

FreeSerif Latin Modern Mono

# CS 224n: Assignment #4

This assignment is split into two sections: *Neural Machine Translation with RNNs* and *Analyzing NMT Systems*. The first is primarily coding and implementation focused, whereas the second entirely consists of written, analysis questions. If you get stuck on the first section, you can always work on the second as the two sections are independent of each other. Note that the NMT system is more complicated than the neural networks we have previously constructed within this class and takes about **4 hours to train on a GPU**. Thus, we strongly recommend you get started early with this assignment. Finally, the notation and implementation of the NMT system is a bit tricky, so if you ever get stuck along the way, please come to Office Hours so that the TAs can support you.

## 1. Neural Machine Translation with RNNs (45 points)

In Machine Translation, our goal is to convert a sentence from the *source* language (e.g. Cherokee) to the *target* language (e.g. English). In this assignment, we will implement a sequence-to-sequence (Seq2Seq) network with attention, to build a Neural Machine Translation (NMT) system. In this section, we describe the **training procedure** for the proposed NMT system, which uses a Bidirectional LSTM Encoder and a Unidirectional LSTM Decoder.

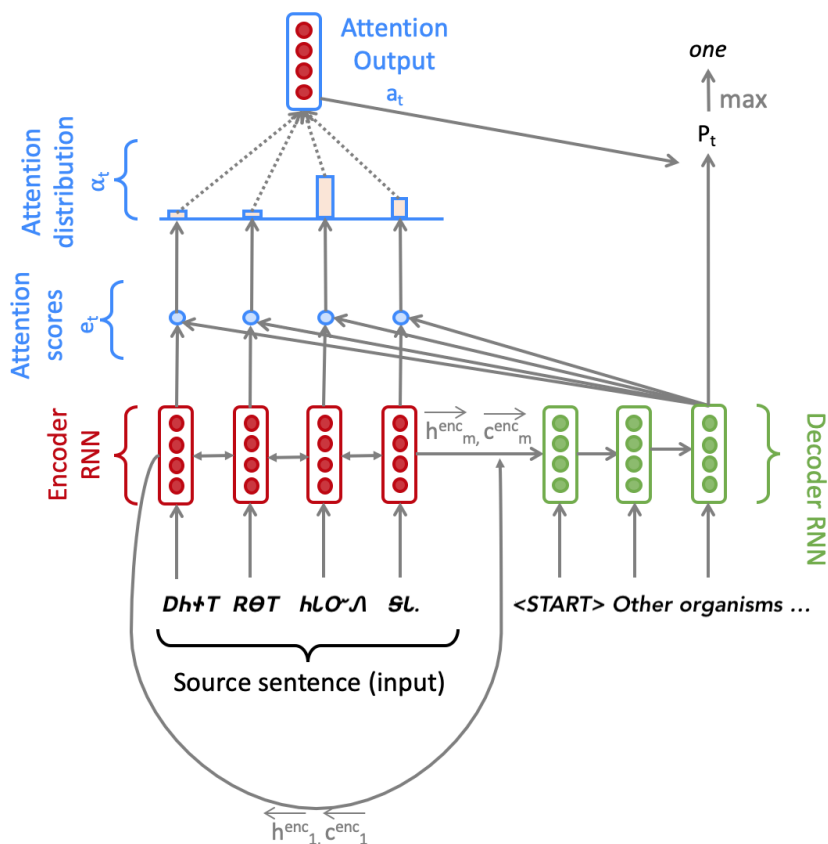


Figure 1: Seq2Seq Model with Multiplicative Attention, shown on the third step of the decoder. Hidden states  $\mathbf{h}_i^{\text{enc}}$  and cell states  $\mathbf{c}_i^{\text{enc}}$  are defined in the next page.

## Model description (training procedure)

Given a sentence in the source language, we look up the subword embeddings from an embeddings matrix, yielding  $\mathbf{x}_1, \dots, \mathbf{x}_m$  ( $\mathbf{x}_i \in \mathbb{R}^{e \times 1}$ ), where  $m$  is the length of the source sentence and  $e$  is the embedding size. We feed these embeddings to the bidirectional encoder, yielding hidden states and cell states for both the forwards ( $\rightarrow$ ) and backwards ( $\leftarrow$ ) LSTMs. The forwards and backwards versions are concatenated to give hidden states  $\mathbf{h}_i^{\text{enc}}$  and cell states  $\mathbf{c}_i^{\text{enc}}$ :

$$\mathbf{h}_i^{\text{enc}} = [\overleftarrow{\mathbf{h}_i^{\text{enc}}}; \overrightarrow{\mathbf{h}_i^{\text{enc}}}] \text{ where } \mathbf{h}_i^{\text{enc}} \in \mathbb{R}^{2h \times 1}, \overleftarrow{\mathbf{h}_i^{\text{enc}}}, \overrightarrow{\mathbf{h}_i^{\text{enc}}} \in \mathbb{R}^{h \times 1} \quad 1 \leq i \leq m \quad (1)$$

$$\mathbf{c}_i^{\text{enc}} = [\overleftarrow{\mathbf{c}_i^{\text{enc}}}; \overrightarrow{\mathbf{c}_i^{\text{enc}}}] \text{ where } \mathbf{c}_i^{\text{enc}} \in \mathbb{R}^{2h \times 1}, \overleftarrow{\mathbf{c}_i^{\text{enc}}}, \overrightarrow{\mathbf{c}_i^{\text{enc}}} \in \mathbb{R}^{h \times 1} \quad 1 \leq i \leq m \quad (2)$$

We then initialize the decoder's first hidden state  $\mathbf{h}_0^{\text{dec}}$  and cell state  $\mathbf{c}_0^{\text{dec}}$  with a linear projection of the encoder's final hidden state and final cell state.<sup>1</sup>

$$\mathbf{h}_0^{\text{dec}} = \mathbf{W}_h [\overleftarrow{\mathbf{h}_1^{\text{enc}}}; \overrightarrow{\mathbf{h}_m^{\text{enc}}}] \text{ where } \mathbf{h}_0^{\text{dec}} \in \mathbb{R}^{h \times 1}, \mathbf{W}_h \in \mathbb{R}^{h \times 2h} \quad (3)$$

$$\mathbf{c}_0^{\text{dec}} = \mathbf{W}_c [\overleftarrow{\mathbf{c}_1^{\text{enc}}}; \overrightarrow{\mathbf{c}_m^{\text{enc}}}] \text{ where } \mathbf{c}_0^{\text{dec}} \in \mathbb{R}^{h \times 1}, \mathbf{W}_c \in \mathbb{R}^{h \times 2h} \quad (4)$$

With the decoder initialized, we must now feed it a target sentence. On the  $t^{\text{th}}$  step, we look up the embedding for the  $t^{\text{th}}$  subword,  $\mathbf{y}_t \in \mathbb{R}^{e \times 1}$ . We then concatenate  $\mathbf{y}_t$  with the *combined-output vector*  $\mathbf{o}_{t-1} \in \mathbb{R}^{h \times 1}$  from the previous timestep (we will explain what this is later down this page!) to produce  $\overline{\mathbf{y}}_t \in \mathbb{R}^{(e+h) \times 1}$ . Note that for the first target subword (i.e. the start token)  $\mathbf{o}_0$  is a zero-vector. We then feed  $\overline{\mathbf{y}}_t$  as input to the decoder.

$$\mathbf{h}_t^{\text{dec}}, \mathbf{c}_t^{\text{dec}} = \text{Decoder}(\overline{\mathbf{y}}_t, \mathbf{h}_{t-1}^{\text{dec}}, \mathbf{c}_{t-1}^{\text{dec}}) \text{ where } \mathbf{h}_t^{\text{dec}} \in \mathbb{R}^{h \times 1}, \mathbf{c}_t^{\text{dec}} \in \mathbb{R}^{h \times 1} \quad (5)$$

$$(6)$$

We then use  $\mathbf{h}_t^{\text{dec}}$  to compute multiplicative attention over  $\mathbf{h}_1^{\text{enc}}, \dots, \mathbf{h}_m^{\text{enc}}$ :

$$\mathbf{e}_{t,i} = (\mathbf{h}_t^{\text{dec}})^T \mathbf{W}_{\text{attProj}} \mathbf{h}_i^{\text{enc}} \text{ where } \mathbf{e}_t \in \mathbb{R}^{m \times 1}, \mathbf{W}_{\text{attProj}} \in \mathbb{R}^{h \times 2h} \quad 1 \leq i \leq m \quad (7)$$

$$\alpha_t = \text{softmax}(\mathbf{e}_t) \text{ where } \alpha_t \in \mathbb{R}^{m \times 1} \quad (8)$$

$$\mathbf{a}_t = \sum_{i=1}^m \alpha_{t,i} \mathbf{h}_i^{\text{enc}} \text{ where } \mathbf{a}_t \in \mathbb{R}^{2h \times 1} \quad (9)$$

$\mathbf{e}_{t,i}$  is a scalar, the  $i$ th element of  $\mathbf{e}_t \in \mathbb{R}^{m \times 1}$ , computed using the hidden state of the decoder at the  $t$ th step,  $\mathbf{h}_t^{\text{dec}} \in \mathbb{R}^{h \times 1}$ , the attention projection  $\mathbf{W}_{\text{attProj}} \in \mathbb{R}^{h \times 2h}$ , and the hidden state of the encoder at the  $i$ th step,  $\mathbf{h}_i^{\text{enc}} \in \mathbb{R}^{2h \times 1}$ .

We now concatenate the attention output  $\mathbf{a}_t$  with the decoder hidden state  $\mathbf{h}_t^{\text{dec}}$  and pass this through a linear layer, tanh, and dropout to attain the *combined-output vector*  $\mathbf{o}_t$ .

$$\mathbf{u}_t = [\mathbf{a}_t; \mathbf{h}_t^{\text{dec}}] \text{ where } \mathbf{u}_t \in \mathbb{R}^{3h \times 1} \quad (10)$$

$$\mathbf{v}_t = \mathbf{W}_u \mathbf{u}_t \text{ where } \mathbf{v}_t \in \mathbb{R}^{h \times 1}, \mathbf{W}_u \in \mathbb{R}^{h \times 3h} \quad (11)$$

$$\mathbf{o}_t = \text{dropout}(\tanh(\mathbf{v}_t)) \text{ where } \mathbf{o}_t \in \mathbb{R}^{h \times 1} \quad (12)$$

---

<sup>1</sup>If it's not obvious, think about why we regard  $[\overleftarrow{\mathbf{h}_1^{\text{enc}}}; \overrightarrow{\mathbf{h}_m^{\text{enc}}}]$  as the 'final hidden state' of the Encoder.

Then, we produce a probability distribution  $\mathbf{P}_t$  over target subwords at the  $t^{\text{th}}$  timestep:

$$\mathbf{P}_t = \text{softmax}(\mathbf{W}_{\text{vocab}} \mathbf{o}_t) \text{ where } \mathbf{P}_t \in \mathbb{R}^{V_t \times 1}, \mathbf{W}_{\text{vocab}} \in \mathbb{R}^{V_t \times h} \quad (13)$$

Here,  $V_t$  is the size of the target vocabulary. Finally, to train the network we then compute the cross entropy loss between  $\mathbf{P}_t$  and  $\mathbf{g}_t$ , where  $\mathbf{g}_t$  is the one-hot vector of the target subword at timestep  $t$ :

$$J_t(\theta) = \text{CrossEntropy}(\mathbf{P}_t, \mathbf{g}_t) \quad (14)$$

Here,  $\theta$  represents all the parameters of the model and  $J_t(\theta)$  is the loss on step  $t$  of the decoder. Now that we have described the model, let's try implementing it for Cherokee to English translation!

## Setting up your Virtual Machine

Follow the instructions in the [CS224n Azure Guide](#) (link also provided on website and Ed) in order to create your VM instance. This should take you approximately 45 minutes. Though you will need the GPU to train your model, we strongly advise that you first develop the code locally and ensure that it runs, before attempting to train it on your VM. GPU time is expensive and limited. It takes approximately **30 minutes to 1 hour** to train the NMT system. We don't want you to accidentally use all your GPU time for debugging your model rather than training and evaluating it. Finally, **make sure that your VM is turned off whenever you are not using it.**

**If your Azure subscription runs out of money, your VM will be temporarily locked and inaccessible. If that happens, please fill out a request form [here](#).**

In order to run the model code on your **local** machine, please run the following command to create the proper virtual environment:

```
conda env create --file local_env.yml
```

Note that this virtual environment **will not** be needed on the VM.

## Implementation and written questions

- (2 points) (coding) In order to apply tensor operations, we must ensure that the sentences in a given batch are of the same length. Thus, we must identify the longest sentence in a batch and pad others to be the same length. Implement the `pad_sents` function in `utils.py`, which shall produce these padded sentences.
- (3 points) (coding) Implement the `__init__` function in `model_embeddings.py` to initialize the necessary source and target embeddings.
- (4 points) (coding) Implement the `__init__` function in `nmt_model.py` to initialize the necessary model embeddings (using the `ModelEmbeddings` class from `model_embeddings.py`) and layers (LSTM, projection, and dropout) for the NMT system.
- (8 points) (coding) Implement the `encode` function in `nmt_model.py`. This function converts the padded source sentences into the tensor  $\mathbf{X}$ , generates  $\mathbf{h}_1^{\text{enc}}, \dots, \mathbf{h}_m^{\text{enc}}$ , and computes the initial state  $\mathbf{h}_0^{\text{dec}}$  and initial cell  $\mathbf{c}_0^{\text{dec}}$  for the Decoder. You can run a non-comprehensive sanity check by executing:

```

131 Requirement already satisfied: asttokens in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
132 Requirement already satisfied: pure-eval in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
133 Requirement already satisfied: executing in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
134 Requirement already satisfied: wheel>=0.26 in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages
135 Requirement already satisfied: grpcio>=1.24.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
136 Requirement already satisfied: markdown>=2.6.8 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
137 Requirement already satisfied: tensorboard-data-server<0.7.0,>=0.6.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
138 Requirement already satisfied: absl-py>=0.4 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
139 Requirement already satisfied: werkzeug>=1.0.1 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
140 Requirement already satisfied: tensorboard-plugin-wit>=1.6.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
141 Requirement already satisfied: google-auth-oauthlib<0.5,>=0.4.1 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
142 Requirement already satisfied: google-auth>=1.6.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
143 Requirement already satisfied: pyasn1-modules>=0.2.1 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
144 Requirement already satisfied: rsa<5,>=3.1.4 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
145 Requirement already satisfied: requests-oauthlib>=0.7.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
146 Requirement already satisfied: pyasn1<0.5.0,>=0.4.6 in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages
147 Requirement already satisfied: oauthlib>=3.0.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages
148 Defaulting to user installation because normal site-packages is not writeable
149 Requirement already satisfied: nltk in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (3.7)
150 Requirement already satisfied: joblib in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages (1.0.1)
151 Requirement already satisfied: click in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages (7.1.2)
152 Requirement already satisfied: tqdm in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (4.64.0)
153 Requirement already satisfied: regex>=2021.8.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (2022.10.31)
154 Defaulting to user installation because normal site-packages is not writeable
155 Running Sanity Check for Question 1d: Encode
156 -----
157 enc_hidden Sanity Checks Passed!
158 dec_init_state[0] Sanity Checks Passed!
159 dec_init_state[1] Sanity Checks Passed!
160 -----
161 All Sanity Checks Passed for Question 1d: Encode!
162 -----
163

```

Figure 2: Result of Sanity check 1d.

```
python sanity_check.py 1d
```

- (e) (8 points) (coding) Implement the decode function in `nmt_model.py`. This function constructs  $\bar{\mathbf{y}}$  and runs the `step` function over every timestep for the input. You can run a non-comprehensive sanity check by executing:

```
python sanity_check.py 1e
```

```

129 Requirement already satisfied: wheel in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages (0.36.2)
130 Requirement already satisfied: zstd in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from pyz-deprecation-shim->itlocal-b=1.5->streamlit->simpletransformers) (2022.2)
131 Requirement already satisfied: asttokens in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from stack-data->ipython>7.23.1->ipykernel>5.1.2->pydeck>0.1.dev5->streamlit->simpletransformers) (2.0.5)
132 Requirement already satisfied: executing in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from stack-data->ipython>7.23.1->ipykernel>5.1.2->pydeck>0.1.dev5->streamlit->simpletransformers) (0.8.4)
133 Requirement already satisfied: pure-eval in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from stack-data->ipython>7.23.1->ipykernel>5.1.2->pydeck>0.1.dev5->streamlit->simpletransformers) (0.1.1)
134 Requirement already satisfied: wheel>=0.26 in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages (from tensorboard->simpletransformers) (0.36.2)
135 Requirement already satisfied: google-auth>=1.6.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simpletransformers) (2.18.0)
136 Requirement already satisfied: markdown>=2.6.8 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simpletransformers) (3.4.3)
137 Requirement already satisfied: tensorboard-plugin-wit>=1.6.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simpletransformers) (1.8.1)
138 Requirement already satisfied: google-auth-oauthlib>=0.5,>=0.4.1 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simpletransformers) (0.4.6)
139 Requirement already satisfied: absl-py>=0.4 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simpletransformers) (1.2.0)
140 Requirement already satisfied: grpcio>=1.24.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simpletransformers) (1.47.0)
141 Requirement already satisfied: werkzeug>=1.0.1 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simpletransformers) (2.2.2)
142 Requirement already satisfied: tensorboard-data-server<0.7.0,>=0.6.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simpletransformers) (0.6.1)
143 Requirement already satisfied: rsa<5,>=3.1.4 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from google-auth>=1.6.3->tensorboard->simpletransformers) (4.9)
144 Requirement already satisfied: pyasn1-modules>=0.2.1 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from google-auth>=1.6.3->tensorboard->simpletransformers) (0.2.8)
145 Requirement already satisfied: requests-oauthlib>=0.7.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from google-auth-oauthlib>=0.5,>=0.4.1->tensorboard->simpletransformers) (1.3.1)
146 Requirement already satisfied: pyasn1<0.5.0,>=0.4.6 in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages (from pyasn1-modules>=0.2.1->google-auth>=1.6.3->tensorboard->simpletransformers) (0.4.8)
147 Requirement already satisfied: oauthlib>=3.0.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from requests-oauthlib>=0.7.0->google-auth-oauthlib>=0.5,>=0.4.1->tensorboard->simpletransformers) (3.1.1)
148 Defaulting to user installation because normal site-packages is not writeable
149 Requirement already satisfied: nltk in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (3.7)
150 Requirement already satisfied: regex>=2021.8.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from nltk) (2022.10.31)
151 Requirement already satisfied: joblib in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages (from nltk) (1.0.1)
152 Requirement already satisfied: tqdm in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from nltk) (4.64.0)
153 Requirement already satisfied: click in /apps/Arch/software/Python/3.9.5-GCCcore-10.3.0/lib/python3.9/site-packages (from nltk) (7.1.2)
154 Defaulting to user installation because normal site-packages is not writeable
155 -----
156 Running Sanity Check for Question 1e: Decode
157 -----
158 torch.Size([23, 5, 2])
159 combined_outputs Sanity Checks Passed!
160 -----
161 All Sanity Checks Passed for Question 1e: Decode!
162 -----
163

```

Figure 3: Result of Sanity check 1e.

- (f) (10 points) (coding) Implement the `step` function in `nmt_model.py`. This function applies the Decoder's LSTM cell for a single timestep, computing the encoding of the target subword  $\mathbf{h}_t^{\text{dec}}$ , the attention scores  $\mathbf{e}_t$ , attention distribution  $\alpha_t$ , the attention output  $\mathbf{a}_t$ , and finally the combined output  $\mathbf{o}_t$ . You can run a non-comprehensive sanity check by executing:

```
python sanity_check.py 1f
```

```

130 Requirement already satisfied: trdata in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from py12-deprecation-shim->streamlit->simltransformers) (1.0.1)
131 Requirement already satisfied: pttokens in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from stack-data->ipython=>7.22.1->ipykernel=>5.1.2->pydeck=>0.1.dev5->simltransformers) (0.1.0)
132 Requirement already satisfied: pure-eval in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from stack-data->ipython=>7.22.1->ipykernel=>5.1.2->pydeck=>0.1.dev5->simltransformers) (0.1.0)
133 Requirement already satisfied: executing in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from stack-data->ipython=>7.22.1->ipykernel=>5.1.2->pydeck=>0.1.dev5->simltransformers) (0.1.0)
134 Requirement already satisfied: tensorboard-plugin-wit==1.6.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simltransformers) (1.6.1)
135 Requirement already satisfied: tensorboard-data-server==0.7.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simltransformers) (0.6.0)
136 Requirement already satisfied: markdown==2.6.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simltransformers) (2.4.2)
137 Requirement already satisfied: grpcio==1.24.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simltransformers) (1.47.0)
138 Requirement already satisfied: wheel==0.26 in /apps/Arch/software/Python/3.9.5-GCCore-10.3.0/lib/python3.9/site-packages (from tensorboard->simltransformers) (0.36.2)
139 Requirement already satisfied: google-auth<3,>=1.6.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simltransformers) (2.10.0)
140 Requirement already satisfied: google-auth-oauthlib==0.5 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simltransformers) (0.4.6)
141 Requirement already satisfied: absl-py==0.4 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simltransformers) (1.2.0)
142 Requirement already satisfied: werkzeug==0.1 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from tensorboard->simltransformers) (2.2.2)
143 Requirement already satisfied: rsa<4,>=3.1.4 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from google-auth<3,>=1.6.3->tensorboard->simltransformers) (4.9)
144 Requirement already satisfied: pyasn1-modules==0.2.1 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from google-auth<3,>=1.6.3->tensorboard->simltransformers) (2.10.0)
145 Requirement already satisfied: requests-oauthlib==0.7.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from google-auth-oauthlib==0.5->=0.4.1->tensorboard->simltransformers) (0.7.0)
146 Requirement already satisfied: pyasn1==0.5.0 in /apps/Arch/software/Python/3.9.5-GCCore-10.3.0/lib/python3.9/site-packages (from pyasn1-modules==0.2.1->google-auth<3,>=1.6.3->tensorboard->simltransformers) (0.5.0)
147 Requirement already satisfied: oauthlib==3.0.0 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from requests-oauthlib==0.7.0->google-auth-oauthlib==0.5->=0.4.1->tensorboard->simltransformers) (3.0.0)
148 Defaulting to user installation because normal site-packages is not writable
149 Requirement already satisfied: nltk in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (3.7)
150 Requirement already satisfied: regex==2021.0.3 in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from nltk (2022.10.31)) (2022.10.31)
151 Requirement already satisfied: click in /apps/Arch/software/Python/3.9.5-GCCore-10.3.0/lib/python3.9/site-packages (from nltk (7.1.2)) (7.1.2)
152 Requirement already satisfied: tqdm in /cephyr/users/layegh/Alvis/.local/lib/python3.9/site-packages (from nltk (4.64.0)) (4.64.0)
153 Requirement already satisfied: joblib in /apps/Arch/software/Python/3.9.5-GCCore-10.3.0/lib/python3.9/site-packages (from nltk (1.0.1)) (1.0.1)
154 Defaulting to user installation because normal site-packages is not writable
155 -----
156 Running Sanity Check for Question 1f: Step
157 -----
158 dec_state[0] Sanity Checks Passed!
159 dec_state[1] Sanity Checks Passed!
160 combined_output Sanity Checks Passed!
161 e_t Sanity Checks Passed!
162 -----
163 All Sanity Checks Passed for Question 1f: Step!
164 -----

```

Figure 4: Result of Sanity check 1f.

(g) (3 points) (written) The `generate_sent_masks()` function in `nmt_model.py` produces a tensor called `enc_masks`. It has shape (batch size, max source sentence length) and contains 1s in positions corresponding to ‘pad’ tokens in the input, and 0s for non-pad tokens. Look at how the masks are used during the attention computation in the `step()` function (lines 295-296).

First explain (in around three sentences) what effect the masks have on the entire attention computation. Then explain (in one or two sentences) why it is necessary to use the masks in this way. The `generate_sent_masks()` function masks the padding tokens in all sentences to value 1 and other tokens to 0. They are then converted to  $-\infty$  (negative infinity) in attention score  $e_t$ . In the calculation of  $\alpha_t$ , they become zero since  $\exp(-\infty) = 0$ . Therefore, they have zero contribution to the attention output  $\mathbf{a}_t$  (weighted average). This is necessary since <pad> tokens are not part of the original source data and do not contain meaningful information and should not be attended to or affect the attention output.

Now it’s time to get things running! Execute the following to generate the necessary vocab file:

```
sh run.sh vocab
```

Or if you are on Windows, use the following command instead. Make sure you execute this in an environment that has python in path. For example, you can run this in the terminal of your IDE or your Anaconda prompt.

```
run.bat vocab
```

As noted earlier, we recommend that you develop the code on your personal computer. Confirm that you are running in the proper conda environment and then execute the following command to train the model on your local machine:

```
sh run.sh train_local
(Windows) run.bat train_local
```

To help with monitoring and debugging, the starter code uses tensorboard to log loss and perplexity during training using TensorBoard<sup>2</sup>. TensorBoard provides tools for logging and visualizing training information from experiments. To open TensorBoard, run the following in your conda environment:

```
tensorboard --logdir=runs
```

<sup>2</sup><https://pytorch.org/docs/stable/tensorboard.html>

You should see a significant decrease in loss during the initial iterations. Once you have ensured that your code does not crash (i.e. let it run till `iter 10` or `iter 20`), power on your VM from the Azure Web Portal. Then read the *Managing Code Deployment to a VM* section of our [Practical Guide to VMs](#) (link also given on website and Ed) for instructions on how to upload your code to the VM.

Next, install necessary packages to your VM by running:

```
pip install -r gpu_requirements.txt
```

Finally, turn to the *Managing Processes on a VM* section of the Practical Guide and follow the instructions to create a new tmux session. Concretely, run the following command to create tmux session called `nmt`.

```
tmux new -s nmt
```

Once your VM is configured and you are in a tmux session, execute:

```
sh run.sh train
(Windows) run.bat train
```

```
398 epoch 6, iter 3010, avg. loss 39.68, avg. ppl 5.60 cum. examples 320, speed 3480.12 words/sec, time elapsed 317.54 sec
399 epoch 6, iter 3020, avg. loss 42.53, avg. ppl 5.82 cum. examples 640, speed 8343.75 words/sec, time elapsed 318.46 sec
400 epoch 6, iter 3030, avg. loss 44.46, avg. ppl 6.13 cum. examples 960, speed 8753.08 words/sec, time elapsed 319.36 sec
401 epoch 6, iter 3040, avg. loss 45.65, avg. ppl 6.26 cum. examples 1280, speed 8697.15 words/sec, time elapsed 320.28 sec
402 epoch 6, iter 3050, avg. loss 40.77, avg. ppl 5.70 cum. examples 1600, speed 8794.23 words/sec, time elapsed 321.13 sec
403 epoch 6, iter 3060, avg. loss 45.16, avg. ppl 6.36 cum. examples 1920, speed 8332.47 words/sec, time elapsed 322.07 sec
404 epoch 6, iter 3070, avg. loss 46.35, avg. ppl 6.57 cum. examples 2240, speed 7760.97 words/sec, time elapsed 323.08 sec
405 epoch 6, iter 3080, avg. loss 42.05, avg. ppl 5.83 cum. examples 2560, speed 8350.15 words/sec, time elapsed 323.99 sec
406 epoch 6, iter 3090, avg. loss 46.47, avg. ppl 6.56 cum. examples 2880, speed 8179.83 words/sec, time elapsed 324.96 sec
407 epoch 6, iter 3100, avg. loss 41.57, avg. ppl 5.82 cum. examples 3200, speed 8073.77 words/sec, time elapsed 325.90 sec
408 epoch 6, iter 3110, avg. loss 45.32, avg. ppl 6.55 cum. examples 3520, speed 8140.35 words/sec, time elapsed 326.85 sec
409 epoch 6, iter 3120, avg. loss 43.26, avg. ppl 5.82 cum. examples 3840, speed 8657.39 words/sec, time elapsed 327.75 sec
410 epoch 6, iter 3130, avg. loss 45.91, avg. ppl 6.26 cum. examples 4160, speed 8479.43 words/sec, time elapsed 328.70 sec
411 epoch 7, iter 3140, avg. loss 41.52, avg. ppl 5.63 cum. examples 4472, speed 8635.67 words/sec, time elapsed 329.57 sec
412 epoch 7, iter 3150, avg. loss 48.39, avg. ppl 6.43 cum. examples 4792, speed 7771.33 words/sec, time elapsed 330.64 sec
413 epoch 7, iter 3160, avg. loss 42.38, avg. ppl 5.73 cum. examples 5112, speed 8866.32 words/sec, time elapsed 331.51 sec
414 epoch 7, iter 3170, avg. loss 45.34, avg. ppl 5.75 cum. examples 5432, speed 9002.49 words/sec, time elapsed 332.44 sec
415 epoch 7, iter 3180, avg. loss 43.79, avg. ppl 6.02 cum. examples 5752, speed 7958.09 words/sec, time elapsed 333.42 sec
416 epoch 7, iter 3190, avg. loss 40.70, avg. ppl 5.59 cum. examples 6072, speed 8349.51 words/sec, time elapsed 334.32 sec
417 epoch 7, iter 3200, cum. loss 41.34, avg. ppl 5.82 cum. examples 6392, speed 7765.94 words/sec, time elapsed 335.29 sec
418 epoch 7, iter 3200, cum. loss 43.63, cum. ppl 6.01 cum. examples 6392
419 begin validation ...
420 validation: iter 3200, dev. ppl 38.009722
421 hit patience 1
422 hit #5 trial
423 early stop!
```

Figure 5: Result of Training.

Once you know your code is running properly, you can detach from session and close your ssh connection to the server. To detach from the session, run:

```
tmux detach
```

You can return to your training model by ssh-ing back into the server and attaching to the tmux session by running:

```
tmux a -t nmt
```

- (h) (3 points) (written) Once your model is done training (**this should take under 1 hour on the VM**), execute the following command to test the model:

```
sh run.sh test
(Windows) run.bat test
```

Please report the model's corpus BLEU Score. It should be larger than 10.



```

13  load test source sentences from [./chr_en_data/test.chr]
14  load test target sentences from [./chr_en_data/test.en]
15  load model from model.bin
16  Corpus BLEU: 12.804304011090135

```

Figure 6: Report of BLEU Score.

- (i) (4 points) (written) In class, we learned about dot product attention, multiplicative attention, and additive attention. As a reminder, dot product attention is  $\mathbf{e}_{t,i} = \mathbf{s}_t^T \mathbf{h}_i$ , multiplicative attention is  $\mathbf{e}_{t,i} = \mathbf{s}_t^T \mathbf{W} \mathbf{h}_i$ , and additive attention is  $\mathbf{e}_{t,i} = \mathbf{v}^T \tanh(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_2 \mathbf{s}_t)$ .
- (2 points) Explain one advantage and one disadvantage of *dot product attention* compared to multiplicative attention.
    - **Advantage:** The *dot product attention* does not require extra parameters which makes it easy to compute.
    - **Disadvantage:** In the *dot product attention* the dimension of query and key vectors must be the same.
  - (2 points) Explain one advantage and one disadvantage of *additive attention* compared to multiplicative attention.
    - **Advantage:** The *additive attention* applies non-linearity which might bring a better performance.
    - **Disadvantage:** The *additive attention* is more computationally expensive to train.

## 2. Analyzing NMT Systems (33 points)

- (a) (3 points) In part 1, we modeled our NMT problem at a subword-level. That is, given a sentence in the source language, we looked up subword components from an embeddings matrix. Alternatively, we could have modeled the NMT problem at the word-level, by looking up whole words from the embeddings matrix. Why might it be important to model our Cherokee-to-English NMT problem at the subword-level vs. the whole word-level? (Hint: Cherokee is a polysynthetic language.)
- Cherokee is a polysynthetic language which means that this language may contain a very large vocabulary. Therefore, a large embedding matrix is required to implement and train an effective word-level NMT model.
- (b) (3 points) Transliteration is the representation of letters or words in the characters of another alphabet or script based on phonetic similarity. For example, the transliteration of (which translates to "do you know") from Cherokee letters to Latin script is tsanvtasgo. In the Cherokee language, "ts-" is a common prefix in many words, but the Cherokee character is "tsa". Using this example, explain why when modeling our Cherokee-to-English NMT problem at the subword-level, training on transliterated Cherokee text may improve performance over training on original Cherokee characters. (Hint: A prefix is a morpheme.) In transliterated Cherokee text, we separate the words into some sub-words, including a pre-fix (morpheme) followed by other sub-words. In other words, we help the model by separating original characters into morphemes. Therefore, training on this dataset might improve the performance.
- (c) (3 points) One challenge of training successful NMT models is lack of language data, particularly for resource-scarce languages like Cherokee. One way of addressing this challenge is with multilingual training, where we train our NMT on multiple languages (including Cherokee). You can read more about multilingual training here: <https://ai.googleblog.com/2019/10/exploring-massively-multilingual.html>. How does multilingual training help in improving NMT performance with low-resource languages?



By training our NMT on multiple languages, we try to transfer the knowledge learned via high-resource languages such as French, English, and ... to low-resource languages like Cherokee when the annotated examples are limited.

- (d) (6 points) Here we present a series of errors we found in the outputs of our NMT model (which is the same as the one you just trained). For each example of a reference (i.e., ‘gold’) English translation, and NMT (i.e., ‘model’) English translation, please:

1. Identify the error in the NMT translation.
2. Provide possible reason(s) why the model may have made the error (either due to a specific linguistic construct or a specific model limitation).
3. Describe one possible way we might alter the NMT system to fix the observed error. There are more than one possible fixes for an error. For example, it could be tweaking the size of the hidden layers or changing the attention mechanism.

Below are the translations that you should analyze as described above. Only analyze the underlined error in each sentence. Rest assured that you don’t need to know Cherokee to answer these questions. You just need to know English! If, however, you would like additional color on the source sentences, feel free to use a resource like <https://www.cherokeedictionary.net/> to look up words.

- i. (2 points) **Source Sentence:** , : .

**Reference Translation:** When she was finished ripping things out, her web looked something like this:

**NMT Translation:** When it was gone out of the web, he said the web in the web. **NMT model lacks the grammar and semantic understanding. Collecting more data to train would be a solution.**

- ii. (2 points) **Source Translation:** , ? .

**Reference Translation:** What’s wrong little tree? the boy asked.

**NMT Translation:** The little little little little little tree? asked him.

- iii. (2 points) **Source Sentence:** “ , ” .

**Reference Translation:** “ ‘Humble,’ ” said Mr. Zuckerman

**NMT Translation:** “It’s not a lot,” said Mr. Zuckerman. **This might be related to the situation where we do not have the word “humble” in the vocabulary. Adding more data can be a solution.**

- (e) (4 points) Now it is time to explore the outputs of the model that you have trained! The test-set translations your model produced in question 1-i should be located in `outputs/test_outputs.txt`.

- i. (2 points) Find a line where the predicted translation is correct for a long (4 or 5 word) sequence of words. Check the training target file (English); does the training file contain that string (almost) verbatim? If so or if not, what does this say about what the MT system learned to do? **in the sentence "But he said unto them, I have seen you in your sight" the sentence "But he said unto them" was verbatim in training data. This can indicate that the MT system might learn based on the target training data too.**
- ii. (2 points) Find a line where the predicted translation starts off correct for a long (4 or 5 word) sequence of words, but then diverges (where the latter part of the sentence seems totally unrelated). What does this say about the model’s decoding behavior? **The sentence "And when I saw him, and he took him in him, he went out." It seems that the decoder forgets the words emitted some time-step ago.**

- (f) (14 points) BLEU score is the most commonly used automatic evaluation metric for NMT systems. It is usually calculated across the entire test set, but here we will consider BLEU defined for a single example.<sup>3</sup> Suppose we have a source sentence  $\mathbf{s}$ , a set of  $k$  reference translations  $\mathbf{r}_1, \dots, \mathbf{r}_k$ , and a

<sup>3</sup>This definition of sentence-level BLEU score matches the `sentence_bleu()` function in the `nltk` Python package. Note

candidate translation  $\mathbf{c}$ . To compute the BLEU score of  $\mathbf{c}$ , we first compute the *modified  $n$ -gram precision*  $p_n$  of  $\mathbf{c}$ , for each of  $n = 1, 2, 3, 4$ , where  $n$  is the  $n$  in  **$n$ -gram**:

$$p_n = \frac{\sum_{\text{ngram} \in \mathbf{c}} \min \left( \max_{i=1, \dots, k} \text{Count}_{\mathbf{r}_i}(\text{ngram}), \text{Count}_{\mathbf{c}}(\text{ngram}) \right)}{\sum_{\text{ngram} \in \mathbf{c}} \text{Count}_{\mathbf{c}}(\text{ngram})} \quad (15)$$

Here, for each of the  $n$ -grams that appear in the candidate translation  $\mathbf{c}$ , we count the maximum number of times it appears in any one reference translation, capped by the number of times it appears in  $\mathbf{c}$  (this is the numerator). We divide this by the number of  $n$ -grams in  $\mathbf{c}$  (denominator).

Next, we compute the *brevity penalty* BP. Let  $\text{len}(\mathbf{c})$  be the length of  $\mathbf{c}$  and let  $\text{len}(\mathbf{r})$  be the length of the reference translation that is closest to  $\text{len}(\mathbf{c})$  (in the case of two equally-close reference translation lengths, choose  $\text{len}(\mathbf{r})$  as the shorter one).

$$BP = \begin{cases} 1 & \text{if } \text{len}(\mathbf{c}) \geq \text{len}(\mathbf{r}) \\ \exp \left( 1 - \frac{\text{len}(\mathbf{r})}{\text{len}(\mathbf{c})} \right) & \text{otherwise} \end{cases} \quad (16)$$

Lastly, the BLEU score for candidate  $\mathbf{c}$  with respect to  $\mathbf{r}_1, \dots, \mathbf{r}_k$  is:

$$BLEU = BP \times \exp \left( \sum_{n=1}^4 \lambda_n \log p_n \right) \quad (17)$$

where  $\lambda_1, \lambda_2, \lambda_3, \lambda_4$  are weights that sum to 1. The log here is natural log.

- i. (5 points) Please consider this example<sup>4</sup>:

Source Sentence  $\mathbf{s}$ : -

Reference Translation  $\mathbf{r}_1$ : *the light shines in the darkness and the darkness has not overcome it*

Reference Translation  $\mathbf{r}_2$ : *and the light shines in the darkness and the darkness did not comprehend it*

NMT Translation  $\mathbf{c}_1$ : and the light shines in the darkness and the darkness can not comprehend

NMT Translation  $\mathbf{c}_2$ : the light shines the darkness has not in the darkness and the trials

Please compute the BLEU scores for  $\mathbf{c}_1$  and  $\mathbf{c}_2$ . Let  $\lambda_i = 0.5$  for  $i \in \{1, 2\}$  and  $\lambda_i = 0$  for  $i \in \{3, 4\}$  (**this means we ignore 3-grams and 4-grams**, i.e., don't compute  $p_3$  or  $p_4$ ). When computing BLEU scores, show your working (i.e., show your computed values for  $p_1, p_2, \text{len}(\mathbf{c}), \text{len}(\mathbf{r})$  and BP). Note that the BLEU scores can be expressed between 0 and 1 or between 0 and 100. The code is using the 0 to 100 scale while in this question we are using the **0 to 1** scale.

Which of the two NMT translations is considered the better translation according to the BLEU Score? Do you agree that it is the better translation? **The BLEU score for  $\mathbf{c}_1$  is 0.877 and the BLEU score for  $\mathbf{c}_2$  is 0.797, which means that the  $\mathbf{c}_1$  performs better.**

- ii. (5 points) Our hard drive was corrupted and we lost Reference Translation  $\mathbf{r}_2$ . Please recompute BLEU scores for  $\mathbf{c}_1$  and  $\mathbf{c}_2$ , this time with respect to  $\mathbf{r}_1$  only. Which of the two NMT translations now receives the higher BLEU score? Do you agree that it is the better translation? **In this situation, the BLEU score for  $\mathbf{c}_1$  is 0.71 and the BLEU score for  $\mathbf{c}_2$  is 0.79 which**

that the NLTK function is sensitive to capitalization. In this question, all text is lowercased, so capitalization is irrelevant.

[http://www.nltk.org/api/nltk.translate.html#nltk.translate.bleu\\_score.sentence\\_bleu](http://www.nltk.org/api/nltk.translate.html#nltk.translate.bleu_score.sentence_bleu)

<sup>4</sup>Due to data availability, many Cherokee sentences with English reference translations are from the Bible. This example is John 1:5. The two reference translations are from the New International Version and the New King James Version translations of the Bible.

means that the  $c_2$  performs better. I do not agree with this due to some syntactic error in its translation.

- iii. (2 points) List two advantages and two disadvantages of BLEU, compared to human evaluation, as an evaluation metric for Machine Translation.

- **Advantages:** (1) It makes the evaluation faster. (2) It is easy to understand.
- **Disadvantages:** (1) It does not consider the semantic relations.

## Submission Instructions

You shall submit this assignment on GradeScope as two submissions – one for “Assignment 4 [coding]” and another for ‘Assignment 4 [written]’:

1. Run the `collect_submission.sh` script on Azure to produce your `assignment4.zip` file. You can use [scp](#) to transfer files between Azure and your local computer.
2. Upload your `assignment4.zip` file to GradeScope to “Assignment 4 [coding]”.
3. Upload your written solutions to GradeScope to “Assignment 4 [written]”. When you submit your assignment, make sure to tag all the pages for each problem according to Gradescope’s submission directions. Points will be deducted if the submission is not correctly tagged.