

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

Getting Started

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

Sentence Transformer

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

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

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

[* \[Semantic Textual Similarity\]\(../../sentence_transformer/usage/semantic_textual_similarity.html\)](#)

[* \[Similarity Calculation\]\(../../sentence_transformer/usage/semantic_textual_similarity.html#similarity-calculation\)](#)

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

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

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

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

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

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

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

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

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

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

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

* [Retrieve & Re-Rank

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

* [Retrieval:

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

* [Re-Ranker:

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

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

* [Pre-trained Bi-Encoders

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

* [Pre-trained Cross-Encoders

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

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

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

* [Agglomerative

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

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

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

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

*

[`paraphrase_mining()`](../../examples/applications/paraphrase-mining/README.html#sentence_transformers.util.paraphrase_mining)

* [Translated Sentence

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

* [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-quantization)

* [Additional

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

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

* [Try it

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

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

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

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

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

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

* [Creating Custom Models](../../sentence_transformer/usage/custom_models.html)

* [Structure of Sentence Transformer

Models](../../sentence_transformer/usage/custom_models.html#structure-of-sentence-transformer-models)

* [Sentence Transformer Model from a Transformers

Model](../../sentence_transformer/usage/custom_models.html#sentence-transformer-model-from-a-transformers-model)

* [Pretrained Models](../../sentence_transformer/pretrained_models.html)

* [Original Models](../../sentence_transformer/pretrained_models.html#original-models)

* [Semantic Search

Models](../../sentence_transformer/pretrained_models.html#semantic-search-models)

* [Multi-QA Models](../../sentence_transformer/pretrained_models.html#multi-qa-models)

* [MSMARCO Passage

Models](../../sentence_transformer/pretrained_models.html#msmarco-passage-models)

* [Multilingual Models](../../sentence_transformer/pretrained_models.html#multilingual-models)

* [Semantic Similarity

Models](../../sentence_transformer/pretrained_models.html#semantic-similarity-models)

* [Bitext Mining](../../sentence_transformer/pretrained_models.html#bitext-mining)

* [Image & Text-Models](../../sentence_transformer/pretrained_models.html#image-text-models)

* [INSTRUCTOR models](../../sentence_transformer/pretrained_models.html#instructor-models)

* [Scientific Similarity

Models](../../sentence_transformer/pretrained_models.html#scientific-similarity-models)

* [Training Overview](../../sentence_transformer/training_overview.html)

* [Why Finetune?](../../sentence_transformer/training_overview.html#why-finetune)

* [Training Components](../../sentence_transformer/training_overview.html#training-components)

* [Dataset](../../sentence_transformer/training_overview.html#dataset)

* [Dataset Format](../../sentence_transformer/training_overview.html#dataset-format)

* [Loss Function](../../sentence_transformer/training_overview.html#loss-function)

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

* [\[Domain Adaptation\]\(../../examples/domain_adaptation/README.html\)](#)

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

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

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

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

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

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

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

* [\[Distributed Training\]\(../../sentence_transformer/training/distributed.html\)](#)

* [\[Comparison\]\(../../sentence_transformer/training/distributed.html#comparison\)](#)

* [\[FSDP\]\(../../sentence_transformer/training/distributed.html#fsdp\)](#)

Cross Encoder

* [\[Usage\]\(../../cross_encoder/usage/usage.html\)](#)

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

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

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

* [\[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](index.html)

* [SentenceTransformer](SentenceTransformer.html)

* [SentenceTransformer](SentenceTransformer.html#id1)

*

[SentenceTransformerModelCardData](SentenceTransformer.html#sentencetransformermodelcarddata)

* [SimilarityFunction](SentenceTransformer.html#similarityfunction)

* [Trainer](trainer.html)

* [SentenceTransformerTrainer](trainer.html#sentencetransformertrainer)

* [Training Arguments](training_args.html)

*

[SentenceTransformerTrainingArguments](training_args.html#sentencetransformertrainingarguments)

* [Losses](losses.html)

* [BatchAllTripletLoss](losses.html#batchalltripletloss)

* [BatchHardSoftMarginTripletLoss](losses.html#batchhardsoftmargintripletloss)

* [BatchHardTripletLoss](losses.html#batchhardtripletloss)

* [BatchSemiHardTripletLoss](losses.html#batchsemihardtripletloss)

* [ContrastiveLoss](losses.html#contrastiveloss)

* [OnlineContrastiveLoss](losses.html#onlinecontrastiveloss)

* [ContrastiveTensionLoss](losses.html#contrastivetensionloss)

*

[ContrastiveTensionLossInBatchNegatives](losses.html#contrastivetensionlossinbatchnegatives)

* [CoSENTLoss](losses.html#cosentloss)

* [AngleELoss](losses.html#angleloss)

* [CosineSimilarityLoss](losses.html#cosinesimilarityloss)

- * [DenoisingAutoEncoderLoss](losses.html#denoisingautoencoderloss)
- * [GISTEmbedLoss](losses.html#gistembedloss)
- * [CachedGISTEmbedLoss](losses.html#cachedgistembedloss)
- * [MSELoss](losses.html#mseloss)
- * [MarginMSELoss](losses.html#marginmseloss)
- * [MatryoshkaLoss](losses.html#matryoshkaloss)
- * [Matryoshka2dLoss](losses.html#matryoshka2dloss)
- * [AdaptiveLayerLoss](losses.html#adaptivelayerloss)
- * [MegaBatchMarginLoss](losses.html#megabatchmarginloss)
- * [MultipleNegativesRankingLoss](losses.html#multiplenegativesrankingloss)
- * [CachedMultipleNegativesRankingLoss](losses.html#cachedmultiplenegativesrankingloss)

*

[MultipleNegativesSymmetricRankingLoss](losses.html#multiplenegativessymmetricrankingloss)

*

[CachedMultipleNegativesSymmetricRankingLoss](losses.html#cachedmultiplenegativessymmetricrankingloss)

- * [SoftmaxLoss](losses.html#softmaxloss)
- * [TripletLoss](losses.html#tripletloss)
- * [Samplers](sampler.html)
- * [BatchSamplers](sampler.html#batchsamplers)
- * [MultiDatasetBatchSamplers](sampler.html#multidatasetbatchsamplers)
- * [Evaluation](evaluation.html)
- * [BinaryClassificationEvaluator](evaluation.html#binaryclassificationevaluator)
- * [EmbeddingSimilarityEvaluator](evaluation.html#embeddingsimilarityevaluator)
- * [InformationRetrievalEvaluator](evaluation.html#informationretrievalevaluator)
- * [NanoBEIREvaluator](evaluation.html#nanobeirevaluator)
- * [MSEEvaluator](evaluation.html#mseevaluator)

- * [ParaphraseMiningEvaluator](evaluation.html#paraphraseminingevaluator)
- * [RerankingEvaluator](evaluation.html#rerankingevaluator)
- * [SentenceEvaluator](evaluation.html#sentenceevaluator)
- * [SequentialEvaluator](evaluation.html#sequentialevaluator)
- * [TranslationEvaluator](evaluation.html#translationevaluator)
- * [TripletEvaluator](evaluation.html#tripletevaluator)
- * [Datasets](datasets.html)
- * [ParallelSentencesDataset](datasets.html#parallelsentencesdataset)
- * [SentenceLabelDataset](datasets.html#sentencelabeldataset)
- * [DenoisingAutoEncoderDataset](datasets.html#denoisingautoencoderdataset)
- * [NoDuplicatesDataLoader](datasets.html#noduplicatesdataloader)
- * Models
- * Main Classes
- * Further Classes
- * [quantization](quantization.html)

*

[`quantize_embeddings()`](quantization.html#sentence_transformers.quantization.quantize_embeddings)

*

[`semantic_search_faiss()`](quantization.html#sentence_transformers.quantization.semantic_search_faiss)

*

[`semantic_search_usearch()`](quantization.html#sentence_transformers.quantization.semantic_search_usearch)

- * [Cross Encoder](../cross_encoder/index.html)
- * [CrossEncoder](../cross_encoder/cross_encoder.html)
- * [CrossEncoder](../cross_encoder/cross_encoder.html#id1)

* [Training Inputs](../cross_encoder/cross_encoder.html#training-inputs)

* [Evaluation](../cross_encoder/evaluation.html)

* [CEBinaryAccuracyEvaluator](../cross_encoder/evaluation.html#cebinaryaccuracyevaluator)

*

[CEBinaryClassificationEvaluator](../cross_encoder/evaluation.html#cebinaryclassificationevaluator)

* [CECorrelationEvaluator](../cross_encoder/evaluation.html#cecorrelationevaluator)

* [CEF1Evaluator](../cross_encoder/evaluation.html#cef1evaluator)

*

[CESoftmaxAccuracyEvaluator](../cross_encoder/evaluation.html#cesoftmaxaccuracyevaluator)

* [CERerankingEvaluator](../cross_encoder/evaluation.html#cererankingevaluator)

* [util](../util.html)

* [Helper Functions](../util.html#module-sentence_transformers.util)

* [community_detection()](../util.html#sentence_transformers.util.community_detection)

* [http_get()](../util.html#sentence_transformers.util.http_get)

* [is_training_available()](../util.html#sentence_transformers.util.is_training_available)

* [mine_hard_negatives()](../util.html#sentence_transformers.util.mine_hard_negatives)

* [normalize_embeddings()](../util.html#sentence_transformers.util.normalize_embeddings)

* [paraphrase_mining()](../util.html#sentence_transformers.util.paraphrase_mining)

* [semantic_search()](../util.html#sentence_transformers.util.semantic_search)

* [truncate_embeddings()](../util.html#sentence_transformers.util.truncate_embeddings)

* [Model Optimization](../util.html#module-sentence_transformers.backend)

*

[export_dynamic_quantized_onnx_model()](../util.html#sentence_transformers.backend.export_dynamic_quantized_onnx_model)

*

[export_optimized_onnx_model()](../util.html#sentence_transformers.backend.export_optimized_onnx_model)

```
[`export_static_quantized_openvino_model()`](../util.html#sentence_transformers.backend.export_static_quantized_openvino_model)
```

- * [Similarity Metrics](../util.html#module-sentence_transformers.util)
- * [`cos_sim()`](../util.html#sentence_transformers.util.cos_sim)
- * [`dot_score()`](../util.html#sentence_transformers.util.dot_score)
- * [`euclidean_sim()`](../util.html#sentence_transformers.util.euclidean_sim)
- * [`manhattan_sim()`](../util.html#sentence_transformers.util.manhattan_sim)
- * [`pairwise_cos_sim()`](../util.html#sentence_transformers.util.pairwise_cos_sim)
- * [`pairwise_dot_score()`](../util.html#sentence_transformers.util.pairwise_dot_score)
- * [`pairwise_euclidean_sim()`](../util.html#sentence_transformers.util.pairwise_euclidean_sim)
- * [`pairwise_manhattan_sim()`](../util.html#sentence_transformers.util.pairwise_manhattan_sim)

___[Sentence Transformers](../../index.html)

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

* [Sentence Transformer](index.html)

* Models

* [Edit on

GitHub](https://github.com/UKPLab/sentence-transformers/blob/master/docs/package_reference/sentence_transformer/models.md)

* * *

Models¶

`sentence_transformers.models` defines different building blocks, that can be

used to create SentenceTransformer networks from scratch. For more details, see [Training Overview](../sentence_transformer/training_overview.html).

Main Classes

```
class sentence_transformers.models.Transformer(_model_name_or_path : str,
_max_seq_length : int | None = None, _model_args : dict[str, Any] | None = None,
_tokenizer_args : dict[str, Any] | None = None, _config_args : dict[str, Any] | None = None,
_cache_dir : str | None = None, _do_lower_case : bool = False, _tokenizer_name_or_path : str |
None = None, _backend : str =
'torch')[source](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers/models/Transformer.py#L31-L551)
```

Hugging Face AutoModel to generate token embeddings. Loads the correct class, e.g. BERT / RoBERTa etc.

Parameters:

model_name_or_path – Hugging Face models name (<<https://huggingface.co/models>>)

max_seq_length – Truncate any inputs longer than max_seq_length

model_args – Keyword arguments passed to the Hugging Face Transformers model

* **tokenizer_args** â€“ Keyword arguments passed to the Hugging Face Transformers tokenizer

* **config_args** â€“ Keyword arguments passed to the Hugging Face Transformers config

* **cache_dir** â€“ Cache dir for Hugging Face Transformers to store/load models

* **do_lower_case** â€“ If true, lowercases the input (independent if the model is cased or not)

* **tokenizer_name_or_path** â€“ Name or path of the tokenizer. When None, then model_name_or_path is used

* **backend** â€“ Backend used for model inference. Can be torch, onnx, or openvino. Default is torch.

Initializes internal Module state, shared by both nn.Module and ScriptModule.

```
_class      _sentence_transformers.models.Pooling(_word_embedding_dimension      :      int_,
_pooling_mode      :      str | None = None_, _pooling_mode_cls_token      :      bool = False_,
_pooling_mode_max_tokens      :      bool = False_, _pooling_mode_mean_tokens      :      bool = True_,
_pooling_mode_mean_sqrt_len_tokens      :      bool = False_, _pooling_mode_weightedmean_tokens      :
bool = False_, _pooling_mode_lasttoken      :      bool = False_, _include_prompt      :      bool =
True_)\[\[source\]\]\(https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transformers\models\Pooling.py#L11-L260\)if•
```

Performs pooling (max or mean) on the token embeddings.

Using pooling, it generates from a variable sized sentence a fixed sized sentence embedding. This layer also allows to use the CLS token if it is returned by the underlying word embedding model. You can concatenate multiple poolings together.

Parameters:

word_embedding_dimension – Dimensions for the word embeddings

pooling_mode – Either `“cls”`, `“lasttoken”`, `“max”`, `“mean”`, `“mean_sqrt_len_tokens”`, or `“weightedmean”`. If set, overwrites the other `pooling_mode_*` settings

pooling_mode_cls_token – Use the first token (CLS token) as text representations

pooling_mode_max_tokens – Use max in each dimension over all tokens.

pooling_mode_mean_tokens – Perform mean-pooling

pooling_mode_mean_sqrt_len_tokens – Perform mean-pooling, but divide by `sqrt(input_length)`.

pooling_mode_weightedmean_tokens – Perform (position) weighted mean pooling. See [SGPT: GPT Sentence Embeddings for Semantic Search](<https://arxiv.org/abs/2202.08904>).

`***pooling_mode_lasttoken**` â€

Perform last token pooling. See [SGPT: GPT Sentence Embeddings for Semantic Search](https://arxiv.org/abs/2202.08904) and [Text and Code Embeddings by Contrastive Pre-Training](https://arxiv.org/abs/2201.10005).

`***include_prompt**` â€ If set to false, the prompt tokens are not included in the pooling. This is useful for reproducing work that does not include the prompt tokens in the pooling like INSTRUCTOR, but otherwise not recommended.

Initializes internal Module state, shared by both nn.Module and ScriptModule.

```
_class _sentence_transformers.models.Dense(_in_features : int_, _out_features : int_, _bias : bool
= True_, _activation_function =Tanh(), _init_weight :
[Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch v2.5\))" | None =
None_, _init_bias : [Tensor](https://pytorch.org/docs/stable/tensors.html#torch.Tensor "(in PyTorch
v2.5\))" | None =
None_)[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transfor
mers\\models\\Dense.py#L14-L94)ïf•
```

Feed-forward function with activation function.

This layer takes a fixed-sized sentence embedding and passes it through a feed-forward layer. Can be used to generate deep averaging networks (DAN).

Parameters:

* **in_features** â€“ Size of the input dimension

* **out_features** â€“ Output size

* **bias** â€“ Add a bias vector

* **activation_function** â€“ Pytorch activation function applied on output

* **init_weight** â€“ Initial value for the matrix of the linear layer

* **init_bias** â€“ Initial value for the bias of the linear layer

Initializes internal Module state, shared by both nn.Module and ScriptModule.

Further Classesïf•

```
_class _sentence_transformers.models.Asym(_sub_modules : dict[str,
list[[Module]](https://pytorch.org/docs/stable/generated/torch.nn.Module.html#torch.nn.Module
"\(in PyTorch v2.5\)"))], _allow_empty_key : bool =
True_)[[source]](https://github.com/UKPLab/sentence-
transformers/blob/master/sentence\_transformers\models\Asym.py#L12-L132)ïf•
```

This model allows to create asymmetric SentenceTransformer models, that apply different models depending on the specified input key.

In the below example, we create two different Dense models for `~query~` and `~doc~`. Text that is passed as `{~query~: ~My query~}` will be passed along along the first Dense model, and text that will be passed as `{~doc~: ~My document~}` will use the other Dense model.

Note, that when you call `encode()`, that only inputs of the same type can be encoded. Mixed-Types cannot be encoded.

Example::

```
word_embedding_model = models.Transformer(model_name) pooling_model =  
models.Pooling(word_embedding_model.get_word_embedding_dimension()) asym_model  
= models.Asym({~query~:  
[models.Dense(word_embedding_model.get_word_embedding_dimension(), 128)],  
~doc~: [models.Dense(word_embedding_model.get_word_embedding_dimension(),  
128)]}) model = SentenceTransformer(modules=[word_embedding_model,  
pooling_model, asym_model])
```

```
model.encode([{~query~: ~Q1~}, {~query~: ~Q2~}])  
model.encode([{~doc~: ~Doc1~}, {~doc~: ~Doc2~}])
```

#You can train it with InputExample like this. Note, that the order must

always be the same: `train_example = InputExample(texts=[{"query": "Train query"}, {"doc": "Document"}], label=1)`

Parameters:

* **sub_modules** Dict in the format `str -> List[models]`. The models in the specified list will be applied for input marked with the respective key.

* **allow_empty_key** If true, inputs without a key can be processed. If false, an exception will be thrown if no key is specified.

```
_class _sentence_transformers.models.Bow(_vocab : list[str], _word_weights :  
dict[str, float] = {}, _unknown_word_weight : float = 1,  
_cumulative_term_frequency : bool =  
True_)[source](https://github.com/UKPLab/sentence-  
transformers/blob/master/sentence\_transformers\models\Bow.py#L16-L96)if•
```

Implements a Bag-of-Words (BoW) model to derive sentence embeddings.

A weighting can be added to allow the generation of tf-idf vectors. The output vector has the size of the vocab.

Initializes internal Module state, shared by both `nn.Module` and `ScriptModule`.

```
_class _sentence_transformers.models.CNN(_in_word_embedding_dimension : int, _out_channels : int = 256, _kernel_sizes : list[int] = [1, 3, 5], _stride_sizes : list[int] | None = None)[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transformers\\models\\CNN.py#L12-L88)if•
```

CNN-layer with multiple kernel-sizes over the word embeddings

Initializes internal Module state, shared by both nn.Module and ScriptModule.

```
_class _sentence_transformers.models.LSTM(_word_embedding_dimension : int, _hidden_dim : int, _num_layers : int = 1, _dropout : float = 0, _bidirectional : bool = True)[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transformers\\models\\LSTM.py#L12-L90)if•
```

Bidirectional LSTM running over word embeddings.

Initializes internal Module state, shared by both nn.Module and ScriptModule.

```
_class _sentence_transformers.models.Normalize[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transformers\\models\\Normalize.py#L7-L22)if•
```


This layer normalizes embeddings to unit length

Initializes internal Module state, shared by both nn.Module and ScriptModule.

```
_class      _sentence_transformers.models.StaticEmbedding(_tokenizer      :      Tokenizer      |
PreTrainedTokenizerFast_,      _embedding_weights      :      ndarray      |
[Tensor])(https://pytorch.org/docs/stable/tensors.html#torch.Tensor "\in PyTorch v2.5\") | None =
None_,      _embedding_dim      :      int      |      None      =      None_,      _**
kwargs_)[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transfo
rmers\models\StaticEmbedding.py#L18-L220)if•
```

Initializes the StaticEmbedding model given a tokenizer. The model is a simple embedding bag model that takes the mean of trained per-token embeddings to compute text embeddings.

Parameters:

* **tokenizer** (`_Tokenizer_` | `_PreTrainedTokenizerFast_`) â€” The tokenizer to be used. Must be a fast tokenizer from `transformers` or `tokenizers`.

* **embedding_weights** (`_np.ndarray_` | `_torch.Tensor_`)(https://pytorch.org/docs/stable/tensors.html#torch.Tensor "\in PyTorch v2.5\") | `_None_`, `_optional_`) â€” Pre-trained embedding weights. Defaults to None.

embedding_dim (`_int_` | `__None__`, `__optional__`) â€” Dimension of the embeddings. Required if `embedding_weights` is not provided. Defaults to `None`.

Example:

```
from sentence_transformers import SentenceTransformer

from sentence_transformers.models import StaticEmbedding

from tokenizers import Tokenizer

# Pre-distilled embeddings:

static_embedding = StaticEmbedding.from_model2vec("minishlab/potion-base-8M")

# or distill your own embeddings:

static_embedding = StaticEmbedding.from_distillation("BAAI/bge-base-en-v1.5", device="cuda")

# or start with randomized embeddings:

tokenizer = Tokenizer.from_pretrained("FacebookAI/xlm-roberta-base")

static_embedding = StaticEmbedding(tokenizer, embedding_dim=512)

model = SentenceTransformer(modules=[static_embedding])

embeddings = model.encode(["What are Pandas?", "The giant panda, also known as the panda bear or simply the panda, is a bear native to south central China."])

similarity = model.similarity(embeddings[0], embeddings[1])

# tensor([[0.8093]]) (If you use potion-base-8M)

# tensor([[0.6234]]) (If you use the distillation method)
```

tensor([[-0.0693]]) (For example, if you use randomized embeddings)

Raises:

ValueError – If the tokenizer is not a fast tokenizer.

ValueError – If neither `embedding_weights` nor `embedding_dim` is provided.

```
_classmethod _from_distillation(_model_name : str, _vocabulary : list[str] | None = None, _device : str | None = None, _pca_dims : int | None = 256, _apply_zipf : bool = True, _use_subword : bool = True) -> StaticEmbedding[[source]](https://github.com/UKPLab/sentence-transformers/blob/master/sentence\_transformers\models\StaticEmbedding.py#L137-L189)if•
```

Creates a `StaticEmbedding` instance from a distillation process using the `model2vec` package.

Parameters:

model_name (_str_) – The name of the model to distill.

* **vocabulary** (_list_ _[_str_ _]_|_None_ __, __optional_) â€“ A list of vocabulary words to use.

Defaults to None.

* **device** (_str_) â€“ The device to run the distillation on (e.g., `~cpu`TM, `~cuda`TM). If not specified, the strongest device is automatically detected. Defaults to None.

* **pca_dims** (_int_ _|_None_ __, __optional_) â€“ The number of dimensions for PCA reduction.

Defaults to 256.

* **apply_zipf** (_bool_) â€“ Whether to apply ZipfTMs law during distillation. Defaults to True.

* **use_subword** (_bool_) â€“ Whether to use subword tokenization. Defaults to True.

Returns:

An instance of StaticEmbedding initialized with the distilled modelTMs

tokenizer and embedding weights.

Return type:

StaticEmbedding

Raises:

ImportError â€“ If the model2vec package is not installed.

`_classmethod _from_model2vec(_model_id_or_path : str_) ->`

`StaticEmbedding`[[source]]([https://github.com/UKPLab/sentence-](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\models\StaticEmbedding.py#L191-L220)

`transformers/blob/master/sentence_transformers\models\StaticEmbedding.py#L191-L220`)if•

Create a StaticEmbedding instance from a model2vec model. This method loads a pre-trained model2vec model and extracts the embedding weights and tokenizer to create a StaticEmbedding instance.

Parameters:

model_id_or_path (_str_) â€“ The identifier or path to the pre-trained model2vec model.

Returns:

An instance of StaticEmbedding initialized with the tokenizer and embedding weights

the model2vec model.

Return type:

StaticEmbedding

Raises:

ImportError â€” If the model2vec package is not installed.

`_class`

`_sentence_transformers.models.WeightedLayerPooling(_word_embedding_dimension_`

`, _num_hidden_layers : int = 12_, _layer_start : int = 4_, _layer_weights`

`=None_)`[\[\[source\]\]\(https://github.com/UKPLab/sentence-](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\models\WeightedLayerPooling.py#L12-L70)

[transformers/blob/master/sentence_transformers\models\WeightedLayerPooling.py#L12-L70\)](https://github.com/UKPLab/sentence-transformers/blob/master/sentence_transformers\models\WeightedLayerPooling.py#L12-L70)[if](#)•

Token embeddings are weighted mean of their different hidden layer representations

Initializes internal Module state, shared by both nn.Module and ScriptModule.

```
_class _sentence_transformers.models.WordEmbeddings(_tokenizer :  
WordTokenizer_, _embedding_weights_ , _update_embeddings : bool = False_,  
_max_seq_length : int =  
1000000_)[[source]](https://github.com/UKPLab/sentence-  
transformers/blob/master/sentence\_transformers\\models\\WordEmbeddings.py#L22-L179)if•
```

Initializes internal Module state, shared by both nn.Module and ScriptModule.

```
_class _sentence_transformers.models.WordWeights(_vocab : list[str]_,  
_word_weights : dict[str, float]_, _unknown_word_weight : float =  
1_)[[source]](https://github.com/UKPLab/sentence-  
transformers/blob/master/sentence\_transformers\\models\\WordWeights.py#L13-L80)if•
```

This model can weight word embeddings, for example, with idf-values.

Initializes the WordWeights class.

Parameters:

* **vocab** (`_List_` [`__str_` `_`]) â€“ Vocabulary of the tokenizer.

* **word_weights** (`_Dict_` [`__str_` `_`, `__float_` `_`]) â€“ Mapping of tokens to a float weight value. Word embeddings are multiplied by this float value. Tokens in `word_weights` must not be equal to the vocab (can contain more or less values).

* **unknown_word_weight** (`_float_` `_`, `__optional_`) â€“ Weight for words in vocab that do not appear in the `word_weights` lookup. These can be, for example, rare words in the vocab where no weight exists. Defaults to 1.

[[Previous](#)](datasets.html "Datasets") [[Next](#)](quantization.html "quantization")

* * *

(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).