

[ ![Logo](../../\_static/logo.png) ](../../index.html)

## Getting Started

- \* [Installation](../../docs/installation.html)

- \* [Install with pip](../../docs/installation.html#install-with-pip)

- \* [Install with Conda](../../docs/installation.html#install-with-conda)

- \* [Install from Source](../../docs/installation.html#install-from-source)

- \* [Editable Install](../../docs/installation.html#editable-install)

- \* [Install PyTorch with CUDA support](../../docs/installation.html#install-pytorch-with-cuda-support)

- \* [Quickstart](../../docs/quickstart.html)

- \* [Sentence Transformer](../../docs/quickstart.html#sentence-transformer)

- \* [Cross Encoder](../../docs/quickstart.html#cross-encoder)

- \* [Next Steps](../../docs/quickstart.html#next-steps)

## Sentence Transformer

- \* [Usage](../../docs/sentence\_transformer/usage/usage.html)

- \* [Computing Embeddings](../computing-embeddings/README.html)

- \* [Initializing a Sentence Transformer Model](../computing-embeddings/README.html#initializing-a-sentence-transformer-model)

- \* [Calculating Embeddings](../computing-embeddings/README.html#calculating-embeddings)

- \* [Prompt Templates](../computing-embeddings/README.html#prompt-templates)

- \* [Input Sequence Length](../computing-embeddings/README.html#id1)

- \* [Multi-Process / Multi-GPU Encoding](../computing-embeddings/README.html#multi-process-multi-gpu-encoding)

[\\* \[Semantic Textual Similarity\]\(../../docs/sentence\\_transformer/usage/semantic\\_textual\\_similarity.html\)](#)

[\\* \[Similarity Calculation\]\(../../docs/sentence\\_transformer/usage/semantic\\_textual\\_similarity.html#similarity-calculation\)](#)

[\\* \[Semantic Search\]\(../semantic-search/README.html\)](#)

[\\* \[Background\]\(../semantic-search/README.html#background\)](#)

[\\* \[Symmetric vs. Asymmetric Semantic Search\]\(../semantic-search/README.html#symmetric-vs-asymmetric-semantic-search\)](#)

[\\* \[Manual Implementation\]\(../semantic-search/README.html#manual-implementation\)](#)

[\\* \[Optimized Implementation\]\(../semantic-search/README.html#optimized-implementation\)](#)

[\\* \[Speed Optimization\]\(../semantic-search/README.html#speed-optimization\)](#)

[\\* \[Elasticsearch\]\(../semantic-search/README.html#elasticsearch\)](#)

[\\* \[Approximate Nearest Neighbor\]\(../semantic-search/README.html#approximate-nearest-neighbor\)](#)

[\\* \[Retrieve & Re-Rank\]\(../semantic-search/README.html#retrieve-re-rank\)](#)

[\\* \[Examples\]\(../semantic-search/README.html#examples\)](#)

[\\* Retrieve & Re-Rank](#)

[\\* Retrieve & Re-Rank Pipeline](#)

[\\* Retrieval: Bi-Encoder](#)

[\\* Re-Ranker: Cross-Encoder](#)

[\\* Example Scripts](#)

[\\* Pre-trained Bi-Encoders \(Retrieval\)](#)

[\\* Pre-trained Cross-Encoders \(Re-Ranker\)](#)

[\\* \[Clustering\]\(../clustering/README.html\)](#)

[\\* \[k-Means\]\(../clustering/README.html#k-means\)](#)

[\\* \[Agglomerative Clustering\]\(../clustering/README.html#agglomerative-clustering\)](#)

\* [Fast Clustering](../clustering/README.html#fast-clustering)

\* [Topic Modeling](../clustering/README.html#topic-modeling)

\* [Paraphrase Mining](../paraphrase-mining/README.html)

\*

[`paraphrase\_mining()`](../paraphrase-mining/README.html#sentence\_transformers.util.paraphrase\_mining)

\* [Translated Sentence Mining](../parallel-sentence-mining/README.html)

\* [Margin Based Mining](../parallel-sentence-mining/README.html#margin-based-mining)

\* [Examples](../parallel-sentence-mining/README.html#examples)

\* [Image Search](../image-search/README.html)

\* [Installation](../image-search/README.html#installation)

\* [Usage](../image-search/README.html#usage)

\* [Examples](../image-search/README.html#examples)

\* [Embedding Quantization](../embedding-quantization/README.html)

\* [Binary Quantization](../embedding-quantization/README.html#binary-quantization)

\* [Scalar (int8) Quantization](../embedding-quantization/README.html#scalar-int8-quantization)

\* [Additional extensions](../embedding-quantization/README.html#additional-extensions)

\* [Demo](../embedding-quantization/README.html#demo)

\* [Try it yourself](../embedding-quantization/README.html#try-it-yourself)

\* [Speeding up Inference](../../docs/sentence\_transformer/usage/efficiency.html)

\* [PyTorch](../../docs/sentence\_transformer/usage/efficiency.html#pytorch)

\* [ONNX](../../docs/sentence\_transformer/usage/efficiency.html#onnx)

\* [OpenVINO](../../docs/sentence\_transformer/usage/efficiency.html#openvino)

\* [Benchmarks](../../docs/sentence\_transformer/usage/efficiency.html#benchmarks)

\* [Creating Custom Models](../../docs/sentence\_transformer/usage/custom\_models.html)

\* [Structure of Sentence Transformer

Models](../../docs/sentence\_transformer/usage/custom\_models.html#structure-of-sentence-transfo

mer-models)

\* [Sentence Transformer Model from a Transformers Model](../../../../docs/sentence\_transformer/usage/custom\_models.html#sentence-transformer-model-from-a-transformers-model)

\* [Pretrained Models](../../../../docs/sentence\_transformer/pretrained\_models.html)

\* [Original Models](../../../../docs/sentence\_transformer/pretrained\_models.html#original-models)

\* [Semantic Search Models](../../../../docs/sentence\_transformer/pretrained\_models.html#semantic-search-models)

\* [Multi-QA Models](../../../../docs/sentence\_transformer/pretrained\_models.html#multi-qa-models)

\* [MSMARCO Passage Models](../../../../docs/sentence\_transformer/pretrained\_models.html#msmarco-passage-models)

\* [Multilingual Models](../../../../docs/sentence\_transformer/pretrained\_models.html#multilingual-models)

\* [Semantic Similarity Models](../../../../docs/sentence\_transformer/pretrained\_models.html#semantic-similarity-models)

\* [Bitext Mining](../../../../docs/sentence\_transformer/pretrained\_models.html#bitext-mining)

\* [Image & Text-Models](../../../../docs/sentence\_transformer/pretrained\_models.html#image-text-models)

\* [INSTRUCTOR models](../../../../docs/sentence\_transformer/pretrained\_models.html#instructor-models)

\* [Scientific Similarity Models](../../../../docs/sentence\_transformer/pretrained\_models.html#scientific-similarity-models)

\* [Training Overview](../../../../docs/sentence\_transformer/training\_overview.html)

\* [Why Finetune?](../../../../docs/sentence\_transformer/training\_overview.html#why-finetune)

\* [Training Components](../../../../docs/sentence\_transformer/training\_overview.html#training-components)

\* [Dataset](../../../../docs/sentence\_transformer/training\_overview.html#dataset)

- \* [Dataset Format](../../../../docs/sentence\_transformer/training\_overview.html#dataset-format)
- \* [Loss Function](../../../../docs/sentence\_transformer/training\_overview.html#loss-function)
- \* [Training Arguments](../../../../docs/sentence\_transformer/training\_overview.html#training-arguments)
- \* [Evaluator](../../../../docs/sentence\_transformer/training\_overview.html#evaluator)
- \* [Trainer](../../../../docs/sentence\_transformer/training\_overview.html#trainer)
- \* [Callbacks](../../../../docs/sentence\_transformer/training\_overview.html#callbacks)
- \* [Multi-Dataset Training](../../../../docs/sentence\_transformer/training\_overview.html#multi-dataset-training)
- \* [Deprecated Training](../../../../docs/sentence\_transformer/training\_overview.html#deprecated-training)
- \* [Best Base Embedding Models](../../../../docs/sentence\_transformer/training\_overview.html#best-base-embedding-models)
- \* [Dataset Overview](../../../../docs/sentence\_transformer/dataset\_overview.html)
- \* [Datasets on the Hugging Face Hub](../../../../docs/sentence\_transformer/dataset\_overview.html#datasets-on-the-hugging-face-hub)
- \* [Pre-existing Datasets](../../../../docs/sentence\_transformer/dataset\_overview.html#pre-existing-datasets)
- \* [Loss Overview](../../../../docs/sentence\_transformer/loss\_overview.html)
- \* [Loss modifiers](../../../../docs/sentence\_transformer/loss\_overview.html#loss-modifiers)
- \* [Distillation](../../../../docs/sentence\_transformer/loss\_overview.html#distillation)
- \* [Commonly used Loss Functions](../../../../docs/sentence\_transformer/loss\_overview.html#commonly-used-loss-functions)
- \* [Custom Loss Functions](../../../../docs/sentence\_transformer/loss\_overview.html#custom-loss-functions)
- \* [Training Examples](../../../../docs/sentence\_transformer/training/examples.html)
- \* [Semantic Textual Similarity](../../../../training/sts/README.html)

- \* [Training data](../../training/sts/README.html#training-data)
- \* [Loss Function](../../training/sts/README.html#loss-function)
- \* [Natural Language Inference](../../training/nli/README.html)
- \* [Data](../../training/nli/README.html#data)
- \* [SoftmaxLoss](../../training/nli/README.html#softmaxloss)
- \* [MultipleNegativesRankingLoss](../../training/nli/README.html#multiplenegativesrankingloss)
- \* [Paraphrase Data](../../training/paraphrases/README.html)
- \* [Pre-Trained Models](../../training/paraphrases/README.html#pre-trained-models)
- \* [Quora Duplicate Questions](../../training/quora\_duplicate\_questions/README.html)
- \* [Training](../../training/quora\_duplicate\_questions/README.html#training)

\*

[MultipleNegativesRankingLoss](../../training/quora\_duplicate\_questions/README.html#multiplenegativesrankingloss)

- \* [Pretrained Models](../../training/quora\_duplicate\_questions/README.html#pretrained-models)
- \* [MS MARCO](../../training/ms\_marco/README.html)
- \* [Bi-Encoder](../../training/ms\_marco/README.html#bi-encoder)
- \* [Matryoshka Embeddings](../../training/matryoshka/README.html)
- \* [Use Cases](../../training/matryoshka/README.html#use-cases)
- \* [Results](../../training/matryoshka/README.html#results)
- \* [Training](../../training/matryoshka/README.html#training)
- \* [Inference](../../training/matryoshka/README.html#inference)
- \* [Code Examples](../../training/matryoshka/README.html#code-examples)
- \* [Adaptive Layers](../../training/adaptive\_layer/README.html)
- \* [Use Cases](../../training/adaptive\_layer/README.html#use-cases)
- \* [Results](../../training/adaptive\_layer/README.html#results)
- \* [Training](../../training/adaptive\_layer/README.html#training)
- \* [Inference](../../training/adaptive\_layer/README.html#inference)

- \* [Code Examples](../../training/adaptive\_layer/README.html#code-examples)
- \* [Multilingual Models](../../training/multilingual/README.html)
- \* [Extend your own models](../../training/multilingual/README.html#extend-your-own-models)
- \* [Training](../../training/multilingual/README.html#training)
- \* [Datasets](../../training/multilingual/README.html#datasets)
- \* [Sources for Training Data](../../training/multilingual/README.html#sources-for-training-data)
- \* [Evaluation](../../training/multilingual/README.html#evaluation)

\* [Available Pre-trained Models](../../training/multilingual/README.html#available-pre-trained-models)

- \* [Usage](../../training/multilingual/README.html#usage)
- \* [Performance](../../training/multilingual/README.html#performance)
- \* [Citation](../../training/multilingual/README.html#citation)
- \* [Model Distillation](../../training/distillation/README.html)
- \* [Knowledge Distillation](../../training/distillation/README.html#knowledge-distillation)

\* [Speed - Performance Trade-Off](../../training/distillation/README.html#speed-performance-trade-off)

- \* [Dimensionality Reduction](../../training/distillation/README.html#dimensionality-reduction)
- \* [Quantization](../../training/distillation/README.html#quantization)
- \* [Augmented SBERT](../../training/data\_augmentation/README.html)
- \* [Motivation](../../training/data\_augmentation/README.html#motivation)

\* [Extend to your own datasets](../../training/data\_augmentation/README.html#extend-to-your-own-datasets)

- \* [Methodology](../../training/data\_augmentation/README.html#methodology)

\* [Scenario 1: Limited or small annotated datasets (few labeled sentence-pairs)](../../training/data\_augmentation/README.html#scenario-1-limited-or-small-annotated-datasets-few-labeled-sentence-pairs)

\* [Scenario 2: No annotated datasets (Only unlabeled

sentence-pairs)](../../training/data\_augmentation/README.html#scenario-2-no-annotated-datasets-only-unlabeled-sentence-pairs)

- \* [Training](../../training/data\_augmentation/README.html#training)
- \* [Citation](../../training/data\_augmentation/README.html#citation)
- \* [Training with Prompts](../../training/prompts/README.html)
- \* [What are Prompts?](../../training/prompts/README.html#what-are-prompts)
  - \* [Why would we train with Prompts?](../../training/prompts/README.html#why-would-we-train-with-prompts)
  - \* [How do we train with Prompts?](../../training/prompts/README.html#how-do-we-train-with-prompts)
- \* [Training with PEFT Adapters](../../training/peft/README.html)
- \* [Compatibility Methods](../../training/peft/README.html#compatibility-methods)
- \* [Adding a New Adapter](../../training/peft/README.html#adding-a-new-adapter)
- \* [Loading a Pretrained Adapter](../../training/peft/README.html#loading-a-pretrained-adapter)
- \* [Training Script](../../training/peft/README.html#training-script)
- \* [Unsupervised Learning](../../unsupervised\_learning/README.html)
- \* [TSDAE](../../unsupervised\_learning/README.html#tsdae)
- \* [SimCSE](../../unsupervised\_learning/README.html#simcse)
- \* [CT](../../unsupervised\_learning/README.html#ct)
  - \* [CT (In-Batch Negative Sampling)](../../unsupervised\_learning/README.html#ct-in-batch-negative-sampling)
  - \* [Masked Language Model (MLM)](../../unsupervised\_learning/README.html#masked-language-model-mlm)
- \* [GenQ](../../unsupervised\_learning/README.html#genq)
- \* [GPL](../../unsupervised\_learning/README.html#gpl)
  - \* [Performance Comparison](../../unsupervised\_learning/README.html#performance-comparison)



- \* [Domain Adaptation](../../domain\_adaptation/README.html)
  - \* [Domain Adaptation vs. Unsupervised Learning](../../domain\_adaptation/README.html#domain-adaptation-vs-unsupervised-learning)
  - \* [Adaptive Pre-Training](../../domain\_adaptation/README.html#adaptive-pre-training)
    - \* [GPL: Generative Pseudo-Labeling](../../domain\_adaptation/README.html#gpl-generative-pseudo-labeling)
- \* [Hyperparameter Optimization](../../training/hpo/README.html)
  - \* [HPO Components](../../training/hpo/README.html#hpo-components)
  - \* [Putting It All Together](../../training/hpo/README.html#putting-it-all-together)
  - \* [Example Scripts](../../training/hpo/README.html#example-scripts)
- \* [Distributed Training](../../docs/sentence\_transformer/training/distributed.html)
  - \* [Comparison](../../docs/sentence\_transformer/training/distributed.html#comparison)
  - \* [FSDP](../../docs/sentence\_transformer/training/distributed.html#fsdp)

## Cross Encoder

- \* [Usage](../../docs/cross\_encoder/usage/usage.html)
- \* Retrieve & Re-Rank
  - \* Retrieve & Re-Rank Pipeline
  - \* Retrieval: Bi-Encoder
  - \* Re-Ranker: Cross-Encoder
  - \* Example Scripts
  - \* Pre-trained Bi-Encoders (Retrieval)
  - \* Pre-trained Cross-Encoders (Re-Ranker)
- \* [Pretrained Models](../../docs/cross\_encoder/pretrained\_models.html)
  - \* [MS MARCO](../../docs/cross\_encoder/pretrained\_models.html#ms-marco)
  - \* [SQuAD (QNLI)](../../docs/cross\_encoder/pretrained\_models.html#squad-qnli)

\* [STSbenchmark](../../docs/cross\_encoder/pretrained\_models.html#stsbenchmark)

\* [Quora Duplicate

Questions](../../docs/cross\_encoder/pretrained\_models.html#quora-duplicate-questions)

\* [NLI](../../docs/cross\_encoder/pretrained\_models.html#nli)

\* [Community Models](../../docs/cross\_encoder/pretrained\_models.html#community-models)

\* [Training Overview](../../docs/cross\_encoder/training\_overview.html)

\* [Training Examples](../../docs/cross\_encoder/training/examples.html)

\* [MS MARCO](../../training/ms\_marco/cross\_encoder\_README.html)

\* [Cross-Encoder](../../training/ms\_marco/cross\_encoder\_README.html#cross-encoder)

\* [Cross-Encoder Knowledge

Distillation](../../training/ms\_marco/cross\_encoder\_README.html#cross-encoder-knowledge-distillation)

## Package Reference

\* [Sentence Transformer](../../docs/package\_reference/sentence\_transformer/index.html)

\*

[SentenceTransformer](../../docs/package\_reference/sentence\_transformer/SentenceTransformer.html)

\*

[SentenceTransformer](../../docs/package\_reference/sentence\_transformer/SentenceTransformer.html#id1)

\*

[SentenceTransformerModelCardData](../../docs/package\_reference/sentence\_transformer/SentenceTransformer.html#sentencetransformermodelcarddata)

\*

[SimilarityFunction](../../docs/package\_reference/sentence\_transformer/SentenceTransformer.html

#similarityfunction)

\* [Trainer](../../docs/package\_reference/sentence\_transformer/trainer.html)

\*

[SentenceTransformerTrainer](../../docs/package\_reference/sentence\_transformer/trainer.html#sentencetransformertrainer)

\* [Training Arguments](../../docs/package\_reference/sentence\_transformer/training\_args.html)

\*

[SentenceTransformerTrainingArguments](../../docs/package\_reference/sentence\_transformer/training\_args.html#sentencetransformertrainingarguments)

\* [Losses](../../docs/package\_reference/sentence\_transformer/losses.html)

\*

[BatchAllTripletLoss](../../docs/package\_reference/sentence\_transformer/losses.html#batchalltripletloss)

\*

[BatchHardSoftMarginTripletLoss](../../docs/package\_reference/sentence\_transformer/losses.html#batchhardsoftmargintripletloss)

\*

[BatchHardTripletLoss](../../docs/package\_reference/sentence\_transformer/losses.html#batchhardtripletloss)

\*

[BatchSemiHardTripletLoss](../../docs/package\_reference/sentence\_transformer/losses.html#batchsemihardtripletloss)

\*

[ContrastiveLoss](../../docs/package\_reference/sentence\_transformer/losses.html#contrastiveloss)

\*

[OnlineContrastiveLoss](../../docs/package\_reference/sentence\_transformer/losses.html#onlinecontrastiveloss)

\*

[ContrastiveTensionLoss](../../docs/package\_reference/sentence\_transformer/losses.html#contrastivetensionloss)

\*

[ContrastiveTensionLossInBatchNegatives](../../docs/package\_reference/sentence\_transformer/losses.html#contrastivetensionlossinbatchnegatives)

\* [CoSENTLoss](../../docs/package\_reference/sentence\_transformer/losses.html#cosentloss)

\* [AngleLoss](../../docs/package\_reference/sentence\_transformer/losses.html#angleloss)

\*

[CosineSimilarityLoss](../../docs/package\_reference/sentence\_transformer/losses.html#cosinesimilarityloss)

\*

[DenoisingAutoEncoderLoss](../../docs/package\_reference/sentence\_transformer/losses.html#denoisingautoencoderloss)

\*

[GISTEmbedLoss](../../docs/package\_reference/sentence\_transformer/losses.html#gistembedloss)

\*

[CachedGISTEmbedLoss](../../docs/package\_reference/sentence\_transformer/losses.html#cachedgistembedloss)

\* [MSELoss](../../docs/package\_reference/sentence\_transformer/losses.html#mseloss)

\*

[MarginMSELoss](../../docs/package\_reference/sentence\_transformer/losses.html#marginmseloss)

\*

[MatryoshkaLoss](../../docs/package\_reference/sentence\_transformer/losses.html#matryoshkaloss)

\*

[Matryoshka2dLoss](../../docs/package\_reference/sentence\_transformer/losses.html#matryoshka2dloss)

\*

[AdaptiveLayerLoss](../../docs/package\_reference/sentence\_transformer/losses.html#adaptivelayerloss)

\*

[MegaBatchMarginLoss](../../docs/package\_reference/sentence\_transformer/losses.html#megabatchmarginloss)

\*

[MultipleNegativesRankingLoss](../../docs/package\_reference/sentence\_transformer/losses.html#multiplenegativesrankingloss)

\*

[CachedMultipleNegativesRankingLoss](../../docs/package\_reference/sentence\_transformer/losses.html#cachedmultiplenegativesrankingloss)

\*

[MultipleNegativesSymmetricRankingLoss](../../docs/package\_reference/sentence\_transformer/losses.html#multiplenegativessymmetricrankingloss)

\*

[CachedMultipleNegativesSymmetricRankingLoss](../../docs/package\_reference/sentence\_transformer/losses.html#cachedmultiplenegativessymmetricrankingloss)

\* [SoftmaxLoss](../../docs/package\_reference/sentence\_transformer/losses.html#softmaxloss)

\* [TripletLoss](../../docs/package\_reference/sentence\_transformer/losses.html#tripletloss)

\* [Samplers](../../docs/package\_reference/sentence\_transformer/sampler.html)

\*

[BatchSamplers](../../docs/package\_reference/sentence\_transformer/sampler.html#batchsamplers)

\*

[MultiDatasetBatchSamplers](../../docs/package\_reference/sentence\_transformer/sampler.html#multidatasetbatchsamplers)

\* [Evaluation](../../docs/package\_reference/sentence\_transformer/evaluation.html)

\*

[BinaryClassificationEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#binaryclassificationevaluator)

\*

[EmbeddingSimilarityEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#embeddingsimilarityevaluator)

\*

[InformationRetrievalEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#informationretrievalevaluator)

\*

[NanoBEIREvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#nanoberevaluator)

\*

[MSEEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#mseevaluator)

\*

[ParaphraseMiningEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#paraphraseminingevaluator)

\*

[RerankingEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#rerankingevaluator)

\*

[SentenceEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#sentenceevaluator)

eevaluator)

\*

[SequentialEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#sequentialevaluator)

\*

[TranslationEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#translationevaluator)

\*

[TripletEvaluator](../../docs/package\_reference/sentence\_transformer/evaluation.html#tripletevaluator)

\* [Datasets](../../docs/package\_reference/sentence\_transformer/datasets.html)

\*

[ParallelSentencesDataset](../../docs/package\_reference/sentence\_transformer/datasets.html#parallelsentencesdataset)

\*

[SentenceLabelDataset](../../docs/package\_reference/sentence\_transformer/datasets.html#sentencelabeldataset)

\*

[DenoisingAutoEncoderDataset](../../docs/package\_reference/sentence\_transformer/datasets.html#denoisingautoencoderdataset)

\*

[NoDuplicatesDataLoader](../../docs/package\_reference/sentence\_transformer/datasets.html#noduplicatesdataloader)

\* [Models](../../docs/package\_reference/sentence\_transformer/models.html)

\*

[Main

Classes](../../docs/package\_reference/sentence\_transformer/models.html#main-classes)

\*

[Further

Classes](../../docs/package\_reference/sentence\_transformer/models.html#further-classes)

\* [quantization](../../docs/package\_reference/sentence\_transformer/quantization.html)

\*

[`quantize\_embeddings()`](../../docs/package\_reference/sentence\_transformer/quantization.html#sentence\_transformers.quantization.quantize\_embeddings)

\*

[`semantic\_search\_faiss()`](../../docs/package\_reference/sentence\_transformer/quantization.html#sentence\_transformers.quantization.semantic\_search\_faiss)

\*

[`semantic\_search\_usearch()`](../../docs/package\_reference/sentence\_transformer/quantization.html#sentence\_transformers.quantization.semantic\_search\_usearch)

\* [Cross Encoder](../../docs/package\_reference/cross\_encoder/index.html)

\* [CrossEncoder](../../docs/package\_reference/cross\_encoder/cross\_encoder.html)

\* [CrossEncoder](../../docs/package\_reference/cross\_encoder/cross\_encoder.html#id1)

\*

[Training

Inputs](../../docs/package\_reference/cross\_encoder/cross\_encoder.html#training-inputs)

\* [Evaluation](../../docs/package\_reference/cross\_encoder/evaluation.html)

\*

[CEBinaryAccuracyEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cebinaryaccuracyevaluator)

\*

[CEBinaryClassificationEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cebinaryclassificationevaluator)

\*

[CECorrelationEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cecorrelationevaluator)

\* [CEF1Evaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cef1evaluator)



\*

[CESoftmaxAccuracyEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cesoftmaxaccuracyevaluator)

\*

[CERerankingEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cererankingevaluator)

\* [util](../../docs/package\_reference/util.html)

\* [Helper Functions](../../docs/package\_reference/util.html#module-sentence\_transformers.util)

\*

[`community\_detection()`](../../docs/package\_reference/util.html#sentence\_transformers.util.community\_detection)

\* [`http\_get()`](../../docs/package\_reference/util.html#sentence\_transformers.util.http\_get)

\*

[`is\_training\_available()`](../../docs/package\_reference/util.html#sentence\_transformers.util.is\_training\_available)

\*

[`mine\_hard\_negatives()`](../../docs/package\_reference/util.html#sentence\_transformers.util.mine\_hard\_negatives)

\*

[`normalize\_embeddings()`](../../docs/package\_reference/util.html#sentence\_transformers.util.normalize\_embeddings)

\*

[`paraphrase\_mining()`](../../docs/package\_reference/util.html#sentence\_transformers.util.paraphrase\_mining)

\*

[`semantic\_search()`](../../docs/package\_reference/util.html#sentence\_transformers.util.semantic\_search)

\*

[`truncate\_embeddings()`](../../docs/package\_reference/util.html#sentence\_transformers.util.truncate\_embeddings)

\*

[Model

Optimization](../../docs/package\_reference/util.html#module-sentence\_transformers.backend)

\*

[`export\_dynamic\_quantized\_onnx\_model()`](../../docs/package\_reference/util.html#sentence\_transformers.backend.export\_dynamic\_quantized\_onnx\_model)

\*

[`export\_optimized\_onnx\_model()`](../../docs/package\_reference/util.html#sentence\_transformers.backend.export\_optimized\_onnx\_model)

\*

[`export\_static\_quantized\_openvino\_model()`](../../docs/package\_reference/util.html#sentence\_transformers.backend.export\_static\_quantized\_openvino\_model)

\* [Similarity Metrics](../../docs/package\_reference/util.html#module-sentence\_transformers.util)

\* [`cos\_sim()`](../../docs/package\_reference/util.html#sentence\_transformers.util.cos\_sim)

\* [`dot\_score()`](../../docs/package\_reference/util.html#sentence\_transformers.util.dot\_score)

\*

[`euclidean\_sim()`](../../docs/package\_reference/util.html#sentence\_transformers.util.euclidean\_sim)

\*

[`manhattan\_sim()`](../../docs/package\_reference/util.html#sentence\_transformers.util.manhattan\_sim)

\*

[`pairwise\_cos\_sim()`](../../docs/package\_reference/util.html#sentence\_transformers.util.pairwise\_cos\_sim)

\*

[`pairwise\_dot\_score()`](../../docs/package\_reference/util.html#sentence\_transformers.util.pairwise\_dot\_score)

\*

[`pairwise\_euclidean\_sim()`](../../docs/package\_reference/util.html#sentence\_transformers.util.pairwise\_euclidean\_sim)

\*

[`pairwise\_manhattan\_sim()`](../../docs/package\_reference/util.html#sentence\_transformers.util.pairwise\_manhattan\_sim)

\_\_[Sentence Transformers](../../index.html)

\* [(../../index.html)]

\* [Usage](../../docs/sentence\_transformer/usage/usage.html)

\* Retrieve & Re-Rank

\*

[

Edit

on

GitHub](https://github.com/UKPLab/sentence-transformers/blob/master/examples/applications/retrieve\_rerank/README.md)

\* \* \*

# Retrieve & Re-Rank

In [Semantic Search](../semantic-search/README.html) we have shown how to use SentenceTransformer to compute embeddings for queries, sentences, and paragraphs and how to use this for semantic search.

For complex search tasks, for example question answering retrieval, the search

can significantly be improved by using **Retrieve & Re-Rank**.

## ## Retrieve & Re-Rank Pipeline

The following pipeline for Information Retrieval / Question Answering

Retrieval works very well. All components are provided and explained in this article:

![[InformationRetrieval]](<https://raw.githubusercontent.com/UKPLab/sentence-transformers/master/docs/img/InformationRetrieval.png>)

Given a search query, we first use a **retrieval system** that retrieves a large list of e.g. 100 possible hits which are potentially relevant for the query. For the retrieval, we can use either lexical search, e.g. with a vector engine like Elasticsearch, or we can use dense retrieval with a bi-encoder. However, the retrieval system might retrieve documents that are not that relevant for the search query. Hence, in a second stage, we use a **re-ranker** based on a **cross-encoder** that scores the relevancy of all candidates for the given search query. The output will be a ranked list of hits we can present to the user.

## ## Retrieval: Bi-Encoder

For the retrieval of the candidate set, we can either use lexical search (e.g. [Elasticsearch](https://www.elastic.co/elasticsearch/)), or we can use a bi-encoder which is implemented in Sentence Transformers.

Lexical search looks for literal matches of the query words in your document collection. It will not recognize synonyms, acronyms or spelling variations. In contrast, semantic search (or dense retrieval) encodes the search query into vector space and retrieves the document embeddings that are close in vector space.

![[SemanticSearch]](<https://raw.githubusercontent.com/UKPLab/sentence-transformers/master/docs/img/SemanticSearch.png>)

Semantic search overcomes the shortcomings of lexical search and can recognize synonym and acronyms. Have a look at the [semantic search article](../semantic-search/README.html) for different options to implement semantic search.

## ## Re-Ranker: Cross-Encoder

The retriever has to be efficient for large document collections with millions of entries. However, it might return irrelevant candidates. A re-ranker based on a Cross-Encoder can substantially improve the final results for the user. The query and a possible document is passed simultaneously to transformer network, which then outputs a single score between 0 and 1 indicating how relevant the document is for the given query.

![[CrossEncoder]](<https://raw.githubusercontent.com/UKPLab/sentence-transformers/master/docs/img/CrossEncoder.png>)

The advantage of Cross-Encoders is the higher performance, as they perform

attention across the query and the document. Scoring thousands or millions of (query, document)-pairs would be rather slow. Hence, we use the retriever to create a set of e.g. 100 possible candidates which are then re-ranked by the Cross-Encoder.

## ## Example Scripts

**[retrieve\_rerank\_simple\_wikipedia.ipynb]**([https://github.com/UKPLab/sentence-transformers/tree/master/examples/applications/retrieve\\_rerank/retrieve\\_rerank\\_simple\\_wikipedia.ipynb](https://github.com/UKPLab/sentence-transformers/tree/master/examples/applications/retrieve_rerank/retrieve_rerank_simple_wikipedia.ipynb)) **[Colab Version]**([https://colab.research.google.com/github/UKPLab/sentence-transformers/blob/master/examples/applications/retrieve\\_rerank/retrieve\\_rerank\\_simple\\_wikipedia.ipynb](https://colab.research.google.com/github/UKPLab/sentence-transformers/blob/master/examples/applications/retrieve_rerank/retrieve_rerank_simple_wikipedia.ipynb)): This script uses the smaller **[Simple English Wikipedia]**([https://simple.wikipedia.org/wiki/Main\\_Page](https://simple.wikipedia.org/wiki/Main_Page)) as document collection to provide answers to user questions / search queries. First, we split all Wikipedia articles into paragraphs and encode them with a bi-encoder. If a new query / question is entered, it is encoded by the same bi-encoder and the paragraphs with the highest cosine-similarity are retrieved (see **[semantic search]**([../semantic-search/README.html](https://github.com/UKPLab/sentence-transformers/blob/master/examples/applications/semantic_search/README.html))). Next, the retrieved candidates are scored by a Cross-Encoder re-ranker and the 5 passages with the highest score from the Cross-Encoder are presented to the user.

**[in\_document\_search\_crossencoder.py]**([https://github.com/UKPLab/sentence-transformers/tree/master/examples/applications/retrieve\\_rerank/in\\_document\\_search\\_crossencoder.py](https://github.com/UKPLab/sentence-transformers/tree/master/examples/applications/retrieve_rerank/in_document_search_crossencoder.py)): **If you only have a small set of paragraphs, we don't do the retrieval stage.** This is for example the case if you want to perform search within a single document. In this example, we take the Wikipedia article about Europe and split it into paragraphs. Then, the search query / question and all paragraphs are scored using the Cross-Encoder re-ranker. The most relevant passages for the query are returned.

## ## Pre-trained Bi-Encoders (Retrieval)if•

The bi-encoder produces embeddings independently for your paragraphs and for your search queries. You can use it like this:

```
from sentence_transformers import SentenceTransformer

model = SentenceTransformer("multi-qa-mpnet-base-dot-v1")

docs = [
    "My first paragraph. That contains information",
    "Python is a programming language.",
]

document_embeddings = model.encode(docs)

query = "What is Python?"

query_embedding = model.encode(query)
```

For more details how to compare the embeddings, see [semantic search](../semantic-search/README.html).

We provide pre-trained models based on:

\* \*\*MS MARCO:\*\* 500k real user queries from Bing search engine. See [MS MARCO models](../../docs/pretrained-models/msmarco-v3.html)

## Pre-trained Cross-Encoders (Re-Ranker)if•

For pre-trained Cross Encoder models, see: [MS MARCO Cross-Encoders](../../docs/pretrained-models/ce-msmarco.html)

[ Previous](../semantic-search/README.html "Semantic Search") [Next](../clustering/README.html "Clustering")

\* \* \*

(C) Copyright 2025.

Built with [Sphinx](https://www.sphinx-doc.org/) using a [theme](https://github.com/readthedocs/sphinx\_rtd\_theme) provided by [Read the Docs](https://readthedocs.org).