# HW3: Neural Machine Translation

Jiafeng Chen     Yufeng Ling     Francisco Rivera

March 7, 2019

## 1   Introduction

In this writeup we consider neural machine translation—given a source sentence in one language, can we generate a target sentence in another? We work with the prevailing encoder-decoder architecture, where one neural network, the *encoder*, maps the source sentence into a numerical tensor—which is supposed to be a latent representation of the source sentence. Another network, the *decoder*, is a language model that takes the encoded sentence as an input, and generates a sentence in the target language.

## 2   Problem Description

In this writeup, we consider the problem of machine translation. Let

$$\boldsymbol{x}_i = [x_{i1}, \dots, x_{iS}] \tag{1}$$

be a source sentence, where each word belongs to some source vocabulary, $x_{it} \in \mathcal{V}_s$. Let

$$\boldsymbol{y}_i = [y_{i1}, \dots, y_{iT}] \tag{2}$$

be a target sentence, with each target word belonging to some target vocabulary, $y_{it} \in \mathcal{V}_t$. Our goal is to learn $p(\boldsymbol{y}_i \mid \boldsymbol{x}_i)$. We often treat the prediction like a language model—i.e. we consider the sequential conditional distributions $p(y_{it} \mid y_{i1}, \dots, y_{i,t-1}, \boldsymbol{x}_i)$.

# 3 Model and Algorithms

## 3.1 Sequence-to-Sequence (Seq2Seq)

In Seq2Seq (Sutskever et al. (2014)), we have a encoder-decoder network architecture. The *encoder*, in the vanilla Seq2Seq implementation, is an LSTM network that takes an embedding of the source sentence $\text{embed}(x_{i1}), \dots, \text{embed}(x_{iS})$ and outputs a list of hidden states $h_{i1}, \dots, h_{iS}$ and cell state $c_{iS}$.

In Seq2Seq, the decoder is another LSTM which is initialized with $h_{iS}$ and $c_{iS}$. At training time, the decoder LSTM takes embeddings of the ground truth target $\text{embed}(y_{it})$ and output probability predictions for $y_{i,t+1}$. These predictions are penalized with the usual cross entropy loss at training time. At prediction time, the decoder LSTM gets passed the start-of-sentence token <s> and outputs predictions for the first word. We then iteratively pass in its top predictions to obtain predictions for future words. For instance, a greedy algorithm to generate a sentence would be taking the top prediction every time and pass the predicted sentence so far into the decoder to obtain the next word—stopping when the end-of-sentence token </s> is the top prediction.[1]

## 3.2 Attention

The attention model (Vaswani et al. (2017)) builds upon the Sequence-to-Sequence model described in section (3.1). Instead of only using the end state of the encoded input as input, we take advantage of a bi-directional RNN and generate target sentence by taking in a weighted input of all the states of the encoder. The intuition is that when we are translating a sentence, the adjacent words of the current word also help inform how we should translate the current word. Such weighted average is aptly named "context". To formalize this, we let the sequential conditional distribution be given by

$$p(y_{it} \mid y_{i1}, \dots, y_{i,t-1}, x_i) = g(y_{i,t-1}, s_t, c_t) \tag{3}$$

---

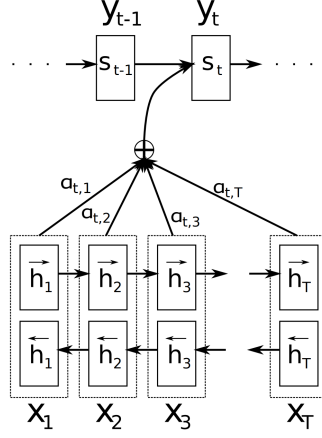[1]This corresponds to beam search with beam size 1.

*Figure 1: Attention Diagram*

where $s_t$ is the hidden state from BiRNN at time $t$ defined by

$$s_t = f(s_{t-1}, y_{i,t-1}, c_t). \tag{4}$$

As mentioned above, the context vector $c_t$ is a weighted average given by

$$c_t = \sum_{j=1}^{T_x} \alpha_{tj} h_j \tag{5}$$

where the annotations $\{h_j\}$ are the concatenated hidden states of the BiRNN and $\{\alpha_{tj}\}$ are the weights, which usually center around the current state $t$. The weights

$$\alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{k=1}^{T_x} \exp(e_t k)} \tag{6}$$

are generated by softmaxing over

$$e_{tk} = a(s_{t-1}, h_k) \tag{7}$$

for some function $a$. In our model, we take $a$ to be the dot product between the two arguments.

## 3.3 Beam Search

After we have trained a model, since the outputs are conditional probabilities over the entire vocabulary, we need to have an algorithm that generates the top $k$ most likely sentences. One can use the greedy algorithm by selecting the local MLE conditional on the previously selected words. However, this can easily run into problems like stucking at a local optimum.

Beam search is proposed based on this idea. Instead of selecting the MLE, we select the top $B$ words with the largest conditional probabilities. The parameter $B$ is known as the "beam width". At every step, we choose the top $B$ words with the highest joint probability conditional on the $B$ previously generated sentenses $\{y_{i,1:t-1}^{(b)}\}_{b=1}^{B}$. That is

$$\arg\max_{y_{it},b} p(y_{it}, y_{i,t-1}^{(b)}, \ldots, y_{i1}^{(b)}, x_i) = \arg\max_{y_{it},b} p(y_{it} \mid y_{i,t-1}^{(b)}, \ldots, y_{i1}^{(b)}, x_i) p(y_{i,t-1}^{(b)}, \ldots, y_{i1}^{(b)}, x_i).$$

In the last step we choose the top $k$ to report as results.

## 3.4 Transformer extension

We also implemented the *transformer* architecture, following Vaswani et al. (2017). Each transformer layer is a self-attention layer,[2] a source-attention layer (only in the decoder), and a fully-connected feedforward network, connected by residual connections[3] and layer normalization.[4] We transform the input via a dense embedding and a *positional embedding*—treating position in a sentence as tokens in a vocabulary, and map to a dense embedding. We concatenate the embeddings and feed into the transformer layers. The model architecture is diagramed in Figure 2.

We trained two transformer models, which we name `Megatron` and `Soundwave`, after beloved *Transformers* antagonists. The `size` parameter is the number of dimensions of the dense embeddings, input to layers, and output from layers (so

---

[2]The attention layers are *multiheaded*, which means that we transform the (query, key, value) inputs by a linear layer before using dot-product attention (and normalization by $\sqrt{d}$. We then concatenate the resulting context vectors and feed into one last linear layer).

[3]A residual connection of a layer adds the input to the layer output. This helps with gradient propagation

[4]Layer normalization standardizes a layer's output along dimension of the dense embedding (i.e. the dimension other than the batch and seqlen dimensions).
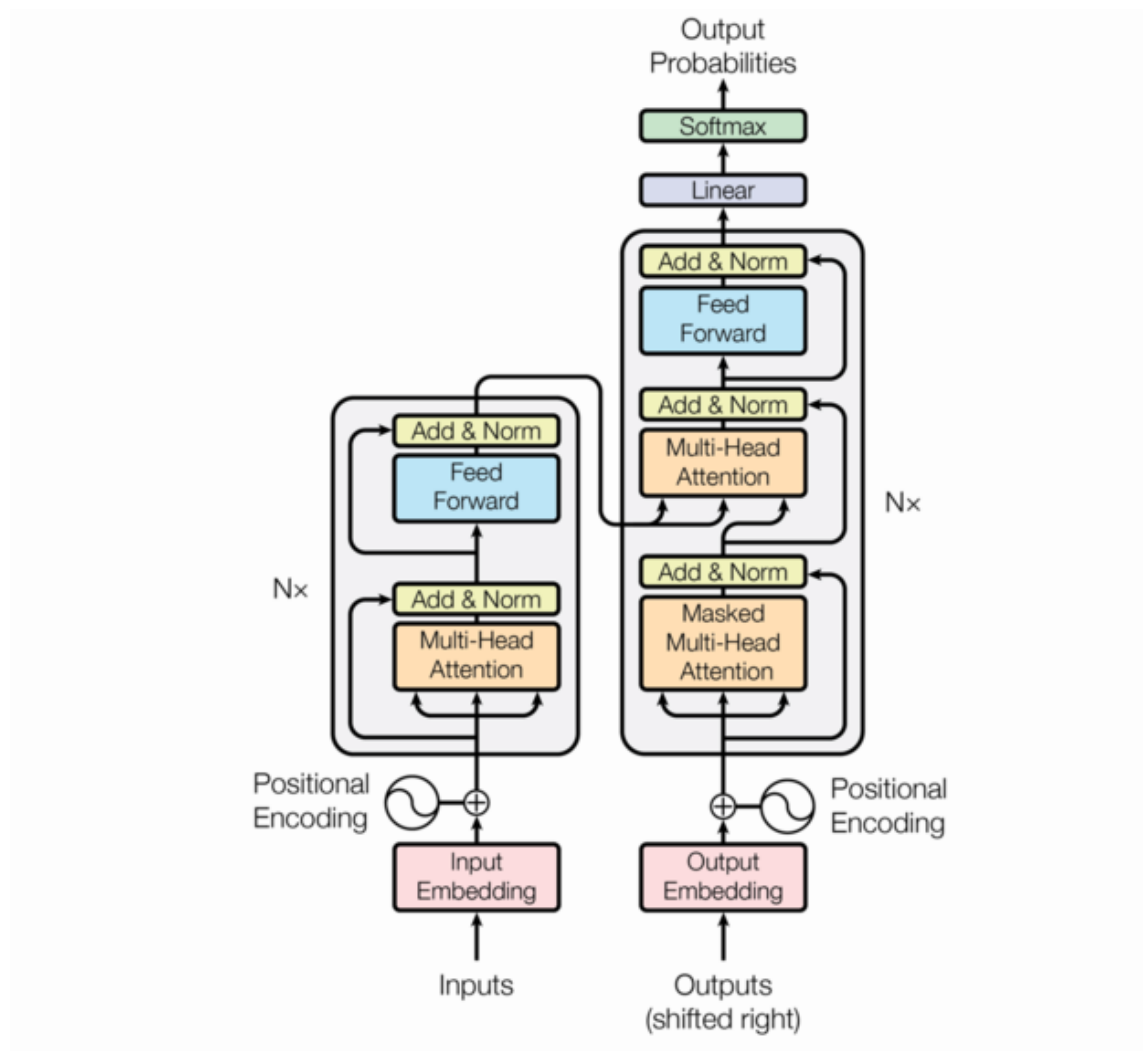
*Figure 2: Diagram in Vaswani et al. (2017)*

that the layers plug into each other). The pos, emb parameters are the number of dimensions in the positional and the word embedding

```
# Megatron
size = 200
pos = 20
emb = 180
head = 4
nlayers = 3
dropout = .2
```

```
# Soundwave
size = 200
pos = 20
emb = 180
head = 3
nlayers = 4
dropout = .1
```

# 4    Experiments

We trained two models and document hyperparameters in Table 1. Results are
displayed in Table 2. The performance of the Sequence-to-Sequence model in
Sutskever et al. (2014) is reasonable, but understandabely weaker due to its sim-
ple structure of combining two LSTMs without any attention layer applied to the
source sentence.

In the attention model proposed in Bahdanau et al. (2014), we used 1/5 of the
number of parameters compared to the paper due to the constraint of training
speed. The main difference is that instead of a 5-layer decoder, we used one layer.
Compared to the previous Sequence-to-Sequence model without attention, the
performance improved significantly as seen in Table 2. We experimented normal-
izing the log-attention weights by factor of $\sqrt{\text{embedding size}}$ before proceeding
to the softmax function to calculate the attention weights; this is because with
more dimensions, dot product tends to be larger, and the softmax tends to be
sharper, which would result in small gradients. However, this resulted in uni-
form weights across all source words and very poor performance. In the properly
trained model, we can observe the attention weights reflecting our intuition. See
Figure 3 and 4 for visualization of the attention weights.

It is also somewhat surprising to us that transformer model in Vaswani et al.
(2017) does not perform well relative to LSTM-based models, contradicting es-
tablished wisdom. We suspect that the reason is due to either poor implemen-
tation/optimization or short sentences and insufficient computing power. At-
tention is powerful when sentences are long, since it allows the model to look
at arbitrarily long dependencies, and attention scales well with GPU hardware

```
Ich  verlor  all  meine  Hoffnung  .  <pad>

<s>  I  lost  all  my  hope  .  </s>

Ich  verlor  all  meine  Hoffnung  .  <pad>

<s>  I  lost  all  my  hope  .  </s>

Ich  verlor  all  meine  Hoffnung  .  <pad>

<s>  I  lost  all  my  hope  .  </s>
```

*Figure 3: Darker background words signify more attention being placed on that word at each point in the translation process.*

and multi-GPU training—we would get neither of these advantages with a small dataset and Google Colab computing power (Training Megatron takes 15 minutes per epoch and training Soundwave takes over 20 minutes, which does not leave room for hyperparameter tuning).
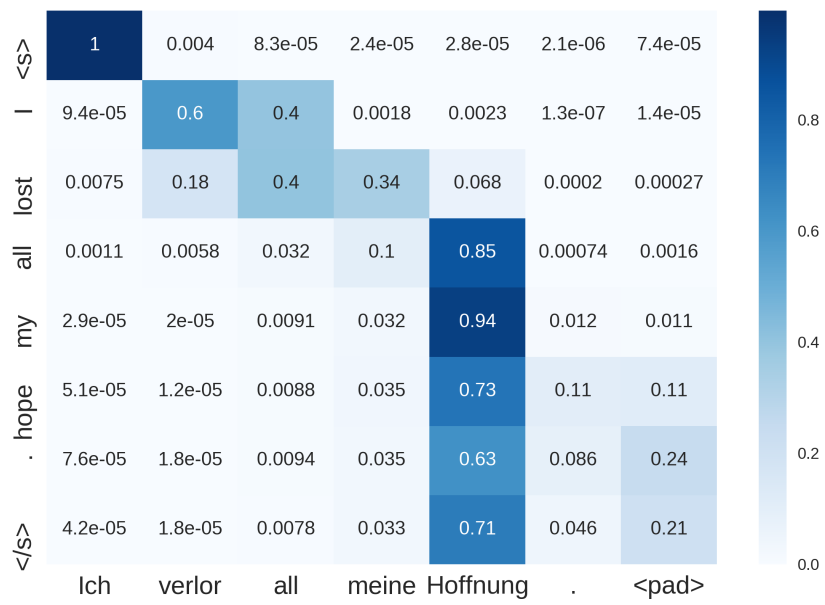
*Figure 4: Attention matrix weights*

| model name | specifications |
|---|---|
| Seq2Seq | 100 embedding, 100 hidden, no dropout |
| Seq2SeqAttn | Bi-directional LSTM, 100 embedding, 100 hidden, no dropout |

*Table 1: Both models used only one layer in the decoder and trained with Adam and learning rate $10^{-3}$ over 10 epochs. For* Seq2Seq, *decrease in validation loss flattened at epoch 8. For* Seq2SeqAttn, *training stopped at epoch 7 after loss starting going up.*

| model name | Loss | Perplexity | BLEU |
|---|---|---|---|
| Seq2Seq | 3.23 | 25.36 | 6.32 |
| Seq2SeqAttn | 2.71 | 15.05 | 17.77 |
| Megatron | 2.7 | 15 | 9.68 |
| Soundwave | 2.8 | 16 | 7.08 |

*Table 2: Performance metrics for different models*

# References

Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.

# A   Model implementation

*Listing 1: Sequence-to-Sequence*

```python
from namedtensor import ntorch
from namedtensor.nn import nn as nnn

class Seq2Seq(nnn.Module):
    def __init__(self, encoder, decoder):
        super().__init__()
        self.encoder = encoder
        self.decoder = decoder
    def forward(self, src, trg):
        encoded = self.encoder(src)
        return self.decoder(encoded, trg)


ce_loss = nnn.CrossEntropyLoss().spec("classes")

# def seq2seq_loss_fn(model, batch):
# length = batch.trg.size('trgSeqlen')
# out = model(batch.src, batch.trg)[{'trgSeqlen': slice(0, length - 1)}]
# trg = batch.trg[{'trgSeqlen': slice(1, length)}]
# mask = (trg != pad_idx_EN).float()

# n = trg.values.numel()
# loss = ce_loss(mask * out, trg)

# return (loss * n - np.log(V_EN) * (n - mask.sum().item())) / mask.sum().item()
```

*Listing 2: LSTM encoder*

```python
from namedtensor import ntorch
from namedtensor.nn import nn as nnn
```

```
3
4  class LSTMEncoder(nnn.Module):
5      def __init__(self, TEXT,
6                   embedding_dim=100,
7                   hidden_dim=150,
8                   num_layers=1,
9                   dropout=0,
10                  history=False,
11                  bidirectional=False):
12          super().__init__()
13          self.history = history
14          pad_idx = TEXT.vocab.stoi['<pad>']
15
16          self.embed = nnn.Embedding(num_embeddings=len(TEXT.vocab),
17                                     embedding_dim=embedding_dim,
18                                     padding_idx=pad_idx)
19
20          self.lstm = nnn.LSTM(input_size=embedding_dim,
21                               hidden_size=hidden_dim,
22                               num_layers=num_layers,
23                               dropout=dropout,
24                               bidirectional=bidirectional) \
25                          .spec("embedding", "srcSeqlen")
26
27      def forward(self, batch_text):
28          embedded = self.embed(batch_text)
29          hidden_states, last_state = self.lstm(embedded)
30          if self.history:
31              return hidden_states, last_state
32          else:
33              return last_state
```

*Listing 3: LSTM decoder*

```
1  import torch
2  from torch import nn
3  from namedtensor import ntorch, NamedTensor
4  from namedtensor.nn import nn as nnn
5  from attention import Attention
6  from utils import *
7
8  device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
9
10
11 class LSTMDecoder(nnn.Module):
12     def __init__(self, TEXT,
13                  embedding_dim=100,
14                  hidden_dim=150,
15                  num_layers=1,
16                  dropout=0):
17         super().__init__()
18
19         pad_idx = TEXT.vocab.stoi['<pad>']
20
21         self.embed = nnn.Embedding(num_embeddings=len(TEXT.vocab),
22                                    embedding_dim=embedding_dim,
23                                    padding_idx=pad_idx)
24
25         self.lstm = nnn.LSTM(input_size=embedding_dim,
```

```python
26                              hidden_size=hidden_dim,
27                              num_layers=num_layers,
28                              dropout=dropout) \
29                  .spec("embedding", "trgSeqlen")
30
31          self.w = nnn.Linear(in_features=hidden_dim,
32                          out_features=len(TEXT.vocab)) \
33                  .spec("embedding", "classes")
34
35      def forward(self, init_state, batch_text):
36          embedded = self.embed(batch_text)
37          hidden_states, (hn, cn) = self.lstm(embedded, init_state)
38          log_probs = self.w(hidden_states)
39          return log_probs
40
41
42  class LSTMDecoderAttn(nnn.Module):
43      def __init__(self, TEXT,
44                  embedding_dim=100,
45                  hidden_dim=150,
46                  num_layers=1,
47                  dropout=0,
48                  attn_normalize=True):
49          super().__init__()
50
51          pad_idx = TEXT.vocab.stoi['<pad>']
52
53          self.embed = nnn.Embedding(num_embeddings=len(TEXT.vocab),
54                                  embedding_dim=embedding_dim,
55                                  padding_idx=pad_idx)
56
57          self.lstm = nnn.LSTM(input_size=embedding_dim + hidden_dim,
58                          hidden_size=hidden_dim,
59                          num_layers=num_layers,
60                          dropout=dropout) \
61                  .spec("embedding", "trgSeqlen")
62
63          self.w = nnn.Linear(in_features=hidden_dim,
64                          out_features=len(TEXT.vocab)) \
65                  .spec("embedding", "classes")
66
67          self.attention = []
68
69          self.attn = Attention(attn_normalize)
70
71      def forward(self, init_state, batch_text):
72          H, (ht, ct) = init_state
73          ht, ct = (unsqueeze(flatten(ht, "layers", "embedding"), "layers"),
74                  unsqueeze(flatten(ct, "layers", "embedding"), "layers"))
75
76          embedded = self.embed(batch_text)
77
78          hidden_states = []
79          attention_weights = []
80          ht_flat = flatten(ht, "layers", "embedding")
81
82          for t in range(embedded.shape["trgSeqlen"]):
83
84              a, context = self.attn(ht_flat, H, "embedding", "srcSeqlen")
```

```
85              context = unsqueeze(context, "trgSeqlen")
86              word_t = embedded[{'trgSeqlen': slice(t, t+1)}]
87              lstm_input = ntorch.cat([word_t, context], "embedding")
88              output, (ht, ct) = self.lstm(lstm_input, (ht, ct))
89              if not self.training:
90                  attention_weights.append(a)
91              ht_flat = flatten(ht, "layers", "embedding")
92              hidden_states.append(ht_flat)
93
94          if not self.training:
95              self.attention = ntorch.stack(attention_weights, 'trgSeqlen')
96          return self.w(ntorch.stack(hidden_states, "trgSeqlen"))
```

*Listing 4: Attention*

```
1  import torch
2  from torch import nn
3  from namedtensor import ntorch, NamedTensor
4  from namedtensor.nn import nn as nnn
5
6  device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
7
8  class Attention(nnn.Module):
9      def __init__(self, normalize=True):
10         super().__init__()
11         self.normalize = normalize
12
13     def forward(self, query, key, embedding_dim, softmax_dim,
14                 mask=None, value=None, dropout=None,
15                 self_attention=False):
16         if self.normalize:
17             log_weights = (query.dot(embedding_dim, key)
18                            / (key.size(embedding_dim) ** .5))
19         else:
20             log_weights = query.dot(embedding_dim, key)
21
22         if mask is not None:
23             log_weights = log_weights.masked_fill_(mask == 0, -1e9)
24
25         if value is None:
26             value = key
27
28         a = log_weights.softmax(softmax_dim)
29
30         if dropout is not None:
31             a = dropout(a)
32
33         return (a, (a * value).sum(softmax_dim))
```

*Listing 5: Beam search*

```
1  import torch
2  from namedtensor import ntorch, NamedTensor
3  from namedtensor.nn import nn as nnn
4
5  device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
6
7  class Beam(nnn.Module):
```

```python
8      def __init__(self, TEXT, beam_size=3, topk=10, max_len=50):
9          super().__init__()
10         self.TEXT = TEXT
11         self.beam_size = beam_size
12         self.topk = topk
13         self.max_len = max_len
14         self.scores = None
15         self.nodes = None
16         self.result = None
17         self.filled = None
18         self.top_scores = None
19         self.top_score_locs = None
20
21     def log_prob(self, model, src):
22         # calculate marginal probability for next word
23         if self.nodes.shape['trgSeqlen'] == 1:
24             # start with SOS
25             log_prob = model(src, self.nodes[{'beam': 0}])[{'trgSeqlen': -1}]
26             return log_prob, log_prob.shape['classes']
27         else:
28             log_prob = [model(src, self.nodes[{'beam': b}])[{'trgSeqlen': -1}]
29                         + self.scores[{'beam': b}]
30                         for b in range(self.beam_size)]
31             return ntorch.cat(log_prob,'classes'), log_prob[0].shape['classes']
32
33     def generate_sentence(self, b, k, vocab_size):
34         sentence_i = self.top_score_locs[{'batch': b, 'classes': k}] / vocab_size
35         prev_sentence = self.nodes[{'batch': b, 'beam': sentence_i}].values[:-1]
36         word = self.top_score_locs[{'batch': b, 'classes': k}].fmod(vocab_size)
37         word = torch.LongTensor([word.values]).to(device)
38         sentence = NamedTensor(torch.cat([prev_sentence, word]), names=('trgSeqlen'))
39
40         return sentence
41
42     def advance(self, batch_size, vocab_size):
43         # increment length and pad using 1
44         increment = ntorch.ones((batch_size, self.beam_size, 1),
45                                 names=('batch', 'beam', 'trgSeqlen')) \
46                     .long().to(device)
47         self.nodes = ntorch.cat([self.nodes, increment], 'trgSeqlen')
48
49         # generate new next word for the entire batch
50         for b in range(batch_size):
51             beam_count = 0
52             new_nodes = ntorch.zeros((self.nodes.shape['trgSeqlen'],
53                                       self.beam_size),
54                                      names=('trgSeqlen', 'beam')) \
55                         .long().to(device)
56             if len(self.result[b]) < self.topk:
57                 for k in range(self.topk):
58                     sentence = self.generate_sentence(b, k, vocab_size)
59                     added_word = sentence[{'trgSeqlen': -1}]
60                     if len(self.result[b]) < self.topk and \
61                         (added_word.values.item() == self.TEXT.vocab.stoi["</s>"]
62                          or self.nodes.shape['trgSeqlen'] == self.max_len + 1):
63                         self.result[b].append(sentence)
64                     else:
65                         new_nodes[{'beam': beam_count}] = sentence
66                         self.scores[{'batch': b, 'beam': beam_count}] = \
```

```
67                         self.top_scores[{'batch': b, 'classes': k}].values.item()
68                     beam_count += 1
69                     if beam_count == self.beam_size:
70                         break
71             self.nodes[{'batch': b}] = new_nodes
72         else:
73             self.filled[b] = True
74
75     def forward(self, model, src):
76         batch_size = src.shape['batch']
77         self.nodes = (ntorch.ones((batch_size, self.beam_size, 1),
78                              names=('batch', 'beam', 'trgSeqlen'))
79                 * self.TEXT.vocab.stoi["<s>"]).long().to(device)
80         self.scores = ntorch.zeros((batch_size, self.beam_size),
81                              names=('batch', 'beam')).to(device)
82         self.result = [[] for _ in range(batch_size)]
83         self.filled = [False] * batch_size
84
85         while (sum(self.filled) < batch_size and
86                self.nodes.shape['trgSeqlen'] <= self.max_len):
87             log_prob, vocab_size = self.log_prob(model, src)
88             self.top_scores, self.top_score_locs = log_prob.topk('classes', self.topk)
89             self.advance(batch_size, vocab_size)
90
91         return self.result
```

## A.1   Transformer implementation

*Listing 6: Encoder Layer*

```
1  from namedtensor.nn import nn as nnn
2  from sublayer_connection import SublayerConnection
3
4  class TransformerEncoderLayer(nnn.Module):
5      def __init__(self, attn, feed_forward, size, dropout=.3):
6          super().__init__()
7          self.attn = attn
8          self.feed_forward = feed_forward
9
10         self.sublayer = nnn.ModuleList([
11                 SublayerConnection(size, dropout, "embedding")
12                 for _ in range(2)])
13         self.size = size
14
15
16     def forward(self, src):
17         a, x = self.sublayer[0](src, lambda src: self.attn(src, src, "embedding",
18             "srcSeqlen"))
18         return self.sublayer[1](x, self.feed_forward)
```

*Listing 7: Decoder Layer*

```
1  from namedtensor.nn import nn as nnn
2  from sublayer_connection import SublayerConnection
3  from namedtensor import ntorch, NamedTensor
4
5
```

```
6  class TransformerDecoderLayer(nnn.Module):
7      def __init__(self, self_attn, src_attn, feed_forward, size, dropout=.3):
8          super().__init__()
9          self.self_attn = self_attn
10         self.src_attn = src_attn
11         self.feed_forward = feed_forward
12
13         self.sublayer = nnn.ModuleList([
14                 SublayerConnection(size, dropout, "embedding")
15                 for _ in range(3)])
16         self.size = size
17
18
19     def forward(self, encoded, trg):
20         trglen = trg.size("trgSeqlen")
21         mask = ntorch.triu(ntorch.ones(trglen, trglen, names=('src','trg'))
22                         .to(trg.values.device),
23                         dims=("src","trg"))
24         a, x = self.sublayer[0](trg, lambda x: self.self_attn(x, x, "embedding",
25             "trgSeqlen", mask=mask))
25         a, x = self.sublayer[1](trg, lambda x: self.src_attn(x, encoded, "embedding",
            "srcSeqlen"))
26
27         return a, self.sublayer[2](x, self.feed_forward)
```

*Listing 8: Encoder*

```
1  from namedtensor.nn import nn as nnn
2  from copy import deepcopy
3  from layernorm import LayerNorm
4  from namedtensor import ntorch, NamedTensor
5
6  MAX_LEN = 20
7
8
9  class TransformerEncoder(nnn.Module):
10     def __init__(self, embed_dim, position_dim, TEXT, layer, nlayers=3):
11         super().__init__()
12         size = layer.size
13         pad_idx = TEXT.vocab.stoi['<pad>']
14         assert embed_dim + position_dim == size, \
15             "Embedding dimension + position_embedding must equal size"
16
17         self.embed = nnn.Embedding(num_embeddings=len(TEXT.vocab),
18                                 embedding_dim=embed_dim,
19                                 padding_idx=pad_idx)
20
21         self.position_embed = nnn.Embedding(num_embeddings=MAX_LEN + 1,
22                                 embedding_dim=position_dim,
23                                 padding_idx=MAX_LEN)
24
25         self.layers = nnn.ModuleList([deepcopy(layer) for _ in range(nlayers)])
26         self.norm = LayerNorm(size, "embedding")
27
28     def forward(self, x):
29         pos = ntorch.ones(*x.shape.values(),
30             names=[*x.shape.keys()]).to(x.values.device)
30         pos_vec = ntorch.arange(x.size("srcSeqlen"), names="srcSeqlen").float()
31         pos_vec[pos_vec > MAX_LEN] = MAX_LEN
```

```
32          embed = self.embed(x)
33          position_embed = self.position_embed((pos *
                pos_vec.to(x.values.device)).long())
34
35
36          x = ntorch.cat([embed, position_embed], "embedding")
37          for layer in self.layers:
38              x = layer(x)
39          return self.norm(x)
```

*Listing 9: Decoder*

```
1   from namedtensor.nn import nn as nnn
2   from copy import deepcopy
3   from layernorm import LayerNorm
4   from namedtensor import ntorch, NamedTensor
5   from position_encoding import PositionalEncoding
6
7   MAX_LEN = 20
8
9   class TransformerDecoder(nnn.Module):
10      def __init__(self, embed_dim, position_dim, TEXT, layer, nlayers=3):
11          super().__init__()
12          size = layer.size
13          pad_idx = TEXT.vocab.stoi['<pad>']
14          assert embed_dim + position_dim == size, \
15              "Embedding dimension + position_embedding must equal size"
16
17          self.embed = nnn.Embedding(num_embeddings=len(TEXT.vocab),
18                                     embedding_dim=embed_dim,
19                                     padding_idx=pad_idx)
20
21          self.position_embed = nnn.Embedding(num_embeddings=MAX_LEN + 1,
22                                     embedding_dim=position_dim,
23                                     padding_idx=MAX_LEN)
24
25          self.layers = nnn.ModuleList([deepcopy(layer) for _ in range(nlayers)])
26          self.norm = LayerNorm(size, "embedding")
27          self.w = nnn.Linear(in_features=size,
28              out_features=len(TEXT.vocab)).spec("embedding", "classes")
29
29      def forward(self, encoded, trg):
30          # attn_weights = []
31          pos = ntorch.ones(*trg.shape.values(),
                names=[*trg.shape.keys()]).to(trg.values.device)
32          pos_vec = ntorch.arange(trg.size("trgSeqlen"), names="trgSeqlen").float()
33          pos_vec[pos_vec > MAX_LEN] = MAX_LEN
34          embed = self.embed(trg)
35          position_embed = self.position_embed((pos *
                pos_vec.to(trg.values.device)).long())
36
37
38          x = ntorch.cat([embed, position_embed], "embedding")
39          for layer in self.layers:
40              a, x = layer(encoded, x)
41              # if not self.training:
42              # attn_weights.append(a)
43          # if not self.training:
44          # self.attn_weights = ntorch.stack(self.attn_weights, "layers")
```

17

```
45          x = self.norm(x)
46          # self.attn_weights = attn_weights
47          return self.w(x)
```

*Listing 10: Sublayer connection*

```
1  import torch
2  from torch import nn
3  from namedtensor import ntorch, NamedTensor
4  from namedtensor.nn import nn as nnn
5  from attention import Attention
6  from layernorm import LayerNorm
7
8  class SublayerConnection(nnn.Module):
9      """
10     A residual connection followed by a layer norm.
11     Note for code simplicity the norm is first as opposed to last.
12     """
13     def __init__(self, size, dropout, dim):
14         super().__init__()
15         self.norm = LayerNorm(size, dim)
16         self.dropout = nnn.Dropout(dropout)
17
18     def forward(self, x, sublayer):
19         "Apply residual connection to any sublayer with the same size."
20         out = sublayer(self.norm(x))
21         if type(out) is tuple:
22             a, context = out
23             return a, x + self.dropout(context)
24         else:
25             return x + self.dropout(out)
```

*Listing 11: Multihead Attention*

```
1  import torch
2  from torch import nn
3  from namedtensor import ntorch, NamedTensor
4  from namedtensor.nn import nn as nnn
5  from attention import Attention
6
7  device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
8
9  class MultiHeadAttention(nnn.Module):
10     """
11     Takes query, key, value tuple and pass through n_head linear layers
12     (with interim_dim_size output dimensions),
13     attend, concatenate, and pass through one final linear layer.
14     The layer maps input_dim_name of input_dim_size to
15     input_dim_name with output_dim_size.
16     """
17     def __init__(self, n_heads, input_dim_name, input_dim_size,
18                  interim_dim_size, output_dim_size, dropout=0, attn=None):
19         super().__init__()
20         self.linears = nnn.ModuleList([nnn.Linear(
21             in_features=input_dim_size,
22             out_features=interim_dim_size).spec(input_dim_name, input_dim_name)
23         for _ in range(n_heads)])
24
```

```
25          if attn is None:
26              self.attn = Attention()
27          else:
28              self.attn = attn
29
30          self.dropout = nnn.Dropout(dropout)
31          self.linear_final = nnn.Linear(
32              in_features=interim_dim_size * n_heads,
33              out_features=output_dim_size).spec(input_dim_name, input_dim_name)
34
35      def forward(self, query, key, embedding_dim, softmax_dim, mask=None, value=None):
36          if value is None:
37              value = key
38          attended = [self.attn(l(query), l(key), embedding_dim,
39              softmax_dim, mask=mask, value=l(value), dropout=self.dropout)
40              for l in self.linears]
41          a = [tup[0] for tup in attended]
42          contexts = ntorch.cat([tup[1] for tup in attended], embedding_dim)
43          return a, self.linear_final(contexts)
```

*Listing 12: Self-attention*

```
1   import torch
2   from torch import nn
3   from namedtensor import ntorch, NamedTensor
4   from namedtensor.nn import nn as nnn
5   from attention import Attention
6
7   class SelfAttention(Attention):
8       def __init__(self, *args):
9           super().__init__(*args)
10      def forward(self, query, key, embedding_dim, softmax_dim, mask=None,
11          value=None, **kwargs):
12          # Source is key,values, Target is query
13          target_name = f"trg{softmax_dim}"
14          trg = key.rename(softmax_dim, target_name)
15          if mask is not None:
16              mask = mask.rename("src", softmax_dim).rename("trg", target_name)
17
18          a, context = super().forward(query, trg, embedding_dim, softmax_dim,
19              mask, value=trg, **kwargs)
20
21          return a, context.rename(target_name, softmax_dim)
```

*Listing 13: Layer norm*

```
1   import torch
2   from torch import nn
3   from namedtensor import ntorch, NamedTensor
4   from namedtensor.nn import nn as nnn
5
6   class LayerNorm(nnn.Module):
7       def __init__(self, features, dim, eps=1e-6):
8           super().__init__()
9           self.a_2 = ntorch.ones(features, names=(dim,))
10          self.b_2 = ntorch.ones(features, names=(dim,))
11          self.register_parameter("layernorma", self.a_2)
12          self.register_parameter("layernormb", self.b_2)
```

```python
13          self.eps = eps
14          self.dim = dim
15
16      def forward(self, x):
17          mean = x.mean(self.dim)
18          std = x.std(self.dim)
19          return self.a_2 * (x - mean) / (std + self.eps) + self.b_2
```