

Multilingual Universal Sentence Encoder for Semantic Retrieval

Yinfei Yang^{a†}, Daniel Cer^{a†}, Amin Ahmad^a, Mandy Guo^a,
Jax Law^a, Noah Constant^a, Gustavo Hernandez Abrego^a, Steve Yuan^b, Chris Tar^a,
Yun-Hsuan Sung^a, Brian Strope^a, Ray Kurzweil^a

^aGoogle AI
Mountain View, CA

^bGoogle
Cambridge, MA

Abstract

We present easy-to-use retrieval focused multilingual sentence embedding models, made available on TensorFlow Hub. The models embed text from *16 languages* into a shared semantic space using a multi-task trained dual-encoder that learns tied cross-lingual representations via translation bridge tasks (Chidambaram et al., 2018). The models achieve *a new state-of-the-art in performance on monolingual and cross-lingual semantic retrieval (SR)*. Competitive performance is obtained on the related tasks of translation pair bitext retrieval (BR) and retrieval question answering (ReQA). On transfer learning tasks, our multilingual embeddings approach, and in some cases exceed, the performance of English only sentence embeddings.

1 Introduction

We introduce three new multilingual members in the *universal sentence encoder (USE)* (Cer et al., 2018) family of sentence embedding models. The models target performance on tasks that involve multilingual semantic similarity and achieve a new state-of-the-art in performance on monolingual and cross-lingual semantic retrieval (SR). One model targets efficient resource usage with a CNN model architecture (Kim, 2014). Another targets accuracy using the Transformer architecture (Vaswani et al., 2017). The third model provides an alternative interface to our multilingual Transformer model for use in retrieval question answering (ReQA). The *16 languages* supported by our multilingual models are given in Table 1.¹

† Corresponding authors:

{yinfeiy, cer}@google.com

¹Language coverage was selected based, in part, on the ease of obtaining data for the tasks used to train our models. Due to character set differences, we treat Simplified Chinese, zh, and Traditional Chinese, zh-tw, prominently used in Taiwan, as two languages within our model.

Languages	Family
Arabic (ar)	Semitic
Chinese (PRC) (zh) Chinese (Taiwan) (zh-tw)	Sino-Tibetan
Dutch(nl) English(en) German (de)	Germanic
French (fr) Italian (it) Portuguese (pt) Spanish (es)	Latin
Japanese (ja)	Japonic
Korean (ko)	Koreanic
Russian (ru) Polish (pl)	Slavic
Thai (th)	Kra-Dai
Turkish (tr)	Turkic

Table 1: Multilingual universal sentence encoder’s supported languages (ISO 639-1). Multilingual sentences are mapped to a shared semantic space.

2 Model Toolkit

Our multilingual models are implemented in TensorFlow (Abadi et al., 2016) and made publicly available on TensorFlow Hub.² Listing 1 illustrates the easy-to-use generation of multilingual sentence embeddings. The models conveniently only rely on TensorFlow without requiring additional libraries or packages. Listing 2 demonstrates using the question answering interface. Responses are encoded with additional context information such that the resulting context aware embeddings have a high dot product similarity score with the questions they answer. This allows for retrieval of indexed candidates using efficient nearest neighbor search.³

3 Encoder Architecture

3.1 Multi-task Dual Encoder Training

Similar to Cer et al. (2018) and Chidambaram et al. (2018), we target broad coverage using a

²<https://www.tensorflow.org/hub/>, Apache 2.0 license, with models available as saved TF graphs.

³Popular efficient search tools include FAISS <https://github.com/facebookresearch/faiss>, Annoy <https://github.com/spotify/annoy>, or FLANN <https://www.cs.ubc.ca/research/flann>.

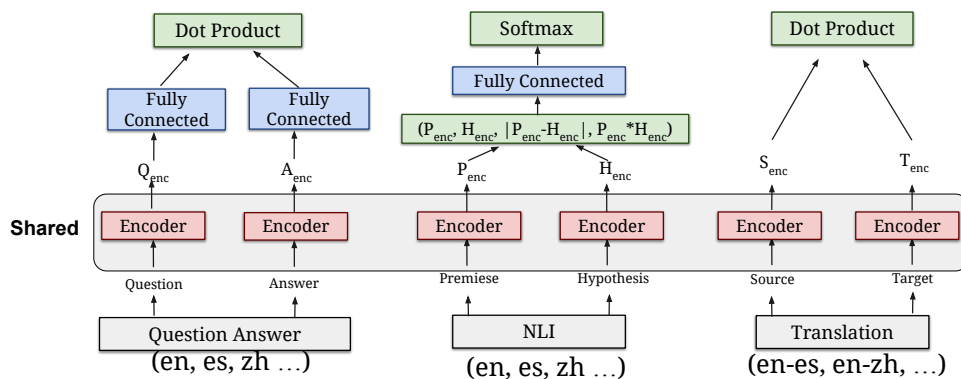


Figure 1: Multilingual universal sentence encoder model training architecture using multi-task training over: (i) retrieval question-answering (ReQA), natural language inference (NLI) and translation ranking. Transformer or CNN based sentence embedding models provide a shared encoder across all tasks.

```
import tensorflow_hub as hub

module = hub.Module("https://tfhub.dev/google/"
    "universal-sentence-encoder-multilingual/1")

multilingual_embeddings = module([
    "Hola Mundo!", "Bonjour le monde!", "Ciao mondo!"
    "Hello World!", "Hallo Welt!", "Hallo Wereld!",
    "你好世界!", "Привет, мир!", "!مرحبا بالعالم!"])
```

Listing 1: Python code mapping multilingual sentences into a shared semantic embedding space.

```
module = hub.Module("https://tfhub.dev/google/"
    "universal-sentence-encoder-multilingual-qa/1")

query_embeddings = module(
    dict(text=["What is your age?"],
        signature="question_encoder", as_dict=True))

candidate_embeddings = module(
    dict(text=["I am 20 years old."],
        context=["I will be 21 next year."],
        signature="response_encoder", as_dict=True))
```

Listing 2: Python code embedding a question and answer for retrieval Question-Answering (ReQA).

multi-task dual-encoder training framework, with a single shared encoder supporting multiple downstream tasks. The training tasks include: a multi-feature question-answer prediction task,⁴ a translation ranking task, and a natural language inference (NLI) task. Additional task specific hidden layers for the question-answering and NLI tasks are added after the shared encoder to provide representational specialization for each type of task. The model training architecture is illustrated at figure 1.

⁴Question-answer prediction is similar to conversational-response prediction (Yang et al., 2018). We treat the question as the conversational input and the answer as the response. For improved answer selection, we provide a bag-of-words (BoW) context feature as an additional input to the answer encoder. For our models, we use the entire paragraph containing the answer as context. The context feature is encoded using a separate DAN encoder.

3.2 SentencePiece

SentencePiece tokenization (Kudo and Richardson, 2018) is used for all of the 16 languages supported by our models.⁵ A single 128k SentencePiece vocabulary is trained from 8 million sentences sampled from our training corpus and balanced across the 16 languages. For validation, the vocab is used to process a development set, separately sampled from the sentence encoding model training corpus. We find the development set character coverage is higher than 99% for all languages, with less than 1% out-of-vocabulary tokens. Each token in the vocab is mapped to a fixed length embedding vector.⁶

3.3 Shared Encoder

Two distinct architectures for the sentence encoding models are provided: (i) transformer (Vaswani et al., 2017), targeted at higher accuracy at the cost of resource consumption; (ii) convolutional neural network (CNN) (Kim, 2014), designed for efficient inference but obtaining reduced accuracy.

Transformer The transformer encoding model embeds sentences using the *encoder* component of the transformer architecture (Vaswani et al., 2017). Bi-directional self-attention is used to compute context-aware representations of tokens in a sentence, taking into account both the ordering and the identity of the tokens. The context-aware token representations are then averaged together to obtain a sentence-level embedding.

CNN The CNN sentence encoding model feeds the input token sequence embeddings into a con-

⁵<https://github.com/google/sentencepiece>

⁶Out-of-vocabulary characters map to an <UNK> token.

Task Name	Task Type	Data Source	Native or Not
Retrieval Question-Answering (ReQA)	Ranking	Web Crawled	Native + MT
Translation Ranking	Ranking	Web Crawled	Native
Natural Language Inference (NLI)	3 way classification	Human Written	Native (en) + MT

Table 2: Training tasks for the multilingual sentence encoder. For better coverage across languages, we combine native text with machine translated (MT) data. For NLI, native data is only used for English (en).

Lang	QA	Translation	NLI
	Native + Translated	Native	Translated
ar	60M	158M	570K
de	75M	517M	570K
en	2.7B	-	570K
es	340M	416M	570K
fr	92M	586M	570K
it	103M	261M	570K
ja	384M	69M	570K
ko	60M	57M	570K
nl	60M	574M	570K
pt	180M	536M	570K
pl	60M	292M	570K
ru	112M	148M	570K
th	60M	70M	570K
tr	69M	415M	570K
zh	1B	112M	570K
zh-t	147M	112M	570K

Table 3: Training examples by task for each of the 16 languages understood by our models.

volutional neural network (Kim, 2014). Similar to the transformer encoder, average pooling is used to turn the token-level embeddings into a fixed-length representation. Sentence embeddings are then obtained by passing the averaged representation through additional feedforward layers.

4 Training and Configuration

4.1 Training Corpus

Training data consists of mined question-answer pairs,⁷ mined translation pairs,⁸ and the Stanford Natural Language Inference (SNLI) corpus (Bowman et al., 2015).⁹ SNLI only contains English data. The number of mined questions-answer pairs also varies across languages with a bias toward a handful of top tier languages. To balance training across languages, we use Google’s translation system to translate SNLI to the other 15 languages.

⁷QA pairs are mined from online forums and QA websites, including Reddit, StackOverflow, and YahooAnswers.

⁸The translation pairs are mined using a system similar to the approach described in Uszkoreit et al. (2010).

⁹MultiNLI (Williams et al., 2018), a more extensive corpus, contains examples from multiple sources but with different licences. Employing SNLI avoids navigating the licensing complexity of using MultiNLI to training public models.

Model	Quora	AskUbuntu	Average
USE _{Trans}	89.1	42.3	65.7
USE _{CNN}	89.2	39.9	64.6
Gillick et al. (2018)	87.5	37.3	62.4

Table 4: MAP@100 on SR (English). Models are compared with the best models from Gillick et al. (2018) that exclude in-domain training data.

We also translate a portion of question-answer pairs to ensure each language has a minimum of 60M training pairs. For each of our datasets, we use 90% of the data for training, and the remaining 10% for development/validation. Table 2 and 3 lists the details of data used for each task / language.

4.2 Model Configuration

Input sentences are truncated to 256 tokens for the CNN model and 100 tokens for transformer. The CNN encoder uses 2 CNN layers with filter width of [1, 2, 3, 5] and 256 filters per width. The Transformer encoder employs 6 transformer layers, with 8 attentions heads, hidden size 512, and filter size 2048. Similar to our prior work (Cer et al., 2018), we configure our models with the intention of making them small and fast enough to be used directly within many downstream applications without the need for model distillation. Model hyperparameters are tuned on development data sampled from the same sources as the training data. We export sentence encoding modules for our two encoder architectures: USE_{Trans} and USE_{CNN}. We also export a larger graph for QA tasks from our Transformer based model that includes QA specific layers and support providing context information from the larger document as USE_{QA Trans+Cxt}.¹⁰

5 Experiments on Retrieval Tasks

In this section we evaluate our multilingual encoding models on semantic retrieval, bitext and

¹⁰While USE_{QA Trans+Cxt} uses the same underlying shared encoder as USE_{Trans} but with additional task specific layers, we anticipate that the models could diverge in the future.

Model	en-es	en-fr	en-ru	en-zh
USE _{Trans}	86.1	83.3	88.9	78.8
USE _{CNN}	85.8	82.7	87.4	79.5
Yang et al. (2019)	89.0	86.1	89.2	87.9

Table 5: P@1 on UN translation pair bitext retrieval (BR). Yang et al. (2019) is a specialized translation retrieval model and the current state-of-the-art.

retrieval question answer tasks.

5.1 Semantic Retrieval (SR)

Following Gillick et al. (2018), we construct semantic retrieval (SR) tasks from the Quora question-pairs (Hoogeveen et al., 2015) and AskUbuntu (Lei et al., 2016) datasets. The SR task is to identify all sentences in the retrieval corpus that are semantically similar to a query sentence.¹¹

For each dataset, we first build a graph connecting each of the positive pairs, and then compute its transitive closure. Each sentence then serves as a test query that should retrieve all of the other sentences it is connected to within the transitive closure. Mean average precision (MAP) is employed to evaluate the models. More details on the constructed datasets can be found in Gillick et al. (2018). Both datasets are English only.

Table 4 shows the MAP@100 on the Quora and AskUbuntu retrieval tasks. We use Gillick et al. (2018) as the baseline model, which is trained using a similar dual encoder architecture. The numbers provided here are for models without focused in-domain training data.¹² Both USE_{CNN} and USE_{Trans} outperform the prior state-of-the-art. USE_{Trans} and USE_{CNN} perform comparably on Quora. However, USE_{Trans} performs notably better than USE_{CNN} on AskUbuntu, suggesting the AskUbuntu data could be more challenging.

5.2 Bitext Retrieval (BR)

Bitext retrieval performance is evaluated on the United Nation (UN) Parallel Corpus (Ziems et al., 2016), containing 86,000 bilingual document pairs matching English (en) documents with their translations in five other languages: French (fr),

¹¹The task is related to paraphrase identification (Dolan et al., 2004) and Semantic Textual Similarity (STS) (Cer et al., 2017), but with the identification of meaning similarity being assessed in the context of a retrieval task.

¹²The model for Quora is trained on Paralex (<http://knowitall.cs.washington.edu/paralex>) and AskUbuntu data. The model for AskUbuntu is trained on Paralex and Quora.

Spanish (es), Russian (ru), Arabic (ar) and Chinese (zh). Document pairs are aligned at the sentence-level, which results in 11.3 million aligned sentence pairs for each language pair.

Table 5 shows sentence-level retrieval precision@1 (P@1) for the proposed models as well as the current state-of-the-art results from Yang et al. (2019), which uses a specialized translation pair retrieval model. USE_{Trans} is generally better than USE_{CNN}, performing lower than the SOTA but not by too much with the exception of en-zh.¹³

Model	SQuAD Dev	SQuAD Train
<i>Paragraph Retrieval</i>		
USE _{QA Trans+Cxt}	63.5	53.3
BM25 (baseline)	61.6	52.4
<i>Sentence Retrieval</i>		
USE _{QA Trans+Cxt}	53.2	43.3
USE _{Trans}	47.1	37.2

Table 6: P@1 for SQuAD ReQA. Models are *not trained* on SQuAD. Dev and Train only refer to the respective sections of the SQuAD dataset.

5.3 Retrieval Question Answering (ReQA)

Similar to the data set construction used for the SR tasks, the SQuAD v1.0 dataset (Rajpurkar et al., 2016) is transformed into a retrieval question answering (ReQA) task.¹⁴ We first break all documents in the dataset into sentences using the sentence splitter distributed with the ReQA evaluation suite.¹⁵ Each question of the (question, answer spans) tuples in the dataset is treated as a query. The task is to retrieve the sentence designated by the tuple answer span. Search is performed on a retrieval corpus consisting of all of the sentences within the corpus. We contrast sentence and paragraph-level retrieval using our models, with the later allowing for comparison against a BM25 baseline (Jones et al., 2000).¹⁶

¹³Performance is degraded from Yang et al. (2019) due to using a single sentencepiece vocabulary to cover 16 languages. Languages like Chinese, Korean, Japanese have much more characters. To ensure the vocab coverage, sentencepiece tends to split the text of these languages into single characters, which increases the difficulty of the task.

¹⁴The retrieval question answering task was suggested by Chen et al. (2017) and then recently explored further by Cakaloglu et al. (2018). However, Cakaloglu et al. (2018)’s use of sampling makes it difficult to directly compare with their results and we provide our own baseline based on BM25.

¹⁵<https://github.com/google/retrieval-qa-eval>

¹⁶BM25 is a strong baseline for text retrieval tasks. Paragraph-level experiments use the BM25 implementa-

Model	en	ar	de	es	fr	it	ja	ko	nl	pt	pl	ru	th	tr	zh / zh-t
<i>Cross-lingual Semantic Retrieval (cl-SR)</i>															
Quora															
USE _{Trans}	89.1	83.1	85.5	86.3	86.7	86.8	85.1	82.5	83.8	86.5	82.1	85.7	85.8	82.5	84.8
USE _{CNN}	89.2	79.9	83.7	85.0	85.0	85.5	82.4	77.6	81.3	85.2	78.3	83.8	83.5	79.9	81.9
LASER		79.7	82.2	83.5	83.1	83.7	-	73.4	82.8	83.6	82.3	82.6	78.6	79.9	-
AskUbuntu															
USE _{Trans}	42.3	38.2	40.0	39.9	39.3	40.2	40.6	40.3	39.5	39.8	38.4	39.6	40.3	37.7	40.1
USE _{CNN}	39.9	33.0	35.0	35.6	35.2	36.1	35.5	35.1	34.5	35.6	32.9	35.2	35.2	32.8	34.6
LASER		24.5	26.1	26.4	26.5	27.0	-	22.0	26.2	26.2	25.7	25.6	23.8	25.0	-
Average															
USE _{Trans}	65.7	60.7	62.8	63.1	63.0	63.5	63.8	62.4	61.7	63.2	60.7	62.7	63.1	60.1	62.5
USE _{CNN}	64.6	56.5	59.4	60.3	60.1	60.8	59.0	56.4	57.9	60.4	55.6	59.5	59.4	56.4	58.3
LASER		52.1	54.2	55.0	54.8	55.4	-	47.7	54.5	54.9	54.0	54.6	51.2	52.5	-
<i>Cross-lingual Retrieval Question Answering (cl-ReQA)</i>															
SQuAD train															
USE _{QA Trans+Cxt}	43.3	33.2	35.2	37.2	37.0	37.0	32.9	31.1	36.6	37.7	34.5	33.2	36.9	32.3	32.7

Table 7: Cross-lingual performance on Quora/AskUbuntu cl-SR (MAP) and SQuAD cl-ReQA (P@1). Queries/questions are machine translated, while retrieval candidates remain in English.

We evaluated ReQA using the SQuAD dev and train sets and without training on the SQuAD data.¹⁷ The sentence and paragraph retrieval P@1 are shown in table 6. For sentence retrieval, we compare encodings produced using context from the text surrounding the retrieval candidate, USE_{QA Trans+Cxt}, to sentence encodings produced without contextual cues, USE_{Trans}. Paragraph retrieval contrasts USE_{QA Trans+Cxt} with BM25.

5.4 Cross-lingual Retrieval

Our English retrieval experiments are extended to explore cross-lingual semantic retrieval (cl-SR) and cross-lingual retrieval question answering (cl-ReQA). SR queries and ReQA questions are machine translated into other languages, while keeping the retrieval candidates in English.¹⁸ Table 7 provides our cross-lingual retrieval results for our transformer and CNN multilingual sentence encoding models. We compare against the state-of-the-art LASER multilingual sentence embedding

tion: <https://github.com/nhirakawa/BM25>, with default parameters. We exclude sentence-level BM25, as BM25 generally performs poorly at this granularity.

¹⁷For sentences, the resulting retrieval task for development set consists of 11,425 questions and 10,248 candidates, and the retrieval task for train set is consists of 87,599 questions and 91,703 candidates. For paragraph retrieval, there are 2,067 retrieval candidates in the development set and 18,896 in the training set. To retrieve paragraphs with our model, we first run sentence retrieval and use the retrieved nearest sentence to select the enclosing paragraph.

¹⁸Poor translations are detected and rejected when the original English text and English back translation have a cosine similarity < 0.5 according our previously released English USE_{Trans} model (Cer et al., 2018).

library (Artetxe and Schwenk, 2019).¹⁹

On both the Quora and AskUbuntu cl-SR tasks, USE_{Trans} outperforms USE_{CNN} and LASER on all datasets, except the Polish (pl) Quora data where LASER achieves slightly better performance.²⁰ USE_{CNN} tends to outperform LASER on Quora and always outperforms LASER by a sizable margin on AskUbuntu. We note that our CNN based model not only outperforms LASER, but also relies on simpler model architecture than LASER’s LSTM based architecture. Given the similar level of performance on Quora between USE_{CNN} and LASER, we suspect the notably better performance on AskUbuntu over LASER is due to differences in the training data provided to encoding models.

6 Experiments on Transfer Tasks

For comparison with prior USE models, English task transfer performance is evaluated on SentEval (Conneau and Kiela, 2018). For sentence classification transfer tasks, the output of the sentence encoders are provided to a task specific DNN. For the pairwise semantic similarity task, the similarity of sentence embeddings u and v is assessed using $-\arccos\left(\frac{uv}{\|u\|\|v\|}\right)$, following Yang et al. (2018). In table 8, our multilingual models show competitive transfer performance when compared to state-of-the-art sentence embedding models. USE_{Trans} performs better than USE_{CNN} on all tasks. Our new

¹⁹<https://github.com/facebookresearch/LASER>

²⁰Results are not presented for LASER on ja and zh due unicode errors.

Model	MR	CR	SUBJ	MPQA	TREC	SST	STS Bench (dev / test)
USE multilingual models							
USE _{CNN}	73.8	83.2	90.1	87.7	96.4	78.1	0.829 / 0.809
USE _{Transformer}	78.1	87.0	92.1	89.9	96.6	80.9	0.837 / 0.825
The state-of-the-art English embedding models							
InferSent (Conneau et al., 2017)	81.1	86.3	92.4	90.2	88.2	84.6	0.801 / 0.758
Skip-Thought LN (Ba et al., 2016)	79.4	83.1	93.7	89.3	–	–	–
Quick-Thought (Logeswaran and Lee, 2018)	82.4	86.0	94.8	90.2	92.4	87.6	–
USE _{DAN} for English (Cer et al., 2018)	72.2	78.5	92.1	86.9	88.1	77.5	0.760 / 0.717
USE _{Transformer} for English (Cer et al., 2018)	82.2	84.2	95.5	88.1	93.2	83.7	0.802 / 0.766

Table 8: Performance on English transfer tasks from SentEval (Conneau and Kiela, 2018).

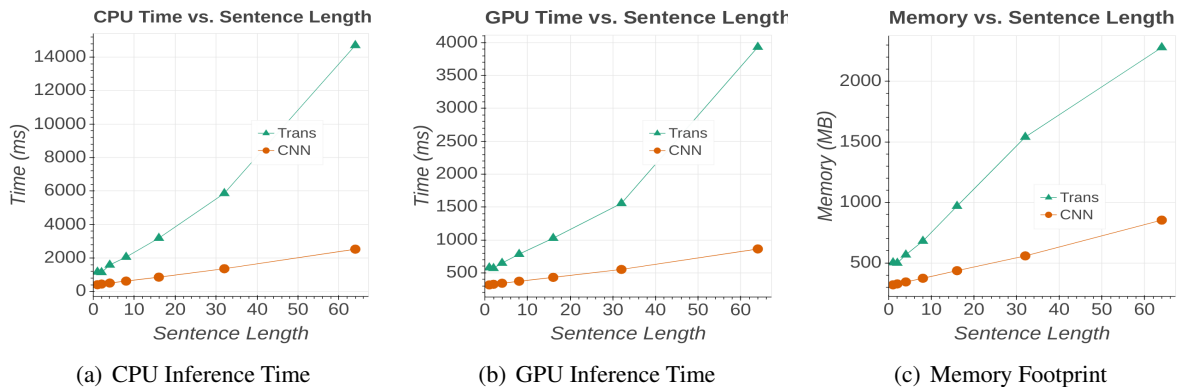


Figure 2: Resource usage for the multilingual Transformer and CNN encoding models.

multilingual USE_{Trans} model outperforms our best previously released English only model, USE_{Trans} for English (Cer et al., 2018), on some tasks.

7 Resource Usage

Figure 2 provides compute and memory usage benchmarks for our models.²¹ Inference times on GPU are 2 to 3 times faster than CPU. Our CNN models have the smallest memory footprint and are the fastest on both CPU and GPU. The memory requirements increase with sentence length, with the Transformer model increasing more than twice as fast as the CNN model.²² While this makes CNNs an attractive choice for efficiently encoding longer texts, this comes with a corresponding drop in accuracy on many retrieval and transfer tasks.

8 Conclusion

Easy-to-use retrieval focused multilingual models for embedding sentence-length text are made avail-

²¹CPU benchmarks are run on Intel(R) Xeon(R) Platinum 8173M CPU @ 2.00GHz. GPU benchmarks were run on an Nvidia v100. Memory footprint was measured on CPU.

²²Transformer models are ultimately governed by a time and space complexity of $O(n^2)$. The benchmarks show for shorter sequence lengths the time and space requirements are dominated by computations that scale linearly with length and have a larger constant factor than the quadratic terms.

able on TensorFlow Hub. Our models embed text from 16 languages into a shared semantic embedding space and *achieve a new state-of-the-art in performance on monolingual and cross-lingual semantic retrieval (SR)*. The models achieve good performance on the related tasks of translation pair bi-text retrieval (BR) and retrieval question answering (ReQA). Monolingual transfer task performance approaches, and in some cases exceeds, English only sentence embedding models. Our models are freely available under an Apache license with additional documentation and tutorial colaboratory notebooks at:

<https://tfhub.dev/s?q=universal-sentence-encoder-multilingual>

Acknowledgments

We thank our teammates from Descartes and other Google groups for their feedback and suggestions. Special thanks goes to Muthu Chidambaram for his early exploration of multilingual training, Taku Kudo for the SentencePiece model support, Chen Chen for the templates used to perform the transfer learning experiments and Mario Guajardo for an early version of the ReQA tutorial Colab.

References

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. *Tensorflow: A system for large-scale machine learning*. In *Proceedings of USENIX OSDI'16*, OSDI'16, pages 265–283.
- Mikel Artetxe and Holger Schwenk. 2019. *Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond*. *Transactions of the Association for Computational Linguistics*, 7:597–610.
- Lei Jimmy Ba, Ryan Kiros, and Geoffrey E. Hinton. 2016. *Layer normalization*. *CoRR*, abs/1607.06450.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. *A large annotated corpus for learning natural language inference*. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642.
- Tolgahan Cakaloglu, Christian Szegedy, and Xiaowei Xu. 2018. *Text embeddings for retrieval from a large knowledge base*. *CoRR*, abs/1810.10176.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. *SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation*. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14.
- Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Brian Strope, and Ray Kurzweil. 2018. *Universal sentence encoder for English*. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 169–174.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. *Reading Wikipedia to answer open-domain questions*. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879.
- Muthuraman Chidambaram, Yinfei Yang, Daniel Cer, Steve Yuan, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. *Learning cross-lingual sentence representations via a multi-task dual-encoder model*. *CoRR*, abs/1810.12836.
- Alexis Conneau and Douwe Kiela. 2018. *SentEval: An evaluation toolkit for universal sentence representations*. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*.
- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. 2017. *Supervised learning of universal sentence representations from natural language inference data*. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 670–680, Copenhagen, Denmark. Association for Computational Linguistics.
- Bill Dolan, Chris Quirk, and Chris Brockett. 2004. *Unsupervised construction of large paraphrase corpora: Exploiting massively parallel news sources*. In *COLING 2004: Proceedings of the 20th International Conference on Computational Linguistics*, pages 350–356.
- Daniel Gillick, Alessandro Presta, and Gaurav Singh Tomar. 2018. *End-to-end retrieval in continuous space*. *CoRR*, abs/1811.08008.
- Doris Hoogeveen, Karin M. Verspoor, and Timothy Baldwin. 2015. *Cquadupstack: A benchmark data set for community question-answering research*. In *Proceedings of the 20th Australasian Document Computing Symposium, ADCS '15*, pages 3:1–3:8.
- K. Sparck Jones, S. Walker, and S. E. Robertson. 2000. *A probabilistic model of information retrieval: Development and comparative experiments*. *Inf. Process. Manage.*, 36(6):779–808.
- Yoon Kim. 2014. *Convolutional neural networks for sentence classification*. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751.
- Taku Kudo and John Richardson. 2018. *SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing*. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71.
- Tao Lei, Hrishikesh Joshi, Regina Barzilay, Tommi Jaakkola, Kateryna Tymoshenko, Alessandro Moschitti, and Lluís Màrquez. 2016. *Semi-supervised question retrieval with gated convolutions*. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1279–1289.
- Lajanugen Logeswaran and Honglak Lee. 2018. *An efficient framework for learning sentence representations*. In *International Conference on Learning Representations (ICLR)*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. *SQuAD: 100,000+ questions for machine comprehension of text*. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392.
- Jakob Uszkoreit, Jay Ponte, Ashok Papat, and Moshe Dubiner. 2010. *Large scale parallel document mining for machine translation*. In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 1101–1109.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. *Attention is all you need*. In *Proceedings of NIPS*, pages 6000–6010.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. *A broad-coverage challenge corpus for sentence understanding through inference*. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122.
- Yinfei Yang, Gustavo Hernández Ábrego, Steve Yuan, Mandy Guo, Qinlan Shen, Daniel Cer, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2019. *Improving multilingual sentence embedding using bi-directional dual encoder with additive margin softmax*. *CoRR*, abs/1902.08564.

Yinfei Yang, Steve Yuan, Daniel Cer, Sheng-yi Kong, Noah Constant, Petr Pilar, Heming Ge, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. [Learning semantic textual similarity from conversations](#). In *Proceedings of The Third Workshop on Representation Learning for NLP*, pages 164–174.

Michał Ziemski, Marcin Junczys-Dowmunt, and Bruno Poulliquen. 2016. [The united nations parallel corpus v1.0](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, pages 3530–3534.