

Hugging Face

- [Models](#)
- [Datasets](#)
- [Spaces](#)
- [Posts](#)
- [Docs](#)
- [Pricing](#)
-
-
-



[BAAI](#)

/

[bge-large-zh-v1.5](#)

[like224](#)

[Feature Extraction](#)[sentence](#)-

[transformers](#)[PyTorch](#)[Safetensors](#)[Transformers](#)[Chinese](#)[bert](#)[sentence-similarity](#)[Inference](#)

[Endpoints](#)[5 papers](#)

[License:mit](#)

[Model card](#)[Files and versions](#)[Community](#)

[15](#)

Deploy

Use in libraries

[Edit model card](#)

FlagEmbedding

[Model List](#) | [FAQ](#) | [Usage](#) | [Evaluation](#) | [Train](#) | [Contact](#) | [Citation](#) | [License](#)

For more details please refer to our Github: [FlagEmbedding](#).

If you are looking for a model that supports more languages, longer texts, and other retrieval methods, you can try using [bge-m3](#).

[English](#) | [中文](#)

FlagEmbedding focuses on retrieval-augmented LLMs, consisting of the following projects currently:

- **Long-Context LLM:** [Activation Beacon](#)
- **Fine-tuning of LM :** [LM-Cocktail](#)
- **Dense Retrieval:** [BGE-M3](#), [LLM Embedder](#), [BGE Embedding](#)
- **Reranker Model:** [BGE Reranker](#)
- **Benchmark:** [C-MTEB](#)

News

- 1/30/2024: Release **BGE-M3**, a new member to BGE model series! M3 stands for Multi-linguality (100+ languages), Multi-granularities (input length up to 8192), Multi-Functionality (unification of dense, lexical, multi-vec/colbert retrieval). It is the first embedding model which supports all three retrieval methods, achieving new SOTA on multi-lingual (MIRACL) and cross-lingual (MKQA) benchmarks. [Technical Report](#) and [Code](#). :fire:
- 1/9/2024: Release [Activation-Beacon](#), an effective, efficient, compatible, and low-cost (training) method to extend the context length of LLM. [Technical Report](#) :fire:
- 12/24/2023: Release **LLaRA**, a LLaMA-7B based dense retriever, leading to state-of-the-art performances on MS MARCO and BEIR. Model and code will be open-sourced. Please stay tuned. [Technical Report](#) :fire:
- 11/23/2023: Release [LM-Cocktail](#), a method to maintain general capabilities during fine-tuning by merging multiple language models. [Technical Report](#) :fire:
- 10/12/2023: Release [LLM-Embedder](#), a unified embedding model to support diverse retrieval augmentation needs for LLMs. [Technical Report](#)
- 09/15/2023: The [technical report](#) and [massive training data](#) of BGE has been released
- 09/12/2023: New models:

- **New reranker model:** release cross-encoder models [BAAI/bge-reranker-base](#) and [BAAI/bge-reranker-large](#), which are more powerful than embedding model. We recommend to use/fine-tune them to re-rank top-k documents returned by embedding models.

- **update embedding model**: release `bge-*/v1.5` embedding model to alleviate the issue of the similarity distribution, and enhance its retrieval ability without instruction.

More

-
-
-
-
-

Model List

`bge` is short for BAAI general embedding.

Model	Language		Description	query instruction for retrieval [1]
BAAI/bge-m3	Multilingual	Inference Fine-tune	Multi-Functionality(dense retrieval, sparse retrieval, multi-vector(colbert)), Multi-Linguality, and Multi-Granularity(8192 tokens)	
BAAI/llm-embedder	English	Inference Fine-tune	a unified embedding model to support diverse retrieval augmentation needs for LLMs	See README
BAAI/bge-reranker-large	Chinese and English	Inference Fine-tune	a cross-encoder model which is more accurate but less efficient [2]	

Model	Language		Description	query instruction for retrieval [1]
BAAI/bge-reranker-base	Chinese and English	Inference Fine-tune	a cross-encoder model which is more accurate but less efficient [2]	
BAAI/bge-large-en-v1.5	English	Inference Fine-tune	version 1.5 with more reasonable similarity distribution	Represent this sentence for searching relevant passages:
BAAI/bge-base-en-v1.5	English	Inference Fine-tune	version 1.5 with more reasonable similarity distribution	Represent this sentence for searching relevant passages:
BAAI/bge-small-en-v1.5	English	Inference Fine-tune	version 1.5 with more reasonable similarity distribution	Represent this sentence for searching relevant passages:

Model	Language		Description	query instruction for retrieval [1]
BAAI/bge-large-zh-v1.5	Chinese	Inference Fine-tune	version 1.5 with more reasonable similarity distribution	为这个句子生成表示以用于检索相关文章：
BAAI/bge-base-zh-v1.5	Chinese	Inference Fine-tune	version 1.5 with more reasonable similarity distribution	为这个句子生成表示以用于检索相关文章：
BAAI/bge-small-zh-v1.5	Chinese	Inference Fine-tune	version 1.5 with more reasonable similarity distribution	为这个句子生成表示以用于检索相关文章：
BAAI/bge-large-en	English	Inference Fine-tune	:trophy: rank 1st in MTEB leaderboard	Represent this sentence for searching relevant passages:
BAAI/bge-base-en	English	Inference Fine-tune	a base-scale model but with similar ability to bge-large-en	Represent this sentence for

Model	Language		Description	query instruction for retrieval [1]
				searching relevant passages:
BAAI/bge-small-en	English	Inference Fine-tune	a small-scale model but with competitive performance	Represent this sentence for searching relevant passages:
BAAI/bge-large-zh	Chinese	Inference Fine-tune	:trophy: rank 1st in C-MTEB benchmark	为这个句子生成表示以用于检索相关文章：
BAAI/bge-base-zh	Chinese	Inference Fine-tune	a base-scale model but with similar ability to bge-large-zh	为这个句子生成表示以用于检索相关文章：
BAAI/bge-small-zh	Chinese	Inference Fine-tune	a small-scale model but with competitive performance	为这个句子生成表示以用于检索相关文章：

[1]: If you need to search the relevant passages to a query, we suggest to add the instruction to the query; in other cases, no instruction is needed, just use the original query directly. In all cases, **no instruction** needs to be added to passages.

[2]: Different from embedding model, reranker uses question and document as input and directly output similarity instead of embedding. To balance the accuracy and time cost, cross-encoder is widely used to re-rank top-k documents retrieved by other simple models. For examples, use bge embedding model to retrieve top 100 relevant documents, and then use bge reranker to re-rank the top 100 document to get the final top-3 results.

All models have been uploaded to Huggingface Hub, and you can see them at <https://huggingface.co/BAAI>. If you cannot open the Huggingface Hub, you also can download the models at <https://model.baai.ac.cn/models>.

Frequently asked questions

1. How to fine-tune bge embedding model?

-
-
-

2. The similarity score between two dissimilar sentences is higher than 0.5

3. When does the query instruction need to be used

Usage

Usage for Embedding Model

Here are some examples for using `bge` models with [FlagEmbedding](#), [Sentence-Transformers](#), [Langchain](#), or [Huggingface Transformers](#).

Using FlagEmbedding

```
pip install -U FlagEmbedding
```

If it doesn't work for you, you can see [FlagEmbedding](#) for more methods to install FlagEmbedding.

```
from FlagEmbedding import FlagModel

sentences_1 = ["样例数据-1", "样例数据-2"]

sentences_2 = ["样例数据-3", "样例数据-4"]
```

```

model = FlagModel('BAAI/bge-large-zh-v1.5',
                  query_instruction_for_retrieval="为这个句子生成表示以用于检索相关文章: ",
                  use_fp16=True) # Setting use_fp16 to True speeds up computation with a slight performance degradation

embeddings_1 = model.encode(sentences_1)

embeddings_2 = model.encode(sentences_2)

similarity = embeddings_1 @ embeddings_2.T

print(similarity)

# for s2p(short query to long passage) retrieval task, suggest to use encode_queries() which will automatically add the instruction to each query

# corpus in retrieval task can still use encode() or encode_corpus(), since they don't need instruction

queries = ['query_1', 'query_2']

passages = ["样例文档-1", "样例文档-2"]

q_embeddings = model.encode_queries(queries)

p_embeddings = model.encode(passages)

scores = q_embeddings @ p_embeddings.T

```

For the value of the argument `query_instruction_for_retrieval`, see [Model List](#).

By default, FlagModel will use all available GPUs when encoding. Please set `os.environ["CUDA_VISIBLE_DEVICES"]` to select specific GPUs. You also can set `os.environ["CUDA_VISIBLE_DEVICES"]=" "` to make all GPUs unavailable.

Using Sentence-Transformers

You can also use the `bge` models with [sentence-transformers](#):

```
pip install -U sentence-transformers

from sentence_transformers import SentenceTransformer

sentences_1 = ["样例数据-1", "样例数据-2"]

sentences_2 = ["样例数据-3", "样例数据-4"]

model = SentenceTransformer('BAAI/bge-large-zh-v1.5')

embeddings_1 = model.encode(sentences_1, normalize_embeddings=True)

embeddings_2 = model.encode(sentences_2, normalize_embeddings=True)

similarity = embeddings_1 @ embeddings_2.T

print(similarity)
```

For s2p(short query to long passage) retrieval task, each short query should start with an instruction (instructions see [Model List](#)). But the instruction is not needed for passages.

```
from sentence_transformers import SentenceTransformer

queries = ['query_1', 'query_2']

passages = ["样例文档-1", "样例文档-2"]

instruction = "为这个句子生成表示以用于检索相关文章: "

model = SentenceTransformer('BAAI/bge-large-zh-v1.5')

q_embeddings = model.encode([instruction+q for q in queries], normalize_embeddings=True)

p_embeddings = model.encode(passages, normalize_embeddings=True)

scores = q_embeddings @ p_embeddings.T
```

Using Langchain

You can use `bge` in langchain like this:

```
from langchain.embeddings import HuggingFaceBgeEmbeddings

model_name = "BAAI/bge-large-en-v1.5"

model_kwargs = {'device': 'cuda'}

encode_kwargs = {'normalize_embeddings': True} # set True to compute cosine similarity

model = HuggingFaceBgeEmbeddings(
    model_name=model_name,
    model_kwargs=model_kwargs,
    encode_kwargs=encode_kwargs,
    query_instruction="为这个句子生成表示以用于检索相关文章: "
)

model.query_instruction = "为这个句子生成表示以用于检索相关文章: "
```

Using HuggingFace Transformers

With the transformers package, you can use the model like this: First, you pass your input through the transformer model, then you select the last hidden state of the first token (i.e., [CLS]) as the sentence embedding.

```
from transformers import AutoTokenizer, AutoModel

import torch

# Sentences we want sentence embeddings for

sentences = ["样例数据-1", "样例数据-2"]
```

```
# Load model from HuggingFace Hub
```

```
tokenizer = AutoTokenizer.from_pretrained('BAAI/bge-large-zh-v1.5')
```

```
model = AutoModel.from_pretrained('BAAI/bge-large-zh-v1.5')
```

```
model.eval()
```

```
# Tokenize sentences
```

```
encoded_input = tokenizer(sentences, padding=True, truncation=True, return_tensors='pt')
```

```
# for s2p(short query to long passage) retrieval task, add an instruction to query (not add instruction for passages)
```

```
# encoded_input = tokenizer([instruction + q for q in queries], padding=True, truncation=True, return_tensors='pt')
```

```
# Compute token embeddings
```

```
with torch.no_grad():
```

```
    model_output = model(**encoded_input)
```

```
    # Perform pooling. In this case, cls pooling.
```

```
    sentence_embeddings = model_output[0][:, 0]
```

```
# normalize embeddings
```

```
sentence_embeddings = torch.nn.functional.normalize(sentence_embeddings, p=2, dim=1)
```

```
print("Sentence embeddings:", sentence_embeddings)
```

Usage for Reranker

Different from embedding model, reranker uses question and document as input and directly output similarity instead of embedding. You can get a relevance score by inputting query and

passage to the reranker. The reranker is optimized based cross-entropy loss, so the relevance score is not bounded to a specific range.

Using FlagEmbedding

```
pip install -U FlagEmbedding
```

Get relevance scores (higher scores indicate more relevance):

```
from FlagEmbedding import FlagReranker
```

```
reranker = FlagReranker('BAAI/bge-reranker-large', use_fp16=True) # Setting use_fp16 to True speeds up computation with a slight performance degradation
```

```
score = reranker.compute_score(['query', 'passage'])
```

```
print(score)
```

```
scores = reranker.compute_score([[['what is panda?', 'hi'], ['what is panda?', 'The giant panda (Ailuropoda melanoleuca), sometimes called a panda bear or simply panda, is a bear species endemic to China.']]])
```

```
print(scores)
```

Using Huggingface transformers

```
import torch
```

```
from transformers import AutoModelForSequenceClassification, AutoTokenizer
```

```
tokenizer = AutoTokenizer.from_pretrained('BAAI/bge-reranker-large')
```

```
model = AutoModelForSequenceClassification.from_pretrained('BAAI/bge-reranker-large')
```

```
model.eval()
```

```
pairs = [['what is panda?', 'hi'], ['what is panda?', 'The giant panda (Ailuropoda melanoleuca), sometimes called a panda bear or simply panda, is a bear species endemic to China.']]
```

```
with torch.no_grad():
```

```
    inputs = tokenizer(pairs, padding=True, truncation=True, return_tensors='pt', max_length=512)
```

```
    scores = model(**inputs, return_dict=True).logits.view(-1, ).float()
```

```
    print(scores)
```

Evaluation

baai-general-embedding models achieve **state-of-the-art performance on both MTEB and C-MTEB leaderboard!** For more details and evaluation tools see our [scripts](#).

- **MTEB:**

Model Name	Dimension	Sequence Length	Average (56)	Retrieval (15)	Clustering (11)	Pair Classification (3)	Reranking (4)	STS (10)
BAAI/bge-large-en-v1.5	1024	512	64.23	54.29	46.08	87.12	60.03	83.11
BAAI/bge-base-en-v1.5	768	512	63.55	53.25	45.77	86.55	58.86	82.4

Model Name	Dimension	Sequence Length	Average (56)	Retrieval (15)	Clustering (11)	Pair Classification (3)	Reranking (4)	STS (10)
BAAI/bge-small-en-v1.5	384	512	62.17	51.68	43.82	84.92	58.36	81.59
bge-large-en	1024	512	63.98	53.9	46.98	85.8	59.48	81.56
bge-base-en	768	512	63.36	53.0	46.32	85.86	58.7	81.84
gte-large	1024	512	63.13	52.22	46.84	85.00	59.13	83.35
gte-base	768	512	62.39	51.14	46.2	84.57	58.61	82.3
e5-large-v2	1024	512	62.25	50.56	44.49	86.03	56.61	82.05
bge-small-en	384	512	62.11	51.82	44.31	83.78	57.97	80.72

Model Name	Dimension	Sequence Length	Average (56)	Retrieval (15)	Clustering (11)	Pair Classification (3)	Reranking (4)	STS (10)
instructor-xl	768	512	61.79	49.26	44.74	86.62	57.29	83.06
e5-base-v2	768	512	61.5	50.29	43.80	85.73	55.91	81.05
gte-small	384	512	61.36	49.46	44.89	83.54	57.7	82.07
text-embedding-ada-002	1536	8192	60.99	49.25	45.9	84.89	56.32	80.97
e5-small-v2	384	512	59.93	49.04	39.92	84.67	54.32	80.39
sentence-t5-xxl	768	512	59.51	42.24	43.72	85.06	56.42	82.63
all-mpnet-base-v2	768	514	57.78	43.81	43.69	83.04	59.36	80.28

Model Name	Dimension	Sequence Length	Average (56)	Retrieval (15)	Clustering (11)	Pair Classification (3)	Reranking (4)	STS (10)
sgpt-bloom-7b1-msmarco	4096	2048	57.59	48.22	38.93	81.9	55.65	77.74

- **C-MTEB:**

We create the benchmark C-MTEB for Chinese text embedding which consists of 31 datasets from 6 tasks. Please refer to [C_MTEB](#) for a detailed introduction.

Model	Embedding dimension	Avg	Retrieval	STS	PairClassification	Classification	Reranking	Clus
BAAI/bge-large-zh-v1.5	1024	64.53	70.46	56.25	81.6	69.13	65.84	48
BAAI/bge-base-zh-v1.5	768	63.13	69.49	53.72	79.75	68.07	65.39	47
BAAI/bge-small-zh-v1.5	512	57.82	61.77	49.11	70.41	63.96	60.92	44
BAAI/bge-large-zh	1024	64.20	71.53	54.98	78.94	68.32	65.11	48

Model	Embedding dimension	Avg	Retrieval	STS	PairClassification	Classification	Reranking	Clus
<u>bge-large-zh-noinstruct</u>	1024	63.53	70.55	53	76.77	68.58	64.91	50
<u>BAAI/bge-base-zh</u>	768	62.96	69.53	54.12	77.5	67.07	64.91	47
<u>multilingual-e5-large</u>	1024	58.79	63.66	48.44	69.89	67.34	56.00	48
<u>BAAI/bge-small-zh</u>	512	58.27	63.07	49.45	70.35	63.64	61.48	45
<u>m3e-base</u>	768	57.10	56.91	50.47	63.99	67.52	59.34	47
<u>m3e-large</u>	1024	57.05	54.75	50.42	64.3	68.2	59.66	48
<u>multilingual-e5-base</u>	768	55.48	61.63	46.49	67.07	65.35	54.35	40

Model	Embedding dimension	Avg	Retrieval	STS	PairClassification	Classification	Reranking	Clus
multilingual-e5-small	384	55.38	59.95	45.27	66.45	65.85	53.86	45
text-embedding-ada-002(OpenAI)	1536	53.02	52.0	43.35	69.56	64.31	54.28	45
luotuo	1024	49.37	44.4	42.78	66.62	61	49.25	44
text2vec-base	768	47.63	38.79	43.41	67.41	62.19	49.45	37
text2vec-large	1024	47.36	41.94	44.97	70.86	60.66	49.16	30

- **Reranking:** See [C_MTEB](#) for evaluation script.

Model	T2Reranking	T2RerankingZh2En*	T2RerankingEn2Zh*	MMarcoReranking	CMedQAv1
text2vec-base-multilingual	64.66	62.94	62.51	14.37	48.46

Model	T2Reranking	T2RerankingZh2En*	T2RerankingEn2Zh*	MMarcoReranking	CMedQAv1
multilingual-e5-small	65.62	60.94	56.41	29.91	67.26
multilingual-e5-large	64.55	61.61	54.28	28.6	67.42
multilingual-e5-base	64.21	62.13	54.68	29.5	66.23
m3e-base	66.03	62.74	56.07	17.51	77.05
m3e-large	66.13	62.72	56.1	16.46	77.76
bge-base-zh-v1.5	66.49	63.25	57.02	29.74	80.47
bge-large-zh-v1.5	65.74	63.39	57.03	28.74	83.45

Model	T2Reranking	T2RerankingZh2En*	T2RerankingEn2Zh*	MMarcoReranking	CMedQAv1
BAAI/bge-reranker-base	67.28	63.95	60.45	35.46	81.26
BAAI/bge-reranker-large	67.6	64.03	61.44	37.16	82.15

* : T2RerankingZh2En and T2RerankingEn2Zh are cross-language retrieval tasks

Train

BAAI Embedding

We pre-train the models using [retromae](#) and train them on large-scale pairs data using contrastive learning. **You can fine-tune the embedding model on your data following our examples.** We also provide a [pre-train example](#). Note that the goal of pre-training is to reconstruct the text, and the pre-trained model cannot be used for similarity calculation directly, it needs to be fine-tuned. More training details for bge see [baai_general_embedding](#).

BGE Reranker

Cross-encoder will perform full-attention over the input pair, which is more accurate than embedding model (i.e., bi-encoder) but more time-consuming than embedding model. Therefore, it can be used to re-rank the top-k documents returned by embedding model. We train the cross-encoder on a multilingual pair data, The data format is the same as embedding model, so you can fine-tune it easily following our [example](#). More details please refer to [./FlagEmbedding/reranker/README.md](#)

Contact

If you have any question or suggestion related to this project, feel free to open an issue or pull request. You also can email Shitao Xiao(stxiao@baai.ac.cn) and Zheng Liu(liuzheng@baai.ac.cn).

Citation

If you find this repository useful, please consider giving a star :star: and citation

```
@misc{bge_embedding,  
    title={C-Pack: Packaged Resources To Advance General Chinese Embedding},  
    author={Shitao Xiao and Zheng Liu and Peitian Zhang and Niklas Muennighoff},  
    year={2023},  
    eprint={2309.07597},  
    archivePrefix={arXiv},  
    primaryClass={cs.CL}  
}
```

License

FlagEmbedding is licensed under the [MIT License](#). The released models can be used for commercial purposes free of charge.

Downloads last month

2,383,668

Safetensors

Model size

326M params

Tensor type

I64

.

F32

Inference API

[Feature Extraction](#)

Compute

This model can be loaded on Inference API (serverless).

JSON OutputMaximize

Spaces using BAAI/bge-large-zh-v1.513



[nyust-eb210/bge-large-zh-v1.5_gradio](#)



[TheresaQWQ/BAAI-bge-large-zh-v1.5](#)



[TtKkk/BAAI-bge-large-zh-v1.5](#)



[sejamenath2023/BAAI-bge-large-zh-v1.5](#)



[Zulelee/langchain-chatchat](#)



[gordonchan/bge-embedding](#)



[mengk/BAAI-bge-large-zh-v1.5-new](#)



[abidlabs/mteb-leaderboard](#)



[isiriai/BAAI-bge-large-zh-v1.5](#)



[TtKkk/BAAI-bge-large-zh-v1.5](#)



[Lisang11/BAAI-bge-large-zh-v1.5](#)



[yulone/bge](#)



[littlezebra/hello-embed](#)

© Hugging Face

[TOS](#)[Privacy](#)[About](#)[Jobs](#)[Models](#)[Datasets](#)[Spaces](#)[Pricing](#)[Docs](#)