

# **4. RNN: Q&A session**

**CS-E4890: Deep Learning**



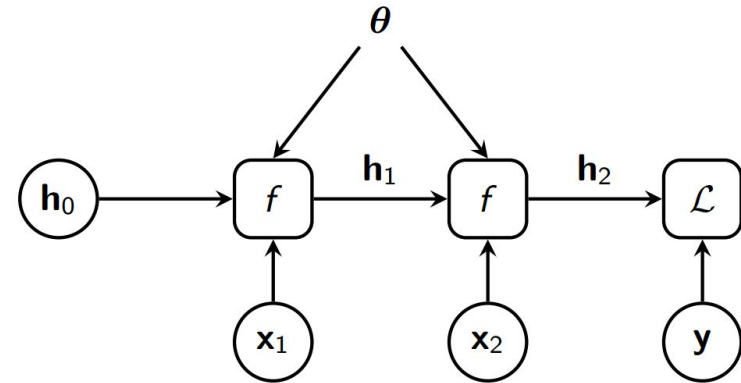
Aalto University  
School of Science

Kalle Kujanpää

24.03.2023

# RNNs: Short Recap

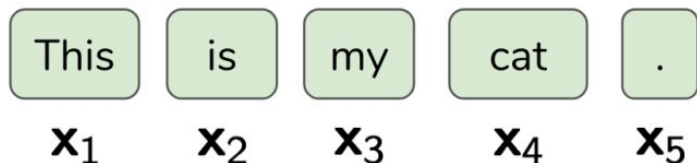
- A family of NNs for handling sequential data with variable length inputs and outputs
- Traditionally used for
  - Natural language processing
  - Speech recognition
  - Time series prediction
  - Reinforcement learning



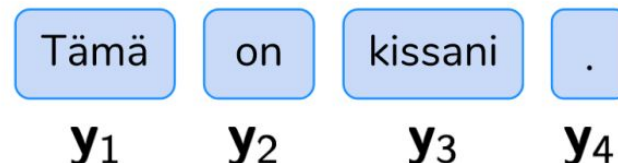
# Neural Machine Translation

- Translate a sentence from source language to target language
- A sequence-to-sequence model, where input and output sequences may be of different lengths

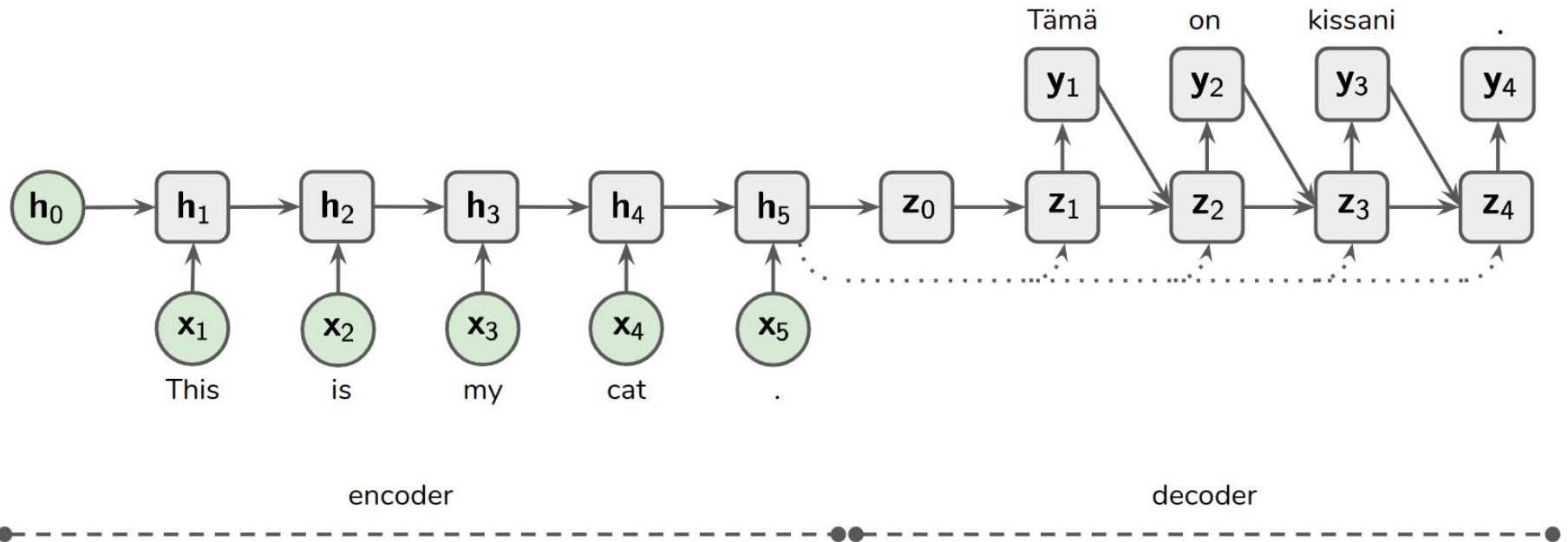
Input: a sequence of words (from the source language)



Output: is a sequence of words (from the target language)

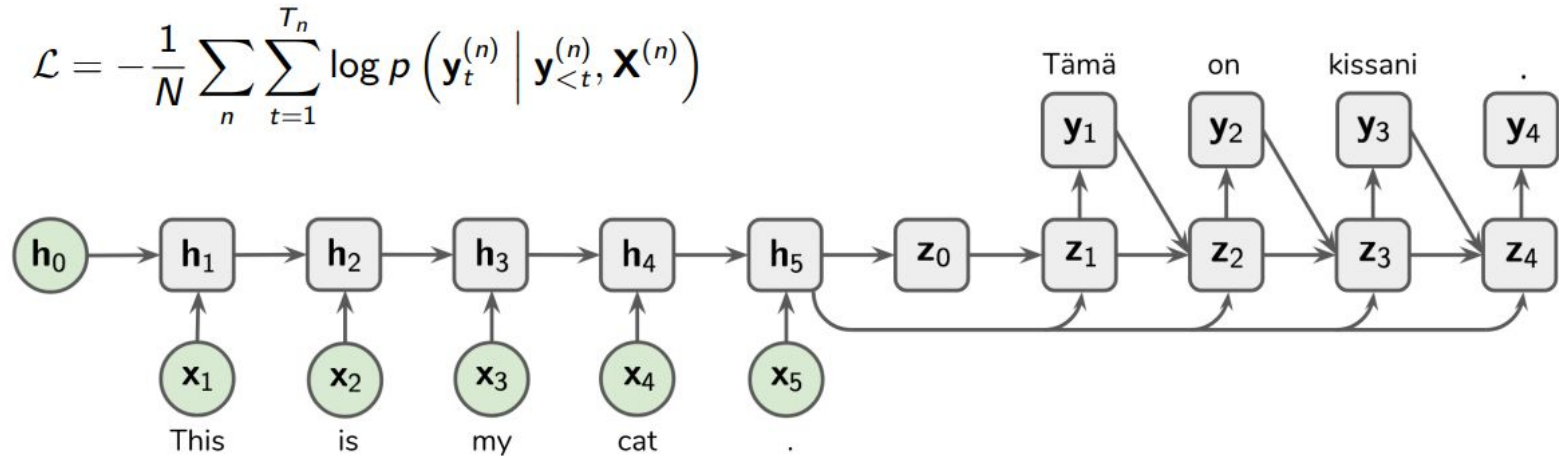


# Sequence-to-sequence model



# Training the Model

- Minimize the negative log-likelihood of the output sequence
- A categorical distribution over the words at each step



# RNNs in PyTorch (from lecture)

- There are two way to build a computational graph with RNNs in PyTorch.
- In simple cases, the whole sequence can be processed with one call:

```
h = torch.zeros(...)  
h = rnn.forward(x, h)
```

- In more difficult cases, you need to build a graph with a for-loop:

```
h = torch.zeros(...)  
for x_t in x:  
    h = rnn.forward(x_t, h)
```

- The initial states of RNNs are often initialized with zeros.

# Word Embeddings (from lecture)

- A simple word representation is one-hot vector. Word  $i$  is represented with vector  $\mathbf{z}$  such that  $z_i = 1, z_{j \neq i} = 0$ .
- Better representation:
  - represent each word  $i$  as a vector  $\mathbf{w}_i$
  - treat all vectors  $\mathbf{w}_i$  as model parameters and tune them in the training procedure
  - this is equivalent to  $\mathbf{W}\mathbf{z}$  where  $\mathbf{W}$  is a matrix of word embeddings (word vectors  $\mathbf{w}_i$  in its columns).
- This is implemented in `torch.nn.Embedding(num_embeddings, embedding_dim)`
  - `num_embeddings` is the size of the dictionary
  - `embedding_dim` is the size of each embedding vector  $\mathbf{w}_i$

# 04\_rnn: Dataset

- French  $\leftrightarrow$  English translation
- Each word is represented by an integer index
  - 4489 French words
  - 2925 English words
  - 8682 Sentences

```
Source sentence: "vous etes depourvues d ambition . EOS"
Sentence as tensor of word indices:
tensor([ 118,  215, 1542,  234, 1209,    5,    1])
Target sentence: "you re unambitious . EOS"
Sentence as tensor of word indices:
tensor([130,  78, 669,   4,   1])
```



# 04\_rnn: 1<sup>st</sup> task – collate

```
def collate(list_of_samples):  
    """Merges a list of samples to form a mini-batch.  
  
    Args:  
        list_of_samples is a list of tuples (src_seq, tgt_seq):  
            src_seq is of shape (src_seq_length,  
            tgt_seq is of shape (tgt_seq_length,  
  
    Returns:  
        src_seqs of shape (max_src_seq_length, batch_size): Tensor of padded source sequences.  
            The sequences should be sorted by length in a decreasing order, that is src_seqs[:,0] should be  
            the longest sequence, and src_seqs[:,-1] should be the shortest.  
        src_seq_lengths: List of lengths of source sequences.  
        tgt_seqs of shape (max_tgt_seq_length, batch_size): Tensor of padded target sequences.  
    """
```

Inputs:

```
[  
    ([1, 2], [3, 4, 5]),  
    ([6, 7, 8], [9, 10]),  
]
```

collate:

Outputs:

src\_seqs:

```
[[6, 1],  
 [7, 2],  
 [8, 0]]
```

tgt\_seqs:

```
[[9, 3],  
 [10, 4],  
 [0, 5]]
```

src\_seq\_lengths:

```
[3, 2]
```

# 04\_rnn: 2<sup>nd</sup> task – encoder

```
def forward(self, pad_seqs, seq_lengths, hidden):
```

```
    """
```

```
    Args:
```

```
        pad_seqs of shape (max_seq_length, batch_size): Padded source sequences.
```

```
        seq_lengths: List of sequence lengths.
```

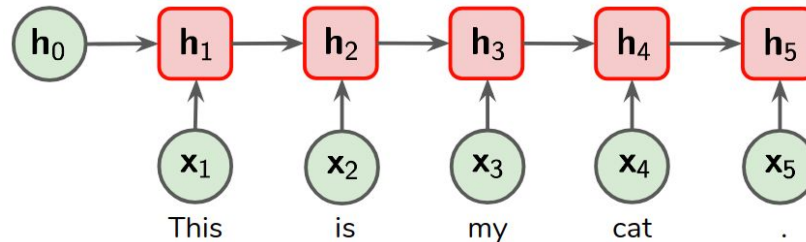
```
        hidden of shape (1, batch_size, hidden_size): Initial states of the GRU.
```

```
    Returns:
```

```
        outputs of shape (max_seq_length, batch_size, hidden_size): Padded outputs of GRU at every step.
```

```
        hidden of shape (1, batch_size, hidden_size): Updated states of the GRU.
```

```
    """
```



# 04\_rnn: 3<sup>rd</sup> task – decoder

```
def forward(self, hidden, pad_tgt_seqs=None, teacher_forcing=False):
```

```
    """
```

Args:

hidden of shape (1, batch\_size, hidden\_size): States of the GRU.

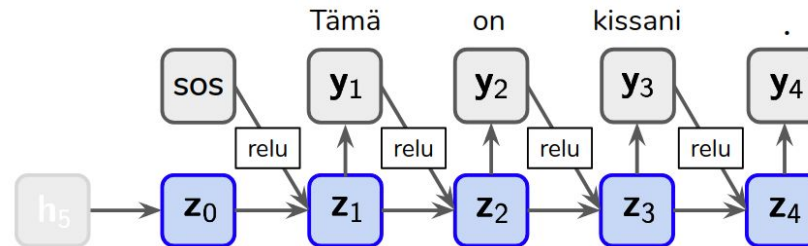
pad\_tgt\_seqs of shape (max\_out\_seq\_length, batch\_size): Tensor of words (word indices) of the target sentence. If None, the output sequence is generated by feeding the decoder's outputs (teacher\_forcing has to be False).

teacher\_forcing (bool): Whether to use teacher forcing or not.

Returns:

outputs of shape (max\_out\_seq\_length, batch\_size, tgt\_dictionary\_size): Tensor of log-probabilities of words in the target language.

hidden of shape (1, batch\_size, hidden\_size): New states of the GRU.



# 04\_rnn: 4<sup>th</sup> task – training loop

- **Loss function: implement log-probabilities**
- **Can be run on a CPU, GPU is faster**
  - I would recommend staying on Jupyter
- **When computing the loss, ignore the padded values**

$$L = -\frac{1}{N} \sum_n \sum_{t=1}^{T_n} \log p \left( \mathbf{y}_t^{(n)} \mid \mathbf{y}_{<t}^{(n)}, \mathbf{X}^{(n)} \right)$$

# 04\_rnn: 5<sup>th</sup> task – translation

```
def translate(encoder, decoder, pad_src_seqs, src_seq_lengths):  
    """Translate sequences from the source language to the target language using the trained model.  
  
    Args:  
        encoder (Encoder): Trained encoder.  
        decoder (Decoder): Trained decoder.  
        pad_src_seqs of shape (max_src_seq_length, batch_size): Padded source sequences.  
        src_seq_lengths: List of source sequence lengths.  
  
    Returns:  
        out_seqs of shape (MAX_LENGTH, batch_size): LongTensor of word indices of the output sequences.  
    """  
    # YOUR CODE HERE  
    raise NotImplementedError()
```