

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

## Getting Started

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

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

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

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

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

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

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

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

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

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

## Sentence Transformer

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

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

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

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

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

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

\* [Multi-Process / Multi-GPU

Encoding](../../applications/computing-embeddings/README.html#multi-process-multi-gpu-encoding)

\* [Semantic Textual

Similarity](../../docs/sentence\_transformer/usage/semantic\_textual\_similarity.html)

\* [Similarity

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

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

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

\* [Symmetric vs. Asymmetric Semantic

Search](../../applications/semantic-search/README.html#symmetric-vs-asymmetric-semantic-search)

\* [Manual

Implementation](../../applications/semantic-search/README.html#manual-implementation)

\* [Optimized

Implementation](../../applications/semantic-search/README.html#optimized-implementation)

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

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

\* [Approximate Nearest

Neighbor](../../applications/semantic-search/README.html#approximate-nearest-neighbor)

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

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

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

\* [Retrieve & Re-Rank

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

\* [Retrieval: Bi-Encoder](../../applications/retrieve\_rerank/README.html#retrieval-bi-encoder)

\* [Re-Ranker:

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

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

\* [Pre-trained Bi-Encoders

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

\* [Pre-trained Cross-Encoders

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

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

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

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

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

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

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

\*

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

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

\* [Margin Based

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

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

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

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

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

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

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

\* [Binary

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

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

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

- \* [\[Demo\]\(../../applications/embedding-quantization/README.html#demo\)](#)
- \* [\[Try it yourself\]\(../../applications/embedding-quantization/README.html#try-it-yourself\)](#)
- \* [\[Speeding up Inference\]\(../../docs/sentence\\_transformer/usage/efficiency.html\)](#)
- \* [\[PyTorch\]\(../../docs/sentence\\_transformer/usage/efficiency.html#pytorch\)](#)
- \* [\[ONNX\]\(../../docs/sentence\\_transformer/usage/efficiency.html#onnx\)](#)
- \* [\[OpenVINO\]\(../../docs/sentence\\_transformer/usage/efficiency.html#openvino\)](#)
- \* [\[Benchmarks\]\(../../docs/sentence\\_transformer/usage/efficiency.html#benchmarks\)](#)
- \* [\[Creating Custom Models\]\(../../docs/sentence\\_transformer/usage/custom\\_models.html\)](#)

[\\* \[Structure of Sentence Transformer Models\]\(../../docs/sentence\\_transformer/usage/custom\\_models.html#structure-of-sentence-transformer-models\)](#)

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

\* [Loss Function](../../docs/sentence\_transformer/training\_overview.html#loss-function)

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

\* [Evaluator](../../docs/sentence\_transformer/training\_overview.html#evaluator)

\* [Trainer](../../docs/sentence\_transformer/training\_overview.html#trainer)

\* [Callbacks](../../docs/sentence\_transformer/training\_overview.html#callbacks)

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

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

\* [Best Base Embedding Models](../../docs/sentence\_transformer/training\_overview.html#best-base-embedding-models)

- \* [Dataset Overview](../../docs/sentence\_transformer/dataset\_overview.html)
- \* [Datasets on the Hugging Face Hub](../../docs/sentence\_transformer/dataset\_overview.html#datasets-on-the-hugging-face-hub)
- \* [Pre-existing Datasets](../../docs/sentence\_transformer/dataset\_overview.html#pre-existing-datasets)
- \* [Loss Overview](../../docs/sentence\_transformer/loss\_overview.html)
- \* [Loss modifiers](../../docs/sentence\_transformer/loss\_overview.html#loss-modifiers)
- \* [Distillation](../../docs/sentence\_transformer/loss\_overview.html#distillation)
- \* [Commonly used Loss Functions](../../docs/sentence\_transformer/loss\_overview.html#commonly-used-loss-functions)
- \* [Custom Loss Functions](../../docs/sentence\_transformer/loss\_overview.html#custom-loss-functions)
- \* [Training Examples](../../docs/sentence\_transformer/training/examples.html)
- \* [Semantic Textual Similarity](../sts/README.html)
- \* [Training data](../sts/README.html#training-data)
- \* [Loss Function](../sts/README.html#loss-function)
- \* [Natural Language Inference](../nli/README.html)
- \* [Data](../nli/README.html#data)
- \* [SoftmaxLoss](../nli/README.html#softmaxloss)
- \* [MultipleNegativesRankingLoss](../nli/README.html#multiplenegativesrankingloss)
- \* [Paraphrase Data](../paraphrases/README.html)
- \* [Pre-Trained Models](../paraphrases/README.html#pre-trained-models)
- \* [Quora Duplicate Questions](../quora\_duplicate\_questions/README.html)
- \* [Training](../quora\_duplicate\_questions/README.html#training)
- \* [MultipleNegativesRankingLoss](../quora\_duplicate\_questions/README.html#multiplenegativesrankingloss)

- \* [Pretrained Models](../quora\_duplicate\_questions/README.html#pretrained-models)
- \* [MS MARCO](../ms\_marco/README.html)
- \* [Bi-Encoder](../ms\_marco/README.html#bi-encoder)
- \* [Matryoshka Embeddings](../matryoshka/README.html)
- \* [Use Cases](../matryoshka/README.html#use-cases)
- \* [Results](../matryoshka/README.html#results)
- \* [Training](../matryoshka/README.html#training)
- \* [Inference](../matryoshka/README.html#inference)
- \* [Code Examples](../matryoshka/README.html#code-examples)
- \* Adaptive Layers
  - \* Use Cases
  - \* Results
  - \* Training
  - \* Inference
  - \* Code Examples
- \* [Multilingual Models](../multilingual/README.html)
- \* [Extend your own models](../multilingual/README.html#extend-your-own-models)
- \* [Training](../multilingual/README.html#training)
- \* [Datasets](../multilingual/README.html#datasets)
- \* [Sources for Training Data](../multilingual/README.html#sources-for-training-data)
- \* [Evaluation](../multilingual/README.html#evaluation)
- \* [Available Pre-trained Models](../multilingual/README.html#available-pre-trained-models)
- \* [Usage](../multilingual/README.html#usage)
- \* [Performance](../multilingual/README.html#performance)
- \* [Citation](../multilingual/README.html#citation)
- \* [Model Distillation](../distillation/README.html)
- \* [Knowledge Distillation](../distillation/README.html#knowledge-distillation)

- \* [Speed - Performance Trade-Off](../distillation/README.html#speed-performance-trade-off)
- \* [Dimensionality Reduction](../distillation/README.html#dimensionality-reduction)
- \* [Quantization](../distillation/README.html#quantization)
- \* [Augmented SBERT](../data\_augmentation/README.html)
- \* [Motivation](../data\_augmentation/README.html#motivation)
  - \* [Extend to your own datasets](../data\_augmentation/README.html#extend-to-your-own-datasets)
  - \* [Methodology](../data\_augmentation/README.html#methodology)
    - \* [Scenario 1: Limited or small annotated datasets (few labeled sentence-pairs)](../data\_augmentation/README.html#scenario-1-limited-or-small-annotated-dataset-s-few-labeled-sentence-pairs)
    - \* [Scenario 2: No annotated datasets (Only unlabeled sentence-pairs)](../data\_augmentation/README.html#scenario-2-no-annotated-datasets-only-unlabeled-sentence-pairs)
  - \* [Training](../data\_augmentation/README.html#training)
  - \* [Citation](../data\_augmentation/README.html#citation)
- \* [Training with Prompts](../prompts/README.html)
  - \* [What are Prompts?](../prompts/README.html#what-are-prompts)
    - \* [Why would we train with Prompts?](../prompts/README.html#why-would-we-train-with-prompts)
    - \* [How do we train with Prompts?](../prompts/README.html#how-do-we-train-with-prompts)
- \* [Training with PEFT Adapters](../peft/README.html)
  - \* [Compatibility Methods](../peft/README.html#compatibility-methods)
  - \* [Adding a New Adapter](../peft/README.html#adding-a-new-adapter)
  - \* [Loading a Pretrained Adapter](../peft/README.html#loading-a-pretrained-adapter)
  - \* [Training Script](../peft/README.html#training-script)
- \* [Unsupervised Learning](../unsupervised\_learning/README.html)



\* [TSDAE](../../unsupervised\_learning/README.html#tsdae)

\* [SimCSE](../../unsupervised\_learning/README.html#simcse)

\* [CT](../../unsupervised\_learning/README.html#ct)

\* [CT (In-Batch Negative Sampling)](../../unsupervised\_learning/README.html#ct-in-batch-negative-sampling)

\* [Masked Language Model (MLM)](../../unsupervised\_learning/README.html#masked-language-model-mlm)

\* [GenQ](../../unsupervised\_learning/README.html#genq)

\* [GPL](../../unsupervised\_learning/README.html#gpl)

\* [Performance Comparison](../../unsupervised\_learning/README.html#performance-comparison)

\* [Domain Adaptation](../../domain\_adaptation/README.html)

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

\* [Adaptive Pre-Training](../../domain\_adaptation/README.html#adaptive-pre-training)

\* [GPL: Generative Pseudo-Labeling](../../domain\_adaptation/README.html#gpl-generative-pseudo-labeling)

\* [Hyperparameter Optimization](../hpo/README.html)

\* [HPO Components](../hpo/README.html#hpo-components)

\* [Putting It All Together](../hpo/README.html#putting-it-all-together)

\* [Example Scripts](../hpo/README.html#example-scripts)

\* [Distributed Training](../../docs/sentence\_transformer/training/distributed.html)

\* [Comparison](../../docs/sentence\_transformer/training/distributed.html#comparison)

\* [FSDP](../../docs/sentence\_transformer/training/distributed.html#fsdp)

Cross Encoder

- \* [Usage](../../docs/cross\_encoder/usage/usage.html)
- \* [Retrieve & Re-Rank](../../applications/retrieve\_rerank/README.html)
  - \* [Retrieve & Re-Rank Pipeline](../../applications/retrieve\_rerank/README.html#retrieve-re-rank-pipeline)
  - \* [Retrieval: Bi-Encoder](../../applications/retrieve\_rerank/README.html#retrieval-bi-encoder)
    - \* [Re-Ranker: Cross-Encoder](../../applications/retrieve\_rerank/README.html#re-ranker-cross-encoder)
  - \* [Example Scripts](../../applications/retrieve\_rerank/README.html#example-scripts)
    - \* [Pre-trained Bi-Encoders (Retrieval)](../../applications/retrieve\_rerank/README.html#pre-trained-bi-encoders-retrieval)
    - \* [Pre-trained Cross-Encoders (Re-Ranker)](../../applications/retrieve\_rerank/README.html#pre-trained-cross-encoders-re-ranker)
- \* [Pretrained Models](../../docs/cross\_encoder/pretrained\_models.html)
  - \* [MS MARCO](../../docs/cross\_encoder/pretrained\_models.html#ms-marco)
  - \* [SQuAD (QNLI)](../../docs/cross\_encoder/pretrained\_models.html#squad-qnli)
  - \* [STSbenchmark](../../docs/cross\_encoder/pretrained\_models.html#stsbenchmark)
    - \* [Quora Duplicate Questions](../../docs/cross\_encoder/pretrained\_models.html#quora-duplicate-questions)
  - \* [NLI](../../docs/cross\_encoder/pretrained\_models.html#nli)
  - \* [Community Models](../../docs/cross\_encoder/pretrained\_models.html#community-models)
- \* [Training Overview](../../docs/cross\_encoder/training\_overview.html)
- \* [Training Examples](../../docs/cross\_encoder/training/examples.html)
- \* [MS MARCO](../ms\_marco/cross\_encoder\_README.html)
  - \* [Cross-Encoder](../ms\_marco/cross\_encoder\_README.html#cross-encoder)
    - \* [Cross-Encoder Knowledge Distillation](../ms\_marco/cross\_encoder\_README.html#cross-encoder-knowledge-distillation)

## Package Reference

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

\*

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

\*

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

\*

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

\*

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

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

\*

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

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

\*

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

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

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

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

\*

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

\*

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

osingautoencoderloss)

\*

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

\*

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

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

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

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

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

\*

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

\*

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

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

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

\*

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

allelsentencesdataset)

\*

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

\*

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

\*

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

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

\*

[Main

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

\*

[Further

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

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

\*

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

\*

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

\*

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

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

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



- \* [CrossEncoder](../../docs/package\_reference/cross\_encoder/cross\_encoder.html#id1)
- \* [Training Inputs](../../docs/package\_reference/cross\_encoder/cross\_encoder.html#training-inputs)
- \* [Evaluation](../../docs/package\_reference/cross\_encoder/evaluation.html)
- \* [CEBinaryAccuracyEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cebinaryaccuracyevaluator)
- \* [CEBinaryClassificationEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cebinaryclassificationevaluator)
- \* [CECorrelationEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cecorrelationevaluator)
- \* [CEF1Evaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cef1evaluator)
- \* [CESoftmaxAccuracyEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cesoftmaxaccuracyevaluator)
- \* [CERerankingEvaluator](../../docs/package\_reference/cross\_encoder/evaluation.html#cererankingevaluator)
- \* [util](../../docs/package\_reference/util.html)
- \* [Helper Functions](../../docs/package\_reference/util.html#module-sentence\_transformers.util)
- \* [community\_detection()](../../docs/package\_reference/util.html#sentence\_transformers.util.community\_detection)
- \* [http\_get()](../../docs/package\_reference/util.html#sentence\_transformers.util.http\_get)

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

[Model

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

\*

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

\*

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

\*

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

nsformers.backend.export\_static\_quantized\_openvino\_model)

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

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

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

\*

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

\*

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

\*

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

\*

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

\*

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

\*

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

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

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

\* [Training Examples](../../docs/sentence\_transformer/training/examples.html)

\* Adaptive Layers

\*

[

Edit

on

GitHub]([https://github.com/UKPLab/sentence-transformers/blob/master/examples/training/adaptive\\_layer/README.md](https://github.com/UKPLab/sentence-transformers/blob/master/examples/training/adaptive_layer/README.md))

\* \* \*

# Adaptive Layers

Embedding models are often encoder models with numerous layers, such as 12 (e.g. [all-mpnet-base-v2](<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>)) or 6 (e.g. [all-MiniLM-L6-v2](<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>)). To get embeddings, every single one of these layers must be traversed. The [2D Matryoshka Sentence Embeddings](<https://arxiv.org/abs/2402.14776v1>) (2DMSE) preprint revisits this concept by proposing an approach to train embedding models that will perform well when only using a selection of all layers. This results in faster inference speeds at relatively low performance costs.

Note

The 2DMSE preprint was later updated and renamed to [ESE: Espresso Sentence Embeddings](<https://arxiv.org/abs/2402.14776>). The Sentence Transformers implementation of Adaptive Layers and Matryoshka2d (Adaptive Layer + Matryoshka Embeddings) are based on the initial preprint, and we accept contributions that implement the updated ESE paper.

## ## Use Cases

The 2DMSE paper mentions that using a few layers of a larger model trained using Adaptive Layers and Matryoshka Representation Learning can outperform a smaller model that was trained like a standard embedding model.

## ## Results

Let's look at the performance that we may be able to expect from an Adaptive Layer embedding model versus a regular embedding model. For this experiment, I have trained two models:

\*

[tomaarsen/mpnet-base-nli-adaptive-layer](https://huggingface.co/tomaarsen/mpnet-base-nli-adaptive-layer): Trained by running [adaptive\_layer\_nli.py](https://github.com/UKPLab/sentence-transformers/tree/master/examples/training/adaptive\_layer/adaptive\_layer\_nli.py) with [microsoft/mpnet-base](https://huggingface.co/microsoft/mpnet-base).

\* [tomaarsen/mpnet-base-nli](https://huggingface.co/tomaarsen/mpnet-base-nli): A near identical model as the former, but using only `MultipleNegativesRankingLoss` rather than `AdaptiveLayerLoss` on top of `MultipleNegativesRankingLoss`. I also use [microsoft/mpnet-base](https://huggingface.co/microsoft/mpnet-base) as the base model.

Both of these models were trained on the AllNLI dataset, which is a concatenation of the [SNLI](https://huggingface.co/datasets/snli) and [MultiNLI](https://huggingface.co/datasets/multi\_nli) datasets. I have

evaluated these models on the

[STSBenchmark](https://huggingface.co/datasets/mteb/stsbenchmark-sts) test set using multiple different embedding dimensions. The results are plotted in the following figure:

![[adaptive\_layer\_results]](https://huggingface.co/tomaarsen/mpnet-base-nli-adaptive-layer/resolve/main/adaptive\_layer\_results.png)

The first figure shows that the Adaptive Layer model stays much more performant when reducing the number of layers in the model. This is also clearly shown in the second figure, which displays that 80% of the performance is preserved when the number of layers is reduced all the way to 1.

Lastly, the third figure shows the expected speedup ratio for GPU & CPU devices in my tests. As you can see, removing half of the layers results in roughly a 2x speedup, at a cost of ~15% performance on STSB (~86 -> ~75 Spearman correlation). When removing even more layers, the performance benefit gets larger for CPUs, and between 5x and 10x speedups are very feasible with a 20% loss in performance.

## ## Training

Training with Adaptive Layer support is quite elementary: rather than applying some loss function on only the last layer, we also apply that same loss function on the pooled embeddings from previous layers. Additionally, we employ a KL-divergence loss that aims to make the embeddings of the non-last layers match that of the last layer. This can be seen as a fascinating

approach of [knowledge distillation](../distillation/README.html#knowledge-distillation), but with the last layer as the teacher model and the prior layers as the student models.

For example, with the 12-layer [microsoft/mpnet-base](https://huggingface.co/microsoft/mpnet-base), it will now be trained such that the model produces meaningful embeddings after each of the 12 layers.

```
from sentence_transformers import SentenceTransformer
from sentence_transformers.losses import CoSENTLoss, AdaptiveLayerLoss

model = SentenceTransformer("microsoft/mpnet-base")

base_loss = CoSENTLoss(model=model)
loss = AdaptiveLayerLoss(model=model, loss=base_loss)
```

\*

**\*\*Reference\*\***

:

```
[`AdaptiveLayerLoss`](../../docs/package_reference/sentence_transformer/losses.html#adaptivelayerloss)
```

Note that training with `AdaptiveLayerLoss` is not notably slower than without using it.

Additionally, this can be combined with the `MatryoshkaLoss` such that the resulting model can be reduced both in the number of layers, but also in the size of the output dimensions. See also the [Matryoshka Embeddings](../matryoshka/README.html) for more information on reducing output dimensions. In Sentence Transformers, the combination of these two losses is called `Matryoshka2dLoss`, and a shorthand is provided for simpler training.

```
from sentence_transformers import SentenceTransformer

from sentence_transformers.losses import CoSENTLoss, Matryoshka2dLoss

model = SentenceTransformer("microsoft/mpnet-base")

base_loss = CoSENTLoss(model=model)

loss = Matryoshka2dLoss(model=model, loss=base_loss, matryoshka_dims=[768, 512, 256, 128,
64])
```

\* **Reference** \*

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

## Inference

After a model has been trained using the Adaptive Layer loss, you can then truncate the model layers to your desired layer count. Note that this requires



doing a bit of surgery on the model itself, and each model is structured a bit differently, so the steps are slightly different depending on the model.

First of all, we will load the model & access the underlying `transformers` model like so:

```
from sentence_transformers import SentenceTransformer
```

```
model = SentenceTransformer("tomaarsen/mpnet-base-nli-adaptive-layer")
```

```
# We can access the underlying model with `model[0].auto_model`
```

```
print(model[0].auto_model)
```

```
MPNetModel(  
  (embeddings): MPNetEmbeddings(  
    (word_embeddings): Embedding(30527, 768, padding_idx=1)  
    (position_embeddings): Embedding(514, 768, padding_idx=1)  
    (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)  
    (dropout): Dropout(p=0.1, inplace=False)  
  )  
  (encoder): MPNetEncoder(  
    (layer): ModuleList(  
      (0-11): 12 x MPNetLayer(  

```

(attention): MPNetAttention(  
 (attn): MPNetSelfAttention(  
 (q): Linear(in\_features=768, out\_features=768, bias=True)  
 (k): Linear(in\_features=768, out\_features=768, bias=True)  
 (v): Linear(in\_features=768, out\_features=768, bias=True)  
 (o): Linear(in\_features=768, out\_features=768, bias=True)  
 (dropout): Dropout(p=0.1, inplace=False)  
 )  
 (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise\_affine=True)  
 (dropout): Dropout(p=0.1, inplace=False)  
)  
(intermediate): MPNetIntermediate(  
 (dense): Linear(in\_features=768, out\_features=3072, bias=True)  
 (intermediate\_act\_fn): GELUActivation()  
)  
(output): MPNetOutput(  
 (dense): Linear(in\_features=3072, out\_features=768, bias=True)  
 (LayerNorm): LayerNorm((768,), eps=1e-05, elementwise\_affine=True)  
 (dropout): Dropout(p=0.1, inplace=False)  
)  
)  
)  
(relative\_attention\_bias): Embedding(32, 12)  
)  
(pooler): MPNetPooler(  
 (dense): Linear(in\_features=768, out\_features=768, bias=True)  
 (activation): Tanh()

)

)

This output will differ depending on the model. We will look for the repeated layers in the encoder. For this MPNet model, this is stored under ``model[0].auto_model.encoder.layer``. Then we can slice the model to only keep the first few layers to speed up the model:

```
new_num_layers = 3
```

```
model[0].auto_model.encoder.layer = model[0].auto_model.encoder.layer[:new_num_layers]
```

Then we can run inference with it using

```
[`SentenceTransformers.encode`](../docs/package_reference/sentence_transformer/SentenceTransformer.html#sentence_transformers.SentenceTransformer.encode).
```

```
from sentence_transformers import SentenceTransformer
```

```
model = SentenceTransformer("tomaarsen/mpnet-base-nli-adaptive-layer")
```

```
new_num_layers = 3
```

```
model[0].auto_model.encoder.layer = model[0].auto_model.encoder.layer[:new_num_layers]
```

```

embeddings = model.encode(
    [
        "The weather is so nice!",
        "It's so sunny outside!",
        "He drove to the stadium.",
    ]
)

# Similarity of the first sentence with the other two
similarities = model.similarity(embeddings[0], embeddings[1:])

# => tensor([[0.7761, 0.1655]])

# compared to tensor([[ 0.7547, -0.0162]]) for the full model

```

As you can see, the similarity between the related sentences is much higher than the unrelated sentence, despite only using 3 layers. Feel free to copy this script locally, modify the `new\_num\_layers`, and observe the difference in similarities.

### ## Code Examples

See the following scripts as examples of how to apply the

[`AdaptiveLayerLoss`](../../docs/package\_reference/sentence\_transformer/losses.html#adaptive-layer-loss)

in practice:

\*\*[adaptive\_layer\_nli.py](https://github.com/UKPLab/sentence-transformers/tree/master/examples/tr

aining/adaptive\_layer/adaptive\_layer\_nli.py)\*\* : This example uses the  
`MultipleNegativesRankingLoss` with `AdaptiveLayerLoss` to train a strong embedding model using  
Natural Language Inference (NLI) data. It is an adaptation of the [NLI](../nli/README.html)  
documentation.

\*

\*\*[adaptive\_layer\_sts.py](https://github.com/UKPLab/sentence-transformers/tree/master/examples/training/adaptive\_layer/adaptive\_layer\_sts.py)\*\* : This example uses the CoSENTLoss with  
AdaptiveLayerLoss to train an embedding model on the training set of the STSBenchmark dataset.  
It is an adaptation of the [STS](../sts/README.html) documentation.

And the following scripts to see how to apply

[`Matryoshka2dLoss`](../../docs/package\_reference/sentence\_transformer/losses.html#matryoshka  
2dloss):

\*

\*\*[2d\_matryoshka\_nli.py](https://github.com/UKPLab/sentence-transformers/tree/master/examples/training/adaptive\_layer/./matryoshka/2d\_matryoshka\_nli.py)\*\* : This example uses the  
`MultipleNegativesRankingLoss` with `Matryoshka2dLoss` to train a strong embedding model using  
Natural Language Inference (NLI) data. It is an adaptation of the [NLI](../nli/README.html)  
documentation.

\*

\*\*[2d\_matryoshka\_sts.py](https://github.com/UKPLab/sentence-transformers/tree/master/examples/training/adaptive\_layer/./matryoshka/2d\_matryoshka\_sts.py)\*\* : This example uses the  
`CoSENTLoss` with `Matryoshka2dLoss` to train an embedding model on the training set of the  
STSBenchmark dataset. It is an adaptation of the [STS](../sts/README.html) documentation.

[ Previous](../matryoshka/README.html "Matryoshka Embeddings") [Next  
(../multilingual/README.html "Multilingual Models")

\* \* \*

(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).