

THUMT: An Open-Source Toolkit for Neural Machine Translation

The Tsinghua Natural Language Processing Group

December, 2017

1 Introduction

Machine translation, which investigates the use of computer to translate human languages automatically, is an important task in natural language processing and artificial intelligence. With the availability of bilingual, machine-readable texts, data-driven approaches to machine translation have gained wide popularity since 1990s. Recent several years have witnessed the rapid development of end-to-end neural machine translation (NMT) (Sutskever et al., 2014; Bahdanau et al., 2015). Capable of learning representations from data, NMT has quickly replaced conventional statistical machine translation (SMT) (Brown et al., 1993; Koehn et al., 2003; Chiang, 2005) to become the new de facto method in practical MT systems (Wu et al., 2016).

On top of [TensorFlow](#), THUMT is an open-source toolkit for neural machine translation developed by the Tsinghua Natural Language Processing Group. THUMT has the following features:

1. *Attention-based translation model.* THUMT implements the standard attention-based encoder-decoder framework for NMT (Bahdanau et al., 2015) as well as the latest Transformer architecture (Vaswani et al., 2017).
2. *Minimum risk training.* Besides standard maximum likelihood estimation (MLE), THUMT also supports minimum risk training (MRT) (Shen et al., 2016) that aims to find a set of model parameters that minimize the expected loss calculated using evaluation metrics such as BLEU (Papineni et al., 2002) on the training data.

2 Installation

2.1 System Requirements

THUMT supports Linux i686 and Mac OSX. The following third-party software is required to install THUMT:

1. Python version 2.7.0 or higher.

2. JRE 1.6 or higher (optional, only used for visualization).

2.2 Installing Prerequisites

We recommend using pip to install the prerequisites of THUMT. The installation starts with python-pip:

```
1 apt-get install python-pip
```

Then, run the following two commands to install argparse and TensorFlow (version $\geq 1.4.0$):

```
1 pip install argparse
2 pip install tensorflow-gpu
```

2.3 Installing THUMT

The source code of THUMT is available both at [the toolkit website](#) (stable release) and GitHub (latest version).

Here is a brief guide on how to install THUMT.

2.3.1 Step1: Unpacking

Unpack the package using the following command:

```
1 tar xvfz THUMT.tar.gz
```

Entering the THUMT folder, you may find two folders (`thumt`, `data`) and three files (`LICENSE`, `UserManual.en.pdf` and `UserManual.zh.pdf`):

1. `thumt`: the source code.
2. `data`: toy training, validation, and test datasets.
3. `docs`: source of the documents.
4. `LICENSE`: license statement.
5. `UserManual.en.pdf`: this document.
6. `UserManual.zh.pdf`: the Chinese version of this document.

2.3.2 Step 2: Modifying Environment Variables

We highly recommend running THUMT on GPU servers. Suppose THUMT runs on NVIDIA GPUs with [the CUDA toolkit](#) version 8.0 installed. Users need to set environment variables to enable the GPU support:

```
1 export PATH=/usr/local/cuda/bin:$PATH
2 export LD_LIBRARY_PATH=/usr/local/cuda/lib64:$LD_LIBRARY_PATH
3 export PYTHONPATH=/PATH/TO/THUMT:$PYTHONPATH
```

To set these environment variable permanently for all future bash sessions, users can simply add the above two lines to the `.bashrc` file in your `$HOME` directory, where `/PATH/TO/THUMT` is the path of THUMT.

3 User Guide

3.1 Data Preparation

Running THUMT involves three types of datasets:

1. *Training set*: a set of parallel sentences used for training NMT models.
2. *Validation set*: a set of source sentences paired with single or multiple target translations used for model selection and hyper-parameter optimization.
3. *Test set*: a set of source sentences paired with single or multiple target translations used for evaluating translation performance on unseen texts.

3.1.1 Training Set

A training set is used for training NMT models. It often consists of two files: one file for source sentences and the other for corresponding target sentences. In the `data` folder, there is an example source file `train.src` of a toy training set that contains seven sentences¹:

```
1 wo hen xihuan yinyue .
2 wo bu xihuan huahua .
3 ni xihuan yinyue me ?
4 shide , wo ye xihuan yiyue .
5 ta ye xihuan yinyue .
6 ta yidian dou bu xihuan yiyue .
7 ta hen xinhuan huahua .
```

The corresponding target file `train.trg` is shown below:

```
1 i like music very much .
2 i do not like painting .
3 do you like music ?
4 yes , i like music too .
5 he also likes music .
6 she does not like music at all .
7 she likes painting very much .
```

Note that each line in the source and target files only contains one tokenized sentence. The source and target sentences with the same line number are translationally equivalent.

¹The Chinese text is romanized for readability.

Our toy training set only contains seven sentence pairs. In practice, NMT often requires millions of sentence pairs to achieve reasonable translation performance.

3.1.2 Validation Set

A validation set is used for model selection and hyper-parameter optimization. During training, THUMT evaluates intermediate models on the validation set periodically. The model that obtains the highest BLEU score on the validation set is chosen as the final learned model. A validation set consists of one source file and one or more target files. There is an example source file `valid.src` of a toy validation set in the `data` folder:

```
1 wo hen xihuan huahua .  
2 wo bu xihuan yinyue .
```

In our toy validation set, there is only one target file `valid.trg` that contains the reference translations of `valid.src`:

```
1 i like painting very much .  
2 i do not like music .
```

Due to the diversity of natural languages, one source sentence often corresponds to multiple reference translations. THUMT also supports multiple references.

3.1.3 Test set

A test set is used for evaluating the learned NMT model on unseen source text. Like a validation set, a test set also contains one source file and one or more target files. It shares the same naming scheme for multiple references with the validation set.

In the `data` folder, there is one source file `test.src`:

```
1 ta xihuan huahua me ?  
2 ta yidian dou bu xihuan huahua .
```

The corresponding target file `test.trg` is shown below:

```
1 does she like painting ?  
2 he does not like painting at all .
```

In our toy data, validation and test sets contain only two sentences. In practice, thousands of sentences are needed for both validation and test sets.

3.2 Training

3.3 Data preprocessing

The most common approach to achieve open vocabularies is to use the Byte Pair Encoding (BPE). The codes of BPE can be found at [Subword-NMT](#).

To encode the training corpora using BPE, you need to generate BPE operations first. The following command will create a file named `bpe32k`, which contains 32k BPE operations. It also outputs two dictionaries named `vocab.src` and `vocab.trg`.

```
1 python subword-nmt/learn_joint_bpe_and_vocab.py
2   --input train.src train.trg -s 32000 -o bpe32k
3   --write-vocabulary vocab.src vocab.trg
```

You also need to encode the training corpora, validation set and test set using the generated BPE operations and dictionaries.

```
1 python subword-nmt/apply_bpe.py
2   --vocabulary vocab.src
3   --vocabulary-threshold 50
4   -c bpe32k < train.src > train.32k.src
5 python subword-nmt/apply_bpe.py
6   --vocabulary vocab.trg
7   --vocabulary-threshold 50
8   -c bpe32k < train.trg > train.32k.trg
9 python subword-nmt/apply_bpe.py
10  --vocabulary vocab.src
11  --vocabulary-threshold 50
12  -c bpe32k < valid.src > valid.32k.src
13 python subword-nmt/apply_bpe.py
14  --vocabulary vocab.trg --vocabulary-threshold 50
15  -c bpe32k < valid.trg > valid.32k.trg
16 python subword-nmt/apply_bpe.py
17  --vocabulary vocab.src
18  --vocabulary-threshold 50
19  -c bpe32k < test.src > test.32k.src
```

The original segmentation can be restored with the following command:

```
1 sed -r 's/(@@ )|(@@ ?$)//g' < input > output
```

3.4 Shuffling Corpora

It is helpful to shuffle training corpora first and we provide `shuffle_corpus.py` script to achieve this functionality, which has the following options

```
1 usage: shuffle_corpus.py [-h] --corpus CORPUS [CORPUS ...]
2                        [--suffix SUFFIX]
3                        [--seed SEED]
4
5 Shuffle corpus
6
7 optional arguments:
8   -h, --help            show this help message and exit
```

```

9  --corpus CORPUS [CORPUS ...]
10                               input corpora
11  --suffix SUFFIX             Suffix of output files
12  --seed SEED                  Random seed

```

The meaning of arguments is shown below:

1. **corpus**: the input corpus.
2. **--suffix** the suffix of output filename.
3. **--seed**: the random seed.
4. **--help** show the help message.

3.4.1 Generating Vocabularies

We have to build vocabularies before training an NMT model. The `build_vocab.py` in the `scripts` folder can be used to create vocabularies. It has the following options:

```

1  usage: build_vocab.py [-h] [--limit LIMIT] [--control CONTROL]
2                        corpus output
3
4  Create vocabulary
5
6  positional arguments:
7    corpus                input corpus
8    output                Output vocabulary name
9
10 optional arguments:
11  -h, --help            show this help message and exit
12  --limit LIMIT         Vocabulary size
13  --control CONTROL     Add control symbols to vocabulary.
14                        Control symbols are separated by comma.

```

The meaning of arguments is shown below:

1. **corpus**: path to source or target corpus.
2. **output**: the name of output vocabulary.
3. **--limit**: control the size of vocabulrary.
4. **--control**: a comma separated string represents the control symbols. For example, the string “<pad>,<eos>,<unk>” will add three control symbols to the beginning of the vocabulary.
5. **--help**: show the help message.

3.4.2 The Python Script for Training

THUMT implements the traditional sequence to sequence model SEQ2SEQ (Sutskever et al., 2014), standard attention-based encoder-decoder framework for neural machine translation RNNSEARCH (Bahdanau et al., 2015) as well as the attention-based Transformer architecture (Vaswani et al., 2017).

The `trainer.py` script in the `bin` folder is used for training NMT models. Its arguments are listed as follows:

```
1 usage: trainer.py [<args>] [-h | --help]
2
3 Training neural machine translation models
4
5 optional arguments:
6   -h, --help            show this help message and exit
7   --input INPUT INPUT  Path of source and target corpus
8   --output OUTPUT      Path to saved models
9   --vocabulary VOCABULARY
10                        Path of source and target vocabulary
11   --validation VALIDATION
12                        Path of validation file
13   --references REFERENCES [REFERENCES ...]
14                        Path of reference files
15   --model MODEL        Name of the model
16   --parameters PARAMETERS
17                        Additional hyper parameters
```

We distinguish between *required* and *optional* arguments. Users must specify the following required arguments to run the `trainer.py` script:

1. `--model`: the name of model architecture. It can be `seq2seq`, `rnnsearch` or `transformer`
2. `--input`: the path of source and target corpora.
3. `--vocabulary`: the path of source and target vocabularies.

The optional arguments of the `trainer.py` script can be omitted in a command. If an optional argument has a default value, the default value will be used in training if the argument is omitted in the command-line argument list. These optional arguments are listed as follows:

1. `--output`: the path to save checkpoints. The default value is `train` directory in the current path.
2. `--validation`: the path to the validation set. The best model will be saved in `train/eval` directory.
3. `--references`: the path to references. Multiple references are also supported.

4. `--parameters`: other optional parameters which will be described in the next subsection. The parameters is a string containing comma-separated `name=value` pairs. Please refer to the document of `tf.contrib.training.parse_values`.
5. `--helpshow` this help message.

3.5 General Parameters

1. `num_threads`: the number of threads used in data processing.
2. `buffer_size`: the buffer size used in data processing.
3. `batch_size`: the batch size used in training stage.
4. `constant_batch_size`: if set to true, then each batch will contains `batch_size` sentences. Otherwise, each batch will contains approximately `batch_size` tokens.
5. `max_length`: the maximum length of training sentences.
6. `train_steps`: the total number of steps in the training stage.
7. `save_checkpoint_steps`: save the checkpoint every N steps.
8. `save_checkpoint_secs`: save the checkpoint every N second.
9. `initializer`: choose different initializing functions. The possible values are `uniform_unit_scaling`, `uniformnormal` and `normal_unit_scaling`.
10. `initializer_gain`: set the parameter of the initializer.
11. `learning_rate`: adjust the learning rate.
12. `learning_rate_decay`: set different learning rate decay function. The possible values are `noam`(used in the transformer architecture)`piecewise_constant`(see `tf.train.piecewise_constant`) and `none`.
13. `learning_rate_boundaries`: see `tf.train.piecewise_constant` for more references.
14. `learning_rate_values`: see `tf.trian.piecewise_constant` for more references.
15. `keep_checkpoint_max`: the maximum number of checkpoints to keep.
16. `keep_top_checkpoint_max`: the maximum number of top performed checkpoints to keep.
17. `eval_steps`: validate the model every N steps.
18. `eval_secs`: validate the model every N seconds.

19. `eval_batch_size`: the batch size used when validating.
20. `decode_alpha`: the length penalty term in the beam search. See (Wu et al., 2016) for more references.
21. `beam_size`: the beam size of beam search.
22. `decode_length`: the maximum length of decoded length, which is the length of source sentence plus this value.
23. `device_list`: the device id of GPUs. For example, set `device_list=[0,1]` will use GPU 0 and 1. Each GPU will process `batch_size` data.

3.6 Parameters specific to Seq2Seq

1. `rnn_cell`: the recurrent unit. Currently, only LSTM is supported.
2. `embedding_size`: the size of source and target embedding.
3. `hidden_size`: the hidden units of RNN layer.
4. `num_hidden_layers`: the number of RNN layers.
5. `dropout`: dropout rate used in the network.
6. `label_smoothing`: the value of label smoothing.
7. `reverse_source`: if set to true, reverse the source sentence automatically (Sutskever et al., 2014).
8. `use_residual`: whether to use residual connections when multi-layered LSTM is employed.

3.7 Parameters specific to RNNsearch

1. `rnn_cell`: the recurrent unit. Currently, only GRU is supported.
2. `embedding_size`: the size of source and target embedding.
3. `hidden_size`: the hidden units of RNN layer.
4. `maxnum`: the hidden units of maxout layer.
5. `dropout`: dropout rate used in the network.
6. `label_smoothing`: the value of label smoothing.

3.8 Parameters specific to Transformer

1. `hidden_size`: the embedding size and hidden size of the network.
2. `filter_size`: the hidden size of Feed-forward layer.
3. `num_encoder_layers`: the number of encoder layers.
4. `num_decoder_layers`: the number of decoder layers.
5. `num_heads`: the number of attention heads used in attention mechanism.
6. `shared_embedding_and_softmax_weights`: whether to share the embedding and softmax weights.
7. `shared_source_target_embedding`: whether to share the source and target embedding.
8. `residual_dropout`: the dropout rate used in residual connection.
9. `attention_dropout`: the dropout rate used in attention mechanism.
10. `relu_dropout`: the dropout rate used in feed forward layer.
11. `label_smoothing`: the value of label smoothing.

3.9 Test

3.9.1 The Python Script for Test

The `translator.py` in the `bin` folder is used to translate unseen test sentences. It has the following arguments:

```
1 usage: translator.py [<args>] [-h | --help]
2
3 Translate using existing NMT models
4
5 optional arguments:
6   -h, --help            show this help message and exit
7   --input INPUT          Path of input file
8   --output OUTPUT        Path of output file
9   --checkpoints CHECKPOINTS [CHECKPOINTS ...]
10                          Path of trained models
11   --vocabulary VOCABULARY VOCABULARY
12                          Path of source and target vocabulary
13   --models MODELS [MODELS ...]
14                          Name of the model
15   --parameters PARAMETERS
16                          Additional hyper parameters
```

Users must specify the following required arguments to run the `translator.py` script:

1. `--input`: the path to input file.
2. `--output`: the path to translated results.
3. `--checkpointsthe` path to trained checkpoints.
4. `--vocabulary`: the path to source and target vocabulary.
5. `--models`: the name of used architecture.

The optional argument of `translator.py` is listed below:

1. `--parameters`: Change parameters used in decoding. Such as `beam_size` or the GPU id.
2. `--help`: Show the help message.

3.9.2 Decoding

The source file of the test set `test.src` in the data folder is

```
1 ta xihuan huahua me ?
2 ta yidian dou bu xihuan huahua .
```

Given a directory to trained models `train`, please run the following command to translate the test set without evaluation:

```
1 python /PATH/TO/THUMT/thumt/bin/translator.py
2 --models transformer
3 --input test.src
4 --output test.trans
5 --checkpoints train
6 --vocabulary vocab.zh vocab.en
7 --parameters=device_list=[0]
```

Note that `test.trans` is the output of THUMT.

3.10 Model Averaging

Model averaging is also called checkpoint ensemble, which is used to average the checkpoints saved during training. The `checkpoint_averaging.py` script is used to average checkpoints, which has the following arguments:

```
1 usage: average.py [<args>] [-h | --help]
2
3 Average checkpoints
4
5 optional arguments:
6 -h, --help            show this help message and exit
7 --path PATH           checkpoint dir
8 --checkpoints CHECKPOINTS
```

9		number of checkpoints to use
10	--output OUTPUT	output path

The meaning of arguments is shown below:

1. **path**: the path to checkpoints, which is the **output** argument specified during training.
2. **--checkpoints**: the number of checkpoints to average. It will load N latest checkpoints.
3. **--output**: the output directory.
4. **--help** show the help message.

3.11 Examples

The following is the command we used to train a base Transformer model. This model shares the source, target and softmax weights. We employed 4 GPUs to accelerate training:

```

1 python /PATH/TO/THUMT/thumt/bin/trainer.py
2   --input train.tok.clean.bpe.32000.en.shuf
3       train.tok.clean.bpe.32000.de.shuf
4   --model transformer --output train/
5   --vocabulary vocab.bpe.32000 vocab.bpe.32000
6   --validation newstest2013.tok.bpe.32000.en
7   --references newstest2013.tok.bpe.32000.de
8   --parameters=batch_size=6250,device_list=[0,1,2,3],
9               eval_steps=5000,train_steps=100000,
10              save_checkpoint_steps=1500,
11              shared_embedding_and_softmax_weights=true,
12              shared_source_target_embedding=true

```

Once the training stage is finished, we use the following command to average checkpoints:

```

1 python /PATH/TO/THUMT/thumt/scripts/checkpoint_averaging.py
2   --path train --checkpoints 5 --output average

```

Finally, we use the `translator.py` to translate the test set.

```

1 python /PATH/TO/THUMT/thumt/bin/translator.py
2   --model transformer
3   --checkpoints average
4   --input newstest2014.tok.bpe.32000.en
5   --vocabulary vocab.bpe.32000 vocab.bpe.32000
6   --output output.txt

```

References

- Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In *Proceedings of ICLR*.
- Brown, P. F., Della Pietra, S. A., Della Pietra, V. J., and Mercer, R. L. (1993). The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*.
- Chiang, D. (2005). A hierarchical phrase-based model for statistical machine translation. In *Proceedings of ACL*.
- Koehn, P., Och, F. J., and Marcu, D. (2003). Statistical phrase-based translation. In *Proceedings of NAACL*.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. (2002). Bleu: A method for automatic evaluation of machine translation. In *Proceedings of ACL*.
- Shen, S., Cheng, Y., He, Z., He, W., Wu, H., Sun, M., and Liu, Y. (2016). Minimum risk training for neural machine translation. In *Proceedings of ACL*.
- Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Proceedings of NIPS*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. *arXiv preprint arXiv:1706.03762*.
- Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Kaiser, L., Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G., Hughes, M., and Dean, J. (2016). Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv:1609.08144v2*.