

Due Date: March 20th, 23:59, 2026

Instructions

- *This assignment is involved – please start well ahead of time.*
- *For all questions, show your work!*
- *Submit your report (PDF) and your code electronically via the course Gradescope page. Your report must contain answers to Problem 2 (all questions).*
- *For open-ended experiments (i.e., experiments that do not have associated test-cases), unless otherwise specified, you do not need to submit code – a report will suffice.*
- *You are encouraged to ask questions on the Piazza. TA for this assignment **Xutong Zhao**.*

Summary: In this assignment, you will implement a sequential language model (a **GRU**) with soft-attention and a self attention model (a **Transformer**).

In problem 1, you will use built-in PyTorch modules to implement a variety of RNN architectures (GRU, Bidirectional GRU Encoder, GRU Encoder Decoder with Soft Attention) along with implementing the Transformer architecture via it's various building blocks.

In problem 2, you will take your implemented modules and perform sentiment analysis on the Yelp Polarity dataset. You will compare the performance of the RNNs to various configurations of the Transformer. This dataset will be pre-processed by a BERT based Hugging Face tokenizer, which outputs the padded sequence with a corresponding mask denoting where the padding is, and a target label. You are provided a google colab allowing you to train your implemented RNN and Transformer models on the dataset. Throughout this assignment, **all sequences will have length 256**, and we will use zero-padding to pad shorter sequences. You will work with mini-batches of data, each with batchsize B elements.

Coding instructions You will be required to use PyTorch to complete all questions. Moreover, this assignment **requires running the models on GPU** (otherwise it will take an incredibly long time); if you don't have access to your own resources (e.g. your own machine, a cluster), please use Google Colab. For some questions, you will be asked to not use certain functions in PyTorch and implement these yourself using primitive functions from `torch`; in that case, the functions in question are explicitly disabled in the tests on Gradescope.

Submission You are expected to submit `gru_solution.py` and `transformer_solution.py` to gradescope for autograding, as well as a PDF report for experimental and open-ended problems.

Problem 1

Implementing a GRU Encoder Decoder with soft attention (15 pts) In this problem, you will be using PyTorch's built-in modules in order to implement a GRU and various architectures that make use of a GRU

1. (6 pts) A Gated Recurrent Unit is a simplified version of a Long-Short Term Memory Network (LSTM). LSTMs learn how to forget and retain important information using input, target, cell, and forget gates, and a running memory cell. GRUs abandons the running memory cell and couples the notion of an input and forgetting gate with an update gate to help model long term dependencies, with reset gates modeling short term dependencies. The GRU is smaller in parameter size, but experimental evidence often shows GRU networks performing similarly to LSTM networks.

The following set of equations defines a portion of the forward pass for a GRU.

$$\begin{aligned}r_t &= \sigma(x_t W_{ir}^T + b_{ir} + h_{t-1} W_{hr}^T + b_{hr}) \\z_t &= \sigma(x_t W_{iz}^T + b_{iz} + h_{t-1} W_{hz}^T + b_{hz}) \\n_t &= \tanh(x_t W_{in}^T + b_{in} + r_t * (h_{t-1} W_{hi}^T + b_{hi})) \\h_t &= (1 - z_t) * n_t + z_t * h_{t-1}.\end{aligned}$$

Where $*$ indicates elementwise product.

In this problem, you will take the given GRU class and implement a GRU using the above set of equations. Be-careful of indexing along sequences, and of appropriately storing each hidden state for the output. For this problem, **you are not allowed** to use the PyTorch `nn.GRU` or `nn.GRUCell` module. Your implementation will be compared against the outputs and backwards pass of Pytorch's `nn.GRU` implementation, on CPU using `torch.float64`. Note however that your implementation must also work with `torch.float32` inputs and on a CUDA enabled device.

2. (3pts) Bidirectional encoders make use of two RNN layers; the 'forward' layer applies an RNN to an input sequence from start to end, while the 'backward' layer applies an RNN to a sequence end to start. You must implement a bidirectional GRU encoder, with dropout on the embedding layer. The two directions of the GRU network will be summed up together.

Due to speed and numerical instability issues, you will use the Pytorch implementation of `nn.GRU` in developing your Encoder and Decoder networks.

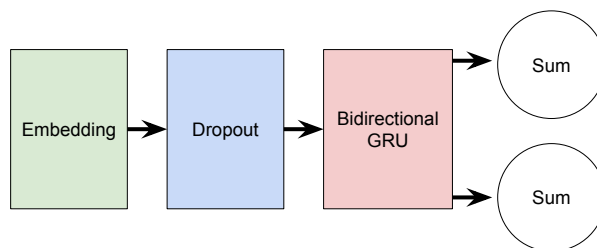


Figure 1: Encoder Network.

3. (4pts) Attention mechanisms help with giving the network access to and concentrating on specific information in the sequence. Attention also helps with improving gradient flow. Now, you will implement a one layer MLP which implements a cross attention mechanism. You **must** make sure that this MLP attention model can perform masking on the attention vector. Also note, the second linear transformation will be implemented as a multiply and a sum.

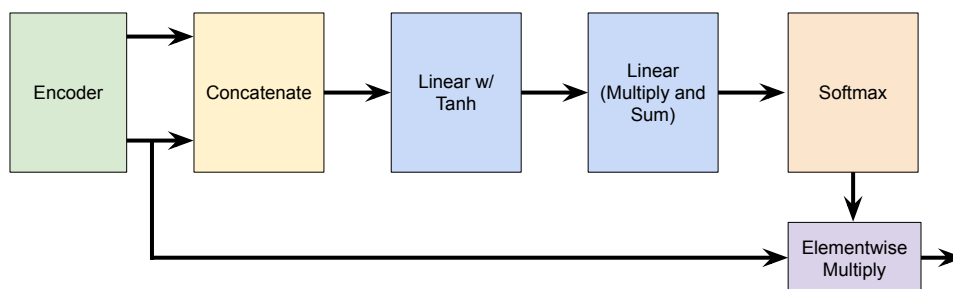


Figure 2: Attention Network.

4. (2pts) The decoder must take in the encoder outputs as input and hidden state, and these must be fed into the attention mechanism. The attended input and the encoder hidden state will then be fed into a GRU layer.

In the file `gru_solution.py`, you are given an `Encoder`, `Attn`, and `DecoderAttn` class containing all the blocks necessary to create this model. In particular, `self.embedding`, located in `Encoder`, is a `nn.Embedding` module that converts sequences of token indices into embeddings, while `self.rnn`, located in both `Encoder` and `DecoderAttn` is a `nn.GRU` module that runs a GRU over a sequence of vectors.

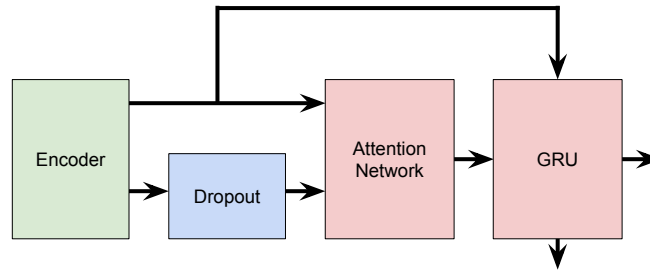


Figure 3: Decoder Network.

Implementing a Transformer (26pts) While typical RNNs “remember” past information by taking their previous hidden state as input at each step, recent years have seen a profusion of methodologies for making use of past information in different ways. The Transformer¹ is one such architecture which uses several self-attention networks (“heads”) in parallel, among other architectural specifics. Implementing a Transformer is a fairly involved process – so we provide most of the boilerplate code and your task is only to implement the multi-head scaled dot-product attention mechanism, as well as the layernorm operation. The attention mechanism here must use padded masking in order to pass the unit tests.

Implementing Layer Normalization (5pts): You will first implement the layer normalization (LayerNorm) technique that we have seen in class. For this assignment, **you are not allowed** to use the PyTorch `nn.LayerNorm` module (nor any function calling `torch.layer_norm`).

As defined in the [layer normalization paper](#), the layernorm operation over a minibatch of inputs x is defined as

$$\text{layernorm}(x) = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \text{weight} + \text{bias}$$

where $\mathbb{E}[x]$ denotes the expectation over x , $\text{Var}[x]$ denotes the variance of x , both of which are only taken over the last dimension of the tensor x here. `weight` and `bias` are learnable affine parameters.

- (5pts) In the file `transformer_solution.py`, implement the `forward()` function of the `LayerNorm` class. Pay extra attention to the lecture slides on the exact details of how $\mathbb{E}[x]$ and $\text{Var}[x]$ are computed. In particular, PyTorch’s function `torch.var` uses an unbiased estimate of the variance by default, defined as the formula on the left-hand side

$$\overline{\text{Var}}(X)_{\text{unbiased}} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad \overline{\text{Var}}(X)_{\text{biased}} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

whereas LayerNorm uses the biased estimate on the right-hand side (where \bar{X} here is the mean estimate). Please refer to the docstrings of this function for more information on input/output signatures.

¹See <https://arxiv.org/abs/1706.03762> for more details.

Implementing the attention mechanism (18pts): You will now implement the core module of the Transformer architecture – the multi-head attention mechanism. Assuming there are m attention heads, the attention vector for the head at index i is given by:

$$\begin{aligned} [\mathbf{q}_1, \dots, \mathbf{q}_m] &= \mathbf{Q}\mathbf{W}_Q + \mathbf{b}_Q & [\mathbf{k}_1, \dots, \mathbf{k}_m] &= \mathbf{K}\mathbf{W}_K + \mathbf{b}_K & [\mathbf{v}_1, \dots, \mathbf{v}_m] &= \mathbf{V}\mathbf{W}_V + \mathbf{b}_V \\ \mathbf{A}_i &= \text{softmax} \left(\frac{\mathbf{q}_i \mathbf{k}_i^\top}{\sqrt{d}} \right) \\ \mathbf{h}_i &= \mathbf{A}_i \mathbf{v}_i \\ A(\mathbf{Q}, \mathbf{K}, \mathbf{V}) &= \text{concat}(\mathbf{h}_1, \dots, \mathbf{h}_m) \mathbf{W}_O + \mathbf{b}_O \end{aligned}$$

Here $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are queries, keys, and values respectively, where all the heads have been concatenated into a single vector (e.g. here $\mathbf{K} \in \mathbb{R}^{T \times md}$, where d is the dimension of a single key vector, and T the length of the sequence). $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ are the corresponding projection matrices (with biases \mathbf{b}), and \mathbf{W}_O is the output projection (with bias \mathbf{b}_O). \mathbf{Q}, \mathbf{K} , and \mathbf{V} are determined by the output of the previous layer in the main network. \mathbf{A}_i are the attention values, which specify which elements of the input sequence each attention head attends to. In this question, **you are not allowed** to use the module `nn.MultiheadAttention` (or any function calling `torch.nn.functional.multi_head_attention_forward`). Please refer to the docstrings of each function for a precise description of what each function is expected to do, and the expected input/output tensors and their shapes.

- (4pts) The equations above require many vector manipulations in order to split and combine head vectors together. For example, the concatenated queries \mathbf{Q} are split into m vectors $[\mathbf{q}_1, \dots, \mathbf{q}_m]$ (one for each head) after an affine projection by \mathbf{W}_Q , and the \mathbf{h}_i 's are then concatenated back for the affine projection with \mathbf{W}_O . In the class `MultiHeadedAttention`, implement the utility functions `split_heads()` and `merge_heads()` to do both of these operations, as well as a transposition for convenience later. For example, for the 1st sequence in the mini-batch:

```
y = split_heads(x)  →  y[0, 1, 2, 3] = x[0, 2, (num_heads=1) * 3]
x = merge_heads(y)  →  x[0, 1, (num_heads=2) * 3] = y[0, 2, 1, 3]
```

These two functions are exactly inverse from one another. Note that in the code, the number of heads m is called `self.num_heads`, and the head dimension d is `self.head_size`. Your functions must handle mini-batches of sequences of vectors, see the docstring for details about the input/output signatures.

- (9pts) In the class `MultiHeadedAttention`, implement the function `get_attention_weights()`, which is responsible for returning \mathbf{A}_i 's (for all the heads at the same time) from \mathbf{q}_i 's and \mathbf{k}_i 's. Concretely, this means taking the softmax over the whole sequence. The softmax is then

$$[\text{softmax}(\mathbf{x})]_\tau = \frac{\exp(x_\tau)}{\sum_i \exp(x_i)}$$

Implement padded masking before the softmax.

4. (2pts) Using the functions you have implemented, complete the function `apply_attention()` in the class `MultiHeadedAttention`, which computes the vectors \mathbf{h}_i 's as a function of \mathbf{q}_i 's, \mathbf{k}_i 's and \mathbf{v}_i 's, and concatenates the head vectors.

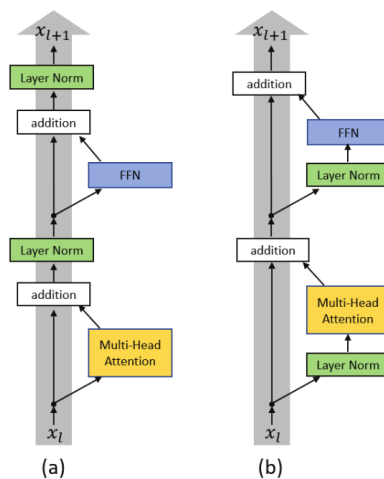
$$\text{apply_attention}(\{\mathbf{q}_i\}_{i=1}^m, \{\mathbf{k}_i\}_{i=1}^m, \{\mathbf{v}_i\}_{i=1}^m) = \text{concat}(\mathbf{h}_1, \dots, \mathbf{h}_m).$$

5. (3pts) Using the functions you have implemented, complete the function `forward()` in the class `MultiHeadedAttention`. You may implement the different affine projections however you want (do not forget the biases), and you can add modules to the `__init__()` function. How many learnable parameters does your module have, as a function of `num_heads` and `head_size`?

The Transformer forward pass (3pts): You now have all building blocks to implement the forward pass of a miniature Transformer model. You are provided a module `PostNormAttentionBlock` which corresponds to a full block with self-attention and a feed-forward neural network, with skip-connections, using the modules `LayerNorm` and `MultiHeadedAttention` you implemented before.

In this part of the exercise, you will fill in the `Transformer` class in `transformer_solution.py`. This module contains all the blocks necessary to create this model. In particular, an embedding layer using input and positional embeddings, `self.transformer` is a `nn.ModuleList` containing the different `Attention Block` layers.

6. (1pts) By taking inspiration from the `PostNormAttentionBlock`, implement the `PreNormAttentionBlock`. You can look at the implementation of the forward function of the `PostNorm` block to complete the forward function in `PreNormAttentionBlock`. See the figure below for a comparison of the post-norm and pre-norm.



(a) Post-LN Transformer layer; (b) Pre-LN Transformer layer.

Figure 4: Image from Ruibin Xiong et al [On Layer Normalization in the Transformer Architecture](#)

7. (2pts) In the class `Transformer`, complete the function `forward()` using the different modules described above.

Problem 2

Training Sequential models (40pts) *Before starting this assignment make sure your code for the previous problems works correctly by uploading it to the Gradescope!*

You will train and evaluate each of the following architectures. For reference, we have provided a *feature-complete* training notebook

(`main.ipynb`) that uses the **ADAMW** optimizer. You are free to modify this notebook as you deem fit. You do not need to submit code for this part of the assignment. However, you are required to create a report that presents the accuracy and training curve comparisons as specified in the following questions. You will train each model for at least 5 epochs, except for the fine-tuning of pre-trained BERT - in that case you need only up to 2 epochs. You are free to train any model for longer if you think that will improve your report.

Note: For each experiment, closely observe the training curves, and report the best validation accuracy score across training iterations (not necessarily the validation score for the last training iteration).

Configurations to run: We have provided 8 experiment configurations for you to run. These configurations span over several neural network architectures. Each run creates a log. Perform the following analysis on the logs.

1. (6pts) You are asked to run 8 experiments. For each of these experiments, plot learning curves (train loss and validation loss/accuracy) over **training iterations**. Figures should have labeled axes and a legend and an explanatory caption.
2. (5pts) Make a table of results summarizing the train and validation performance for each experiment, indicating the architecture, and including total wall clock time it took to train the architecture, and how long it takes to evaluate the model.² Sort by experiment number, and make sure to refer to these experiment numbers in bold script for easy reference. Bold the best validation result. The table should have an explanatory caption, and appropriate column and/or row headers. Any shorthand or symbols in the table should be explained in the caption.
3. (2pts) Among the 8 configurations, which would you use if you were most concerned with wall-clock time? With generalization performance?
4. (3pts) Compare experiments **1** and **2**. What was the impact of having dropout?
5. (3pts) Compare experiments **2** and **3**. Are there benefits of having Encoder-Decoder network instead of just an Encoder? Speculate as to why or why not.
6. (3pts) Compare experiments **3** and **4**. What was the impact of training the GRU network with soft-attention?

²You can also make the table in LaTeX; for convenience you can use tools like **LaTeX table generator** to generate tables online and get the corresponding LaTeX code.

7. (3pts) In experiments **5**, **6**, **7** and **8**, you trained or finetuned a Transformer. Given the recent high profile Transformer based models, are the results as you expected? Speculate as to why or why not.
8. (3pts) For each of the experiment configurations above, measure the average steady-state GPU memory usage. Comment about the GPU memory footprints of each model, discussing reasons behind increased or decreased memory consumption where applicable. One way to measure GPU usage would be via monitoring `nvidia-smi`. If you are using Google Colab, use the resources tab to monitor GPU usage.
9. (7pts) Comment on the learning curve behavior of the various models you trained, under different hyper-parameter settings. Did a particular class of models over/underfit more easily than the others? Exhibited instabilities in training? Can you make an informed guess of the various steps a practitioner can take to prevent over/under-fitting or instability in this case?
10. (5 pts) In this assignment you implemented Cross Attention in the form of Soft Attention, and then you implemented Self Attention. Give a high level overview of what an attention mechanism is, and go into each attention mechanism (Self, Cross), and their differences. Write two paragraphs total.