[![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-tra
nsformer-model)
* [Calculating

ings) [Prompt

Embeddings] (.../../examples/applications/computing-embeddings/README.html # calculating-embeddings) + (.../examples/applications/computing-embeddings/README.html # calculating-embeddings/README.html # calculating-embeddings/README

Templates](../../examples/applications/computing-embeddings/README.html#prompt-templates)

Sequence [Input Length](../../examples/applications/computing-embeddings/README.html#id1) Multi-GPU [Multi-Process / Encoding](../../examples/applications/computing-embeddings/README.html#multi-process-multi-gp u-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 Asymmetric Semantic VS. Search](../../examples/applications/semantic-search/README.html#symmetric-vs-asymmetric-sema ntic-search) [Manual Implementation](../../examples/applications/semantic-search/README.html#manual-implementation ) [Optimized Implementation](../../examples/applications/semantic-search/README.html#optimized-implementati on) [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-neighbo r) & [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-retrie val) [Pre-trained Cross-Encoders (Re-Ranker)](../../examples/applications/retrieve rerank/README.html#pre-trained-cross-encodersre-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\_tra nsformers.util.paraphrase\_mining) [Translated Sentence

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

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

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

)

ation)

s)

- \* [OpenVINO](../sentence\_transformer/usage/efficiency.html#openvino)
- \* [Benchmarks](../sentence\_transformer/usage/efficiency.html#benchmarks)
- \* [Creating Custom Models](../sentence transformer/usage/custom models.html)
  - [Structure Transformer of Sentence

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-ga-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)

[Multiple Negatives Ranking Loss] (../../examples/training/nli/README.html # multiple negatives ranking loss) (.../examples/training/nli/README.html # multiple negatives ranking loss) (.../examples/training/nli/README.html # multiple negatives ranking loss) (.../examples/training/nli/README.html # multiple negatives/training/nli/README.html # mul

ss)

- \* [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#m ultiplenegativesrankingloss)

\* [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/multilingua	al/RE	ADME.ht	tml#datase	ets)		
		*	[Sources		for	Training
Data](//examples/training/multilingual/READM	E.htn	nl#source	es-for-trair	ing-data	)	
* [Evaluation](//examples/training/multiling	ual/R	EADME.	html#eval	uation)		
		*	[A	vailable	!	Pre-trained
Models](//examples/training/multilingual/READ	DME.I	html#ava	ilable-pre-	trained-r	nodels)	
* [Usage](//examples/training/multilingual/l	READ	OME.htm	l#usage)			
* [Performance](//examples/training/multili	ngua	I/READIV	1E.html#pe	erforman	ce)	
* [Citation](//examples/training/multilingual	/REA	DME.htn	nl#citation	)		
* [Model Distillation](//examples/training/dis-	tillatio	on/READ	ME.html)			
* [Knowledge Distillation](//examples/traini	ng/di	stillation/	README.	.html#kn	owledge-	distillation)
	*		[Speed	-	Pe	erformance
Trade-Off](//examples/training/distillation/REA	DME	.html#sp	eed-perfoi	mance-t	rade-off)	
				*	[Dim	ensionality
Reduction](//examples/training/distillation/REA	ADME	.html#di	mensional	ity-reduc	tion)	
* [Quantization](//examples/training/distilla	tion/F	README	.html#qua	ntization	)	
* [Augmented SBERT](//examples/training/o	data_	augment	tation/REA	DME.htr	nl)	
* [Motivation](//examples/training/data_aug	gmen	tation/RE	EADME.ht	ml#motiv	ation)	
	*	[Exte	end	to	your	own
datasets](//examples/training/data_augmentat	ion/R	EADME.	.html#exte	nd-to-yo	ur-own-d	atasets)
* [Methodology](//examples/training/data_a	augm	entation/	/README	.html#me	ethodolog	jy)
* [Scenario 1: Limited	or	small a	annotated	datase	ets (fev	v labeled
sentence-pairs)](//examples/training/data_aug	ment	ation/RE	ADME.htr	nl#scena	rio-1-limi	ted-or-sm
all-annotated-datasets-few-labeled-sentence-pai	rs)					
* [Scenario 2: I	No	annotat	ted dat	asets	(Only	unlabeled
sentence-pairs)](//examples/training/data_aug	ment	ation/RE	ADME.htr	nl#scena	rio-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 with train we Prompts?](../../examples/training/prompts/README.html#why-would-we-train-with-prompts) [How do train with we 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) Pretrained [Loading а 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 Model Language (MLM)](../../examples/unsupervised\_learning/README.html#masked-language-model-mlm) \* [GenQ](../../examples/unsupervised\_learning/README.html#geng) \* [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 Unsupervised VS. Learning](../../examples/domain\_adaptation/README.html#domain-adaptation-vs-unsupervised-lear ning) \* [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](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)

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

[Pre-trained

**Bi-Encoders** 

[Pre-trained Cross-Encoders

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

- \* Pretrained Models
  - \* MS MARCO
  - \* SQuAD (QNLI)
  - \* STSbenchmark
  - \* Quora Duplicate Questions
  - \* NLI
  - \* Community Models
- \* [Training Overview](training\_overview.html)
- \* [Training Examples](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-knowled ge-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/SentenceTransformer.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\_arg s.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#batchsemihar dtripletloss)

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

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

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

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

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

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

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

\* [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)

ŧ

\*

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

\*

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

\*

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

r

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

•

[CachedMultipleNegativesSymmetricRankingLoss](../package\_reference/sentence\_transformer/loss es.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#multidatase tbatchsamplers)

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

•

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

\*

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

dingsimilarityevaluator) [InformationRetrievalEvaluator](../package\_reference/sentence\_transformer/evaluation.html#informa tionretrievalevaluator) [NanoBEIREvaluator](../package\_reference/sentence\_transformer/evaluation.html#nanobeirevaluat or) \* [MSEEvaluator](../package\_reference/sentence\_transformer/evaluation.html#mseevaluator) [ParaphraseMiningEvaluator](../package\_reference/sentence\_transformer/evaluation.html#paraphra seminingevaluator) [RerankingEvaluator](../package\_reference/sentence\_transformer/evaluation.html#rerankingevaluat or) [SentenceEvaluator](../package\_reference/sentence\_transformer/evaluation.html#sentenceevaluato r) [SequentialEvaluator](../package reference/sentence transformer/evaluation.html#sequentialevalua tor) [TranslationEvaluator](../package\_reference/sentence\_transformer/evaluation.html#translationevalu ator)

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

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

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

ncesdataset)

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

ataset)

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

ngautoencoderdataset)

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

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

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

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

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

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

.

\*

\* [Evaluation](../package\_reference/cross\_encoder/evaluation.html) [CEBinaryAccuracyEvaluator](../package\_reference/cross\_encoder/evaluation.html#cebinaryaccura cyevaluator) [CEBinaryClassificationEvaluator](../package\_reference/cross\_encoder/evaluation.html#cebinarycla ssificationevaluator) [CECorrelationEvaluator](../package reference/cross encoder/evaluation.html#cecorrelationevaluat or) \* [CEF1Evaluator](../package\_reference/cross\_encoder/evaluation.html#cef1evaluator) [CESoftmaxAccuracyEvaluator](../package\_reference/cross\_encoder/evaluation.html#cesoftmaxacc uracyevaluator) [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\_det ection) \* [`http\_get()`](../package\_reference/util.html#sentence\_transformers.util.http\_get)

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

able)

\*

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

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

[`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\_transformer s.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_sco
`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)  * Pretrained Models  * [ Edit or
GitHub](https://github.com/UKPLab/sentence-transformers/blob/master/docs/cross_encoder/pretrain
ed_models.md)
* *
<sup>‡</sup> Pretrained Modelsï <i>f</i> •
Ve have released various pre-trained Cross Encoder models via our [Cross
Encoder Hugging Face organization](https://huggingface.co/models?author=cross-

encoder). Additionally, numerous community CrossEncoder models have been publicly released on the Hugging Face Hub.

Each of these models can be easily downloaded and used like so:

from sentence\_transformers import CrossEncoder import torch

model = CrossEncoder("cross-encoder/ms-marco-MiniLM-L-6-v2",
default\_activation\_function=torch.nn.Sigmoid())

scores = model.predict([

("How many people live in Berlin?", "Berlin had a population of 3,520,031 registered inhabitants in an area of 891.82 square kilometers."),

("How many people live in Berlin?", "Berlin is well known for its museums."),

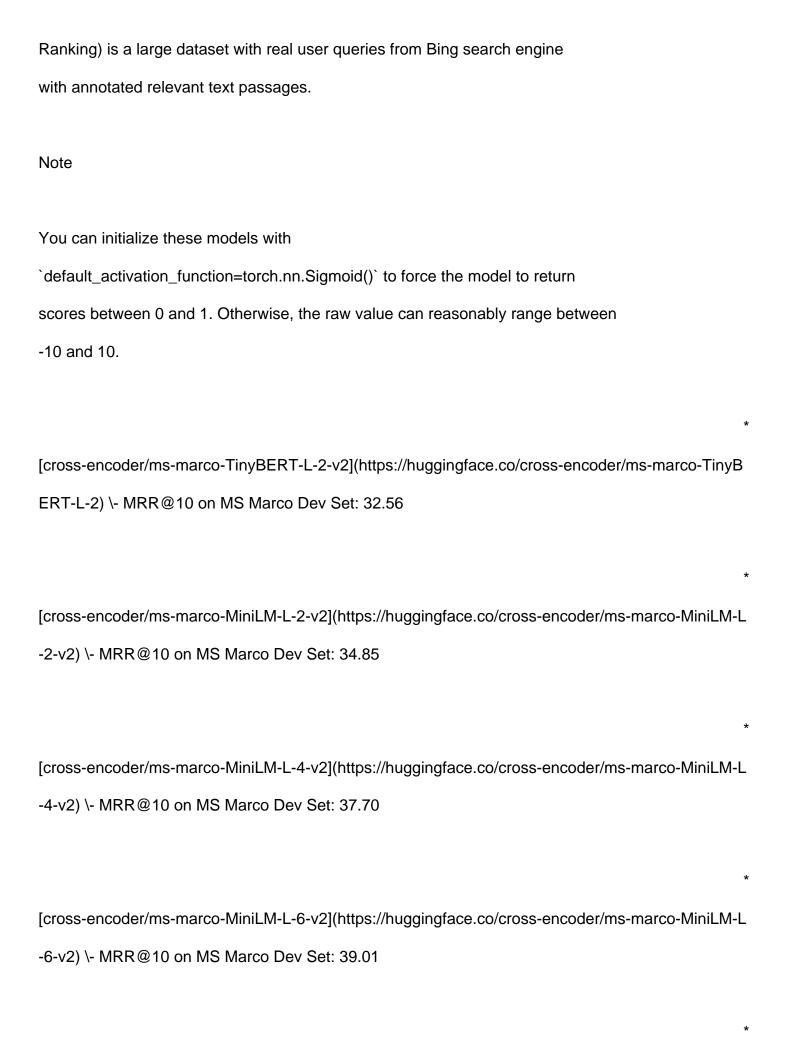
# => array([0.9998173, 0.01312432], dtype=float32)

Cross-Encoders require text pairs as inputs and output a score 0…1 (if the Sigmoid activation function is used). They do not work for individual sentences and they don't compute embeddings for individual texts.

## MS MARCOïf•

])

[MS MARCO Passage Retrieval](https://github.com/microsoft/MSMARCO-Passage-



[cross-encoder/ms-marco-MiniLM-L-12-v2](https://huggingface.co/cross-encoder/ms-marco-MiniLM-L-12-v2) \- MRR@10 on MS Marco Dev Set: 39.02

For details on the usage, see [Retrieve & Re-

Rank](../../examples/applications/retrieve\_rerank/README.html) or [MS MARCO Cross-Encoders](../pretrained-models/ce-msmarco.html).

## SQuAD (QNLI)ïf•

QNLI is based on the [SQuAD dataset](https://rajpurkar.github.io/SQuAD-explorer/) ([HF](https://huggingface.co/datasets/rajpurkar/squad)) and was introduced by the [GLUE Benchmark](https://arxiv.org/abs/1804.07461)
([HF](https://huggingface.co/datasets/nyu-mll/glue)). Given a passage from Wikipedia, annotators created questions that are answerable by that passage.

[cross-encoder/qnli-distilroberta-base](https://huggingface.co/cross-encoder/qnli-distilroberta-base)
\- Accuracy on QNLI dev set: 90.96

\* [cross-encoder/qnli-electra-base](https://huggingface.co/cross-encoder/qnli-electra-base) \Accuracy on QNLI dev set: 93.21

## STSbenchmarkïf•

The following models can be used like this:

from sentence\_transformers import CrossEncoder

model = CrossEncoder("cross-encoder/stsb-roberta-base")
scores = model.predict([("It's a wonderful day outside.", "It's so sunny today!"), ("It's a wonderful day outside.", "He drove to work earlier.")])

# => array([0.60443085, 0.00240758], dtype=float32)

They return a score 0â€l1 indicating the semantic similarity of the given sentence pair.

\* [cross-encoder/stsb-TinyBERT-L-4](https://huggingface.co/cross-encoder/stsb-TinyBERT-L-4) \STSbenchmark test performance: 85.50

[cross-encoder/stsb-distilroberta-base](https://huggingface.co/cross-encoder/stsb-distilroberta-base)
\- STSbenchmark test performance: 87.92

- \* [cross-encoder/stsb-roberta-base](https://huggingface.co/cross-encoder/stsb-roberta-base) \STSbenchmark test performance: 90.17
- \* [cross-encoder/stsb-roberta-large](https://huggingface.co/cross-encoder/stsb-roberta-large) \STSbenchmark test performance: 91.47

## Quora Duplicate Questionsïf•

These models have been trained on the [Quora duplicate questions dataset](https://huggingface.co/datasets/sentence-transformers/quoraduplicates). They can used like the STSb models and give a score 0â€11 indicating the probability that two questions are duplicate questions.

[cross-encoder/quora-distilroberta-base](https://huggingface.co/cross-encoder/quora-distilroberta-ba se) \- Average Precision dev set: 87.48

\* [cross-encoder/quora-roberta-base](https://huggingface.co/cross-encoder/quora-roberta-base) \-

Average Precision dev set: 87.80

\* [cross-encoder/quora-roberta-large](https://huggingface.co/cross-encoder/quora-roberta-large) \-

Average Precision dev set: 87.91

Note

The model don't work for question similarity. The question \_How to learn Java and How to learn Python will get a low score, as these guestions are not duplicates. For question similarity, the respective bi-encoder trained on the Quora dataset yields much more meaningful results.

## NLIïf•

Given two sentences, are these contradicting each other, entailing one the other or are these neutral? The following models were trained on the [SNLI](https://huggingface.co/datasets/stanfordnlp/snli) and

[MultiNLI](https://huggingface.co/datasets/nyu-mll/multi_nli) datasets.
* [cross-encoder/nli-deberta-v3-base](https://huggingface.co/cross-encoder/nli-deberta-v3-base) \- Accuracy on MNLI mismatched set: 90.04
* [cross-encoder/nli-deberta-base](https://huggingface.co/cross-encoder/nli-deberta-base) \- Accuracy on MNLI mismatched set: 88.08
* [cross-encoder/nli-deberta-v3-xsmall](https://huggingface.co/cross-encoder/nli-deberta-v3-xsmall) \- Accuracy on MNLI mismatched set: 87.77
* [cross-encoder/nli-deberta-v3-small](https://huggingface.co/cross-encoder/nli-deberta-v3-small) \- Accuracy on MNLI mismatched set: 87.55
* [cross-encoder/nli-roberta-base](https://huggingface.co/cross-encoder/nli-roberta-base) \- Accuracy on MNLI mismatched set: 87.47
* [cross-encoder/nli-MiniLM2-L6-H768](https://huggingface.co/cross-encoder/nli-MiniLM2-L6-H768)  \- Accuracy on MNLI mismatched set: 86.89
* [cross-encoder/nli-distilroberta-base](https://huggingface.co/cross-encoder/nli-distilroberta-base)  \- Accuracy on MNLI mismatched set: 83.98
from sentence_transformers import CrossEncoder

```
model = CrossEncoder("cross-encoder/nli-deberta-v3-base")
  scores = model.predict([
    ("A man is eating pizza", "A man eats something"),
      ("A black race car starts up in front of a crowd of people.", "A man is driving down a lonely
road."),
  ])
  # Convert scores to labels
  label_mapping = ["contradiction", "entailment", "neutral"]
  labels = [label_mapping[score_max] for score_max in scores.argmax(axis=1)]
  # => ['entailment', 'contradiction']
## Community Modelsif•
Some notable models from the Community include:
 * [BAAI/bge-reranker-base](https://huggingface.co/BAAI/bge-reranker-base)
 * [BAAI/bge-reranker-large](https://huggingface.co/BAAI/bge-reranker-large)
 * [BAAI/bge-reranker-v2-m3](https://huggingface.co/BAAI/bge-reranker-v2-m3)
 * [BAAI/bge-reranker-v2-gemma](https://huggingface.co/BAAI/bge-reranker-v2-gemma)
```

[BAAI/bge-reranker-v2-minicpm-layerwise](https://huggingface.co/BAAI/bge-reranker-v2-minicpm-la

yerwise)
* [jinaai/jina-reranker-v1-tiny-en](https://huggingface.co/jinaai/jina-reranker-v1-tiny-en)
* [jinaai/jina-reranker-v1-turbo-en](https://huggingface.co/jinaai/jina-reranker-v1-turbo-en)
[mixedbread-ai/mxbai-rerank-xsmall-v1](https://huggingface.co/mixedbread-ai/mxbai-rerank-xsmall-v1)
[mixedbread-ai/mxbai-rerank-base-v1](https://huggingface.co/mixedbread-ai/mxbai-rerank-base-v1)
[mixedbread-ai/mxbai-rerank-large-v1](https://huggingface.co/mixedbread-ai/mxbai-rerank-large-v1)
[maidalun1020/bce-reranker-base_v1](https://huggingface.co/maidalun1020/bce-reranker-base_v1)
[ Previous](usage/usage.html "Usage") [Next ](training_overview.html "Training Overview")
* * *
(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).