

MTEB: Massive Text Embedding Benchmark

Niklas Muennighoff¹, Nouamane Tazi¹, Loïc Magne¹, Nils Reimers^{2*}

¹Hugging Face ²cohere.ai

¹firstname@hf.co ²info@nils-reimers.de

Abstract

Text embeddings are commonly evaluated on a small set of datasets from a single task not covering their possible applications to other tasks. It is unclear whether state-of-the-art embeddings on semantic textual similarity (STS) can be equally well applied to other tasks like clustering or reranking. This makes progress in the field difficult to track, as various models are constantly being proposed without proper evaluation. To solve this problem, we introduce the Massive Text Embedding Benchmark (MTEB). MTEB spans 8 embedding tasks covering a total of 58 datasets and 112 languages. Through the benchmarking of 33 models on MTEB, we establish the most comprehensive benchmark of text embeddings to date. We find that no particular text embedding method dominates across all tasks. This suggests that the field has yet to converge on a universal text embedding method and scale it up sufficiently to provide state-of-the-art results on all embedding tasks. MTEB comes with open-source code and a public leaderboard at <https://github.com/embeddings-benchmark/mteb>.

1 Introduction

Natural language embeddings power a variety of use cases from clustering and topic representation (Aggarwal and Zhai, 2012; Angelov, 2020) to search systems and text mining (Huang et al., 2020; Zhu et al., 2021; Nayak, 2019) to feature representations for downstream models (Saharia et al., 2022; Borgeaud et al., 2022). Using generative language models or cross-encoders for these applications is often intractable, as they may require exponentially more computations (Reimers and Gurevych, 2019).

However, the evaluation regime of current text embedding models rarely covers the breadth of

their possible use cases. For example, SimCSE (Gao et al., 2021b) or SBERT (Reimers and Gurevych, 2019) solely evaluate on STS and classification tasks, leaving open questions about the transferability of the embedding models to search or clustering tasks. STS is known to poorly correlate with other real-world use cases (Neelakantan et al., 2022; Wang et al., 2021). Further, evaluating embedding methods on many tasks requires implementing multiple evaluation pipelines. Implementation details like pre-processing or hyperparameters may influence the results making it unclear whether performance improvements simply come from a favorable evaluation pipeline. This leads to the “blind” application of these models to new use cases in industry or requires incremental work to reevaluate them on different tasks.

The Massive Text Embedding Benchmark (MTEB) aims to provide clarity on how models perform on a variety of embedding tasks and thus serves as the gateway to finding universal text embeddings applicable to a variety of tasks. MTEB consists of 58 datasets covering 112 languages from 8 embedding tasks: Bitext mining, classification, clustering, pair classification, reranking, retrieval, STS and summarization. MTEB software is available open-source¹ enabling evaluation of any embedding model by adding less than 10 lines of code. Datasets and the MTEB leaderboard are available on the Hugging Face Hub².

We evaluate over 30 models on MTEB with additional speed and memory benchmarking to provide a holistic view of the state of text embedding models. We cover both models available open-source as well as models accessible via APIs, such as the OpenAI Embeddings endpoint. We find there to be no single best solution, with different models dom-

¹<https://github.com/embeddings-benchmark/mteb>

²<https://huggingface.co/spaces/mteb/leaderboard>

*Most of the work done while at Hugging Face. Correspondence to n.muennighoff@gmail.com.

inating different tasks. Our benchmarking sheds light on the weaknesses and strengths of individual models, such as SimCSE’s (Gao et al., 2021b) low performance on clustering and retrieval despite its strong performance on STS. We hope our work makes selecting the right embedding model easier and simplifies future embedding research.

2 Related Work

2.1 Benchmarks

Benchmarks, such as (Super)GLUE (Wang et al., 2018, 2019) or Big-BENCH (Srivastava et al., 2022), and evaluation frameworks (Gao et al., 2021a) play a key role in driving NLP progress. Yearly released SemEval datasets (Agirre et al., 2012, 2013, 2014, 2015, 2016) are commonly used as the go-to benchmark for text embeddings. SemEval datasets correspond to the task of semantic textual similarity (STS) requiring models to embed similar sentences with geometrically close embeddings. Due to the limited expressivity of a single SemEval dataset, SentEval (Conneau and Kiela, 2018) aggregates multiple STS datasets. SentEval focuses on fine-tuning classifiers on top of embeddings. It lacks tasks like retrieval or clustering, where embeddings are directly compared without additional classifiers. Further, the toolkit was proposed in 2018 and thus does not provide easy support for recent trends like text embeddings from transformers (Reimers and Gurevych, 2019). Due to the insufficiency of STS benchmarking, USEB (Wang et al., 2021) was introduced consisting mostly of reranking tasks. Consequently, it does not cover tasks like retrieval or classification. Meanwhile, the recently released BEIR Benchmark (Thakur et al., 2021) has become the standard for the evaluation of embeddings for zero-shot information retrieval.

MTEB unifies datasets from different embedding tasks into a common, accessible evaluation framework. MTEB incorporates SemEval datasets (STS11 - STS22) and BEIR alongside a variety of other datasets from various tasks to provide a holistic performance review of text embedding models.

2.2 Embedding Models

Text embedding models like Glove (Pennington et al., 2014) lack context awareness and are thus commonly labeled as Word Embedding Models. They consist of a layer mapping each input word to a vector often followed by an averaging layer to provide a final embedding invariant of input length.

Transformers (Vaswani et al., 2017) inject context awareness into language models via self-attention and form the foundation of most recent embedding models. BERT (Devlin et al., 2018) uses the transformer architecture and performs large-scale self-supervised pre-training. The resulting model can directly be used to produce text embeddings via an averaging operation alike Glove. Building on InferSent (Conneau et al., 2017), SBERT (Reimers and Gurevych, 2019) demonstrated it to be beneficial to perform additional fine-tuning of the transformer for competitive embedding performance. Most recent fine-tuned embedding models use a contrastive loss objective to perform supervised fine-tuning on positive and negative text pairs (Gao et al., 2021b; Wang et al., 2021; Ni et al., 2021b; Muennighoff, 2022). Due to the large variety of available pre-trained transformers (Wolf et al., 2020), there is an at least equally large variety of potential text embedding models to be explored. This leads to confusion about which model provides practitioners with the best performance for their embedding use case.

We benchmark both word embedding and transformer models on MTEB quantifying gains provided by often much slower context aware models.

3 The MTEB Benchmark

3.1 Desiderata

MTEB is built on a set of desiderata: **(a) Diversity:** MTEB aims to provide an understanding of the usability of embedding models in various use cases. The benchmark comprises 8 different tasks, with up to 15 datasets each. Of the 58 total datasets in MTEB, 10 are multilingual, covering 112 different languages. Sentence-level and paragraph-level datasets are included to contrast performance on short and long texts. **(b) Simplicity:** MTEB provides a simple API for plugging in any model that given a list of texts can produce a vector for each list item with a consistent shape. This makes it possible to benchmark a diverse set of models. **(c) Extensibility:** New datasets for existing tasks can be benchmarked in MTEB via a single file that specifies the task and a Hugging Face dataset name where the data has been uploaded (Lhoest et al., 2021). New tasks require implementing a task interface for loading the data and an evaluator for benchmarking. We welcome dataset, task or metric contributions from the community via pull requests to continue the development of MTEB. **(d) Repro-**

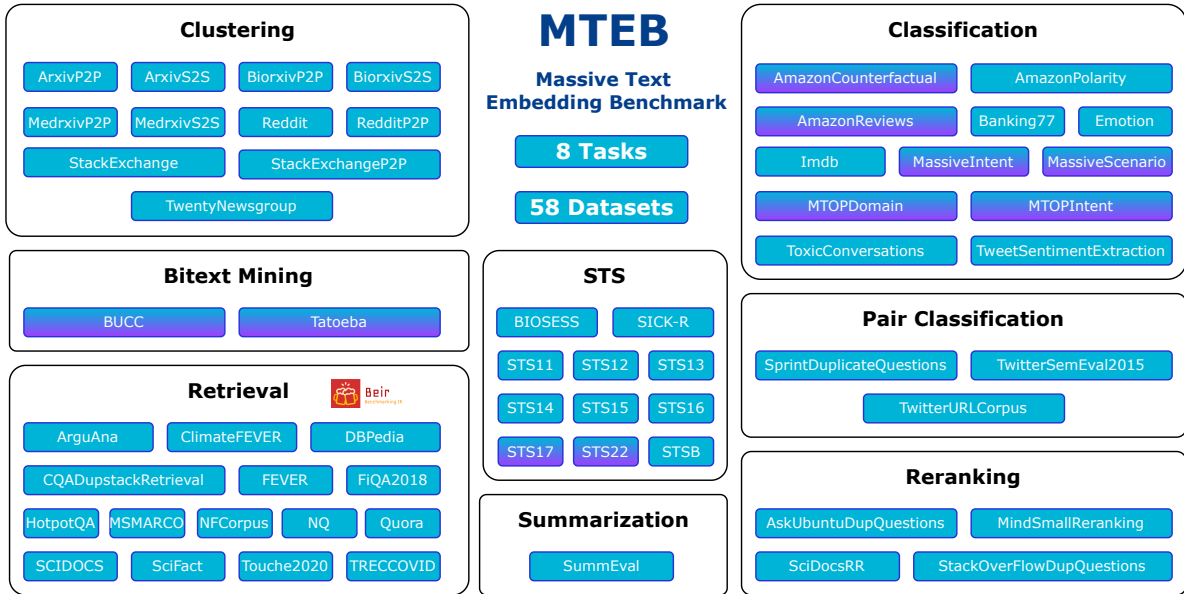


Figure 1: An overview of tasks and datasets in MTEB. Multilingual datasets are marked with a purple shade.

ducibility: Through versioning at a dataset and software level, we aim to make it easy to reproduce results in MTEB. JSON files corresponding to all results available in this paper have been made available together with the MTEB benchmark³.

3.2 Tasks and Evaluation

Figure 1 provides an overview of tasks and datasets available in MTEB. Dataset statistics are available in Table 2. The benchmark consists of the following 8 task types:

Bitext Mining Inputs are two sets of sentences from two different languages. For each sentence in the first set, the best match in the second set needs to be found. The matches are commonly translations. The provided model is used to embed each sentence and the closest pairs are found via cosine similarity. F1 serves as the main metric for bitext mining. Accuracy, precision and recall are also computed.

Classification A train and test set are embedded with the provided model. The train set embeddings are used to train a logistic regression classifier with 100 maximum iterations, which is scored on the test set. The main metric is accuracy with average precision and f1 additionally provided.

Clustering Given a set of sentences or paragraphs, the goal is to group them into meaningful clusters. A mini-batch k-means model with batch size 32 and k equal to the number of different labels (Pedregosa et al., 2011) is trained on the embedded texts. The model is scored using v-measure (Rosenberg and Hirschberg, 2007). V-measure does not depend on the cluster label, thus the permutation of labels does not affect the score.

Pair Classification A pair of text inputs is provided and a label needs to be assigned. Labels are typically binary variables denoting duplicate or paraphrase pairs. The two texts are embedded and their distance is computed with various metrics (cosine similarity, dot product, euclidean distance, manhattan distance). Using the best binary threshold accuracy, average precision, f1, precision and recall are computed. The average precision score based on cosine similarity is the main metric.

Reranking Inputs are a query and a list of relevant and irrelevant reference texts. The aim is to rank the results according to their relevance to the query. The model is used to embed the references which are then compared to the query using cosine similarity. The resulting ranking is scored for each query and averaged across all queries. Metrics are mean MRR@k and MAP with the latter being the main metric.

³<https://huggingface.co/datasets/mteb/results>

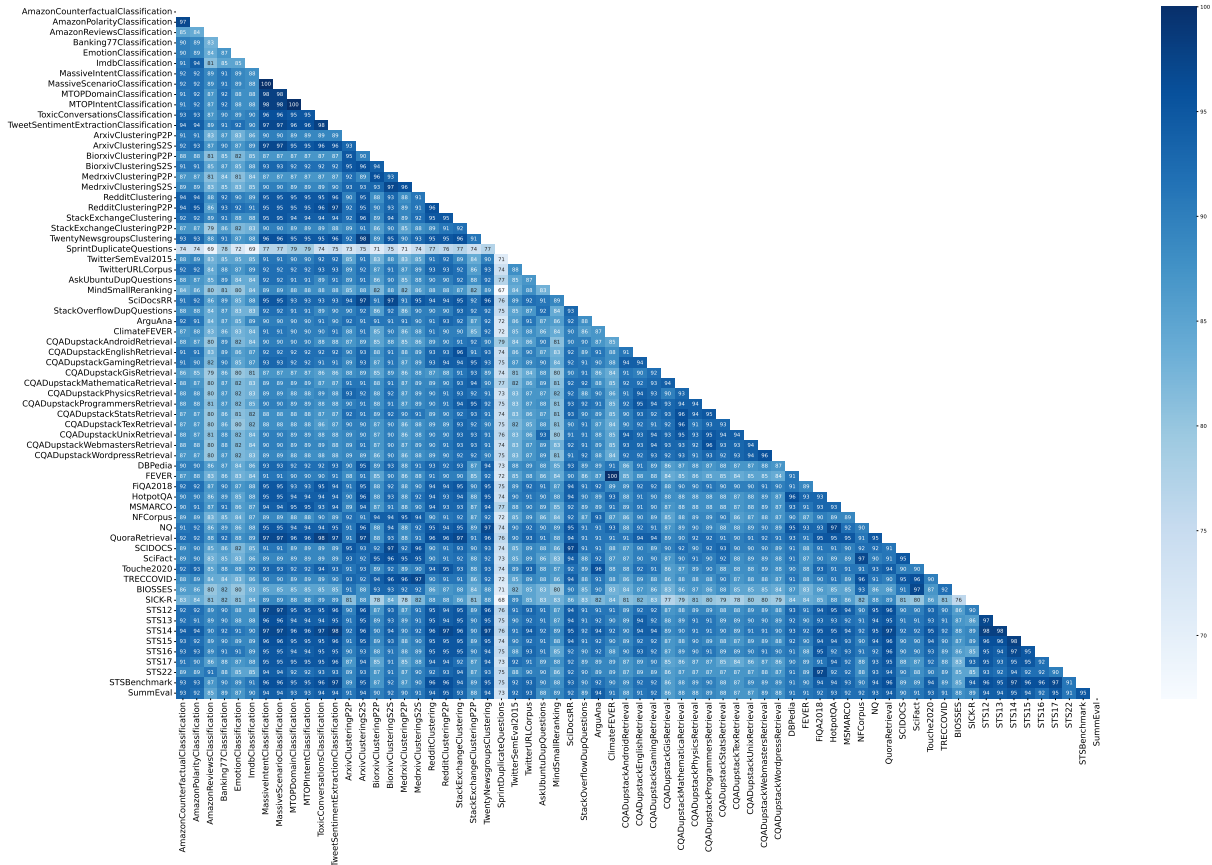


Figure 2: Similarity of MTEB datasets. We use the best model on MTEB STS (ST5-XXL, see Table 1) to embed 100 samples for each dataset. Cosine similarities between the averaged embeddings are computed and visualized.

Retrieval Each dataset consists of a corpus, queries and a mapping for each query to relevant documents from the corpus. The aim is to find these relevant documents. The provided model is used to embed all queries and all corpus documents and similarity scores are computed using cosine similarity. After ranking the corpus documents for each query based on the scores, nDCG@k, MRR@k, MAP@k, precision@k and recall@k are computed for several values of k . nDCG@10 serves as the main metric. MTEB reuses datasets and evaluation from BEIR (Thakur et al., 2021).

Semantic Textual Similarity (STS) Given a sentence pair the aim is to determine their similarity. Labels are continuous scores with higher numbers indicating more similar sentences. The provided model is used to embed the sentences and their similarity is computed using various distance metrics. Distances are benchmarked with ground truth similarities using Pearson and Spearman correlations. Spearman correlation based on cosine similarity serves as the main metric (Reimers et al., 2016).

Summarization A set of human-written and machine-generated summaries are provided. The aim is to score the machine summaries. The provided model is first used to embed all summaries. For each machine summary embedding, distances to all human summary embeddings are computed. The closest score (e.g. highest cosine similarity) is kept and used as the model’s score of a single machine-generated summary. Pearson and Spearman correlations with ground truth human assessments of the machine-generated summaries are computed. Like for STS, Spearman correlation based on cosine similarity serves as the main metric (Reimers et al., 2016).

3.3 Datasets

To further the diversity of MTEB, datasets of varying text lengths are included. All datasets are grouped into three categories:

Sentence to sentence (S2S) A sentence is compared with another sentence. An example of S2S are all current STS tasks in MTEB, where the similarity between two sentences is assessed.

Num. Datasets (→)	Class. 12	Clust. 11	PairClass. 3	Rerank. 4	Retr. 15	STS 10	Summ. 1	Avg. 56
<i>Self-supervised methods</i>								
Glove	57.29	27.73	70.92	43.29	21.62	61.85	28.87	41.97
Komninos	57.65	26.57	72.94	44.75	21.22	62.47	30.49	42.06
BERT	61.66	30.12	56.33	43.44	10.59	54.36	29.82	38.33
SimCSE-BERT-unsup	62.50	29.04	70.33	46.47	20.29	74.33	31.15	45.45
<i>Supervised methods</i>								
SimCSE-BERT-sup	67.32	33.43	73.68	47.54	21.82	79.12	23.31	48.72
coCondenser-msmarco	64.71	37.64	81.74	51.84	32.96	76.47	29.50	52.35
Contriever	66.68	41.10	82.53	53.14	41.88	76.51	30.36	56.00
SPECTER	52.37	34.06	61.37	48.10	15.88	61.02	27.66	40.28
LaBSE	62.71	29.55	78.87	48.42	18.99	70.80	31.05	45.21
LASER2	53.65	15.28	68.86	41.44	7.93	55.32	26.80	33.63
MiniLM-L6	63.06	42.35	82.37	58.04	41.95	78.90	30.81	56.26
MiniLM-L12	63.21	41.81	82.41	<u>58.44</u>	42.69	79.80	27.90	56.53
MiniLM-L12-multilingual	64.30	37.14	78.45	53.62	32.45	78.92	30.67	52.44
MPNet	65.07	<u>43.69</u>	83.04	59.36	43.81	80.28	27.49	57.78
MPNet-multilingual	67.91	38.40	80.81	53.80	35.34	80.73	31.57	54.71
OpenAI Ada Similarity	70.44	37.52	76.86	49.02	18.36	78.60	26.94	49.52
SGPT-125M-nli	61.46	30.95	71.78	47.56	20.90	74.71	30.26	45.97
SGPT-5.8B-nli	70.14	36.98	77.03	52.33	32.34	80.53	30.38	53.74
SGPT-125M-msmarco	60.72	35.79	75.23	50.58	37.04	73.41	28.90	51.23
SGPT-1.3B-msmarco	66.52	39.92	79.58	54.00	44.49	75.74	25.44	56.11
SGPT-2.7B-msmarco	67.13	39.83	80.65	54.67	46.54	76.83	27.87	57.12
SGPT-5.8B-msmarco	68.13	40.35	82.00	56.56	50.25	78.10	24.75	58.81
SGPT-BLOOM-7.1B-msmarco	66.19	38.93	81.90	55.65	48.21	77.74	24.99	57.44
GTR-Base	65.25	38.63	83.85	54.23	44.67	77.07	29.67	56.19
GTR-Large	67.14	41.60	85.33	55.36	47.42	78.19	29.50	58.28
GTR-XL	67.11	41.51	86.13	55.96	47.96	77.80	30.21	58.42
GTR-XXL	67.41	42.42	<u>86.12</u>	<u>56.65</u>	<u>48.48</u>	78.38	30.64	<u>58.97</u>
ST5-Base	69.81	40.21	85.17	53.09	33.63	81.14	<u>31.39</u>	55.27
ST5-Large	72.31	41.65	84.97	54.00	36.71	<u>81.83</u>	29.64	57.06
ST5-XL	<u>72.84</u>	42.34	86.06	54.71	38.47	81.66	29.91	57.87
ST5-XXL	73.42	43.71	85.06	56.43	42.24	82.63	30.08	59.51

Table 1: Average of the main metric (see Section 3.2) per task per model on MTEB English subsets.

Paragraph to paragraph (P2P) A paragraph is compared with another paragraph. MTEB imposes no limit on the input length, leaving it up to the models to truncate if necessary. Several clustering tasks are framed as both S2S and P2P tasks. The former only compare titles, while the latter include both title and content. For ArxivClustering, for example, abstracts are concatenated to the title in the P2P setting.

Sentence to paragraph (S2P) A few retrieval datasets are mixed in a S2P setting. Here a query is a single sentence, while documents are long paragraphs consisting of multiple sentences.

Similarities across 56 MTEB datasets are visualized in Figure 2. Several datasets rely on

the same corpora, such as ClimateFEVER and FEVER, resulting in a score of 1. Clusters of similar datasets can be seen among CQADupstack variations and STS datasets. S2S and P2P variations of the same dataset tend to also be similar. Scientific datasets, such as SciDocsRR, SciFact, ArxivClustering, show high similarities among each other even when coming from different tasks (Reranking, Retrieval and Clustering in this case).

4 Results

4.1 Models

We evaluate on the test splits of all datasets except for MSMARCO, where the dev split is used following Thakur et al. (2021). We benchmark models claiming state-of-the-art results on various embed-

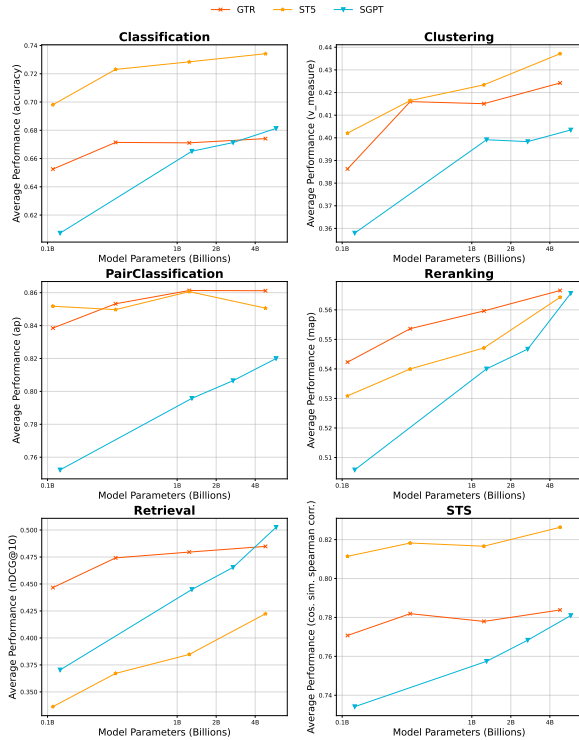


Figure 3: MTEB performance scales with model size. The smallest SGPT variant underperforms similar-sized GTR and ST5 variants. This may be due to the bias-only fine-tuning SGPT employs, which catches up with full fine-tuning only as model size and thus the number of bias parameters increases (Muennighoff, 2022).

ding tasks leading to a high representation of transformers (Vaswani et al., 2017). We group models into self-supervised and supervised methods.

Self-supervised methods (a) **Transformer-based BERT** (Devlin et al., 2018) is trained using self-supervised mask and sentence prediction tasks. By taking the mean across the sequence length (mean-pooling) the model can directly be used to produce text embeddings. SimCSE-Unsup (Gao et al., 2021b) uses BERT as a foundation and performs additional self-supervised training. (b) **Non-transformer:** Komninos (Komninos and Manandhar, 2016) and Glove (Pennington et al., 2014) are two word embedding models that directly map words to vectors. Hence, their embeddings lack context awareness, but provide significant speed-ups.

Supervised methods The original transformer model (Vaswani et al., 2017) consists of an encoder and decoder network. Subsequent transformers often train only encoders like BERT (Devlin et al.,

2018) or decoders like GPT (Radford et al., 2019).

(a) **Transformer encoder methods** coCondenser (Gao and Callan, 2021), Contriever (Izacard et al., 2021), LaBSE (Feng et al., 2020) and SimCSE-BERT-sup (Gao et al., 2021b) are based on the pre-trained BERT model (Devlin et al., 2018). coCondenser and Contriever add a self-supervised stage prior to supervised fine-tuning for a total of three training stages. LaBSE uses BERT to perform additional pre-training on parallel data to produce a competitive bitext mining model. SPECTER (Cohan et al., 2020a) relies on the pre-trained SciBERT (Beltagy et al., 2019) variant instead and fine-tunes on citation graphs. GTR (Ni et al., 2021b) and ST5 (Ni et al., 2021a) are based on the encoder part of the T5 model (Rafael et al., 2020) and only differ in their fine-tuning datasets. After additional self-supervised training, ST5 does contrastive fine-tuning on NLI (Ni et al., 2021a; Gao et al., 2021b) being geared towards STS tasks. Meanwhile, GTR fine-tunes on MS-MARCO and focuses on retrieval tasks. MPNet and MiniLM correspond to fine-tuned embedding models (Reimers and Gurevych, 2019) of the pre-trained MPNet (Song et al., 2020) and MiniLM (Wang et al., 2020) models using diverse datasets to target any embedding use case.

(b) **Transformer decoder methods** SGPT Bi-Encoders (Muennighoff, 2022) perform contrastive fine-tuning of <0.1% of pre-trained parameters using weighted-mean pooling. Similar to ST5 and GTR, SGPT-nli models are geared towards STS, while SGPT-msmarco models towards retrieval. SGPT-msmarco models embed queries and documents for retrieval with different special tokens to help the model distinguish their role. For non-retrieval tasks, we use its query representations. We benchmark publicly available SGPT models based on GPT-NeoX (Andonian et al., 2021), GPT-J (Wang and Komatsuzaki, 2021) and BLOOM (Scao et al., 2022). Alternatively, cpt-text (Neelakantan et al., 2022) passes pre-trained GPT decoders through a two-stage process using last token pooling to provide embeddings from decoders. We benchmark their models via the OpenAI Embeddings API⁴.

(c) **Non-transformer** LASER (Heffernan et al., 2022) is the only context aware non-transformer model we benchmark, relying on an LSTM

⁴<https://beta.openai.com/docs/guides/embeddings>

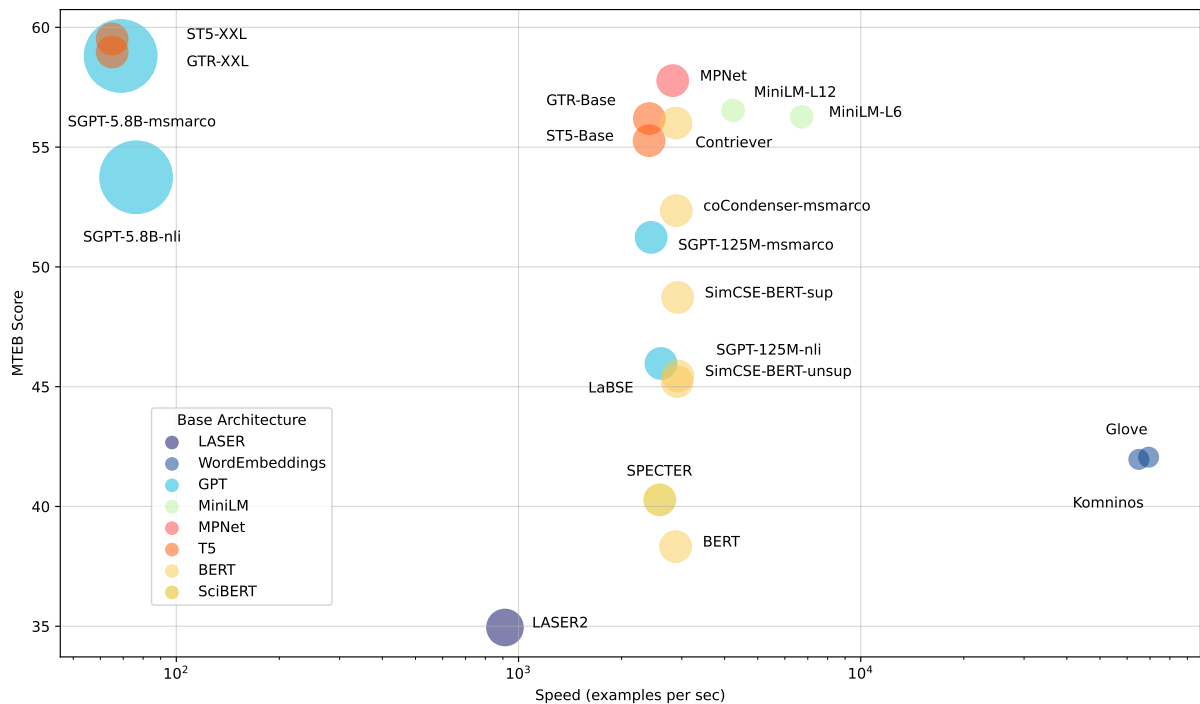


Figure 4: Performance, speed, and size of produced embeddings (size of the circles) of different embedding models. Embedding sizes range from 1.2 kB (Glove / Komninos) to 16.4 kB (SGPT-5.8B) per example. Speed was benchmarked on STS15 using 1x Nvidia A100 80GB with CUDA 11.6.

(Hochreiter and Schmidhuber, 1997) instead. Similar to LaBSE, the model trains on parallel data and focuses on bitext mining applications.

4.2 Analysis

Based on the results in Table 1, we observe that there is considerable variability between tasks. No model claims the state-of-the-art in all seven English tasks. There is even more variability in the results per dataset present in the appendix. Further, there remains a large gap between self-supervised and supervised methods. Self-supervised large language models have been able to close this gap in many natural language generation tasks (Chowdhery et al., 2022). However, they appear to still require supervised fine-tuning for competitive embedding performance.

We find that performance strongly correlates with model size, see Figure 3. A majority of MTEB tasks are dominated by multi-billion parameter models. However, these come at a significant cost as we investigate in Section 4.3.

Classification ST5 models dominate the classification task across most datasets, as can be seen in detail in the full results in the appendix. ST5-XXL has the highest average performance, 3% ahead of

the best non-ST5 model, OpenAI Ada Similarity.

Clustering Despite being almost 50x smaller, the MPNet embedding model is on par with the ST5-XXL state-of-the-art on Clustering. This may be due to the large variety of datasets MPNet (and MiniLM) has been fine-tuned on. Clustering requires coherent distances between a large number of embeddings. Models like SimCSE-sup or SGPT-nli, which are only fine-tuned on a single dataset, NLI, may produce incoherent embeddings when encountering topics unseen during fine-tuning. Relatedly, we find that the query embeddings of SGPT-msmarco and the Ada Search endpoint are competitive with SGPT-nli and the Ada Similarity endpoint, respectively. We refer to the public leaderboard⁵ for Ada Search results. This could be due to the MSMARCO dataset being significantly larger than NLI. Thus, while the OpenAI docs recommend using the similarity embeddings for clustering use cases⁶, the retrieval query embeddings may be the better choice in some cases.

⁵<https://huggingface.co/spaces/mteb/leaderboard>

⁶<https://beta.openai.com/docs/guides/embeddings/similarity-embeddings>

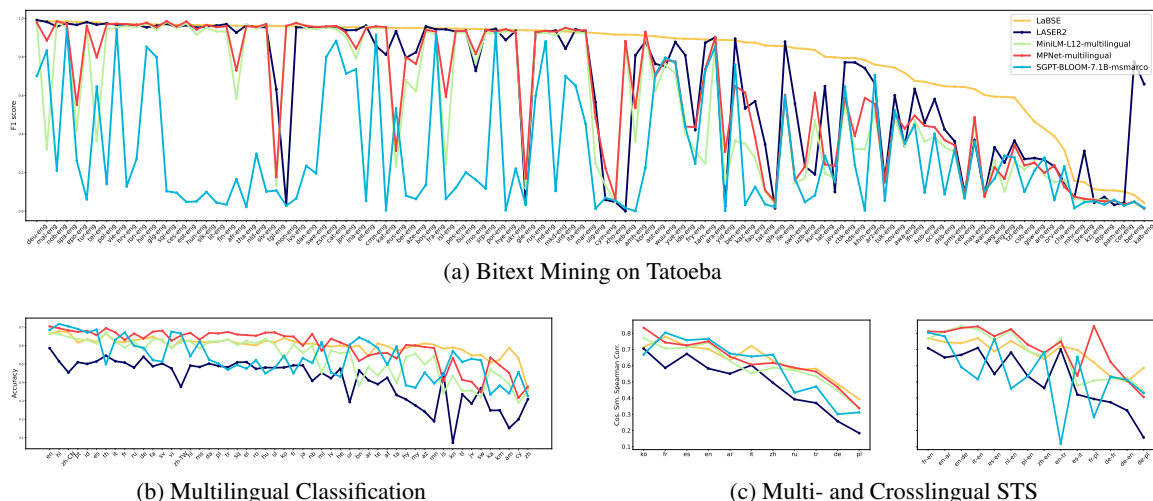


Figure 5: MTEB multilingual performance. Bitext mining is dominated by LaBSE, while classification and STS results are mixed. SGPT-BLOOM-7B1-msmarco tends to perform well on the languages BLOOM has been pre-trained on, such as Chinese, French and Portuguese.

Pair Classification GTR-XL and GTR-XXL have the strongest performance. Pair classification is closest to STS in its framing, yet models rank significantly differently on the two tasks. This highlights the importance of benchmarking on a diverse set of tasks to avoid blindly reusing a model for a different task.

Reranking MPNet and MiniLM models perform strongly on reranking tasks. On SciDocsRR (Cohan et al., 2020a) they perform far better than bigger models, which is likely due to parts of SciDocsRR being included in their training data. Our scale of experiments and that of model pre-training make controlling for data contamination challenging. Thus, we ignore overlap of MTEB datasets with model training datasets in MTEB scores. As long as enough datasets are averaged, we believe these effects to be insignificant.

Retrieval SGPT-5.8B-msmarco is the best embedding model on the BEIR subset in MTEB as well as on the full BEIR benchmark (Thakur et al., 2021; Muennighoff, 2022). The even larger 7.1B SGPT model making use of BLOOM (Scao et al., 2022) performs significantly weaker, which is likely due to the multilinguality of BLOOM. Models geared towards STS (SimCSE, ST5, SGPT-nli) perform badly on retrieval tasks. Retrieval tasks are unique in that there are two distinct types of texts: Queries and documents (“asymmetric”), while other tasks only have a single type of text (“symmetric”). On the QuoraRetrieval dataset, which has been shown to be largely symmetric

(Muennighoff, 2022), the playing field is more even with SGPT-5.8B-nli outperforming SGPT-5.8B-msmarco, see Table 11.

STS & Summarization Retrieval models (GTR, SGPT-msmarco) perform badly on STS, while ST5-XXL has the highest performance. This highlights the bifurcation of the field into separate embedding models for retrieval (asymmetric) and similarity (symmetric) use cases (Muennighoff, 2022).

4.3 Efficiency

We investigate the latency-performance trade-off of models in Figure 4. The graph allows for significant elimination of model candidates in the model selection process. It brings model selection down to three clusters:

Maximum speed Word Embedding models offer maximum speed with Glove taking the lead on both performance and speed, thus making the choice simple in this case.

Maximum performance If latency is less important than performance, the left-hand side of the graph offers a cluster of highly performant, but slow models. Depending on the task at hand, GTR-XXL, ST5-XXL or SGPT-5.8B may be the right choice, see Section 4.2. SGPT-5.8B comes with the additional caveat of its high-dimensional embeddings requiring more storage.

Speed and performance The fine-tuned MPNet and MiniLM models lead the middle cluster making the choice easy.

4.4 Multilinguality

MTEB comes with 10 multilingual datasets across bitext mining, classification and STS tasks. We investigate performance on these in Figure 5. Tabular results can be found in Tables 12, 13 and 14.

Bitext Mining LaBSE (Feng et al., 2020) performs strongly across a wide array of languages in bitext mining. Meanwhile, LASER2 shows high variance across different languages. While there are additional language-specific LASER2 models available for some of the languages we benchmark, we use the default multilingual LASER2 model for all languages. This is to provide a fair one-to-one comparison of models. In practice, however, the high variance of LASER2’s performance may be resolved by mixing its model variants. MPNet, MiniLM and SGPT-BLOOM-7B1-msmarco perform poorly on languages they have not been pre-trained on, such as German for the latter.

Classification & STS On multilingual classification and STS, the multilingual MPNet provides the overall strongest performance. It outperforms the slightly faster multilingual MiniLM on almost all languages. Both models have been trained on the same languages, thus bringing decision-making down to performance vs speed. SGPT-BLOOM-7B1-msmarco provides state-of-the-art performance on languages like Hindi, Portuguese, Chinese or French, which the model has seen extensively during pre-training. It also performs competitively on languages like Russian or Japanese that unintentionally leaked into its pre-training data (Muennighoff et al., 2022). However, it is not much ahead of the much cheaper MPNet. LASER2 performs consistently worse than other models.

5 Conclusion

In this work, we presented the Massive Text Embedding Benchmark (MTEB). Consisting of 8 text embedding tasks with up to 15 datasets each and covering 112 languages, MTEB aims to provide reliable embedding performance estimates. By open-sourcing MTEB alongside a leaderboard, we provide a foundation for further pushing the state-of-the-art of available text embeddings.

To introduce MTEB, we have conducted the most comprehensive benchmarking of text embeddings to date. Through the course of close to 5,000 experiments on over 30 different models, we have set up solid baselines for future research to build

on. We found model performance on different tasks to vary strongly with no model claiming state-of-the-art on all tasks. Our studies on scaling behavior, model efficiency and multilinguality revealed various intricacies of models that should ease the decision-making process for future research or industry applications of text embeddings.

We welcome task, dataset or metric contributions to the MTEB codebase⁷ as well as additions to the leaderboard via our automatic submission format⁸.

⁷<https://github.com/embeddings-benchmark/mteb>

⁸<https://huggingface.co/spaces/mteb/leaderboard>

Acknowledgments

This work was granted access to the HPC resources of Institut du développement et des ressources en informatique scientifique (IDRIS) du Centre national de la recherche scientifique (CNRS) under the allocation 2021-A0101012475 made by Grand équipement national de calcul intensif (GENCI). In particular, all the evaluations and data processing ran on the Jean Zay cluster of IDRIS, and we want to thank the IDRIS team for responsive support throughout the project, in particular Rémi Lacroix.

We thank Douwe Kiela, Teven Le Scao and Nandan Thakur for feedback and suggestions.

References

- Charu C Aggarwal and ChengXiang Zhai. 2012. A survey of text clustering algorithms. In *Mining text data*, pages 77–128. Springer.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Inigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, et al. 2015. Semeval-2015 task 2: Semantic textual similarity, english, spanish and pilot on interpretability. In *Proceedings of the 9th international workshop on semantic evaluation (SemEval 2015)*, pages 252–263.
- Eneko Agirre, Carmen Banea, Claire Cardie, Daniel M Cer, Mona T Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. 2014. Semeval-2014 task 10: Multilingual semantic textual similarity. In *SemEval@ COLING*, pages 81–91.
- Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, Aitor Gonzalez Agirre, Rada Mihalcea, German Rigau Claramunt, and Janyce Wiebe. 2016. Semeval-2016 task 1: Semantic textual similarity, monolingual and cross-lingual evaluation. In *SemEval-2016. 10th International Workshop on Semantic Evaluation; 2016 Jun 16-17; San Diego, CA. Stroudsburg (PA): ACL; 2016. p. 497-511*. ACL (Association for Computational Linguistics).
- Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. 2012. Semeval-2012 task 6: A pilot on semantic textual similarity. In ** SEM 2012: The First Joint Conference on Lexical and Computational Semantics—Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 385–393.
- Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. 2013. * sem 2013 shared task: Semantic textual similarity. In *Second joint conference on lexical and computational semantics (* SEM)*, volume 1: proceedings of the Main conference and the shared task: semantic textual similarity, pages 32–43.
- Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, et al. 2023. Santa-coder: don’t reach for the stars! *arXiv preprint arXiv:2301.03988*.
- Alex Andonian, Quentin Anthony, Stella Biderman, Sid Black, Preetham Gali, Leo Gao, Eric Hallahan, Josh Levy-Kramer, Connor Leahy, Lucas Nestler, Kip Parker, Michael Pieler, Shivanshu Purohit, Tri Songz, Phil Wang, and Samuel Weinbach. 2021. *GPT-NeoX: Large scale autoregressive language modeling in pytorch*.
- Dimo Angelov. 2020. Top2vec: Distributed representations of topics. *arXiv preprint arXiv:2008.09470*.
- Akari Asai, Jungo Kasai, Jonathan H Clark, Kenton Lee, Eunsol Choi, and Hannaneh Hajishirzi. 2020. Xor qa: Cross-lingual open-retrieval question answering. *arXiv preprint arXiv:2010.11856*.
- Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. Scibert: A pretrained language model for scientific text. *arXiv preprint arXiv:1903.10676*.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International Conference on Machine Learning*, pages 2206–2240. PMLR.
- Micael Carvalho, Rémi Cadène, David Picard, Laure Soulier, Nicolas Thome, and Matthieu Cord. 2018. Cross-modal retrieval in the cooking context: Learning semantic text-image embeddings. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 35–44.
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. *Efficient intent detection with dual sentence encoders*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.

- Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S Weld. 2020a. Specter: Document-level representation learning using citation-informed transformers. *arXiv preprint arXiv:2004.07180*.
- Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S. Weld. 2020b. Specter: Document-level representation learning using citation-informed transformers.
- Alexis Conneau and Douwe Kiela. 2018. Senteval: An evaluation toolkit for universal sentence representations. *arXiv preprint arXiv:1803.05449*.
- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic Barrault, and Antoine Bordes. 2017. Supervised learning of universal sentence representations from natural language inference data. *arXiv preprint arXiv:1705.02364*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2020. Summeval: Re-evaluating summarization evaluation.
- Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2020. Language-agnostic bert sentence embedding. *arXiv preprint arXiv:2007.01852*.
- Jack FitzGerald, Christopher Hench, Charith Peris, Scott Mackie, Kay Rottmann, Ana Sanchez, Aaron Nash, Liam Urbach, Vishesh Kakarala, Richa Singh, Swetha Ranganath, Laurie Crist, Misha Britan, Wouter Leeuwis, Gokhan Tur, and Prem Nataraajan. 2022. Massive: A 1m-example multilingual natural language understanding dataset with 51 typologically-diverse languages.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. 2021a. A framework for few-shot language model evaluation. *Version v0. 0.1. Sept*.
- Luyu Gao and Jamie Callan. 2021. Unsupervised corpus aware language model pre-training for dense passage retrieval. *arXiv preprint arXiv:2108.05540*.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021b. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*.
- Gregor Geigle, Nils Reimers, Andreas Rücklé, and Iryna Gurevych. 2021. Tweac: Transformer with extendable qa agent classifiers.
- Kevin Heffernan, Onur Çelebi, and Holger Schwenk. 2022. Bitext mining using distilled sentence representations for low-resource languages. *arXiv preprint arXiv:2205.12654*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padmanabhan, Giuseppe Ottaviano, and Linjun Yang. 2020. Embedding-based retrieval in facebook search. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2553–2561.
- Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. Code-searchnet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436*.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Towards unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*.
- Alexandros Komninos and Suresh Manandhar. 2016. Dependency based embeddings for sentence classification tasks. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, pages 1490–1500.
- Wuwei Lan, Siyu Qiu, Hua He, and Wei Xu. 2017. A continuously growing dataset of sentential paraphrases. In *Proceedings of The 2017 Conference on Empirical Methods on Natural Language Processing (EMNLP)*, pages 1235–1245. Association for Computational Linguistics.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, et al. 2021. Datasets: A community library for natural language processing. *arXiv preprint arXiv:2109.02846*.
- Haoran Li, Abhinav Arora, Shuohui Chen, Anchit Gupta, Sonal Gupta, and Yashar Mehdad. 2020. Mtop: A comprehensive multilingual task-oriented semantic parsing benchmark.
- Xueqing Liu, Chi Wang, Yue Leng, and ChengXiang Zhai. 2018. Linkso: a dataset for learning to retrieve similar question answer pairs on software development forums. In *Proceedings of the 4th ACM SIGSOFT International Workshop on NLP for Software Engineering*, pages 2–5.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.

- Julian McAuley and Jure Leskovec. 2013. [Hidden factors and hidden topics: Understanding rating dimensions with review text](#). RecSys '13, New York, NY, USA. Association for Computing Machinery.
- Niklas Muennighoff. 2020. Vilio: State-of-the-art visio-linguistic models applied to hateful memes. *arXiv preprint arXiv:2012.07788*.
- Niklas Muennighoff. 2022. Sgpt: Gpt sentence embeddings for semantic search. *arXiv preprint arXiv:2202.08904*.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- Pandu Nayak. 2019. [Understanding searches better than ever before](#).
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. 2022. Text and code embeddings by contrastive pre-training. *arXiv preprint arXiv:2201.10005*.
- Jianmo Ni, Gustavo Hernández Ábrego, Noah Constant, Ji Ma, Keith B Hall, Daniel Cer, and Yinfei Yang. 2021a. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models. *arXiv preprint arXiv:2108.08877*.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Ábrego, Ji Ma, Vincent Y Zhao, Yi Luan, Keith B Hall, Ming-Wei Chang, et al. 2021b. Large dual encoders are generalizable retrievers. *arXiv preprint arXiv:2112.07899*.
- Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*.
- James O’Neill, Polina Rozenshtein, Ryuichi Kiryo, Motoko Kubota, and Danushka Bollegala. 2021. [I wish i would have loved this one, but i didn’t – a multilingual dataset for counterfactual detection in product reviews](#).
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67.
- Nils Reimers, Philip Beyer, and Iryna Gurevych. 2016. Task-oriented intrinsic evaluation of semantic textual similarity. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 87–96.
- Nils Reimers and Iryna Gurevych. 2019. Sentencebert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- Facebook Research. [Tatoeba multilingual test set](#).
- Andrew Rosenberg and Julia Hirschberg. 2007. V-measure: A conditional entropy-based external cluster evaluation measure. pages 410–420.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*.
- Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. [CARER: Contextualized affect representations for emotion recognition](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3687–3697, Brussels, Belgium. Association for Computational Linguistics.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Darsh Shah, Tao Lei, Alessandro Moschitti, Salvatore Romeo, and Preslav Nakov. 2018. [Adversarial domain adaptation for duplicate question detection](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1056–1063, Brussels, Belgium. Association for Computational Linguistics.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tiejun Liu. 2020. Mpnet: Masked and permuted pre-training for language understanding. *Advances in Neural Information Processing Systems*, 33:16857–16867.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta,

- Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.
- Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*.
- Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. **Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models**.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.
- Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>.
- Kexin Wang, Nils Reimers, and Iryna Gurevych. 2021. Tsdae: Using transformer-based sequential denoising auto-encoder for unsupervised sentence embedding learning. *arXiv preprint arXiv:2104.06979*.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in Neural Information Processing Systems*, 33:5776–5788.
- Samuel Weinbach, Marco Bellagente, Constantin Eichenberg, Andrew Dai, Robert Baldock, Souradeep Nanda, Björn Deiseroth, Koen Oostermeijer, Hannah Teufel, and Andres Felipe Cruz-Salinas. 2022. M-vader: A model for diffusion with multimodal context. *arXiv preprint arXiv:2212.02936*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45.
- Fangzhao Wu, Ying Qiao, Jiun-Hung Chen, Chuhan Wu, Tao Qi, Jianxun Lian, Danyang Liu, Xing Xie, Jianfeng Gao, Winnie Wu, et al. 2020. Mind: A large-scale dataset for news recommendation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3597–3606.
- Wei Xu, Chris Callison-Burch, and William B Dolan. 2015. Semeval-2015 task 1: Paraphrase and semantic similarity in twitter (pit). In *Proceedings of the 9th international workshop on semantic evaluation (SemEval 2015)*, pages 1–11.
- Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamaloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. 2022. Making a miracle: Multilingual information retrieval across a continuum of languages. *arXiv preprint arXiv:2210.09984*.
- Jeffrey Zhu, Mingqin Li, Jason Li, and Cassandra Oduola. 2021. **Bing delivers more contextualized search using quantized transformer inference on nvidia gpus in azure**.
- Pierre Zweigenbaum, Serge Sharoff, and Reinhard Rapp. 2016. Towards preparation of the second bucc shared task: Detecting parallel sentences in comparable corpora. In *Proceedings of the Ninth Workshop on Building and Using Comparable Corpora. European Language Resources Association (ELRA), Portoroz, Slovenia*, pages 38–43.
- Pierre Zweigenbaum, Serge Sharoff, and Reinhard Rapp. 2017. Overview of the second bucc shared task: Spotting parallel sentences in comparable corpora. In *Proceedings of the 10th Workshop on Building and Using Comparable Corpora*, pages 60–67.
- Pierre Zweigenbaum, Serge Sharoff, and Reinhard Rapp. 2018. Overview of the third bucc shared task: Spotting parallel sentences in comparable corpora. In *Proceedings of 11th workshop on building and using comparable corpora*, pages 39–42.

A Datasets

Table 2 provides a summary along with statistics of all MTEB tasks. In the following, we give a brief description of each dataset included in MTEB.

A.1 Clustering

ArxivClusteringS2S, **ArxivClusteringP2P**, **BiorxivClusteringS2S**, **BiorxivClusteringP2P**, **MedrxivClusteringP2P**, **MedrxivClusteringS2S** These datasets are custom-made for MTEB using the public APIs from arXiv⁹ and bioRxiv/medRxiv¹⁰. For S2S datasets, the input text is simply the title of the paper, while for P2P the input text is the concatenation of the title and the abstract. The cluster labels are generated using categories given to the papers by humans. For bioRxiv and medRxiv this category is unique, but for arXiv multiple categories can be given to a single paper so we only use the first one. For bioRxiv and medRxiv there is only one level of category (e.g. biochemistry, genetics, microbiology, etc.) hence we only perform clustering based on that label. For arXiv there is a main category and secondary category: for example "cs.AI" means the main category is Computer Science and the sub-category is AI, math.AG means the main category is Mathematics and the sub-category is Algebraic Geometry etc. Hence, we create three types of splits:

(a) Main category clustering Articles are only clustered based on the main category (Math, Physics, Computer Science etc.). This split evaluates coarse clustering capacity of a model.

(b) Secondary category clustering within the same main category Articles are clustered based on their secondary category, but within a given main category, for example only Math papers that need to be clustered into Algebraic Geometry, Functional Analysis, Numerical Analysis etc. This split evaluates fine-grained clustering capacity of a model, as differentiating some sub-categories can be very difficult.

(c) Secondary category clustering Articles are clustered based on their secondary category for all main categories, so the labels can be Number Theory, Computational Complexity, Astrophysics of Galaxies etc. These splits evaluate fine-grained

clustering capacity, as well as multi-scale capacities i.e. is a model able to both separate Maths from Physics as well as Probability from Algebraic Topology at the same time.

For every dataset, split and strategy, we select subsets of all labels and then sample articles from those labels. This yields splits with a varying amount and size of clusters.

RedditClustering (Geigle et al., 2021): Clustering of titles from 199 subreddits. Clustering of 25 splits, each with 10-50 classes, and each class with 100 - 1000 sentences

RedditClusteringP2P Dataset created for MTEB using available data from Reddit posts¹¹. The task consists of clustering the concatenation of title+post according to their subreddit. It contains 10 splits, with 10 and 100 clusters per split and 1,000 to 100,000 posts.

StackExchangeClustering (Geigle et al., 2021) Clustering of titles from 121 stackexchanges. Clustering of 25 splits, each with 10-50 classes, and each class with 100-1000 sentences.

StackExchangeClusteringP2P Dataset created for MTEB using available data from StackExchange posts¹². The task consists of clustering the concatenation of title and post according to their subreddit. It contains 10 splits, with 10 to 100 clusters and 5,000 to 10,000 posts per split.

TwentyNewsgroupsClustering¹³ Clustering of the 20 Newsgroups dataset, given titles of article the goal is to find the newsgroup (20 in total). Contains 10 splits, each with 20 classes, with each split containing between 1,000 and 10,000 titles.

A.2 Classification

AmazonCounterfactual (O’Neill et al., 2021) A collection of Amazon customer reviews annotated for counterfactual detection pair classification. For each review the label is either "counterfactual" or "not-counterfactual". This is a multilingual dataset with 4 available languages.

¹¹<https://huggingface.co/datasets/sentence-transformers/reddit-title-body>

¹²https://huggingface.co/datasets/flax-sentence-embeddings/stackexchange_title_body_jsonl

¹³https://scikit-learn.org/0.19/datasets/twenty_newsgroups.html

⁹<https://arxiv.org/help/api/>

¹⁰<https://api.biorxiv.org/>

Name	Type	Categ.	#Lang.	Train Samples	Dev Samples	Test Samples	Train avg. chars	Dev avg. chars	Test avg. chars
BUCC	BitextMining	s2s	4	0	0	641684	0	0	101.3
Tatoeba	BitextMining	s2s	112	0	0	2000	0	0	39.4
AmazonCounterfactualClassification	Classification	s2s	4	4018	335	670	107.3	109.2	106.1
AmazonPolarityClassification	Classification	p2p	1	3600000	0	400000	431.6	0	431.4
AmazonReviewsClassification	Classification	s2s	6	1200000	30000	30000	160.5	159.2	160.4
Banking77Classification	Classification	s2s	1	10003	0	3080	59.5	0	54.2
EmotionClassification	Classification	s2s	1	16000	2000	2000	96.8	95.3	96.6
ImdbClassification	Classification	p2p	1	25000	0	25000	1325.1	0	1293.8
MassiveIntentClassification	Classification	s2s	51	11514	2033	2974	35.0	34.8	34.6
MassiveScenarioClassification	Classification	s2s	51	11514	2033	2974	35.0	34.8	34.6
MTOPDomainClassification	Classification	s2s	6	15667	2235	4386	36.6	36.5	36.8
MTOPIntentClassification	Classification	s2s	6	15667	2235	4386	36.6	36.5	36.8
ToxicConversationsClassification	Classification	s2s	1	50000	0	50000	298.8	0	296.6
TweetSentimentExtractionClassification	Classification	s2s	1	27481	0	3534	68.3	0	67.8
ArxivClusteringP2P	Clustering	p2p	1	0	0	732723	0	0	1009.9
ArxivClusteringS2S	Clustering	s2s	1	0	0	732723	0	0	74.0
BiorxivClusteringP2P	Clustering	p2p	1	0	0	75000	0	0	1666.2
BiorxivClusteringS2S	Clustering	s2s	1	0	0	75000	0	0	101.6
MedrxivClusteringP2P	Clustering	p2p	1	0	0	37500	0	0	1981.2
MedrxivClusteringS2S	Clustering	s2s	1	0	0	37500	0	0	114.7
RedditClustering	Clustering	s2s	1	0	420464	420464	0	64.7	64.7
RedditClusteringP2P	Clustering	p2p	1	0	0	459399	0	0	727.7
StackExchangeClustering	Clustering	s2s	1	0	417060	373850	0	56.8	57.0
StackExchangeClusteringP2P	Clustering	p2p	1	0	0	75000	0	0	1090.7
TwentyNewsgroupsClustering	Clustering	s2s	1	0	0	59545	0	0	32.0
SprintDuplicateQuestions	PairClassification	s2s	1	0	101000	101000	0	65.2	67.9
TwitterSemEval2015	PairClassification	s2s	1	0	0	16777	0	0	38.3
TwitterURLCorpus	PairClassification	s2s	1	0	0	51534	0	0	79.5
AskUbuntuDupQuestions	Reranking	s2s	1	0	0	2255	0	0	52.5
MindSmallReranking	Reranking	s2s	1	231530	0	107968	69.0	0	70.9
SciDocsRR	Reranking	s2s	1	0	19594	19599	0	69.4	69.0
StackOverflowDupQuestions	Reranking	s2s	1	23018	3467	3467	49.6	49.8	49.8
ArguAna	Retrieval	p2p	1	0	0	10080	0	0	1052.9
ClimateFEVER	Retrieval	s2p	1	0	0	5418128	0	0	539.1
CQADupstackAndroidRetrieval	Retrieval	s2p	1	0	0	23697	0	0	578.7
CQADupstackEnglishRetrieval	Retrieval	s2p	1	0	0	41791	0	0	467.1
CQADupstackGamingRetrieval	Retrieval	s2p	1	0	0	46896	0	0	474.7
CQADupstackGisRetrieval	Retrieval	s2p	1	0	0	38522	0	0	991.1
CQADupstackMathematicaRetrieval	Retrieval	s2p	1	0	0	17509	0	0	1103.7
CQADupstackPhysicsRetrieval	Retrieval	s2p	1	0	0	39355	0	0	799.4
CQADupstackProgrammersRetrieval	Retrieval	s2p	1	0	0	33052	0	0	1030.2
CQADupstackStatsRetrieval	Retrieval	s2p	1	0	0	42921	0	0	1041.0
CQADupstackTexRetrieval	Retrieval	s2p	1	0	0	71090	0	0	1246.9
CQADupstackUnixRetrieval	Retrieval	s2p	1	0	0	48454	0	0	984.7
CQADupstackWebmastersRetrieval	Retrieval	s2p	1	0	0	17911	0	0	689.8
CQADupstackWordpressRetrieval	Retrieval	s2p	1	0	0	49146	0	0	1111.9
DBPedia	Retrieval	s2p	1	0	4635989	4636322	0	310.2	310.1
FEVER	Retrieval	s2p	1	0	0	5423234	0	0	538.6
FiQA2018	Retrieval	s2p	1	0	0	58286	0	0	760.4
HotpotQA	Retrieval	s2p	1	0	0	5240734	0	0	288.6
MSMARCO	Retrieval	s2p	1	0	8848803	8841866	0	336.6	336.8
MSMARCOv2	Retrieval	s2p	1	138641342	138368101	0	341.4	342.0	0
NFCorpus	Retrieval	s2p	1	0	0	3956	0	0	1462.7
NQ	Retrieval	s2p	1	0	0	2684920	0	0	492.7
QuoraRetrieval	Retrieval	s2s	1	0	0	532931	0	0	62.9
SCIDOCS	Retrieval	s2p	1	0	0	26657	0	0	1161.9
SciFact	Retrieval	s2p	1	0	0	5483	0	0	1422.3
Touche2020	Retrieval	s2p	1	0	0	382594	0	0	1720.1
TRECCOVID	Retrieval	s2p	1	0	0	171382	0	0	1117.4
BIOSESSE	STS	s2s	1	200	200	200	156.6	156.6	156.6
SICK-R	STS	s2s	1	19854	19854	19854	46.1	46.1	46.1
STS12	STS	s2s	1	4468	0	6216	100.7	0	64.7
STS13	STS	s2s	1	0	0	3000	0	0	54.0
STS14	STS	s2s	1	0	0	7500	0	0	54.3
STS15	STS	s2s	1	0	0	6000	0	0	57.7
STS16	STS	s2s	1	0	0	2372	0	0	65.3
STS17	STS	s2s	11	0	0	500	0	0	43.3
STS22	STS	p2p	18	0	0	8060	0	0	1992.8
STSBenchmark	STS	s2s	1	11498	3000	2758	57.6	64.0	53.6
SummEval	Summarization	p2p	1	0	0	2800	0	0	359.8

Table 2: Tasks in MTEB

AmazonPolarity (McAuley and Leskovec, 2013) A collection of Amazon customer reviews annotated for polarity classification. For each review the label is either "positive" or "negative".

AmazonReviews (McAuley and Leskovec, 2013) A collection of Amazon reviews designed to aid research in multilingual text classification. For each review the label is the score given by the review between 0 and 4 (1-5 stars). This is a

multilingual dataset with 6 available languages.

Banking77 (Casanueva et al., 2020) Dataset composed of online banking queries annotated with their corresponding intents. For each user query the label is an intent among 77 intents like 'activate_my_card', 'apple_pay', 'bank_transfer', etc.

Emotion (Saravia et al., 2018) Dataset of English Twitter messages with six basic emotions: anger, fear, joy, love, sadness, and surprise.

Imdb (Maas et al., 2011) Large movie review dataset with labels being positive or negative.

MassiveIntent (FitzGerald et al., 2022) A collection of Amazon Alexa virtual assistant utterances annotated with the associated intent. For each user utterance the label is one of 60 intents like 'play_music', 'alarm_set', etc. This is a multilingual dataset with 51 available languages.

MassiveScenario (FitzGerald et al., 2022) A collection of Amazon Alexa virtual assistant utterances annotated with the associated intent. For each user utterance the label is a theme among 60 scenarios like 'music', 'weather', etc. This is a multilingual dataset with 51 available languages.

MTOPTopDomain / MTOPTopIntent Multilingual sentence datasets from the MTOPTop (Li et al., 2020) benchmark. We refer to their paper for details.

ToxicConversations Dataset from Kaggle competition¹⁴. Collection of comments from the Civil Comments platform together with annotations if the comment is toxic or not.

TweetSentimentExtraction Dataset from Kaggle competition¹⁵. Sentiment classification of tweets as neutral, positive or negative.

A.3 Pair Classification

SprintDuplicateQuestions (Shah et al., 2018): Collection of questions from the Sprint community. The goal is to classify a pair of sentences as duplicates or not.

TwitterSemEval2015 (Xu et al., 2015) Paraphrase-Pairs of Tweets from the SemEval 2015 workshop. The goal is to classify a pair of tweets as paraphrases or not.

¹⁴<https://www.kaggle.com/competitions/jigsaw-unintended-bias-in-toxicity-classification>

¹⁵<https://www.kaggle.com/competitions/tweet-sentiment-extraction>

TwitterURLCorpus (Lan et al., 2017) Paraphrase-Pairs of Tweets. The goal is to classify a pair of tweets as paraphrases or not.

A.4 Bitext Mining

BUCC (Zweigenbaum et al., 2016, 2017, 2018) BUCC provides big set of sentences (~ 10-70k each) for English, French, Russian, German and Chinese, along with associated pairs annotation. The annotated pairs here corresponds to a pairs of translated sentences, i.e. a sentence and its translation in the other language.

Tatoeba (Research) Tatoeba provides sets of sentences (1000 sentences each) for 112 languages with annotated associated pairs. Each pair is one sentence and its translation in another language.

A.5 Reranking

AskUbuntuDupQuestions¹⁶ Questions from AskUbuntu with manual annotations marking pairs of questions as similar or dissimilar.

MindSmall (Wu et al., 2020) Large-scale English Dataset for News Recommendation Research. Ranking news article titles given the title of a news article. The idea is to recommend other news from the one you are reading.

SciDocsRR (Cohan et al., 2020b) Ranking of related scientific papers based on their title.

StackOverflowDupQuestions (Liu et al., 2018) Stack Overflow Duplicate Questions Task for questions with the tags Java, JavaScript and Python, ranking questions as duplicates or not.

A.6 Semantic Textual Similarity (STS)

STS12, STS13, STS14, STS15, STS16, STS17, STS22, STSBenchmark (Agirre et al., 2012, 2013)^{17,18,19,20} Original STS benchmark, with scores from 0 to 5. The selection of sentences includes text from image captions, news headlines and user forums. In total they contain between 1,000 and 20,000 sentences. STS12 - STS16 and

¹⁶<https://github.com/taolei87/askubuntu>
¹⁷<https://alt.qcri.org/semeval2014/tas-k10/>

¹⁸<https://alt.qcri.org/semeval2015/tas-k2/>

¹⁹<https://alt.qcri.org/semeval2016/tas-k1/>

²⁰<https://competitions.codalab.org/competitions/33835>

STSBenchmark are monolingual English benchmarks. STS17 and STS22 contain crosslingual pairs of sentences, where the goal is to assess the similarity of two sentences in different languages. STS17 has 11 language pairs (among Korean, Arabic, English, French, German, Turkish, Spanish, Italian and Dutch) and STS22 has 18 language pairs (among Arabic, English, French, German, Turkish, Spanish, Polish, Italian, Russian and Chinese).

BIOSESSE²¹ Contains 100 sentence pairs from the biomedical field.

SICK-R (Agirre et al., 2014) Sentences Involving Compositional Knowledge (SICK) contains a large number of sentence pairs (10 000) that are lexically, syntactically and semantically rich.

A.7 Summarization

SummEval (Fabbri et al., 2020) Summaries generated by recent summarization models trained on CNN or DailyMail alongside human annotations.

A.8 Retrieval

We refer to the BEIR paper (Thakur et al., 2021), which contains description of each dataset. For MTEB, we include all publicly available datasets: **ArguAna**, **ClimateFEVER**, **CQADupstack**, **DBPedia**, **FEVER**, **FiQA2018**, **HotpotQA**, **MS-MARCO**, **NFCorpus**, **NQ**, **Quora**, **SCIDOCS**, **SciFact**, **Touche2020**, **TRECCOVID**.

B Limitations of MTEB

While MTEB aims to be a diverse benchmark to provide holistic performance reviews, the benchmark has its limitations. We list them here:

1. Long document datasets MTEB covers multiple text lengths (S2S, P2P, S2P), but very long documents are still missing. The longest datasets in MTEB have a few hundred words, and longer text sizes could be relevant for use cases like retrieval.

2. Task imbalance Tasks in MTEB have a different amount of datasets with summarization consisting of only a single dataset. This means MTEB average scores, which are computed over all datasets, are biased towards tasks with many datasets, notably retrieval, classification and clustering. As MTEB grows, we hope to add more datasets to currently underrepresented tasks like summarization or pair classification.

²¹<https://tabilab.cmpe.boun.edu.tr/BIOSESSE/DataSet.html>

3. Multilinguality MTEB contains multilingual classification, STS and bitext mining datasets. However, retrieval and clustering are English-only. SGPT-BLOOM-7B1-msmarco is geared towards multilingual retrieval datasets and due to the lack thereof cannot be comprehensively benchmarked in MTEB. Further, MTEB does not contain any code datasets that could be used to benchmark code models (Neelakantan et al., 2022; Allal et al., 2023). It should be easy to extend MTEB with datasets, such as CodeSearchNet (Husain et al., 2019), TyDI QA (Clark et al., 2020), XOR QA (Asai et al., 2020) or MIRACL (Zhang et al., 2022).

4. Additional modalities Text embeddings are commonly used as input features for downstream models, such as in our classification task. This can involve other modalities, notably image content (Carvalho et al., 2018; Tan and Bansal, 2019; Muennighoff, 2020; Nichol et al., 2021; Saharia et al., 2022; Weinbach et al., 2022). We have focused solely on natural language applications and leave extensive benchmarking of text embeddings as inputs for other modalities to future work.

C Examples

Tables 3-9 provide examples for each dataset for each task. For retrieval datasets, we refer to the BEIR paper (Thakur et al., 2021).

D Correlations

Figure 6 provides correlation heatmaps for model performance and MTEB tasks.

E Models

Table 10 provides publicly available model checkpoints used for MTEB evaluation.

F Additional results

Tables 11 until the end provide results on individual datasets of MTEB. The results are additionally available in json format on the Hugging Face Hub²² and can be inspected on the leaderboard²³.

²²<https://huggingface.co/datasets/mteb/results>

²³<https://huggingface.co/spaces/mteb/leaderboard>

Dataset	Text	Label
AmazonCounterfactualClassification	In person it looks as though it would have cost a lot more.	counterfactual
AmazonPolarityClassification	an absolute masterpiece I am quite sure any of you actually taking the time to read this have played the game at least once, and heard at least a few of the tracks here. And whether you were aware of it or not, Mitsuda's music contributed greatly to the...	positive
AmazonReviewsClassification	solo llega una unidad cuando te obligan a comprar dos Te obligan a comprar dos unidades y te llega solo una y no hay forma de reclamar, una autentica estafa, no compreis!!	0
Banking77Classification	What currencies is an exchange rate calculated in?	exchange_rate
EmotionClassification	i feel so inhibited in someone elses kitchen like im painting on someone elses picture	sadness
ImdbClassification	When I first saw a glimpse of this movie, I quickly noticed the actress who was playing the role of Lucille Ball. Rachel York's portrayal of Lucy is absolutely awful. Lucille Ball was an astounding comedian with incredible talent. To think about a legend like Lucille Ball being portrayed the way she was in the movie is horrendous. I cannot believe...	negative
MassiveIntentClassification	réveille-moi à neuf heures du matin le vendredi	alarm_set
MassiveScenarioClassification	tell me the artist of this song	music
MTOPDomainClassification	Maricopa County weather forecast for this week	weather
MTOPIntentClassification	what ingredients do is have left	GET_INFO_RECIPES
ToxicConversationsClassification	The guy's a damn cop, so what do you expect?	toxic
TweetSentimentExtractionClassification	I really really like the song Love Story by Taylor Swift	positive

Table 3: Classification examples

Dataset	Text	Cluster
ArxivClusteringP2P	Finite groups of rank two which do not involve $Qd(p)$. Let $p > 3$ be a prime. We show that if G is a finite group with p -rank equal to 2, then G involves $Qd(p)$ if and only if G p' -involves $Qd(p)$. This allows us to use a version of Glauberman's ZJ-theorem to give a more direct construction of finite group actions on mod- p homotopy spheres. We give an example to illustrate that the above conclusion does not hold for $p \leq 3$.	math
ArxivClusteringS2S	Vertical shift and simultaneous Diophantine approximation on polynomial curves	math
BiorxivClusteringP2P	Innate Immune sensing of Influenza A viral RNA through IFI16 promotes pyroptotic cell death Programmed cell death pathways are triggered by various stresses or stimuli, including viral infections. The mechanism underlying the regulation of these pathways upon Influenza A virus IAV infection is not well characterized. We report that a cytosolic DNA sensor IFI16 is...	immunology
BiorxivClusteringS2S	Association of CDH11 with ASD revealed by matched-gene co-expression analysis and mouse behavioral	neuroscience
MedrxivClusteringP2P	Temporal trends in the incidence of haemophagocytic lymphohistiocytosis: a nationwide cohort study from England 2003-2018. Haemophagocytic lymphohistiocytosis (HLH) is rare, results in high mortality and is increasingly being diagnosed. Little is known about what is driving the apparent rise in the incidence of this disease. Using national linked electronic health data from hospital admissions and death certification cases of HLH that were diagnosed in England between 1/1/2003 and 31/12/2018 were identified using a previously validated approach. We calculated incidence...	infectious diseases
MedrxivClusteringS2S	Current and Lifetime Somatic Symptom Burden Among Transition-aged Young Adults on the Autism Spectrum	psychiatry and clinical psychology
RedditClustering	Could anyone tell me what breed my bicolor kitten is?	r/cats
RedditClusteringP2P	Headaches after working out? Hey guys! I've been diagnosed with adhd since I was seven. I just recently got re-diagnosed (22f) and I've been out on a different medication, adderall I was normally taking vyvanse but because of cost and no insurance adderall was more affordable. I've noticed that if I take adderall and workout...	r/ADHD
StackExchangeClustering	Does this property characterize a space as Hausdorff?	math.stackexchange.com
StackExchangeClusteringP2P	Google play services error DEBUG: Application is pausing, which disconnects the RTMP client. I am having this issue from past day with Google Play Services Unity. What happens is, when I install app directly ot device via Unity, the Google Play Services work fine but when I upload it as beta to play store console and install it via that then it starts to give " DEBUG: Application is pausing, which disconnects the RTMP client" error. I have a proper SHA1 key.	unity
TwentyNewsgroupsClustering	Commercial mining activities on the moon	14

Table 4: Clustering examples

Dataset	Sentence 1	Sentence 2	Label
SprintDuplicateQuestions	Franklin U722 USB modem signal strength	How do I know if my Franklin U772 USB Modem has a weak signal ?	1
TwitterSemEval2015	All the home alones watching 8 mile", "All the home alones watching 8 mile	The last rap battle in 8 Mile nevr gets old ahah	0
TwitterURLCorpus	How the metaphors we use to describe discovery affect men and women in the sciences	Light Bulbs or Seeds ? How Metaphors for Ideas Influence Judgments About Genius	0

Table 5: Pair classification examples. Labels are binary.

Dataset	Query	Positive	Negative
AskUbuntuDupQuestions	change the application icon theme but not changing the panel icons	change folder icons in ubuntu-mono-dark theme	change steam tray icon back to default
MindSmallReranking	Man accused in probe of Giuliani associates is freed on bail	Studies show these are the best and worst states for your retirement	There are 14 cheap days to fly left in 2019: When are they and what deals can you score?
SciDocsRR	Discovering social circles in ego networks	Benchmarks for testing community detection algorithms on directed and weighted graphs with overlapping communities.	Improving www proxies performance with greedy-dual-size-frequency caching policy
StackOverflowDupQuestions	Java launch error selection does not contain a main type	Error: Selection does not contain a main type	Selection Sort in Java

Table 6: Reranking examples

Dataset	Sentence 1	Sentence 2	Score
BIOSSES	It has recently been shown that Craf is essential for Kras G12D-induced NSCLC.	It has recently become evident that Craf is essential for the onset of Kras-driven non-small cell lung cancer.	4.0
SICK-R	A group of children is playing in the house and there is no man standing in the background	A group of kids is playing in a yard and an old man is standing in the background	3.2
STS12	Nationally, the federal Centers for Disease Control and Prevention recorded 4,156 cases of West Nile, including 284 deaths.	There were 293 human cases of West Nile in Indiana in 2002, including 11 deaths statewide.	1.7
STS13	this frame has to do with people (the residents) residing in locations , sometimes with a co-resident .	inhabit or live in ; be an inhabitant of ;	2.8
STS14	then the captain was gone.	then the captain came back.	0.8
STS15	you 'll need to check the particular policies of each publisher to see what is allowed and what is not allowed.	if you need to publish the book and you have found one publisher that allows it.	3.0
STS16	you do not need to worry.	you don 't have to worry.	5.0
STS17	La gente muestra su afecto el uno por el otro.	A women giving something to other lady.	1.4
STS22	El secretario general de la Asociación Gremial de los Trabajadores del Subte y Premetro de Metrodelegados, Beto Pianelli, dijo que el Gobierno porteño debe convocar "inmediatamente" a licitación para la compra de nuevos trenes y retirar los que quedan en circulación...	En diálogo con el servicio informativo de la Radio Pública, el ministro de Salud de la Nación, Ginés González García, habló sobre el avance del coronavirus en la Argentina y se manifestó a favor de prorrogar la cuarentena obligatoria dis-puesta por...	1
STSBenchmark	A man is playing the cello.	A man seated is playing the cello.	4.25

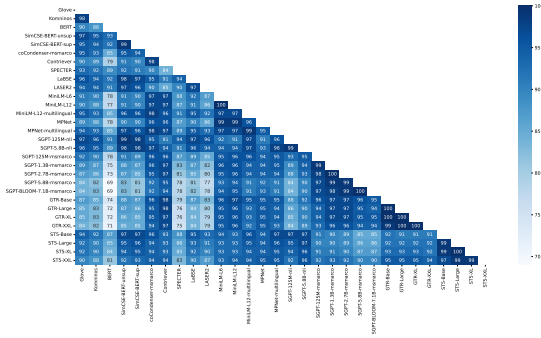
Table 7: STS examples. Scores are continuous between 0 and 5 (included).

Dataset	First set sentence	Second set sentence
BUCC	Morales remporte l'élection présidentielle de 2005 à la majorité absolue.	Morales went on to win the 2005 presidential election with an absolute majority.
Tatoeba	Chi le ha detto che Tom l'ha fatto?	Who told you that Tom did that?

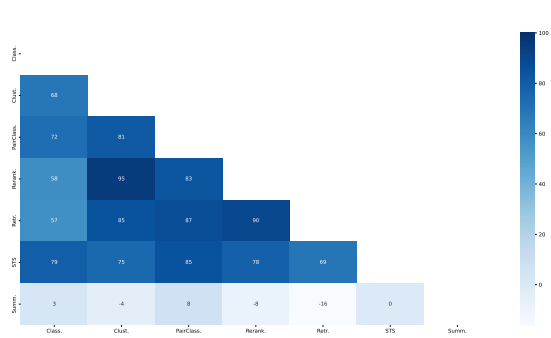
Table 8: Bitext mining examples

Dataset	Human Summary	Machine Summary	Relevance
SummEval	V. Stiviano must pay back \$2.6 million in gifts from Donald Sterling. Sterling's wife claimed the ex-Clippers used the couple's money for the gifts. The items included a Ferrari, two Bentleys and a Range Rover.	donald sterling , nba team last year . sterling 's wife sued for \$ 2.6 million in gifts . sterling says he is the former female companion who has lost the . sterling has ordered v. stiviano to pay back \$ 2.6 m in gifts after his wife sued . sterling also includes a \$ 391 easter bunny costume , \$ 299 and a \$ 299 .	1.7

Table 9: Summarization example



(a) Model correlation based on all results



(b) Task correlation based on average task results

Figure 6: Pearson correlations across model and task results. **Left:** Size variants of the same architecture show high correlations. **Right:** Performance on clustering and reranking correlates strongest, while summarization and classification show weaker correlation with other tasks.

Model	Public Checkpoint
Glove	https://huggingface.co/sentence-transformers/average_word_embeddings_glove.6B.300d
Komninos	https://huggingface.co/sentence-transformers/average_word_embeddings_komninos
BERT	https://huggingface.co/bert-base-uncased
SimCSE-BERT-unsup	https://huggingface.co/princeton-nlp/unsup-simcse-bert-base-uncased
SimCSE-BERT-sup	https://huggingface.co/princeton-nlp/sup-simcse-bert-base-uncased
coCondenser-msmarco	https://huggingface.co/sentence-transformers/msmarco-bert-co-condensor
Contriever	https://huggingface.co/nthakur/contriever-base-msmarco
SPECTER	https://huggingface.co/sentence-transformers/allenai-specter
LaBSE	https://huggingface.co/sentence-transformers/LaBSE
LASER2	https://github.com/facebookresearch/LASER
MiniLM-L6	https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2
MiniLM-L12	https://huggingface.co/sentence-transformers/all-MiniLM-L12-v2
MiniLM-L12-multilingual	https://huggingface.co/sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2
MPNet	https://huggingface.co/sentence-transformers/all-mpnet-base-v2
MPNet-multilingual	https://huggingface.co/sentence-transformers/paraphrase-multilingual-mpnet-base-v2
MiniLM-L12-multilingual	https://huggingface.co/sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2
SGPT-125M-nli	https://huggingface.co/Muennighoff/SGPT-125M-weightedmean-nli-bitfit
SGPT-5.8B-nli	https://huggingface.co/Muennighoff/SGPT-5.8B-weightedmean-nli-bitfit
SGPT-125M-msmarco	https://huggingface.co/Muennighoff/SGPT-125M-weightedmean-msmarco-specb-bitfit
SGPT-1.3B-msmarco	https://huggingface.co/Muennighoff/SGPT-1.3B-weightedmean-msmarco-specb-bitfit
SGPT-2.7B-msmarco	https://huggingface.co/Muennighoff/SGPT-2.7B-weightedmean-msmarco-specb-bitfit
SGPT-5.8B-msmarco	https://huggingface.co/Muennighoff/SGPT-5.8B-weightedmean-msmarco-specb-bitfit
SGPT-BLOOM-7.1B-msmarco	https://huggingface.co/bigscience/sgpt-bloom-7b1-msmarco
SGPT-BLOOM-1.7B-nli	https://huggingface.co/bigscience-data/sgpt-bloom-1b7-nli
GTR-Base	https://huggingface.co/sentence-transformers/gtr-t5-base
GTR-Large	https://huggingface.co/sentence-transformers/gtr-t5-large
GTR-XL	https://huggingface.co/sentence-transformers/gtr-t5-xl
GTR-XXL	https://huggingface.co/sentence-transformers/gtr-t5-xxl
ST5-Base	https://huggingface.co/sentence-transformers/sentence-t5-base
ST5-Large	https://huggingface.co/sentence-transformers/sentence-t5-large
ST5-XL	https://huggingface.co/sentence-transformers/sentence-t5-xl
ST5-XXL	https://huggingface.co/sentence-transformers/sentence-t5-xxl

Table 10: Publicly available model links used for evaluation

Dataset	Language	LASER2	LaBSE	MiniLM-L12-multilingual	MPNet-multilingual	SGPT-BLOOM-7.1B-msmarco
BUCC	de-en	99.21	99.35	97.11	98.59	54.00
BUCC	fr-en	98.39	98.72	94.99	96.89	97.06
BUCC	ru-en	97.62	97.78	95.06	96.44	45.30
BUCC	zh-en	97.70	99.16	95.63	97.56	97.96
Tatoeba	sqi-eng	97.22	96.76	98.17	98.57	10.38
Tatoeba	fry-eng	42.07	89.31	31.13	43.54	24.62
Tatoeba	kur-eng	19.09	83.59	46.94	61.44	8.26
Tatoeba	tur-eng	98.03	98.00	95.08	96.17	6.15
Tatoeba	deu-eng	99.07	99.20	97.02	97.73	70.10
Tatoeba	nld-eng	95.35	96.07	94.58	95.50	29.74
Tatoeba	ron-eng	96.52	96.92	95.30	96.43	27.23
Tatoeba	ang-eng	25.22	59.28	10.24	16.72	28.76
Tatoeba	ido-eng	80.86	89.42	40.25	43.91	43.91
Tatoeba	jav-eng	9.95	79.77	17.04	23.39	15.02
Tatoeba	isl-eng	94.32	94.75	24.07	59.25	6.29
Tatoeba	slv-eng	95.40	96.03	96.92	97.08	10.14
Tatoeba	cym-eng	5.85	92.00	13.25	22.31	6.97
Tatoeba	kaz-eng	53.30	87.49	34.89	61.49	3.32
Tatoeba	est-eng	96.43	96.55	97.33	98.40	4.76
Tatoeba	heb-eng	0.00	91.53	86.88	88.26	1.69
Tatoeba	gla-eng	1.52	85.66	3.61	4.72	2.09
Tatoeba	mar-eng	92.93	92.65	92.38	93.83	45.53
Tatoeba	lat-eng	64.81	80.07	19.47	24.25	28.76
Tatoeba	bel-eng	79.54	95.00	67.73	79.94	8.03
Tatoeba	pmo-eng	36.23	64.57	30.70	34.19	31.94
Tatoeba	gle-eng	4.20	93.80	11.62	16.85	3.26
Tatoeba	pes-eng	93.13	94.70	92.59	93.47	12.13
Tatoeba	nob-eng	95.77	98.40	97.73	98.53	21.07
Tatoeba	bul-eng	93.57	94.58	92.65	93.52	20.09
Tatoeba	chk-eng	77.17	79.44	55.37	58.68	64.63
Tatoeba	hun-eng	95.20	96.55	91.58	94.18	5.07
Tatoeba	uig-eng	56.49	92.40	24.39	48.35	1.27
Tatoeba	rus-eng	92.58	93.75	91.87	92.92	59.84
Tatoeba	spa-eng	97.33	98.40	95.42	97.00	94.48
Tatoeba	hye-eng	88.72	94.09	93.28	94.38	0.50
Tatoeba	tel-eng	96.72	97.86	36.40	79.73	64.62
Tatoeba	afr-eng	92.59	96.18	58.22	72.96	16.62
Tatoeba	mon-eng	3.42	95.91	95.04	96.14	2.85
Tatoeba	arz-eng	66.16	76.00	51.26	55.69	70.66
Tatoeba	hrv-eng	96.72	96.95	95.98	97.00	12.79
Tatoeba	nor-eng	60.02	74.38	47.99	50.23	52.23
Tatoeba	gsw-eng	27.52	46.50	25.74	25.12	21.03
Tatoeba	nds-eng	77.13	79.42	32.16	38.88	23.92
Tatoeba	ukr-eng	93.52	93.97	92.82	92.67	22.06
Tatoeba	uzb-eng	23.20	84.23	17.14	23.19	4.71
Tatoeba	lit-eng	96.20	96.47	93.16	95.37	4.49
Tatoeba	ina-eng	93.93	95.37	79.13	84.32	73.67
Tatoeba	lfn-eng	63.39	67.54	47.02	49.56	44.85
Tatoeba	zsm-eng	95.41	95.62	95.31	95.80	79.95
Tatoeba	ita-eng	94.32	92.72	93.05	93.76	65.04
Tatoeba	cmn-eng	85.62	95.10	94.93	95.83	91.45
Tatoeba	lvs-eng	95.33	95.88	97.87	97.53	6.55
Tatoeba	gle-eng	96.14	96.82	94.00	95.32	79.86
Tatoeba	ceb-eng	9.93	64.42	8.05	7.39	6.64
Tatoeba	bre-eng	31.2	15.07	5.56	6.42	4.67
Tatoeba	ben-eng	89.43	88.55	36.48	64.90	75.98
Tatoeba	swg-eng	33.10	59.36	26.31	22.80	16.89
Tatoeba	arq-eng	26.63	42.69	18.60	19.84	27.75
Tatoeba	kab-eng	65.88	4.31	1.16	1.41	1.69
Tatoeba	fra-eng	94.28	94.86	91.72	93.12	91.44
Tatoeba	por-eng	94.54	94.14	92.13	93.02	92.62
Tatoeba	tat-eng	34.74	85.92	10.25	10.89	3.59
Tatoeba	oci-eng	58.13	65.81	38.57	43.49	40.17
Tatoeba	pol-eng	97.32	97.22	94.28	96.95	14.09
Tatoeba	war-eng	8.25	60.29	7.25	7.42	10.38
Tatoeba	aze-eng	82.41	94.93	62.10	76.36	6.32
Tatoeba	vie-eng	96.73	97.20	95.12	97.23	94.20
Tatoeba	nno-eng	72.75	94.48	76.34	81.41	16.28
Tatoeba	cha-eng	14.86	31.77	15.98	12.59	23.26
Tatoeba	mhr-eng	6.86	15.74	6.89	7.57	1.56
Tatoeba	dan-eng	95.22	95.71	94.80	96.17	23.52
Tatoeba	ell-eng	96.20	95.35	95.43	94.93	5.34
Tatoeba	amh-eng	80.82	91.47	36.21	53.49	0.03
Tatoeba	pam-eng	3.24	10.73	5.41	5.39	5.85
Tatoeba	hsb-eng	45.75	67.11	36.10	44.32	9.68
Tatoeba	srp-eng	93.64	94.43	92.24	94.12	11.69
Tatoeba	epo-eng	96.61	98.20	41.73	55.12	26.20
Tatoeba	kzj-eng	4.46	11.33	6.24	5.88	5.17
Tatoeba	awa-eng	33.74	71.70	35.43	42.83	35.01
Tatoeba	fac-eng	57.04	87.40	27.51	38.24	12.61
Tatoeba	mal-eng	98.16	98.45	32.20	88.46	83.30
Tatoeba	ile-eng	87.88	85.58	57.71	60.36	59.59
Tatoeba	bos-eng	95.86	94.92	93.27	94.02	13.65
Tatoeba	cor-eng	4.45	10.11	3.42	3.53	2.83
Tatoeba	cat-eng	95.80	95.38	94.42	96.05	88.31
Tatoeba	eus-eng	93.32	95.01	23.18	31.33	53.38
Tatoeba	yue-eng	87.75	89.58	71.45	77.58	77.03
Tatoeba	swe-eng	95.31	95.63	94.42	95.45	19.53
Tatoeba	dtp-eng	7.39	10.85	5.69	5.03	3.41
Tatoeba	kat-eng	81.16	95.02	95.44	95.46	0.42
Tatoeba	jpn-eng	93.78	95.38	90.41	92.51	71.36
Tatoeba	ceb-eng	27.03	53.57	21.56	23.73	10.03
Tatoeba	sho-eng	4.68	91.55	4.52	6.53	5.51
Tatoeba	orv-eng	23.24	38.93	15.10	23.77	5.79
Tatoeba	ind-eng	92.98	93.66	92.74	93.50	88.04
Tatoeba	tuk-eng	16.35	75.27	15.16	14.91	5.48
Tatoeba	max-eng	36.96	63.26	45.25	48.77	36.14
Tatoeba	swb-eng	55.66	84.50	14.48	16.02	16.74
Tatoeba	hin-eng	95.32	96.87	97.62	97.75	85.23
Tatoeba	dsb-eng	42.34	64.81	33.43	36.85	8.78
Tatoeba	ber-eng	77.63	8.40	4.43	4.88	4.92
Tatoeba	tam-eng	87.32	89.0	24.64	73.60	72.76
Tatoeba	slk-eng	95.82	96.5	95.15	96.62	9.98
Tatoeba	tgl-eng	63.19	96.02	13.09	17.67	10.70
Tatoeba	ast-eng	76.35	90.68	62.17	70.08	71.13
Tatoeba	mkl-eng	93.63	93.6	91.00	93.02	10.47
Tatoeba	khm-eng	74.19	78.37	32.11	58.80	0.37
Tatoeba	ces-eng	95.52	96.68	95.12	95.73	9.55
Tatoeba	tzl-eng	36.56	58.88	25.46	34.21	27.82
Tatoeba	urd-eng	84.23	93.22	94.57	95.12	70.10
Tatoeba	ara-eng	90.14	88.80	87.93	90.19	85.37
Tatoeba	kor-eng	87.97	90.95	92.52	93.07	22.39
Tatoeba	yid-eng	2.49	88.79	14.38	30.73	0.16
Tatoeba	fin-eng	96.98	96.37	93.10	95.92	3.41
Tatoeba	tha-eng	96.38	96.14	96.72	95.99	2.22
Tatoeba	wuu-eng	75.09	90.18	76.00	78.25	79.58
Average	mix	67.42	81.75	57.98	63.38	31.08

Table 12: Multilingual bitext mining results. Scores are f1.

Dataset	Language	LASER2	LaBSE	MimLM-L12-multilingual	MPNet-multilingual	SGPT-BLOOM-7.1B-msmarco
AmazonCounterfactualClassification	de	67.82	73.17	68.35	69.95	61.35
AmazonCounterfactualClassification	ja	68.76	76.42	63.45	69.79	58.23
AmazonReviewsClassification	de	31.07	39.92	35.91	39.52	29.70
AmazonReviewsClassification	es	32.72	39.39	37.49	39.99	35.97
AmazonReviewsClassification	fr	31.12	38.52	35.30	39.00	35.92
AmazonReviewsClassification	ja	28.94	36.44	33.24	36.64	27.64
AmazonReviewsClassification	zh	30.89	36.45	35.26	37.74	32.63
MassiveIntentClassification	af	38.01	56.12	45.88	52.32	47.85
MassiveIntentClassification	am	12.70	55.71	36.75	41.55	33.30
MassiveIntentClassification	ar	37.16	50.86	45.14	51.43	59.25
MassiveIntentClassification	az	19.98	58.97	47.42	56.98	45.24
MassiveIntentClassification	bn	42.51	58.22	35.34	48.79	61.59
MassiveIntentClassification	cy	17.33	50.16	26.12	27.87	44.92
MassiveIntentClassification	da	45.61	58.25	57.73	62.77	51.23
MassiveIntentClassification	de	44.79	56.21	50.71	59.57	56.10
MassiveIntentClassification	el	46.71	57.03	58.70	62.62	46.13
MassiveIntentClassification	es	45.44	58.32	59.66	64.43	66.35
MassiveIntentClassification	fa	45.01	62.33	61.02	65.34	51.20
MassiveIntentClassification	fi	45.94	60.12	57.54	62.28	45.33
MassiveIntentClassification	fr	46.13	60.47	60.25	64.82	66.95
MassiveIntentClassification	he	42.55	56.55	52.51	58.21	43.18
MassiveIntentClassification	hi	40.20	59.40	58.37	62.77	63.54
MassiveIntentClassification	hu	42.77	59.52	60.41	63.87	44.73
MassiveIntentClassification	hy	28.07	56.20	51.60	57.74	38.13
MassiveIntentClassification	id	45.81	61.12	59.85	65.43	64.06
MassiveIntentClassification	is	39.86	54.90	30.83	37.05	44.35
MassiveIntentClassification	it	48.25	59.83	59.61	64.68	60.77
MassiveIntentClassification	ja	45.30	63.11	60.89	63.74	61.22
MassiveIntentClassification	ju	24.30	50.98	32.37	36.49	50.94
MassiveIntentClassification	ka	22.70	48.35	43.03	49.85	33.84
MassiveIntentClassification	km	22.48	48.55	40.04	45.47	37.34
MassiveIntentClassification	kn	4.32	56.24	40.98	50.63	53.54
MassiveIntentClassification	ko	44.26	60.99	50.30	61.82	53.36
MassiveIntentClassification	lv	39.75	57.10	54.68	61.29	46.50
MassiveIntentClassification	ml	41.33	57.91	42.41	54.34	58.27
MassiveIntentClassification	mn	16.20	58.50	51.77	56.59	40.28
MassiveIntentClassification	ms	43.23	58.60	54.76	60.70	59.65
MassiveIntentClassification	my	25.37	57.35	52.01	57.09	37.42
MassiveIntentClassification	nb	37.74	57.91	55.50	62.60	49.41
MassiveIntentClassification	nl	45.00	59.37	59.31	63.57	52.09
MassiveIntentClassification	pl	44.99	59.71	59.43	64.30	50.48
MassiveIntentClassification	pt	48.55	60.16	61.27	64.89	66.69
MassiveIntentClassification	ro	44.30	57.92	58.39	62.80	50.53
MassiveIntentClassification	ru	44.29	60.67	59.04	63.26	58.32
MassiveIntentClassification	sl	44.72	59.37	57.36	63.51	47.74
MassiveIntentClassification	sq	46.12	58.03	56.59	62.49	48.94
MassiveIntentClassification	sv	45.95	59.66	59.43	64.73	50.79
MassiveIntentClassification	sw	31.89	51.62	29.57	31.95	49.81
MassiveIntentClassification	ta	29.63	55.04	36.77	50.17	56.40
MassiveIntentClassification	te	36.03	58.32	40.22	52.82	49.41
MassiveIntentClassification	th	43.39	56.58	58.97	61.11	44.43
MassiveIntentClassification	tl	29.73	55.28	33.67	38.83	50.21
MassiveIntentClassification	tr	43.93	60.91	59.90	64.54	46.56
MassiveIntentClassification	ur	26.11	56.70	52.80	56.37	56.75
MassiveIntentClassification	vi	44.33	56.67	56.61	59.68	64.53
MassiveIntentClassification	zh-CN	40.62	63.86	61.99	65.33	67.07
MassiveIntentClassification	zh-TW	32.93	59.51	58.77	62.35	62.89
MassiveScenarioClassification	af	47.10	63.39	53.64	59.67	51.47
MassiveScenarioClassification	am	17.70	62.02	41.89	48.97	34.87
MassiveScenarioClassification	ar	45.21	57.72	51.74	57.78	65.21
MassiveScenarioClassification	az	28.20	63.48	52.06	61.53	45.58
MassiveScenarioClassification	bn	50.52	61.84	41.17	54.53	67.30
MassiveScenarioClassification	cy	22.58	56.13	31.72	35.26	46.29
MassiveScenarioClassification	da	54.87	65.24	66.87	71.00	53.52
MassiveScenarioClassification	de	54.34	62.39	57.40	67.34	61.74
MassiveScenarioClassification	el	55.47	64.58	66.14	68.81	48.96
MassiveScenarioClassification	es	52.77	63.61	65.04	70.42	73.34
MassiveScenarioClassification	fa	52.50	67.46	65.86	69.88	53.17
MassiveScenarioClassification	fi	52.63	64.58	63.75	67.60	44.69
MassiveScenarioClassification	fr	54.32	65.10	66.06	70.69	72.91
MassiveScenarioClassification	he	52.41	63.23	65.16	67.78	63.10
MassiveScenarioClassification	hi	47.37	64.40	65.21	67.92	69.27
MassiveScenarioClassification	hu	53.43	65.82	66.56	70.30	45.16
MassiveScenarioClassification	hy	33.57	61.25	56.11	63.02	38.73
MassiveScenarioClassification	id	54.38	65.84	66.16	70.73	70.13
MassiveScenarioClassification	is	49.78	61.94	37.52	44.16	44.21
MassiveScenarioClassification	it	54.84	64.09	65.00	69.73	65.57
MassiveScenarioClassification	ja	54.12	67.72	66.50	69.69	65.76
MassiveScenarioClassification	ju	32.71	58.29	38.60	44.20	54.79
MassiveScenarioClassification	ka	26.92	53.38	50.66	57.30	32.99
MassiveScenarioClassification	km	27.23	56.18	45.96	53.14	39.34
MassiveScenarioClassification	kn	10.06	61.74	45.73	56.08	60.50
MassiveScenarioClassification	ko	52.01	67.26	55.66	68.52	55.69
MassiveScenarioClassification	lv	44.82	61.87	59.80	66.28	44.35
MassiveScenarioClassification	ml	49.10	62.26	47.69	60.13	65.53
MassiveScenarioClassification	mn	21.51	62.60	57.07	60.85	38.72
MassiveScenarioClassification	ms	53.60	65.63	61.71	65.81	64.99
MassiveScenarioClassification	my	29.72	62.94	59.10	63.03	36.84
MassiveScenarioClassification	nb	43.90	64.29	64.25	70.24	51.80
MassiveScenarioClassification	nl	53.33	65.16	65.52	70.37	56.32
MassiveScenarioClassification	pl	52.92	64.56	65.04	68.99	49.98
MassiveScenarioClassification	pt	53.41	63.28	65.79	70.09	71.46
MassiveScenarioClassification	ro	50.48	62.41	64.17	67.95	53.69
MassiveScenarioClassification	ru	51.84	65.25	65.24	69.92	61.60
MassiveScenarioClassification	sl	51.29	64.25	64.01	70.81	48.04
MassiveScenarioClassification	sq	55.65	64.54	64.31	69.63	50.06
MassiveScenarioClassification	sv	54.64	66.01	67.14	71.60	51.73
MassiveScenarioClassification	sw	42.04	58.36	34.86	37.29	54.22
MassiveScenarioClassification	ta	36.72	59.08	42.62	55.96	62.77
MassiveScenarioClassification	te	42.08	64.13	46.46	58.81	62.59
MassiveScenarioClassification	th	52.15	64.34	67.01	69.44	45.18
MassiveScenarioClassification	tl	37.34	60.23	37.37	43.99	52.06
MassiveScenarioClassification	tr	52.56	65.43	66.55	70.4	47.21
MassiveScenarioClassification	ur	32.60	61.52	60.43	62.9	64.26
MassiveScenarioClassification	vi	50.97	61.05	60.72	65.71	70.61
MassiveScenarioClassification	zh-CN	50.22	70.85	67.44	71.23	73.95
MassiveScenarioClassification	zh-TW	42.32	67.08	65.70	68.73	70.30
MTOPODomainClassification	de	74.08	86.95	79.20	85.73	82.05
MTOPODomainClassification	es	73.47	84.07	83.04	86.96	93.55
MTOPODomainClassification	fr	72.26	84.14	78.63	81.21	90.98
MTOPODomainClassification	hi	72.95	85.11	81.36	84.76	89.33
MTOPODomainClassification	th	72.68	81.24	79.99	82.51	60.49
MTOPIntentClassification	de	51.62	63.42	54.23	61.27	61.92
MTOPIntentClassification	es	52.75	64.44	60.28	66.59	74.49
MTOPIntentClassification	fr	50.12	62.01	54.05	59.76	69.12
MTOPIntentClassification	hi	45.55	62.58	59.90	62.37	64.85
MTOPIntentClassification	th	50.07	64.61	61.96	64.80	49.36
Average	mix	42.85	60.77	54.87	60.39	54.4

Table 13: Multilingual classification results. Scores are accuracy.

Dataset	Language	Komninos	LASER2	LaBSE	MiniLM-L12-multilingual	MPNet-multilingual	SGPT-BLOOM-7.1B-msmarco
STS17	ko-ko	2.54	70.52	71.32	77.03	83.41	66.89
STS17	ar-ar	13.78	67.47	69.07	79.16	79.10	76.42
STS17	en-ar	9.08	65.05	74.51	81.22	80.85	78.07
STS17	en-de	-3.11	66.66	73.85	84.22	83.28	59.10
STS17	en-tr	-0.45	70.05	72.07	76.74	74.90	11.80
STS17	es-en	-8.18	55.30	65.71	84.44	86.11	78.22
STS17	es-es	48.23	79.67	80.83	85.56	85.14	86.00
STS17	fr-en	5.81	70.82	76.98	76.59	81.17	80.46
STS17	it-en	3.64	70.98	76.99	82.35	84.24	51.58
STS17	nl-en	-0.44	68.12	75.22	81.71	82.51	45.85
STS22	de	33.04	25.69	48.58	44.64	46.70	30.05
STS22	es	48.53	54.92	63.18	56.56	59.91	65.41
STS22	pl	12.47	18.34	39.30	33.74	33.65	31.13
STS22	tr	47.38	36.97	58.15	53.39	56.30	47.14
STS22	ar	32.42	42.57	57.67	46.2	52.19	58.67
STS22	ru	19.44	39.24	57.49	57.08	58.74	43.36
STS22	zh	4.78	49.41	63.02	58.75	61.75	66.78
STS22	fr	49.43	58.61	77.95	70.55	74.30	80.38
STS22	de-en	28.65	32.35	50.14	52.65	50.81	51.16
STS22	es-en	26.97	54.34	71.86	67.33	70.26	75.06
STS22	it	57.77	60.31	72.22	55.22	60.65	65.65
STS22	pl-en	45.55	53.63	69.41	69.02	73.07	53.31
STS22	zh-en	14.05	46.19	64.02	65.71	67.96	68.45
STS22	es-it	41.10	42.21	69.69	47.67	53.70	65.50
STS22	de-fr	14.77	37.41	53.28	51.73	62.34	53.28
STS22	de-pl	11.21	15.67	58.69	44.22	40.53	43.05
STS22	fr-pl	39.44	39.44	61.98	50.71	84.52	28.17
Average	mix	22.14	51.55	65.67	64.23	67.71	57.81

Table 14: Multilingual STS Results. Scores are Spearman correlations of cosine similarities.