# Notebook

February 21, 2019

```python
In [1]: # import modules & set up logging
        import gensim, logging
        import smart_open, os
        logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.
        import datetime
        import pandas as pd
        import multiprocessing

        # fichier incltu dans le projet
        import save_notebook
```

```
D:\Outil\Anaconda\envs\majeure-ml-env\lib\site-packages\gensim\utils.py:1197: UserWarning: det
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

# 1 Déclaration données

```python
In [2]: now = str(datetime.datetime.now()).replace(" ","")
```

```python
In [3]: word_embedding = "word2vec"
```

# 2 Prepare data

```python
In [4]: filenames =  os.listdir("../wikipedia/data")
```

```python
In [5]: #Créé un fichier ou chaque ligne continent tout un fichier
        # path="../wikipedia/data"
        path="../wikipedia/data/"
        with open('./data/wikipedia_informatic.txt', 'w+',encoding="utf8" ) as out_file:

            for fname in filenames:
        #         print(fname)
                if "ipynb_checkpoints" in fname:
                    continue
                try:
                    with open(path + fname, encoding="utf8") as in_file:
                        out_file.write(in_file.read().replace("\n",""))
```

```
            except:
                continue
```

In [6]: 
```python
# On lit et on tokenize le fichier
with open('./data/wikipedia_informatic.txt', 'r', encoding="utf8") as f:
    wiki_vocab = f.readlines()
wiki_vocab = [x.strip() for x in wiki_vocab]

wiki_vocab_tokenized = []
# for line in wiki_vocab:
#     print(gensim.utils.simple_preprocess(line))
# wiki_vocab_tokenized.append(gensim.utils.simple_preprocess(str(wiki_vocab)))
```

In [7]: 
```python
wiki_vocab_tokenized = gensim.utils.simple_preprocess(str(wiki_vocab))
```

## 3 Create model

In [ ]: 
```python
# build vocabulary and train model
model = gensim.models.Word2Vec(
    [wiki_vocab_tokenized],
    size=150,
    seed=1234,
    window=10,
    min_count=2,
    workers=multiprocessing.cpu_count())

date_before_learning = datetime.datetime.now()
model.train([wiki_vocab_tokenized], total_examples=len(wiki_vocab_tokenized), epochs=20
time_training = datetime.datetime.now() - date_before_learning
```

```
2019-02-21 17:16:40,363 : WARNING : consider setting layer size to a multiple of 4 for greater
2019-02-21 17:16:40,364 : INFO : collecting all words and their counts
2019-02-21 17:16:40,365 : INFO : PROGRESS: at sentence #0, processed 0 words, keeping 0 word ty
2019-02-21 17:16:41,604 : INFO : collected 134060 word types from a corpus of 5386246 raw words
2019-02-21 17:16:41,606 : INFO : Loading a fresh vocabulary
2019-02-21 17:16:41,995 : INFO : min_count=2 retains 62658 unique words (46% of original 134060
2019-02-21 17:16:41,997 : INFO : min_count=2 leaves 5314844 word corpus (98% of original 5386246
2019-02-21 17:16:42,233 : INFO : deleting the raw counts dictionary of 134060 items
2019-02-21 17:16:42,237 : INFO : sample=0.001 downsamples 28 most-common words
2019-02-21 17:16:42,238 : INFO : downsampling leaves estimated 4091074 word corpus (77.0% of pr
2019-02-21 17:16:42,477 : INFO : estimated required memory for 62658 words and 150 dimensions:
2019-02-21 17:16:42,479 : INFO : resetting layer weights
2019-02-21 17:16:43,719 : INFO : training model with 4 workers on 62658 vocabulary and 150 feat
2019-02-21 17:16:43,728 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:43,729 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:43,730 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:43,775 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:43,776 : INFO : EPOCH - 1 : training on 5386246 raw words (10000 effective wor
```

```
2019-02-21 17:16:43,786 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:43,787 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:43,788 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:43,821 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:43,823 : INFO : EPOCH - 2 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:16:43,831 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:43,833 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:43,833 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:43,864 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:43,866 : INFO : EPOCH - 3 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:16:43,875 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:43,876 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:43,877 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:43,906 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:43,907 : INFO : EPOCH - 4 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:16:43,915 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:43,917 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:43,918 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:43,946 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:43,947 : INFO : EPOCH - 5 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:16:43,948 : INFO : training on a 26931230 raw words (50000 effective words) took
2019-02-21 17:16:43,949 : WARNING : under 10 jobs per worker: consider setting a smaller `batch
2019-02-21 17:16:43,977 : WARNING : Effective 'alpha' higher than previous training cycles
2019-02-21 17:16:43,978 : INFO : training model with 4 workers on 62658 vocabulary and 150 feat
2019-02-21 17:16:43,989 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:43,990 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:43,991 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,020 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,021 : INFO : EPOCH - 1 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:16:44,022 : WARNING : EPOCH - 1 : supplied example count (1) did not equal expect
2019-02-21 17:16:44,034 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,035 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,036 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,066 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,067 : INFO : EPOCH - 2 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:16:44,067 : WARNING : EPOCH - 2 : supplied example count (1) did not equal expect
2019-02-21 17:16:44,077 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,079 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,080 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,108 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,109 : INFO : EPOCH - 3 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:16:44,109 : WARNING : EPOCH - 3 : supplied example count (1) did not equal expect
2019-02-21 17:16:44,121 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,122 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,123 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,153 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,153 : INFO : EPOCH - 4 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:16:44,154 : WARNING : EPOCH - 4 : supplied example count (1) did not equal expect
```

```
2019-02-21 17:16:44,164 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,164 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,165 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,196 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,200 : INFO : EPOCH - 5 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,201 : WARNING : EPOCH - 5 : supplied example count (1) did not equal expect
2019-02-21 17:16:44,216 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,217 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,218 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,247 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,248 : INFO : EPOCH - 6 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,249 : WARNING : EPOCH - 6 : supplied example count (1) did not equal expect
2019-02-21 17:16:44,262 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,264 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,264 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,293 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,294 : INFO : EPOCH - 7 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,295 : WARNING : EPOCH - 7 : supplied example count (1) did not equal expect
2019-02-21 17:16:44,306 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,307 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,307 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,336 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,337 : INFO : EPOCH - 8 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,338 : WARNING : EPOCH - 8 : supplied example count (1) did not equal expect
2019-02-21 17:16:44,346 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,347 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,348 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,377 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,377 : INFO : EPOCH - 9 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,378 : WARNING : EPOCH - 9 : supplied example count (1) did not equal expect
2019-02-21 17:16:44,391 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,393 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,394 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,423 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,424 : INFO : EPOCH - 10 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:44,425 : WARNING : EPOCH - 10 : supplied example count (1) did not equal expe
2019-02-21 17:16:44,439 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,441 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,442 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,470 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,472 : INFO : EPOCH - 11 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:44,473 : WARNING : EPOCH - 11 : supplied example count (1) did not equal expe
2019-02-21 17:16:44,484 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,487 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,488 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,516 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,518 : INFO : EPOCH - 12 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:44,521 : WARNING : EPOCH - 12 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:44,531 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,532 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,534 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,562 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,563 : INFO : EPOCH - 13 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,564 : WARNING : EPOCH - 13 : supplied example count (1) did not equal expec
2019-02-21 17:16:44,575 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,576 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,577 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,607 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,608 : INFO : EPOCH - 14 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,608 : WARNING : EPOCH - 14 : supplied example count (1) did not equal expec
2019-02-21 17:16:44,620 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,622 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,623 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,650 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,651 : INFO : EPOCH - 15 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,652 : WARNING : EPOCH - 15 : supplied example count (1) did not equal expec
2019-02-21 17:16:44,662 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,663 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,664 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,692 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,693 : INFO : EPOCH - 16 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,693 : WARNING : EPOCH - 16 : supplied example count (1) did not equal expec
2019-02-21 17:16:44,704 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,706 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,707 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,735 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,736 : INFO : EPOCH - 17 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,736 : WARNING : EPOCH - 17 : supplied example count (1) did not equal expec
2019-02-21 17:16:44,747 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,748 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,749 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,778 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,779 : INFO : EPOCH - 18 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,780 : WARNING : EPOCH - 18 : supplied example count (1) did not equal expec
2019-02-21 17:16:44,792 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,793 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,795 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,823 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,824 : INFO : EPOCH - 19 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,824 : WARNING : EPOCH - 19 : supplied example count (1) did not equal expec
2019-02-21 17:16:44,836 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,837 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,838 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,865 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,866 : INFO : EPOCH - 20 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:44,867 : WARNING : EPOCH - 20 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:44,877 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,879 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,881 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,907 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,908 : INFO : EPOCH - 21 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:44,908 : WARNING : EPOCH - 21 : supplied example count (1) did not equal expe
2019-02-21 17:16:44,918 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,919 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,921 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,949 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,950 : INFO : EPOCH - 22 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:44,950 : WARNING : EPOCH - 22 : supplied example count (1) did not equal expe
2019-02-21 17:16:44,961 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:44,964 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:44,965 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:44,991 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:44,991 : INFO : EPOCH - 23 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:44,992 : WARNING : EPOCH - 23 : supplied example count (1) did not equal expe
2019-02-21 17:16:45,001 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,002 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,003 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,031 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,032 : INFO : EPOCH - 24 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:45,033 : WARNING : EPOCH - 24 : supplied example count (1) did not equal expe
2019-02-21 17:16:45,044 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,045 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,047 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,074 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,075 : INFO : EPOCH - 25 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:45,076 : WARNING : EPOCH - 25 : supplied example count (1) did not equal expe
2019-02-21 17:16:45,090 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,092 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,093 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,125 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,126 : INFO : EPOCH - 26 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:45,127 : WARNING : EPOCH - 26 : supplied example count (1) did not equal expe
2019-02-21 17:16:45,141 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,142 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,142 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,171 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,173 : INFO : EPOCH - 27 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:45,173 : WARNING : EPOCH - 27 : supplied example count (1) did not equal expe
2019-02-21 17:16:45,193 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,206 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,206 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,222 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,222 : INFO : EPOCH - 28 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:45,223 : WARNING : EPOCH - 28 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:45,233 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,235 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,237 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,263 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,264 : INFO : EPOCH - 29 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,265 : WARNING : EPOCH - 29 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,276 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,277 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,277 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,305 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,307 : INFO : EPOCH - 30 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,308 : WARNING : EPOCH - 30 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,318 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,320 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,321 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,346 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,347 : INFO : EPOCH - 31 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,348 : WARNING : EPOCH - 31 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,358 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,360 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,362 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,389 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,390 : INFO : EPOCH - 32 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,391 : WARNING : EPOCH - 32 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,401 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,403 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,404 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,433 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,434 : INFO : EPOCH - 33 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,434 : WARNING : EPOCH - 33 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,443 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,444 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,445 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,473 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,474 : INFO : EPOCH - 34 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,474 : WARNING : EPOCH - 34 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,484 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,489 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,490 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,514 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,515 : INFO : EPOCH - 35 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,516 : WARNING : EPOCH - 35 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,527 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,528 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,528 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,556 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,557 : INFO : EPOCH - 36 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,557 : WARNING : EPOCH - 36 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:45,567 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,572 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,573 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,598 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,599 : INFO : EPOCH - 37 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,599 : WARNING : EPOCH - 37 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,610 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,611 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,612 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,639 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,640 : INFO : EPOCH - 38 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,640 : WARNING : EPOCH - 38 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,652 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,653 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,654 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,700 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,705 : INFO : EPOCH - 39 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,706 : WARNING : EPOCH - 39 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,717 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,718 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,718 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,745 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,746 : INFO : EPOCH - 40 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,747 : WARNING : EPOCH - 40 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,759 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,760 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,761 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,788 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,789 : INFO : EPOCH - 41 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,789 : WARNING : EPOCH - 41 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,799 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,800 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,800 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,826 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,827 : INFO : EPOCH - 42 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,828 : WARNING : EPOCH - 42 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,848 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,849 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,849 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,877 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,878 : INFO : EPOCH - 43 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,878 : WARNING : EPOCH - 43 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,889 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,890 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,891 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,918 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,919 : INFO : EPOCH - 44 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,920 : WARNING : EPOCH - 44 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:45,936 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,937 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,938 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:45,968 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:45,969 : INFO : EPOCH - 45 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:45,970 : WARNING : EPOCH - 45 : supplied example count (1) did not equal expec
2019-02-21 17:16:45,979 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:45,980 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:45,981 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,008 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,009 : INFO : EPOCH - 46 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,009 : WARNING : EPOCH - 46 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,019 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,020 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,021 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,048 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,049 : INFO : EPOCH - 47 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,050 : WARNING : EPOCH - 47 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,060 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,061 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,061 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,089 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,089 : INFO : EPOCH - 48 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,090 : WARNING : EPOCH - 48 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,101 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,102 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,103 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,129 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,130 : INFO : EPOCH - 49 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,131 : WARNING : EPOCH - 49 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,141 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,142 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,143 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,171 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,173 : INFO : EPOCH - 50 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,174 : WARNING : EPOCH - 50 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,189 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,192 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,193 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,219 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,220 : INFO : EPOCH - 51 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,221 : WARNING : EPOCH - 51 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,229 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,230 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,231 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,259 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,260 : INFO : EPOCH - 52 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,260 : WARNING : EPOCH - 52 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:46,271 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,271 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,272 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,299 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,300 : INFO : EPOCH - 53 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,300 : WARNING : EPOCH - 53 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,310 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,311 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,311 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,339 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,340 : INFO : EPOCH - 54 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,341 : WARNING : EPOCH - 54 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,350 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,351 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,352 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,380 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,381 : INFO : EPOCH - 55 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,381 : WARNING : EPOCH - 55 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,392 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,393 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,394 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,421 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,456 : INFO : EPOCH - 56 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,457 : WARNING : EPOCH - 56 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,468 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,470 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,471 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,500 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,501 : INFO : EPOCH - 57 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,504 : WARNING : EPOCH - 57 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,517 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,518 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,519 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,547 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,547 : INFO : EPOCH - 58 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,548 : WARNING : EPOCH - 58 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,558 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,560 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,562 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,586 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,587 : INFO : EPOCH - 59 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,588 : WARNING : EPOCH - 59 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,601 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,603 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,604 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,629 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,629 : INFO : EPOCH - 60 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,630 : WARNING : EPOCH - 60 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:46,640 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,643 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,644 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,669 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,670 : INFO : EPOCH - 61 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:16:46,671 : WARNING : EPOCH - 61 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,680 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,681 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,682 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,709 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,710 : INFO : EPOCH - 62 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:16:46,711 : WARNING : EPOCH - 62 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,719 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,721 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,721 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,748 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,749 : INFO : EPOCH - 63 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:16:46,750 : WARNING : EPOCH - 63 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,761 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,762 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,763 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,790 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,792 : INFO : EPOCH - 64 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:16:46,792 : WARNING : EPOCH - 64 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,804 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,805 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,805 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,833 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,834 : INFO : EPOCH - 65 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:16:46,834 : WARNING : EPOCH - 65 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,845 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,846 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,847 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,874 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,875 : INFO : EPOCH - 66 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:16:46,876 : WARNING : EPOCH - 66 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,885 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,886 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,887 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,912 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,914 : INFO : EPOCH - 67 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:16:46,915 : WARNING : EPOCH - 67 : supplied example count (1) did not equal expec
2019-02-21 17:16:46,926 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,927 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,928 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,954 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,955 : INFO : EPOCH - 68 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:16:46,955 : WARNING : EPOCH - 68 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:46,965 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:46,967 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:46,970 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:46,995 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:46,996 : INFO : EPOCH - 69 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:46,997 : WARNING : EPOCH - 69 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,006 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,007 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,008 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,034 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,035 : INFO : EPOCH - 70 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,036 : WARNING : EPOCH - 70 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,050 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,051 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,052 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,078 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,087 : INFO : EPOCH - 71 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,088 : WARNING : EPOCH - 71 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,096 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,096 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,097 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,124 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,125 : INFO : EPOCH - 72 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,126 : WARNING : EPOCH - 72 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,137 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,138 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,139 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,165 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,167 : INFO : EPOCH - 73 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,168 : WARNING : EPOCH - 73 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,183 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,186 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,187 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,211 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,212 : INFO : EPOCH - 74 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,212 : WARNING : EPOCH - 74 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,222 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,223 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,224 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,250 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,252 : INFO : EPOCH - 75 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,252 : WARNING : EPOCH - 75 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,259 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,261 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,262 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,290 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,290 : INFO : EPOCH - 76 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,291 : WARNING : EPOCH - 76 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:47,305 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,306 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,306 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,332 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,371 : INFO : EPOCH - 77 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,371 : WARNING : EPOCH - 77 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,381 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,382 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,382 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,409 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,410 : INFO : EPOCH - 78 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,411 : WARNING : EPOCH - 78 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,421 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,422 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,423 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,448 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,449 : INFO : EPOCH - 79 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,450 : WARNING : EPOCH - 79 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,463 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,464 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,464 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,491 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,492 : INFO : EPOCH - 80 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,492 : WARNING : EPOCH - 80 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,502 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,503 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,504 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,530 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,531 : INFO : EPOCH - 81 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,532 : WARNING : EPOCH - 81 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,542 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,543 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,544 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,570 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,571 : INFO : EPOCH - 82 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,571 : WARNING : EPOCH - 82 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,581 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,582 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,584 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,610 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,611 : INFO : EPOCH - 83 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,611 : WARNING : EPOCH - 83 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,618 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,619 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,620 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,648 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,649 : INFO : EPOCH - 84 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,650 : WARNING : EPOCH - 84 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:47,669 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,670 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,671 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,696 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,698 : INFO : EPOCH - 85 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,698 : WARNING : EPOCH - 85 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,709 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,710 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,711 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,736 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,737 : INFO : EPOCH - 86 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,738 : WARNING : EPOCH - 86 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,747 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,748 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,749 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,776 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,777 : INFO : EPOCH - 87 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,778 : WARNING : EPOCH - 87 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,789 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,790 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,791 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,816 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,817 : INFO : EPOCH - 88 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,818 : WARNING : EPOCH - 88 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,827 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,828 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,829 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,856 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,857 : INFO : EPOCH - 89 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,857 : WARNING : EPOCH - 89 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,867 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,868 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,869 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,895 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,896 : INFO : EPOCH - 90 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,897 : WARNING : EPOCH - 90 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,908 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,917 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,918 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,936 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,937 : INFO : EPOCH - 91 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,938 : WARNING : EPOCH - 91 : supplied example count (1) did not equal expec
2019-02-21 17:16:47,947 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,948 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,948 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:47,975 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:47,976 : INFO : EPOCH - 92 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:47,977 : WARNING : EPOCH - 92 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:16:47,986 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:47,988 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:47,988 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,015 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,015 : INFO : EPOCH - 93 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:48,016 : WARNING : EPOCH - 93 : supplied example count (1) did not equal expec
2019-02-21 17:16:48,025 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,026 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,027 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,053 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,054 : INFO : EPOCH - 94 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:48,055 : WARNING : EPOCH - 94 : supplied example count (1) did not equal expec
2019-02-21 17:16:48,064 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,065 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,065 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,091 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,092 : INFO : EPOCH - 95 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:48,092 : WARNING : EPOCH - 95 : supplied example count (1) did not equal expec
2019-02-21 17:16:48,100 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,101 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,103 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,129 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,130 : INFO : EPOCH - 96 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:48,131 : WARNING : EPOCH - 96 : supplied example count (1) did not equal expec
2019-02-21 17:16:48,141 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,142 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,143 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,169 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,170 : INFO : EPOCH - 97 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:48,171 : WARNING : EPOCH - 97 : supplied example count (1) did not equal expec
2019-02-21 17:16:48,182 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,184 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,184 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,214 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,215 : INFO : EPOCH - 98 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:48,216 : WARNING : EPOCH - 98 : supplied example count (1) did not equal expec
2019-02-21 17:16:48,226 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,227 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,227 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,253 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,254 : INFO : EPOCH - 99 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:16:48,255 : WARNING : EPOCH - 99 : supplied example count (1) did not equal expec
2019-02-21 17:16:48,268 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,270 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,271 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,296 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,297 : INFO : EPOCH - 100 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:48,298 : WARNING : EPOCH - 100 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:48,307 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,308 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,309 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,337 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,338 : INFO : EPOCH - 101 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:48,338 : WARNING : EPOCH - 101 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,350 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,351 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,352 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,377 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,378 : INFO : EPOCH - 102 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:48,379 : WARNING : EPOCH - 102 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,389 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,390 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,392 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,422 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,423 : INFO : EPOCH - 103 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:48,423 : WARNING : EPOCH - 103 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,435 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,437 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,437 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,465 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,471 : INFO : EPOCH - 104 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:48,472 : WARNING : EPOCH - 104 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,482 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,486 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,486 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,511 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,512 : INFO : EPOCH - 105 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:48,513 : WARNING : EPOCH - 105 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,523 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,524 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,525 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,551 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,552 : INFO : EPOCH - 106 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:48,553 : WARNING : EPOCH - 106 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,564 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,565 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,566 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,594 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,594 : INFO : EPOCH - 107 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:48,595 : WARNING : EPOCH - 107 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,605 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,605 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,606 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,631 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,632 : INFO : EPOCH - 108 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:48,633 : WARNING : EPOCH - 108 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:48,645 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,647 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,648 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,673 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,674 : INFO : EPOCH - 109 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,675 : WARNING : EPOCH - 109 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,684 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,686 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,686 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,713 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,714 : INFO : EPOCH - 110 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,715 : WARNING : EPOCH - 110 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,726 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,727 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,728 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,757 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,758 : INFO : EPOCH - 111 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,758 : WARNING : EPOCH - 111 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,768 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,769 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,770 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,796 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,797 : INFO : EPOCH - 112 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,798 : WARNING : EPOCH - 112 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,810 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,811 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,812 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,837 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,838 : INFO : EPOCH - 113 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,838 : WARNING : EPOCH - 113 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,848 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,850 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,850 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,877 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,878 : INFO : EPOCH - 114 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,879 : WARNING : EPOCH - 114 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,888 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,889 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,889 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,915 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,916 : INFO : EPOCH - 115 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,917 : WARNING : EPOCH - 115 : supplied example count (1) did not equal expe
2019-02-21 17:16:48,927 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,928 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,928 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,954 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,955 : INFO : EPOCH - 116 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,957 : WARNING : EPOCH - 116 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:48,971 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:48,972 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:48,973 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:48,998 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:48,999 : INFO : EPOCH - 117 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:48,999 : WARNING : EPOCH - 117 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,010 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,011 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,012 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,037 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,038 : INFO : EPOCH - 118 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,039 : WARNING : EPOCH - 118 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,052 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,054 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,055 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,079 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,080 : INFO : EPOCH - 119 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,081 : WARNING : EPOCH - 119 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,091 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,091 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,092 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,117 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,118 : INFO : EPOCH - 120 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,119 : WARNING : EPOCH - 120 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,127 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,129 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,129 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,157 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,158 : INFO : EPOCH - 121 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,158 : WARNING : EPOCH - 121 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,168 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,169 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,170 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,196 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,197 : INFO : EPOCH - 122 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,198 : WARNING : EPOCH - 122 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,208 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,209 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,210 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,237 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,238 : INFO : EPOCH - 123 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,239 : WARNING : EPOCH - 123 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,252 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,254 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,254 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,280 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,290 : INFO : EPOCH - 124 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,291 : WARNING : EPOCH - 124 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:49,301 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,303 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,303 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,330 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,331 : INFO : EPOCH - 125 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,332 : WARNING : EPOCH - 125 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,342 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,343 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,344 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,370 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,371 : INFO : EPOCH - 126 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,372 : WARNING : EPOCH - 126 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,381 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,382 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,383 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,420 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,433 : INFO : EPOCH - 127 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,433 : WARNING : EPOCH - 127 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,444 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,446 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,446 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,471 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,472 : INFO : EPOCH - 128 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,473 : WARNING : EPOCH - 128 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,482 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,482 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,483 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,509 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,510 : INFO : EPOCH - 129 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,510 : WARNING : EPOCH - 129 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,519 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,520 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,521 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,547 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,549 : INFO : EPOCH - 130 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,550 : WARNING : EPOCH - 130 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,571 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,574 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,576 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,599 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,600 : INFO : EPOCH - 131 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,601 : WARNING : EPOCH - 131 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,613 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,614 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,615 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,642 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,643 : INFO : EPOCH - 132 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:49,644 : WARNING : EPOCH - 132 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:49,655 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,656 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,656 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,682 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,683 : INFO : EPOCH - 133 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:49,684 : WARNING : EPOCH - 133 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,692 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,694 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,694 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,721 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,721 : INFO : EPOCH - 134 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:49,722 : WARNING : EPOCH - 134 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,731 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,732 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,733 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,759 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,760 : INFO : EPOCH - 135 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:49,761 : WARNING : EPOCH - 135 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,772 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,773 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,773 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,798 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,799 : INFO : EPOCH - 136 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:49,800 : WARNING : EPOCH - 136 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,811 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,812 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,813 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,840 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,850 : INFO : EPOCH - 137 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:49,851 : WARNING : EPOCH - 137 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,861 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,862 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,863 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,889 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,889 : INFO : EPOCH - 138 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:49,890 : WARNING : EPOCH - 138 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,899 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,900 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,901 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,926 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,927 : INFO : EPOCH - 139 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:49,928 : WARNING : EPOCH - 139 : supplied example count (1) did not equal expe
2019-02-21 17:16:49,936 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,938 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,938 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:49,964 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:49,965 : INFO : EPOCH - 140 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:49,966 : WARNING : EPOCH - 140 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:49,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:49,977 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:49,978 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,003 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,004 : INFO : EPOCH - 141 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:50,005 : WARNING : EPOCH - 141 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,014 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,015 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,015 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,040 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,041 : INFO : EPOCH - 142 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:50,042 : WARNING : EPOCH - 142 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,052 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,054 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,055 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,082 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,085 : INFO : EPOCH - 143 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:50,086 : WARNING : EPOCH - 143 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,095 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,096 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,096 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,123 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,124 : INFO : EPOCH - 144 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:50,125 : WARNING : EPOCH - 144 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,133 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,135 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,136 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,162 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,163 : INFO : EPOCH - 145 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:50,164 : WARNING : EPOCH - 145 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,179 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,181 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,182 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,208 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,208 : INFO : EPOCH - 146 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:50,209 : WARNING : EPOCH - 146 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,218 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,219 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,220 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,244 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,245 : INFO : EPOCH - 147 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:50,246 : WARNING : EPOCH - 147 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,256 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,257 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,257 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,284 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,285 : INFO : EPOCH - 148 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:50,286 : WARNING : EPOCH - 148 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:50,295 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,296 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,297 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,322 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,323 : INFO : EPOCH - 149 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,323 : WARNING : EPOCH - 149 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,333 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,333 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,334 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,359 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,360 : INFO : EPOCH - 150 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,360 : WARNING : EPOCH - 150 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,373 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,374 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,374 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,403 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,410 : INFO : EPOCH - 151 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,411 : WARNING : EPOCH - 151 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,421 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,422 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,423 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,448 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,449 : INFO : EPOCH - 152 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,450 : WARNING : EPOCH - 152 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,460 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,461 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,462 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,488 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,489 : INFO : EPOCH - 153 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,490 : WARNING : EPOCH - 153 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,501 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,503 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,503 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,529 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,530 : INFO : EPOCH - 154 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,531 : WARNING : EPOCH - 154 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,540 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,545 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,547 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,569 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,570 : INFO : EPOCH - 155 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,570 : WARNING : EPOCH - 155 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,580 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,581 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,581 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,607 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,608 : INFO : EPOCH - 156 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,608 : WARNING : EPOCH - 156 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:50,619 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,621 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,621 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,646 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,647 : INFO : EPOCH - 157 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,648 : WARNING : EPOCH - 157 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,658 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,659 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,659 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,686 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,687 : INFO : EPOCH - 158 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,688 : WARNING : EPOCH - 158 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,696 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,698 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,698 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,724 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,725 : INFO : EPOCH - 159 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,726 : WARNING : EPOCH - 159 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,734 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,736 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,736 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,762 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,764 : INFO : EPOCH - 160 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,765 : WARNING : EPOCH - 160 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,774 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,776 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,777 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,803 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,804 : INFO : EPOCH - 161 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,805 : WARNING : EPOCH - 161 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,814 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,815 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,815 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,841 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,842 : INFO : EPOCH - 162 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,842 : WARNING : EPOCH - 162 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,851 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,852 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,853 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,877 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,878 : INFO : EPOCH - 163 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,879 : WARNING : EPOCH - 163 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,891 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,892 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,893 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,919 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,920 : INFO : EPOCH - 164 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:50,921 : WARNING : EPOCH - 164 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:50,931 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,940 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,941 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,959 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,960 : INFO : EPOCH - 165 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:50,961 : WARNING : EPOCH - 165 : supplied example count (1) did not equal expe
2019-02-21 17:16:50,970 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:50,971 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:50,971 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:50,996 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:50,996 : INFO : EPOCH - 166 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:50,997 : WARNING : EPOCH - 166 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,007 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,009 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,011 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,036 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,037 : INFO : EPOCH - 167 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:51,039 : WARNING : EPOCH - 167 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,048 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,050 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,050 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,076 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,077 : INFO : EPOCH - 168 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:51,077 : WARNING : EPOCH - 168 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,086 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,088 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,088 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,113 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,114 : INFO : EPOCH - 169 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:51,114 : WARNING : EPOCH - 169 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,123 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,125 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,125 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,152 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,153 : INFO : EPOCH - 170 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:51,154 : WARNING : EPOCH - 170 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,202 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,204 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,205 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,231 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,236 : INFO : EPOCH - 171 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:51,237 : WARNING : EPOCH - 171 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,246 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,247 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,248 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,274 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,275 : INFO : EPOCH - 172 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:51,276 : WARNING : EPOCH - 172 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:51,286 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,287 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,287 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,312 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,313 : INFO : EPOCH - 173 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:51,313 : WARNING : EPOCH - 173 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,323 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,323 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,324 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,350 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,351 : INFO : EPOCH - 174 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:51,352 : WARNING : EPOCH - 174 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,364 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,365 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,366 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,399 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,400 : INFO : EPOCH - 175 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:51,400 : WARNING : EPOCH - 175 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,410 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,411 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,412 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,436 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,437 : INFO : EPOCH - 176 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:51,438 : WARNING : EPOCH - 176 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,446 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,447 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,448 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,473 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,474 : INFO : EPOCH - 177 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:51,475 : WARNING : EPOCH - 177 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,484 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,498 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,498 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,513 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,514 : INFO : EPOCH - 178 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:51,515 : WARNING : EPOCH - 178 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,524 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,526 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,526 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,552 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,553 : INFO : EPOCH - 179 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:51,554 : WARNING : EPOCH - 179 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,562 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,564 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,564 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,589 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,590 : INFO : EPOCH - 180 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:51,591 : WARNING : EPOCH - 180 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:51,600 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,602 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,602 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,627 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,629 : INFO : EPOCH - 181 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:51,629 : WARNING : EPOCH - 181 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,638 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,639 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,639 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,670 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,671 : INFO : EPOCH - 182 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:51,672 : WARNING : EPOCH - 182 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,684 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,685 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,686 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,713 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,714 : INFO : EPOCH - 183 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:51,715 : WARNING : EPOCH - 183 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,726 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,727 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,728 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,753 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,760 : INFO : EPOCH - 184 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:51,761 : WARNING : EPOCH - 184 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,771 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,772 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,772 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,797 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,798 : INFO : EPOCH - 185 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:51,798 : WARNING : EPOCH - 185 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,810 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,811 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,812 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,838 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,840 : INFO : EPOCH - 186 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:51,840 : WARNING : EPOCH - 186 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,850 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,850 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,852 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,877 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,878 : INFO : EPOCH - 187 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:51,879 : WARNING : EPOCH - 187 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,887 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,888 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,890 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,928 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,929 : INFO : EPOCH - 188 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:51,930 : WARNING : EPOCH - 188 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:51,945 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,946 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,946 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:51,971 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:51,972 : INFO : EPOCH - 189 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:51,973 : WARNING : EPOCH - 189 : supplied example count (1) did not equal expe
2019-02-21 17:16:51,982 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:51,986 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:51,987 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,013 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,014 : INFO : EPOCH - 190 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:52,015 : WARNING : EPOCH - 190 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,024 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,025 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,026 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,052 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,055 : INFO : EPOCH - 191 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:52,056 : WARNING : EPOCH - 191 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,065 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,066 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,067 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,092 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,093 : INFO : EPOCH - 192 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:52,094 : WARNING : EPOCH - 192 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,103 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,104 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,105 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,129 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,130 : INFO : EPOCH - 193 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:52,131 : WARNING : EPOCH - 193 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,141 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,142 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,145 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,171 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,172 : INFO : EPOCH - 194 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:52,173 : WARNING : EPOCH - 194 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,188 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,189 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,189 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,215 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,216 : INFO : EPOCH - 195 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:52,217 : WARNING : EPOCH - 195 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,227 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,228 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,228 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,254 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,255 : INFO : EPOCH - 196 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:52,255 : WARNING : EPOCH - 196 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:52,264 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,265 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,266 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,292 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,294 : INFO : EPOCH - 197 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,294 : WARNING : EPOCH - 197 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,306 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,307 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,307 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,335 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,336 : INFO : EPOCH - 198 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,337 : WARNING : EPOCH - 198 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,346 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,347 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,348 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,374 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,375 : INFO : EPOCH - 199 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,375 : WARNING : EPOCH - 199 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,385 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,386 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,387 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,413 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,413 : INFO : EPOCH - 200 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,414 : WARNING : EPOCH - 200 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,424 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,425 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,426 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,451 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,452 : INFO : EPOCH - 201 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,453 : WARNING : EPOCH - 201 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,462 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,463 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,463 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,488 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,489 : INFO : EPOCH - 202 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,490 : WARNING : EPOCH - 202 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,499 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,500 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,501 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,527 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,528 : INFO : EPOCH - 203 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,528 : WARNING : EPOCH - 203 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,539 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,540 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,541 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,567 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,568 : INFO : EPOCH - 204 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,569 : WARNING : EPOCH - 204 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:52,578 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,579 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,580 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,606 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,607 : INFO : EPOCH - 205 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,608 : WARNING : EPOCH - 205 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,617 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,618 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,619 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,644 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,644 : INFO : EPOCH - 206 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,645 : WARNING : EPOCH - 206 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,655 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,656 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,657 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,683 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,684 : INFO : EPOCH - 207 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,685 : WARNING : EPOCH - 207 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,723 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,727 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,728 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,728 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,729 : INFO : EPOCH - 208 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,729 : WARNING : EPOCH - 208 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,740 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,741 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,742 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,767 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,768 : INFO : EPOCH - 209 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,769 : WARNING : EPOCH - 209 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,778 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,779 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,780 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,806 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,807 : INFO : EPOCH - 210 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,808 : WARNING : EPOCH - 210 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,817 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,818 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,819 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,850 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,852 : INFO : EPOCH - 211 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,852 : WARNING : EPOCH - 211 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,864 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,878 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,878 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,891 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,892 : INFO : EPOCH - 212 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,893 : WARNING : EPOCH - 212 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:52,903 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,904 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,904 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,932 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,932 : INFO : EPOCH - 213 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,933 : WARNING : EPOCH - 213 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,944 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,945 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,945 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:52,970 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:52,971 : INFO : EPOCH - 214 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:52,971 : WARNING : EPOCH - 214 : supplied example count (1) did not equal expe
2019-02-21 17:16:52,979 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:52,981 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:52,981 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,008 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,009 : INFO : EPOCH - 215 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:53,009 : WARNING : EPOCH - 215 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,018 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,020 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,020 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,047 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,049 : INFO : EPOCH - 216 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:53,049 : WARNING : EPOCH - 216 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,060 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,061 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,062 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,087 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,088 : INFO : EPOCH - 217 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:53,089 : WARNING : EPOCH - 217 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,098 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,099 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,100 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,127 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,128 : INFO : EPOCH - 218 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:53,128 : WARNING : EPOCH - 218 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,138 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,139 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,140 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,167 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,168 : INFO : EPOCH - 219 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:53,169 : WARNING : EPOCH - 219 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,180 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,182 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,182 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,209 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,210 : INFO : EPOCH - 220 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:53,210 : WARNING : EPOCH - 220 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:53,220 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,221 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,221 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,246 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,247 : INFO : EPOCH - 221 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,248 : WARNING : EPOCH - 221 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,259 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,260 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,262 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,287 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,288 : INFO : EPOCH - 222 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,288 : WARNING : EPOCH - 222 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,298 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,299 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,300 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,326 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,327 : INFO : EPOCH - 223 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,328 : WARNING : EPOCH - 223 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,337 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,338 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,339 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,364 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,365 : INFO : EPOCH - 224 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,365 : WARNING : EPOCH - 224 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,375 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,376 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,376 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,403 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,403 : INFO : EPOCH - 225 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,404 : WARNING : EPOCH - 225 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,418 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,420 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,420 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,447 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,451 : INFO : EPOCH - 226 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,452 : WARNING : EPOCH - 226 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,461 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,462 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,462 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,488 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,489 : INFO : EPOCH - 227 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,489 : WARNING : EPOCH - 227 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,499 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,500 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,500 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,526 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,527 : INFO : EPOCH - 228 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,527 : WARNING : EPOCH - 228 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:53,538 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,539 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,540 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,567 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,568 : INFO : EPOCH - 229 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,569 : WARNING : EPOCH - 229 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,578 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,579 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,579 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,605 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,606 : INFO : EPOCH - 230 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,606 : WARNING : EPOCH - 230 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,615 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,616 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,617 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,641 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,642 : INFO : EPOCH - 231 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,643 : WARNING : EPOCH - 231 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,656 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,657 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,658 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,685 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,686 : INFO : EPOCH - 232 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,687 : WARNING : EPOCH - 232 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,696 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,697 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,698 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,725 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,726 : INFO : EPOCH - 233 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,727 : WARNING : EPOCH - 233 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,736 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,738 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,739 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,764 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,765 : INFO : EPOCH - 234 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,766 : WARNING : EPOCH - 234 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,776 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,777 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,778 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,803 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,804 : INFO : EPOCH - 235 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,804 : WARNING : EPOCH - 235 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,813 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,814 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,815 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,839 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,840 : INFO : EPOCH - 236 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,841 : WARNING : EPOCH - 236 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:53,850 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,851 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,852 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,878 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,879 : INFO : EPOCH - 237 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,879 : WARNING : EPOCH - 237 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,890 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,891 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,892 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,918 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,919 : INFO : EPOCH - 238 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,920 : WARNING : EPOCH - 238 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,930 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,931 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,932 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,959 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:53,961 : INFO : EPOCH - 239 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:53,961 : WARNING : EPOCH - 239 : supplied example count (1) did not equal expe
2019-02-21 17:16:53,971 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:53,972 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:53,973 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:53,999 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,002 : INFO : EPOCH - 240 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,003 : WARNING : EPOCH - 240 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,015 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,016 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,017 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,042 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,043 : INFO : EPOCH - 241 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,043 : WARNING : EPOCH - 241 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,055 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,056 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,057 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,082 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,083 : INFO : EPOCH - 242 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,083 : WARNING : EPOCH - 242 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,092 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,093 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,094 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,120 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,121 : INFO : EPOCH - 243 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,122 : WARNING : EPOCH - 243 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,132 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,133 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,133 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,159 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,160 : INFO : EPOCH - 244 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,160 : WARNING : EPOCH - 244 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:54,171 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,172 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,175 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,200 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,201 : INFO : EPOCH - 245 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,201 : WARNING : EPOCH - 245 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,212 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,213 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,214 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,238 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,239 : INFO : EPOCH - 246 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,239 : WARNING : EPOCH - 246 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,251 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,253 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,254 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,279 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,279 : INFO : EPOCH - 247 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,280 : WARNING : EPOCH - 247 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,293 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,295 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,295 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,322 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,323 : INFO : EPOCH - 248 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,324 : WARNING : EPOCH - 248 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,336 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,337 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,338 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,364 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,365 : INFO : EPOCH - 249 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,366 : WARNING : EPOCH - 249 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,376 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,377 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,378 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,403 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,404 : INFO : EPOCH - 250 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,405 : WARNING : EPOCH - 250 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,414 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,415 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,416 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,441 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,441 : INFO : EPOCH - 251 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,442 : WARNING : EPOCH - 251 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,450 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,451 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,452 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,477 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,478 : INFO : EPOCH - 252 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,478 : WARNING : EPOCH - 252 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:54,488 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,489 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,490 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,514 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,515 : INFO : EPOCH - 253 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,515 : WARNING : EPOCH - 253 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,525 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,527 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,528 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,552 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,553 : INFO : EPOCH - 254 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,554 : WARNING : EPOCH - 254 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,575 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,577 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,578 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,600 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,601 : INFO : EPOCH - 255 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,602 : WARNING : EPOCH - 255 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,614 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,615 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,616 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,642 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,643 : INFO : EPOCH - 256 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,644 : WARNING : EPOCH - 256 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,655 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,657 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,658 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,682 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,683 : INFO : EPOCH - 257 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,684 : WARNING : EPOCH - 257 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,693 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,694 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,695 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,720 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,721 : INFO : EPOCH - 258 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,722 : WARNING : EPOCH - 258 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,731 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,732 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,733 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,757 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,758 : INFO : EPOCH - 259 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,759 : WARNING : EPOCH - 259 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,770 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,771 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,771 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,796 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,797 : INFO : EPOCH - 260 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:54,797 : WARNING : EPOCH - 260 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:54,807 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,808 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,809 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,833 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,834 : INFO : EPOCH - 261 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,835 : WARNING : EPOCH - 261 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,851 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,858 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,859 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,876 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,877 : INFO : EPOCH - 262 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,878 : WARNING : EPOCH - 262 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,889 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,890 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,891 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,915 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,916 : INFO : EPOCH - 263 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,917 : WARNING : EPOCH - 263 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,926 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,927 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,927 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,952 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,953 : INFO : EPOCH - 264 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,954 : WARNING : EPOCH - 264 : supplied example count (1) did not equal expe
2019-02-21 17:16:54,964 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:54,964 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:54,965 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:54,990 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:54,991 : INFO : EPOCH - 265 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:54,992 : WARNING : EPOCH - 265 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,000 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,001 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,002 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,026 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,027 : INFO : EPOCH - 266 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:55,028 : WARNING : EPOCH - 266 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,037 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,039 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,039 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,063 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,064 : INFO : EPOCH - 267 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:55,065 : WARNING : EPOCH - 267 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,075 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,076 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,077 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,104 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,105 : INFO : EPOCH - 268 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:55,106 : WARNING : EPOCH - 268 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:55,117 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,118 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,119 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,145 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,146 : INFO : EPOCH - 269 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:55,147 : WARNING : EPOCH - 269 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,157 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,158 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,159 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,185 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,186 : INFO : EPOCH - 270 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:55,187 : WARNING : EPOCH - 270 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,200 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,202 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,203 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,228 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,229 : INFO : EPOCH - 271 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:55,229 : WARNING : EPOCH - 271 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,241 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,242 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,242 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,269 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,270 : INFO : EPOCH - 272 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:55,270 : WARNING : EPOCH - 272 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,279 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,280 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,281 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,306 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,307 : INFO : EPOCH - 273 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:55,307 : WARNING : EPOCH - 273 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,315 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,316 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,317 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,342 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,343 : INFO : EPOCH - 274 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:55,343 : WARNING : EPOCH - 274 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,353 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,354 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,355 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,387 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,388 : INFO : EPOCH - 275 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:55,390 : WARNING : EPOCH - 275 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,400 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,402 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,403 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,429 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,430 : INFO : EPOCH - 276 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:55,431 : WARNING : EPOCH - 276 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:55,442 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,443 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,444 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,470 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,471 : INFO : EPOCH - 277 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,471 : WARNING : EPOCH - 277 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,481 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,482 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,483 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,508 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,509 : INFO : EPOCH - 278 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,509 : WARNING : EPOCH - 278 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,518 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,520 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,520 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,544 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,545 : INFO : EPOCH - 279 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,545 : WARNING : EPOCH - 279 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,554 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,556 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,556 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,580 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,581 : INFO : EPOCH - 280 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,581 : WARNING : EPOCH - 280 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,591 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,592 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,593 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,616 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,617 : INFO : EPOCH - 281 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,618 : WARNING : EPOCH - 281 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,627 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,628 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,629 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,655 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,656 : INFO : EPOCH - 282 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,657 : WARNING : EPOCH - 282 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,672 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,674 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,675 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,701 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,706 : INFO : EPOCH - 283 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,707 : WARNING : EPOCH - 283 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,716 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,717 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,719 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,745 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,746 : INFO : EPOCH - 284 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,748 : WARNING : EPOCH - 284 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:55,759 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,760 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,761 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,786 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,787 : INFO : EPOCH - 285 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,788 : WARNING : EPOCH - 285 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,797 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,798 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,799 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,825 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,826 : INFO : EPOCH - 286 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,827 : WARNING : EPOCH - 286 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,837 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,838 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,838 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,863 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,864 : INFO : EPOCH - 287 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,864 : WARNING : EPOCH - 287 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,873 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,875 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,876 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,901 : INFO : EPOCH - 288 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,902 : WARNING : EPOCH - 288 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,911 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,912 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,912 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,936 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,937 : INFO : EPOCH - 289 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,938 : WARNING : EPOCH - 289 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,950 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,952 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,952 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:55,979 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:55,987 : INFO : EPOCH - 290 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:55,988 : WARNING : EPOCH - 290 : supplied example count (1) did not equal expe
2019-02-21 17:16:55,998 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:55,999 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:55,999 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,025 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,026 : INFO : EPOCH - 291 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,027 : WARNING : EPOCH - 291 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,037 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,039 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,040 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,066 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,067 : INFO : EPOCH - 292 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,068 : WARNING : EPOCH - 292 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:56,077 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,078 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,079 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,105 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,106 : INFO : EPOCH - 293 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,106 : WARNING : EPOCH - 293 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,116 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,117 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,119 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,144 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,145 : INFO : EPOCH - 294 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,145 : WARNING : EPOCH - 294 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,153 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,154 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,155 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,180 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,181 : INFO : EPOCH - 295 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,182 : WARNING : EPOCH - 295 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,196 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,197 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,197 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,223 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,225 : INFO : EPOCH - 296 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,226 : WARNING : EPOCH - 296 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,233 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,235 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,236 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,261 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,263 : INFO : EPOCH - 297 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,263 : WARNING : EPOCH - 297 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,275 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,276 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,276 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,304 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,305 : INFO : EPOCH - 298 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,305 : WARNING : EPOCH - 298 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,315 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,316 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,317 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,343 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,344 : INFO : EPOCH - 299 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,345 : WARNING : EPOCH - 299 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,356 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,357 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,358 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,384 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,385 : INFO : EPOCH - 300 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,385 : WARNING : EPOCH - 300 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:56,395 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,396 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,397 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,422 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,423 : INFO : EPOCH - 301 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,423 : WARNING : EPOCH - 301 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,434 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,435 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,436 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,461 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,462 : INFO : EPOCH - 302 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,463 : WARNING : EPOCH - 302 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,472 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,473 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,497 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,498 : INFO : EPOCH - 303 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,499 : WARNING : EPOCH - 303 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,510 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,510 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,511 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,536 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,540 : INFO : EPOCH - 304 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,541 : WARNING : EPOCH - 304 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,550 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,551 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,552 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,577 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,578 : INFO : EPOCH - 305 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,578 : WARNING : EPOCH - 305 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,588 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,589 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,590 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,616 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,617 : INFO : EPOCH - 306 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,619 : WARNING : EPOCH - 306 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,629 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,630 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,631 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,657 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,659 : INFO : EPOCH - 307 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,659 : WARNING : EPOCH - 307 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,669 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,671 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,671 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,698 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,699 : INFO : EPOCH - 308 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:56,700 : WARNING : EPOCH - 308 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:56,710 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,711 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,711 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,740 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,741 : INFO : EPOCH - 309 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:56,742 : WARNING : EPOCH - 309 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,751 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,752 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,753 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,779 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,780 : INFO : EPOCH - 310 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:56,781 : WARNING : EPOCH - 310 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,791 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,792 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,793 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,820 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,822 : INFO : EPOCH - 311 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:56,822 : WARNING : EPOCH - 311 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,832 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,833 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,835 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,858 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,860 : INFO : EPOCH - 312 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:56,861 : WARNING : EPOCH - 312 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,869 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,870 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,871 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,896 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,896 : INFO : EPOCH - 313 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:56,897 : WARNING : EPOCH - 313 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,907 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,908 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,909 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,933 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,934 : INFO : EPOCH - 314 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:56,935 : WARNING : EPOCH - 314 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,944 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,945 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,946 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:56,970 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:56,972 : INFO : EPOCH - 315 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:56,973 : WARNING : EPOCH - 315 : supplied example count (1) did not equal expe
2019-02-21 17:16:56,984 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:56,986 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:56,987 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,012 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,013 : INFO : EPOCH - 316 : training on 5386246 raw words (10000 effective 
2019-02-21 17:16:57,014 : WARNING : EPOCH - 316 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:57,023 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,024 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,024 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,051 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,052 : INFO : EPOCH - 317 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,053 : WARNING : EPOCH - 317 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,063 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,065 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,066 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,090 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,092 : INFO : EPOCH - 318 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,092 : WARNING : EPOCH - 318 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,102 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,103 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,104 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,130 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,131 : INFO : EPOCH - 319 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,131 : WARNING : EPOCH - 319 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,145 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,146 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,146 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,172 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,173 : INFO : EPOCH - 320 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,173 : WARNING : EPOCH - 320 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,187 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,189 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,190 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,215 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,216 : INFO : EPOCH - 321 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,216 : WARNING : EPOCH - 321 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,226 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,227 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,227 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,252 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,253 : INFO : EPOCH - 322 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,254 : WARNING : EPOCH - 322 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,263 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,264 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,265 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,289 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,290 : INFO : EPOCH - 323 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,291 : WARNING : EPOCH - 323 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,301 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,302 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,303 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,327 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,328 : INFO : EPOCH - 324 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,328 : WARNING : EPOCH - 324 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:57,338 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,339 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,340 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,366 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,367 : INFO : EPOCH - 325 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:57,367 : WARNING : EPOCH - 325 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,391 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,392 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,393 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,419 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,420 : INFO : EPOCH - 326 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:57,420 : WARNING : EPOCH - 326 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,432 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,433 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,434 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,459 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,460 : INFO : EPOCH - 327 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:57,460 : WARNING : EPOCH - 327 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,473 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,473 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,474 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,500 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,502 : INFO : EPOCH - 328 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:57,502 : WARNING : EPOCH - 328 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,511 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,512 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,513 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,539 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,540 : INFO : EPOCH - 329 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:57,540 : WARNING : EPOCH - 329 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,552 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,553 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,554 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,579 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,580 : INFO : EPOCH - 330 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:57,581 : WARNING : EPOCH - 330 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,590 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,591 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,592 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,617 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,619 : INFO : EPOCH - 331 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:57,620 : WARNING : EPOCH - 331 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,630 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,631 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,631 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,658 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,659 : INFO : EPOCH - 332 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:57,660 : WARNING : EPOCH - 332 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:57,668 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,669 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,669 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,695 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,696 : INFO : EPOCH - 333 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,697 : WARNING : EPOCH - 333 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,707 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,708 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,708 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,734 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,735 : INFO : EPOCH - 334 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,736 : WARNING : EPOCH - 334 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,745 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,746 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,746 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,771 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,772 : INFO : EPOCH - 335 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,772 : WARNING : EPOCH - 335 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,782 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,783 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,784 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,809 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,810 : INFO : EPOCH - 336 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,810 : WARNING : EPOCH - 336 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,821 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,823 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,823 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,850 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,851 : INFO : EPOCH - 337 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,852 : WARNING : EPOCH - 337 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,862 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,863 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,864 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,890 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,891 : INFO : EPOCH - 338 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,892 : WARNING : EPOCH - 338 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,902 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,904 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,906 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,930 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,931 : INFO : EPOCH - 339 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,932 : WARNING : EPOCH - 339 : supplied example count (1) did not equal expe
2019-02-21 17:16:57,944 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,945 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,946 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:57,972 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:57,974 : INFO : EPOCH - 340 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:57,974 : WARNING : EPOCH - 340 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:57,984 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:57,985 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:57,986 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,011 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,012 : INFO : EPOCH - 341 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:58,012 : WARNING : EPOCH - 341 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,022 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,024 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,024 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,050 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,051 : INFO : EPOCH - 342 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:58,052 : WARNING : EPOCH - 342 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,061 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,062 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,063 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,089 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,090 : INFO : EPOCH - 343 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:58,091 : WARNING : EPOCH - 343 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,101 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,102 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,103 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,128 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,129 : INFO : EPOCH - 344 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:58,129 : WARNING : EPOCH - 344 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,140 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,141 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,141 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,166 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,167 : INFO : EPOCH - 345 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:58,168 : WARNING : EPOCH - 345 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,179 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,181 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,181 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,210 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,211 : INFO : EPOCH - 346 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:58,213 : WARNING : EPOCH - 346 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,224 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,225 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,225 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,251 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,256 : INFO : EPOCH - 347 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:58,257 : WARNING : EPOCH - 347 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,265 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,266 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,267 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,292 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,293 : INFO : EPOCH - 348 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:58,294 : WARNING : EPOCH - 348 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:58,307 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,309 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,310 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,335 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,336 : INFO : EPOCH - 349 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:58,337 : WARNING : EPOCH - 349 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,348 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,349 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,350 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,376 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,377 : INFO : EPOCH - 350 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:58,377 : WARNING : EPOCH - 350 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,387 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,388 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,389 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,414 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,414 : INFO : EPOCH - 351 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:58,415 : WARNING : EPOCH - 351 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,425 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,426 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,427 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,453 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,454 : INFO : EPOCH - 352 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:58,454 : WARNING : EPOCH - 352 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,464 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,465 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,466 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,492 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,493 : INFO : EPOCH - 353 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:58,494 : WARNING : EPOCH - 353 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,505 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,506 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,507 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,532 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,533 : INFO : EPOCH - 354 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:58,534 : WARNING : EPOCH - 354 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,544 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,545 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,545 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,572 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,573 : INFO : EPOCH - 355 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:58,573 : WARNING : EPOCH - 355 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,585 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,587 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,588 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,614 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,615 : INFO : EPOCH - 356 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:58,616 : WARNING : EPOCH - 356 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:58,626 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,627 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,627 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,652 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,653 : INFO : EPOCH - 357 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,654 : WARNING : EPOCH - 357 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,663 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,664 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,664 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,690 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,691 : INFO : EPOCH - 358 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,691 : WARNING : EPOCH - 358 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,701 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,702 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,703 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,727 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,728 : INFO : EPOCH - 359 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,729 : WARNING : EPOCH - 359 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,738 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,739 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,740 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,765 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,766 : INFO : EPOCH - 360 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,767 : WARNING : EPOCH - 360 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,775 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,776 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,777 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,804 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,810 : INFO : EPOCH - 361 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,811 : WARNING : EPOCH - 361 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,822 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,822 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,823 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,848 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,849 : INFO : EPOCH - 362 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,850 : WARNING : EPOCH - 362 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,859 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,860 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,861 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,887 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,888 : INFO : EPOCH - 363 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,888 : WARNING : EPOCH - 363 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,897 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,898 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,899 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,925 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,926 : INFO : EPOCH - 364 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,927 : WARNING : EPOCH - 364 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:58,936 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,937 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,938 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:58,964 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:58,965 : INFO : EPOCH - 365 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:58,965 : WARNING : EPOCH - 365 : supplied example count (1) did not equal expe
2019-02-21 17:16:58,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:58,977 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:58,978 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,003 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,004 : INFO : EPOCH - 366 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:59,005 : WARNING : EPOCH - 366 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,014 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,014 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,015 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,041 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,042 : INFO : EPOCH - 367 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:59,043 : WARNING : EPOCH - 367 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,058 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,059 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,059 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,085 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,093 : INFO : EPOCH - 368 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:59,094 : WARNING : EPOCH - 368 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,103 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,104 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,104 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,130 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,131 : INFO : EPOCH - 369 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:59,132 : WARNING : EPOCH - 369 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,141 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,142 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,143 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,169 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,170 : INFO : EPOCH - 370 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:59,171 : WARNING : EPOCH - 370 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,183 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,184 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,185 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,209 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,210 : INFO : EPOCH - 371 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:59,211 : WARNING : EPOCH - 371 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,220 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,221 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,221 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,245 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,246 : INFO : EPOCH - 372 : training on 5386246 raw words (10000 effective u
2019-02-21 17:16:59,247 : WARNING : EPOCH - 372 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:59,257 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,258 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,259 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,283 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,284 : INFO : EPOCH - 373 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:59,285 : WARNING : EPOCH - 373 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,294 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,295 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,296 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,321 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,323 : INFO : EPOCH - 374 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:59,323 : WARNING : EPOCH - 374 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,332 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,333 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,333 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,359 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,366 : INFO : EPOCH - 375 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:59,366 : WARNING : EPOCH - 375 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,375 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,376 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,377 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,401 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,402 : INFO : EPOCH - 376 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:59,403 : WARNING : EPOCH - 376 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,411 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,413 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,414 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,438 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,438 : INFO : EPOCH - 377 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:59,439 : WARNING : EPOCH - 377 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,448 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,450 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,450 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,489 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,505 : INFO : EPOCH - 378 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:59,506 : WARNING : EPOCH - 378 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,515 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,516 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,517 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,542 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,543 : INFO : EPOCH - 379 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:59,544 : WARNING : EPOCH - 379 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,554 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,555 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,556 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,581 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,582 : INFO : EPOCH - 380 : training on 5386246 raw words (10000 effective w
2019-02-21 17:16:59,583 : WARNING : EPOCH - 380 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:59,593 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,594 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,594 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,620 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,621 : INFO : EPOCH - 381 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,622 : WARNING : EPOCH - 381 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,632 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,632 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,633 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,659 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,667 : INFO : EPOCH - 382 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,667 : WARNING : EPOCH - 382 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,682 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,684 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,685 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,711 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,712 : INFO : EPOCH - 383 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,712 : WARNING : EPOCH - 383 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,723 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,724 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,725 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,750 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,751 : INFO : EPOCH - 384 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,752 : WARNING : EPOCH - 384 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,761 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,762 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,763 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,789 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,790 : INFO : EPOCH - 385 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,790 : WARNING : EPOCH - 385 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,799 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,800 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,801 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,825 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,826 : INFO : EPOCH - 386 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,827 : WARNING : EPOCH - 386 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,836 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,837 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,838 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,863 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,864 : INFO : EPOCH - 387 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,865 : WARNING : EPOCH - 387 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,874 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,875 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,876 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,901 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,902 : INFO : EPOCH - 388 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,903 : WARNING : EPOCH - 388 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:16:59,912 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,913 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,914 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,940 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,941 : INFO : EPOCH - 389 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,942 : WARNING : EPOCH - 389 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,951 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,953 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,953 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:16:59,979 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:16:59,980 : INFO : EPOCH - 390 : training on 5386246 raw words (10000 effective
2019-02-21 17:16:59,981 : WARNING : EPOCH - 390 : supplied example count (1) did not equal expe
2019-02-21 17:16:59,991 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:16:59,992 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:16:59,993 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,021 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,021 : INFO : EPOCH - 391 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,022 : WARNING : EPOCH - 391 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,031 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,033 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,034 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,058 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,059 : INFO : EPOCH - 392 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,060 : WARNING : EPOCH - 392 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,069 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,070 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,070 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,095 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,096 : INFO : EPOCH - 393 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,097 : WARNING : EPOCH - 393 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,108 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,109 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,110 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,135 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,136 : INFO : EPOCH - 394 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,136 : WARNING : EPOCH - 394 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,146 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,147 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,147 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,173 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,174 : INFO : EPOCH - 395 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,175 : WARNING : EPOCH - 395 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,189 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,190 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,191 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,216 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,217 : INFO : EPOCH - 396 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,218 : WARNING : EPOCH - 396 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:00,227 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,228 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,228 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,255 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,256 : INFO : EPOCH - 397 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,256 : WARNING : EPOCH - 397 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,266 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,268 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,269 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,293 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,294 : INFO : EPOCH - 398 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,294 : WARNING : EPOCH - 398 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,305 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,306 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,307 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,332 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,333 : INFO : EPOCH - 399 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,333 : WARNING : EPOCH - 399 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,343 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,343 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,344 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,369 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,370 : INFO : EPOCH - 400 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,371 : WARNING : EPOCH - 400 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,379 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,380 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,381 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,405 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,406 : INFO : EPOCH - 401 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,407 : WARNING : EPOCH - 401 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,416 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,417 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,418 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,442 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,443 : INFO : EPOCH - 402 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,444 : WARNING : EPOCH - 402 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,453 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,453 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,454 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,478 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,479 : INFO : EPOCH - 403 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,480 : WARNING : EPOCH - 403 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,492 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,493 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,494 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,520 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,524 : INFO : EPOCH - 404 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,525 : WARNING : EPOCH - 404 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:00,533 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,534 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,535 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,560 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,561 : INFO : EPOCH - 405 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,561 : WARNING : EPOCH - 405 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,571 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,572 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,573 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,603 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,605 : INFO : EPOCH - 406 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,606 : WARNING : EPOCH - 406 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,620 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,623 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,626 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,649 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,650 : INFO : EPOCH - 407 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,651 : WARNING : EPOCH - 407 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,659 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,660 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,662 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,687 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,688 : INFO : EPOCH - 408 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,689 : WARNING : EPOCH - 408 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,706 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,707 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,708 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,732 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,733 : INFO : EPOCH - 409 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,733 : WARNING : EPOCH - 409 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,743 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,744 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,744 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,770 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,771 : INFO : EPOCH - 410 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,772 : WARNING : EPOCH - 410 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,783 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,785 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,786 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,810 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,811 : INFO : EPOCH - 411 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,812 : WARNING : EPOCH - 411 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,821 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,822 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,823 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,849 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,850 : INFO : EPOCH - 412 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:00,851 : WARNING : EPOCH - 412 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:00,863 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,864 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,864 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,890 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,891 : INFO : EPOCH - 413 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,891 : WARNING : EPOCH - 413 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,901 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,902 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,902 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,927 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,928 : INFO : EPOCH - 414 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,928 : WARNING : EPOCH - 414 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,939 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,940 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,941 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:00,966 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:00,967 : INFO : EPOCH - 415 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:00,968 : WARNING : EPOCH - 415 : supplied example count (1) did not equal expe
2019-02-21 17:17:00,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:00,978 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:00,978 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,004 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,005 : INFO : EPOCH - 416 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:01,006 : WARNING : EPOCH - 416 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,015 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,016 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,016 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,041 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,042 : INFO : EPOCH - 417 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:01,043 : WARNING : EPOCH - 417 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,054 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,054 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,055 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,080 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,081 : INFO : EPOCH - 418 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:01,081 : WARNING : EPOCH - 418 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,091 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,105 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,106 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,119 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,120 : INFO : EPOCH - 419 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:01,120 : WARNING : EPOCH - 419 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,129 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,130 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,131 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,155 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,156 : INFO : EPOCH - 420 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:01,157 : WARNING : EPOCH - 420 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:01,166 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,167 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,169 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,197 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,197 : INFO : EPOCH - 421 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:01,198 : WARNING : EPOCH - 421 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,208 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,209 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,210 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,235 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,236 : INFO : EPOCH - 422 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:01,237 : WARNING : EPOCH - 422 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,247 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,248 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,249 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,274 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,275 : INFO : EPOCH - 423 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:01,276 : WARNING : EPOCH - 423 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,287 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,288 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,288 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,312 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,313 : INFO : EPOCH - 424 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:01,314 : WARNING : EPOCH - 424 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,324 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,325 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,326 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,351 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,351 : INFO : EPOCH - 425 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:01,352 : WARNING : EPOCH - 425 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,361 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,363 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,363 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,390 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,391 : INFO : EPOCH - 426 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:01,392 : WARNING : EPOCH - 426 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,401 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,402 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,403 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,428 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,429 : INFO : EPOCH - 427 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:01,429 : WARNING : EPOCH - 427 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,439 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,440 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,441 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,464 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,466 : INFO : EPOCH - 428 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:01,466 : WARNING : EPOCH - 428 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:01,475 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,476 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,477 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,503 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,504 : INFO : EPOCH - 429 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,505 : WARNING : EPOCH - 429 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,515 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,516 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,517 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,542 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,543 : INFO : EPOCH - 430 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,544 : WARNING : EPOCH - 430 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,552 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,556 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,557 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,583 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,584 : INFO : EPOCH - 431 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,585 : WARNING : EPOCH - 431 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,593 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,594 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,595 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,621 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,622 : INFO : EPOCH - 432 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,623 : WARNING : EPOCH - 432 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,636 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,637 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,637 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,664 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,671 : INFO : EPOCH - 433 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,672 : WARNING : EPOCH - 433 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,681 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,682 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,683 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,708 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,709 : INFO : EPOCH - 434 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,710 : WARNING : EPOCH - 434 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,719 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,720 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,721 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,746 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,747 : INFO : EPOCH - 435 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,747 : WARNING : EPOCH - 435 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,758 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,759 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,760 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,786 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,787 : INFO : EPOCH - 436 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,788 : WARNING : EPOCH - 436 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:01,798 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,799 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,800 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,824 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,825 : INFO : EPOCH - 437 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,826 : WARNING : EPOCH - 437 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,837 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,838 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,839 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,863 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,864 : INFO : EPOCH - 438 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,864 : WARNING : EPOCH - 438 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,874 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,875 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,875 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,901 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,904 : INFO : EPOCH - 439 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,905 : WARNING : EPOCH - 439 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,914 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,915 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,916 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,942 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,943 : INFO : EPOCH - 440 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,944 : WARNING : EPOCH - 440 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,955 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,956 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,956 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:01,982 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:01,983 : INFO : EPOCH - 441 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:01,984 : WARNING : EPOCH - 441 : supplied example count (1) did not equal expe
2019-02-21 17:17:01,994 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:01,995 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:01,996 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,022 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,023 : INFO : EPOCH - 442 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,024 : WARNING : EPOCH - 442 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,032 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,034 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,034 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,059 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,060 : INFO : EPOCH - 443 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,060 : WARNING : EPOCH - 443 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,069 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,070 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,071 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,097 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,097 : INFO : EPOCH - 444 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,098 : WARNING : EPOCH - 444 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:02,109 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,110 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,111 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,136 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,137 : INFO : EPOCH - 445 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,137 : WARNING : EPOCH - 445 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,146 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,147 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,148 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,172 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,173 : INFO : EPOCH - 446 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,174 : WARNING : EPOCH - 446 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,192 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,193 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,193 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,221 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,224 : INFO : EPOCH - 447 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,225 : WARNING : EPOCH - 447 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,236 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,238 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,238 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,263 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,264 : INFO : EPOCH - 448 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,265 : WARNING : EPOCH - 448 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,276 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,277 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,277 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,303 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,304 : INFO : EPOCH - 449 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,304 : WARNING : EPOCH - 449 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,313 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,314 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,315 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,340 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,341 : INFO : EPOCH - 450 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,342 : WARNING : EPOCH - 450 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,352 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,353 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,354 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,380 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,380 : INFO : EPOCH - 451 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,381 : WARNING : EPOCH - 451 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,391 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,392 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,393 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,418 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,419 : INFO : EPOCH - 452 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:02,419 : WARNING : EPOCH - 452 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:02,429 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,430 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,431 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,456 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,457 : INFO : EPOCH - 453 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,458 : WARNING : EPOCH - 453 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,473 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,474 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,475 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,501 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,510 : INFO : EPOCH - 454 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,511 : WARNING : EPOCH - 454 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,520 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,521 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,522 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,546 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,547 : INFO : EPOCH - 455 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,547 : WARNING : EPOCH - 455 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,556 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,558 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,558 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,583 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,584 : INFO : EPOCH - 456 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,585 : WARNING : EPOCH - 456 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,594 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,594 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,595 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,619 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,620 : INFO : EPOCH - 457 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,621 : WARNING : EPOCH - 457 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,630 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,631 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,632 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,656 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,657 : INFO : EPOCH - 458 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,658 : WARNING : EPOCH - 458 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,666 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,667 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,668 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,692 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,693 : INFO : EPOCH - 459 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,693 : WARNING : EPOCH - 459 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,702 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,703 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,703 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,727 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,728 : INFO : EPOCH - 460 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,729 : WARNING : EPOCH - 460 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:02,739 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,740 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,740 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,767 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,768 : INFO : EPOCH - 461 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,769 : WARNING : EPOCH - 461 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,785 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,786 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,787 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,812 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,812 : INFO : EPOCH - 462 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,813 : WARNING : EPOCH - 462 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,823 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,824 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,825 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,850 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,851 : INFO : EPOCH - 463 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,852 : WARNING : EPOCH - 463 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,864 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,865 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,865 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,890 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,891 : INFO : EPOCH - 464 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,892 : WARNING : EPOCH - 464 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,901 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,903 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,903 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,928 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,929 : INFO : EPOCH - 465 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,930 : WARNING : EPOCH - 465 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,940 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,941 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,942 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:02,966 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:02,968 : INFO : EPOCH - 466 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:02,968 : WARNING : EPOCH - 466 : supplied example count (1) did not equal expe
2019-02-21 17:17:02,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:02,978 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:02,978 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,003 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,004 : INFO : EPOCH - 467 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,004 : WARNING : EPOCH - 467 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,013 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,014 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,015 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,041 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,042 : INFO : EPOCH - 468 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,042 : WARNING : EPOCH - 468 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:03,053 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,055 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,063 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,080 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,081 : INFO : EPOCH - 469 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:03,082 : WARNING : EPOCH - 469 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,091 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,092 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,094 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,118 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,119 : INFO : EPOCH - 470 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:03,120 : WARNING : EPOCH - 470 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,129 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,129 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,130 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,156 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,157 : INFO : EPOCH - 471 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:03,158 : WARNING : EPOCH - 471 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,168 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,169 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,170 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,196 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,197 : INFO : EPOCH - 472 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:03,198 : WARNING : EPOCH - 472 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,208 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,210 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,210 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,236 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,237 : INFO : EPOCH - 473 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:03,238 : WARNING : EPOCH - 473 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,246 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,247 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,248 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,273 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,274 : INFO : EPOCH - 474 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:03,275 : WARNING : EPOCH - 474 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,284 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,285 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,286 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,311 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,311 : INFO : EPOCH - 475 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:03,312 : WARNING : EPOCH - 475 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,322 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,323 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,324 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,349 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,350 : INFO : EPOCH - 476 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:03,351 : WARNING : EPOCH - 476 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:03,359 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,360 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,361 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,386 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,387 : INFO : EPOCH - 477 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,388 : WARNING : EPOCH - 477 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,399 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,400 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,401 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,426 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,427 : INFO : EPOCH - 478 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,427 : WARNING : EPOCH - 478 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,438 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,439 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,440 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,464 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,465 : INFO : EPOCH - 479 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,466 : WARNING : EPOCH - 479 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,476 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,477 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,477 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,503 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,503 : INFO : EPOCH - 480 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,504 : WARNING : EPOCH - 480 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,514 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,515 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,516 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,540 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,541 : INFO : EPOCH - 481 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,542 : WARNING : EPOCH - 481 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,552 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,553 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,553 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,578 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,579 : INFO : EPOCH - 482 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,582 : WARNING : EPOCH - 482 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,590 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,592 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,592 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,619 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,620 : INFO : EPOCH - 483 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,621 : WARNING : EPOCH - 483 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,632 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,634 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,635 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,660 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,661 : INFO : EPOCH - 484 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:03,662 : WARNING : EPOCH - 484 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:03,673 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,674 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,674 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,699 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,700 : INFO : EPOCH - 485 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:03,701 : WARNING : EPOCH - 485 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,711 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,712 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,712 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,737 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,738 : INFO : EPOCH - 486 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:03,739 : WARNING : EPOCH - 486 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,747 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,750 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,752 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,777 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,778 : INFO : EPOCH - 487 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:03,779 : WARNING : EPOCH - 487 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,791 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,792 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,792 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,817 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,818 : INFO : EPOCH - 488 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:03,819 : WARNING : EPOCH - 488 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,827 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,829 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,829 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,854 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,855 : INFO : EPOCH - 489 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:03,856 : WARNING : EPOCH - 489 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,866 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,867 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,868 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,894 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,895 : INFO : EPOCH - 490 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:03,896 : WARNING : EPOCH - 490 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,906 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,906 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,907 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,932 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,933 : INFO : EPOCH - 491 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:03,934 : WARNING : EPOCH - 491 : supplied example count (1) did not equal expe
2019-02-21 17:17:03,943 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,944 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,945 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:03,970 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:03,971 : INFO : EPOCH - 492 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:03,972 : WARNING : EPOCH - 492 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:03,981 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:03,982 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:03,983 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,008 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,009 : INFO : EPOCH - 493 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:04,009 : WARNING : EPOCH - 493 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,021 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,022 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,023 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,048 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,049 : INFO : EPOCH - 494 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:04,049 : WARNING : EPOCH - 494 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,058 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,059 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,060 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,085 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,086 : INFO : EPOCH - 495 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:04,087 : WARNING : EPOCH - 495 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,096 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,097 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,097 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,122 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,123 : INFO : EPOCH - 496 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:04,123 : WARNING : EPOCH - 496 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,132 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,133 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,134 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,159 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,160 : INFO : EPOCH - 497 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:04,161 : WARNING : EPOCH - 497 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,170 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,171 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,173 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,199 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,201 : INFO : EPOCH - 498 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:04,204 : WARNING : EPOCH - 498 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,213 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,214 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,215 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,242 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,243 : INFO : EPOCH - 499 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:04,244 : WARNING : EPOCH - 499 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,255 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,256 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,258 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,283 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,284 : INFO : EPOCH - 500 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:04,284 : WARNING : EPOCH - 500 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:04,294 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,295 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,296 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,345 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,348 : INFO : EPOCH - 501 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,349 : WARNING : EPOCH - 501 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,357 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,359 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,360 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,385 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,386 : INFO : EPOCH - 502 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,386 : WARNING : EPOCH - 502 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,396 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,398 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,398 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,423 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,424 : INFO : EPOCH - 503 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,425 : WARNING : EPOCH - 503 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,434 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,435 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,435 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,460 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,461 : INFO : EPOCH - 504 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,461 : WARNING : EPOCH - 504 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,472 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,473 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,498 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,499 : INFO : EPOCH - 505 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,500 : WARNING : EPOCH - 505 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,508 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,509 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,510 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,550 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,560 : INFO : EPOCH - 506 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,560 : WARNING : EPOCH - 506 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,569 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,571 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,572 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,595 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,596 : INFO : EPOCH - 507 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,597 : WARNING : EPOCH - 507 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,608 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,609 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,609 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,635 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,636 : INFO : EPOCH - 508 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,637 : WARNING : EPOCH - 508 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:04,645 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,647 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,649 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,673 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,674 : INFO : EPOCH - 509 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,675 : WARNING : EPOCH - 509 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,683 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,685 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,686 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,711 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,712 : INFO : EPOCH - 510 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,712 : WARNING : EPOCH - 510 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,721 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,723 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,723 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,746 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,747 : INFO : EPOCH - 511 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,748 : WARNING : EPOCH - 511 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,758 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,759 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,759 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,785 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,786 : INFO : EPOCH - 512 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,789 : WARNING : EPOCH - 512 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,797 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,798 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,799 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,825 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,827 : INFO : EPOCH - 513 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,827 : WARNING : EPOCH - 513 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,837 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,838 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,838 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,862 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,863 : INFO : EPOCH - 514 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,864 : WARNING : EPOCH - 514 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,873 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,874 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,875 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,901 : INFO : EPOCH - 515 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,902 : WARNING : EPOCH - 515 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,911 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,912 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,912 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,938 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,939 : INFO : EPOCH - 516 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,939 : WARNING : EPOCH - 516 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:04,948 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,951 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,951 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:04,977 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:04,977 : INFO : EPOCH - 517 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:04,978 : WARNING : EPOCH - 517 : supplied example count (1) did not equal expe
2019-02-21 17:17:04,987 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:04,988 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:04,988 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,012 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,013 : INFO : EPOCH - 518 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,014 : WARNING : EPOCH - 518 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,024 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,026 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,027 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,052 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,054 : INFO : EPOCH - 519 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,055 : WARNING : EPOCH - 519 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,063 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,064 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,065 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,091 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,092 : INFO : EPOCH - 520 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,093 : WARNING : EPOCH - 520 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,102 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,104 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,104 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,128 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,129 : INFO : EPOCH - 521 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,130 : WARNING : EPOCH - 521 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,139 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,140 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,142 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,166 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,167 : INFO : EPOCH - 522 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,168 : WARNING : EPOCH - 522 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,179 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,180 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,180 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,206 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,207 : INFO : EPOCH - 523 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,208 : WARNING : EPOCH - 523 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,217 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,219 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,219 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,244 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,245 : INFO : EPOCH - 524 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,245 : WARNING : EPOCH - 524 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:05,254 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,255 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,256 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,281 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,282 : INFO : EPOCH - 525 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:05,283 : WARNING : EPOCH - 525 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,292 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,293 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,294 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,319 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,321 : INFO : EPOCH - 526 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:05,322 : WARNING : EPOCH - 526 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,331 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,333 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,334 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,358 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,363 : INFO : EPOCH - 527 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:05,363 : WARNING : EPOCH - 527 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,373 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,374 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,375 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,400 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,401 : INFO : EPOCH - 528 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:05,402 : WARNING : EPOCH - 528 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,411 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,412 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,413 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,438 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,439 : INFO : EPOCH - 529 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:05,440 : WARNING : EPOCH - 529 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,447 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,449 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,450 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,475 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,475 : INFO : EPOCH - 530 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:05,476 : WARNING : EPOCH - 530 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,486 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,487 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,488 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,512 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,513 : INFO : EPOCH - 531 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:05,514 : WARNING : EPOCH - 531 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,524 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,525 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,526 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,550 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,551 : INFO : EPOCH - 532 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:05,552 : WARNING : EPOCH - 532 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:05,561 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,562 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,562 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,589 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,591 : INFO : EPOCH - 533 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,592 : WARNING : EPOCH - 533 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,601 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,602 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,603 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,629 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,630 : INFO : EPOCH - 534 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,631 : WARNING : EPOCH - 534 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,640 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,641 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,642 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,667 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,668 : INFO : EPOCH - 535 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,669 : WARNING : EPOCH - 535 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,678 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,679 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,680 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,705 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,706 : INFO : EPOCH - 536 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,707 : WARNING : EPOCH - 536 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,716 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,717 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,718 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,743 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,744 : INFO : EPOCH - 537 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,745 : WARNING : EPOCH - 537 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,755 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,755 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,756 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,782 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,783 : INFO : EPOCH - 538 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,784 : WARNING : EPOCH - 538 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,794 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,794 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,795 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,820 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,821 : INFO : EPOCH - 539 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,822 : WARNING : EPOCH - 539 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,831 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,832 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,832 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,858 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,859 : INFO : EPOCH - 540 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:05,859 : WARNING : EPOCH - 540 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:05,870 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,871 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,872 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,897 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,905 : INFO : EPOCH - 541 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:05,906 : WARNING : EPOCH - 541 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,915 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,916 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,916 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,941 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,942 : INFO : EPOCH - 542 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:05,943 : WARNING : EPOCH - 542 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,951 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,952 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,953 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:05,978 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:05,979 : INFO : EPOCH - 543 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:05,980 : WARNING : EPOCH - 543 : supplied example count (1) did not equal expe
2019-02-21 17:17:05,988 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:05,990 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:05,990 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,015 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,016 : INFO : EPOCH - 544 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,018 : WARNING : EPOCH - 544 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,026 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,027 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,028 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,054 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,055 : INFO : EPOCH - 545 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,056 : WARNING : EPOCH - 545 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,065 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,066 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,067 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,091 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,092 : INFO : EPOCH - 546 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,093 : WARNING : EPOCH - 546 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,102 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,103 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,104 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,130 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,131 : INFO : EPOCH - 547 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,131 : WARNING : EPOCH - 547 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,142 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,143 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,143 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,169 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,170 : INFO : EPOCH - 548 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,170 : WARNING : EPOCH - 548 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:06,181 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,184 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,185 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,210 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,211 : INFO : EPOCH - 549 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,212 : WARNING : EPOCH - 549 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,224 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,225 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,225 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,250 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,251 : INFO : EPOCH - 550 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,252 : WARNING : EPOCH - 550 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,260 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,261 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,262 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,288 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,289 : INFO : EPOCH - 551 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,289 : WARNING : EPOCH - 551 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,298 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,299 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,301 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,324 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,325 : INFO : EPOCH - 552 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,326 : WARNING : EPOCH - 552 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,335 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,336 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,337 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,362 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,363 : INFO : EPOCH - 553 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,364 : WARNING : EPOCH - 553 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,372 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,374 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,374 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,400 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,403 : INFO : EPOCH - 554 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,404 : WARNING : EPOCH - 554 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,413 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,414 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,415 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,441 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,442 : INFO : EPOCH - 555 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,442 : WARNING : EPOCH - 555 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,451 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,452 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,453 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,478 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,479 : INFO : EPOCH - 556 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,480 : WARNING : EPOCH - 556 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:06,490 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,491 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,493 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,517 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,518 : INFO : EPOCH - 557 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,519 : WARNING : EPOCH - 557 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,528 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,529 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,530 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,554 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,555 : INFO : EPOCH - 558 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,556 : WARNING : EPOCH - 558 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,564 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,565 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,565 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,590 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,591 : INFO : EPOCH - 559 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,592 : WARNING : EPOCH - 559 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,601 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,601 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,602 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,627 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,628 : INFO : EPOCH - 560 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,628 : WARNING : EPOCH - 560 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,638 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,639 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,639 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,665 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,666 : INFO : EPOCH - 561 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,666 : WARNING : EPOCH - 561 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,676 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,680 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,681 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,703 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,705 : INFO : EPOCH - 562 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,705 : WARNING : EPOCH - 562 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,715 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,723 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,723 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,742 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,742 : INFO : EPOCH - 563 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,743 : WARNING : EPOCH - 563 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,753 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,753 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,754 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,780 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,781 : INFO : EPOCH - 564 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:06,782 : WARNING : EPOCH - 564 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:06,793 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,795 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,796 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,821 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,822 : INFO : EPOCH - 565 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:06,823 : WARNING : EPOCH - 565 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,833 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,834 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,835 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,860 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,861 : INFO : EPOCH - 566 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:06,862 : WARNING : EPOCH - 566 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,871 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,872 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,873 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,898 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,898 : INFO : EPOCH - 567 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:06,899 : WARNING : EPOCH - 567 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,908 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,910 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,910 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,935 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,936 : INFO : EPOCH - 568 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:06,936 : WARNING : EPOCH - 568 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,947 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,948 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,948 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:06,973 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:06,974 : INFO : EPOCH - 569 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:06,975 : WARNING : EPOCH - 569 : supplied example count (1) did not equal expe
2019-02-21 17:17:06,983 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:06,984 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:06,985 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,010 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,011 : INFO : EPOCH - 570 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,012 : WARNING : EPOCH - 570 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,024 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,025 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,026 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,051 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,052 : INFO : EPOCH - 571 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,053 : WARNING : EPOCH - 571 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,064 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,066 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,066 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,091 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,091 : INFO : EPOCH - 572 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,093 : WARNING : EPOCH - 572 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:07,102 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,103 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,104 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,128 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,129 : INFO : EPOCH - 573 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,130 : WARNING : EPOCH - 573 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,139 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,140 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,141 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,165 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,166 : INFO : EPOCH - 574 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,167 : WARNING : EPOCH - 574 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,177 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,178 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,179 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,206 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,207 : INFO : EPOCH - 575 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,208 : WARNING : EPOCH - 575 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,217 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,218 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,219 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,244 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,244 : INFO : EPOCH - 576 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,245 : WARNING : EPOCH - 576 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,255 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,256 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,257 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,283 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,284 : INFO : EPOCH - 577 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,284 : WARNING : EPOCH - 577 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,294 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,295 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,295 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,322 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,323 : INFO : EPOCH - 578 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,323 : WARNING : EPOCH - 578 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,333 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,334 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,335 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,358 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,359 : INFO : EPOCH - 579 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,360 : WARNING : EPOCH - 579 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,369 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,370 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,371 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,397 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,398 : INFO : EPOCH - 580 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,399 : WARNING : EPOCH - 580 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:07,408 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,409 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,410 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,435 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,435 : INFO : EPOCH - 581 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,435 : WARNING : EPOCH - 581 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,445 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,447 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,448 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,473 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,474 : INFO : EPOCH - 582 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,475 : WARNING : EPOCH - 582 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,493 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,497 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,498 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,521 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,521 : INFO : EPOCH - 583 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,522 : WARNING : EPOCH - 583 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,535 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,536 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,537 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,562 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,573 : INFO : EPOCH - 584 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,573 : WARNING : EPOCH - 584 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,584 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,586 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,587 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,614 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,615 : INFO : EPOCH - 585 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,616 : WARNING : EPOCH - 585 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,625 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,626 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,626 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,653 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,654 : INFO : EPOCH - 586 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,655 : WARNING : EPOCH - 586 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,668 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,669 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,671 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,696 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,696 : INFO : EPOCH - 587 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,697 : WARNING : EPOCH - 587 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,707 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,708 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,708 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,735 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,736 : INFO : EPOCH - 588 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:07,737 : WARNING : EPOCH - 588 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:07,747 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,748 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,749 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,775 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,776 : INFO : EPOCH - 589 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:07,777 : WARNING : EPOCH - 589 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,785 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,787 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,787 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,813 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,814 : INFO : EPOCH - 590 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:07,815 : WARNING : EPOCH - 590 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,825 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,826 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,827 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,853 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,855 : INFO : EPOCH - 591 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:07,855 : WARNING : EPOCH - 591 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,865 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,867 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,868 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,893 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,894 : INFO : EPOCH - 592 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:07,895 : WARNING : EPOCH - 592 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,903 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,904 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,905 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,930 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,930 : INFO : EPOCH - 593 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:07,931 : WARNING : EPOCH - 593 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,941 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,942 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,943 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:07,968 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:07,969 : INFO : EPOCH - 594 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:07,970 : WARNING : EPOCH - 594 : supplied example count (1) did not equal expe
2019-02-21 17:17:07,978 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:07,980 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:07,980 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,006 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,007 : INFO : EPOCH - 595 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,007 : WARNING : EPOCH - 595 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,016 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,017 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,018 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,043 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,043 : INFO : EPOCH - 596 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,044 : WARNING : EPOCH - 596 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:08,054 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,055 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,055 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,080 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,081 : INFO : EPOCH - 597 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,082 : WARNING : EPOCH - 597 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,090 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,091 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,092 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,117 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,118 : INFO : EPOCH - 598 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,118 : WARNING : EPOCH - 598 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,127 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,128 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,128 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,155 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,156 : INFO : EPOCH - 599 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,157 : WARNING : EPOCH - 599 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,166 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,175 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,176 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,195 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,196 : INFO : EPOCH - 600 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,197 : WARNING : EPOCH - 600 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,211 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,212 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,213 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,237 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,238 : INFO : EPOCH - 601 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,239 : WARNING : EPOCH - 601 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,246 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,248 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,250 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,275 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,275 : INFO : EPOCH - 602 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,276 : WARNING : EPOCH - 602 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,285 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,286 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,287 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,312 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,312 : INFO : EPOCH - 603 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,313 : WARNING : EPOCH - 603 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,323 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,324 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,325 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,350 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,351 : INFO : EPOCH - 604 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,352 : WARNING : EPOCH - 604 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:08,360 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,362 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,362 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,387 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,388 : INFO : EPOCH - 605 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,389 : WARNING : EPOCH - 605 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,397 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,398 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,399 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,424 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,425 : INFO : EPOCH - 606 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,426 : WARNING : EPOCH - 606 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,437 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,438 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,439 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,464 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,465 : INFO : EPOCH - 607 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,466 : WARNING : EPOCH - 607 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,478 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,479 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,480 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,505 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,506 : INFO : EPOCH - 608 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,506 : WARNING : EPOCH - 608 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,517 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,518 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,519 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,544 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,545 : INFO : EPOCH - 609 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,546 : WARNING : EPOCH - 609 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,556 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,557 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,557 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,582 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,583 : INFO : EPOCH - 610 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,584 : WARNING : EPOCH - 610 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,592 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,594 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,594 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,624 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,625 : INFO : EPOCH - 611 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,645 : WARNING : EPOCH - 611 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,658 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,660 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,660 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,685 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,686 : INFO : EPOCH - 612 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:08,686 : WARNING : EPOCH - 612 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:08,695 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,696 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,697 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,722 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,725 : INFO : EPOCH - 613 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:08,725 : WARNING : EPOCH - 613 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,734 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,735 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,736 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,761 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,778 : INFO : EPOCH - 614 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:08,779 : WARNING : EPOCH - 614 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,791 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,792 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,793 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,818 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,818 : INFO : EPOCH - 615 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:08,819 : WARNING : EPOCH - 615 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,827 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,828 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,829 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,854 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,854 : INFO : EPOCH - 616 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:08,855 : WARNING : EPOCH - 616 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,865 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,866 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,868 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,892 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,892 : INFO : EPOCH - 617 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:08,893 : WARNING : EPOCH - 617 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,902 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,903 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,904 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,929 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,930 : INFO : EPOCH - 618 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:08,930 : WARNING : EPOCH - 618 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,940 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,941 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,941 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:08,966 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:08,967 : INFO : EPOCH - 619 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:08,968 : WARNING : EPOCH - 619 : supplied example count (1) did not equal expe
2019-02-21 17:17:08,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:08,977 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:08,978 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,003 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,004 : INFO : EPOCH - 620 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,005 : WARNING : EPOCH - 620 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:09,013 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,014 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,014 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,040 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,041 : INFO : EPOCH - 621 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,042 : WARNING : EPOCH - 621 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,051 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,052 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,053 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,078 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,079 : INFO : EPOCH - 622 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,080 : WARNING : EPOCH - 622 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,090 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,090 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,091 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,116 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,117 : INFO : EPOCH - 623 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,117 : WARNING : EPOCH - 623 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,126 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,127 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,128 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,152 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,153 : INFO : EPOCH - 624 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,153 : WARNING : EPOCH - 624 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,162 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,163 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,164 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,191 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,192 : INFO : EPOCH - 625 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,192 : WARNING : EPOCH - 625 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,202 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,204 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,204 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,229 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,230 : INFO : EPOCH - 626 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,230 : WARNING : EPOCH - 626 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,240 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,241 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,241 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,266 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,267 : INFO : EPOCH - 627 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,268 : WARNING : EPOCH - 627 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,277 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,278 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,279 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,305 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,306 : INFO : EPOCH - 628 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,306 : WARNING : EPOCH - 628 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:09,319 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,320 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,321 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,347 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,350 : INFO : EPOCH - 629 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,351 : WARNING : EPOCH - 629 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,360 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,361 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,361 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,387 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,388 : INFO : EPOCH - 630 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,388 : WARNING : EPOCH - 630 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,397 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,398 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,398 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,424 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,425 : INFO : EPOCH - 631 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,426 : WARNING : EPOCH - 631 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,435 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,436 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,437 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,462 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,463 : INFO : EPOCH - 632 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,464 : WARNING : EPOCH - 632 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,472 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,474 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,474 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,499 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,501 : INFO : EPOCH - 633 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,501 : WARNING : EPOCH - 633 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,509 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,511 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,511 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,535 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,536 : INFO : EPOCH - 634 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,537 : WARNING : EPOCH - 634 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,545 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,546 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,547 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,571 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,572 : INFO : EPOCH - 635 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,573 : WARNING : EPOCH - 635 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,587 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,588 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,588 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,613 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,614 : INFO : EPOCH - 636 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:09,615 : WARNING : EPOCH - 636 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:09,625 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,627 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,627 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,651 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,652 : INFO : EPOCH - 637 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,653 : WARNING : EPOCH - 637 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,664 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,665 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,666 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,690 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,691 : INFO : EPOCH - 638 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,691 : WARNING : EPOCH - 638 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,702 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,702 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,704 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,727 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,728 : INFO : EPOCH - 639 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,729 : WARNING : EPOCH - 639 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,737 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,738 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,739 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,764 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,764 : INFO : EPOCH - 640 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,765 : WARNING : EPOCH - 640 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,775 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,776 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,777 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,801 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,802 : INFO : EPOCH - 641 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,803 : WARNING : EPOCH - 641 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,811 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,812 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,813 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,837 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,838 : INFO : EPOCH - 642 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,839 : WARNING : EPOCH - 642 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,847 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,847 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,848 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,874 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,874 : INFO : EPOCH - 643 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,875 : WARNING : EPOCH - 643 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,883 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,884 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,885 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,911 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,912 : INFO : EPOCH - 644 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:09,913 : WARNING : EPOCH - 644 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:09,921 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,922 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,923 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,948 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,949 : INFO : EPOCH - 645 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:09,949 : WARNING : EPOCH - 645 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,958 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,959 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,960 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:09,985 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:09,986 : INFO : EPOCH - 646 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:09,986 : WARNING : EPOCH - 646 : supplied example count (1) did not equal expe
2019-02-21 17:17:09,994 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:09,995 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:09,996 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,020 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,021 : INFO : EPOCH - 647 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:10,022 : WARNING : EPOCH - 647 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,030 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,031 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,032 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,056 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,057 : INFO : EPOCH - 648 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:10,057 : WARNING : EPOCH - 648 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,067 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,068 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,068 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,091 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,093 : INFO : EPOCH - 649 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:10,093 : WARNING : EPOCH - 649 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,103 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,103 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,104 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,129 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,130 : INFO : EPOCH - 650 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:10,130 : WARNING : EPOCH - 650 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,139 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,141 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,142 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,167 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,170 : INFO : EPOCH - 651 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:10,171 : WARNING : EPOCH - 651 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,183 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,185 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,186 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,211 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,212 : INFO : EPOCH - 652 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:10,213 : WARNING : EPOCH - 652 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:10,222 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,223 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,223 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,249 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,250 : INFO : EPOCH - 653 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,251 : WARNING : EPOCH - 653 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,258 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,260 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,261 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,285 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,286 : INFO : EPOCH - 654 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,286 : WARNING : EPOCH - 654 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,296 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,296 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,299 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,322 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,323 : INFO : EPOCH - 655 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,324 : WARNING : EPOCH - 655 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,332 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,333 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,334 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,359 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,360 : INFO : EPOCH - 656 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,360 : WARNING : EPOCH - 656 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,369 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,370 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,371 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,396 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,397 : INFO : EPOCH - 657 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,397 : WARNING : EPOCH - 657 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,407 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,408 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,409 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,435 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,436 : INFO : EPOCH - 658 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,436 : WARNING : EPOCH - 658 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,445 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,446 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,446 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,471 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,472 : INFO : EPOCH - 659 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,472 : WARNING : EPOCH - 659 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,480 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,481 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,482 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,508 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,509 : INFO : EPOCH - 660 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,509 : WARNING : EPOCH - 660 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:10,519 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,520 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,520 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,545 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,546 : INFO : EPOCH - 661 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:10,546 : WARNING : EPOCH - 661 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,555 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,556 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,556 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,582 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,583 : INFO : EPOCH - 662 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:10,583 : WARNING : EPOCH - 662 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,593 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,594 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,594 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,620 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,621 : INFO : EPOCH - 663 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:10,621 : WARNING : EPOCH - 663 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,631 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,632 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,633 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,659 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,660 : INFO : EPOCH - 664 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:10,661 : WARNING : EPOCH - 664 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,668 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,671 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,671 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,697 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,698 : INFO : EPOCH - 665 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:10,699 : WARNING : EPOCH - 665 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,707 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,708 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,709 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,734 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,736 : INFO : EPOCH - 666 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:10,736 : WARNING : EPOCH - 666 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,745 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,745 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,746 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,771 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,772 : INFO : EPOCH - 667 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:10,772 : WARNING : EPOCH - 667 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,781 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,781 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,782 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,807 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,808 : INFO : EPOCH - 668 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:10,809 : WARNING : EPOCH - 668 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:10,818 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,819 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,819 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,844 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,845 : INFO : EPOCH - 669 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,845 : WARNING : EPOCH - 669 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,855 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,856 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,857 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,882 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,883 : INFO : EPOCH - 670 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,883 : WARNING : EPOCH - 670 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,894 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,895 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,896 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,922 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,923 : INFO : EPOCH - 671 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,923 : WARNING : EPOCH - 671 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,933 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,937 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,937 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:10,959 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:10,961 : INFO : EPOCH - 672 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:10,962 : WARNING : EPOCH - 672 : supplied example count (1) did not equal expe
2019-02-21 17:17:10,978 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:10,984 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:10,985 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,004 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,005 : INFO : EPOCH - 673 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:11,006 : WARNING : EPOCH - 673 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,014 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,017 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,018 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,041 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,042 : INFO : EPOCH - 674 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:11,043 : WARNING : EPOCH - 674 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,054 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,055 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,055 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,079 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,080 : INFO : EPOCH - 675 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:11,081 : WARNING : EPOCH - 675 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,089 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,090 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,091 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,115 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,116 : INFO : EPOCH - 676 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:11,117 : WARNING : EPOCH - 676 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:11,127 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,128 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,129 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,154 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,155 : INFO : EPOCH - 677 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:11,156 : WARNING : EPOCH - 677 : supplied example count (1) did not equal exp
2019-02-21 17:17:11,164 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,165 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,166 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,191 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,192 : INFO : EPOCH - 678 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:11,193 : WARNING : EPOCH - 678 : supplied example count (1) did not equal exp
2019-02-21 17:17:11,204 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,205 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,206 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,231 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,232 : INFO : EPOCH - 679 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:11,233 : WARNING : EPOCH - 679 : supplied example count (1) did not equal exp
2019-02-21 17:17:11,241 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,242 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,243 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,269 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,270 : INFO : EPOCH - 680 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:11,271 : WARNING : EPOCH - 680 : supplied example count (1) did not equal exp
2019-02-21 17:17:11,280 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,281 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,282 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,307 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,308 : INFO : EPOCH - 681 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:11,309 : WARNING : EPOCH - 681 : supplied example count (1) did not equal exp
2019-02-21 17:17:11,317 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,318 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,319 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,343 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,344 : INFO : EPOCH - 682 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:11,344 : WARNING : EPOCH - 682 : supplied example count (1) did not equal exp
2019-02-21 17:17:11,355 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,355 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,356 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,381 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,382 : INFO : EPOCH - 683 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:11,382 : WARNING : EPOCH - 683 : supplied example count (1) did not equal exp
2019-02-21 17:17:11,392 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,392 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,393 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,418 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,419 : INFO : EPOCH - 684 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:11,420 : WARNING : EPOCH - 684 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:11,428 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,429 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,430 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,454 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,455 : INFO : EPOCH - 685 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:11,456 : WARNING : EPOCH - 685 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,464 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,465 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,466 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,491 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,493 : INFO : EPOCH - 686 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:11,493 : WARNING : EPOCH - 686 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,502 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,503 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,504 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,530 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,530 : INFO : EPOCH - 687 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:11,531 : WARNING : EPOCH - 687 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,540 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,541 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,542 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,567 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,568 : INFO : EPOCH - 688 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:11,569 : WARNING : EPOCH - 688 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,578 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,579 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,580 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,603 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,604 : INFO : EPOCH - 689 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:11,605 : WARNING : EPOCH - 689 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,614 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,615 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,616 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,640 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,641 : INFO : EPOCH - 690 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:11,642 : WARNING : EPOCH - 690 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,651 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,652 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,653 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,677 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,678 : INFO : EPOCH - 691 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:11,679 : WARNING : EPOCH - 691 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,687 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,688 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,689 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,713 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,714 : INFO : EPOCH - 692 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:11,714 : WARNING : EPOCH - 692 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:11,723 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,724 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,725 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,750 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,751 : INFO : EPOCH - 693 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:11,752 : WARNING : EPOCH - 693 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,760 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,761 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,762 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,788 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,789 : INFO : EPOCH - 694 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:11,789 : WARNING : EPOCH - 694 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,798 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,800 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,801 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,825 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,826 : INFO : EPOCH - 695 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:11,827 : WARNING : EPOCH - 695 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,837 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,838 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,838 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,863 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,864 : INFO : EPOCH - 696 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:11,864 : WARNING : EPOCH - 696 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,873 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,874 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,876 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,900 : INFO : EPOCH - 697 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:11,901 : WARNING : EPOCH - 697 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,910 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,911 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,912 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,937 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,938 : INFO : EPOCH - 698 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:11,939 : WARNING : EPOCH - 698 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,947 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,947 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,948 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:11,973 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:11,974 : INFO : EPOCH - 699 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:11,975 : WARNING : EPOCH - 699 : supplied example count (1) did not equal expe
2019-02-21 17:17:11,983 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:11,984 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:11,985 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,010 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,010 : INFO : EPOCH - 700 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,011 : WARNING : EPOCH - 700 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:12,024 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,025 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,026 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,050 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,051 : INFO : EPOCH - 701 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,051 : WARNING : EPOCH - 701 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,060 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,065 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,065 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,088 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,088 : INFO : EPOCH - 702 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,089 : WARNING : EPOCH - 702 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,101 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,102 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,103 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,127 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,128 : INFO : EPOCH - 703 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,129 : WARNING : EPOCH - 703 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,139 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,140 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,141 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,165 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,166 : INFO : EPOCH - 704 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,167 : WARNING : EPOCH - 704 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,178 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,180 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,181 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,206 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,206 : INFO : EPOCH - 705 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,207 : WARNING : EPOCH - 705 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,215 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,216 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,217 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,242 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,243 : INFO : EPOCH - 706 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,243 : WARNING : EPOCH - 706 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,252 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,253 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,254 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,278 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,279 : INFO : EPOCH - 707 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,280 : WARNING : EPOCH - 707 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,289 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,290 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,290 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,315 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,316 : INFO : EPOCH - 708 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:12,317 : WARNING : EPOCH - 708 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:12,326 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,326 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,327 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,353 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,353 : INFO : EPOCH - 709 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,354 : WARNING : EPOCH - 709 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,363 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,364 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,365 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,389 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,390 : INFO : EPOCH - 710 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,391 : WARNING : EPOCH - 710 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,399 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,400 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,401 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,425 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,426 : INFO : EPOCH - 711 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,427 : WARNING : EPOCH - 711 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,437 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,438 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,439 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,463 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,464 : INFO : EPOCH - 712 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,465 : WARNING : EPOCH - 712 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,473 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,474 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,475 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,499 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,500 : INFO : EPOCH - 713 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,501 : WARNING : EPOCH - 713 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,509 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,510 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,511 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,536 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,537 : INFO : EPOCH - 714 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,537 : WARNING : EPOCH - 714 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,545 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,546 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,547 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,572 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,573 : INFO : EPOCH - 715 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,573 : WARNING : EPOCH - 715 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,586 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,587 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,588 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,635 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,640 : INFO : EPOCH - 716 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,645 : WARNING : EPOCH - 716 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:12,654 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,655 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,656 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,682 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,683 : INFO : EPOCH - 717 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:12,683 : WARNING : EPOCH - 717 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,693 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,694 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,695 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,719 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,720 : INFO : EPOCH - 718 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:12,721 : WARNING : EPOCH - 718 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,731 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,732 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,732 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,757 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,758 : INFO : EPOCH - 719 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:12,759 : WARNING : EPOCH - 719 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,768 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,769 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,769 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,795 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,795 : INFO : EPOCH - 720 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:12,796 : WARNING : EPOCH - 720 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,806 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,808 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,838 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,838 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,839 : INFO : EPOCH - 721 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:12,840 : WARNING : EPOCH - 721 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,851 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,852 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,853 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,877 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,878 : INFO : EPOCH - 722 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:12,879 : WARNING : EPOCH - 722 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,889 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,891 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,894 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,916 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,917 : INFO : EPOCH - 723 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:12,918 : WARNING : EPOCH - 723 : supplied example count (1) did not equal expe
2019-02-21 17:17:12,927 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,936 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,936 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,954 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,955 : INFO : EPOCH - 724 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:12,955 : WARNING : EPOCH - 724 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:12,964 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:12,965 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:12,966 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:12,991 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:12,992 : INFO : EPOCH - 725 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:12,993 : WARNING : EPOCH - 725 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,001 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,003 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,003 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,027 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,029 : INFO : EPOCH - 726 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:13,029 : WARNING : EPOCH - 726 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,037 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,039 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,039 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,064 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,065 : INFO : EPOCH - 727 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:13,066 : WARNING : EPOCH - 727 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,074 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,076 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,076 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,100 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,101 : INFO : EPOCH - 728 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:13,101 : WARNING : EPOCH - 728 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,110 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,111 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,111 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,136 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,136 : INFO : EPOCH - 729 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:13,137 : WARNING : EPOCH - 729 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,146 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,147 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,147 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,173 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,173 : INFO : EPOCH - 730 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:13,175 : WARNING : EPOCH - 730 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,185 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,186 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,187 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,213 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,214 : INFO : EPOCH - 731 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:13,215 : WARNING : EPOCH - 731 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,225 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,226 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,226 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,251 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,252 : INFO : EPOCH - 732 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:13,252 : WARNING : EPOCH - 732 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:13,261 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,262 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,263 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,287 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,288 : INFO : EPOCH - 733 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,289 : WARNING : EPOCH - 733 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,297 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,298 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,299 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,323 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,324 : INFO : EPOCH - 734 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,325 : WARNING : EPOCH - 734 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,334 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,335 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,335 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,359 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,359 : INFO : EPOCH - 735 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,360 : WARNING : EPOCH - 735 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,368 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,369 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,370 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,395 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,396 : INFO : EPOCH - 736 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,396 : WARNING : EPOCH - 736 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,405 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,406 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,407 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,431 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,432 : INFO : EPOCH - 737 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,433 : WARNING : EPOCH - 737 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,445 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,446 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,447 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,472 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,479 : INFO : EPOCH - 738 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,480 : WARNING : EPOCH - 738 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,490 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,491 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,492 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,517 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,518 : INFO : EPOCH - 739 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,518 : WARNING : EPOCH - 739 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,526 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,528 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,528 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,554 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,555 : INFO : EPOCH - 740 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,556 : WARNING : EPOCH - 740 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:13,565 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,566 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,567 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,590 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,591 : INFO : EPOCH - 741 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,592 : WARNING : EPOCH - 741 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,601 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,602 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,602 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,626 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,627 : INFO : EPOCH - 742 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,628 : WARNING : EPOCH - 742 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,637 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,638 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,639 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,663 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,664 : INFO : EPOCH - 743 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,665 : WARNING : EPOCH - 743 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,674 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,675 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,675 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,699 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,700 : INFO : EPOCH - 744 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,701 : WARNING : EPOCH - 744 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,712 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,715 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,716 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,739 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,740 : INFO : EPOCH - 745 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,740 : WARNING : EPOCH - 745 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,749 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,771 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,772 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,776 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,777 : INFO : EPOCH - 746 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,778 : WARNING : EPOCH - 746 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,788 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,788 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,789 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,816 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,818 : INFO : EPOCH - 747 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,818 : WARNING : EPOCH - 747 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,827 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,828 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,828 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,853 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,854 : INFO : EPOCH - 748 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:13,854 : WARNING : EPOCH - 748 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:13,864 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,865 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,865 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,889 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,890 : INFO : EPOCH - 749 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:13,891 : WARNING : EPOCH - 749 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,901 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,902 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,902 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,927 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,928 : INFO : EPOCH - 750 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:13,932 : WARNING : EPOCH - 750 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,945 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,946 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,947 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:13,973 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:13,974 : INFO : EPOCH - 751 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:13,974 : WARNING : EPOCH - 751 : supplied example count (1) did not equal expe
2019-02-21 17:17:13,989 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:13,991 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:13,992 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,017 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,018 : INFO : EPOCH - 752 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:14,018 : WARNING : EPOCH - 752 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,027 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,028 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,028 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,055 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,056 : INFO : EPOCH - 753 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:14,057 : WARNING : EPOCH - 753 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,065 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,067 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,069 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,093 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,094 : INFO : EPOCH - 754 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:14,095 : WARNING : EPOCH - 754 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,104 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,105 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,106 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,131 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,133 : INFO : EPOCH - 755 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:14,133 : WARNING : EPOCH - 755 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,142 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,143 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,144 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,170 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,171 : INFO : EPOCH - 756 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:14,172 : WARNING : EPOCH - 756 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:14,184 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,186 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,188 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,211 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,212 : INFO : EPOCH - 757 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:14,212 : WARNING : EPOCH - 757 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,223 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,223 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,224 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,249 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,250 : INFO : EPOCH - 758 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:14,251 : WARNING : EPOCH - 758 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,260 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,261 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,261 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,287 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,289 : INFO : EPOCH - 759 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:14,290 : WARNING : EPOCH - 759 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,299 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,301 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,302 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,330 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,331 : INFO : EPOCH - 760 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:14,331 : WARNING : EPOCH - 760 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,345 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,346 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,346 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,373 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,375 : INFO : EPOCH - 761 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:14,375 : WARNING : EPOCH - 761 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,398 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,401 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,402 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,427 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,428 : INFO : EPOCH - 762 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:14,429 : WARNING : EPOCH - 762 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,442 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,443 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,444 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,471 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,472 : INFO : EPOCH - 763 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:14,473 : WARNING : EPOCH - 763 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,489 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,490 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,491 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,516 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,517 : INFO : EPOCH - 764 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:14,517 : WARNING : EPOCH - 764 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:14,526 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,528 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,528 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,554 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,557 : INFO : EPOCH - 765 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:14,558 : WARNING : EPOCH - 765 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,571 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,572 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,573 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,600 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,601 : INFO : EPOCH - 766 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:14,602 : WARNING : EPOCH - 766 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,611 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,613 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,613 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,638 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,639 : INFO : EPOCH - 767 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:14,640 : WARNING : EPOCH - 767 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,650 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,652 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,652 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,680 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,681 : INFO : EPOCH - 768 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:14,682 : WARNING : EPOCH - 768 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,692 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,693 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,694 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,719 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,720 : INFO : EPOCH - 769 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:14,721 : WARNING : EPOCH - 769 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,729 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,731 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,731 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,760 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,760 : INFO : EPOCH - 770 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:14,761 : WARNING : EPOCH - 770 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,776 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,777 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,779 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,806 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,807 : INFO : EPOCH - 771 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:14,809 : WARNING : EPOCH - 771 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,842 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,844 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,847 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,875 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,876 : INFO : EPOCH - 772 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:14,878 : WARNING : EPOCH - 772 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:14,890 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,908 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,910 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,922 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,924 : INFO : EPOCH - 773 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:14,924 : WARNING : EPOCH - 773 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,941 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,943 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,945 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:14,971 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:14,972 : INFO : EPOCH - 774 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:14,975 : WARNING : EPOCH - 774 : supplied example count (1) did not equal expe
2019-02-21 17:17:14,988 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:14,989 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:14,991 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,014 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,015 : INFO : EPOCH - 775 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:15,016 : WARNING : EPOCH - 775 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,026 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,028 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,028 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,054 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,055 : INFO : EPOCH - 776 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:15,056 : WARNING : EPOCH - 776 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,071 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,072 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,073 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,097 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,099 : INFO : EPOCH - 777 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:15,099 : WARNING : EPOCH - 777 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,113 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,115 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,116 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,140 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,141 : INFO : EPOCH - 778 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:15,142 : WARNING : EPOCH - 778 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,158 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,188 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,189 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,189 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,191 : INFO : EPOCH - 779 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:15,192 : WARNING : EPOCH - 779 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,221 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,224 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,225 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,248 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,248 : INFO : EPOCH - 780 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:15,249 : WARNING : EPOCH - 780 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:15,262 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,263 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,265 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,288 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,289 : INFO : EPOCH - 781 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,289 : WARNING : EPOCH - 781 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,304 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,305 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,306 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,330 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,331 : INFO : EPOCH - 782 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,332 : WARNING : EPOCH - 782 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,347 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,355 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,356 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,373 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,374 : INFO : EPOCH - 783 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,375 : WARNING : EPOCH - 783 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,387 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,389 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,389 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,416 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,417 : INFO : EPOCH - 784 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,418 : WARNING : EPOCH - 784 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,430 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,437 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,438 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,457 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,458 : INFO : EPOCH - 785 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,459 : WARNING : EPOCH - 785 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,484 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,487 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,489 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,511 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,512 : INFO : EPOCH - 786 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,513 : WARNING : EPOCH - 786 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,525 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,526 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,527 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,553 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,555 : INFO : EPOCH - 787 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,555 : WARNING : EPOCH - 787 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,572 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,574 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,574 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,599 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,601 : INFO : EPOCH - 788 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,602 : WARNING : EPOCH - 788 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:15,613 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,616 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,617 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,641 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,642 : INFO : EPOCH - 789 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,642 : WARNING : EPOCH - 789 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,657 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,658 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,658 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,683 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,684 : INFO : EPOCH - 790 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,685 : WARNING : EPOCH - 790 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,696 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,698 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,699 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,724 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,726 : INFO : EPOCH - 791 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,727 : WARNING : EPOCH - 791 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,740 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,741 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,741 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,767 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,780 : INFO : EPOCH - 792 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,781 : WARNING : EPOCH - 792 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,792 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,793 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,794 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,819 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,821 : INFO : EPOCH - 793 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,822 : WARNING : EPOCH - 793 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,835 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,836 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,836 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,862 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,863 : INFO : EPOCH - 794 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,864 : WARNING : EPOCH - 794 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,880 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,896 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,896 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,906 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,908 : INFO : EPOCH - 795 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,908 : WARNING : EPOCH - 795 : supplied example count (1) did not equal expe
2019-02-21 17:17:15,924 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,927 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,928 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,951 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,952 : INFO : EPOCH - 796 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,953 : WARNING : EPOCH - 796 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:15,966 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:15,980 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:15,981 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:15,992 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:15,993 : INFO : EPOCH - 797 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:15,993 : WARNING : EPOCH - 797 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,007 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,008 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,010 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,035 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,039 : INFO : EPOCH - 798 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:16,039 : WARNING : EPOCH - 798 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,055 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,056 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,057 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,084 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,090 : INFO : EPOCH - 799 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:16,090 : WARNING : EPOCH - 799 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,103 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,104 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,104 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,130 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,131 : INFO : EPOCH - 800 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:16,132 : WARNING : EPOCH - 800 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,144 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,145 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,145 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,172 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,174 : INFO : EPOCH - 801 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:16,174 : WARNING : EPOCH - 801 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,196 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,206 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,206 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,223 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,224 : INFO : EPOCH - 802 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:16,225 : WARNING : EPOCH - 802 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,238 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,240 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,242 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,268 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,269 : INFO : EPOCH - 803 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:16,270 : WARNING : EPOCH - 803 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,283 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,284 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,284 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,310 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,311 : INFO : EPOCH - 804 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:16,312 : WARNING : EPOCH - 804 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:16,326 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,334 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,335 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,355 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,355 : INFO : EPOCH - 805 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,356 : WARNING : EPOCH - 805 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,382 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,384 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,385 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,409 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,411 : INFO : EPOCH - 806 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,412 : WARNING : EPOCH - 806 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,425 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,426 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,427 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,452 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,453 : INFO : EPOCH - 807 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,453 : WARNING : EPOCH - 807 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,463 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,464 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,465 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,489 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,490 : INFO : EPOCH - 808 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,491 : WARNING : EPOCH - 808 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,501 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,502 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,503 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,527 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,528 : INFO : EPOCH - 809 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,528 : WARNING : EPOCH - 809 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,539 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,539 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,540 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,565 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,566 : INFO : EPOCH - 810 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,566 : WARNING : EPOCH - 810 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,576 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,577 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,578 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,603 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,604 : INFO : EPOCH - 811 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,604 : WARNING : EPOCH - 811 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,613 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,615 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,617 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,640 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,641 : INFO : EPOCH - 812 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,641 : WARNING : EPOCH - 812 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:16,651 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,659 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,660 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,678 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,679 : INFO : EPOCH - 813 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,679 : WARNING : EPOCH - 813 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,689 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,690 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,690 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,715 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,716 : INFO : EPOCH - 814 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,717 : WARNING : EPOCH - 814 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,724 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,725 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,725 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,750 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,751 : INFO : EPOCH - 815 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,751 : WARNING : EPOCH - 815 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,759 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,760 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,761 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,787 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,788 : INFO : EPOCH - 816 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,788 : WARNING : EPOCH - 816 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,797 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,798 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,800 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,823 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,824 : INFO : EPOCH - 817 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,825 : WARNING : EPOCH - 817 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,833 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,835 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,835 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,860 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,861 : INFO : EPOCH - 818 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,862 : WARNING : EPOCH - 818 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,872 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,873 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,874 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,901 : INFO : EPOCH - 819 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,901 : WARNING : EPOCH - 819 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,910 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,911 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,911 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,937 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,938 : INFO : EPOCH - 820 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:16,938 : WARNING : EPOCH - 820 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:16,946 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,947 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,948 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:16,976 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:16,977 : INFO : EPOCH - 821 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:16,977 : WARNING : EPOCH - 821 : supplied example count (1) did not equal expe
2019-02-21 17:17:16,987 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:16,988 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:16,989 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,014 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,015 : INFO : EPOCH - 822 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:17,015 : WARNING : EPOCH - 822 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,024 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,025 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,025 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,049 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,051 : INFO : EPOCH - 823 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:17,051 : WARNING : EPOCH - 823 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,090 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,091 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,091 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,116 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,117 : INFO : EPOCH - 824 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:17,117 : WARNING : EPOCH - 824 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,130 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,130 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,132 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,159 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,160 : INFO : EPOCH - 825 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:17,161 : WARNING : EPOCH - 825 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,183 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,195 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,213 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,214 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,215 : INFO : EPOCH - 826 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:17,216 : WARNING : EPOCH - 826 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,230 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,233 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,234 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,257 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,258 : INFO : EPOCH - 827 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:17,259 : WARNING : EPOCH - 827 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,273 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,274 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,274 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,299 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,301 : INFO : EPOCH - 828 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:17,301 : WARNING : EPOCH - 828 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:17,313 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,314 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,315 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,339 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,340 : INFO : EPOCH - 829 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:17,341 : WARNING : EPOCH - 829 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,353 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,354 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,355 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,379 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,381 : INFO : EPOCH - 830 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:17,381 : WARNING : EPOCH - 830 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,395 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,396 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,397 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,418 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,419 : INFO : EPOCH - 831 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:17,419 : WARNING : EPOCH - 831 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,432 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,434 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,434 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,460 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,461 : INFO : EPOCH - 832 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:17,462 : WARNING : EPOCH - 832 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,476 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,477 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,477 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,503 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,508 : INFO : EPOCH - 833 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:17,509 : WARNING : EPOCH - 833 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,521 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,522 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,522 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,547 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,549 : INFO : EPOCH - 834 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:17,550 : WARNING : EPOCH - 834 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,561 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,562 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,563 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,589 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,590 : INFO : EPOCH - 835 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:17,590 : WARNING : EPOCH - 835 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,605 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,614 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,614 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,632 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,633 : INFO : EPOCH - 836 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:17,634 : WARNING : EPOCH - 836 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:17,645 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,646 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,647 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,672 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,673 : INFO : EPOCH - 837 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:17,673 : WARNING : EPOCH - 837 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,687 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,689 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,691 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,713 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,730 : INFO : EPOCH - 838 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:17,731 : WARNING : EPOCH - 838 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,742 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,743 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,744 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,770 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,773 : INFO : EPOCH - 839 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:17,774 : WARNING : EPOCH - 839 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,789 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,823 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,824 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,824 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,825 : INFO : EPOCH - 840 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:17,827 : WARNING : EPOCH - 840 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,841 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,841 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,842 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,867 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,869 : INFO : EPOCH - 841 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:17,869 : WARNING : EPOCH - 841 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,882 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,883 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,885 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,908 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,909 : INFO : EPOCH - 842 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:17,910 : WARNING : EPOCH - 842 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,921 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,922 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,922 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,948 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,949 : INFO : EPOCH - 843 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:17,951 : WARNING : EPOCH - 843 : supplied example count (1) did not equal expe
2019-02-21 17:17:17,963 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:17,964 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:17,965 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:17,990 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:17,991 : INFO : EPOCH - 844 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:17,992 : WARNING : EPOCH - 844 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:18,005 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,005 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,006 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,030 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,031 : INFO : EPOCH - 845 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,032 : WARNING : EPOCH - 845 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,042 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,046 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,047 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,069 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,071 : INFO : EPOCH - 846 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,071 : WARNING : EPOCH - 846 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,083 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,084 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,085 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,111 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,112 : INFO : EPOCH - 847 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,112 : WARNING : EPOCH - 847 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,122 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,122 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,123 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,150 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,150 : INFO : EPOCH - 848 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,150 : WARNING : EPOCH - 848 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,159 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,160 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,160 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,185 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,186 : INFO : EPOCH - 849 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,187 : WARNING : EPOCH - 849 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,197 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,199 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,199 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,225 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,226 : INFO : EPOCH - 850 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,226 : WARNING : EPOCH - 850 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,235 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,236 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,237 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,260 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,261 : INFO : EPOCH - 851 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,262 : WARNING : EPOCH - 851 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,271 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,273 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,273 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,297 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,298 : INFO : EPOCH - 852 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,299 : WARNING : EPOCH - 852 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:18,309 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,310 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,311 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,336 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,337 : INFO : EPOCH - 853 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,338 : WARNING : EPOCH - 853 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,345 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,347 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,347 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,372 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,372 : INFO : EPOCH - 854 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,373 : WARNING : EPOCH - 854 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,384 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,385 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,387 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,411 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,412 : INFO : EPOCH - 855 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,413 : WARNING : EPOCH - 855 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,423 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,425 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,425 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,450 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,451 : INFO : EPOCH - 856 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,451 : WARNING : EPOCH - 856 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,460 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,462 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,463 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,487 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,488 : INFO : EPOCH - 857 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,488 : WARNING : EPOCH - 857 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,498 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,498 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,499 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,524 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,525 : INFO : EPOCH - 858 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,525 : WARNING : EPOCH - 858 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,534 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,535 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,536 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,561 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,562 : INFO : EPOCH - 859 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,563 : WARNING : EPOCH - 859 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,572 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,572 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,573 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,598 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,599 : INFO : EPOCH - 860 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:18,600 : WARNING : EPOCH - 860 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:18,608 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,609 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,610 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,634 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,635 : INFO : EPOCH - 861 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:18,635 : WARNING : EPOCH - 861 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,644 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,646 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,647 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,673 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,674 : INFO : EPOCH - 862 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:18,675 : WARNING : EPOCH - 862 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,684 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,685 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,685 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,710 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,711 : INFO : EPOCH - 863 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:18,711 : WARNING : EPOCH - 863 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,722 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,723 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,723 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,751 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,753 : INFO : EPOCH - 864 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:18,754 : WARNING : EPOCH - 864 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,763 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,763 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,764 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,790 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,791 : INFO : EPOCH - 865 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:18,791 : WARNING : EPOCH - 865 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,800 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,802 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,802 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,827 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,828 : INFO : EPOCH - 866 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:18,829 : WARNING : EPOCH - 866 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,838 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,839 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,840 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,864 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,865 : INFO : EPOCH - 867 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:18,865 : WARNING : EPOCH - 867 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,873 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,875 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,875 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,901 : INFO : EPOCH - 868 : training on 5386246 raw words (10000 effective v
2019-02-21 17:17:18,901 : WARNING : EPOCH - 868 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:18,911 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,911 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,912 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,937 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,939 : INFO : EPOCH - 869 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:18,939 : WARNING : EPOCH - 869 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,947 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,948 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,949 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:18,975 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:18,978 : INFO : EPOCH - 870 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:18,979 : WARNING : EPOCH - 870 : supplied example count (1) did not equal expe
2019-02-21 17:17:18,987 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:18,989 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:18,990 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,015 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,016 : INFO : EPOCH - 871 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,017 : WARNING : EPOCH - 871 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,026 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,027 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,028 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,053 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,054 : INFO : EPOCH - 872 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,055 : WARNING : EPOCH - 872 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,063 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,064 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,064 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,088 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,089 : INFO : EPOCH - 873 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,090 : WARNING : EPOCH - 873 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,097 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,098 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,099 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,123 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,124 : INFO : EPOCH - 874 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,125 : WARNING : EPOCH - 874 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,134 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,135 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,135 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,160 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,161 : INFO : EPOCH - 875 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,161 : WARNING : EPOCH - 875 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,171 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,172 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,173 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,199 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,200 : INFO : EPOCH - 876 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,200 : WARNING : EPOCH - 876 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:19,209 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,210 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,211 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,237 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,237 : INFO : EPOCH - 877 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,238 : WARNING : EPOCH - 877 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,254 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,255 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,256 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,280 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,281 : INFO : EPOCH - 878 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,282 : WARNING : EPOCH - 878 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,290 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,292 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,292 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,317 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,318 : INFO : EPOCH - 879 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,319 : WARNING : EPOCH - 879 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,327 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,328 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,328 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,353 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,354 : INFO : EPOCH - 880 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,355 : WARNING : EPOCH - 880 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,363 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,364 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,365 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,389 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,390 : INFO : EPOCH - 881 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,391 : WARNING : EPOCH - 881 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,400 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,401 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,401 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,426 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,427 : INFO : EPOCH - 882 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,427 : WARNING : EPOCH - 882 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,436 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,438 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,438 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,463 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,464 : INFO : EPOCH - 883 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,465 : WARNING : EPOCH - 883 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,474 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,475 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,476 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,499 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,500 : INFO : EPOCH - 884 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,501 : WARNING : EPOCH - 884 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:19,512 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,513 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,514 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,539 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,540 : INFO : EPOCH - 885 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,541 : WARNING : EPOCH - 885 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,556 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,557 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,558 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,583 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,584 : INFO : EPOCH - 886 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,585 : WARNING : EPOCH - 886 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,595 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,596 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,597 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,621 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,622 : INFO : EPOCH - 887 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,623 : WARNING : EPOCH - 887 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,632 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,633 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,634 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,658 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,659 : INFO : EPOCH - 888 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,659 : WARNING : EPOCH - 888 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,668 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,669 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,669 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,693 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,694 : INFO : EPOCH - 889 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,695 : WARNING : EPOCH - 889 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,705 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,706 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,706 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,730 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,731 : INFO : EPOCH - 890 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,731 : WARNING : EPOCH - 890 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,741 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,742 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,742 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,766 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,767 : INFO : EPOCH - 891 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,768 : WARNING : EPOCH - 891 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,776 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,777 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,778 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,802 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,803 : INFO : EPOCH - 892 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,804 : WARNING : EPOCH - 892 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:19,813 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,815 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,816 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,842 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,843 : INFO : EPOCH - 893 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,843 : WARNING : EPOCH - 893 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,853 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,853 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,854 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,878 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,879 : INFO : EPOCH - 894 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,880 : WARNING : EPOCH - 894 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,890 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,891 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,891 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,917 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,917 : INFO : EPOCH - 895 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,918 : WARNING : EPOCH - 895 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,926 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,927 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,928 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,953 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,954 : INFO : EPOCH - 896 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,955 : WARNING : EPOCH - 896 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,963 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,964 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:19,965 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:19,989 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:19,990 : INFO : EPOCH - 897 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:19,990 : WARNING : EPOCH - 897 : supplied example count (1) did not equal expe
2019-02-21 17:17:19,998 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:19,999 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,001 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,024 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,025 : INFO : EPOCH - 898 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:20,026 : WARNING : EPOCH - 898 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,035 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,036 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,036 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,060 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,061 : INFO : EPOCH - 899 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:20,062 : WARNING : EPOCH - 899 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,071 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,073 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,073 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,097 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,098 : INFO : EPOCH - 900 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:20,098 : WARNING : EPOCH - 900 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:20,107 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,108 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,118 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,134 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,135 : INFO : EPOCH - 901 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:20,135 : WARNING : EPOCH - 901 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,143 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,144 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,145 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,169 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,170 : INFO : EPOCH - 902 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:20,171 : WARNING : EPOCH - 902 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,183 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,184 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,185 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,209 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,210 : INFO : EPOCH - 903 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:20,211 : WARNING : EPOCH - 903 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,219 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,220 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,221 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,245 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,246 : INFO : EPOCH - 904 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:20,247 : WARNING : EPOCH - 904 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,256 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,257 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,257 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,281 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,282 : INFO : EPOCH - 905 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:20,283 : WARNING : EPOCH - 905 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,291 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,292 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,292 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,317 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,318 : INFO : EPOCH - 906 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:20,319 : WARNING : EPOCH - 906 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,327 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,327 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,328 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,352 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,353 : INFO : EPOCH - 907 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:20,354 : WARNING : EPOCH - 907 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,362 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,363 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,364 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,388 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,389 : INFO : EPOCH - 908 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:20,390 : WARNING : EPOCH - 908 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:20,398 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,399 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,400 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,425 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,425 : INFO : EPOCH - 909 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:20,426 : WARNING : EPOCH - 909 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,435 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,436 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,436 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,461 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,462 : INFO : EPOCH - 910 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:20,462 : WARNING : EPOCH - 910 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,472 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,473 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,498 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,499 : INFO : EPOCH - 911 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:20,499 : WARNING : EPOCH - 911 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,508 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,509 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,510 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,538 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,539 : INFO : EPOCH - 912 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:20,539 : WARNING : EPOCH - 912 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,547 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,548 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,549 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,574 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,575 : INFO : EPOCH - 913 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:20,575 : WARNING : EPOCH - 913 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,586 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,587 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,588 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,612 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,613 : INFO : EPOCH - 914 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:20,614 : WARNING : EPOCH - 914 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,622 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,623 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,623 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,648 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,649 : INFO : EPOCH - 915 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:20,650 : WARNING : EPOCH - 915 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,658 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,659 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,660 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,685 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,686 : INFO : EPOCH - 916 : training on 5386246 raw words (10000 effective 
2019-02-21 17:17:20,687 : WARNING : EPOCH - 916 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:20,697 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,699 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,700 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,723 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,724 : INFO : EPOCH - 917 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:20,725 : WARNING : EPOCH - 917 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,739 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,742 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,743 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,765 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,767 : INFO : EPOCH - 918 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:20,767 : WARNING : EPOCH - 918 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,777 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,778 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,779 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,804 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,805 : INFO : EPOCH - 919 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:20,805 : WARNING : EPOCH - 919 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,813 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,814 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,815 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,838 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,839 : INFO : EPOCH - 920 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:20,840 : WARNING : EPOCH - 920 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,848 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,849 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,850 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,875 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,876 : INFO : EPOCH - 921 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:20,877 : WARNING : EPOCH - 921 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,885 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,887 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,887 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,912 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,913 : INFO : EPOCH - 922 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:20,913 : WARNING : EPOCH - 922 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,922 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,923 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,923 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,947 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,948 : INFO : EPOCH - 923 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:20,949 : WARNING : EPOCH - 923 : supplied example count (1) did not equal expe
2019-02-21 17:17:20,958 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,959 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,960 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:20,984 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:20,985 : INFO : EPOCH - 924 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:20,985 : WARNING : EPOCH - 924 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:20,995 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:20,996 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:20,997 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,022 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,031 : INFO : EPOCH - 925 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,032 : WARNING : EPOCH - 925 : supplied example count (1) did not equal exp
2019-02-21 17:17:21,040 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,041 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,042 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,068 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,069 : INFO : EPOCH - 926 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,069 : WARNING : EPOCH - 926 : supplied example count (1) did not equal exp
2019-02-21 17:17:21,078 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,079 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,080 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,104 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,105 : INFO : EPOCH - 927 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,105 : WARNING : EPOCH - 927 : supplied example count (1) did not equal exp
2019-02-21 17:17:21,113 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,115 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,116 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,140 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,141 : INFO : EPOCH - 928 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,142 : WARNING : EPOCH - 928 : supplied example count (1) did not equal exp
2019-02-21 17:17:21,150 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,152 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,153 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,178 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,179 : INFO : EPOCH - 929 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,180 : WARNING : EPOCH - 929 : supplied example count (1) did not equal exp
2019-02-21 17:17:21,193 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,194 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,195 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,219 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,227 : INFO : EPOCH - 930 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,228 : WARNING : EPOCH - 930 : supplied example count (1) did not equal exp
2019-02-21 17:17:21,239 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,240 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,242 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,264 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,265 : INFO : EPOCH - 931 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,266 : WARNING : EPOCH - 931 : supplied example count (1) did not equal exp
2019-02-21 17:17:21,277 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,278 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,278 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,303 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,304 : INFO : EPOCH - 932 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,305 : WARNING : EPOCH - 932 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:21,315 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,322 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,323 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,341 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,361 : INFO : EPOCH - 933 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,361 : WARNING : EPOCH - 933 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,378 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,379 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,379 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,404 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,405 : INFO : EPOCH - 934 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,405 : WARNING : EPOCH - 934 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,414 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,415 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,416 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,440 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,441 : INFO : EPOCH - 935 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,441 : WARNING : EPOCH - 935 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,450 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,451 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,452 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,476 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,476 : INFO : EPOCH - 936 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,477 : WARNING : EPOCH - 936 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,487 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,488 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,488 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,512 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,513 : INFO : EPOCH - 937 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,513 : WARNING : EPOCH - 937 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,522 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,523 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,524 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,548 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,549 : INFO : EPOCH - 938 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,550 : WARNING : EPOCH - 938 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,558 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,559 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,560 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,584 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,585 : INFO : EPOCH - 939 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,586 : WARNING : EPOCH - 939 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,594 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,595 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,596 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,621 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,622 : INFO : EPOCH - 940 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,622 : WARNING : EPOCH - 940 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:21,630 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,632 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,633 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,668 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,669 : INFO : EPOCH - 941 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,669 : WARNING : EPOCH - 941 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,678 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,679 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,679 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,722 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,792 : INFO : EPOCH - 942 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,793 : WARNING : EPOCH - 942 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,812 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,836 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,839 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,840 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,840 : INFO : EPOCH - 943 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,841 : WARNING : EPOCH - 943 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,878 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,899 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,915 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,915 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,916 : INFO : EPOCH - 944 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,917 : WARNING : EPOCH - 944 : supplied example count (1) did not equal expe
2019-02-21 17:17:21,944 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:21,955 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:21,964 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:21,969 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:21,980 : INFO : EPOCH - 945 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:21,986 : WARNING : EPOCH - 945 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,007 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,013 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,034 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,035 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,036 : INFO : EPOCH - 946 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:22,036 : WARNING : EPOCH - 946 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,088 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,093 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,094 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,114 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,134 : INFO : EPOCH - 947 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:22,134 : WARNING : EPOCH - 947 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,168 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,196 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,196 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,197 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,198 : INFO : EPOCH - 948 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:22,198 : WARNING : EPOCH - 948 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:22,239 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,251 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,252 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,265 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,267 : INFO : EPOCH - 949 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:22,268 : WARNING : EPOCH - 949 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,305 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,314 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,315 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,333 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,339 : INFO : EPOCH - 950 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:22,339 : WARNING : EPOCH - 950 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,363 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,369 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,369 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,390 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,398 : INFO : EPOCH - 951 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:22,399 : WARNING : EPOCH - 951 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,432 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,459 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,460 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,460 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,461 : INFO : EPOCH - 952 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:22,461 : WARNING : EPOCH - 952 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,495 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,527 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,528 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,529 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,529 : INFO : EPOCH - 953 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:22,530 : WARNING : EPOCH - 953 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,560 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,574 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,575 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,585 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,587 : INFO : EPOCH - 954 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:22,587 : WARNING : EPOCH - 954 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,608 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,634 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,635 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,636 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,636 : INFO : EPOCH - 955 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:22,637 : WARNING : EPOCH - 955 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,786 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,813 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,814 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,814 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,815 : INFO : EPOCH - 956 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:22,815 : WARNING : EPOCH - 956 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:22,871 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,917 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,917 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,918 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,919 : INFO : EPOCH - 957 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:22,919 : WARNING : EPOCH - 957 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,944 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,960 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:22,961 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:22,969 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:22,973 : INFO : EPOCH - 958 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:22,974 : WARNING : EPOCH - 958 : supplied example count (1) did not equal expe
2019-02-21 17:17:22,997 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:22,999 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,001 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,023 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,025 : INFO : EPOCH - 959 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,025 : WARNING : EPOCH - 959 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,051 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,060 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,060 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,078 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,087 : INFO : EPOCH - 960 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,088 : WARNING : EPOCH - 960 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,108 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,113 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,113 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,134 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,136 : INFO : EPOCH - 961 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,137 : WARNING : EPOCH - 961 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,158 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,184 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,185 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,185 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,186 : INFO : EPOCH - 962 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,186 : WARNING : EPOCH - 962 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,224 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,225 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,225 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,251 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,252 : INFO : EPOCH - 963 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,253 : WARNING : EPOCH - 963 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,320 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,322 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,322 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,347 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,348 : INFO : EPOCH - 964 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,349 : WARNING : EPOCH - 964 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:23,358 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,358 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,359 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,385 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,386 : INFO : EPOCH - 965 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,386 : WARNING : EPOCH - 965 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,395 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,396 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,397 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,422 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,429 : INFO : EPOCH - 966 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,430 : WARNING : EPOCH - 966 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,439 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,440 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,441 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,466 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,467 : INFO : EPOCH - 967 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,468 : WARNING : EPOCH - 967 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,476 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,477 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,478 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,503 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,504 : INFO : EPOCH - 968 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,505 : WARNING : EPOCH - 968 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,518 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,519 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,520 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,546 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,547 : INFO : EPOCH - 969 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,548 : WARNING : EPOCH - 969 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,565 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,585 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,586 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,592 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,593 : INFO : EPOCH - 970 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,593 : WARNING : EPOCH - 970 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,606 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,611 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,612 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,633 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,634 : INFO : EPOCH - 971 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,635 : WARNING : EPOCH - 971 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,651 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,653 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,654 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,678 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,679 : INFO : EPOCH - 972 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,680 : WARNING : EPOCH - 972 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:23,693 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,695 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,696 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,720 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,723 : INFO : EPOCH - 973 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,724 : WARNING : EPOCH - 973 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,739 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,742 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,743 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,766 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,767 : INFO : EPOCH - 974 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,768 : WARNING : EPOCH - 974 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,781 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,782 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,783 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,807 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,819 : INFO : EPOCH - 975 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,820 : WARNING : EPOCH - 975 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,837 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,840 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,841 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,864 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,865 : INFO : EPOCH - 976 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,867 : WARNING : EPOCH - 976 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,878 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,879 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,879 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,904 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,905 : INFO : EPOCH - 977 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,906 : WARNING : EPOCH - 977 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,915 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,916 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,917 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:23,942 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:23,961 : INFO : EPOCH - 978 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:23,962 : WARNING : EPOCH - 978 : supplied example count (1) did not equal expe
2019-02-21 17:17:23,978 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:23,979 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:23,979 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,004 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,005 : INFO : EPOCH - 979 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:24,006 : WARNING : EPOCH - 979 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,014 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,015 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,016 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,041 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,042 : INFO : EPOCH - 980 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:24,042 : WARNING : EPOCH - 980 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:24,060 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,066 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,066 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,092 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,097 : INFO : EPOCH - 981 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:24,103 : WARNING : EPOCH - 981 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,123 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,138 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,139 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,150 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,151 : INFO : EPOCH - 982 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:24,152 : WARNING : EPOCH - 982 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,169 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,173 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,174 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,196 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,203 : INFO : EPOCH - 983 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:24,209 : WARNING : EPOCH - 983 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,224 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,226 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,227 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,251 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,253 : INFO : EPOCH - 984 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:24,253 : WARNING : EPOCH - 984 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,264 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,266 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,266 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,292 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,294 : INFO : EPOCH - 985 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:24,294 : WARNING : EPOCH - 985 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,307 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,308 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,309 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,335 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,336 : INFO : EPOCH - 986 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:24,337 : WARNING : EPOCH - 986 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,348 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,349 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,350 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,375 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,376 : INFO : EPOCH - 987 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:24,377 : WARNING : EPOCH - 987 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,390 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,392 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,393 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,417 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,420 : INFO : EPOCH - 988 : training on 5386246 raw words (10000 effective w
2019-02-21 17:17:24,420 : WARNING : EPOCH - 988 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:24,432 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,433 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,434 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,460 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,462 : INFO : EPOCH - 989 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:24,462 : WARNING : EPOCH - 989 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,474 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,476 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,477 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,502 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,503 : INFO : EPOCH - 990 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:24,504 : WARNING : EPOCH - 990 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,517 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,524 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,525 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,544 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,545 : INFO : EPOCH - 991 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:24,546 : WARNING : EPOCH - 991 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,565 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,566 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,567 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,592 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,593 : INFO : EPOCH - 992 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:24,593 : WARNING : EPOCH - 992 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,606 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,607 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,608 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,634 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,635 : INFO : EPOCH - 993 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:24,636 : WARNING : EPOCH - 993 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,648 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,649 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,650 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,675 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,676 : INFO : EPOCH - 994 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:24,676 : WARNING : EPOCH - 994 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,686 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,687 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,688 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,713 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,714 : INFO : EPOCH - 995 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:24,715 : WARNING : EPOCH - 995 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,728 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,730 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,731 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,756 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,757 : INFO : EPOCH - 996 : training on 5386246 raw words (10000 effective u
2019-02-21 17:17:24,758 : WARNING : EPOCH - 996 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:17:24,767 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,768 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,769 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,794 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,795 : INFO : EPOCH - 997 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:24,796 : WARNING : EPOCH - 997 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,806 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,807 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,808 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,833 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,834 : INFO : EPOCH - 998 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:24,834 : WARNING : EPOCH - 998 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,841 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,843 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,844 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,870 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,874 : INFO : EPOCH - 999 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:24,875 : WARNING : EPOCH - 999 : supplied example count (1) did not equal expe
2019-02-21 17:17:24,886 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,888 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,888 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,914 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,915 : INFO : EPOCH - 1000 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:24,916 : WARNING : EPOCH - 1000 : supplied example count (1) did not equal exp
2019-02-21 17:17:24,928 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,929 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,931 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,954 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,955 : INFO : EPOCH - 1001 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:24,956 : WARNING : EPOCH - 1001 : supplied example count (1) did not equal exp
2019-02-21 17:17:24,965 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:24,966 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:24,966 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:24,991 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:24,992 : INFO : EPOCH - 1002 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:24,993 : WARNING : EPOCH - 1002 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,003 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,005 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,005 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,031 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,034 : INFO : EPOCH - 1003 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,035 : WARNING : EPOCH - 1003 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,042 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,043 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,044 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,070 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,071 : INFO : EPOCH - 1004 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,071 : WARNING : EPOCH - 1004 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:25,080 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,081 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,082 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,107 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,117 : INFO : EPOCH - 1005 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,117 : WARNING : EPOCH - 1005 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,129 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,133 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,134 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,156 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,166 : INFO : EPOCH - 1006 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,167 : WARNING : EPOCH - 1006 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,184 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,190 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,193 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,219 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,220 : INFO : EPOCH - 1007 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,221 : WARNING : EPOCH - 1007 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,233 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,234 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,234 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,259 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,261 : INFO : EPOCH - 1008 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,261 : WARNING : EPOCH - 1008 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,270 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,272 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,272 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,298 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,299 : INFO : EPOCH - 1009 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,300 : WARNING : EPOCH - 1009 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,309 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,310 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,310 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,336 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,337 : INFO : EPOCH - 1010 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,337 : WARNING : EPOCH - 1010 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,346 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,347 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,348 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,372 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,373 : INFO : EPOCH - 1011 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,374 : WARNING : EPOCH - 1011 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,384 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,385 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,386 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,411 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,412 : INFO : EPOCH - 1012 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,413 : WARNING : EPOCH - 1012 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:25,430 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,434 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,435 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,456 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,457 : INFO : EPOCH - 1013 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,457 : WARNING : EPOCH - 1013 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,473 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,474 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,497 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,498 : INFO : EPOCH - 1014 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,499 : WARNING : EPOCH - 1014 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,507 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,509 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,510 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,534 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,535 : INFO : EPOCH - 1015 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,536 : WARNING : EPOCH - 1015 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,546 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,547 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,548 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,573 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,574 : INFO : EPOCH - 1016 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,574 : WARNING : EPOCH - 1016 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,585 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,586 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,587 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,611 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,612 : INFO : EPOCH - 1017 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,612 : WARNING : EPOCH - 1017 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,620 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,622 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,622 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,649 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,650 : INFO : EPOCH - 1018 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,650 : WARNING : EPOCH - 1018 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,659 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,660 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,661 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,685 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,686 : INFO : EPOCH - 1019 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,687 : WARNING : EPOCH - 1019 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,695 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,697 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,697 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,723 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,724 : INFO : EPOCH - 1020 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,725 : WARNING : EPOCH - 1020 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:25,734 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,735 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,736 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,760 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,761 : INFO : EPOCH - 1021 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,762 : WARNING : EPOCH - 1021 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,774 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,776 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,777 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,801 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,802 : INFO : EPOCH - 1022 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,803 : WARNING : EPOCH - 1022 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,811 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,812 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,813 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,837 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,838 : INFO : EPOCH - 1023 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,839 : WARNING : EPOCH - 1023 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,848 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,849 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,850 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,874 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,875 : INFO : EPOCH - 1024 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,875 : WARNING : EPOCH - 1024 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,883 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,885 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,885 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,911 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,912 : INFO : EPOCH - 1025 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,913 : WARNING : EPOCH - 1025 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,922 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,923 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,923 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,948 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,949 : INFO : EPOCH - 1026 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,950 : WARNING : EPOCH - 1026 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,957 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,958 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,959 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:25,984 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:25,985 : INFO : EPOCH - 1027 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:25,985 : WARNING : EPOCH - 1027 : supplied example count (1) did not equal exp
2019-02-21 17:17:25,994 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:25,995 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:25,996 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,021 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,022 : INFO : EPOCH - 1028 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,023 : WARNING : EPOCH - 1028 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:26,031 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,032 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,033 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,058 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,059 : INFO : EPOCH - 1029 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,060 : WARNING : EPOCH - 1029 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,069 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,070 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,070 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,096 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,097 : INFO : EPOCH - 1030 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,098 : WARNING : EPOCH - 1030 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,106 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,107 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,108 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,134 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,135 : INFO : EPOCH - 1031 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,135 : WARNING : EPOCH - 1031 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,144 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,145 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,146 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,170 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,171 : INFO : EPOCH - 1032 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,172 : WARNING : EPOCH - 1032 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,185 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,186 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,186 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,211 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,212 : INFO : EPOCH - 1033 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,213 : WARNING : EPOCH - 1033 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,221 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,222 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,223 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,248 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,249 : INFO : EPOCH - 1034 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,250 : WARNING : EPOCH - 1034 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,258 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,259 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,260 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,285 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,315 : INFO : EPOCH - 1035 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,315 : WARNING : EPOCH - 1035 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,326 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,326 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,327 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,351 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,352 : INFO : EPOCH - 1036 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,353 : WARNING : EPOCH - 1036 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:17:26,363 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,365 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,367 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,390 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,409 : INFO : EPOCH - 1037 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,410 : WARNING : EPOCH - 1037 : supplied example count (1) did not equal exp
2019-02-21 17:17:26,424 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,425 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,426 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,450 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,451 : INFO : EPOCH - 1038 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,452 : WARNING : EPOCH - 1038 : supplied example count (1) did not equal exp
2019-02-21 17:17:26,461 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,462 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,463 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,487 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,488 : INFO : EPOCH - 1039 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,489 : WARNING : EPOCH - 1039 : supplied example count (1) did not equal exp
2019-02-21 17:17:26,497 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,498 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,499 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,524 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,525 : INFO : EPOCH - 1040 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,525 : WARNING : EPOCH - 1040 : supplied example count (1) did not equal exp
2019-02-21 17:17:26,534 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,535 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,536 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,560 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,561 : INFO : EPOCH - 1041 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,562 : WARNING : EPOCH - 1041 : supplied example count (1) did not equal exp
2019-02-21 17:17:26,571 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,572 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,573 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,597 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,598 : INFO : EPOCH - 1042 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,599 : WARNING : EPOCH - 1042 : supplied example count (1) did not equal exp
2019-02-21 17:17:26,608 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,609 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,610 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,634 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,635 : INFO : EPOCH - 1043 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,636 : WARNING : EPOCH - 1043 : supplied example count (1) did not equal exp
2019-02-21 17:17:26,651 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,653 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,653 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,678 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,679 : INFO : EPOCH - 1044 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,680 : WARNING : EPOCH - 1044 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:26,691 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,692 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,693 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,717 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,718 : INFO : EPOCH - 1045 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,719 : WARNING : EPOCH - 1045 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,728 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,729 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,729 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,754 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,755 : INFO : EPOCH - 1046 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,756 : WARNING : EPOCH - 1046 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,765 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,766 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,767 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,791 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,792 : INFO : EPOCH - 1047 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,793 : WARNING : EPOCH - 1047 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,803 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,804 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,804 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,828 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,829 : INFO : EPOCH - 1048 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,830 : WARNING : EPOCH - 1048 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,839 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,839 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,840 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,865 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,865 : INFO : EPOCH - 1049 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,866 : WARNING : EPOCH - 1049 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,876 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,877 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,878 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,902 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,903 : INFO : EPOCH - 1050 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,904 : WARNING : EPOCH - 1050 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,914 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,921 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,922 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,939 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,940 : INFO : EPOCH - 1051 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,941 : WARNING : EPOCH - 1051 : supplied example count (1) did not equal ex
2019-02-21 17:17:26,950 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,951 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,952 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:26,976 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:26,977 : INFO : EPOCH - 1052 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:26,978 : WARNING : EPOCH - 1052 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:17:26,987 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:26,988 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:26,988 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,012 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,014 : INFO : EPOCH - 1053 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,014 : WARNING : EPOCH - 1053 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,023 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,024 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,025 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,050 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,051 : INFO : EPOCH - 1054 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,051 : WARNING : EPOCH - 1054 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,059 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,060 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,061 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,086 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,087 : INFO : EPOCH - 1055 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,088 : WARNING : EPOCH - 1055 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,095 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,097 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,097 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,122 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,123 : INFO : EPOCH - 1056 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,124 : WARNING : EPOCH - 1056 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,132 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,134 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,134 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,159 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,160 : INFO : EPOCH - 1057 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,160 : WARNING : EPOCH - 1057 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,168 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,170 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,171 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,198 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,199 : INFO : EPOCH - 1058 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,200 : WARNING : EPOCH - 1058 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,208 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,210 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,211 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,236 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,237 : INFO : EPOCH - 1059 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,238 : WARNING : EPOCH - 1059 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,246 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,248 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,248 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,274 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,275 : INFO : EPOCH - 1060 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,275 : WARNING : EPOCH - 1060 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:27,284 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,286 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,286 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,311 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,312 : INFO : EPOCH - 1061 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,313 : WARNING : EPOCH - 1061 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,321 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,322 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,323 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,347 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,348 : INFO : EPOCH - 1062 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,349 : WARNING : EPOCH - 1062 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,357 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,358 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,359 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,384 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,385 : INFO : EPOCH - 1063 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,386 : WARNING : EPOCH - 1063 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,394 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,395 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,396 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,420 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,421 : INFO : EPOCH - 1064 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,422 : WARNING : EPOCH - 1064 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,430 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,431 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,432 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,456 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,457 : INFO : EPOCH - 1065 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,457 : WARNING : EPOCH - 1065 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,465 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,467 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,468 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,492 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,493 : INFO : EPOCH - 1066 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,494 : WARNING : EPOCH - 1066 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,503 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,504 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,505 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,530 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,531 : INFO : EPOCH - 1067 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,532 : WARNING : EPOCH - 1067 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,540 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,541 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,542 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,567 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,568 : INFO : EPOCH - 1068 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,569 : WARNING : EPOCH - 1068 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:27,576 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,578 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,579 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,604 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,605 : INFO : EPOCH - 1069 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,605 : WARNING : EPOCH - 1069 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,613 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,615 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,615 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,640 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,641 : INFO : EPOCH - 1070 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,642 : WARNING : EPOCH - 1070 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,653 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,654 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,655 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,679 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,680 : INFO : EPOCH - 1071 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,681 : WARNING : EPOCH - 1071 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,689 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,691 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,691 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,716 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,717 : INFO : EPOCH - 1072 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,718 : WARNING : EPOCH - 1072 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,728 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,729 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,730 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,756 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,757 : INFO : EPOCH - 1073 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,758 : WARNING : EPOCH - 1073 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,766 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,767 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,768 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,792 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,793 : INFO : EPOCH - 1074 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,794 : WARNING : EPOCH - 1074 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,803 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,805 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,806 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,829 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,830 : INFO : EPOCH - 1075 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,831 : WARNING : EPOCH - 1075 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,839 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,840 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,841 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,865 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,867 : INFO : EPOCH - 1076 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,867 : WARNING : EPOCH - 1076 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:27,876 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,877 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,878 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,904 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,905 : INFO : EPOCH - 1077 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,906 : WARNING : EPOCH - 1077 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,914 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,915 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,916 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,940 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,941 : INFO : EPOCH - 1078 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,942 : WARNING : EPOCH - 1078 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,950 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,951 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,952 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:27,977 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:27,978 : INFO : EPOCH - 1079 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:27,979 : WARNING : EPOCH - 1079 : supplied example count (1) did not equal exp
2019-02-21 17:17:27,988 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:27,989 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:27,989 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,014 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,015 : INFO : EPOCH - 1080 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,016 : WARNING : EPOCH - 1080 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,024 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,025 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,025 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,051 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,052 : INFO : EPOCH - 1081 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,053 : WARNING : EPOCH - 1081 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,062 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,063 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,063 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,088 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,089 : INFO : EPOCH - 1082 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,090 : WARNING : EPOCH - 1082 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,099 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,100 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,101 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,125 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,126 : INFO : EPOCH - 1083 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,127 : WARNING : EPOCH - 1083 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,136 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,137 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,138 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,162 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,163 : INFO : EPOCH - 1084 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,164 : WARNING : EPOCH - 1084 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:28,173 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,174 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,175 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,199 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,200 : INFO : EPOCH - 1085 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,200 : WARNING : EPOCH - 1085 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,209 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,210 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,211 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,235 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,236 : INFO : EPOCH - 1086 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,237 : WARNING : EPOCH - 1086 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,246 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,247 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,247 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,272 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,274 : INFO : EPOCH - 1087 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,274 : WARNING : EPOCH - 1087 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,283 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,284 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,284 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,309 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,309 : INFO : EPOCH - 1088 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,310 : WARNING : EPOCH - 1088 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,320 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,321 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,322 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,346 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,347 : INFO : EPOCH - 1089 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,348 : WARNING : EPOCH - 1089 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,357 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,358 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,358 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,383 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,384 : INFO : EPOCH - 1090 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,385 : WARNING : EPOCH - 1090 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,394 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,395 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,396 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,420 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,421 : INFO : EPOCH - 1091 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,422 : WARNING : EPOCH - 1091 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,431 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,432 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,433 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,458 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,459 : INFO : EPOCH - 1092 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,459 : WARNING : EPOCH - 1092 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:28,468 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,469 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,470 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,494 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,495 : INFO : EPOCH - 1093 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,496 : WARNING : EPOCH - 1093 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,504 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,532 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,533 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,534 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,535 : INFO : EPOCH - 1094 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,536 : WARNING : EPOCH - 1094 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,545 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,546 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,547 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,572 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,573 : INFO : EPOCH - 1095 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,574 : WARNING : EPOCH - 1095 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,593 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,595 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,595 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,620 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,621 : INFO : EPOCH - 1096 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,622 : WARNING : EPOCH - 1096 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,631 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,632 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,633 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,656 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,674 : INFO : EPOCH - 1097 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,675 : WARNING : EPOCH - 1097 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,684 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,685 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,686 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,710 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,711 : INFO : EPOCH - 1098 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,712 : WARNING : EPOCH - 1098 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,720 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,721 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,722 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,746 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,747 : INFO : EPOCH - 1099 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,748 : WARNING : EPOCH - 1099 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,756 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,757 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,758 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,782 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,783 : INFO : EPOCH - 1100 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,783 : WARNING : EPOCH - 1100 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:28,792 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,793 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,794 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,818 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,819 : INFO : EPOCH - 1101 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,820 : WARNING : EPOCH - 1101 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,829 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,830 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,831 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,854 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,855 : INFO : EPOCH - 1102 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,856 : WARNING : EPOCH - 1102 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,865 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,866 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,867 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,892 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,893 : INFO : EPOCH - 1103 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,894 : WARNING : EPOCH - 1103 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,904 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,905 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,906 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,930 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,931 : INFO : EPOCH - 1104 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,932 : WARNING : EPOCH - 1104 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,940 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,941 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,942 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:28,966 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:28,967 : INFO : EPOCH - 1105 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:28,968 : WARNING : EPOCH - 1105 : supplied example count (1) did not equal exp
2019-02-21 17:17:28,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:28,977 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:28,978 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,003 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,004 : INFO : EPOCH - 1106 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,005 : WARNING : EPOCH - 1106 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,013 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,015 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,016 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,040 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,041 : INFO : EPOCH - 1107 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,042 : WARNING : EPOCH - 1107 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,051 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,052 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,053 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,077 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,078 : INFO : EPOCH - 1108 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,079 : WARNING : EPOCH - 1108 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:29,087 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,089 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,089 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,114 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,115 : INFO : EPOCH - 1109 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,116 : WARNING : EPOCH - 1109 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,124 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,125 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,126 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,151 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,152 : INFO : EPOCH - 1110 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,153 : WARNING : EPOCH - 1110 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,160 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,162 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,163 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,189 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,190 : INFO : EPOCH - 1111 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,191 : WARNING : EPOCH - 1111 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,205 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,207 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,207 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,233 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,234 : INFO : EPOCH - 1112 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,234 : WARNING : EPOCH - 1112 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,242 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,244 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,244 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,270 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,271 : INFO : EPOCH - 1113 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,271 : WARNING : EPOCH - 1113 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,279 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,280 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,282 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,306 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,307 : INFO : EPOCH - 1114 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,308 : WARNING : EPOCH - 1114 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,316 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,318 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,318 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,342 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,343 : INFO : EPOCH - 1115 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,344 : WARNING : EPOCH - 1115 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,353 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,354 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,355 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,380 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,381 : INFO : EPOCH - 1116 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,382 : WARNING : EPOCH - 1116 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:29,390 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,391 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,392 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,416 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,417 : INFO : EPOCH - 1117 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,418 : WARNING : EPOCH - 1117 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,426 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,427 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,428 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,452 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,453 : INFO : EPOCH - 1118 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,454 : WARNING : EPOCH - 1118 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,462 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,463 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,464 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,488 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,489 : INFO : EPOCH - 1119 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,490 : WARNING : EPOCH - 1119 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,498 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,499 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,501 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,533 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,534 : INFO : EPOCH - 1120 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,535 : WARNING : EPOCH - 1120 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,542 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,544 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,544 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,569 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,570 : INFO : EPOCH - 1121 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,570 : WARNING : EPOCH - 1121 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,579 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,581 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,582 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,607 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,608 : INFO : EPOCH - 1122 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,608 : WARNING : EPOCH - 1122 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,618 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,619 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,620 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,646 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,648 : INFO : EPOCH - 1123 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,649 : WARNING : EPOCH - 1123 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,657 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,658 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,659 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,683 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,684 : INFO : EPOCH - 1124 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,685 : WARNING : EPOCH - 1124 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:29,693 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,695 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,695 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,720 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,721 : INFO : EPOCH - 1125 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,721 : WARNING : EPOCH - 1125 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,734 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,735 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,736 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,760 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,761 : INFO : EPOCH - 1126 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,762 : WARNING : EPOCH - 1126 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,773 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,774 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,775 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,800 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,801 : INFO : EPOCH - 1127 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,802 : WARNING : EPOCH - 1127 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,812 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,813 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,814 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,837 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,838 : INFO : EPOCH - 1128 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,839 : WARNING : EPOCH - 1128 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,850 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,851 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,852 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,877 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,878 : INFO : EPOCH - 1129 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,878 : WARNING : EPOCH - 1129 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,888 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,890 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,890 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,915 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,916 : INFO : EPOCH - 1130 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,917 : WARNING : EPOCH - 1130 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,929 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,931 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,931 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,956 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,957 : INFO : EPOCH - 1131 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,958 : WARNING : EPOCH - 1131 : supplied example count (1) did not equal exp
2019-02-21 17:17:29,969 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:29,971 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:29,971 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:29,996 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:29,997 : INFO : EPOCH - 1132 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:29,998 : WARNING : EPOCH - 1132 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:30,009 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,010 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,011 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,036 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,037 : INFO : EPOCH - 1133 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,037 : WARNING : EPOCH - 1133 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,049 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,050 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,051 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,074 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,075 : INFO : EPOCH - 1134 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,076 : WARNING : EPOCH - 1134 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,085 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,086 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,087 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,111 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,112 : INFO : EPOCH - 1135 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,112 : WARNING : EPOCH - 1135 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,120 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,122 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,122 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,147 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,148 : INFO : EPOCH - 1136 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,149 : WARNING : EPOCH - 1136 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,157 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,158 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,159 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,183 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,186 : INFO : EPOCH - 1137 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,186 : WARNING : EPOCH - 1137 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,199 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,202 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,203 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,225 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,226 : INFO : EPOCH - 1138 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,227 : WARNING : EPOCH - 1138 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,235 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,237 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,238 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,262 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,264 : INFO : EPOCH - 1139 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,264 : WARNING : EPOCH - 1139 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,273 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,274 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,275 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,300 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,301 : INFO : EPOCH - 1140 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,302 : WARNING : EPOCH - 1140 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:30,312 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,313 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,314 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,339 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,340 : INFO : EPOCH - 1141 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,341 : WARNING : EPOCH - 1141 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,350 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,351 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,352 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,377 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,378 : INFO : EPOCH - 1142 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,379 : WARNING : EPOCH - 1142 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,388 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,389 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,390 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,414 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,415 : INFO : EPOCH - 1143 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,416 : WARNING : EPOCH - 1143 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,424 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,425 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,426 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,450 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,451 : INFO : EPOCH - 1144 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,452 : WARNING : EPOCH - 1144 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,461 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,462 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,463 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,487 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,488 : INFO : EPOCH - 1145 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,489 : WARNING : EPOCH - 1145 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,498 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,506 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,507 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,525 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,526 : INFO : EPOCH - 1146 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,527 : WARNING : EPOCH - 1146 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,535 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,537 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,538 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,562 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,563 : INFO : EPOCH - 1147 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,564 : WARNING : EPOCH - 1147 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,572 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,573 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,574 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,599 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,600 : INFO : EPOCH - 1148 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,600 : WARNING : EPOCH - 1148 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:30,608 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,610 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,610 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,635 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,636 : INFO : EPOCH - 1149 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,637 : WARNING : EPOCH - 1149 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,645 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,646 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,647 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,671 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,672 : INFO : EPOCH - 1150 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,672 : WARNING : EPOCH - 1150 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,681 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,682 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,683 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,708 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,709 : INFO : EPOCH - 1151 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,709 : WARNING : EPOCH - 1151 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,718 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,719 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,720 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,745 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,747 : INFO : EPOCH - 1152 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,747 : WARNING : EPOCH - 1152 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,757 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,758 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,759 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,785 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,786 : INFO : EPOCH - 1153 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,787 : WARNING : EPOCH - 1153 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,822 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,824 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,825 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,825 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,826 : INFO : EPOCH - 1154 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,827 : WARNING : EPOCH - 1154 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,836 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,837 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,838 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,862 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,863 : INFO : EPOCH - 1155 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,864 : WARNING : EPOCH - 1155 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,872 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,875 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,875 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,899 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,900 : INFO : EPOCH - 1156 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,901 : WARNING : EPOCH - 1156 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:30,910 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,911 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,912 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,937 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,938 : INFO : EPOCH - 1157 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,939 : WARNING : EPOCH - 1157 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,948 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,949 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,950 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:30,975 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:30,981 : INFO : EPOCH - 1158 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:30,982 : WARNING : EPOCH - 1158 : supplied example count (1) did not equal exp
2019-02-21 17:17:30,991 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:30,992 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:30,992 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,017 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,018 : INFO : EPOCH - 1159 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,018 : WARNING : EPOCH - 1159 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,028 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,029 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,029 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,053 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,054 : INFO : EPOCH - 1160 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,055 : WARNING : EPOCH - 1160 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,064 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,065 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,066 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,090 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,091 : INFO : EPOCH - 1161 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,092 : WARNING : EPOCH - 1161 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,102 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,104 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,104 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,129 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,130 : INFO : EPOCH - 1162 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,131 : WARNING : EPOCH - 1162 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,139 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,140 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,141 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,165 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,166 : INFO : EPOCH - 1163 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,167 : WARNING : EPOCH - 1163 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,176 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,177 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,177 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,201 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,202 : INFO : EPOCH - 1164 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,202 : WARNING : EPOCH - 1164 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:31,211 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,212 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,213 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,238 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,239 : INFO : EPOCH - 1165 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,239 : WARNING : EPOCH - 1165 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,248 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,250 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,250 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,275 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,276 : INFO : EPOCH - 1166 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,276 : WARNING : EPOCH - 1166 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,288 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,289 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,290 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,315 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,323 : INFO : EPOCH - 1167 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,323 : WARNING : EPOCH - 1167 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,332 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,333 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,334 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,358 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,359 : INFO : EPOCH - 1168 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,360 : WARNING : EPOCH - 1168 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,369 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,370 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,371 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,396 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,397 : INFO : EPOCH - 1169 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,398 : WARNING : EPOCH - 1169 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,408 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,409 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,410 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,434 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,435 : INFO : EPOCH - 1170 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,436 : WARNING : EPOCH - 1170 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,445 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,446 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,447 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,472 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,473 : INFO : EPOCH - 1171 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,473 : WARNING : EPOCH - 1171 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,482 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,483 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,483 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,507 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,508 : INFO : EPOCH - 1172 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,509 : WARNING : EPOCH - 1172 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:31,518 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,519 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,520 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,543 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,544 : INFO : EPOCH - 1173 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,545 : WARNING : EPOCH - 1173 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,556 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,557 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,558 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,581 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,582 : INFO : EPOCH - 1174 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,582 : WARNING : EPOCH - 1174 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,595 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,597 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,598 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,622 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,623 : INFO : EPOCH - 1175 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,624 : WARNING : EPOCH - 1175 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,635 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,636 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,636 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,661 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,662 : INFO : EPOCH - 1176 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,663 : WARNING : EPOCH - 1176 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,671 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,673 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,673 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,698 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,699 : INFO : EPOCH - 1177 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,699 : WARNING : EPOCH - 1177 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,707 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,708 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,709 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,734 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,735 : INFO : EPOCH - 1178 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,735 : WARNING : EPOCH - 1178 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,745 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,746 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,747 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,772 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,773 : INFO : EPOCH - 1179 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,774 : WARNING : EPOCH - 1179 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,786 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,787 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,788 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,813 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,814 : INFO : EPOCH - 1180 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,815 : WARNING : EPOCH - 1180 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:31,827 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,828 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,830 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,854 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,855 : INFO : EPOCH - 1181 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,856 : WARNING : EPOCH - 1181 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,864 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,866 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,867 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,892 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,893 : INFO : EPOCH - 1182 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,893 : WARNING : EPOCH - 1182 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,903 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,904 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,905 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,929 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,930 : INFO : EPOCH - 1183 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,931 : WARNING : EPOCH - 1183 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,940 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,941 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,941 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:31,965 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:31,966 : INFO : EPOCH - 1184 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:31,967 : WARNING : EPOCH - 1184 : supplied example count (1) did not equal exp
2019-02-21 17:17:31,975 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:31,976 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:31,977 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,002 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,003 : INFO : EPOCH - 1185 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,004 : WARNING : EPOCH - 1185 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,012 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,013 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,014 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,038 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,039 : INFO : EPOCH - 1186 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,040 : WARNING : EPOCH - 1186 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,048 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,050 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,050 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,075 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,076 : INFO : EPOCH - 1187 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,076 : WARNING : EPOCH - 1187 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,085 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,086 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,087 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,112 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,113 : INFO : EPOCH - 1188 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,114 : WARNING : EPOCH - 1188 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:32,122 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,123 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,124 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,149 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,150 : INFO : EPOCH - 1189 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,150 : WARNING : EPOCH - 1189 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,159 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,160 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,161 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,187 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,188 : INFO : EPOCH - 1190 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,189 : WARNING : EPOCH - 1190 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,200 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,201 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,201 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,226 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,227 : INFO : EPOCH - 1191 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,228 : WARNING : EPOCH - 1191 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,236 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,237 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,238 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,262 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,263 : INFO : EPOCH - 1192 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,264 : WARNING : EPOCH - 1192 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,273 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,274 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,275 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,299 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,300 : INFO : EPOCH - 1193 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,301 : WARNING : EPOCH - 1193 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,309 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,310 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,310 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,334 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,335 : INFO : EPOCH - 1194 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,336 : WARNING : EPOCH - 1194 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,344 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,345 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,346 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,370 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,371 : INFO : EPOCH - 1195 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,372 : WARNING : EPOCH - 1195 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,381 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,382 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,383 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,407 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,408 : INFO : EPOCH - 1196 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,408 : WARNING : EPOCH - 1196 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:32,417 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,425 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,426 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,444 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,445 : INFO : EPOCH - 1197 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,446 : WARNING : EPOCH - 1197 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,454 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,455 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,456 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,482 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,483 : INFO : EPOCH - 1198 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,484 : WARNING : EPOCH - 1198 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,493 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,494 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,495 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,519 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,521 : INFO : EPOCH - 1199 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,521 : WARNING : EPOCH - 1199 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,529 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,531 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,531 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,557 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,558 : INFO : EPOCH - 1200 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,559 : WARNING : EPOCH - 1200 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,567 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,568 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,569 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,595 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,596 : INFO : EPOCH - 1201 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,596 : WARNING : EPOCH - 1201 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,605 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,606 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,607 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,632 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,633 : INFO : EPOCH - 1202 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,634 : WARNING : EPOCH - 1202 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,641 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,643 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,644 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,671 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,671 : INFO : EPOCH - 1203 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,672 : WARNING : EPOCH - 1203 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,681 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,682 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,682 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,707 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,708 : INFO : EPOCH - 1204 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,709 : WARNING : EPOCH - 1204 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:32,719 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,720 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,720 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,745 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,746 : INFO : EPOCH - 1205 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,746 : WARNING : EPOCH - 1205 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,756 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,757 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,757 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,782 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,783 : INFO : EPOCH - 1206 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,783 : WARNING : EPOCH - 1206 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,792 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,793 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,794 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,819 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,820 : INFO : EPOCH - 1207 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,820 : WARNING : EPOCH - 1207 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,828 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,830 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,831 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,855 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,856 : INFO : EPOCH - 1208 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,857 : WARNING : EPOCH - 1208 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,865 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,867 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,867 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,892 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,893 : INFO : EPOCH - 1209 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,894 : WARNING : EPOCH - 1209 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,902 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,903 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,904 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,928 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,929 : INFO : EPOCH - 1210 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,930 : WARNING : EPOCH - 1210 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,939 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,940 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,940 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:32,966 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:32,967 : INFO : EPOCH - 1211 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:32,968 : WARNING : EPOCH - 1211 : supplied example count (1) did not equal exp
2019-02-21 17:17:32,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:32,977 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:32,978 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,003 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,004 : INFO : EPOCH - 1212 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,005 : WARNING : EPOCH - 1212 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:33,013 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,014 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,015 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,039 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,040 : INFO : EPOCH - 1213 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,041 : WARNING : EPOCH - 1213 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,050 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,051 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,052 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,075 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,076 : INFO : EPOCH - 1214 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,077 : WARNING : EPOCH - 1214 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,085 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,086 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,087 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,113 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,114 : INFO : EPOCH - 1215 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,115 : WARNING : EPOCH - 1215 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,123 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,124 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,125 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,150 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,151 : INFO : EPOCH - 1216 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,151 : WARNING : EPOCH - 1216 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,160 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,161 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,162 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,188 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,189 : INFO : EPOCH - 1217 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,190 : WARNING : EPOCH - 1217 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,202 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,203 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,204 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,229 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,269 : INFO : EPOCH - 1218 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,270 : WARNING : EPOCH - 1218 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,279 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,280 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,280 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,306 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,307 : INFO : EPOCH - 1219 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,308 : WARNING : EPOCH - 1219 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,317 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,318 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,319 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,343 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,344 : INFO : EPOCH - 1220 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,345 : WARNING : EPOCH - 1220 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:33,353 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,354 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,354 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,378 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,379 : INFO : EPOCH - 1221 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,380 : WARNING : EPOCH - 1221 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,388 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,389 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,390 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,414 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,415 : INFO : EPOCH - 1222 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,416 : WARNING : EPOCH - 1222 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,424 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,425 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,426 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,449 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,459 : INFO : EPOCH - 1223 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,461 : WARNING : EPOCH - 1223 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,498 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,499 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,499 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,523 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,524 : INFO : EPOCH - 1224 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,525 : WARNING : EPOCH - 1224 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,538 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,540 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,541 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,564 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,565 : INFO : EPOCH - 1225 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,566 : WARNING : EPOCH - 1225 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,574 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,575 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,575 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,600 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,601 : INFO : EPOCH - 1226 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,602 : WARNING : EPOCH - 1226 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,610 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,611 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,612 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,636 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,637 : INFO : EPOCH - 1227 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,637 : WARNING : EPOCH - 1227 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,648 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,649 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,650 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,674 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,675 : INFO : EPOCH - 1228 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,676 : WARNING : EPOCH - 1228 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:33,684 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,685 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,686 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,710 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,711 : INFO : EPOCH - 1229 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,712 : WARNING : EPOCH - 1229 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,720 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,721 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,722 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,747 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,748 : INFO : EPOCH - 1230 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,748 : WARNING : EPOCH - 1230 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,759 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,760 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,762 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,785 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,786 : INFO : EPOCH - 1231 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,786 : WARNING : EPOCH - 1231 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,801 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,804 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,805 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,825 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,826 : INFO : EPOCH - 1232 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,827 : WARNING : EPOCH - 1232 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,837 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,838 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,839 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,862 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,863 : INFO : EPOCH - 1233 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,864 : WARNING : EPOCH - 1233 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,872 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,873 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,874 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,898 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,899 : INFO : EPOCH - 1234 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,900 : WARNING : EPOCH - 1234 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,908 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,909 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,910 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,935 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,936 : INFO : EPOCH - 1235 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,937 : WARNING : EPOCH - 1235 : supplied example count (1) did not equal exp
2019-02-21 17:17:33,945 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,946 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,947 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:33,971 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:33,972 : INFO : EPOCH - 1236 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:33,973 : WARNING : EPOCH - 1236 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:33,981 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:33,982 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:33,983 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,007 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,008 : INFO : EPOCH - 1237 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,009 : WARNING : EPOCH - 1237 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,018 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,019 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,020 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,048 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,049 : INFO : EPOCH - 1238 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,050 : WARNING : EPOCH - 1238 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,058 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,068 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,069 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,084 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,085 : INFO : EPOCH - 1239 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,086 : WARNING : EPOCH - 1239 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,094 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,095 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,096 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,121 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,122 : INFO : EPOCH - 1240 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,123 : WARNING : EPOCH - 1240 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,132 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,134 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,134 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,158 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,159 : INFO : EPOCH - 1241 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,160 : WARNING : EPOCH - 1241 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,168 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,169 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,170 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,194 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,195 : INFO : EPOCH - 1242 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,195 : WARNING : EPOCH - 1242 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,206 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,207 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,208 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,232 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,234 : INFO : EPOCH - 1243 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,234 : WARNING : EPOCH - 1243 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,243 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,244 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,245 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,270 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,271 : INFO : EPOCH - 1244 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,271 : WARNING : EPOCH - 1244 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:34,285 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,286 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,287 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,312 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,320 : INFO : EPOCH - 1245 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,320 : WARNING : EPOCH - 1245 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,329 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,330 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,331 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,355 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,356 : INFO : EPOCH - 1246 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,357 : WARNING : EPOCH - 1246 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,366 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,368 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,368 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,393 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,394 : INFO : EPOCH - 1247 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,395 : WARNING : EPOCH - 1247 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,404 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,405 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,406 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,430 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,431 : INFO : EPOCH - 1248 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,432 : WARNING : EPOCH - 1248 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,439 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,441 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,441 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,466 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,467 : INFO : EPOCH - 1249 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,468 : WARNING : EPOCH - 1249 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,476 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,477 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,478 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,502 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,503 : INFO : EPOCH - 1250 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,504 : WARNING : EPOCH - 1250 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,511 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,513 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,513 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,538 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,539 : INFO : EPOCH - 1251 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,542 : WARNING : EPOCH - 1251 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,551 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,552 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,552 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,576 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,585 : INFO : EPOCH - 1252 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,585 : WARNING : EPOCH - 1252 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:34,593 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,595 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,595 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,621 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,621 : INFO : EPOCH - 1253 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,622 : WARNING : EPOCH - 1253 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,630 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,631 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,632 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,657 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,658 : INFO : EPOCH - 1254 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,659 : WARNING : EPOCH - 1254 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,668 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,669 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,669 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,693 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,694 : INFO : EPOCH - 1255 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,695 : WARNING : EPOCH - 1255 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,704 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,705 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,705 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,730 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,731 : INFO : EPOCH - 1256 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,731 : WARNING : EPOCH - 1256 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,740 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,741 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,742 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,767 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,768 : INFO : EPOCH - 1257 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,768 : WARNING : EPOCH - 1257 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,776 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,778 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,779 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,804 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,805 : INFO : EPOCH - 1258 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,805 : WARNING : EPOCH - 1258 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,813 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,815 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,815 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,841 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,842 : INFO : EPOCH - 1259 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,843 : WARNING : EPOCH - 1259 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,851 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,852 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,853 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,880 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,881 : INFO : EPOCH - 1260 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,881 : WARNING : EPOCH - 1260 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:34,891 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,892 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,892 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,917 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,918 : INFO : EPOCH - 1261 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,919 : WARNING : EPOCH - 1261 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,927 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,929 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,930 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,953 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,955 : INFO : EPOCH - 1262 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,955 : WARNING : EPOCH - 1262 : supplied example count (1) did not equal exp
2019-02-21 17:17:34,964 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:34,965 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:34,965 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:34,990 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:34,991 : INFO : EPOCH - 1263 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:34,992 : WARNING : EPOCH - 1263 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,001 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,002 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,003 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,027 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,028 : INFO : EPOCH - 1264 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,029 : WARNING : EPOCH - 1264 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,037 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,039 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,040 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,065 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,066 : INFO : EPOCH - 1265 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,067 : WARNING : EPOCH - 1265 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,075 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,076 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,077 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,102 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,103 : INFO : EPOCH - 1266 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,104 : WARNING : EPOCH - 1266 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,113 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,114 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,114 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,140 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,141 : INFO : EPOCH - 1267 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,142 : WARNING : EPOCH - 1267 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,150 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,151 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,152 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,177 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,178 : INFO : EPOCH - 1268 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,179 : WARNING : EPOCH - 1268 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:35,190 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,191 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,192 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,218 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,219 : INFO : EPOCH - 1269 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,220 : WARNING : EPOCH - 1269 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,229 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,230 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,231 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,255 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,256 : INFO : EPOCH - 1270 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,257 : WARNING : EPOCH - 1270 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,266 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,267 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,268 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,291 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,292 : INFO : EPOCH - 1271 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,293 : WARNING : EPOCH - 1271 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,303 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,304 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,305 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,329 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,330 : INFO : EPOCH - 1272 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,330 : WARNING : EPOCH - 1272 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,338 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,339 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,340 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,365 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,366 : INFO : EPOCH - 1273 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,367 : WARNING : EPOCH - 1273 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,375 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,376 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,377 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,402 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,402 : INFO : EPOCH - 1274 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,403 : WARNING : EPOCH - 1274 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,411 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,413 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,414 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,438 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,439 : INFO : EPOCH - 1275 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,440 : WARNING : EPOCH - 1275 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,448 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,450 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,450 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,476 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,477 : INFO : EPOCH - 1276 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,477 : WARNING : EPOCH - 1276 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:35,486 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,487 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,488 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,512 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,513 : INFO : EPOCH - 1277 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,514 : WARNING : EPOCH - 1277 : supplied example count (1) did not equal ex
2019-02-21 17:17:35,522 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,523 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,524 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,549 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,550 : INFO : EPOCH - 1278 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,550 : WARNING : EPOCH - 1278 : supplied example count (1) did not equal ex
2019-02-21 17:17:35,559 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,560 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,561 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,585 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,588 : INFO : EPOCH - 1279 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,613 : WARNING : EPOCH - 1279 : supplied example count (1) did not equal ex
2019-02-21 17:17:35,622 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,623 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,624 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,649 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,650 : INFO : EPOCH - 1280 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,651 : WARNING : EPOCH - 1280 : supplied example count (1) did not equal ex
2019-02-21 17:17:35,662 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,663 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,664 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,688 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,717 : INFO : EPOCH - 1281 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,717 : WARNING : EPOCH - 1281 : supplied example count (1) did not equal ex
2019-02-21 17:17:35,725 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,726 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,727 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,752 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,753 : INFO : EPOCH - 1282 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,754 : WARNING : EPOCH - 1282 : supplied example count (1) did not equal ex
2019-02-21 17:17:35,762 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,763 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,764 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,788 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,789 : INFO : EPOCH - 1283 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,790 : WARNING : EPOCH - 1283 : supplied example count (1) did not equal ex
2019-02-21 17:17:35,798 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,800 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,800 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,825 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,826 : INFO : EPOCH - 1284 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,826 : WARNING : EPOCH - 1284 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:17:35,835 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,836 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,837 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,861 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,862 : INFO : EPOCH - 1285 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,863 : WARNING : EPOCH - 1285 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,871 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,873 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,873 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,898 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,899 : INFO : EPOCH - 1286 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,900 : WARNING : EPOCH - 1286 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,908 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,909 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,910 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,934 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,936 : INFO : EPOCH - 1287 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,937 : WARNING : EPOCH - 1287 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,945 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,946 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,947 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:35,971 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:35,987 : INFO : EPOCH - 1288 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:35,988 : WARNING : EPOCH - 1288 : supplied example count (1) did not equal exp
2019-02-21 17:17:35,996 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:35,997 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:35,998 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,023 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,024 : INFO : EPOCH - 1289 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,025 : WARNING : EPOCH - 1289 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,034 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,035 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,035 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,059 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,060 : INFO : EPOCH - 1290 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,061 : WARNING : EPOCH - 1290 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,071 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,072 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,073 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,097 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,098 : INFO : EPOCH - 1291 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,099 : WARNING : EPOCH - 1291 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,108 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,108 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,109 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,140 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,147 : INFO : EPOCH - 1292 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,148 : WARNING : EPOCH - 1292 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:36,161 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,162 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,163 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,198 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,224 : INFO : EPOCH - 1293 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,228 : WARNING : EPOCH - 1293 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,238 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,241 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,242 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,268 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,272 : INFO : EPOCH - 1294 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,272 : WARNING : EPOCH - 1294 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,281 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,294 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,296 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,309 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,310 : INFO : EPOCH - 1295 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,310 : WARNING : EPOCH - 1295 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,319 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,321 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,322 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,348 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,349 : INFO : EPOCH - 1296 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,350 : WARNING : EPOCH - 1296 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,358 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,359 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,360 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,386 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,387 : INFO : EPOCH - 1297 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,387 : WARNING : EPOCH - 1297 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,397 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,398 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,399 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,425 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,426 : INFO : EPOCH - 1298 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,427 : WARNING : EPOCH - 1298 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,434 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,436 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,436 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,461 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,462 : INFO : EPOCH - 1299 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,463 : WARNING : EPOCH - 1299 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,473 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,473 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,499 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,500 : INFO : EPOCH - 1300 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,501 : WARNING : EPOCH - 1300 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:36,509 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,515 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,516 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,536 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,537 : INFO : EPOCH - 1301 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,538 : WARNING : EPOCH - 1301 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,546 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,556 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,557 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,574 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,575 : INFO : EPOCH - 1302 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,575 : WARNING : EPOCH - 1302 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,587 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,587 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,588 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,613 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,614 : INFO : EPOCH - 1303 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,614 : WARNING : EPOCH - 1303 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,623 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,624 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,625 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,651 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,652 : INFO : EPOCH - 1304 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,653 : WARNING : EPOCH - 1304 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,661 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,662 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,663 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,687 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,688 : INFO : EPOCH - 1305 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,688 : WARNING : EPOCH - 1305 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,697 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,698 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,699 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,723 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,734 : INFO : EPOCH - 1306 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,734 : WARNING : EPOCH - 1306 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,743 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,744 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,744 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,773 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,774 : INFO : EPOCH - 1307 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,776 : WARNING : EPOCH - 1307 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,795 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,799 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,799 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,822 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,823 : INFO : EPOCH - 1308 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,824 : WARNING : EPOCH - 1308 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:36,833 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,834 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,835 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,860 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,869 : INFO : EPOCH - 1309 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,870 : WARNING : EPOCH - 1309 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,879 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,880 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,881 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,906 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,907 : INFO : EPOCH - 1310 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,907 : WARNING : EPOCH - 1310 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,916 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,918 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,919 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,944 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,945 : INFO : EPOCH - 1311 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,945 : WARNING : EPOCH - 1311 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,954 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,955 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,956 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:36,981 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:36,982 : INFO : EPOCH - 1312 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:36,982 : WARNING : EPOCH - 1312 : supplied example count (1) did not equal exp
2019-02-21 17:17:36,990 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:36,992 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:36,993 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,018 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,019 : INFO : EPOCH - 1313 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,020 : WARNING : EPOCH - 1313 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,029 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,030 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,030 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,056 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,057 : INFO : EPOCH - 1314 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,057 : WARNING : EPOCH - 1314 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,067 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,068 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,068 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,093 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,094 : INFO : EPOCH - 1315 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,095 : WARNING : EPOCH - 1315 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,104 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,105 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,105 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,131 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,132 : INFO : EPOCH - 1316 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,133 : WARNING : EPOCH - 1316 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:37,141 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,142 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,143 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,168 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,169 : INFO : EPOCH - 1317 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,170 : WARNING : EPOCH - 1317 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,180 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,181 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,182 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,207 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,208 : INFO : EPOCH - 1318 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,208 : WARNING : EPOCH - 1318 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,217 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,218 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,219 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,243 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,244 : INFO : EPOCH - 1319 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,245 : WARNING : EPOCH - 1319 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,254 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,255 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,256 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,281 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,282 : INFO : EPOCH - 1320 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,283 : WARNING : EPOCH - 1320 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,292 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,293 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,294 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,319 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,320 : INFO : EPOCH - 1321 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,320 : WARNING : EPOCH - 1321 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,329 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,331 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,331 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,356 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,357 : INFO : EPOCH - 1322 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,358 : WARNING : EPOCH - 1322 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,367 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,368 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,369 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,393 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,396 : INFO : EPOCH - 1323 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,397 : WARNING : EPOCH - 1323 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,405 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,406 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,407 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,432 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,433 : INFO : EPOCH - 1324 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,434 : WARNING : EPOCH - 1324 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:37,442 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,444 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,445 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,469 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,470 : INFO : EPOCH - 1325 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,471 : WARNING : EPOCH - 1325 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,480 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,481 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,482 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,506 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,507 : INFO : EPOCH - 1326 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,507 : WARNING : EPOCH - 1326 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,516 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,517 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,518 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,543 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,544 : INFO : EPOCH - 1327 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,544 : WARNING : EPOCH - 1327 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,552 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,553 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,554 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,579 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,580 : INFO : EPOCH - 1328 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,581 : WARNING : EPOCH - 1328 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,589 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,590 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,591 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,616 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,618 : INFO : EPOCH - 1329 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,619 : WARNING : EPOCH - 1329 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,627 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,628 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,629 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,654 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,655 : INFO : EPOCH - 1330 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,656 : WARNING : EPOCH - 1330 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,665 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,666 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,666 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,691 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,692 : INFO : EPOCH - 1331 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,692 : WARNING : EPOCH - 1331 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,701 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,702 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,704 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,731 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,732 : INFO : EPOCH - 1332 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,733 : WARNING : EPOCH - 1332 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:37,741 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,743 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,743 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,768 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,769 : INFO : EPOCH - 1333 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,770 : WARNING : EPOCH - 1333 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,779 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,780 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,781 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,806 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,807 : INFO : EPOCH - 1334 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,807 : WARNING : EPOCH - 1334 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,816 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,817 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,818 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,843 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,844 : INFO : EPOCH - 1335 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,845 : WARNING : EPOCH - 1335 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,854 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,855 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,856 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,881 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,882 : INFO : EPOCH - 1336 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,883 : WARNING : EPOCH - 1336 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,892 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,892 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,893 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,918 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,919 : INFO : EPOCH - 1337 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,920 : WARNING : EPOCH - 1337 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,928 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,929 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,930 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:37,956 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:37,969 : INFO : EPOCH - 1338 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:37,969 : WARNING : EPOCH - 1338 : supplied example count (1) did not equal exp
2019-02-21 17:17:37,977 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:37,979 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:37,980 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,006 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,007 : INFO : EPOCH - 1339 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,008 : WARNING : EPOCH - 1339 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,017 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,018 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,018 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,043 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,044 : INFO : EPOCH - 1340 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,045 : WARNING : EPOCH - 1340 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:38,053 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,055 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,056 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,081 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,090 : INFO : EPOCH - 1341 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,091 : WARNING : EPOCH - 1341 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,100 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,101 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,102 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,126 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,127 : INFO : EPOCH - 1342 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,128 : WARNING : EPOCH - 1342 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,137 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,138 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,138 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,163 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,164 : INFO : EPOCH - 1343 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,165 : WARNING : EPOCH - 1343 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,174 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,186 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,187 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,201 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,202 : INFO : EPOCH - 1344 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,202 : WARNING : EPOCH - 1344 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,214 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,215 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,216 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,240 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,241 : INFO : EPOCH - 1345 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,242 : WARNING : EPOCH - 1345 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,251 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,258 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,271 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,288 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,290 : INFO : EPOCH - 1346 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,290 : WARNING : EPOCH - 1346 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,310 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,311 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,312 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,337 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,338 : INFO : EPOCH - 1347 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,338 : WARNING : EPOCH - 1347 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,348 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,349 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,350 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,374 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,375 : INFO : EPOCH - 1348 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,376 : WARNING : EPOCH - 1348 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:38,385 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,386 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,387 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,413 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,414 : INFO : EPOCH - 1349 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,415 : WARNING : EPOCH - 1349 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,424 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,425 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,426 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,451 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,452 : INFO : EPOCH - 1350 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,453 : WARNING : EPOCH - 1350 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,462 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,463 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,463 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,488 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,489 : INFO : EPOCH - 1351 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,490 : WARNING : EPOCH - 1351 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,499 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,500 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,500 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,525 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,526 : INFO : EPOCH - 1352 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,526 : WARNING : EPOCH - 1352 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,537 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,538 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,539 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,564 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,565 : INFO : EPOCH - 1353 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,566 : WARNING : EPOCH - 1353 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,574 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,575 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,576 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,602 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,603 : INFO : EPOCH - 1354 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,603 : WARNING : EPOCH - 1354 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,612 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,613 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,614 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,639 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,640 : INFO : EPOCH - 1355 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,640 : WARNING : EPOCH - 1355 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,649 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,650 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,651 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,675 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,676 : INFO : EPOCH - 1356 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,676 : WARNING : EPOCH - 1356 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:38,685 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,686 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,687 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,711 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,712 : INFO : EPOCH - 1357 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,713 : WARNING : EPOCH - 1357 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,722 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,722 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,723 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,747 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,748 : INFO : EPOCH - 1358 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,749 : WARNING : EPOCH - 1358 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,758 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,759 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,760 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,785 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,786 : INFO : EPOCH - 1359 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,787 : WARNING : EPOCH - 1359 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,795 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,796 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,797 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,822 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,823 : INFO : EPOCH - 1360 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,824 : WARNING : EPOCH - 1360 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,833 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,834 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,834 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,859 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,860 : INFO : EPOCH - 1361 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,861 : WARNING : EPOCH - 1361 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,870 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,870 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,871 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,896 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,897 : INFO : EPOCH - 1362 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,898 : WARNING : EPOCH - 1362 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,907 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,908 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,909 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,940 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,949 : INFO : EPOCH - 1363 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,950 : WARNING : EPOCH - 1363 : supplied example count (1) did not equal exp
2019-02-21 17:17:38,966 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:38,967 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:38,968 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:38,993 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:38,994 : INFO : EPOCH - 1364 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:38,995 : WARNING : EPOCH - 1364 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:39,004 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,005 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,005 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,030 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,031 : INFO : EPOCH - 1365 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,032 : WARNING : EPOCH - 1365 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,040 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,041 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,041 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,066 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,067 : INFO : EPOCH - 1366 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,068 : WARNING : EPOCH - 1366 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,076 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,077 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,079 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,104 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,105 : INFO : EPOCH - 1367 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,105 : WARNING : EPOCH - 1367 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,114 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,115 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,116 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,141 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,142 : INFO : EPOCH - 1368 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,142 : WARNING : EPOCH - 1368 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,150 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,152 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,153 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,178 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,179 : INFO : EPOCH - 1369 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,180 : WARNING : EPOCH - 1369 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,190 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,192 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,192 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,217 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,218 : INFO : EPOCH - 1370 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,219 : WARNING : EPOCH - 1370 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,227 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,229 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,230 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,254 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,256 : INFO : EPOCH - 1371 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,256 : WARNING : EPOCH - 1371 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,267 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,268 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,269 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,293 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,294 : INFO : EPOCH - 1372 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,295 : WARNING : EPOCH - 1372 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:39,304 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,305 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,306 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,331 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,332 : INFO : EPOCH - 1373 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,333 : WARNING : EPOCH - 1373 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,341 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,342 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,342 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,367 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,368 : INFO : EPOCH - 1374 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,368 : WARNING : EPOCH - 1374 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,377 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,378 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,379 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,404 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,405 : INFO : EPOCH - 1375 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,405 : WARNING : EPOCH - 1375 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,414 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,415 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,416 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,439 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,440 : INFO : EPOCH - 1376 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,441 : WARNING : EPOCH - 1376 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,449 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,451 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,451 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,477 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,478 : INFO : EPOCH - 1377 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,478 : WARNING : EPOCH - 1377 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,487 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,489 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,489 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,515 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,516 : INFO : EPOCH - 1378 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,517 : WARNING : EPOCH - 1378 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,524 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,526 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,526 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,552 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,553 : INFO : EPOCH - 1379 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,553 : WARNING : EPOCH - 1379 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,562 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,563 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,564 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,588 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,589 : INFO : EPOCH - 1380 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,590 : WARNING : EPOCH - 1380 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:39,598 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,600 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,600 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,625 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,626 : INFO : EPOCH - 1381 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,627 : WARNING : EPOCH - 1381 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,636 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,637 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,637 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,663 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,664 : INFO : EPOCH - 1382 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,665 : WARNING : EPOCH - 1382 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,673 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,674 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,675 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,699 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,700 : INFO : EPOCH - 1383 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,701 : WARNING : EPOCH - 1383 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,709 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,711 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,711 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,736 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,737 : INFO : EPOCH - 1384 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,738 : WARNING : EPOCH - 1384 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,746 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,747 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,748 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,773 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,774 : INFO : EPOCH - 1385 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,775 : WARNING : EPOCH - 1385 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,784 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,785 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,785 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,809 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,811 : INFO : EPOCH - 1386 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,811 : WARNING : EPOCH - 1386 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,820 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,821 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,822 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,846 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,847 : INFO : EPOCH - 1387 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,848 : WARNING : EPOCH - 1387 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,856 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,857 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,858 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,883 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,884 : INFO : EPOCH - 1388 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,885 : WARNING : EPOCH - 1388 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:39,896 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,897 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,898 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,922 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,923 : INFO : EPOCH - 1389 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,924 : WARNING : EPOCH - 1389 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,933 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,935 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,935 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,960 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,961 : INFO : EPOCH - 1390 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,961 : WARNING : EPOCH - 1390 : supplied example count (1) did not equal exp
2019-02-21 17:17:39,971 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:39,972 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:39,973 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:39,997 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:39,998 : INFO : EPOCH - 1391 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:39,999 : WARNING : EPOCH - 1391 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,009 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,009 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,010 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,035 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,036 : INFO : EPOCH - 1392 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,036 : WARNING : EPOCH - 1392 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,044 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,046 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,046 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,071 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,072 : INFO : EPOCH - 1393 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,073 : WARNING : EPOCH - 1393 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,081 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,082 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,083 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,107 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,108 : INFO : EPOCH - 1394 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,109 : WARNING : EPOCH - 1394 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,119 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,120 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,120 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,144 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,145 : INFO : EPOCH - 1395 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,146 : WARNING : EPOCH - 1395 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,154 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,155 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,156 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,181 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,182 : INFO : EPOCH - 1396 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,184 : WARNING : EPOCH - 1396 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:40,193 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,194 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,195 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,219 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,220 : INFO : EPOCH - 1397 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,221 : WARNING : EPOCH - 1397 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,231 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,232 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,233 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,258 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,259 : INFO : EPOCH - 1398 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,259 : WARNING : EPOCH - 1398 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,269 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,271 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,271 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,299 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,304 : INFO : EPOCH - 1399 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,304 : WARNING : EPOCH - 1399 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,314 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,315 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,316 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,341 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,342 : INFO : EPOCH - 1400 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,343 : WARNING : EPOCH - 1400 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,352 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,353 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,353 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,378 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,379 : INFO : EPOCH - 1401 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,379 : WARNING : EPOCH - 1401 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,388 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,389 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,390 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,415 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,432 : INFO : EPOCH - 1402 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,433 : WARNING : EPOCH - 1402 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,442 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,443 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,444 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,469 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,470 : INFO : EPOCH - 1403 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,471 : WARNING : EPOCH - 1403 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,478 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,480 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,480 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,505 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,506 : INFO : EPOCH - 1404 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,506 : WARNING : EPOCH - 1404 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:40,514 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,524 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,524 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,541 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,542 : INFO : EPOCH - 1405 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,543 : WARNING : EPOCH - 1405 : supplied example count (1) did not equal ex
2019-02-21 17:17:40,550 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,552 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,552 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,578 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,579 : INFO : EPOCH - 1406 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,580 : WARNING : EPOCH - 1406 : supplied example count (1) did not equal ex
2019-02-21 17:17:40,589 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,590 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,590 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,615 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,616 : INFO : EPOCH - 1407 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,617 : WARNING : EPOCH - 1407 : supplied example count (1) did not equal ex
2019-02-21 17:17:40,625 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,627 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,628 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,652 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,653 : INFO : EPOCH - 1408 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,654 : WARNING : EPOCH - 1408 : supplied example count (1) did not equal ex
2019-02-21 17:17:40,665 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,665 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,666 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,691 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,691 : INFO : EPOCH - 1409 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,692 : WARNING : EPOCH - 1409 : supplied example count (1) did not equal ex
2019-02-21 17:17:40,705 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,706 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,707 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,731 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,732 : INFO : EPOCH - 1410 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,733 : WARNING : EPOCH - 1410 : supplied example count (1) did not equal ex
2019-02-21 17:17:40,741 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,743 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,743 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,769 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,770 : INFO : EPOCH - 1411 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,770 : WARNING : EPOCH - 1411 : supplied example count (1) did not equal ex
2019-02-21 17:17:40,779 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,780 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,781 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,806 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,807 : INFO : EPOCH - 1412 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,807 : WARNING : EPOCH - 1412 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:17:40,816 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,817 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,818 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,843 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,844 : INFO : EPOCH - 1413 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,844 : WARNING : EPOCH - 1413 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,855 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,856 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,856 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,882 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,883 : INFO : EPOCH - 1414 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,883 : WARNING : EPOCH - 1414 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,891 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,892 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,893 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,918 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,918 : INFO : EPOCH - 1415 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,919 : WARNING : EPOCH - 1415 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,926 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,928 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,929 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,954 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,955 : INFO : EPOCH - 1416 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,955 : WARNING : EPOCH - 1416 : supplied example count (1) did not equal exp
2019-02-21 17:17:40,963 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:40,965 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:40,965 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:40,990 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:40,991 : INFO : EPOCH - 1417 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:40,992 : WARNING : EPOCH - 1417 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,002 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,003 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,004 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,029 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,029 : INFO : EPOCH - 1418 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,030 : WARNING : EPOCH - 1418 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,039 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,040 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,041 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,066 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,067 : INFO : EPOCH - 1419 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,068 : WARNING : EPOCH - 1419 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,076 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,077 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,078 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,103 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,104 : INFO : EPOCH - 1420 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,105 : WARNING : EPOCH - 1420 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:41,114 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,115 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,116 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,140 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,142 : INFO : EPOCH - 1421 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,142 : WARNING : EPOCH - 1421 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,151 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,152 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,153 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,178 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,179 : INFO : EPOCH - 1422 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,180 : WARNING : EPOCH - 1422 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,191 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,192 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,192 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,217 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,218 : INFO : EPOCH - 1423 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,219 : WARNING : EPOCH - 1423 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,230 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,231 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,232 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,256 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,257 : INFO : EPOCH - 1424 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,258 : WARNING : EPOCH - 1424 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,267 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,268 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,269 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,293 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,293 : INFO : EPOCH - 1425 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,294 : WARNING : EPOCH - 1425 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,304 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,305 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,306 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,331 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,332 : INFO : EPOCH - 1426 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,333 : WARNING : EPOCH - 1426 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,341 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,351 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,352 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,369 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,370 : INFO : EPOCH - 1427 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,370 : WARNING : EPOCH - 1427 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,380 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,381 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,382 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,406 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,407 : INFO : EPOCH - 1428 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,408 : WARNING : EPOCH - 1428 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:41,419 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,420 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,421 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,447 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,448 : INFO : EPOCH - 1429 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,448 : WARNING : EPOCH - 1429 : supplied example count (1) did not equal ex
2019-02-21 17:17:41,460 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,461 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,462 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,487 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,488 : INFO : EPOCH - 1430 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,489 : WARNING : EPOCH - 1430 : supplied example count (1) did not equal ex
2019-02-21 17:17:41,498 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,499 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,499 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,524 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,525 : INFO : EPOCH - 1431 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,526 : WARNING : EPOCH - 1431 : supplied example count (1) did not equal ex
2019-02-21 17:17:41,534 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,536 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,536 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,561 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,562 : INFO : EPOCH - 1432 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,563 : WARNING : EPOCH - 1432 : supplied example count (1) did not equal ex
2019-02-21 17:17:41,572 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,573 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,573 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,599 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,600 : INFO : EPOCH - 1433 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,600 : WARNING : EPOCH - 1433 : supplied example count (1) did not equal ex
2019-02-21 17:17:41,613 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,614 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,615 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,639 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,640 : INFO : EPOCH - 1434 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,641 : WARNING : EPOCH - 1434 : supplied example count (1) did not equal ex
2019-02-21 17:17:41,650 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,652 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,653 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,676 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,678 : INFO : EPOCH - 1435 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,678 : WARNING : EPOCH - 1435 : supplied example count (1) did not equal ex
2019-02-21 17:17:41,688 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,689 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,690 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,715 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,716 : INFO : EPOCH - 1436 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,717 : WARNING : EPOCH - 1436 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:17:41,726 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,727 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,728 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,753 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,754 : INFO : EPOCH - 1437 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,755 : WARNING : EPOCH - 1437 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,763 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,765 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,765 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,790 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,792 : INFO : EPOCH - 1438 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,792 : WARNING : EPOCH - 1438 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,803 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,804 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,804 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,828 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,830 : INFO : EPOCH - 1439 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,830 : WARNING : EPOCH - 1439 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,839 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,840 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,841 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,866 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,867 : INFO : EPOCH - 1440 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,867 : WARNING : EPOCH - 1440 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,878 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,879 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,880 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,904 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,905 : INFO : EPOCH - 1441 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,906 : WARNING : EPOCH - 1441 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,915 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,916 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,917 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,941 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,942 : INFO : EPOCH - 1442 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,943 : WARNING : EPOCH - 1442 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,951 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,952 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,953 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:41,978 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:41,979 : INFO : EPOCH - 1443 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:41,980 : WARNING : EPOCH - 1443 : supplied example count (1) did not equal exp
2019-02-21 17:17:41,989 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:41,990 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:41,990 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,014 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,015 : INFO : EPOCH - 1444 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,016 : WARNING : EPOCH - 1444 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:42,025 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,026 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,026 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,051 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,052 : INFO : EPOCH - 1445 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,052 : WARNING : EPOCH - 1445 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,061 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,062 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,063 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,088 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,089 : INFO : EPOCH - 1446 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,090 : WARNING : EPOCH - 1446 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,098 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,100 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,101 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,126 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,127 : INFO : EPOCH - 1447 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,128 : WARNING : EPOCH - 1447 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,136 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,137 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,138 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,163 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,164 : INFO : EPOCH - 1448 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,164 : WARNING : EPOCH - 1448 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,173 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,174 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,175 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,201 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,202 : INFO : EPOCH - 1449 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,203 : WARNING : EPOCH - 1449 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,213 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,214 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,214 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,239 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,240 : INFO : EPOCH - 1450 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,241 : WARNING : EPOCH - 1450 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,250 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,251 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,252 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,276 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,277 : INFO : EPOCH - 1451 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,278 : WARNING : EPOCH - 1451 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,287 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,288 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,289 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,314 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,315 : INFO : EPOCH - 1452 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,315 : WARNING : EPOCH - 1452 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:42,323 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,324 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,325 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,350 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,351 : INFO : EPOCH - 1453 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,352 : WARNING : EPOCH - 1453 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,360 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,362 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,372 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,388 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,389 : INFO : EPOCH - 1454 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,390 : WARNING : EPOCH - 1454 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,398 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,399 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,400 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,424 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,425 : INFO : EPOCH - 1455 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,426 : WARNING : EPOCH - 1455 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,435 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,436 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,437 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,461 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,462 : INFO : EPOCH - 1456 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,463 : WARNING : EPOCH - 1456 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,473 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,473 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,498 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,499 : INFO : EPOCH - 1457 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,500 : WARNING : EPOCH - 1457 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,508 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,509 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,509 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,534 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,535 : INFO : EPOCH - 1458 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,536 : WARNING : EPOCH - 1458 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,544 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,545 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,546 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,572 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,583 : INFO : EPOCH - 1459 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,584 : WARNING : EPOCH - 1459 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,595 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,596 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,597 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,621 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,622 : INFO : EPOCH - 1460 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,623 : WARNING : EPOCH - 1460 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:42,637 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,638 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,639 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,664 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,665 : INFO : EPOCH - 1461 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,666 : WARNING : EPOCH - 1461 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,676 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,677 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,678 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,703 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,704 : INFO : EPOCH - 1462 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,705 : WARNING : EPOCH - 1462 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,722 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,739 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,740 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,741 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,741 : INFO : EPOCH - 1463 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,742 : WARNING : EPOCH - 1463 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,751 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,753 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,753 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,778 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,779 : INFO : EPOCH - 1464 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,780 : WARNING : EPOCH - 1464 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,789 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,789 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,790 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,814 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,815 : INFO : EPOCH - 1465 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,816 : WARNING : EPOCH - 1465 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,824 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,825 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,826 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,850 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,851 : INFO : EPOCH - 1466 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,852 : WARNING : EPOCH - 1466 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,861 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,862 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,863 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,887 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,888 : INFO : EPOCH - 1467 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,889 : WARNING : EPOCH - 1467 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,898 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,907 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,907 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,924 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,926 : INFO : EPOCH - 1468 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,926 : WARNING : EPOCH - 1468 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:42,935 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,936 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,937 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,962 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:42,963 : INFO : EPOCH - 1469 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:42,964 : WARNING : EPOCH - 1469 : supplied example count (1) did not equal exp
2019-02-21 17:17:42,972 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:42,973 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:42,974 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:42,999 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,000 : INFO : EPOCH - 1470 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,001 : WARNING : EPOCH - 1470 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,009 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,010 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,011 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,036 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,037 : INFO : EPOCH - 1471 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,038 : WARNING : EPOCH - 1471 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,046 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,047 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,048 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,073 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,074 : INFO : EPOCH - 1472 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,075 : WARNING : EPOCH - 1472 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,084 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,085 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,086 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,110 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,111 : INFO : EPOCH - 1473 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,112 : WARNING : EPOCH - 1473 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,121 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,122 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,123 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,147 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,148 : INFO : EPOCH - 1474 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,149 : WARNING : EPOCH - 1474 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,159 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,160 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,161 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,186 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,194 : INFO : EPOCH - 1475 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,195 : WARNING : EPOCH - 1475 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,205 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,207 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,208 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,232 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,234 : INFO : EPOCH - 1476 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,234 : WARNING : EPOCH - 1476 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:43,243 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,244 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,245 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,269 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,270 : INFO : EPOCH - 1477 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,271 : WARNING : EPOCH - 1477 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,280 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,281 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,281 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,306 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,307 : INFO : EPOCH - 1478 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,308 : WARNING : EPOCH - 1478 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,317 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,318 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,319 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,343 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,344 : INFO : EPOCH - 1479 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,345 : WARNING : EPOCH - 1479 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,353 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,355 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,355 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,380 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,381 : INFO : EPOCH - 1480 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,382 : WARNING : EPOCH - 1480 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,389 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,391 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,391 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,416 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,417 : INFO : EPOCH - 1481 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,418 : WARNING : EPOCH - 1481 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,427 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,429 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,430 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,454 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,455 : INFO : EPOCH - 1482 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,455 : WARNING : EPOCH - 1482 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,469 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,470 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,471 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,495 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,496 : INFO : EPOCH - 1483 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,497 : WARNING : EPOCH - 1483 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,507 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,508 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,509 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,535 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,536 : INFO : EPOCH - 1484 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,536 : WARNING : EPOCH - 1484 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:43,544 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,545 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,546 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,570 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,571 : INFO : EPOCH - 1485 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,572 : WARNING : EPOCH - 1485 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,581 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,582 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,583 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,607 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,608 : INFO : EPOCH - 1486 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,609 : WARNING : EPOCH - 1486 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,618 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,619 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,620 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,644 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,645 : INFO : EPOCH - 1487 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,646 : WARNING : EPOCH - 1487 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,654 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,656 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,657 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,682 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,683 : INFO : EPOCH - 1488 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,684 : WARNING : EPOCH - 1488 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,692 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,694 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,695 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,720 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,721 : INFO : EPOCH - 1489 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,721 : WARNING : EPOCH - 1489 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,732 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,739 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,740 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,758 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,759 : INFO : EPOCH - 1490 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,760 : WARNING : EPOCH - 1490 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,768 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,770 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,770 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,795 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,797 : INFO : EPOCH - 1491 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,797 : WARNING : EPOCH - 1491 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,807 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,808 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,809 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,834 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,835 : INFO : EPOCH - 1492 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,836 : WARNING : EPOCH - 1492 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:43,845 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,846 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,847 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,871 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,872 : INFO : EPOCH - 1493 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,873 : WARNING : EPOCH - 1493 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,881 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,882 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,883 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,908 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,909 : INFO : EPOCH - 1494 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,910 : WARNING : EPOCH - 1494 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,917 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,918 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,919 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,944 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,945 : INFO : EPOCH - 1495 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,945 : WARNING : EPOCH - 1495 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,954 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,956 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,956 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:43,982 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:43,983 : INFO : EPOCH - 1496 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:43,984 : WARNING : EPOCH - 1496 : supplied example count (1) did not equal exp
2019-02-21 17:17:43,993 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:43,995 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:43,996 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,021 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,022 : INFO : EPOCH - 1497 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,023 : WARNING : EPOCH - 1497 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,031 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,033 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,033 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,057 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,059 : INFO : EPOCH - 1498 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,059 : WARNING : EPOCH - 1498 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,068 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,069 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,070 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,094 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,095 : INFO : EPOCH - 1499 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,096 : WARNING : EPOCH - 1499 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,104 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,105 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,106 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,130 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,131 : INFO : EPOCH - 1500 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,132 : WARNING : EPOCH - 1500 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:44,139 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,141 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,142 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,166 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,167 : INFO : EPOCH - 1501 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,168 : WARNING : EPOCH - 1501 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,178 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,179 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,180 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,204 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,205 : INFO : EPOCH - 1502 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,206 : WARNING : EPOCH - 1502 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,214 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,215 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,216 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,240 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,241 : INFO : EPOCH - 1503 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,242 : WARNING : EPOCH - 1503 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,250 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,252 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,252 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,277 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,278 : INFO : EPOCH - 1504 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,279 : WARNING : EPOCH - 1504 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,287 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,288 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,288 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,313 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,314 : INFO : EPOCH - 1505 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,315 : WARNING : EPOCH - 1505 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,323 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,324 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,325 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,349 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,350 : INFO : EPOCH - 1506 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,351 : WARNING : EPOCH - 1506 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,359 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,360 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,361 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,386 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,387 : INFO : EPOCH - 1507 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,388 : WARNING : EPOCH - 1507 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,396 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,398 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,398 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,423 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,425 : INFO : EPOCH - 1508 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,425 : WARNING : EPOCH - 1508 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:44,433 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,435 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,436 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,461 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,462 : INFO : EPOCH - 1509 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,463 : WARNING : EPOCH - 1509 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,472 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,473 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,498 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,499 : INFO : EPOCH - 1510 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,499 : WARNING : EPOCH - 1510 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,507 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,508 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,509 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,534 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,535 : INFO : EPOCH - 1511 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,536 : WARNING : EPOCH - 1511 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,544 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,546 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,546 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,570 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,571 : INFO : EPOCH - 1512 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,572 : WARNING : EPOCH - 1512 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,581 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,582 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,583 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,608 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,609 : INFO : EPOCH - 1513 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,609 : WARNING : EPOCH - 1513 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,618 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,619 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,620 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,645 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,646 : INFO : EPOCH - 1514 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,647 : WARNING : EPOCH - 1514 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,655 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,656 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,657 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,681 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,682 : INFO : EPOCH - 1515 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,682 : WARNING : EPOCH - 1515 : supplied example count (1) did not equal exp
2019-02-21 17:17:44,691 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,692 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,692 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,717 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,718 : INFO : EPOCH - 1516 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,719 : WARNING : EPOCH - 1516 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:44,728 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,729 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,730 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,755 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,756 : INFO : EPOCH - 1517 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,757 : WARNING : EPOCH - 1517 : supplied example count (1) did not equal ex
2019-02-21 17:17:44,765 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,767 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,768 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,793 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,794 : INFO : EPOCH - 1518 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,795 : WARNING : EPOCH - 1518 : supplied example count (1) did not equal ex
2019-02-21 17:17:44,804 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,805 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,806 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,830 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,831 : INFO : EPOCH - 1519 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,831 : WARNING : EPOCH - 1519 : supplied example count (1) did not equal ex
2019-02-21 17:17:44,840 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,841 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,841 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,866 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,867 : INFO : EPOCH - 1520 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,868 : WARNING : EPOCH - 1520 : supplied example count (1) did not equal ex
2019-02-21 17:17:44,876 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,877 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,878 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,902 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,903 : INFO : EPOCH - 1521 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,904 : WARNING : EPOCH - 1521 : supplied example count (1) did not equal ex
2019-02-21 17:17:44,912 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,913 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,914 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,939 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,944 : INFO : EPOCH - 1522 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,959 : WARNING : EPOCH - 1522 : supplied example count (1) did not equal ex
2019-02-21 17:17:44,967 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:44,968 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:44,969 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:44,995 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:44,996 : INFO : EPOCH - 1523 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:44,997 : WARNING : EPOCH - 1523 : supplied example count (1) did not equal ex
2019-02-21 17:17:45,006 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,007 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,008 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,032 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,033 : INFO : EPOCH - 1524 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,034 : WARNING : EPOCH - 1524 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:17:45,041 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,042 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,043 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,068 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,093 : INFO : EPOCH - 1525 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,093 : WARNING : EPOCH - 1525 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,103 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,104 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,105 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,129 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,130 : INFO : EPOCH - 1526 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,131 : WARNING : EPOCH - 1526 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,138 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,139 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,140 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,165 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,166 : INFO : EPOCH - 1527 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,167 : WARNING : EPOCH - 1527 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,175 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,176 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,176 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,201 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,203 : INFO : EPOCH - 1528 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,203 : WARNING : EPOCH - 1528 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,212 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,213 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,214 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,239 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,240 : INFO : EPOCH - 1529 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,240 : WARNING : EPOCH - 1529 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,249 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,251 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,252 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,276 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,277 : INFO : EPOCH - 1530 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,278 : WARNING : EPOCH - 1530 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,286 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,288 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,288 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,314 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,315 : INFO : EPOCH - 1531 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,315 : WARNING : EPOCH - 1531 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,323 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,324 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,325 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,350 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,351 : INFO : EPOCH - 1532 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,351 : WARNING : EPOCH - 1532 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:45,359 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,360 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,361 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,386 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,387 : INFO : EPOCH - 1533 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,387 : WARNING : EPOCH - 1533 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,395 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,397 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,398 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,421 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,422 : INFO : EPOCH - 1534 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,423 : WARNING : EPOCH - 1534 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,431 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,433 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,433 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,457 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,458 : INFO : EPOCH - 1535 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,459 : WARNING : EPOCH - 1535 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,468 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,469 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,470 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,494 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,495 : INFO : EPOCH - 1536 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,496 : WARNING : EPOCH - 1536 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,505 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,506 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,506 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,532 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,533 : INFO : EPOCH - 1537 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,534 : WARNING : EPOCH - 1537 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,542 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,543 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,544 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,569 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,570 : INFO : EPOCH - 1538 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,571 : WARNING : EPOCH - 1538 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,580 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,581 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,581 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,606 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,607 : INFO : EPOCH - 1539 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,608 : WARNING : EPOCH - 1539 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,617 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,618 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,619 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,644 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,645 : INFO : EPOCH - 1540 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,646 : WARNING : EPOCH - 1540 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:45,654 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,655 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,656 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,681 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,682 : INFO : EPOCH - 1541 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,682 : WARNING : EPOCH - 1541 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,691 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,692 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,693 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,718 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,719 : INFO : EPOCH - 1542 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,719 : WARNING : EPOCH - 1542 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,728 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,729 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,730 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,754 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,755 : INFO : EPOCH - 1543 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,756 : WARNING : EPOCH - 1543 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,765 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,766 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,766 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,790 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,791 : INFO : EPOCH - 1544 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,794 : WARNING : EPOCH - 1544 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,803 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,804 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,805 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,830 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,837 : INFO : EPOCH - 1545 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,838 : WARNING : EPOCH - 1545 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,847 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,848 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,849 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,873 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,874 : INFO : EPOCH - 1546 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,875 : WARNING : EPOCH - 1546 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,884 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,886 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,886 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,911 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,912 : INFO : EPOCH - 1547 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,913 : WARNING : EPOCH - 1547 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,922 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,923 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,923 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,948 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,949 : INFO : EPOCH - 1548 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,949 : WARNING : EPOCH - 1548 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:45,958 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,958 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,959 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:45,984 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:45,985 : INFO : EPOCH - 1549 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:45,985 : WARNING : EPOCH - 1549 : supplied example count (1) did not equal exp
2019-02-21 17:17:45,994 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:45,995 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:45,996 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,020 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,021 : INFO : EPOCH - 1550 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,022 : WARNING : EPOCH - 1550 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,030 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,031 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,033 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,057 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,058 : INFO : EPOCH - 1551 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,058 : WARNING : EPOCH - 1551 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,067 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,069 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,069 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,094 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,096 : INFO : EPOCH - 1552 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,096 : WARNING : EPOCH - 1552 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,106 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,107 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,108 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,134 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,136 : INFO : EPOCH - 1553 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,137 : WARNING : EPOCH - 1553 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,146 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,147 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,148 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,171 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,173 : INFO : EPOCH - 1554 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,173 : WARNING : EPOCH - 1554 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,186 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,187 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,188 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,213 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,214 : INFO : EPOCH - 1555 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,215 : WARNING : EPOCH - 1555 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,224 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,225 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,225 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,250 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,251 : INFO : EPOCH - 1556 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,252 : WARNING : EPOCH - 1556 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:46,261 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,262 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,263 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,287 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,288 : INFO : EPOCH - 1557 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,289 : WARNING : EPOCH - 1557 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,298 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,299 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,300 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,324 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,325 : INFO : EPOCH - 1558 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,326 : WARNING : EPOCH - 1558 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,335 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,336 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,336 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,360 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,362 : INFO : EPOCH - 1559 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,362 : WARNING : EPOCH - 1559 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,371 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,372 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,373 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,398 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,399 : INFO : EPOCH - 1560 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,399 : WARNING : EPOCH - 1560 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,408 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,409 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,410 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,435 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,436 : INFO : EPOCH - 1561 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,436 : WARNING : EPOCH - 1561 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,444 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,446 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,447 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,470 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,471 : INFO : EPOCH - 1562 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,472 : WARNING : EPOCH - 1562 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,482 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,484 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,484 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,509 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,510 : INFO : EPOCH - 1563 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,511 : WARNING : EPOCH - 1563 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,519 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,520 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,521 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,545 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,546 : INFO : EPOCH - 1564 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,547 : WARNING : EPOCH - 1564 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:46,555 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,556 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,557 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,583 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,584 : INFO : EPOCH - 1565 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,584 : WARNING : EPOCH - 1565 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,593 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,594 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,595 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,619 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,620 : INFO : EPOCH - 1566 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,620 : WARNING : EPOCH - 1566 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,629 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,631 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,631 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,656 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,657 : INFO : EPOCH - 1567 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,658 : WARNING : EPOCH - 1567 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,666 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,667 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,668 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,693 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,694 : INFO : EPOCH - 1568 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,695 : WARNING : EPOCH - 1568 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,704 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,705 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,706 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,730 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,731 : INFO : EPOCH - 1569 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,732 : WARNING : EPOCH - 1569 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,740 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,742 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,743 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,768 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,769 : INFO : EPOCH - 1570 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,770 : WARNING : EPOCH - 1570 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,779 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,780 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,781 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,806 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,807 : INFO : EPOCH - 1571 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,808 : WARNING : EPOCH - 1571 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,817 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,818 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,818 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,843 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,844 : INFO : EPOCH - 1572 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,844 : WARNING : EPOCH - 1572 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:46,854 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,855 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,856 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,880 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,881 : INFO : EPOCH - 1573 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,882 : WARNING : EPOCH - 1573 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,890 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,892 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,892 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,916 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,917 : INFO : EPOCH - 1574 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,918 : WARNING : EPOCH - 1574 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,926 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,928 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,928 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,953 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,954 : INFO : EPOCH - 1575 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,954 : WARNING : EPOCH - 1575 : supplied example count (1) did not equal exp
2019-02-21 17:17:46,963 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:46,964 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:46,965 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:46,989 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:46,990 : INFO : EPOCH - 1576 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:46,991 : WARNING : EPOCH - 1576 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,000 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,001 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,002 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,025 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,026 : INFO : EPOCH - 1577 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,027 : WARNING : EPOCH - 1577 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,035 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,037 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,037 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,061 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,062 : INFO : EPOCH - 1578 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,063 : WARNING : EPOCH - 1578 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,072 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,084 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,085 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,098 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,099 : INFO : EPOCH - 1579 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,100 : WARNING : EPOCH - 1579 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,108 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,110 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,110 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,136 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,137 : INFO : EPOCH - 1580 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,138 : WARNING : EPOCH - 1580 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:47,147 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,148 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,148 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,173 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,174 : INFO : EPOCH - 1581 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,174 : WARNING : EPOCH - 1581 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,186 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,188 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,188 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,213 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,214 : INFO : EPOCH - 1582 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,214 : WARNING : EPOCH - 1582 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,223 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,224 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,224 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,248 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,250 : INFO : EPOCH - 1583 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,251 : WARNING : EPOCH - 1583 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,260 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,260 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,261 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,286 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,291 : INFO : EPOCH - 1584 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,304 : WARNING : EPOCH - 1584 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,316 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,317 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,317 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,342 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,350 : INFO : EPOCH - 1585 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,351 : WARNING : EPOCH - 1585 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,361 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,362 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,363 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,388 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,389 : INFO : EPOCH - 1586 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,389 : WARNING : EPOCH - 1586 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,398 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,399 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,400 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,426 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,441 : INFO : EPOCH - 1587 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,442 : WARNING : EPOCH - 1587 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,453 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,454 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,454 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,479 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,480 : INFO : EPOCH - 1588 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,480 : WARNING : EPOCH - 1588 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:47,489 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,490 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,490 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,515 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,516 : INFO : EPOCH - 1589 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,517 : WARNING : EPOCH - 1589 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,525 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,526 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,526 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,551 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,552 : INFO : EPOCH - 1590 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,553 : WARNING : EPOCH - 1590 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,561 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,562 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,562 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,587 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,588 : INFO : EPOCH - 1591 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,589 : WARNING : EPOCH - 1591 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,597 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,598 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,599 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,624 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,625 : INFO : EPOCH - 1592 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,625 : WARNING : EPOCH - 1592 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,636 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,637 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,638 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,661 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,676 : INFO : EPOCH - 1593 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,677 : WARNING : EPOCH - 1593 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,684 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,686 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,687 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,711 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,712 : INFO : EPOCH - 1594 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,713 : WARNING : EPOCH - 1594 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,722 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,723 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,724 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,748 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,749 : INFO : EPOCH - 1595 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,750 : WARNING : EPOCH - 1595 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,758 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,759 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,760 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,784 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,785 : INFO : EPOCH - 1596 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,786 : WARNING : EPOCH - 1596 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:47,794 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,795 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,796 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,820 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,822 : INFO : EPOCH - 1597 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,822 : WARNING : EPOCH - 1597 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,831 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,832 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,833 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,857 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,859 : INFO : EPOCH - 1598 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,859 : WARNING : EPOCH - 1598 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,867 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,869 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,870 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,894 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,895 : INFO : EPOCH - 1599 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,896 : WARNING : EPOCH - 1599 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,905 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,906 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,906 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,931 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,932 : INFO : EPOCH - 1600 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,933 : WARNING : EPOCH - 1600 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,941 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,942 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,943 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:47,969 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:47,970 : INFO : EPOCH - 1601 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:47,970 : WARNING : EPOCH - 1601 : supplied example count (1) did not equal exp
2019-02-21 17:17:47,979 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:47,980 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:47,981 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,005 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,006 : INFO : EPOCH - 1602 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,007 : WARNING : EPOCH - 1602 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,016 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,017 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,018 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,041 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,042 : INFO : EPOCH - 1603 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,043 : WARNING : EPOCH - 1603 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,052 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,053 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,053 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,078 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,079 : INFO : EPOCH - 1604 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,080 : WARNING : EPOCH - 1604 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:48,088 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,089 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,090 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,115 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,116 : INFO : EPOCH - 1605 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,116 : WARNING : EPOCH - 1605 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,125 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,126 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,127 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,151 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,152 : INFO : EPOCH - 1606 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,153 : WARNING : EPOCH - 1606 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,161 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,162 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,163 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,189 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,190 : INFO : EPOCH - 1607 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,191 : WARNING : EPOCH - 1607 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,204 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,205 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,205 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,231 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,237 : INFO : EPOCH - 1608 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,238 : WARNING : EPOCH - 1608 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,247 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,248 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,248 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,274 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,275 : INFO : EPOCH - 1609 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,275 : WARNING : EPOCH - 1609 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,285 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,286 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,286 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,310 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,312 : INFO : EPOCH - 1610 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,312 : WARNING : EPOCH - 1610 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,321 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,322 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,323 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,347 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,348 : INFO : EPOCH - 1611 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,349 : WARNING : EPOCH - 1611 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,357 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,358 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,358 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,383 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,384 : INFO : EPOCH - 1612 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,385 : WARNING : EPOCH - 1612 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:48,393 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,394 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,395 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,419 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,420 : INFO : EPOCH - 1613 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,421 : WARNING : EPOCH - 1613 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,429 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,430 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,431 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,456 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,459 : INFO : EPOCH - 1614 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,460 : WARNING : EPOCH - 1614 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,469 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,470 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,470 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,495 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,500 : INFO : EPOCH - 1615 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,501 : WARNING : EPOCH - 1615 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,510 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,511 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,512 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,537 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,538 : INFO : EPOCH - 1616 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,538 : WARNING : EPOCH - 1616 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,547 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,548 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,549 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,573 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,574 : INFO : EPOCH - 1617 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,575 : WARNING : EPOCH - 1617 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,585 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,587 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,587 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,613 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,614 : INFO : EPOCH - 1618 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,615 : WARNING : EPOCH - 1618 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,622 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,624 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,625 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,650 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,651 : INFO : EPOCH - 1619 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,652 : WARNING : EPOCH - 1619 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,661 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,662 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,662 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,686 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,687 : INFO : EPOCH - 1620 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,688 : WARNING : EPOCH - 1620 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:48,697 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,698 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,698 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,723 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,724 : INFO : EPOCH - 1621 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,725 : WARNING : EPOCH - 1621 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,737 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,738 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,739 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,764 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,774 : INFO : EPOCH - 1622 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,774 : WARNING : EPOCH - 1622 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,782 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,783 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,784 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,810 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,811 : INFO : EPOCH - 1623 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,811 : WARNING : EPOCH - 1623 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,821 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,822 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,822 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,847 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,848 : INFO : EPOCH - 1624 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,848 : WARNING : EPOCH - 1624 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,856 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,857 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,858 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,883 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,884 : INFO : EPOCH - 1625 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,884 : WARNING : EPOCH - 1625 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,892 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,893 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,895 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,919 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,920 : INFO : EPOCH - 1626 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,921 : WARNING : EPOCH - 1626 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,929 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,930 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,931 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,955 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,956 : INFO : EPOCH - 1627 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,957 : WARNING : EPOCH - 1627 : supplied example count (1) did not equal exp
2019-02-21 17:17:48,966 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:48,967 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:48,967 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:48,992 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:48,993 : INFO : EPOCH - 1628 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:48,993 : WARNING : EPOCH - 1628 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:49,002 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,003 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,004 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,029 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,030 : INFO : EPOCH - 1629 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,031 : WARNING : EPOCH - 1629 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,040 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,041 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,041 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,066 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,067 : INFO : EPOCH - 1630 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,068 : WARNING : EPOCH - 1630 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,079 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,080 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,081 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,105 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,106 : INFO : EPOCH - 1631 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,107 : WARNING : EPOCH - 1631 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,116 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,118 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,118 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,143 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,144 : INFO : EPOCH - 1632 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,145 : WARNING : EPOCH - 1632 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,152 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,154 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,155 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,180 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,180 : INFO : EPOCH - 1633 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,181 : WARNING : EPOCH - 1633 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,192 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,193 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,194 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,218 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,219 : INFO : EPOCH - 1634 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,220 : WARNING : EPOCH - 1634 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,229 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,229 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,230 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,255 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,256 : INFO : EPOCH - 1635 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,257 : WARNING : EPOCH - 1635 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,266 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,267 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,267 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,292 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,294 : INFO : EPOCH - 1636 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,295 : WARNING : EPOCH - 1636 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:49,304 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,305 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,305 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,331 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,332 : INFO : EPOCH - 1637 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,333 : WARNING : EPOCH - 1637 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,342 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,343 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,344 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,369 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,369 : INFO : EPOCH - 1638 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,370 : WARNING : EPOCH - 1638 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,379 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,380 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,381 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,406 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,407 : INFO : EPOCH - 1639 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,407 : WARNING : EPOCH - 1639 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,417 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,418 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,418 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,442 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,444 : INFO : EPOCH - 1640 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,444 : WARNING : EPOCH - 1640 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,453 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,454 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,455 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,480 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,481 : INFO : EPOCH - 1641 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,481 : WARNING : EPOCH - 1641 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,490 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,491 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,492 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,517 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,518 : INFO : EPOCH - 1642 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,519 : WARNING : EPOCH - 1642 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,526 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,527 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,528 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,553 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,554 : INFO : EPOCH - 1643 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,555 : WARNING : EPOCH - 1643 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,567 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,568 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,569 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,596 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,607 : INFO : EPOCH - 1644 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,608 : WARNING : EPOCH - 1644 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:49,620 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,622 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,623 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,648 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,667 : INFO : EPOCH - 1645 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,668 : WARNING : EPOCH - 1645 : supplied example count (1) did not equal ex
2019-02-21 17:17:49,677 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,679 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,680 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,704 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,705 : INFO : EPOCH - 1646 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,705 : WARNING : EPOCH - 1646 : supplied example count (1) did not equal ex
2019-02-21 17:17:49,714 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,716 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,717 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,741 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,742 : INFO : EPOCH - 1647 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,743 : WARNING : EPOCH - 1647 : supplied example count (1) did not equal ex
2019-02-21 17:17:49,751 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,752 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,752 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,776 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,777 : INFO : EPOCH - 1648 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,778 : WARNING : EPOCH - 1648 : supplied example count (1) did not equal ex
2019-02-21 17:17:49,786 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,787 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,788 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,812 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,813 : INFO : EPOCH - 1649 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,814 : WARNING : EPOCH - 1649 : supplied example count (1) did not equal ex
2019-02-21 17:17:49,822 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,823 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,824 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,848 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,849 : INFO : EPOCH - 1650 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,850 : WARNING : EPOCH - 1650 : supplied example count (1) did not equal ex
2019-02-21 17:17:49,859 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,860 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,861 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,888 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,889 : INFO : EPOCH - 1651 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,890 : WARNING : EPOCH - 1651 : supplied example count (1) did not equal ex
2019-02-21 17:17:49,898 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,904 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,906 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,924 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,925 : INFO : EPOCH - 1652 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,925 : WARNING : EPOCH - 1652 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:17:49,934 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,935 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,936 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:49,963 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:49,964 : INFO : EPOCH - 1653 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:49,965 : WARNING : EPOCH - 1653 : supplied example count (1) did not equal exp
2019-02-21 17:17:49,985 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:49,986 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:49,987 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,012 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,013 : INFO : EPOCH - 1654 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,013 : WARNING : EPOCH - 1654 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,021 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,023 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,024 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,049 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,050 : INFO : EPOCH - 1655 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,050 : WARNING : EPOCH - 1655 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,058 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,061 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,062 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,091 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,096 : INFO : EPOCH - 1656 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,097 : WARNING : EPOCH - 1656 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,138 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,148 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,149 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,168 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,169 : INFO : EPOCH - 1657 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,170 : WARNING : EPOCH - 1657 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,183 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,192 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,193 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,211 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,212 : INFO : EPOCH - 1658 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,213 : WARNING : EPOCH - 1658 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,222 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,223 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,223 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,249 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,250 : INFO : EPOCH - 1659 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,251 : WARNING : EPOCH - 1659 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,259 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,260 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,260 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,285 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,286 : INFO : EPOCH - 1660 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,287 : WARNING : EPOCH - 1660 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:50,295 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,297 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,297 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,322 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,323 : INFO : EPOCH - 1661 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,324 : WARNING : EPOCH - 1661 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,334 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,336 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,336 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,362 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,363 : INFO : EPOCH - 1662 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,364 : WARNING : EPOCH - 1662 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,372 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,373 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,374 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,399 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,400 : INFO : EPOCH - 1663 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,401 : WARNING : EPOCH - 1663 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,409 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,410 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,411 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,436 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,437 : INFO : EPOCH - 1664 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,438 : WARNING : EPOCH - 1664 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,446 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,448 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,448 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,474 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,475 : INFO : EPOCH - 1665 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,476 : WARNING : EPOCH - 1665 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,485 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,486 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,486 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,511 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,512 : INFO : EPOCH - 1666 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,513 : WARNING : EPOCH - 1666 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,521 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,522 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,523 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,548 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,549 : INFO : EPOCH - 1667 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,550 : WARNING : EPOCH - 1667 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,557 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,558 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,559 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,584 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,585 : INFO : EPOCH - 1668 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,586 : WARNING : EPOCH - 1668 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:50,594 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,595 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,596 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,622 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,623 : INFO : EPOCH - 1669 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,623 : WARNING : EPOCH - 1669 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,633 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,634 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,635 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,659 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,660 : INFO : EPOCH - 1670 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,661 : WARNING : EPOCH - 1670 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,670 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,671 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,671 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,697 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,698 : INFO : EPOCH - 1671 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,698 : WARNING : EPOCH - 1671 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,707 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,708 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,709 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,735 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,737 : INFO : EPOCH - 1672 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,738 : WARNING : EPOCH - 1672 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,746 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,747 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,748 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,774 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,775 : INFO : EPOCH - 1673 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,775 : WARNING : EPOCH - 1673 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,784 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,786 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,786 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,812 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,813 : INFO : EPOCH - 1674 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,814 : WARNING : EPOCH - 1674 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,822 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,823 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,824 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,849 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,850 : INFO : EPOCH - 1675 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,851 : WARNING : EPOCH - 1675 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,859 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,860 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,861 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,886 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,887 : INFO : EPOCH - 1676 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,888 : WARNING : EPOCH - 1676 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:50,896 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,898 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,899 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,924 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,925 : INFO : EPOCH - 1677 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,926 : WARNING : EPOCH - 1677 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,934 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,935 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,936 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:50,961 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:50,962 : INFO : EPOCH - 1678 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:50,963 : WARNING : EPOCH - 1678 : supplied example count (1) did not equal exp
2019-02-21 17:17:50,975 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:50,976 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:50,976 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,001 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,010 : INFO : EPOCH - 1679 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,011 : WARNING : EPOCH - 1679 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,019 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,019 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,020 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,046 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,047 : INFO : EPOCH - 1680 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,047 : WARNING : EPOCH - 1680 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,056 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,057 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,058 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,083 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,084 : INFO : EPOCH - 1681 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,084 : WARNING : EPOCH - 1681 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,093 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,094 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,095 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,119 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,120 : INFO : EPOCH - 1682 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,121 : WARNING : EPOCH - 1682 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,130 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,131 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,132 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,156 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,157 : INFO : EPOCH - 1683 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,158 : WARNING : EPOCH - 1683 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,167 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,168 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,169 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,196 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,196 : INFO : EPOCH - 1684 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,197 : WARNING : EPOCH - 1684 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:51,207 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,208 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,208 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,233 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,234 : INFO : EPOCH - 1685 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,235 : WARNING : EPOCH - 1685 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,247 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,248 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,248 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,273 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,283 : INFO : EPOCH - 1686 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,284 : WARNING : EPOCH - 1686 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,293 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,294 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,295 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,320 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,321 : INFO : EPOCH - 1687 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,321 : WARNING : EPOCH - 1687 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,329 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,331 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,332 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,356 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,357 : INFO : EPOCH - 1688 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,358 : WARNING : EPOCH - 1688 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,368 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,369 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,370 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,393 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,394 : INFO : EPOCH - 1689 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,395 : WARNING : EPOCH - 1689 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,403 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,405 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,405 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,430 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,431 : INFO : EPOCH - 1690 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,432 : WARNING : EPOCH - 1690 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,440 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,441 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,442 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,466 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,467 : INFO : EPOCH - 1691 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,468 : WARNING : EPOCH - 1691 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,477 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,478 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,479 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,503 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,504 : INFO : EPOCH - 1692 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,505 : WARNING : EPOCH - 1692 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:51,516 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,517 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,518 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,542 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,543 : INFO : EPOCH - 1693 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,544 : WARNING : EPOCH - 1693 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,553 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,554 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,555 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,580 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,581 : INFO : EPOCH - 1694 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,581 : WARNING : EPOCH - 1694 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,591 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,592 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,593 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,618 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,619 : INFO : EPOCH - 1695 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,619 : WARNING : EPOCH - 1695 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,630 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,631 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,632 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,657 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,658 : INFO : EPOCH - 1696 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,659 : WARNING : EPOCH - 1696 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,668 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,669 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,669 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,695 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,696 : INFO : EPOCH - 1697 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,696 : WARNING : EPOCH - 1697 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,705 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,706 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,707 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,731 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,732 : INFO : EPOCH - 1698 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,733 : WARNING : EPOCH - 1698 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,742 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,743 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,744 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,767 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,769 : INFO : EPOCH - 1699 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,769 : WARNING : EPOCH - 1699 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,780 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,781 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,782 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,807 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,817 : INFO : EPOCH - 1700 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,818 : WARNING : EPOCH - 1700 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:51,826 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,828 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,829 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,853 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,854 : INFO : EPOCH - 1701 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,855 : WARNING : EPOCH - 1701 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,864 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,865 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,866 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,890 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,891 : INFO : EPOCH - 1702 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,892 : WARNING : EPOCH - 1702 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,901 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,902 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,902 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,927 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,928 : INFO : EPOCH - 1703 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,929 : WARNING : EPOCH - 1703 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,936 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,938 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,939 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:51,963 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:51,964 : INFO : EPOCH - 1704 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:51,965 : WARNING : EPOCH - 1704 : supplied example count (1) did not equal exp
2019-02-21 17:17:51,974 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:51,975 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:51,976 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,001 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,002 : INFO : EPOCH - 1705 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,002 : WARNING : EPOCH - 1705 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,011 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,012 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,013 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,038 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,039 : INFO : EPOCH - 1706 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,040 : WARNING : EPOCH - 1706 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,052 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,053 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,053 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,078 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,090 : INFO : EPOCH - 1707 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,090 : WARNING : EPOCH - 1707 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,099 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,100 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,101 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,126 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,127 : INFO : EPOCH - 1708 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,128 : WARNING : EPOCH - 1708 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:52,136 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,137 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,138 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,162 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,191 : INFO : EPOCH - 1709 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,191 : WARNING : EPOCH - 1709 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,204 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,205 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,206 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,230 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,232 : INFO : EPOCH - 1710 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,233 : WARNING : EPOCH - 1710 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,241 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,242 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,244 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,267 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,268 : INFO : EPOCH - 1711 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,269 : WARNING : EPOCH - 1711 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,278 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,279 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,279 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,305 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,332 : INFO : EPOCH - 1712 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,333 : WARNING : EPOCH - 1712 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,341 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,353 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,354 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,367 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,368 : INFO : EPOCH - 1713 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,369 : WARNING : EPOCH - 1713 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,376 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,378 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,378 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,403 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,405 : INFO : EPOCH - 1714 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,405 : WARNING : EPOCH - 1714 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,413 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,414 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,415 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,440 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,441 : INFO : EPOCH - 1715 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,442 : WARNING : EPOCH - 1715 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,450 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,451 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,452 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,476 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,477 : INFO : EPOCH - 1716 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,478 : WARNING : EPOCH - 1716 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:52,486 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,487 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,488 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,513 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,514 : INFO : EPOCH - 1717 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,514 : WARNING : EPOCH - 1717 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,523 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,524 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,525 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,549 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,550 : INFO : EPOCH - 1718 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,551 : WARNING : EPOCH - 1718 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,559 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,560 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,561 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,586 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,587 : INFO : EPOCH - 1719 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,588 : WARNING : EPOCH - 1719 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,596 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,597 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,598 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,624 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,625 : INFO : EPOCH - 1720 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,626 : WARNING : EPOCH - 1720 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,635 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,636 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,636 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,662 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,663 : INFO : EPOCH - 1721 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,664 : WARNING : EPOCH - 1721 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,672 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,674 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,674 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,698 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,700 : INFO : EPOCH - 1722 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,700 : WARNING : EPOCH - 1722 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,709 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,710 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,710 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,735 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,736 : INFO : EPOCH - 1723 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,737 : WARNING : EPOCH - 1723 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,745 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,746 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,747 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,771 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,772 : INFO : EPOCH - 1724 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,773 : WARNING : EPOCH - 1724 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:52,781 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,783 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,783 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,807 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,809 : INFO : EPOCH - 1725 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,809 : WARNING : EPOCH - 1725 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,818 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,819 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,820 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,844 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,845 : INFO : EPOCH - 1726 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,846 : WARNING : EPOCH - 1726 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,857 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,858 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,859 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,883 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,884 : INFO : EPOCH - 1727 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,885 : WARNING : EPOCH - 1727 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,894 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,901 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,902 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,919 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,920 : INFO : EPOCH - 1728 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,921 : WARNING : EPOCH - 1728 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,930 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,931 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,932 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,956 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,957 : INFO : EPOCH - 1729 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,958 : WARNING : EPOCH - 1729 : supplied example count (1) did not equal exp
2019-02-21 17:17:52,966 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:52,967 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:52,968 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:52,992 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:52,993 : INFO : EPOCH - 1730 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:52,994 : WARNING : EPOCH - 1730 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,002 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,003 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,004 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,029 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,030 : INFO : EPOCH - 1731 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,030 : WARNING : EPOCH - 1731 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,039 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,040 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,041 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,065 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,066 : INFO : EPOCH - 1732 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,067 : WARNING : EPOCH - 1732 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:53,075 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,076 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,077 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,101 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,102 : INFO : EPOCH - 1733 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,103 : WARNING : EPOCH - 1733 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,111 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,112 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,113 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,138 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,139 : INFO : EPOCH - 1734 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,140 : WARNING : EPOCH - 1734 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,149 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,150 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,151 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,177 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,178 : INFO : EPOCH - 1735 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,179 : WARNING : EPOCH - 1735 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,191 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,192 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,192 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,217 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,218 : INFO : EPOCH - 1736 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,218 : WARNING : EPOCH - 1736 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,228 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,229 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,229 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,253 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,254 : INFO : EPOCH - 1737 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,255 : WARNING : EPOCH - 1737 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,262 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,263 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,264 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,288 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,289 : INFO : EPOCH - 1738 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,290 : WARNING : EPOCH - 1738 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,299 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,300 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,301 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,325 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,326 : INFO : EPOCH - 1739 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,326 : WARNING : EPOCH - 1739 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,335 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,336 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,337 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,361 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,362 : INFO : EPOCH - 1740 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,363 : WARNING : EPOCH - 1740 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:53,371 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,372 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,373 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,397 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,398 : INFO : EPOCH - 1741 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,399 : WARNING : EPOCH - 1741 : supplied example count (1) did not equal ex
2019-02-21 17:17:53,407 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,408 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,409 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,433 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,434 : INFO : EPOCH - 1742 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,435 : WARNING : EPOCH - 1742 : supplied example count (1) did not equal ex
2019-02-21 17:17:53,443 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,444 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,445 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,469 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,470 : INFO : EPOCH - 1743 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,471 : WARNING : EPOCH - 1743 : supplied example count (1) did not equal ex
2019-02-21 17:17:53,479 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,480 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,481 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,505 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,506 : INFO : EPOCH - 1744 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,507 : WARNING : EPOCH - 1744 : supplied example count (1) did not equal ex
2019-02-21 17:17:53,516 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,518 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,518 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,542 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,543 : INFO : EPOCH - 1745 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,544 : WARNING : EPOCH - 1745 : supplied example count (1) did not equal ex
2019-02-21 17:17:53,552 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,553 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,554 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,579 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,580 : INFO : EPOCH - 1746 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,580 : WARNING : EPOCH - 1746 : supplied example count (1) did not equal ex
2019-02-21 17:17:53,589 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,590 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,590 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,616 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,617 : INFO : EPOCH - 1747 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,618 : WARNING : EPOCH - 1747 : supplied example count (1) did not equal ex
2019-02-21 17:17:53,626 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,627 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,627 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,651 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,652 : INFO : EPOCH - 1748 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,653 : WARNING : EPOCH - 1748 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:17:53,661 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,663 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,663 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,688 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,689 : INFO : EPOCH - 1749 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,689 : WARNING : EPOCH - 1749 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,697 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,698 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,699 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,723 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,724 : INFO : EPOCH - 1750 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,725 : WARNING : EPOCH - 1750 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,733 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,735 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,736 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,760 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,761 : INFO : EPOCH - 1751 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,762 : WARNING : EPOCH - 1751 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,772 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,773 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,773 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,798 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,800 : INFO : EPOCH - 1752 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,801 : WARNING : EPOCH - 1752 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,809 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,810 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,811 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,835 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,840 : INFO : EPOCH - 1753 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,841 : WARNING : EPOCH - 1753 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,850 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,851 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,852 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,876 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,878 : INFO : EPOCH - 1754 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,879 : WARNING : EPOCH - 1754 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,887 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,888 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,889 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,915 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,916 : INFO : EPOCH - 1755 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,917 : WARNING : EPOCH - 1755 : supplied example count (1) did not equal exp
2019-02-21 17:17:53,926 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,927 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,927 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,953 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,954 : INFO : EPOCH - 1756 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,954 : WARNING : EPOCH - 1756 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:53,962 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:53,964 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:53,964 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:53,989 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:53,990 : INFO : EPOCH - 1757 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:53,991 : WARNING : EPOCH - 1757 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,000 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,001 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,002 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,025 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,026 : INFO : EPOCH - 1758 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,027 : WARNING : EPOCH - 1758 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,036 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,037 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,038 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,062 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,063 : INFO : EPOCH - 1759 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,064 : WARNING : EPOCH - 1759 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,075 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,076 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,077 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,102 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,103 : INFO : EPOCH - 1760 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,104 : WARNING : EPOCH - 1760 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,113 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,122 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,124 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,140 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,141 : INFO : EPOCH - 1761 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,142 : WARNING : EPOCH - 1761 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,150 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,151 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,152 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,178 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,179 : INFO : EPOCH - 1762 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,180 : WARNING : EPOCH - 1762 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,191 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,192 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,194 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,217 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,218 : INFO : EPOCH - 1763 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,219 : WARNING : EPOCH - 1763 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,228 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,229 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,230 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,256 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,257 : INFO : EPOCH - 1764 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,258 : WARNING : EPOCH - 1764 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:54,266 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,267 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,268 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,292 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,293 : INFO : EPOCH - 1765 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,294 : WARNING : EPOCH - 1765 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,304 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,305 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,305 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,330 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,331 : INFO : EPOCH - 1766 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,332 : WARNING : EPOCH - 1766 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,340 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,341 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,342 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,367 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,370 : INFO : EPOCH - 1767 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,371 : WARNING : EPOCH - 1767 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,379 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,381 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,381 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,406 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,407 : INFO : EPOCH - 1768 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,408 : WARNING : EPOCH - 1768 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,417 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,418 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,419 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,444 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,445 : INFO : EPOCH - 1769 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,446 : WARNING : EPOCH - 1769 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,454 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,456 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,457 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,480 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,481 : INFO : EPOCH - 1770 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,482 : WARNING : EPOCH - 1770 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,490 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,492 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,493 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,518 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,519 : INFO : EPOCH - 1771 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,519 : WARNING : EPOCH - 1771 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,558 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,560 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,561 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,586 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,587 : INFO : EPOCH - 1772 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,587 : WARNING : EPOCH - 1772 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:54,596 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,597 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,598 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,623 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,624 : INFO : EPOCH - 1773 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,625 : WARNING : EPOCH - 1773 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,644 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,674 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,675 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,676 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,677 : INFO : EPOCH - 1774 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,678 : WARNING : EPOCH - 1774 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,687 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,688 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,688 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,712 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,713 : INFO : EPOCH - 1775 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,714 : WARNING : EPOCH - 1775 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,722 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,723 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,724 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,749 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,750 : INFO : EPOCH - 1776 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,750 : WARNING : EPOCH - 1776 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,758 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,758 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,758 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,783 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,784 : INFO : EPOCH - 1777 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,785 : WARNING : EPOCH - 1777 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,792 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,793 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,794 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,819 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,821 : INFO : EPOCH - 1778 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,821 : WARNING : EPOCH - 1778 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,829 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,831 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,832 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,856 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,857 : INFO : EPOCH - 1779 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,858 : WARNING : EPOCH - 1779 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,868 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,869 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,869 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,894 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,895 : INFO : EPOCH - 1780 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,896 : WARNING : EPOCH - 1780 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:54,907 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,908 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,909 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,933 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,934 : INFO : EPOCH - 1781 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,935 : WARNING : EPOCH - 1781 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,949 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,951 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,951 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:54,976 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:54,977 : INFO : EPOCH - 1782 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:54,977 : WARNING : EPOCH - 1782 : supplied example count (1) did not equal exp
2019-02-21 17:17:54,988 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:54,989 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:54,990 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,015 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,016 : INFO : EPOCH - 1783 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,017 : WARNING : EPOCH - 1783 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,025 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,026 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,027 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,051 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,052 : INFO : EPOCH - 1784 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,053 : WARNING : EPOCH - 1784 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,061 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,062 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,063 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,087 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,088 : INFO : EPOCH - 1785 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,089 : WARNING : EPOCH - 1785 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,098 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,099 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,099 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,124 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,125 : INFO : EPOCH - 1786 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,126 : WARNING : EPOCH - 1786 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,134 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,135 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,135 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,160 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,161 : INFO : EPOCH - 1787 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,162 : WARNING : EPOCH - 1787 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,170 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,171 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,172 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,198 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,199 : INFO : EPOCH - 1788 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,200 : WARNING : EPOCH - 1788 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:55,209 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,211 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,211 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,237 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,238 : INFO : EPOCH - 1789 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,238 : WARNING : EPOCH - 1789 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,247 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,248 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,248 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,273 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,274 : INFO : EPOCH - 1790 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,275 : WARNING : EPOCH - 1790 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,285 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,286 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,286 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,312 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,313 : INFO : EPOCH - 1791 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,314 : WARNING : EPOCH - 1791 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,321 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,323 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,323 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,349 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,350 : INFO : EPOCH - 1792 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,350 : WARNING : EPOCH - 1792 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,357 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,359 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,359 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,385 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,386 : INFO : EPOCH - 1793 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,387 : WARNING : EPOCH - 1793 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,395 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,396 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,396 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,421 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,422 : INFO : EPOCH - 1794 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,422 : WARNING : EPOCH - 1794 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,432 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,433 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,433 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,458 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,458 : INFO : EPOCH - 1795 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,459 : WARNING : EPOCH - 1795 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,472 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,473 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,498 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,505 : INFO : EPOCH - 1796 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,506 : WARNING : EPOCH - 1796 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:55,514 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,515 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,516 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,540 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,541 : INFO : EPOCH - 1797 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,541 : WARNING : EPOCH - 1797 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,551 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,552 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,552 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,577 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,578 : INFO : EPOCH - 1798 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,579 : WARNING : EPOCH - 1798 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,587 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,589 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,590 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,615 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,616 : INFO : EPOCH - 1799 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,616 : WARNING : EPOCH - 1799 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,624 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,626 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,627 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,652 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,653 : INFO : EPOCH - 1800 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,653 : WARNING : EPOCH - 1800 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,661 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,663 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,663 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,689 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,690 : INFO : EPOCH - 1801 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,690 : WARNING : EPOCH - 1801 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,699 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,700 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,701 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,727 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,727 : INFO : EPOCH - 1802 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,728 : WARNING : EPOCH - 1802 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,740 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,741 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,742 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,767 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,775 : INFO : EPOCH - 1803 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,776 : WARNING : EPOCH - 1803 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,785 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,786 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,787 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,811 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,812 : INFO : EPOCH - 1804 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,813 : WARNING : EPOCH - 1804 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:55,821 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,822 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,823 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,848 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,849 : INFO : EPOCH - 1805 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,849 : WARNING : EPOCH - 1805 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,858 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,859 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,860 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,883 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,885 : INFO : EPOCH - 1806 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,885 : WARNING : EPOCH - 1806 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,893 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,895 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,896 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,920 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,921 : INFO : EPOCH - 1807 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,922 : WARNING : EPOCH - 1807 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,931 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,932 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,932 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,956 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,957 : INFO : EPOCH - 1808 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,958 : WARNING : EPOCH - 1808 : supplied example count (1) did not equal exp
2019-02-21 17:17:55,967 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:55,968 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:55,969 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:55,993 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:55,994 : INFO : EPOCH - 1809 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:55,994 : WARNING : EPOCH - 1809 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,005 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,006 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,007 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,032 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,044 : INFO : EPOCH - 1810 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,044 : WARNING : EPOCH - 1810 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,053 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,054 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,055 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,079 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,081 : INFO : EPOCH - 1811 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,081 : WARNING : EPOCH - 1811 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,089 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,090 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,091 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,116 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,117 : INFO : EPOCH - 1812 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,118 : WARNING : EPOCH - 1812 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:56,126 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,127 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,128 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,152 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,153 : INFO : EPOCH - 1813 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,153 : WARNING : EPOCH - 1813 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,162 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,163 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,164 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,191 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,192 : INFO : EPOCH - 1814 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,192 : WARNING : EPOCH - 1814 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,202 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,203 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,204 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,228 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,229 : INFO : EPOCH - 1815 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,230 : WARNING : EPOCH - 1815 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,238 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,239 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,240 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,264 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,265 : INFO : EPOCH - 1816 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,266 : WARNING : EPOCH - 1816 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,275 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,276 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,276 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,301 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,306 : INFO : EPOCH - 1817 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,307 : WARNING : EPOCH - 1817 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,315 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,316 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,317 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,342 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,343 : INFO : EPOCH - 1818 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,344 : WARNING : EPOCH - 1818 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,352 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,353 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,354 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,380 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,381 : INFO : EPOCH - 1819 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,381 : WARNING : EPOCH - 1819 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,390 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,391 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,392 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,416 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,417 : INFO : EPOCH - 1820 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,418 : WARNING : EPOCH - 1820 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:56,426 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,428 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,428 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,453 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,454 : INFO : EPOCH - 1821 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,455 : WARNING : EPOCH - 1821 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,463 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,464 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,465 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,490 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,491 : INFO : EPOCH - 1822 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,491 : WARNING : EPOCH - 1822 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,502 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,503 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,504 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,527 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,528 : INFO : EPOCH - 1823 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,529 : WARNING : EPOCH - 1823 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,540 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,541 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,542 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,567 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,568 : INFO : EPOCH - 1824 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,569 : WARNING : EPOCH - 1824 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,578 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,581 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,583 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,605 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,606 : INFO : EPOCH - 1825 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,607 : WARNING : EPOCH - 1825 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,619 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,620 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,620 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,644 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,645 : INFO : EPOCH - 1826 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,646 : WARNING : EPOCH - 1826 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,655 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,656 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,656 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,682 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,683 : INFO : EPOCH - 1827 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,683 : WARNING : EPOCH - 1827 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,692 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,693 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,694 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,719 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,720 : INFO : EPOCH - 1828 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,720 : WARNING : EPOCH - 1828 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:56,729 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,730 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,730 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,755 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,756 : INFO : EPOCH - 1829 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,773 : WARNING : EPOCH - 1829 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,792 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,793 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,794 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,819 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,820 : INFO : EPOCH - 1830 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,821 : WARNING : EPOCH - 1830 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,829 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,831 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,831 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,856 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,857 : INFO : EPOCH - 1831 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,858 : WARNING : EPOCH - 1831 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,867 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,868 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,869 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,910 : INFO : EPOCH - 1832 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,911 : WARNING : EPOCH - 1832 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,927 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,928 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,928 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,952 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,953 : INFO : EPOCH - 1833 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,954 : WARNING : EPOCH - 1833 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,963 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:56,964 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:56,965 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:56,988 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:56,989 : INFO : EPOCH - 1834 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:56,990 : WARNING : EPOCH - 1834 : supplied example count (1) did not equal exp
2019-02-21 17:17:56,998 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,000 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,001 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,024 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,025 : INFO : EPOCH - 1835 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,025 : WARNING : EPOCH - 1835 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,034 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,035 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,036 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,060 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,061 : INFO : EPOCH - 1836 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,062 : WARNING : EPOCH - 1836 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:57,072 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,074 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,075 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,099 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,099 : INFO : EPOCH - 1837 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,100 : WARNING : EPOCH - 1837 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,108 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,121 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,121 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,136 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,137 : INFO : EPOCH - 1838 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,138 : WARNING : EPOCH - 1838 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,146 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,147 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,148 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,173 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,174 : INFO : EPOCH - 1839 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,175 : WARNING : EPOCH - 1839 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,186 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,187 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,188 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,213 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,214 : INFO : EPOCH - 1840 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,214 : WARNING : EPOCH - 1840 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,224 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,225 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,225 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,250 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,251 : INFO : EPOCH - 1841 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,251 : WARNING : EPOCH - 1841 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,259 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,261 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,261 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,286 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,287 : INFO : EPOCH - 1842 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,288 : WARNING : EPOCH - 1842 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,297 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,297 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,298 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,323 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,324 : INFO : EPOCH - 1843 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,324 : WARNING : EPOCH - 1843 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,333 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,334 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,335 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,360 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,361 : INFO : EPOCH - 1844 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,362 : WARNING : EPOCH - 1844 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:57,371 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,372 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,372 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,397 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,398 : INFO : EPOCH - 1845 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,399 : WARNING : EPOCH - 1845 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,407 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,408 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,409 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,434 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,435 : INFO : EPOCH - 1846 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,435 : WARNING : EPOCH - 1846 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,444 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,445 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,446 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,471 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,472 : INFO : EPOCH - 1847 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,473 : WARNING : EPOCH - 1847 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,482 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,483 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,483 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,508 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,509 : INFO : EPOCH - 1848 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,510 : WARNING : EPOCH - 1848 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,518 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,519 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,520 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,545 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,546 : INFO : EPOCH - 1849 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,546 : WARNING : EPOCH - 1849 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,555 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,556 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,557 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,581 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,583 : INFO : EPOCH - 1850 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,583 : WARNING : EPOCH - 1850 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,592 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,592 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,593 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,618 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,619 : INFO : EPOCH - 1851 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,619 : WARNING : EPOCH - 1851 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,632 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,634 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,634 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,659 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,668 : INFO : EPOCH - 1852 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,669 : WARNING : EPOCH - 1852 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:57,678 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,679 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,680 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,704 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,705 : INFO : EPOCH - 1853 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,706 : WARNING : EPOCH - 1853 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,715 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,716 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,717 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,741 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,742 : INFO : EPOCH - 1854 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,743 : WARNING : EPOCH - 1854 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,751 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,752 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,752 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,777 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,778 : INFO : EPOCH - 1855 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,778 : WARNING : EPOCH - 1855 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,787 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,788 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,789 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,815 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,816 : INFO : EPOCH - 1856 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,816 : WARNING : EPOCH - 1856 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,825 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,826 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,826 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,851 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,852 : INFO : EPOCH - 1857 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,852 : WARNING : EPOCH - 1857 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,860 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,862 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,863 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,887 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,888 : INFO : EPOCH - 1858 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,889 : WARNING : EPOCH - 1858 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,899 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,900 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,901 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,926 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,931 : INFO : EPOCH - 1859 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,932 : WARNING : EPOCH - 1859 : supplied example count (1) did not equal exp
2019-02-21 17:17:57,941 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,942 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,942 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:57,967 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:57,968 : INFO : EPOCH - 1860 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:57,969 : WARNING : EPOCH - 1860 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:57,977 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:57,978 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:57,979 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,003 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,004 : INFO : EPOCH - 1861 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,005 : WARNING : EPOCH - 1861 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,014 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,015 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,016 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,041 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,042 : INFO : EPOCH - 1862 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,042 : WARNING : EPOCH - 1862 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,051 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,053 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,054 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,078 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,079 : INFO : EPOCH - 1863 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,079 : WARNING : EPOCH - 1863 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,088 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,089 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,090 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,114 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,115 : INFO : EPOCH - 1864 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,116 : WARNING : EPOCH - 1864 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,124 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,125 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,125 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,149 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,150 : INFO : EPOCH - 1865 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,151 : WARNING : EPOCH - 1865 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,159 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,160 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,161 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,187 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,188 : INFO : EPOCH - 1866 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,188 : WARNING : EPOCH - 1866 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,201 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,212 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,212 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,227 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,228 : INFO : EPOCH - 1867 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,229 : WARNING : EPOCH - 1867 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,237 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,238 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,239 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,263 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,264 : INFO : EPOCH - 1868 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,265 : WARNING : EPOCH - 1868 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:58,273 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,274 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,275 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,298 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,300 : INFO : EPOCH - 1869 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,300 : WARNING : EPOCH - 1869 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,308 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,310 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,311 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,335 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,336 : INFO : EPOCH - 1870 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,337 : WARNING : EPOCH - 1870 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,346 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,347 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,347 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,372 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,373 : INFO : EPOCH - 1871 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,374 : WARNING : EPOCH - 1871 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,382 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,383 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,384 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,410 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,411 : INFO : EPOCH - 1872 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,412 : WARNING : EPOCH - 1872 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,420 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,421 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,422 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,447 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,448 : INFO : EPOCH - 1873 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,448 : WARNING : EPOCH - 1873 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,457 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,457 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,458 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,484 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,485 : INFO : EPOCH - 1874 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,485 : WARNING : EPOCH - 1874 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,494 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,495 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,496 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,519 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,520 : INFO : EPOCH - 1875 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,521 : WARNING : EPOCH - 1875 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,530 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,531 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,533 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,556 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,557 : INFO : EPOCH - 1876 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,558 : WARNING : EPOCH - 1876 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:58,566 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,568 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,568 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,593 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,594 : INFO : EPOCH - 1877 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,595 : WARNING : EPOCH - 1877 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,604 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,605 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,605 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,630 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,631 : INFO : EPOCH - 1878 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,632 : WARNING : EPOCH - 1878 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,640 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,641 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,642 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,666 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,667 : INFO : EPOCH - 1879 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,668 : WARNING : EPOCH - 1879 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,676 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,677 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,678 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,703 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,704 : INFO : EPOCH - 1880 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,705 : WARNING : EPOCH - 1880 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,714 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,714 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,715 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,739 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,741 : INFO : EPOCH - 1881 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,742 : WARNING : EPOCH - 1881 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,751 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,753 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,753 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,778 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,779 : INFO : EPOCH - 1882 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,779 : WARNING : EPOCH - 1882 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,788 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,789 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,790 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,815 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,816 : INFO : EPOCH - 1883 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,817 : WARNING : EPOCH - 1883 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,825 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,826 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,827 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,851 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,852 : INFO : EPOCH - 1884 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,853 : WARNING : EPOCH - 1884 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:58,861 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,863 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,864 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,887 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,888 : INFO : EPOCH - 1885 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,889 : WARNING : EPOCH - 1885 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,898 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,899 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,900 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,924 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,925 : INFO : EPOCH - 1886 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,926 : WARNING : EPOCH - 1886 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,934 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,935 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,936 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,960 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,961 : INFO : EPOCH - 1887 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,962 : WARNING : EPOCH - 1887 : supplied example count (1) did not equal exp
2019-02-21 17:17:58,970 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:58,972 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:58,973 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:58,997 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:58,998 : INFO : EPOCH - 1888 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:58,999 : WARNING : EPOCH - 1888 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,007 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,009 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,010 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,034 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,035 : INFO : EPOCH - 1889 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,036 : WARNING : EPOCH - 1889 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,046 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,047 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,048 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,072 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,083 : INFO : EPOCH - 1890 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,102 : WARNING : EPOCH - 1890 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,115 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,116 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,116 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,141 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,142 : INFO : EPOCH - 1891 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,143 : WARNING : EPOCH - 1891 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,152 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,160 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,161 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,178 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,179 : INFO : EPOCH - 1892 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,179 : WARNING : EPOCH - 1892 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:59,189 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,190 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,191 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,217 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,219 : INFO : EPOCH - 1893 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,225 : WARNING : EPOCH - 1893 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,247 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,248 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,249 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,273 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,274 : INFO : EPOCH - 1894 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,275 : WARNING : EPOCH - 1894 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,283 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,284 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,284 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,308 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,310 : INFO : EPOCH - 1895 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,310 : WARNING : EPOCH - 1895 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,319 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,320 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,320 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,344 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,345 : INFO : EPOCH - 1896 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,346 : WARNING : EPOCH - 1896 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,354 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,355 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,356 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,381 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,383 : INFO : EPOCH - 1897 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,384 : WARNING : EPOCH - 1897 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,393 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,394 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,394 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,419 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,423 : INFO : EPOCH - 1898 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,424 : WARNING : EPOCH - 1898 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,432 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,433 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,434 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,459 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,460 : INFO : EPOCH - 1899 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,460 : WARNING : EPOCH - 1899 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,469 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,470 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,471 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,496 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,497 : INFO : EPOCH - 1900 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,498 : WARNING : EPOCH - 1900 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:59,506 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,507 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,508 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,533 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,534 : INFO : EPOCH - 1901 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,535 : WARNING : EPOCH - 1901 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,543 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,544 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,545 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,569 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,570 : INFO : EPOCH - 1902 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,571 : WARNING : EPOCH - 1902 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,579 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,580 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,581 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,605 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,606 : INFO : EPOCH - 1903 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,607 : WARNING : EPOCH - 1903 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,616 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,617 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,618 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,642 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,643 : INFO : EPOCH - 1904 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,644 : WARNING : EPOCH - 1904 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,654 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,655 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,656 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,681 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,682 : INFO : EPOCH - 1905 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,683 : WARNING : EPOCH - 1905 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,692 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,693 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,694 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,718 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,719 : INFO : EPOCH - 1906 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,720 : WARNING : EPOCH - 1906 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,728 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,730 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,730 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,755 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,756 : INFO : EPOCH - 1907 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,757 : WARNING : EPOCH - 1907 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,766 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,767 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,768 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,792 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,794 : INFO : EPOCH - 1908 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,794 : WARNING : EPOCH - 1908 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:17:59,804 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,805 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,806 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,831 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,832 : INFO : EPOCH - 1909 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,832 : WARNING : EPOCH - 1909 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,841 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,842 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,843 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,868 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,869 : INFO : EPOCH - 1910 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,870 : WARNING : EPOCH - 1910 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,879 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,880 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,881 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,905 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,908 : INFO : EPOCH - 1911 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,909 : WARNING : EPOCH - 1911 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,918 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,919 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,919 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,944 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,951 : INFO : EPOCH - 1912 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,952 : WARNING : EPOCH - 1912 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,960 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,962 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,963 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:17:59,988 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:17:59,989 : INFO : EPOCH - 1913 : training on 5386246 raw words (10000 effective
2019-02-21 17:17:59,989 : WARNING : EPOCH - 1913 : supplied example count (1) did not equal exp
2019-02-21 17:17:59,997 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:17:59,997 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:17:59,998 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,023 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,024 : INFO : EPOCH - 1914 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,025 : WARNING : EPOCH - 1914 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,035 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,035 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,036 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,060 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,061 : INFO : EPOCH - 1915 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,062 : WARNING : EPOCH - 1915 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,074 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,075 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,076 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,100 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,101 : INFO : EPOCH - 1916 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,102 : WARNING : EPOCH - 1916 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:00,110 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,112 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,112 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,137 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,138 : INFO : EPOCH - 1917 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,139 : WARNING : EPOCH - 1917 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,148 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,149 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,149 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,174 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,175 : INFO : EPOCH - 1918 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,175 : WARNING : EPOCH - 1918 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,187 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,188 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,189 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,214 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,215 : INFO : EPOCH - 1919 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,216 : WARNING : EPOCH - 1919 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,224 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,226 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,226 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,252 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,253 : INFO : EPOCH - 1920 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,253 : WARNING : EPOCH - 1920 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,262 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,263 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,263 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,288 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,289 : INFO : EPOCH - 1921 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,290 : WARNING : EPOCH - 1921 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,298 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,300 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,301 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,325 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,326 : INFO : EPOCH - 1922 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,327 : WARNING : EPOCH - 1922 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,335 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,336 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,337 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,362 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,363 : INFO : EPOCH - 1923 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,363 : WARNING : EPOCH - 1923 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,370 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,372 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,373 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,398 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,399 : INFO : EPOCH - 1924 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,399 : WARNING : EPOCH - 1924 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:00,407 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,409 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,410 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,434 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,435 : INFO : EPOCH - 1925 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,436 : WARNING : EPOCH - 1925 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,445 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,446 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,447 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,471 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,472 : INFO : EPOCH - 1926 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,473 : WARNING : EPOCH - 1926 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,481 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,483 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,484 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,509 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,510 : INFO : EPOCH - 1927 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,511 : WARNING : EPOCH - 1927 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,523 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,526 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,527 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,551 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,552 : INFO : EPOCH - 1928 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,552 : WARNING : EPOCH - 1928 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,561 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,562 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,563 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,588 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,589 : INFO : EPOCH - 1929 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,590 : WARNING : EPOCH - 1929 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,599 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,600 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,601 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,625 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,626 : INFO : EPOCH - 1930 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,627 : WARNING : EPOCH - 1930 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,639 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,640 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,640 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,664 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,665 : INFO : EPOCH - 1931 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,666 : WARNING : EPOCH - 1931 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,676 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,677 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,677 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,701 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,702 : INFO : EPOCH - 1932 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,702 : WARNING : EPOCH - 1932 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:00,710 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,712 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,712 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,736 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,737 : INFO : EPOCH - 1933 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,738 : WARNING : EPOCH - 1933 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,745 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,746 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,747 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,772 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,773 : INFO : EPOCH - 1934 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,773 : WARNING : EPOCH - 1934 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,782 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,783 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,784 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,808 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,810 : INFO : EPOCH - 1935 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,810 : WARNING : EPOCH - 1935 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,819 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,820 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,820 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,845 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,846 : INFO : EPOCH - 1936 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,847 : WARNING : EPOCH - 1936 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,855 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,856 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,856 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,881 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,882 : INFO : EPOCH - 1937 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,883 : WARNING : EPOCH - 1937 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,891 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,892 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,893 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,918 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,919 : INFO : EPOCH - 1938 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,919 : WARNING : EPOCH - 1938 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,928 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,929 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,930 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,954 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,955 : INFO : EPOCH - 1939 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,955 : WARNING : EPOCH - 1939 : supplied example count (1) did not equal exp
2019-02-21 17:18:00,964 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:00,965 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:00,966 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:00,990 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:00,991 : INFO : EPOCH - 1940 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:00,991 : WARNING : EPOCH - 1940 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:01,000 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,000 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,001 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,026 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,027 : INFO : EPOCH - 1941 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,027 : WARNING : EPOCH - 1941 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,036 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,037 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,037 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,063 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,064 : INFO : EPOCH - 1942 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,064 : WARNING : EPOCH - 1942 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,073 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,074 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,074 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,099 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,100 : INFO : EPOCH - 1943 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,100 : WARNING : EPOCH - 1943 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,109 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,110 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,110 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,134 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,135 : INFO : EPOCH - 1944 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,135 : WARNING : EPOCH - 1944 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,144 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,145 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,145 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,169 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,170 : INFO : EPOCH - 1945 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,171 : WARNING : EPOCH - 1945 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,182 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,185 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,185 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,210 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,210 : INFO : EPOCH - 1946 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,211 : WARNING : EPOCH - 1946 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,219 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,221 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,221 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,246 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,247 : INFO : EPOCH - 1947 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,247 : WARNING : EPOCH - 1947 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,255 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,257 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,257 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,287 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,293 : INFO : EPOCH - 1948 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,294 : WARNING : EPOCH - 1948 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:01,306 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,308 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,309 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,334 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,342 : INFO : EPOCH - 1949 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,342 : WARNING : EPOCH - 1949 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,351 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,352 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,353 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,378 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,379 : INFO : EPOCH - 1950 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,379 : WARNING : EPOCH - 1950 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,387 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,389 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,389 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,415 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,416 : INFO : EPOCH - 1951 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,416 : WARNING : EPOCH - 1951 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,450 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,453 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,454 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,455 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,455 : INFO : EPOCH - 1952 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,456 : WARNING : EPOCH - 1952 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,465 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,466 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,467 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,492 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,492 : INFO : EPOCH - 1953 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,492 : WARNING : EPOCH - 1953 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,504 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,505 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,506 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,532 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,533 : INFO : EPOCH - 1954 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,534 : WARNING : EPOCH - 1954 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,543 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,544 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,545 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,570 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,571 : INFO : EPOCH - 1955 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,572 : WARNING : EPOCH - 1955 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,594 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,607 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,613 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,613 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,614 : INFO : EPOCH - 1956 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,615 : WARNING : EPOCH - 1956 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:01,624 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,625 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,626 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,652 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,653 : INFO : EPOCH - 1957 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,654 : WARNING : EPOCH - 1957 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,667 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,668 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,669 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,694 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,695 : INFO : EPOCH - 1958 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,696 : WARNING : EPOCH - 1958 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,704 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,706 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,706 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,732 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,733 : INFO : EPOCH - 1959 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,734 : WARNING : EPOCH - 1959 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,742 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,744 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,744 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,769 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,770 : INFO : EPOCH - 1960 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,770 : WARNING : EPOCH - 1960 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,779 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,780 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,781 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,805 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,806 : INFO : EPOCH - 1961 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,807 : WARNING : EPOCH - 1961 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,815 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,817 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,817 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,842 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,843 : INFO : EPOCH - 1962 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,844 : WARNING : EPOCH - 1962 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,853 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,854 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,855 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,879 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,880 : INFO : EPOCH - 1963 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,881 : WARNING : EPOCH - 1963 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,893 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,894 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,895 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,921 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,922 : INFO : EPOCH - 1964 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,922 : WARNING : EPOCH - 1964 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:01,935 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,938 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,939 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:01,964 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:01,966 : INFO : EPOCH - 1965 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:01,967 : WARNING : EPOCH - 1965 : supplied example count (1) did not equal exp
2019-02-21 17:18:01,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:01,978 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:01,979 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,003 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,004 : INFO : EPOCH - 1966 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,005 : WARNING : EPOCH - 1966 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,016 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,017 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,018 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,044 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,045 : INFO : EPOCH - 1967 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,046 : WARNING : EPOCH - 1967 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,055 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,056 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,056 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,082 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,083 : INFO : EPOCH - 1968 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,084 : WARNING : EPOCH - 1968 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,093 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,094 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,095 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,119 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,120 : INFO : EPOCH - 1969 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,121 : WARNING : EPOCH - 1969 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,130 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,131 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,132 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,157 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,158 : INFO : EPOCH - 1970 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,158 : WARNING : EPOCH - 1970 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,167 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,168 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,169 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,195 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,196 : INFO : EPOCH - 1971 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,197 : WARNING : EPOCH - 1971 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,208 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,209 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,210 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,234 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,235 : INFO : EPOCH - 1972 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,236 : WARNING : EPOCH - 1972 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:02,244 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,245 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,246 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,271 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,272 : INFO : EPOCH - 1973 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,273 : WARNING : EPOCH - 1973 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,282 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,283 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,284 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,308 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,309 : INFO : EPOCH - 1974 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,310 : WARNING : EPOCH - 1974 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,319 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,320 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,321 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,345 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,346 : INFO : EPOCH - 1975 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,347 : WARNING : EPOCH - 1975 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,355 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,356 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,357 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,381 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,382 : INFO : EPOCH - 1976 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,383 : WARNING : EPOCH - 1976 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,391 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,392 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,393 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,417 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,418 : INFO : EPOCH - 1977 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,419 : WARNING : EPOCH - 1977 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,428 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,429 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,430 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,454 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,455 : INFO : EPOCH - 1978 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,456 : WARNING : EPOCH - 1978 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,467 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,468 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,468 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,493 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,502 : INFO : EPOCH - 1979 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,503 : WARNING : EPOCH - 1979 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,512 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,513 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,513 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,538 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,539 : INFO : EPOCH - 1980 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,540 : WARNING : EPOCH - 1980 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:02,550 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,551 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,552 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,578 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,579 : INFO : EPOCH - 1981 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,580 : WARNING : EPOCH - 1981 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,588 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,589 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,590 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,615 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,616 : INFO : EPOCH - 1982 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,617 : WARNING : EPOCH - 1982 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,626 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,627 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,628 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,660 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,670 : INFO : EPOCH - 1983 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,671 : WARNING : EPOCH - 1983 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,691 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,692 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,693 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,719 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,720 : INFO : EPOCH - 1984 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,721 : WARNING : EPOCH - 1984 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,730 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,731 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,732 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,757 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,758 : INFO : EPOCH - 1985 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,759 : WARNING : EPOCH - 1985 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,768 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,769 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,770 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,795 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,796 : INFO : EPOCH - 1986 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,797 : WARNING : EPOCH - 1986 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,806 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,807 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,808 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,832 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,833 : INFO : EPOCH - 1987 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,834 : WARNING : EPOCH - 1987 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,842 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,843 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,844 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,868 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,869 : INFO : EPOCH - 1988 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,870 : WARNING : EPOCH - 1988 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:02,881 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,882 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,882 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,907 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,908 : INFO : EPOCH - 1989 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,909 : WARNING : EPOCH - 1989 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,917 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,921 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,922 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,943 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,944 : INFO : EPOCH - 1990 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,944 : WARNING : EPOCH - 1990 : supplied example count (1) did not equal exp
2019-02-21 17:18:02,952 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:02,954 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:02,955 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:02,980 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:02,991 : INFO : EPOCH - 1991 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:02,992 : WARNING : EPOCH - 1991 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,000 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,002 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,002 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,027 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,028 : INFO : EPOCH - 1992 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,029 : WARNING : EPOCH - 1992 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,039 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,040 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,041 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,067 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,076 : INFO : EPOCH - 1993 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,077 : WARNING : EPOCH - 1993 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,094 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,095 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,096 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,121 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,122 : INFO : EPOCH - 1994 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,123 : WARNING : EPOCH - 1994 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,135 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,136 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,137 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,161 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,162 : INFO : EPOCH - 1995 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,163 : WARNING : EPOCH - 1995 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,173 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,179 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,179 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,200 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,205 : INFO : EPOCH - 1996 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,206 : WARNING : EPOCH - 1996 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:03,217 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,219 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,219 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,244 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,245 : INFO : EPOCH - 1997 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,246 : WARNING : EPOCH - 1997 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,255 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,256 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,256 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,281 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,282 : INFO : EPOCH - 1998 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,282 : WARNING : EPOCH - 1998 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,290 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,292 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,292 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,317 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,318 : INFO : EPOCH - 1999 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,319 : WARNING : EPOCH - 1999 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,327 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,328 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,329 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,354 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,355 : INFO : EPOCH - 2000 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,356 : WARNING : EPOCH - 2000 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,364 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,365 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,366 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,390 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,391 : INFO : EPOCH - 2001 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,392 : WARNING : EPOCH - 2001 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,400 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,401 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,402 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,426 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,427 : INFO : EPOCH - 2002 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,427 : WARNING : EPOCH - 2002 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,435 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,437 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,438 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,462 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,463 : INFO : EPOCH - 2003 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,464 : WARNING : EPOCH - 2003 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,471 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,473 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,473 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,498 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,499 : INFO : EPOCH - 2004 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,500 : WARNING : EPOCH - 2004 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:03,509 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,510 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,511 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,535 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,538 : INFO : EPOCH - 2005 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,539 : WARNING : EPOCH - 2005 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,550 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,551 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,552 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,576 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,580 : INFO : EPOCH - 2006 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,581 : WARNING : EPOCH - 2006 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,590 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,591 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,592 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,617 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,643 : INFO : EPOCH - 2007 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,644 : WARNING : EPOCH - 2007 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,652 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,653 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,654 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,679 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,681 : INFO : EPOCH - 2008 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,681 : WARNING : EPOCH - 2008 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,690 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,691 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,692 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,717 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,718 : INFO : EPOCH - 2009 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,718 : WARNING : EPOCH - 2009 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,730 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,731 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,731 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,755 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,756 : INFO : EPOCH - 2010 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,756 : WARNING : EPOCH - 2010 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,765 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,772 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,772 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,793 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,794 : INFO : EPOCH - 2011 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,794 : WARNING : EPOCH - 2011 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,804 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,812 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,813 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,831 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,832 : INFO : EPOCH - 2012 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,833 : WARNING : EPOCH - 2012 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:03,841 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,842 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,842 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,868 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,890 : INFO : EPOCH - 2013 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,891 : WARNING : EPOCH - 2013 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,902 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,903 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,904 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,929 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,937 : INFO : EPOCH - 2014 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,937 : WARNING : EPOCH - 2014 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,947 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,948 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,949 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:03,975 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:03,976 : INFO : EPOCH - 2015 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:03,976 : WARNING : EPOCH - 2015 : supplied example count (1) did not equal exp
2019-02-21 17:18:03,985 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:03,986 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:03,987 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,013 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,014 : INFO : EPOCH - 2016 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,014 : WARNING : EPOCH - 2016 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,025 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,026 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,027 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,052 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,053 : INFO : EPOCH - 2017 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,054 : WARNING : EPOCH - 2017 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,062 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,063 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,064 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,088 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,089 : INFO : EPOCH - 2018 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,090 : WARNING : EPOCH - 2018 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,099 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,100 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,100 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,124 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,126 : INFO : EPOCH - 2019 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,126 : WARNING : EPOCH - 2019 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,134 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,135 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,136 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,160 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,161 : INFO : EPOCH - 2020 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,162 : WARNING : EPOCH - 2020 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:04,173 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,174 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,175 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,199 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,200 : INFO : EPOCH - 2021 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,201 : WARNING : EPOCH - 2021 : supplied example count (1) did not equal ex
2019-02-21 17:18:04,210 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,212 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,214 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,237 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,238 : INFO : EPOCH - 2022 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,238 : WARNING : EPOCH - 2022 : supplied example count (1) did not equal ex
2019-02-21 17:18:04,249 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,250 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,251 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,276 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,276 : INFO : EPOCH - 2023 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,277 : WARNING : EPOCH - 2023 : supplied example count (1) did not equal ex
2019-02-21 17:18:04,286 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,287 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,288 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,312 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,313 : INFO : EPOCH - 2024 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,314 : WARNING : EPOCH - 2024 : supplied example count (1) did not equal ex
2019-02-21 17:18:04,322 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,323 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,323 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,348 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,349 : INFO : EPOCH - 2025 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,350 : WARNING : EPOCH - 2025 : supplied example count (1) did not equal ex
2019-02-21 17:18:04,357 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,359 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,359 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,384 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,385 : INFO : EPOCH - 2026 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,385 : WARNING : EPOCH - 2026 : supplied example count (1) did not equal ex
2019-02-21 17:18:04,394 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,395 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,395 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,419 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,420 : INFO : EPOCH - 2027 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,421 : WARNING : EPOCH - 2027 : supplied example count (1) did not equal ex
2019-02-21 17:18:04,431 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,432 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,432 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,457 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,458 : INFO : EPOCH - 2028 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,459 : WARNING : EPOCH - 2028 : supplied example count (1) did not equal ex
```

```
2019-02-21 17:18:04,468 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,468 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,469 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,493 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,496 : INFO : EPOCH - 2029 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,497 : WARNING : EPOCH - 2029 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,505 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,507 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,507 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,532 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,533 : INFO : EPOCH - 2030 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,534 : WARNING : EPOCH - 2030 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,543 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,544 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,545 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,570 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,571 : INFO : EPOCH - 2031 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,572 : WARNING : EPOCH - 2031 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,580 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,582 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,583 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,607 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,607 : INFO : EPOCH - 2032 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,608 : WARNING : EPOCH - 2032 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,617 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,618 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,619 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,644 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,649 : INFO : EPOCH - 2033 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,650 : WARNING : EPOCH - 2033 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,658 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,659 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,660 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,685 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,686 : INFO : EPOCH - 2034 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,686 : WARNING : EPOCH - 2034 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,697 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,699 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,699 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,724 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,736 : INFO : EPOCH - 2035 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,737 : WARNING : EPOCH - 2035 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,746 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,747 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,747 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,772 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,773 : INFO : EPOCH - 2036 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,774 : WARNING : EPOCH - 2036 : supplied example count (1) did not equal exp
```

```
2019-02-21 17:18:04,783 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,784 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,785 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,810 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,811 : INFO : EPOCH - 2037 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,811 : WARNING : EPOCH - 2037 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,820 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:18:04,821 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:18:04,822 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:18:04,846 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:18:04,847 : INFO : EPOCH - 2038 : training on 5386246 raw words (10000 effective
2019-02-21 17:18:04,848 : WARNING : EPOCH - 2038 : supplied example count (1) did not equal exp
2019-02-21 17:18:04,856 : INFO : worker thread finished; awaiting finish of 3 more threads
```

```
In [ ]: model.save(str("model/" + now.replace(".","-").replace(":","-") + ".model"))
```

## 4 Test

```
In [ ]: result = model.wv.most_similar(positive="microsoft",topn=10)
```

```
In [ ]: print(result)
```

## 5 Save

```
In [ ]: import save_notebook
```

```
In [ ]: name_notebook_exported = save_notebook.save_notebook("word2vec_with_gensim.ipynb")
```

```
In [ ]: def write_result(word_embedding, time_training,name_notebook_exported, fname ):
            if not os.path.isfile(fname):
                f=open(fname, "a+")
                f.write("Nombre de mots;Model de word embedding;temps d'apprentissage;Notebook
                f.close
            f=open(fname, "a+")
            f.write("\n" + str( str(len(wiki_vocab_tokenized)) + ";" +word_embedding + ";"+ str
            f.close
```

```
In [ ]: write_result(word_embedding, time_training,name_notebook_exported, "resultats.csv" )
```

```
In [ ]: name_notebook_exported
```