systems 30 (2017).

[102] Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. 2024. Sfrembedding-mistral: enhance text retrieval with transfer learning. Salesforce AI Research Blog 3 (2024).

[103] Tomas Mikolov. 2013. Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781 (2013).

[104] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. Advances in neural information processing systems (2013).

[105] Gabriel de Souza P Moreira, Radek Osmulski, Mengyao Xu, Ronay Ak, Benedikt Schifferer, and Even Oldridge. 2024. NV-Retriever: Improving text embedding models with effective hard-negative mining. arXiv preprint arXiv:2407.15831 (2024).

[106] John X. Morris and Alexander M. Rush. [n. d.]. Contextual Document Embeddings. ([n. d.]). https://doi.org/10.48550/arXiv.2410.02525 arXiv:2410.02525 [cs]

[107] Niklas Muennighoff. 2022. SGPT: GPT Sentence Embeddings for Semantic Search. CoRR abs/2202.08904 (2022). arXiv:2202.08904 https://arxiv.org/abs/2202.08904

[108] Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. 2024. Generative representational instruction tuning. arXiv preprint arXiv:2402.09906 (2024).

[109] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. MTEB: Massive text embedding benchmark. arXiv preprint arXiv:2210.07316 (2022).

[110] Andrew Neeser, Kaylen Latimer, Aadyant Khatri, Chris Latimer, and Naren Ramakrishnan. [n. d.]. QuOTE: Question-Oriented Text Embeddings. ([n. d.]). https://doi.org/10.48550/arXiv.2502.10976 arXiv:2502.10976 [cs]

[111] Jianmo Ni, Gustavo Hernandez Abrego, Noah Constant, Ji Ma, Keith Hall, Daniel Cer, and Yinfei Yang. 2022. Sentence-T5: Scalable Sentence Encoders from Pre-trained Text-to-Text Models. In Findings of the Association for Computational Linguistics: ACL 2022, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, Dublin, Ireland, 1864–1874. https://doi.org/10.18653/v1/2022.findings-acl.146

[112] Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Ábrego, Ji Ma, Vincent Y Zhao, Yi Luan, Keith B Hall, Ming-Wei Chang, et al. 2021. Large dual encoders are generalizable retrievers. arXiv preprint arXiv:2112.07899 (2021).

[113] Zhijie Nie, Richong Zhang, and Zhanyu Wu. 2024. A Text is Worth Several Tokens: Text Embedding from LLMs Secretly Aligns Well with The Key Tokens. arXiv preprint arXiv:2406.17378 (2024).

[114] Zach Nussbaum, John X Morris, Brandon Duderstadt, and Andriy Mulyar. 2024. Nomic embed: Training a reproducible long context text embedder. arXiv preprint arXiv:2402.01613 (2024).

[115] Hanseok Oh, Hyunji Lee, Seonghyeon Ye, Haebin Shin, Hansol Jang, Changwook Jun, and Minjoon Seo. 2024. INSTRUCTIR: A Benchmark for Instruction Following of Information Retrieval Models. arXiv:2402.14334 (2024). https://doi.org/10.48550/arXiv.2402.14334 arXiv:2402.14334

[116] OpenAI. 2023. GPT-4 Technical Report. ArXiv abs/2303.08774 (2023).

[117] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. Advances in Neural Information Processing Systems 35 (2022), 27730–27744.

[118] Minsu Park, Seyeon Choi, Chanyeol Choi, Junseong Kim, and Jy yong Sohn. 2024. Improving Multi-lingual Alignment Through Soft Contrastive Learning. ArXiv abs/2405.16155 (2024). https://api.semanticscholar.org/CorpusID:270063220

[119] Wenjun Peng, Derong Xu, Tong Xu, Jianjin Zhang, and Enhong Chen. 2023. Are gpt embeddings useful for ads and recommendation?. In International Conference on Knowledge Science, Engineering and Management. Springer, 151–162.

[120] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 1532–1543.

[121] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. CoRR abs/1802.05365 (2018). arXiv:1802.05365 http://arxiv.org/abs/1802.05365

[122] Alina Petukhova, Joao P Matos-Carvalho, and Nuno Fachada. 2024. Text clustering with LLM embeddings. arXiv preprint arXiv:2403.15112 (2024).

[123] Alina Petukhova, Joao P Matos-Carvalho, and Nuno Fachada. 2025. Text clustering with large language model embeddings. International Journal of Cognitive Computing in Engineering 6 (2025), 100–108.

[124] Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? arXiv preprint arXiv:1906.01502 (2019).

[125] Jacob Portes, Alexander Trott, Sam Havens, Daniel King, Abhinav Venigalla, Moin Nadeem, Nikhil Sardana, Daya Khudia, and Jonathan Frankle. 2024. MosaicBERT: a bidirectional encoder optimized for fast pretraining. Advances in Neural Information Processing Systems 36 (2024).

[126] Rafał Poświata, Sławomir Dadas, and Michał Perełkiewicz. 2024. PL-MTEB: Polish Massive Text Embedding Benchmark. arXiv preprint arXiv:2405.10138 (2024).

[127] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2020. RocketQA: An optimized training approach to dense passage retrieval for open-domain question answering. arXiv preprint arXiv:2010.08191 (2020).

[128] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In International conference on machine learning. PMLR, 8748–8763.

[129] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).

[130] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. arXiv preprint arXiv:2112.11446 (2021).

[131] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of machine learning research 21, 140 (2020), 1–67.

[132] Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2022. Parallel context windows for large language models. arXiv preprint arXiv:2212.10947 (2022).

[133] Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. RocketQAv2: A Joint Training Method for Dense Passage Retrieval and Passage Re-ranking. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, 2825–2835. https://doi.org/10.18653/v1/2021.emnlp-main.224

[134] Paul K Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. 2023. Audiopalm: A large language model that can speak and listen. arXiv preprint arXiv:2306.12925 (2023).

[135] Shahnewaz Karim Sakib and Anindya Bijoy Das. 2024. Challenging Fairness: A Comprehensive Exploration of Bias in LLM-Based Recommendations. arXiv preprint arXiv:2409.10825 (2024).

[136] Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. arXiv preprint arXiv:2110.08207 (2021).

[137] Fabian David Schmidt, Philipp Borchert, Ivan Vulić, and Goran Glavaš. 2024. Self-Distillation for Model Stacking Unlocks Cross-Lingual NLU in 200+ Languages. arXiv preprint arXiv:2406.12739 (2024).

[138] Ninglu Shao, Shitao Xiao, Zheng Liu, and Peitian Zhang. 2024. Extensible Embedding: A Flexible Multipler For LLM's Context Length. arXiv preprint arXiv:2402.11577 (2024).

[139] Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Xiaolong Huang, Binxing Jiao, Linjun Yang, and Daxin Jiang. 2022. LexMAE: Lexicon-Bottlenecked Pretraining for Large-Scale Retrieval. In The Eleventh International Conference on Learning Representations.

[140] Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Guodong Long, Kai Zhang, and Daxin Jiang. 2023. Unifier: A unified retriever for large-scale retrieval. In Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. 4787–4799.

[141] Sonish Sivarajkumar, Haneef Ahamed Mohammad, David Oniani, Kirk Roberts, William Hersh, Hongfang Liu, Daqing He, Shyam Visweswaran, and Yanshan Wang. 2024. Clinical information retrieval: A literature review. Journal of Healthcare Informatics Research (2024), 1–40.

[142] Tingyu Song, Guo Gan, Mingsheng Shang, and Yilun Zhao. 2025. IFIR: A Comprehensive Benchmark for Evaluating Instruction-Following in Expert-Domain Information Retrieval. arXiv:2503.04644 (2025). https://doi.org/10.48550/arXiv.2503.04644 arXiv:2503.04644

[143] Jacob Mitchell Springer, Suhas Kotha, Daniel Fried, Graham Neubig, and Aditi Raghunathan. 2024. Repetition Improves Language Model Embeddings. arXiv preprint arXiv:2402.15449 (2024).

[144] Saba Sturua, Isabelle Mohr, Mohammad Kalim Akram, Michael Günther, Bo Wang, Markus Krimmel, Feng Wang, Georgios Mastrapas, Andreas Koukounas, Nan Wang, and Han Xiao. 2024. Jina-Embeddings-v3: Multilingual Embeddings With Task LoRA. arXiv:2409.10173 http://arxiv.org/abs/2409.10173

[145] Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen tau Yih, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. 2023. One Embedder, Any Task: Instruction-Finetuned Text Embeddings. arXiv:2212.09741 [cs.CL] https://arxiv.org/abs/2212.09741

[146] Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen-tau Yih, Noah A Smith, Luke Zettlemoyer, and Tao Yu. 2023. One Embedder, Any Task: Instruction-Finetuned Text Embeddings. In Findings of the Association for Computational Linguistics: ACL 2023. 1102–1121.

[147] Hongjin Su, Howard Yen, Mengzhou Xia, Weijia Shi, Niklas Muennighoff, Han-yu Wang, Haisu Liu, Quan Shi, Zachary S. Siegel, Michael Tang, Ruoxi Sun, Jinsung Yoon, Sercan O. Arik, Danqi Chen, and Tao Yu. 2025. BRIGHT: A Realistic and Challenging Benchmark for Reasoning-Intensive Retrieval. arXiv:2407.12883 (2025). https://doi.org/10.48550/arXiv.2407.12883 arXiv:2407.12883

[148] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. Neurocomputing 568 (2024), 127063.

[149] Hao Sun, Yunyi Shen, and Jean-Francois Ton. 2025. Rethinking reward modeling in preference-based large language model alignment. In The Thirteenth International Conference on Learning Representations.

[150] Hao Sun, Yunyi Shen, Jean-Francois Ton, and Mihaela van der Schaar. 2025. Reusing Embeddings: Reproducible Reward Model Research in Large Language Model Alignment without GPUs. arXiv preprint arXiv:2502.04357 (2025).

[151] Weiwei Sun, Zhengliang Shi, Jiulong Wu, Lingyong Yan, Xinyu Ma, Yiding Liu, Min Cao, Dawei Yin, and Zhaochun Ren. 2024. MAIR: A Massive Benchmark for Evaluating Instructed Retrieval. arXiv:2410.10127 (2024). https://doi.org/10.48550/arXiv.2410.10127 arXiv:2410.10127

[152] Zeinab Sadat Taghavi, Ali Modarressi, Yunpu Ma, and Hinrich Schütze. [n. d.]. ImpliRet: Benchmarking the Implicit Fact Retrieval Challenge. ([n. d.]). https://doi.org/10.48550/arXiv.2506.14407 arXiv:2506.14407 [cs]

[153] Yixuan Tang and Yi Yang. 2025. Finmteb: Finance massive text embedding benchmark. arXiv preprint arXiv:2502.10990 (2025).

[154] Chongyang Tao, Chang Liu, Tao Shen, Can Xu, Xiubo Geng, Binxing Jiao, and Daxin Jiang. 2024. ADAM: Dense Retrieval Distillation with Adaptive Dark Examples. In Findings of the Association for Computational Linguistics ACL 2024. 11639–11651.

[155] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford Alpaca: An Instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca.

[156] Atula Tejaswi, Yoonsang Lee, Sujay Sanghavi, and Eunsol Choi. [n. d.]. RARe: Retrieval Augmented Retrieval with In-Context Examples. ([n. d.]). https://doi.org/10.48550/arXiv.2410.20088 arXiv:2410.20088 [cs]

[157] Raghuveer Thirukovalluru and Bhuwan Dhingra. 2024. GenEOL: Harnessing the Generative Power of LLMs for Training-Free Sentence Embeddings. arXiv:2410.14635 (2024). arXiv:2410.14635 http://arxiv.org/abs/2410.14635

[158] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications. arXiv preprint arXiv:2201.08239 (2022).

[159] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023).

[160] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023).

[161] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023).

[162] Hayato Tsukagoshi, Ryohei Sasano, and Koichi Takeda. 2021. Defsent: Sentence embeddings using definition sentences. arXiv preprint arXiv:2105.04339 (2021).

[163] Oleg Vasilyev, Randy Sawaya, and John Bohannon. 2024. Preserving Multilingual Quality While Tuning Query Encoder on English Only. arXiv preprint arXiv:2407.00923 (2024).

[164] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing systems 30 (2017).

[165] Voyage-AI. 2024. voyage-large-2-instruct: Instruction-tuned and rank 1 on MTEB.

[166] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Simlm: Pre-training with representation bottleneck for dense passage retrieval. arXiv preprint arXiv:2207.02578 (2022).

[167] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. arXiv preprint arXiv:2212.03533 (2022).

[168] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. Improving text embeddings with large language models. arXiv preprint arXiv:2401.00368 (2024).

[169] Suyuchen Wang, Ivan Kobyzev, Peng Lu, Mehdi Rezagholizadeh, and Bang Liu. 2024. Resonance rope: Improving context length generalization of large language models. arXiv preprint arXiv:2403.00071 (2024).

[170] Xiaoyue Wang, Jianyou Wang, Weili Cao, Kaicheng Wang, Ramamohan Paturi, and Leon Bergen. 2024. BIRCO: A Benchmark of Information Retrieval Tasks with Complex Objectives. arXiv:2402.14151 (2024). https://doi.org/10.48550/arXiv.2402.14151 arXiv:2402.14151

[171] Silvan Wehrli, Bert Arnrich, and Christopher Irrgang. 2024. German text embedding clustering benchmark. arXiv preprint arXiv:2401.02709 (2024).

[172] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. arXiv preprint arXiv:2109.01652 (2021).

[173] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent Abilities of Large Language Models. Trans. Mach. Learn. Res. 2022 (2022). https://openreview.net/forum?id=yzkSU5zdwD

[174] Orion Weller, Benjamin Chang, Sean MacAvaney, Kyle Lo, Arman Cohan, Benjamin Van Durme, Dawn Lawrie, and Luca Soldaini. 2024. FollowIR: Evaluating and Teaching Information Retrieval Models to Follow Instructions. arXiv:2403.15246 (2024). https://doi.org/10.48550/arXiv.2403.15246 arXiv:2403.15246

[175] Orion Weller, Benjamin Chang, Eugene Yang, Mahsa Yarmohammadi, Sam Barham, Sean MacAvaney, Arman Cohan, Luca Soldaini, Benjamin Van Durme, and Dawn Lawrie. 2025. mFollowIR: A Multilingual Benchmark for Instruction Following in Retrieval. arXiv:2501.19264 (2025). https://doi.org/10.48550/arXiv.2501.19264 arXiv:2501.19264

[176] BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. arXiv preprint arXiv:2211.05100 (2022).

[177] Chenghao Xiao, G Thomas Hudson, and Noura Al Moubayed. 2024. Rar-b: Reasoning as retrieval benchmark. arXiv preprint arXiv:2404.06347 (2024).

[178] Shitao Xiao, Zheng Liu, Yingxia Shao, and Zhao Cao. 2022. RetroMAE: Pre-training retrieval-oriented language models via masked auto-encoder. arXiv preprint arXiv:2205.12035 (2022).

[179] Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. 2023. C-pack: Packaged resources to advance general chinese embedding. arXiv preprint arXiv:2309.07597 (2023).

[180] Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. 2024. C-pack: Packed resources for general chinese embeddings. In Proceedings of the 47th international ACM SIGIR conference on research and development in information retrieval. 641–649.

[181] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. Approximate nearest neighbor negative contrastive learning for dense text retrieval. arXiv preprint arXiv:2007.00808 (2020).

[182] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2023. WizardLM: Empowering large pre-trained language models to follow complex instructions. In The Twelfth International Conference on Learning Representations.

[183] Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. 2024. A survey on knowledge distillation of large language models. arXiv preprint arXiv:2402.13116 (2024).

[184] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2020. mT5: A massively multilingual pre-trained text-to-text transformer. arXiv preprint arXiv:2010.11934 (2020).

[185] Ruiran Yan, Zheng Liu, and Defu Lian. [n. d.]. O1 Embedder: Let Retrievers Think Before Action. ([n. d.]). https://doi.org/10.48550/arXiv.2502.07555 arXiv:2502.07555 [cs]

[186] Dongkeun Yoon, Joel Jang, Sungdong Kim, Seungone Kim, Sheikh Shafayat, and Minjoon Seo. 2024. Langbridge: Multilingual reasoning without multilingual supervision. arXiv preprint arXiv:2401.10695 (2024).

[187] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. Coca: Contrastive captioners are image-text foundation models. arXiv preprint arXiv:2205.01917 (2022).

[188] Munazza Zaib, Quan Z Sheng, and Wei Emma Zhang. 2020. A short survey of pre-trained language models for conversational ai-a new age in nlp. In Proceedings of the Australasian computer science week multiconference. 1–4.

[189] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. arXiv preprint arXiv:2210.02414 (2022).

[190] Bowen Zhang, Kehua Chang, and Chunping Li. 2024. Simple Techniques for Enhancing Sentence Embeddings in Generative Language Models. arXiv preprint arXiv:2404.03921 (2024).

[191] Hang Zhang, Xin Li, and Lidong Bing. 2023. Video-llama: An instruction-tuned audio-visual language model for video understanding. arXiv preprint arXiv:2306.02858 (2023).

[192] Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Dou, and Jian-Yun Nie. 2023. Retrieve anything to augment large language models. arXiv preprint arXiv:2310.07554 (2023).

[193] Xin Zhang, Yanzhao Zhang, Dingkun Long, Wen Xie, Ziqi Dai, Jialong Tang, Huan Lin, Baosong Yang, Pengjun Xie, Fei Huang, et al. 2024. mGTE: Generalized Long-Context Text Representation and Reranking Models for Multilingual Text Retrieval. arXiv preprint arXiv:2407.19669 (2024).

[194] Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. [n. d.]. Qwen3 Embedding: Advancing Text Embedding and Reranking Through Foundation Models. ([n. d.]). https://doi.org/10.48550/arXiv.2506.05176 arXiv:2506.05176 [cs]

[195] Zihuai Zhao, Wenqi Fan, Jiatong Li, Yunqing Liu, Xiaowei Mei, Yiqi Wang, Zhen Wen, Fei Wang, Xiangyu Zhao, Jiliang Tang, et al. 2024. Recommender systems in the era of large language models (llms). IEEE Transactions on Knowledge and Data Engineering (2024).

[196] Lucia Zheng, Neel Guha, Javokhir Arifov, Sarah Zhang, Michal Skreta, Christopher D. Manning, Peter Henderson, and Daniel E. Ho. 2025. A Reasoning-Focused Legal Retrieval Benchmark. In Proceedings of the Symposium on Computer Science and Law on ZZZ. ACM, Munich Germany, 169–193. https://doi.org/10.1145/3709025.3712219

[197] Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, LILI YU, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. LIMA: Less Is More for Alignment. In Thirty-seventh Conference on Neural Information Processing Systems. https://openreview.net/forum?id=KBMOKmX2he

[198] Junjie Zhou, Zheng Liu, Shitao Xiao, Bo Zhao, and Yongping Xiong. 2024. VISTA: Visualized Text Embedding For Universal Multi-Modal Retrieval. arXiv preprint arXiv:2406.04292 (2024).

[199] Jianqun Zhou, Yuanlei Zheng, Wei Chen, Qianqian Zheng, Hui Su, Wei Zhang, Rui Meng, and Xiaoyu Shen. 2025. Beyond Content Relevance: Evaluating Instruction Following in Retrieval Models. arXiv:2410.23841 (2025). https://doi.org/10.48550/arXiv.2410.23841 arXiv:2410.23841

[200] Yucheng Zhou, Tao Shen, Xiubo Geng, Chongyang Tao, Can Xu, Guodong Long, Binxing Jiao, and Daxin Jiang. 2022. Towards robust ranker for text retrieval. arXiv preprint arXiv:2206.08063 (2022).

[201] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592 (2023).

[202] Dawei Zhu, Liang Wang, Nan Yang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. 2024. LongEmbed: Extending Embedding Models for Long Context Retrieval. arXiv preprint arXiv:2404.12096 (2024).

[203] 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, Yaser Al-Onaizan, Mohit

Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 4375–4391. https://aclanthology.org/2024.emnlp-main.250

[204] Erfan Zinvandi, Morteza Alikhani, Mehran Sarmadi, Zahra Pourbahman, Sepehr Arvin, Reza Kazemi, and Arash Amini. 2025. Famteb: Massive text embedding benchmark in persian language. arXiv preprint arXiv:2502.11571 (2025).

[205] Husein Zolkepli, Aisyah Razak, Kamarul Adha, and Ariff Nazhan. 2024. Multi-Lingual Malaysian Embedding: Leveraging Large Language Models for Semantic Representations. arXiv preprint arXiv:2402.03053 (2024).