

1. Attention Mechanisms for Sequence Modelling

Sequence-to-Sequence is a powerful paradigm of formulating machine learning problems. Broadly, as long as we can formulate a problem as a mapping from a sequence of inputs to a sequence of outputs, we can use sequence-to-sequence models to solve it. For example, in machine translation, we can formulate the problem as a mapping from a sequence of words in one language to a sequence of words in another language. While some RNN architectures we previously covered possess the capability to maintain a memory of the previous inputs/outputs, to compute output, the memory states need to encompass information of many previous states, which can be difficult especially when performing tasks with long-term dependencies.

To understand the limitations of vanilla RNN architectures, we consider the task of changing the case of a sentence, given a prompt token. For example, given a mixed case sequence like “<U> I am a student”, the model should identify this as an upper-case task based on token <U>, and convert it to “I AM A STUDENT”. Similarly, given “<L> I am a student”, the lower-case task is to convert it to “i am a student”.

We can formulate this task as a character-level sequence-to-sequence problem, where the input sequence is the mixed case sentence, and the output sequence is the desired case sentence. In this exercise, we use an encoder-decoder architecture to solve the task. The encoder is a vanilla RNN that takes the input sequence as input, and outputs a sequence of hidden states. The decoder is also a vanilla RNN that takes the last hidden state from the encoder as input, and outputs the desired case sentence.

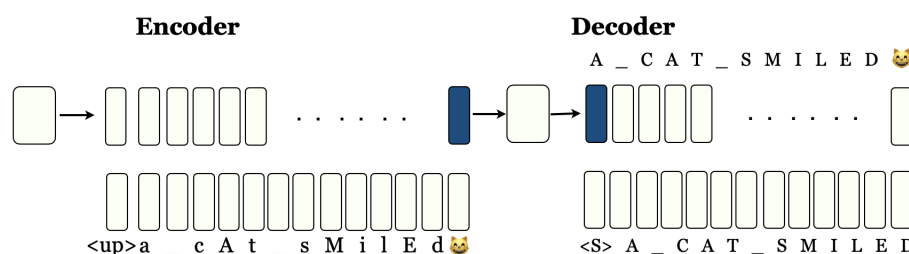


Figure 1: String Manipulation as a Sequence-to-Sequence Problem

(a) What information do RNNs store?

It is important to understand how information propagates through RNNs. Particularly in the context of sequence-to-sequence models, we want to understand what information is stored in the hidden states, and what information is stored in the weights (encoder & decoder). To understand this, we consider the different components of the RNN architecture.

- **Input sequence:** The input sequence is a sequence of T tokens.
- **Encoder Weights:** The shared learnable parameters of the encoder, $W_{\text{enc}} \in \mathbb{R}^{d \times h}$
- **Bottleneck Activations:** The encoder hidden state at time T , that is passes to decoder.
- **Output sequence:** The output sequence is a sequence of T vectors (might be different length).
- **Decoder Weights:** The shared learnable parameters of the decoder, $W_{\text{dec}} \in \mathbb{R}^{d \times h}$

Consider the following questions in the context of these modules:

- Which of these components change during inference?
- When performing gradient based updates, how are the decoder weights trained? How is gradient propagated through the encoder?
- During training, what is the role of the input/output sequence?

(b) Information Bottleneck in RNNs

Consider the architecture shown in Figure 1. This is a simple encoder-decoder architecture with a single hidden layer in the encoder and decoder. The encoder takes the input sequence as input, and outputs a sequence of hidden states. The decoder takes the last hidden state from the encoder as input, and outputs the desired case sentence. **What information needs to be stored in the hidden state to perform the upper-case/lower-case task? Are there any limitations to this architecture?**

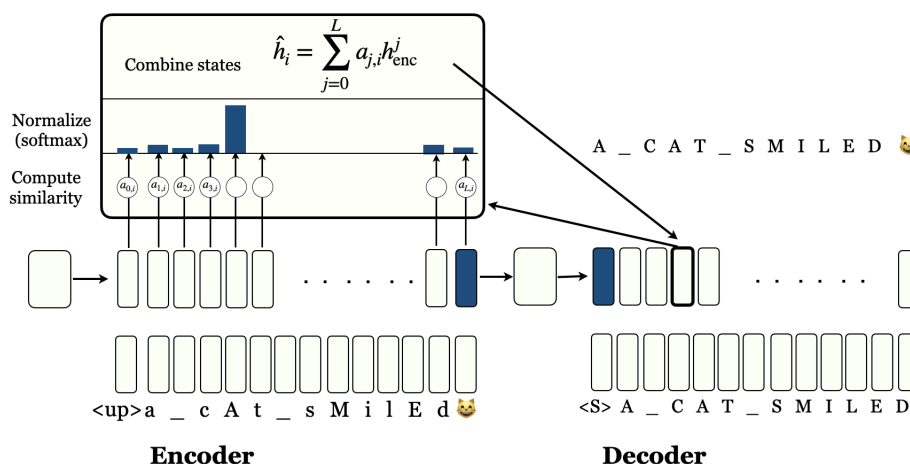


Figure 2: Attention Mechanism for Sequence Modelling with RNNs.

(c) Introducing Attention

Instead of storing all the information in the hidden state, we can use attention to selectively store information. The idea of attention is to query the encoder hidden states with a query vector, and use the resulting attention weights to compute a weighted sum of the hidden states. This weighted sum is then used as the input to the readout layer that computes the output token at each time step of the decoder. **How would you modify the encoder-decoder architecture to incorporate attention?**

(d) Attention & RNNs

How does adding attention allow the model to bypass the information bottleneck? In particular, what information in the following modules would allow the model to perform the capitalization task ?

- **Encoder Weights**
- **Attention Scores**
- **Bottleneck Activations**
- **Decoder Weights**

(e) Positional Encoding

As noted above, the hidden state of the encoder at position t should contain information about the position of token in the input sequence. To incorporate this information, we can add a *positional encoding* to the input tokens. **For sequences of length T discuss how you would add positional encoding to the input sequence.**

2. Query-Key-Value Mechanics in Self-Attention

Self-attention is the core building block for the Transformer model, which has kickstarted the amazing progress in recent deep learning foundation models. The concept of self-attention is not restricted to Transformer exclusively though. In this question, you will be studying the detailed mechanics of the Query-Key-Value interaction in a self-attention block, in the context of RNN. Here we will be studying the case where the encoder only takes in 3 inputs, with the variables defined as:

$$h_1 = \begin{bmatrix} 5 \\ 1 \end{bmatrix} \quad h_2 = \begin{bmatrix} 4 \\ 3 \end{bmatrix} \quad h_3 = \begin{bmatrix} -2 \\ 6 \end{bmatrix} \quad W_q = \begin{bmatrix} 1 & 0 \\ 2 & 9 \end{bmatrix} \quad W_k = \begin{bmatrix} 0 & 1 \\ 6 & 0 \end{bmatrix} \quad W_v = \begin{bmatrix} 4 & 3 \\ 9 & 1 \end{bmatrix}$$

where h_i denotes the hidden states at timestep i at some arbitrary self-attention layer. Each q_i, k_i, v_i is determined by $q_i = W_q h_i, k_i = W_k h_i, v_i = W_v h_i$.

(a) **Compute** $\hat{\alpha}_{2,1}, \hat{\alpha}_{2,2}, \hat{\alpha}_{2,3}$.

(b) Fig 3 shows a rough sketch of how self-attention is performed in one timestep of a Transformer block.

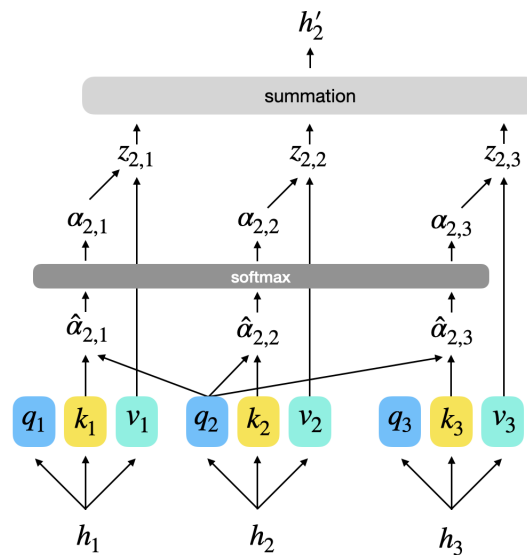


Figure 3: Self-attention of a timestep in Encoder

Write out these operations in equations, ie. $h'_2 = \text{SelfAttention}(h_1, h_2, h_3)$, what is *SelfAttention*? You can define intermediate variables instead of expressing everything in one line.

(c) For simplicity, let's use **argmax** instead of the softmax layer in the diagram. What is h'_2 in this case?

(d) In practice, it is common to have multiple self-attention operations happening in one layer. Each self-attention block is referred to as a *head*, and thus the entire block is typically called *Multi-head Self-Attention (MSA)* in papers. **What's the benefit of having multiple heads?** (*Hint: why do we want multiple kernel filters in ConvNets?*)

Contributors:

- Kumar Krishna Agrawal.
- Kevin Li.