# University of Central Missouri

# Department of Computer Science & Cybersecurity

## CS5760 Natural Language Processing

## Fall 2025

## Homework 4.

## Student name:

## Kopparthi Sai Siva Shankara Vara Prasad

## Student ID: 700765221

## Submission Requirements:

- Once finished your assignment push your source code to your repo (GitHub) and explain the work through the ReadMe file properly. Make sure you add your student info in the ReadMe file.
- Submit your GitHub link on the Bright Space.
- Comment your code appropriately ***IMPORTANT.***
- Any submission after provided deadline is considered as a late submission.

**Part I. Short Answer**

**1) RNN Families & Use-Cases (Many-to-X)**
a) Map each task to the most suitable RNN I/O pattern and explain why (1–2 lines each):
• next-word prediction

• sentiment of a sentence

 • NER

• machine translation
(Choose from: one-to-many, many-to-one, many-to-many aligned, many-to-many unaligned.)

**Answer)**

**RNN Patterns for Tasks**

- **Next-word prediction** is best suited to a **one-to-many** architecture because the goal is to generate a sequence of outputs from a single input token or prompt. For example, given an initial input word or character, the model needs to generate a sequence of subsequent tokens, making predictions step-by-step. This pattern captures sequential data and can be seen in text generation systems where a single input such as "Once" unfolds into a longer sentence or paragraph. The model learns contextual dependencies and gradually unfolds information.
- **Sentiment analysis of a sentence** follows a **many-to-one** RNN pattern. Here, the input consists of multiple tokens (a sentence or document), and the output is a single label or continuous value representing the sentiment—often positive, neutral, or negative. This setup allows the RNN to process each token in the sentence, associate meaning to the context, and finally condense everything into a single output. This structure is powerful because the sequential nature of RNNs considers how words influence one another, allowing nuanced sentiment detection.
- **Named Entity Recognition (NER)** operates using a **many-to-many aligned** pattern because it requires tagging each token of an input sequence with an output label (e.g., PERSON, LOCATION, etc.). Since every input word is associated with a corresponding output tag, alignment of input and output lengths is essential. The RNN processes the full sequence and emits labels simultaneously, capturing both local and contextual patterns in the input text while maintaining a one-to-one mapping between tokens and labels.
- **Machine translation** uses a **many-to-many unaligned** architecture because the input and output sequence lengths usually differ, and the structure between input and output tokens is not directly aligned. For instance, translating "How are you?" into Spanish becomes "¿Cómo estás?", where word order and structure differ. This requires an encoder-decoder model, where the encoder reads the input sentence into a context vector, and the decoder generates the translated sentence step-by-step. This separation helps the model handle long-range dependencies and cross-lingual challenges.

b) In one sentence, explain how "unrolling" over time enables BPTT and weight sharing.

**Answer)**

*Unrolling* an RNN over time allows us to treat the temporal sequence as a deep structure where each timestep corresponds to a layer. This unrolled representation is essential for **Backpropagation Through Time (BPTT)**, as it allows gradients to flow backward through time steps and share weights across them. This shared structure significantly reduces the number of parameters and enables efficient learning, though it also introduces challenges like vanishing gradients.

c) Give one advantage and one limitation of weight sharing across time in RNNs.

**Answer)**

Weight sharing in RNNs offers the tremendous advantage of dramatically reducing model size since the same parameters are reused at every time step. It also allows the model to generalize patterns across different input positions because it doesn't rely on position-specific weights. However, the downside is that this can make learning long-term dependencies difficult, especially when early timesteps have little influence on later predictions due to the vanishing gradient phenomenon.

**2) Vanishing Gradients & Remedies**
a) Describe the vanishing gradient problem in RNNs and how it affects long-range dependencies.

**Answer)**

The **vanishing gradient problem** occurs in RNNs when gradients shrink exponentially as they propagate backward through time, especially over long sequences. This happens due to repeated multiplication by small values (like small derivatives from activation functions), effectively erasing the gradient information. As a result, it's hard for the model to learn dependencies that span many time steps (e.g., linking a pronoun to a noun several words earlier). This fundamentally limits the model's ability to handle complex contextual patterns and long-term relationships.

b) List two *architectural* solutions and briefly how each helps gradient flow.

**Answer)**

Two architectural solutions include:

- **LSTM (Long Short-Term Memory)**: LSTMs mitigate vanishing gradients through their cell state, which can carry information across long sequences with minimal modification,

thanks to gating mechanisms. These gates enable the controlled flow and storage of information and preserve gradients during backpropagation.

- **GRU (Gated Recurrent Unit)**: GRUs simplify the LSTM by combining the forget and input gates into a single update gate. Their simpler structure lowers computational cost while still effectively addressing vanishing gradients by allowing the network to preserve information across many time steps using gated mechanisms.

c) Give one *training* technique (not an architecture change) that can mitigate the issue and why.

**Answer)**

A training technique that combats vanishing gradients is **gradient clipping**, where the gradients are scaled down if they exceed a predefined threshold. This ensures stability during updates, especially in recurrent structures, and prevents exploding gradients, which indirectly helps by not amplifying small gradients to the point where they become negligible over time.

**3) LSTM Gates & Cell State**
a) Explain the roles of the forget, input, and output gates. For each, name its activation and purpose.

**Answer)**

In an LSTM, the **forget gate**—using a sigmoid activation—decides what information to discard from the cell state so old, irrelevant information is removed. The **input gate**, also controlled by a sigmoid function, determines what new information from the current input and hidden state should be added to the cell. The **output gate** regulates how much of the cell state influences the current output by applying sigmoid activation followed by a tanh function for the output vector. Together, these gates form memory that can adaptively store and forget information.

b) Why is the LSTM cell state often described as providing a "linear path" for gradients?

**Answer)**

The LSTM cell state serves as a nearly linear path for gradients over time, meaning there are minimal disruptions or transformations applied to it from step to step. This property creates a highway-like flow for gradients to pass back unobstructed, which is critical for learning long-range dependencies in sequences by avoiding exponential decay of gradients.

c) In one or two sentences, contrast "what to remember" vs. "what to expose" in LSTMs.

**Answer)**

*What to remember* refers to the choice of storing information in the cell state through the input and forget gates—this determines which parts of past information are important. *What to expose* is determined by the output gate and involves controlling which parts of the stored information influence the current output. This separation allows LSTMs to filter information effectively through time for efficient learning.

## 4) Self-Attention

a) Define Query (Q), Key (K), and Value (V) in the context of self-attention.

**Answer)**

In self-attention, the **Query (Q)** represents the vector the model uses to search for relevant information, the **Key (K)** represents labels or markers that match against queries, and the **Value (V)** carries the actual information to be used in downstream tasks. For each position in the sequence, the self-attention mechanism looks at all other positions via Q-K similarity and uses those weighted similarities to combine values.

b) Write the formula for dot-product attention.

**Answer)**

The scaled dot-product attention formula is:

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

where $QK^T$ computes similarity, softmax normalizes these into weights, and $V$ applies those weights to extract relevant information.

c) Why do we divide by $\sqrt{d_k}$?

**Answer)**

Dividing by $\sqrt{d_k}$ (where $d_k$ is the Key dimension) stabilizes softmax computation by preventing large values in the dot product, which could result in a sharp, uninformative softmax distribution. This helps maintain smooth gradient flow and consistent learning across different dimensionalities.

**5) Multi-Head Attention & Residual Connections**

a) Why do Transformers use multi-head attention instead of single-head attention?

**Answer)**

Transformers use **multi-head attention** to allow the model to attend to different parts of the input in different ways. For example, one head might focus on syntactic structure, another on semantic roles, and another on positional relationships. Having multiple perspectives enriches the contextual understanding significantly and allows deeper modeling of linguistic patterns than a single attention head.

b) What is the purpose of Add & Norm (Residual + LayerNorm)? Explain two benefits.

**Answer)**

*Add & Norm*, composed of residual connections and layer normalization, is essential to ensure stable and efficient training. Residual connections ease gradient flow by adding inputs directly to outputs of sublayers, reducing vanishing gradient problems and improving deep network learning. Layer normalization ensures that training remains stable by normalizing activations throughout the network (especially across batch sizes of 1, like in machine translation).

c) Describe one example of linguistic relation that different heads might capture (e.g., coreference, syntax).

**Answer)**

For example, one head might capture **coreference resolution**, connecting "Barack Obama" to "he" later in the text. Another might focus on **syntactic dependencies**, determining which word is the subject of a verb. This division of labor makes the model highly expressive.

**6) Encoder–Decoder with Masked Attention**

a) Why does the decoder use masked self-attention? What problem does it prevent?

**Answer)**

Masked self-attention ensures that during sequence generation, the model cannot look ahead at future tokens. This prevents "cheating" where the decoder accidentally uses future context during training. Each token can only attend to previous tokens, preserving the autoregressive nature of generation and modeling true prediction context.

b) What is the difference between encoder self-attention and encoder–decoder cross-attention?

**Answer)**

**Encoder self-attention** computes attention over the entire source sequence to build contextualized token embeddings. **Encoder–decoder cross-attention**, on the other hand, allows the decoder to attend to encoder outputs, enabling it to incorporate source context while generating target sequences. This mechanism is essential for translation-based tasks.

c) During inference (no teacher forcing), how does the model generate tokens step by step?

**Answer)**

During inference, token generation happens one step at a time without teacher forcing. The model generates the first token, feeds it back as input to predict the next one, and repeats this process until a stopping condition like an EOS token is reached. This process helps the model learn real-world sequence generation.

**Part II: Programming**

**Q1. Character-Level RNN Language Model ("hello" toy & beyond)**
Goal: Train a tiny character-level RNN to predict the next character given previous characters.

Data (toy to start):

- Start with a small toy corpus you create (e.g., several "hello…", "help…", short words/sentences).

- Then expand to a short plain-text file of ~50–200 KB (any public-domain text of your choice).

Model:

- Embedding → RNN (Vanilla RNN or GRU or LSTM) → Linear → Softmax over characters.
- Hidden size 64–256; sequence length 50–100; batch size 64; train 5–20 epochs.

Train:

- Teacher forcing (use the true previous char as input during training).
- Cross-entropy loss; Adam optimizer.

Report:

1. Training/validation loss curves.
2. Sample 3 temperature-controlled generations (e.g., $\tau = 0.7$, 1.0, 1.2) for 200–400 chars each.
3. A 3–5 sentence reflection: what changes when you vary sequence length, hidden size, and temperature?
   (Connect to slides: embedding, sampling loop, teacher forcing, tradeoffs)

**OUTPUT:**

**Entire data is written in GitHub Readme file. Here is the Screenshot of the output.**

## Q2. Mini Transformer Encoder for Sentences

Task: Build a mini Transformer Encoder (NOT full decoder) to process a batch of sentences.

Steps:

1. Use a small dataset (e.g., 10 short sentences of your choice).
2. Tokenize and embed the text.
3. Add sinusoidal positional encoding.
4. Implement:
   o Self-attention layer
   o Multi-head attention (2 or 4 heads)
   o Feed-forward layer
   o Add & Norm
5. Show:
   o Input tokens
   o Final contextual embeddings
   o Attention heatmap between words (visual or printed)

## OUTPUT:

**Entire data is written in GitHub Readme file. Here is the Screenshot of the output.**

## Q3. Implement Scaled Dot-Product Attention

Goal: Implement the attention function from your slides:

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Requirements:

- Write a function in PyTorch or TensorFlow to compute attention.
- Test it using random Q, K, V inputs.
- Print:
    - Attention weight matrix
    - Output vectors
    - Softmax stability check (before and after scaling)

## OUTPUT:

**Entire data is written in GitHub Readme file. Here is the Screenshot of the output.**