

[ ![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)

[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)

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)

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

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

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

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

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)

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

\* [\[Sentence Transformer Model from a Transformers](#)

[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\_overview.html)

\* [Loss modifiers](../loss\_overview.html#loss-modifiers)

\* [Distillation](../loss\_overview.html#distillation)

\* [Commonly used Loss Functions](../loss\_overview.html#commonly-used-loss-functions)

\* [Custom Loss Functions](../loss\_overview.html#custom-loss-functions)

\* [Training Examples](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#multiplenegativesrankin  
gloss)

\* [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-labelin

g)

- \* [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
- \* Comparison
- \* 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)
- r)
  - \* [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/SentenceTran

sformer.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#onlinecontrastiv

eloss)

\*

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

\*

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

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

\* [AngleELoss](../../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#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)

\* [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#binary classificationevaluator)

\*

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

eddingsimilarityevaluator)

\*

[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#paraphraseminingevaluator)

\*

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

\*

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

\*

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

\*

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

\*

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

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

\*

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

\*

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

\*

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

\*

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

- \* [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)

\* [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)]

\* [Training Examples](examples.html)

\* Distributed Training

\*

[

Edit

on

GitHub](https://github.com/UKPLab/sentence-transformers/blob/master/docs/sentence\_transformer/training/distributed.rst)

\* \* \*

## # Distributed Training

Sentence Transformers implements two forms of distributed training: Data Parallel (DP) and Distributed Data Parallel (DDP). Read the [Data Parallelism documentation](https://huggingface.co/docs/transformers/en/perf\_train\_gpu\_many#data-parallelism) on Hugging Face for more details on these strategies. Some of the key differences include:

1. DDP is generally faster than DP because it has to communicate less data.
2. With DP, GPU 0 does the bulk of the work, while with DDP, the work is distributed more evenly across all GPUs.
3. DDP allows for training across multiple machines, while DP is limited to a single machine.

In short, **DDP is generally recommended**. You can use DDP by running your normal training scripts with `torchrun` or `accelerate`. For example, if you have a script called `train_script.py`, you can run it with DDP using the following command:

Via `torchrun`

\* [torchrun documentation](https://pytorch.org/docs/stable/elastic/run.html)

```
torchrun --nproc_per_node=4 train_script.py
```

Via `accelerate`

\* [accelerate documentation](https://huggingface.co/docs/accelerate/en/index)

```
accelerate launch --num_processes 4 train_script.py
```

## Note

When performing distributed training, you have to wrap your code in a `main` function and call it with `if \_\_name\_\_ == "\_\_main\_\_":`. This is because each process will run the entire script, so you don't want to run the same code multiple times. Here is an example of how to do this:

```
from sentence_transformers import SentenceTransformer,
SentenceTransformerTrainingArguments, SentenceTransformerTrainer

# Other imports here


def main():

    # Your training code here
```

```
if __name__ == "__main__":  
    main()
```

### Note

When using an [Evaluator](../training\_overview.html#evaluator), the evaluator only runs on the first device unlike the training and evaluation datasets, which are shared across all devices.

## ## Comparison of

The following table shows the speedup of DDP over DP and no parallelism given a certain hardware setup.

- \* Hardware: a `p3.8xlarge` AWS instance, i.e. 4x V100 GPUs
- \* Model being trained: [microsoft/mpnet-base](https://huggingface.co/microsoft/mpnet-base) (133M parameters)

\* Maximum sequence length: 384 (following [all-mpnet-base-v2](https://huggingface.co/sentence-transformers/all-mpnet-base-v2))

- \* Training datasets: MultiNLI, SNLI and STSB (note: these have short texts)

Losses: [\[`SoftmaxLoss`\]\(../../package\\_reference/sentence\\_transformer/losses.html#sentence\\_transformers.Losses\)](https://pytorch.org/docs/stable/nn.functional.html#torch.nn.functional.softmax)

osses.SoftmaxLoss "sentence\_transformers.losses.SoftmaxLoss") for MultiNLI and SNLI,  
[`CosineSimilarityLoss`](../package\_reference/sentence\_transformer/losses.html#sentence\_transformers.losses.CosineSimilarityLoss "sentence\_transformers.losses.CosineSimilarityLoss") for STSB

\* Batch size per device: 32

Strategy | Launcher | Samples per Second

---|---|---

No Parallelism | `CUDA\_VISIBLE\_DEVICES=0 python train\_script.py` | 2724

Data Parallel (DP) | `python train\_script.py` (DP is used by default when launching a script with  
`python`) | 3675 (1.349x speedup)

**\*\*Distributed Data Parallel (DDP)\*\*** | `torchrun --nproc\_per\_node=4 train\_script.py` or `accelerate  
launch --num\_processes 4 train\_script.py` | **\*\*6980 (2.562x speedup)\*\***

## FSDP

Fully Sharded Data Parallelism (FSDP) is another distributed training strategy that is not fully supported by Sentence Transformers. It is a more advanced version of DDP that is particularly useful for very large models. Note that in the previous comparison, FSDP reaches 5782 samples per second (2.122x speedup), i.e. **\*\*worse than DDP\*\***. FSDP only makes sense with very large models. If you want to use FSDP with Sentence Transformers, you have to be aware of the following limitations:

\* You can't use the `evaluator` functionality with FSDP.

\* You have to save the trained model with

``trainer.accelerator.state.fsdps_plugin.set_state_dict_type("FULL_STATE_DICT")`` followed with ``trainer.save_model("output")``.

\* You have to use ``fsdp=["full_shard", "auto_wrap"]`` and ``fsdp_config={"transformer_layer_cls_to_wrap": "BertLayer"}`` in your ``SentenceTransformerTrainingArguments``, where ``BertLayer`` is the repeated layer in the encoder that houses the multi-head attention and feed-forward layers, so e.g. ``BertLayer`` or ``MPNetLayer``.

Read the [FSDP

documentation](https://huggingface.co/docs/accelerate/en/usage\_guides/fsdp) by

Accelerate for more details.

[ Previous](../examples/training/hpo/README.html "Hyperparameter Optimization") [Next ](../cross\_encoder/usage/usage.html "Usage")

\* \* \*

(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).