

[ ![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](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)

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

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

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

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

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

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

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

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

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

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

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

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

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

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

- \* [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)
          - \* [Margin Based

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-quantiz

ation)

- \* [Additional

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

s)

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

- \* [Try it

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

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

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

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

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

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

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

- \* [Structure of Sentence Transformer

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

Model](usage/custom\_models.html#sentence-transformer-model-from-a-transformers-model)

- \* [Pretrained Models](pretrained\_models.html)
- \* [Original Models](pretrained\_models.html#original-models)
- \* [Semantic Search Models](pretrained\_models.html#semantic-search-models)
- \* [Multi-QA Models](pretrained\_models.html#multi-qa-models)
- \* [MSMARCO Passage Models](pretrained\_models.html#msmarco-passage-models)
- \* [Multilingual Models](pretrained\_models.html#multilingual-models)
- \* [Semantic Similarity Models](pretrained\_models.html#semantic-similarity-models)
- \* [Bitext Mining](pretrained\_models.html#bitext-mining)
- \* [Image & Text-Models](pretrained\_models.html#image-text-models)
- \* [INSTRUCTOR models](pretrained\_models.html#instructor-models)
- \* [Scientific Similarity Models](pretrained\_models.html#scientific-similarity-models)
- \* [Training Overview](training\_overview.html)
- \* [Why Finetune?](training\_overview.html#why-finetune)
- \* [Training Components](training\_overview.html#training-components)
- \* [Dataset](training\_overview.html#dataset)
  - \* [Dataset Format](training\_overview.html#dataset-format)
- \* [Loss Function](training\_overview.html#loss-function)
- \* [Training Arguments](training\_overview.html#training-arguments)
- \* [Evaluator](training\_overview.html#evaluator)
- \* [Trainer](training\_overview.html#trainer)
  - \* [Callbacks](training\_overview.html#callbacks)
- \* [Multi-Dataset Training](training\_overview.html#multi-dataset-training)
- \* [Deprecated Training](training\_overview.html#deprecated-training)
- \* [Best Base Embedding Models](training\_overview.html#best-base-embedding-models)
- \* [Dataset Overview](dataset\_overview.html)

- \* [Datasets on the Hugging Face Hub](dataset\_overview.html#datasets-on-the-hugging-face-hub)
- \* [Pre-existing Datasets](dataset\_overview.html#pre-existing-datasets)
- \* Loss Overview
  - \* Loss modifiers
  - \* Distillation
  - \* Commonly used Loss Functions
  - \* Custom Loss Functions
- \* [Training Examples](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#multiplenegativesrankingloss)

\* [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](training/distributed.html)
- \* [Comparison](training/distributed.html#comparison)
- \* [FSDP](training/distributed.html#fsdp)

## Cross Encoder

- \* [Usage](../cross\_encoder/usage/usage.html)
- \* [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)
    - \* [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](../package\_reference/sentence\_transformer/index.html)
  - \* [SentenceTransformer](../package\_reference/sentence\_transformer/SentenceTransformer.html)
- \*
- [SentenceTransformer](../package\_reference/sentence\_transformer/SentenceTransformer.html#id1)
- \*
- [SentenceTransformerModelCardData](../package\_reference/sentence\_transformer/SentenceTransformerModelCardData.html#sentencetransformermodelcarddata)
- \*
- [SimilarityFunction](../package\_reference/sentence\_transformer/SentenceTransformer.html#similarityfunction)
- \* [Trainer](../package\_reference/sentence\_transformer/trainer.html)
- \*
- [SentenceTransformerTrainer](../package\_reference/sentence\_transformer/trainer.html#sentencetransformertrainer)

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

\*

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

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

\*

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

\*

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

\*

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

\*

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

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

\*

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

\*

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

\*

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

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

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

\*

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

\*

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

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

\*

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

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

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

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

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

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

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

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

\*

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

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

\*

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

\*

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

\*

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

\*

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

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

\*

[ParaphraseMiningEvaluator](../package\_reference/sentence\_transformer/evaluation.html#paraphrase-mining-evaluator)

\*

[RerankingEvaluator](../package\_reference/sentence\_transformer/evaluation.html#reranking-evaluator)

\*

[SentenceEvaluator](../package\_reference/sentence\_transformer/evaluation.html#sentence-evaluator)

\*

[SequentialEvaluator](../package\_reference/sentence\_transformer/evaluation.html#sequential-evaluator)

\*

[TranslationEvaluator](../package\_reference/sentence\_transformer/evaluation.html#translation-evaluator)

\* [TripletEvaluator](../package\_reference/sentence\_transformer/evaluation.html#triplet-evaluator)

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

\*

[ParallelSentencesDataset](../package\_reference/sentence\_transformer/datasets.html#parallel-sentences-dataset)

\*

[SentenceLabelDataset](../package\_reference/sentence\_transformer/datasets.html#sentence-label-dataset)

\*

[DenoisingAutoEncoderDataset](../package\_reference/sentence\_transformer/datasets.html#denoising-auto-encoder-dataset)

\*

[NoDuplicatesDataLoader](../package\_reference/sentence\_transformer/datasets.html#no-duplicates-dataloader)

ataloader)

- \* [Models](../package\_reference/sentence\_transformer/models.html)
- \* [Main Classes](../package\_reference/sentence\_transformer/models.html#main-classes)
- \* [Further Classes](../package\_reference/sentence\_transformer/models.html#further-classes)
- \* [quantization](../package\_reference/sentence\_transformer/quantization.html)

\*

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

\*

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

\*

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

- \* [Cross Encoder](../package\_reference/cross\_encoder/index.html)
- \* [CrossEncoder](../package\_reference/cross\_encoder/cross\_encoder.html)
- \* [CrossEncoder](../package\_reference/cross\_encoder/cross\_encoder.html#id1)
- \* [Training Inputs](../package\_reference/cross\_encoder/cross\_encoder.html#training-inputs)
- \* [Evaluation](../package\_reference/cross\_encoder/evaluation.html)

\*

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

\*

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

\*

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



or)

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

\*

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

\*

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

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

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

\*

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

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

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

\*

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

\*

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

\*

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

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

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

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

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

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

\*

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

\*

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

\*

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

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

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

\* [(../index.html)

\* Loss Overview

\* [ Edit on

GitHub](https://github.com/UKPLab/sentence-transformers/blob/master/docs/sentence\_transformer/loss\_overview.md)

\* \* \*

## # Loss Overview¶

Loss functions play a critical role in the performance of your fine-tuned model. Sadly, there is no “one size fits all” loss function. Ideally, this table should help narrow down your choice of loss function(s) by matching them to your data formats.

### Note

You can often convert one training data format into another, allowing more loss functions to be viable for your scenario. For example, `(sentence\_A, sentence\_B)` pairs with `class` labels can be converted into `(anchor, positive, negative)` triplets by sampling sentences with the same or different

classes.

Inputs | Labels | Appropriate Loss Functions

---|---|---

<code>`single</code>	<code>sentences`</code>		<code>`class`</code>	
[ `BatchAllTripletLoss` ](../package_reference/sentence_transformer/losses.html#batchalltripletloss)				
[ `BatchHardSoftMarginTripletLoss` ](../package_reference/sentence_transformer/losses.html#batchhardsoftmargintripletloss)				
[ `BatchHardTripletLoss` ](../package_reference/sentence_transformer/losses.html#batchhardtripletloss)				
[ `BatchSemiHardTripletLoss` ](../package_reference/sentence_transformer/losses.html#batchsemihardtripletloss)				
<code>`single</code>	<code>sentences`</code>		<code>`none`</code>	
[ `ContrastiveTensionLoss` ](../package_reference/sentence_transformer/losses.html#contrastivetensionloss)				
[ `DenoisingAutoEncoderLoss` ](../package_reference/sentence_transformer/losses.html#denoisingautoencoderloss)				
<code>`(anchor,</code>	<code>anchor)</code>	<code>pairs`</code>		<code>`none`</code>
[ `ContrastiveTensionLossInBatchNegatives` ](../package_reference/sentence_transformer/losses.html#contrastivetensionlossinbatchnegatives)				
<code>`(damaged_sentence,</code>	<code>original_sentence)</code>	<code>pairs`</code>		<code>`none`</code>
[ `DenoisingAutoEncoderLoss` ](../package_reference/sentence_transformer/losses.html#denoisingautoencoderloss)				
<code>`(sentence_A,</code>	<code>sentence_B)</code>	<code>pairs`</code>		<code>`class`</code>
[ `SoftmaxLoss` ](../package_reference/sentence_transformer/losses.html#softmaxloss)				
<code>`(anchor,</code>	<code>positive)</code>	<code>pairs`</code>		<code>`none`</code>
[ `MultipleNegativesRankingLoss` ](../package_reference/sentence_transformer/losses.html#multiple				

negativesrankingloss)

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

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

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

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

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

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

`(anchor, positive/negative) pairs` | `1 if positive, 0 if negative` |

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

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

`(sentence\_A, sentence\_B) pairs` | `float similarity score` |

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

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

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

`(anchor, positive, negative) triplets` | `none` |

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

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

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

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

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

`(anchor, positive, negative\_1, ..., negative\_n)` | `none` |

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

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

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

## ## Loss modifiers¶

These loss functions can be seen as `_loss modifiers_` : they work on top of standard loss functions, but apply those loss functions in different ways to try and instil useful properties into the trained embedding model.

For example, models trained with

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

produce embeddings whose size can be truncated without notable losses in performance, and models trained with

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

still perform well when you remove model layers for faster inference.

Texts | Labels | Appropriate Loss Functions

---|---|---

`any` | `any` |

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

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

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

## ## Distillation¶

These loss functions are specifically designed to be used when distilling the knowledge from one model into another. For example, when finetuning a small model to behave more like a larger & stronger one, or when finetuning a model to become multi-lingual.

Texts | Labels | Appropriate Loss Functions

---|---|---

`sentence` | `model` sentence embeddings` |

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

`sentence\_1, sentence\_2, ..., sentence\_N` | `model` sentence embeddings` |

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

`(query, passage\_one, passage\_two) triplets` | `gold\_sim(query, passage\_one) - gold\_sim(query, passage\_two)` |

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

## ## Commonly used Loss Functions¶

In practice, not all loss functions get used equally often. The most common scenarios are:

\* `(anchor, positive) pairs` without any labels:

[`MultipleNegativesRankingLoss`](../package\_reference/sentence\_transformer/losses.html#multiple-negativesrankingloss) is commonly used to train the top performing embedding models. This data is often relatively cheap to obtain, and the models are generally very performant. [`CachedMultipleNegativesRankingLoss`](../package\_reference/sentence\_transformer/losses.html#cachedmultiplenegativesrankingloss) is often used to increase the batch size, resulting in superior performance.

\* `(sentence_A, sentence_B) pairs` with a `float similarity score`:`

[`CosineSimilarityLoss`](../package\_reference/sentence\_transformer/losses.html#cosinesimilarityloss) is traditionally used a lot, though more recently [`CoSENTLoss`](../package\_reference/sentence\_transformer/losses.html#cosentloss) and [`AnglELoss`](../package\_reference/sentence\_transformer/losses.html#angleloss) are used as drop-in replacements with superior performance.

## ## Custom Loss Functions¶

Advanced users can create and train with their own loss functions. Custom loss functions only have a few requirements:

\* They must be a subclass of [`torch.nn.Module`](https://pytorch.org/docs/stable/generated/torch.nn.Module.html#torch.nn.Module "(in PyTorch v2.5)").

\* They must have ``model`` as the first argument in the constructor.

\* They must implement a ``forward`` method that accepts ``sentence_features`` and ``labels``. The former is a list of tokenized batches, one element for each column. These tokenized batches can be



fed directly to the `model` being trained to produce embeddings. The latter is an optional tensor of labels. The method must return a single loss value.

To get full support with the automatic model card generation, you may also wish to implement:

- \* a `get_config_dict` method that returns a dictionary of loss parameters.
- \* a `citation` property so your work gets cited in all models that train with the loss.

[ [Previous](#)](dataset\_overview.html "Dataset Overview") [ [Next](#)](training/examples.html "Training Examples")

\* \* \*

(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).