

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

## Getting Started

- \* [Installation](../installation.html)
- \* [Install with pip](../installation.html#install-with-pip)
- \* [Install with Conda](../installation.html#install-with-conda)
- \* [Install from Source](../installation.html#install-from-source)
- \* [Editable Install](../installation.html#editable-install)
- \* [Install PyTorch with CUDA support](../installation.html#install-pytorch-with-cuda-support)
- \* [Quickstart](../quickstart.html)
- \* [Sentence Transformer](../quickstart.html#sentence-transformer)
- \* [Cross Encoder](../quickstart.html#cross-encoder)
- \* [Next Steps](../quickstart.html#next-steps)

## Sentence Transformer

- \* [Usage](../sentence\_transformer/usage/usage.html)
- \* [Computing Embeddings](../examples/applications/computing-embeddings/README.html)
  - \* [Initializing a Sentence Transformer Model](../examples/applications/computing-embeddings/README.html#initializing-a-sentence-transformer-model)
  - \* [Calculating Embeddings](../examples/applications/computing-embeddings/README.html#calculating-embeddings)
  - \* [Prompt Templates](../examples/applications/computing-embeddings/README.html#prompt-templates)

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

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

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

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

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

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

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

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

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

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

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

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

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

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

\* [Retrieve & Re-Rank](../../examples/applications/retrieve\_rerank/README.html)

\* [Retrieve & Re-Rank

Pipeline](../../examples/applications/retrieve\_rerank/README.html#retrieve-re-rank-pipeline)

\* [Retrieval:

Bi-Encoder](../../examples/applications/retrieve\_rerank/README.html#retrieval-bi-encoder)

\* [Re-Ranker:

Cross-Encoder](../../examples/applications/retrieve\_rerank/README.html#re-ranker-cross-encoder)

\* [Example Scripts](../../examples/applications/retrieve\_rerank/README.html#example-scripts)

\* [Pre-trained Bi-Encoders

(Retrieval)](../../examples/applications/retrieve\_rerank/README.html#pre-trained-bi-encoders-retrieval)

\* [Pre-trained Cross-Encoders

(Re-Ranker)](../../examples/applications/retrieve\_rerank/README.html#pre-trained-cross-encoders-re-ranker)

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

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

\* [Agglomerative

Clustering](../../examples/applications/clustering/README.html#agglomerative-clustering)

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

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

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

\*

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

\* [Translated Sentence

Mining](../../examples/applications/parallel-sentence-mining/README.html)

Mining](../../examples/applications/parallel-sentence-mining/README.html#margin-based-mining)

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

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

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

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

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

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

\* [Binary

Quantization](../../examples/applications/embedding-quantization/README.html#binary-quantization)

\* [Scalar (int8)

Quantization](../../examples/applications/embedding-quantization/README.html#scalar-int8-quantization)

\* [Additional

extensions](../../examples/applications/embedding-quantization/README.html#additional-extensions)

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

\* [Try it

yourself](../../examples/applications/embedding-quantization/README.html#try-it-yourself)

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

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

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

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

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

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

Models](../sentence\_transformer/usage/custom\_models.html#structure-of-sentence-transformer-models)

- \* [Sentence Transformer Model from a Transformers Model](../sentence\_transformer/usage/custom\_models.html#sentence-transformer-model-from-a-transformers-model)
- \* [Pretrained Models](../sentence\_transformer/pretrained\_models.html)
- \* [Original Models](../sentence\_transformer/pretrained\_models.html#original-models)
  - \* [Semantic Search Models](../sentence\_transformer/pretrained\_models.html#semantic-search-models)
  - \* [Multi-QA Models](../sentence\_transformer/pretrained\_models.html#multi-qa-models)
    - \* [MSMARCO Passage Models](../sentence\_transformer/pretrained\_models.html#msmarco-passage-models)
  - \* [Multilingual Models](../sentence\_transformer/pretrained\_models.html#multilingual-models)
    - \* [Semantic Similarity Models](../sentence\_transformer/pretrained\_models.html#semantic-similarity-models)
  - \* [Bitext Mining](../sentence\_transformer/pretrained\_models.html#bitext-mining)
  - \* [Image & Text-Models](../sentence\_transformer/pretrained\_models.html#image-text-models)
  - \* [INSTRUCTOR models](../sentence\_transformer/pretrained\_models.html#instructor-models)
    - \* [Scientific Similarity Models](../sentence\_transformer/pretrained\_models.html#scientific-similarity-models)
- \* [Training Overview](../sentence\_transformer/training\_overview.html)
  - \* [Why Finetune?](../sentence\_transformer/training\_overview.html#why-finetune)
  - \* [Training Components](../sentence\_transformer/training\_overview.html#training-components)
  - \* [Dataset](../sentence\_transformer/training\_overview.html#dataset)
    - \* [Dataset Format](../sentence\_transformer/training\_overview.html#dataset-format)
  - \* [Loss Function](../sentence\_transformer/training\_overview.html#loss-function)
  - \* [Training Arguments](../sentence\_transformer/training\_overview.html#training-arguments)

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

- \* [Paraphrase Data](../../examples/training/paraphrases/README.html)
- \* [Pre-Trained Models](../../examples/training/paraphrases/README.html#pre-trained-models)
- \* [Quora Duplicate Questions](../../examples/training/quora\_duplicate\_questions/README.html)
- \* [Training](../../examples/training/quora\_duplicate\_questions/README.html#training)

\*

[MultipleNegativesRankingLoss](../../examples/training/quora\_duplicate\_questions/README.html#multiple-negatives-ranking-loss)

\*

[Pretrained

Models](../../examples/training/quora\_duplicate\_questions/README.html#pretrained-models)

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

\*

[Extend your own

models](../../examples/training/multilingual/README.html#extend-your-own-models)

- \* [Training](../../examples/training/multilingual/README.html#training)

- \* [Datasets](../../examples/training/multilingual/README.html#datasets)
- \* [Sources for Training Data](../../examples/training/multilingual/README.html#sources-for-training-data)
- \* [Evaluation](../../examples/training/multilingual/README.html#evaluation)
- \* [Available Pre-trained Models](../../examples/training/multilingual/README.html#available-pre-trained-models)
- \* [Usage](../../examples/training/multilingual/README.html#usage)
- \* [Performance](../../examples/training/multilingual/README.html#performance)
- \* [Citation](../../examples/training/multilingual/README.html#citation)
- \* [Model Distillation](../../examples/training/distillation/README.html)
- \* [Knowledge Distillation](../../examples/training/distillation/README.html#knowledge-distillation)
- \* [Speed - Performance Trade-Off](../../examples/training/distillation/README.html#speed-performance-trade-off)
- \* [Dimensionality Reduction](../../examples/training/distillation/README.html#dimensionality-reduction)
- \* [Quantization](../../examples/training/distillation/README.html#quantization)
- \* [Augmented SBERT](../../examples/training/data\_augmentation/README.html)
- \* [Motivation](../../examples/training/data\_augmentation/README.html#motivation)
- \* [Extend to your own datasets](../../examples/training/data\_augmentation/README.html#extend-to-your-own-datasets)
- \* [Methodology](../../examples/training/data\_augmentation/README.html#methodology)
  - \* [Scenario 1: Limited or small annotated datasets (few labeled sentence-pairs)](../../examples/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)](../../examples/training/data\_augmentation/README.html#scenario-2-no-annotated-datasets-only-unlabeled-sentence-pairs)



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

\* [Domain Adaptation vs. Unsupervised Learning](../../examples/domain\_adaptation/README.html#domain-adaptation-vs-unsupervised-learning)

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

## Cross Encoder

- \* [Usage](../cross\_encoder/usage/usage.html)
- \* [Retrieve & Re-Rank Pipeline](../../examples/applications/retrieve\_rerank/README.html#retrieve-re-rank-pipeline)
- \* [Retrieval: Bi-Encoder](../../examples/applications/retrieve\_rerank/README.html#retrieval-bi-encoder)
- \* [Re-Ranker: Cross-Encoder](../../examples/applications/retrieve\_rerank/README.html#re-ranker-cross-encoder)
- \* [Example Scripts](../../examples/applications/retrieve\_rerank/README.html#example-scripts)
- \* [Pre-trained Bi-Encoders (Retrieval)](../../examples/applications/retrieve\_rerank/README.html#pre-trained-bi-encoders-retrieval)

val)

\* [Pre-trained Cross-Encoders

(Re-Ranker)](../../examples/applications/retrieve\_rerank/README.html#pre-trained-cross-encoders-re-ranker)

\* [Pretrained Models](../cross\_encoder/pretrained\_models.html)

\* [MS MARCO](../cross\_encoder/pretrained\_models.html#ms-marco)

\* [SQuAD (QNLI)](../cross\_encoder/pretrained\_models.html#squad-qnli)

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

\* [Quora Duplicate

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

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

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

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

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

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

\*

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

\* [Cross-Encoder Knowledge

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

## Package Reference

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

\* [SentenceTransformer](sentence\_transformer/SentenceTransformer.html)

\* [SentenceTransformer](sentence\_transformer/SentenceTransformer.html#id1)

\*

[SentenceTransformerModelCardData](sentence\_transformer/SentenceTransformer.html#sentence-transformer-model-card-data)

- \* [SimilarityFunction](sentence\_transformer/SentenceTransformer.html#similarity-function)
- \* [Trainer](sentence\_transformer/trainer.html)
- \* [SentenceTransformerTrainer](sentence\_transformer/trainer.html#sentence-transformer-trainer)
- \* [Training Arguments](sentence\_transformer/training\_args.html)

\*

[SentenceTransformerTrainingArguments](sentence\_transformer/training\_args.html#sentence-transformer-training-arguments)

- \* [Losses](sentence\_transformer/losses.html)
- \* [BatchAllTripletLoss](sentence\_transformer/losses.html#batch-all-triplet-loss)

\*

[BatchHardSoftMarginTripletLoss](sentence\_transformer/losses.html#batch-hard-soft-margin-triplet-loss)

- \* [BatchHardTripletLoss](sentence\_transformer/losses.html#batch-hard-triplet-loss)
- \* [BatchSemiHardTripletLoss](sentence\_transformer/losses.html#batch-semi-hard-triplet-loss)
- \* [ContrastiveLoss](sentence\_transformer/losses.html#contrastive-loss)
- \* [OnlineContrastiveLoss](sentence\_transformer/losses.html#online-contrastive-loss)
- \* [ContrastiveTensionLoss](sentence\_transformer/losses.html#contrastive-tension-loss)

\*

[ContrastiveTensionLossInBatchNegatives](sentence\_transformer/losses.html#contrastive-tension-loss-in-batch-negatives)

- \* [CoSENTLoss](sentence\_transformer/losses.html#co-sent-loss)
- \* [AngleELoss](sentence\_transformer/losses.html#angle-e-loss)
- \* [CosineSimilarityLoss](sentence\_transformer/losses.html#cosine-similarity-loss)
- \* [DenoisingAutoEncoderLoss](sentence\_transformer/losses.html#denoising-auto-encoder-loss)
- \* [GISTEmbedLoss](sentence\_transformer/losses.html#gist-embed-loss)

- \* [CachedGISTEmbedLoss](sentence\_transformer/losses.html#cachedgistembedloss)
- \* [MSELoss](sentence\_transformer/losses.html#mseloss)
- \* [MarginMSELoss](sentence\_transformer/losses.html#marginmseloss)
- \* [MatryoshkaLoss](sentence\_transformer/losses.html#matryoshkaloss)
- \* [Matryoshka2dLoss](sentence\_transformer/losses.html#matryoshka2dloss)
- \* [AdaptiveLayerLoss](sentence\_transformer/losses.html#adaptivelayerloss)
- \* [MegaBatchMarginLoss](sentence\_transformer/losses.html#megabatchmarginloss)

\*

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

\*

[CachedMultipleNegativesRankingLoss](sentence\_transformer/losses.html#cachedmultiplenegative  
srankingloss)

\*

[MultipleNegativesSymmetricRankingLoss](sentence\_transformer/losses.html#multiplenegativessym  
metricrankingloss)

\*

[CachedMultipleNegativesSymmetricRankingLoss](sentence\_transformer/losses.html#cachedmultipl  
enegativessymmetricrankingloss)

- \* [SoftmaxLoss](sentence\_transformer/losses.html#softmaxloss)

- \* [TripletLoss](sentence\_transformer/losses.html#tripletloss)

- \* [Samplers](sentence\_transformer/sampler.html)

- \* [BatchSamplers](sentence\_transformer/sampler.html#batchsamplers)

- \* [MultiDatasetBatchSamplers](sentence\_transformer/sampler.html#multidatasetbatchsamplers)

- \* [Evaluation](sentence\_transformer/evaluation.html)

\*

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

\*

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

\*

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

\* [NanoBEIREvaluator](sentence\_transformer/evaluation.html#nanobeirevaluator)

\* [MSEEvaluator](sentence\_transformer/evaluation.html#mseevaluator)

\*

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

\* [RerankingEvaluator](sentence\_transformer/evaluation.html#rerankingevaluator)

\* [SentenceEvaluator](sentence\_transformer/evaluation.html#sentenceevaluator)

\* [SequentialEvaluator](sentence\_transformer/evaluation.html#sequentialevaluator)

\* [TranslationEvaluator](sentence\_transformer/evaluation.html#translationevaluator)

\* [TripletEvaluator](sentence\_transformer/evaluation.html#tripletevaluator)

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

\* [ParallelSentencesDataset](sentence\_transformer/datasets.html#parallelsentencesdataset)

\* [SentenceLabelDataset](sentence\_transformer/datasets.html#sentencelabeldataset)

\*

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

\* [NoDuplicatesDataLoader](sentence\_transformer/datasets.html#noduplicatesdataloader)

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

\* [Main Classes](sentence\_transformer/models.html#main-classes)

\* [Further Classes](sentence\_transformer/models.html#further-classes)

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

\*

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

\*

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

\*

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

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

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

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

\* [Training Inputs](cross\_encoder/cross\_encoder.html#training-inputs)

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

\* [CEBinaryAccuracyEvaluator](cross\_encoder/evaluation.html#cebinaryaccuracyevaluator)

\*

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

\* [CECorrelationEvaluator](cross\_encoder/evaluation.html#cecorrelationevaluator)

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

\* [CESoftmaxAccuracyEvaluator](cross\_encoder/evaluation.html#cesoftmaxaccuracyevaluator)

\* [CERerankingEvaluator](cross\_encoder/evaluation.html#cererankingevaluator)

\* util

\* Helper Functions

\* `community\_detection()`

\* `http\_get()`

\* `is\_training\_available()`

\* `mine\_hard\_negatives()`

\* `normalize\_embeddings()`

\* `paraphrase\_mining()`

\* `semantic\_search()`

\* ``truncate_embeddings()``

## \* Model Optimization

\* ``export_dynamic_quantized_onnx_model()``

\* ``export_optimized_onnx_model()``

\* ``export_static_quantized_openvino_model()``

## \* Similarity Metrics

\* ``cos_sim()``

\* ``dot_score()``

\* ``euclidean_sim()``

\* ``manhattan_sim()``

\* ``pairwise_cos_sim()``

\* ``pairwise_dot_score()``

\* ``pairwise_euclidean_sim()``

\* ``pairwise_manhattan_sim()``

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

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

\* util

\* [ Edit on

GitHub](https://github.com/UKPLab/sentence-transformers/blob/master/docs/package\_reference/util  
.md)

\* \* \*

# util



`sentence\_transformers.util` defines different helpful functions to work with text embeddings.

## ## Helper Functions

```
sentence_transformers.util.community_detection(_embeddings :  
[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "\ (in PyTorch v2.5\)" ) | ndarray_,  
_threshold : float = 0.75_, _min_community_size : int = 10_, _batch_size : int = 1024_,  
_show_progress_bar : bool = False_) ->  
list[list[int]]
```

(source: [https://github.com/UKPLab/sentence-transformers/blob/master/sentence\\_transformers/util.py#L1151-L1250](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers/util.py#L1151-L1250))

Function for Fast Community Detection.

Finds in the embeddings all communities, i.e. embeddings that are close (closer than threshold). Returns only communities that are larger than min\_community\_size. The communities are returned in decreasing order. The first element in each list is the central point in the community.

Parameters:

\* **embeddings** ([\_torch.Tensor\_](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "\ (in PyTorch v2.5\)" ) \_or\_ \_numpy.ndarray\_) – The input embeddings.

\* **threshold** (\_float\_) â€“ The threshold for determining if two embeddings are close. Defaults to 0.75.

\* **min\_community\_size** (\_int\_) â€“ The minimum size of a community to be considered. Defaults to 10.

\* **batch\_size** (\_int\_) â€“ The batch size for computing cosine similarity scores. Defaults to 1024.

\* **show\_progress\_bar** (\_bool\_) â€“ Whether to show a progress bar during computation. Defaults to False.

Returns:

A list of communities, where each community is represented as a list of indices.

Return type:

List[List[int]]

sentence\_transformers.util.http\_get(\_url : str\_, \_path : str\_) ->

None[[source]]([https://github.com/UKPLab/sentence-](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\\util.py#L1033-L1067)

[transformers/blob/master/sentence\\_transformers\\util.py#L1033-L1067](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\\util.py#L1033-L1067))if•

Downloads a URL to a given path on disk.

Parameters:

\* **url** (\_str\_) â€“ The URL to download.

\* **path** (\_str\_) â€“ The path to save the downloaded file.

Raises:

**requests.HTTPError** â€“ If the HTTP request returns a non-200 status code.

Returns:

None

sentence\_transformers.util.is\_training\_available() ->

bool[[source]](<https://github.com/UKPLab/sentence->

transformers/blob/master/sentence\_transformers\\util.py#L1488-L1493)if•

Returns True if we have the required dependencies for training Sentence

Transformers models, i.e. Huggingface datasets and Huggingface accelerate.

```
sentence_transformers.util.mine_hard_negatives(_dataset : Dataset_, _model :  
[SentenceTransformer])(sentence_transformer/SentenceTransformer.html#sentence_transformers.S  
entenceTransformer "sentence_transformers.SentenceTransformer")_, _anchor_column_name : str  
| None = None_, _positive_column_name : str | None = None_, _corpus : list[str] | None = None_,  
_cross_encoder :  
[CrossEncoder](cross_encoder/cross_encoder.html#sentence_transformers.cross_encoder.CrossE  
ncoder "sentence_transformers.cross_encoder.CrossEncoder") | None = None_, _range_min : int =  
0_, _range_max : int | None = None_, _max_score : float | None = None_, _min_score : float | None  
= None_, _margin : float | None = None_, _num_negatives : int = 3_, _sampling_strategy :  
Literal['random', 'top'] = 'top'_, _as_triplets : bool = True_, _batch_size : int = 32_, _faiss_batch_size  
: int = 16384_, _use_faiss : bool = False_, _use_multi_process : list[str] | bool = False_, _verbose :  
bool = True_) ->  
Dataset[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transfor  
mers\\util.py#L520-L1030)if•
```

Add hard negatives to a dataset of (anchor, positive) pairs to create (anchor,  
positive, negative) triplets or (anchor, positive, negative\_1, ...,  
negative\_n) tuples.

Hard negative mining is a technique to improve the quality of a dataset by adding hard negatives, which are texts that may appear similar to the anchor, but are not. Using hard negatives can improve the performance of models trained on the dataset.

This function uses a SentenceTransformer model to embed the sentences in the dataset, and then finds the closest matches to each anchor sentence in the dataset. It then samples negatives from the closest matches, optionally using a CrossEncoder model to rescore the candidates.

You can influence the candidate negative selection in various ways:

- \* **range\_min** : Minimum rank of the closest matches to consider as negatives: useful to skip the most similar texts to avoid marking texts as negative that are actually positives.

- \* **range\_max** : Maximum rank of the closest matches to consider as negatives: useful to limit the number of candidates to sample negatives from. A lower value makes processing faster, but may result in less candidate negatives that satisfy the margin or max\_score conditions.

- \* **max\_score** : Maximum score to consider as a negative: useful to skip candidates that are too similar to the anchor.

- \* **min\_score** : Minimum score to consider as a negative: useful to skip candidates that are too dissimilar to the anchor.

- \* **margin** : Margin for hard negative mining: useful to skip candidates negatives whose similarity to the anchor is within a certain margin of the positive pair. A value of 0 can be used to enforce that

the negative is always further away from the anchor than the positive.

\* **sampling\_strategy** : Sampling strategy for negatives: `â€œtopâ€` or `â€œrandomâ€`. `â€œtopâ€` will always sample the top n candidates as negatives, while `â€œrandomâ€` will sample n negatives randomly from the candidates that satisfy the margin or max\_score conditions.

## Example

```
>>> from sentence_transformers.util import mine_hard_negatives
>>> from sentence_transformers import SentenceTransformer
>>> from datasets import load_dataset
>>> # Load a Sentence Transformer model
>>> model = SentenceTransformer("all-MiniLM-L6-v2")
>>>
>>> # Load a dataset to mine hard negatives from
>>> dataset = load_dataset("sentence-transformers/natural-questions", split="train")
>>> dataset
Dataset({
  features: ['query', 'answer'],
  num_rows: 100231
})
>>> dataset = mine_hard_negatives(
...     dataset=dataset,
...     model=model,
...     range_min=10,
```

$$\dots)$$

Batches:

17.83it/s]

Batches:

99.60it/s]

## Querying

FAISS

index:

â-â-â-â-â-â-â-â- 784/784 [00:00<00:00, 884.99it/s]

Count	100,231	431,255	431,255
-------	---------	---------	---------

Mean	0.6866	0.4289	0.2804
Median	0.7010	0.4193	0.2740
Std	0.1125	0.0754	0.0999
Min	0.0303	0.1720	0.1001
25%	0.6221	0.3747	0.1991
50%	0.7010	0.4193	0.2740
75%	0.7667	0.4751	0.3530
Max	0.9584	0.7743	0.7003

Skipped 1289492 potential negatives (25.23%) due to the margin of 0.1.

Skipped 39 potential negatives (0.00%) due to the maximum score of 0.8.

Could not find enough negatives for 69900 samples (13.95%). Consider adjusting the range\_max, range\_min, margin and max\_score parameters if you'd like to find more valid negatives.

```
>>> # Note: The minimum similarity difference is 0.1001 due to our margin of 0.1
```

```
>>> dataset
```

```
Dataset({
  features: ['query', 'answer', 'negative'],
  num_rows: 431255
})
```

```
>>> dataset[0]
```

```
{
  'query': 'when did richmond last play in a preliminary final',
  'answer': "Richmond Football Club Richmond began 2017 with 5 straight wins, a feat it had not
achieved since 1995. A series of close losses hampered the Tigers throughout the middle of the
season, including a 5-point loss to the Western Bulldogs, 2-point loss to Fremantle, and a 3-point
loss to the Giants. Richmond ended the season strongly with convincing victories over Fremantle
and St Kilda in the final two rounds, elevating the club to 3rd on the ladder. Richmond's first final of
the season against the Cats at the MCG attracted a record qualifying final crowd of 95,028; the
```



Tigers won by 51 points. Having advanced to the first preliminary finals for the first time since 2001, Richmond defeated Greater Western Sydney by 36 points in front of a crowd of 94,258 to progress to the Grand Final against Adelaide, their first Grand Final appearance since 1982. The attendance was 100,021, the largest crowd to a grand final since 1986. The Crows led at quarter time and led by as many as 13, but the Tigers took over the game as it progressed and scored seven straight goals at one point. They eventually would win by 48 points – 16.12 (108) to Adelaide's 8.12 (60) – to end their 37-year flag drought.[22] Dustin Martin also became the first player to win a Premiership medal, the Brownlow Medal and the Norm Smith Medal in the same season, while Damien Hardwick was named AFL Coaches Association Coach of the Year. Richmond's jump from 13th to premiers also marked the biggest jump from one AFL season to the next.",

'negative': "2018 NRL Grand Final The 2018 NRL Grand Final was the conclusive and premiership-deciding game of the 2018 National Rugby League season and was played on Sunday September 30 at Sydney's ANZ Stadium.[1] The match was contested between minor premiers the Sydney Roosters and defending premiers the Melbourne Storm. In front of a crowd of 82,688, Sydney won the match 21–6 to claim their 14th premiership title and their first since 2013. Roosters five-eighth Luke Keary was awarded the Clive Churchill Medal as the game's official man of the match."

}

```
>>> dataset.push_to_hub("natural-questions-hard-negatives", "triplet-all")
```

Parameters:

**dataset** (`_Dataset_`) – A dataset containing (anchor, positive) pairs.

\*

**\*\*model\*\***

```
([SentenceTransformer_](sentence_transformer/SentenceTransformer.html#sentence_transformer
s.SentenceTransformer "sentence_transformers.SentenceTransformer")) â€œ A
```

SentenceTransformer model to use for embedding the sentences.

\* **\*\*anchor\_column\_name\*\*** (\_str\_ \_\_, \_\_optional\_) â€œ The column name in dataset that contains the anchor/query. Defaults to None, in which case the first column in dataset will be used.

\* **\*\*positive\_column\_name\*\*** (\_str\_ \_\_, \_\_optional\_) â€œ The column name in dataset that contains the positive candidates. Defaults to None, in which case the second column in dataset will be used.

\* **\*\*corpus\*\*** (\_List\_ \_[\_str\_ \_] \_\_, \_\_optional\_) â€œ A list containing documents as strings that will be used as candidate negatives in addition to the second column in dataset. Defaults to None, in which case the second column in dataset will exclusively be used as the negative candidate corpus.

\*

**\*\*cross\_encoder\*\***

```
([CrossEncoder_](cross_encoder/cross_encoder.html#sentence_transformers.cross_encoder.Cros
sEncoder "sentence_transformers.cross_encoder.CrossEncoder") __, __optional_) â€œ A
```

CrossEncoder model to use for rescoring the candidates. Defaults to None.

\* **\*\*range\_min\*\*** (\_int\_) â€œ Minimum rank of the closest matches to consider as negatives. Defaults to 0.

\* **\*\*range\_max\*\*** (\_int\_ \_\_, \_\_optional\_) â€œ Maximum rank of the closest matches to consider as negatives. Defaults to None.

\* **\*\*max\_score\*\*** (\_float\_ \_\_, \_\_optional\_) â€œ Maximum score to consider as a negative. Defaults to

None.

\* **min\_score** (\_float\_ \_\_, \_\_optional\_) â€œ Minimum score to consider as a negative. Defaults to None.

\* **margin** (\_float\_ \_\_, \_\_optional\_) â€œ Margin for hard negative mining. Defaults to None.

\* **num\_negatives** (\_int\_) â€œ Number of negatives to sample. Defaults to 3.

\* **sampling\_strategy** (\_Literal\_ \_[\_\_" random"\_\_, \_\_" top"\_\_]\_) â€œ Sampling strategy for negatives: â€œtopâ€œ or â€œrandomâ€œ. Defaults to â€œtopâ€œ.

\* **as\_triplets** (\_bool\_) â€œ If True, returns up to num\_negatives (anchor, positive, negative) triplets for each input sample. If False, returns 1 (anchor, positive, negative\_1, â€œ, negative\_n) tuple for each input sample. Defaults to True.

\* **batch\_size** (\_int\_) â€œ Batch size for encoding the dataset. Defaults to 32.

\* **faiss\_batch\_size** (\_int\_) â€œ Batch size for FAISS top-k search. Defaults to 16384.

\* **use\_faiss** (\_bool\_) â€œ Whether to use FAISS for similarity search. May be recommended for large datasets. Defaults to False.

\* **use\_multi\_process** (\_bool\_ \_[\_\_List\_ \_[\_\_str\_ \_\_]\_\_, \_\_optional\_]\_) â€œ Whether to use multi-GPU/CPU processing. If True, uses all GPUs if CUDA is available, and 4 CPU processes if itâ€™s not available. You can also pass a list of PyTorch devices like [â€œcuda:0â€œ, â€œcuda:1â€œ, â€œiâ€œ] or [â€œcpuâ€œ, â€œcpuâ€œ, â€œcpuâ€œ, â€œcpuâ€œ].

\* **verbose** (`_bool_`) – Whether to print statistics and logging. Defaults to True.

Returns:

A dataset containing (anchor, positive, negative) triplets or (anchor, positive, negative\_1, ..., negative\_n) tuples.

Return type:

Dataset

`sentence_transformers.util.normalize_embeddings(_embeddings :  
[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in  
PyTorch v2.5\))"_) ->  
[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in  
PyTorch v2.5\))")` [[source]] ([https://github.com/UKPLab/sentence-  
transformers/blob/master/sentence\\_transformers\util.py#L264-L274](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\util.py#L264-L274)) if •

Normalizes the embeddings matrix, so that each sentence embedding has unit length.

Parameters:

**\*\*embeddings\*\*** (`_Tensor_`) – The input embeddings matrix.

Returns:

The normalized embeddings matrix.

Return type:

Tensor

```
sentence_transformers.util.paraphrase_mining(_model, sentences: list[str], show_progress_bar:
bool = False, batch_size: int = 32, query_chunk_size: int = 5000, corpus_chunk_size: int = 100000,
max_pairs: int = 500000, top_k: int = 100, score_function: ~typing.Callable[[_torch.Tensor,
~torch.Tensor], ~torch.Tensor] = <function cos_sim>) -> list[list[float |
int]]
```

[[source]]([https://github.com/UKPLab/sentence-transformers/blob/master/sentence\\_transformers\\util.py#L317-L359](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\\util.py#L317-L359))if•

Given a list of sentences / texts, this function performs paraphrase mining.

It compares all sentences against all other sentences and returns a list with the pairs that have the highest cosine similarity score.

Parameters:

**\* \*\*model\*\***  
([\_SentenceTransformer\_]([sentence\\_transformer/SentenceTransformer.html#sentence\\_transformers.SentenceTransformer](https://www.sbert.net/docs/sentence_transformer/SentenceTransformer.html#sentence_transformers.SentenceTransformer) "sentence\_transformers.SentenceTransformer")) â€“ SentenceTransformer model for embedding computation

**\* \*\*sentences\*\*** (\_List\_ \_[\_\_str\_ \_]) â€“ A list of strings (texts or sentences)

**\* \*\*show\_progress\_bar\*\*** (\_bool\_ \_\_, \_\_optional\_\_) â€“ Plotting of a progress bar. Defaults to False.

**\* \*\*batch\_size\*\*** (\_int\_ \_\_, \_\_optional\_\_) â€“ Number of texts that are encoded simultaneously by the model. Defaults to 32.

**\* \*\*query\_chunk\_size\*\*** (\_int\_ \_\_, \_\_optional\_\_) â€“ Search for most similar pairs for #query\_chunk\_size at the same time. Decrease, to lower memory footprint (increases run-time). Defaults to 5000.

**\* \*\*corpus\_chunk\_size\*\*** (\_int\_ \_\_, \_\_optional\_\_) â€“ Compare a sentence simultaneously against #corpus\_chunk\_size other sentences. Decrease, to lower memory footprint (increases run-time). Defaults to 100000.

\* **max\_pairs** (\_int\_ \_\_, \_\_optional\_) â€“ Maximal number of text pairs returned. Defaults to 500000.

\* **top\_k** (\_int\_ \_\_, \_\_optional\_) â€“ For each sentence, we retrieve up to top\_k other sentences. Defaults to 100.

\* **score\_function** (\_Callable\_ \_[[\_Tensor\_ \_\_, \_\_Tensor\_ \_]\_\_, \_\_Tensor\_ \_]\_\_, \_\_optional\_) â€“ Function for computing scores. By default, cosine similarity. Defaults to cos\_sim.

Returns:

Returns a list of triplets with the format [score, id1, id2]

Return type:

List[List[Union[float, int]]]

sentence\_transformers.util.semantic\_search(\_query\_embeddings: ~torch.Tensor,  
corpus\_embeddings: ~torch.Tensor, query\_chunk\_size: int = 100, corpus\_chunk\_size: int = 500000,  
top\_k: int = 10, score\_function: ~typing.Callable[[\_Tensor, \_Tensor], ~torch.Tensor] =  
<function cos\_sim>\_) -> list[list[dict[str, int | float]]][[source]]([https://github.com/UKPLab/sentence-transformers/blob/master/sentence\\_transformers](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers)

ers\util.py#L440-L517)if•

This function performs a cosine similarity search between a list of query embeddings and a list of corpus embeddings. It can be used for Information Retrieval / Semantic Search for corpora up to about 1 Million entries.

Parameters:

\* \*\*query\_embeddings\*\* (\_Tensor\_) â€“ A 2 dimensional tensor with the query embeddings.

\* \*\*corpus\_embeddings\*\* (\_Tensor\_) â€“ A 2 dimensional tensor with the corpus embeddings.

\* \*\*query\_chunk\_size\*\* (\_int\_ \_\_, \_\_optional\_\_) â€“ Process 100 queries simultaneously. Increasing that value increases the speed, but requires more memory. Defaults to 100.

\* \*\*corpus\_chunk\_size\*\* (\_int\_ \_\_, \_\_optional\_\_) â€“ Scans the corpus 100k entries at a time. Increasing that value increases the speed, but requires more memory. Defaults to 500000.

\* \*\*top\_k\*\* (\_int\_ \_\_, \_\_optional\_\_) â€“ Retrieve top k matching entries. Defaults to 10.

\* \*\*score\_function\*\* (\_Callable\_ \_[[[\_Tensor\_ \_\_, \_\_Tensor\_ \_]\_\_], \_\_Tensor\_ \_]\_\_, \_\_optional\_\_) â€“ Function for computing scores. By default, cosine similarity.

Returns:



A list with one entry for each query. Each entry is a list of dictionaries with the keys `corpus_id` and `score`, sorted by decreasing cosine similarity scores.

Return type:

`List[List[Dict[str, Union[int, float]]]]`

`sentence_transformers.util.truncate_embeddings(_embeddings : ndarray, _truncate_dim : int | None)` ->

`ndarray`[[source]]([https://github.com/UKPLab/sentence-transformers/blob/master/sentence\\_transformers/util.py#L285-L314](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers/util.py#L285-L314)) if •

`sentence_transformers.util.truncate_embeddings(_embeddings : [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\))", _truncate_dim : int | None) -> [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\))"`

Truncates the embeddings matrix.

Parameters:

\* **\*\*embeddings\*\*** (`_Union_` `_[_np.ndarray_`  
`_[_torch.Tensor_]`(<https://pytorch.org/docs/stable/tensors.html#torch.Tensor> `"\in PyTorch v2.5\")`  
`_]`) â€œ Embeddings to truncate.

\* **\*\*truncate\_dim\*\*** (`_Optional_` `_[_int_ _]`) â€œ The dimension to truncate sentence embeddings to. None does no truncation.

## Example

```
>>> from sentence_transformers import SentenceTransformer
>>> from sentence_transformers.util import truncate_embeddings
>>> model = SentenceTransformer("tomaarsen/mpnet-base-nli-matryoshka")
>>> embeddings = model.encode(["It's so nice outside!", "Today is a beautiful day.", "He drove to
work earlier"])
>>> embeddings.shape
(3, 768)
>>> model.similarity(embeddings, embeddings)
tensor([[1.0000, 0.8100, 0.1426],
        [0.8100, 1.0000, 0.2121],
        [0.1426, 0.2121, 1.0000]])
>>> truncated_embeddings = truncate_embeddings(embeddings, 128)
>>> truncated_embeddings.shape
>>> model.similarity(truncated_embeddings, truncated_embeddings)
```

```
tensor([[1.0000, 0.8092, 0.1987],
        [0.8092, 1.0000, 0.2716],
        [0.1987, 0.2716, 1.0000]])
```

Returns:

Truncated embeddings.

Return type:

```
Union[np.ndarray,
      [torch.Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "\in
      PyTorch v2.5\"))]
```

## Model Optimization

```
sentence_transformers.backend.export_dynamic_quantized_onnx_model(_model :
[SentenceTransformer](sentence\_transformer/SentenceTransformer.html#sentence\_transformers.S
entenceTransformer "sentence_transformers.SentenceTransformer"), _quantization_config :
QuantizationConfig | Literal['arm64', 'avx2', 'avx512', 'avx512_vnni'], _model_name_or_path : str,
_push_to_hub : bool = False, _create_pr : bool = False, _file_suffix : str | None = None) ->
None[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transforme
```

rs\\backend.py#L117-L198)if•

Export a quantized ONNX model from a SentenceTransformer model.

This function applies dynamic quantization, i.e. without a calibration dataset. Each of the default quantization configurations quantize the model to int8, allowing for faster inference on CPUs, but are likely slower on GPUs.

See <[https://sbert.net/docs/sentence\\_transformer/usage/efficiency.html](https://sbert.net/docs/sentence_transformer/usage/efficiency.html)> for more information & benchmarks.

Parameters:

\* **model** ([SentenceTransformer]) (sentence\_transformer/SentenceTransformer.html#sentence\_transformer.SentenceTransformer "sentence\_transformers.SentenceTransformer") â€œ The SentenceTransformer model to be quantized. Must be loaded with backend=â€•onnxâ€•.

\* **quantization\_config** (\_QuantizationConfig) â€œ The quantization configuration.

\* **model\_name\_or\_path** (\_str\_) â€œ The path or Hugging Face Hub repository name where the quantized model will be saved.

\* **push\_to\_hub** (\_bool\_, \_\_optional\_) â€œ Whether to push the quantized model to the Hugging

Face Hub. Defaults to False.

\* \*\*create\_pr\*\* (\_bool\_ \_\_, \_\_optional\_\_) â€œ Whether to create a pull request when pushing to the Hugging Face Hub. Defaults to False.

\* **file\_suffix** (`_str_ | __None_`, `__optional_`) â€” The suffix to add to the quantized model file name. Defaults to None.

Raises:

**\*\*ImportError\*\*** â€” If the required packages optimum and onnxruntime are not installed.

\* **ValueError** â€” If the provided model is not a valid SentenceTransformer model loaded with backend=â€•onnxâ€•.

**ValueError** – If the provided `quantization_config` is not valid.

Returns:

None

```
sentence_transformers.backend.export_optimized_onnx_model(_model :  
[SentenceTransformer](sentence_transformer/SentenceTransformer.html#sentence_transformers.S
```

```

entenceTransformer("sentence_transformers.SentenceTransformer"), _optimization_config :
OptimizationConfig | Literal['O1', 'O2', 'O3', 'O4'], _model_name_or_path : str, _push_to_hub :
bool = False, _create_pr : bool = False, _file_suffix : str | None = None) ->
None[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transforme
rs\\backend.py#L28-L114)if•

```

Export an optimized ONNX model from a SentenceTransformer model.

The O1-O4 optimization levels are defined by Optimum and are documented here:

<[https://huggingface.co/docs/optimum/main/en/onnxruntime/usage\\_guides/optimization](https://huggingface.co/docs/optimum/main/en/onnxruntime/usage_guides/optimization)>

The optimization levels are:

- \* O1: basic general optimizations.
- \* O2: basic and extended general optimizations, transformers-specific fusions.
- \* O3: same as O2 with GELU approximation.
- \* O4: same as O3 with mixed precision (fp16, GPU-only)

See <[https://sbert.net/docs/sentence\\_transformer/usage/efficiency.html](https://sbert.net/docs/sentence_transformer/usage/efficiency.html)> for  
more information & benchmarks.

Parameters:

\*

**\*\*model\*\***

([\_SentenceTransformer\_](sentence\_transformer/SentenceTransformer.html#sentence\_transformer  
s.SentenceTransformer "sentence\_transformers.SentenceTransformer")) â€œ The  
SentenceTransformer model to be optimized. Must be loaded with backend=â€•onnxâ€•.

\* **\*\*optimization\_config\*\*** (\_OptimizationConfig\_ | \_Literal\_ [\_ " O1" \_\_, \_\_ " O2" \_\_, \_\_ " O3" \_\_, \_\_ "  
O4" \_\_]\_ ) â€œ The optimization configuration or level.

\* **\*\*model\_name\_or\_path\*\*** (\_str\_) â€œ The path or Hugging Face Hub repository name where the  
optimized model will be saved.

\* **\*\*push\_to\_hub\*\*** (\_bool\_ \_\_, \_\_optional\_) â€œ Whether to push the optimized model to the Hugging  
Face Hub. Defaults to False.

\* **\*\*create\_pr\*\*** (\_bool\_ \_\_, \_\_optional\_) â€œ Whether to create a pull request when pushing to the  
Hugging Face Hub. Defaults to False.

\* **\*\*file\_suffix\*\*** (\_str\_ | \_None\_ \_\_, \_\_optional\_) â€œ The suffix to add to the optimized model file  
name. Defaults to None.

Raises:

\* **\*\*ImportError\*\*** â€œ If the required packages optimum and onnxruntime are not installed.

\* \*\*ValueError\*\* â€“ If the provided model is not a valid SentenceTransformer model loaded with backend=â€•onnxâ€•.

\* \*\*ValueError\*\* â€“ If the provided optimization\_config is not valid.

Returns:

None

```
sentence_transformers.backend.export_static_quantized_openvino_model(_model :  
[SentenceTransformer](sentence_transformer/SentenceTransformer.html#sentence_transformers.S  
entenceTransformer "sentence_transformers.SentenceTransformer"), _quantization_config :  
OVQuantizationConfig | dict | None_, _model_name_or_path : str_, _dataset_name : str | None =  
None_, _dataset_config_name : str | None = None_, _dataset_split : str | None = None_,  
_column_name : str | None = None_, _push_to_hub : bool = False_, _create_pr : bool = False_,  
_file_suffix : str = 'qint8_quantized') ->  
None[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transforme  
rs\\backend.py#L201-L309)if•
```

Export a quantized OpenVINO model from a SentenceTransformer model.

This function applies Post-Training Static Quantization (PTQ) using a calibration dataset, which calibrates quantization constants without requiring



model retraining. Each default quantization configuration converts the model to int8 precision, enabling faster inference while maintaining accuracy.

See <[https://sbert.net/docs/sentence\\_transformer/usage/efficiency.html](https://sbert.net/docs/sentence_transformer/usage/efficiency.html)> for more information & benchmarks.

Parameters:

```

*
**model**

([_SentenceTransformer_](sentence_transformer/SentenceTransformer.html#sentence_transformer
s.SentenceTransformer "sentence_transformers.SentenceTransformer"))  â€œ The
SentenceTransformer model to be quantized. Must be loaded with backend=â€•openvinoâ€•.

* **quantization_config** (_OVQuantizationConfig_ |__dict_ |__None_)  â€œ The quantization
configuration. If None, default values are used.

* **model_name_or_path** (_str_)  â€œ The path or Hugging Face Hub repository name where the
quantized model will be saved.

* **dataset_name** (_str_ __, __optional_)  â€œ The name of the dataset to load for calibration. If not
specified, the sst2 subset of the glue dataset will be used by default.

* **dataset_config_name** (_str_ __, __optional_)  â€œ The specific configuration of the dataset to
load.
```

\* **dataset\_split** (\_str\_ \_\_, \_\_optional\_) â€“ The split of the dataset to load (e.g., â€˜trainâ€™™, â€˜testâ€™™). Defaults to None.

\* **column\_name** (\_str\_ \_\_, \_\_optional\_) â€“ The column name in the dataset to use for calibration. Defaults to None.

\* **push\_to\_hub** (\_bool\_ \_\_, \_\_optional\_) â€“ Whether to push the quantized model to the Hugging Face Hub. Defaults to False.

\* **create\_pr** (\_bool\_ \_\_, \_\_optional\_) â€“ Whether to create a pull request when pushing to the Hugging Face Hub. Defaults to False.

\* **file\_suffix** (\_str\_ \_\_, \_\_optional\_) â€“ The suffix to add to the quantized model file name. Defaults to qint8\_quantized.

Raises:

\* **ImportError** â€“ If the required packages optimum and openvino are not installed.

\* **ValueError** â€“ If the provided model is not a valid SentenceTransformer model loaded with backend=â€•openvinoâ€•.

\* **ValueError** â€“ If the provided quantization\_config is not valid.

Returns:

None

## ## Similarity Metrics

```
def cos_sim(_a : list | ndarray | Tensor, _b : list | ndarray | Tensor) -> Tensor:
    """Computes the cosine similarity between two tensors.
    Parameters:
        _a (Union[list, np.ndarray, Tensor]) â€œ The first tensor.
        _b (Union[list, np.ndarray, Tensor]) â€œ The second tensor.
    Returns:
        A Tensor representing the cosine similarity between the two tensors.
    """
    # Convert inputs to tensors if they are not already
    a = torch.tensor(_a) if isinstance(_a, (list, ndarray)) else _a
    b = torch.tensor(_b) if isinstance(_b, (list, ndarray)) else _b

    # Compute the dot product and the norms
    dot_product = torch.dot(a, b)
    norm_a = torch.norm(a)
    norm_b = torch.norm(b)

    # Compute the cosine similarity
    cos_sim = dot_product / (norm_a * norm_b)

    return cos_sim
```

Computes the cosine similarity between two tensors.

Parameters:

**a** (Union[list, np.ndarray, Tensor]) â€œ The first tensor.

**b** (Union[list, np.ndarray, Tensor]) â€œ The second tensor.

Returns:

Matrix with  $\text{res}[i][j] = \text{cos\_sim}(a[i], b[j])$

Return type:

Tensor

`sentence_transformers.util.dot_score(_a : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)"_) , _b : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)"_)) -> [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)"_)[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transformers\\util.py#L128-L142)if•`

Computes the dot-product  $\text{dot\_prod}(a[i], b[j])$  for all  $i$  and  $j$ .

Parameters:

**a** (`_Union_` [`_list_`, `_np.ndarray_`, `_Tensor_`]) â€œ The first tensor.

**b** (`_Union_` [`_list_`, `_np.ndarray_`, `_Tensor_`]) â€œ The second tensor.

Returns:

Matrix with  $\text{res}[i][j] = \text{dot\_prod}(a[i], b[j])$

Return type:

Tensor

`sentence_transformers.util.euclidean_sim(_a : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5)"), _b : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5)"))`  
-> `[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5)")`[[source]]([https://github.com/UKPLab/sentence-transformers/blob/master/sentence\\_transformers\\util.py#L196-L210](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\\util.py#L196-L210))if•

Computes the euclidean similarity (i.e., negative distance) between two tensors.

Parameters:

**a** (\_Union\_ [\_list\_ \_\_, \_\_np.ndarray\_ \_\_, \_\_Tensor\_ \_]) â€œ The first tensor.

**b** (\_Union\_ [\_list\_ \_\_, \_\_np.ndarray\_ \_\_, \_\_Tensor\_ \_]) â€œ The second tensor.

Returns:

Matrix with  $\text{res}[i][j] = -\text{euclidean\_distance}(a[i], b[j])$

Return type:

Tensor

`sentence_transformers.util.manhattan_sim(_a : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)" ), _b : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)" ))`  
-> `[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)" )`[[source]]([https://github.com/UKPLab/sentence-transformers/blob/master/sentence\\_transformers\\util.py#L162-L176](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\\util.py#L162-L176))if•

Computes the manhattan similarity (i.e., negative distance) between two tensors.

Parameters:

**a** (`_Union_` [`__list_`, `__np.ndarray_`, `__Tensor_`]) â€” The first tensor.

**b** (`_Union_` [`__list_`, `__np.ndarray_`, `__Tensor_`]) â€” The second tensor.

Returns:

Matrix with `res[i][j] = -manhattan_distance(a[i], b[j])`

Return type:

Tensor

`sentence_transformers.util.pairwise_cos_sim(_a :`

[Tensor](<https://pytorch.org/docs/stable/tensors.html#torch.Tensor> "(in PyTorch v2.5\))", \_b :

[Tensor](<https://pytorch.org/docs/stable/tensors.html#torch.Tensor> "(in PyTorch v2.5\))") ->

[Tensor](<https://pytorch.org/docs/stable/tensors.html#torch.Tensor> "(in

PyTorch v2.5\")[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transformers/util.py#L111-L125)if•

Computes the pairwise cosine similarity `cos_sim(a[i], b[i])`.

Parameters:

**a** (`Union[_List_, np.ndarray_, Tensor_]`) â€œ The first tensor.`

**b** (`Union[_List_, np.ndarray_, Tensor_]`) â€œ The second tensor.`

Returns:

Vector with `res[i] = cos_sim(a[i], b[i])`

Return type:

Tensor



```

sentence_transformers.util.pairwise_dot_score(_a :
[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in
PyTorch v2.5\))", _b :
[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in
PyTorch v2.5\))") ->
[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in
PyTorch v2.5\))")[[source]](https://github.com/UKPLab/sentence-
transformers/blob/master/sentence_transformers\util.py#L145-L159)if•

```

Computes the pairwise dot-product `dot_prod(a[i], b[i])`.

Parameters:

**a** (`_Union_` [`list`, `np.ndarray`, `Tensor`]) â€œ The first tensor.

**b** (`_Union_` [`list`, `np.ndarray`, `Tensor`]) â€œ The second tensor.

Returns:

Vector with `res[i] = dot_prod(a[i], b[i])`

Return type:

Tensor

`sentence_transformers.util.pairwise_euclidean_sim(_a : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)" ), _b : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)" ))`[[source]]([https://github.com/UKPLab/sentence-transformers/blob/master/sentence\\_transformers\\util.py#L213-L227](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\\util.py#L213-L227))if•

Computes the euclidean distance (i.e., negative distance) between pairs of tensors.

Parameters:

**a** (`_Union_` [`list`, `np.ndarray`, `Tensor`]) â€œ The first tensor.

**b** (`_Union_` [`list`, `np.ndarray`, `Tensor`]) â€œ The second tensor.

Returns:

Vector with `res[i] = -euclidean_distance(a[i], b[i])`

Return type:

Tensor

`sentence_transformers.util.pairwise_manhattan_sim(_a : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)" ), _b : list | ndarray | [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\)" ))`[[source]]([https://github.com/UKPLab/sentence-transformers/blob/master/sentence\\_transformers\\util.py#L179-L193](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\\util.py#L179-L193))*if*•

Computes the manhattan similarity (i.e., negative distance) between pairs of tensors.

Parameters:

**a** (`_Union_` [`list`, `np.ndarray`, `Tensor`]) â€œ The first tensor.

**b** (`_Union_` [`list`, `np.ndarray`, `Tensor`]) â€œ The second tensor.

Returns:

Vector with  $\text{res}[i] = -\text{manhattan\_distance}(a[i], b[i])$

Return type:

Tensor

[ Previous](cross\_encoder/evaluation.html "Evaluation")

\* \* \*

(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).