

Deep Learning for Sequential Data (II)

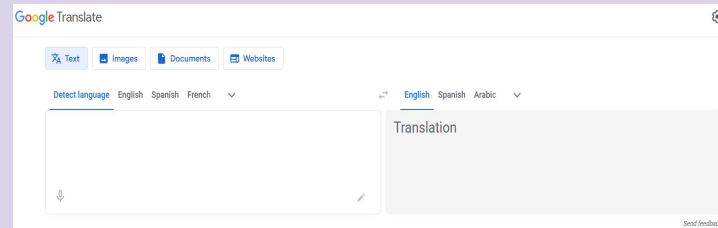
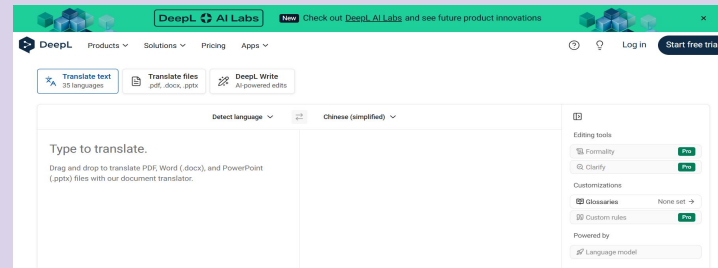
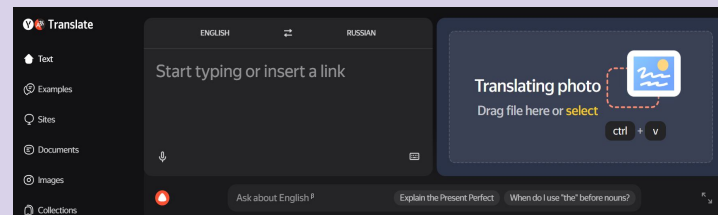
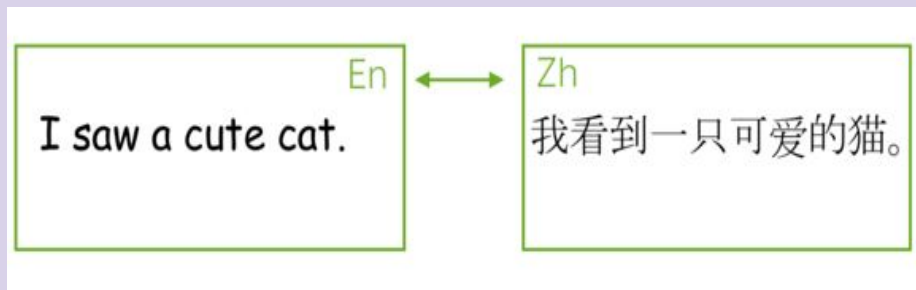
Seq2Seq and Transformers

Learning Outcomes

By the end of this week, you'll be able to:

- Understand Seq2Seq and how they might be applied to solve problems such as machine translation, language modeling, and image captioning.
- Understand local and global attention mechanisms in deep models.
- Understand Transformer and BERT, two latest SOTA models for NLP.
- Understand the self-attention mechanism and how it is used in Transformer.

Use Case



Use Case

Meta AI Research The Latest AI Research About Get Llama

initialization, language modeling, and back translation.

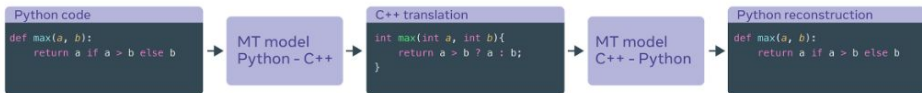
Cross-lingual masked model pretraining



Denoising auto-encoding



Back - translation



This graphic shows how TransCoder leverages the three principles of unsupervised machine translation.

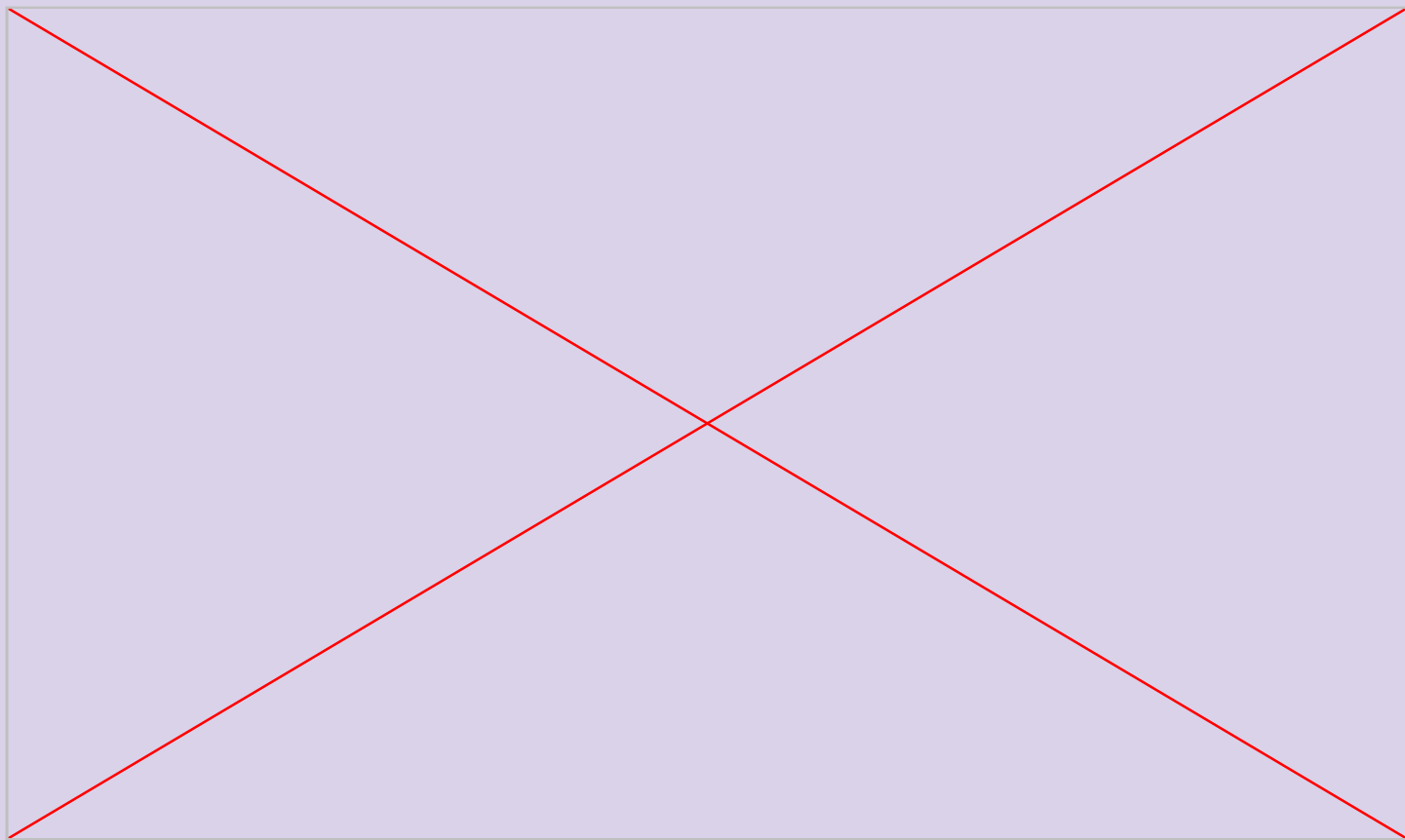
Python input

```
def SumOfKSubArray(arr, n, k):  
    Sum = 0  
    S = deque()  
    G = deque()  
    for i in range(k):  
        while (len(S) > 0 and arr[S[-1]] >= arr[i]):  
            S.pop()  
        while (len(G) > 0 and arr[G[-1]] <= arr[i]):  
            G.pop()  
        S.append(i)  
        G.append(i)  
    for i in range(k, n):  
        Sum += arr[S[0]] + arr[G[0]]  
        while (len(S) > 0 and S[0] <= i - k):  
            S.popleft()  
        while (len(G) > 0 and G[0] <= i - k):  
            G.popleft()  
        while (len(S) > 0 and arr[S[-1]] >= arr[i]):  
            S.pop()  
        while (len(G) > 0 and arr[G[-1]] <= arr[i]):  
            G.pop()  
        S.append(i)  
        G.append(i)  
    Sum += arr[S[0]] + arr[G[0]]  
    return Sum
```

C++ unsupervised translation

```
int SumOfKSubArray(int arr[], int n, int k){  
    int Sum = 0;  
    deque<int> S;  
    deque<int> G;  
    for(int i = 0; i < k; i++){  
        while((!S.empty() && arr[S.back()] >= arr[i])  
            S.pop_back();  
        while((!G.empty() && arr[G.back()] <= arr[i])  
            G.pop_back();  
        S.push_back(i);  
        G.push_back(i);  
    }  
    for(int i = k; i < n; i++){  
        Sum += arr[S.front()] + arr[G.front()];  
        while((!S.empty() && S.front() <= i - k)  
            S.pop_front();  
        while((!G.empty() && G.front() <= i - k)  
            G.pop_front();  
        while((!S.empty() && arr[S.back()] >= arr[i])  
            S.pop_back();  
        while((!G.empty() && arr[G.back()] <= arr[i])  
            G.pop_back();  
        S.push_back(i);  
        G.push_back(i);  
    }  
    Sum += arr[S.front()] + arr[G.front()];  
    return Sum;  
}
```

Sequence and Sequence



Sequence and Sequence

Source sequence:

Я видел котю на мате <eos>
"I" "saw" "cat" "on" "mat"

Target sequence:

I saw a cat on a mat <eos>
previous tokens we want the model to predict this

← one training example

← one step for this example

Model prediction: $p(* | \text{I saw a, Я ... <eos>})$

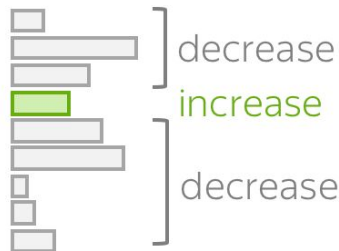


← cat →

Target

0
0
0
0
1
0
0
0
0
0
0

Loss = $-\log(p(\text{cat})) \rightarrow \min$



Encoder

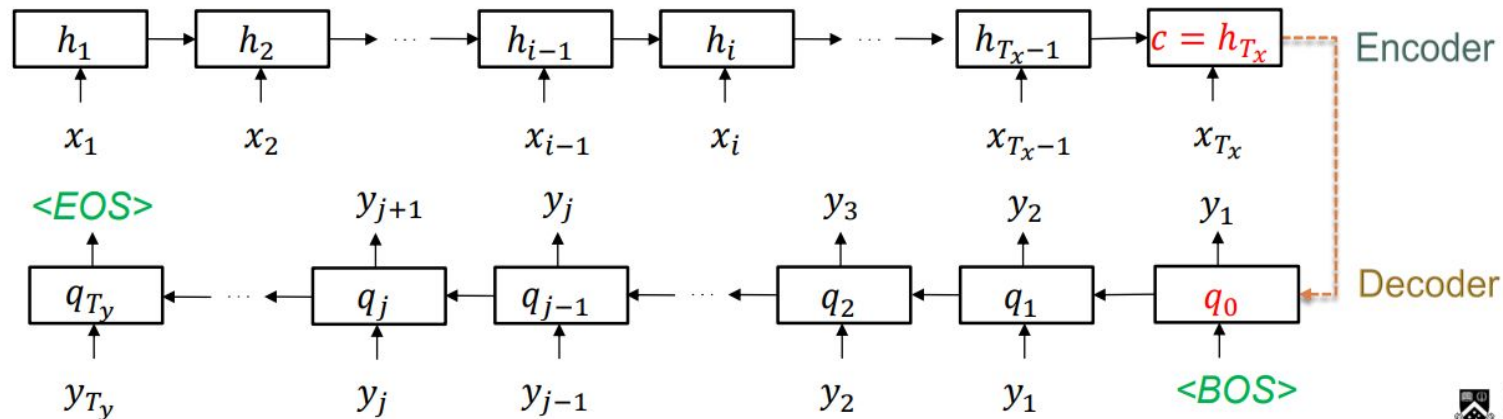
- Produces the context vector $\mathbf{c} = \mathbf{h}_{T_x}$ of the input sequence
- Context vector \mathbf{c} summarizes input sequence $[\mathbf{x}_1, \dots, \mathbf{x}_{T_x}]$.

Decoder

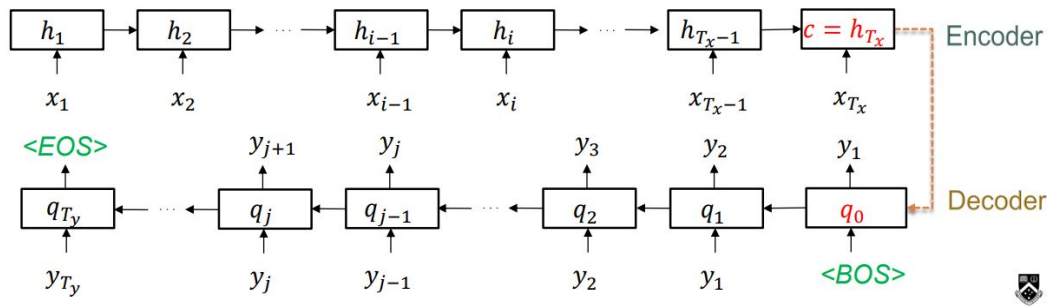
- Decodes the encode \mathbf{c} to the output sequence

Special symbols

- $\langle \text{EOS} \rangle$ signifies the end of a sequence
- $\langle \text{BOS} \rangle$ signifies the beginning of a sequence



Step 1: Dataset Preparation



```
dataset = datasets.load_dataset("bentrevett/multi30k")
```

{'en': 'Two young, White males are outside near many bushes.',

'de': 'Zwei junge weiße Männer sind im Freien in der Nähe vieler Büsche.'}

Statistics

train

(en) 29000 sentences, 377534 words, 13.0 words/sent
(de) 29000 sentences, 360706 words, 12.4 words/sent
(fr) 29000 sentences, 409845 words, 14.1 words/sent
(cs) 29000 sentences, 297212 words, 10.2 words/sent

val

(en) 1014 sentences, 13308 words, 13.1 words/sent
(de) 1014 sentences, 12828 words, 12.7 words/sent
(fr) 1014 sentences, 14381 words, 14.2 words/sent
(cs) 1014 sentences, 10342 words, 10.2 words/sent

test_2016_flickr

(en) 1000 sentences, 12968 words, 13.0 words/sent
(de) 1000 sentences, 12103 words, 12.1 words/sent
(fr) 1000 sentences, 13988 words, 14.0 words/sent
(cs) 1000 sentences, 10497 words, 10.5 words/sent

test_2017_flickr

(en) 1000 sentences, 11376 words, 11.4 words/sent
(de) 1000 sentences, 10758 words, 10.8 words/sent
(fr) 1000 sentences, 12596 words, 12.6 words/sent

test_2017_mscoco

(en) 461 sentences, 5239 words, 11.4 words/sent
(de) 461 sentences, 5158 words, 11.2 words/sent
(fr) 461 sentences, 5710 words, 12.4 words/sent

Step 2: Tokenization

Text Analytics and Language Models

□ Text normalization

○ Expanding contractions

- **Contractions** are shortened version of words or syllables
- e.g., isn't → is not, you're → you are
- Exist extensively and pose a problem to text analytics

○ Lemmatization

- removing word **affixes** to get to a **base form** of the *root* word.
- e.g. cars → car, running → run, is → be

○ Removing special characters and symbols

- e.g. !, .

○ Removing stop words

- e.g., a, and

Distributed representation of words

$$\text{employees} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 10.109 \\ -0.542 \\ 0.349 \\ 0.271 \\ 0.487 \end{pmatrix}$$

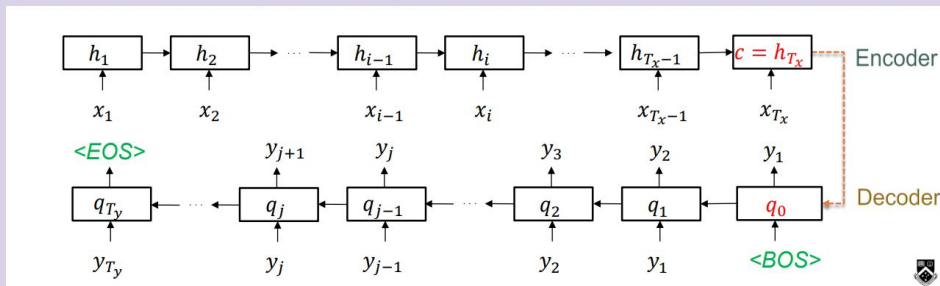


```
import spacy
```

1. A tokenizer is used to turn a string into a list of tokens

2. "good morning!" becomes ["good", "morning", "!"].

Step 2: Tokenization

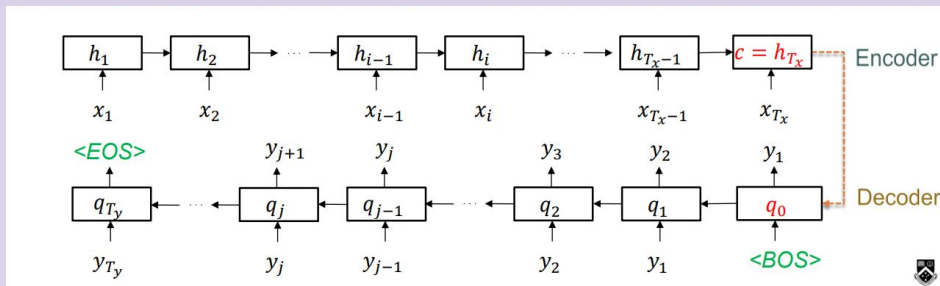


```
import spacy
```

```
1 en_nlp = spacy.load("en_core_web_sm")  
2 de_nlp = spacy.load("de_core_news_sm")
```

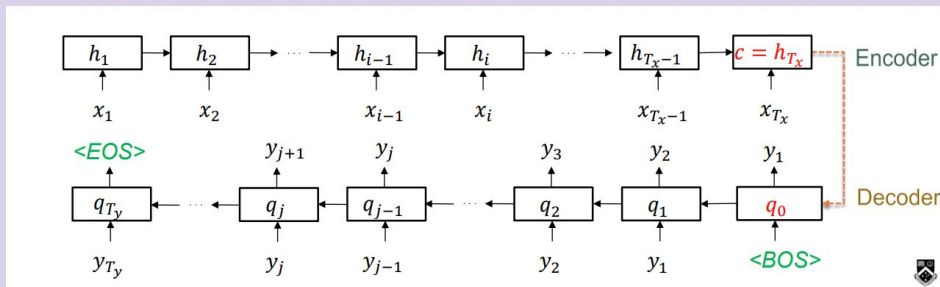
```
1 string = "What a lovely day it is today!"  
2  
3 [token.text for token in en_nlp.tokenizer(string)]  
  
['What', 'a', 'lovely', 'day', 'it', 'is', 'today', '!']
```

Step 2: Tokenization



```
1 max_length = 1_000
2 lower = True
3 sos_token = "<sos>"
4 eos_token = "<eos>"
5
6 fn_kwargs = {
7     "en_nlp": en_nlp,
8     "de_nlp": de_nlp,
9     "max_length": max_length,
10    "lower": lower,
11    "sos_token": sos_token,
12    "eos_token": eos_token,
13 }
14
15 train_data = train_data.map(tokenize_example, fn_kwargs=fn_kwargs)
16 valid_data = valid_data.map(tokenize_example, fn_kwargs=fn_kwargs)
17 test_data = test_data.map(tokenize_example, fn_kwargs=fn_kwargs)
```

Step 2: Tokenization



```
{'en': 'Two young, White males are outside near many bushes.',  
'de': 'Zwei junge weiße Männer sind im Freien in der Nähe vieler Büsche.',  
'en_tokens': ['<sos>',  
'two',  
'young',  
,  
'white',  
'males',  
'are',  
'outside',  
'near',  
'many',  
'bushes',  
,  
'<eos>'],  
'de_tokens': ['<sos>',  
'zwei',  
'junge',  
'weiße',  
'männer',  
'sind',  
'im',  
'freien',  
'in',  
'der',  
'nähe',  
'vieler',  
'büsche',  
,  
'<eos>']}]}
```

Step 2.1 : Vocabulary

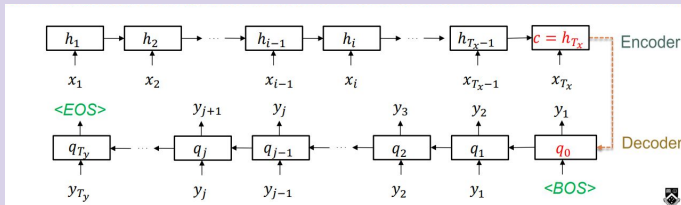
Vocabularies

Next, we'll build the *vocabulary* for the source and target languages. The vocabulary is used to associate each unique token in our dataset with an index (an integer), e.g. "hello" = 1, "world" = 2, "bye" = 3, "hates" = 4, etc

```
{'en': 'Two young, White males are outside near many bushes.',  
 'de': 'Zwei junge weiße Männer sind im Freien in der Nähe vieler Büsche.',  
 'en_tokens': ['<sos>',  
               'two',  
               'young',  
               ',',  
               'white',  
               'males',  
               'are',  
               'outside',  
               'near',  
               'many',  
               'bushes',  
               '.',  
               '<eos>'],  
 'de_tokens': ['<sos>',  
               'zwei',  
               'junge',  
               'weiße',  
               'männer',  
               'sind',  
               'im',  
               'freien',  
               'in',  
               'der',  
               'nähe',  
               'vieler',  
               'büsche',  
               '.',  
               '<eos>']}]
```

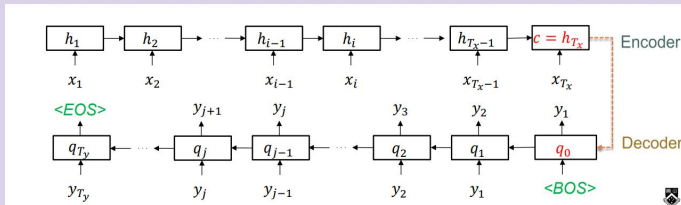
```
1 min_freq = 2  
2 unk_token = "<unk>"  
3 pad_token = "<pad>"  
4  
5 special_tokens = [  
6     unk_token,  
7     pad_token,  
8     sos_token,  
9     eos_token,  
10 ]  
11  
12 en_vocab = torchtext.vocab.build_vocab_from_iterator(  
13     train_data["en_tokens"],  
14     min_freq=min_freq,  
15     specials=special_tokens,  
16 )  
17  
18 de_vocab = torchtext.vocab.build_vocab_from_iterator(  
19     train_data["de_tokens"],  
20     min_freq=min_freq,  
21     specials=special_tokens,  
22 )
```

Step 2.2: Encoder



```
1 class Encoder(nn.Module):
2     def __init__(self, input_dim, embedding_dim, hidden_dim, n_layers, dropout):
3         super().__init__()
4         self.hidden_dim = hidden_dim
5         self.n_layers = n_layers
6         self.embedding = nn.Embedding(input_dim, embedding_dim)
7         self.rnn = nn.LSTM(embedding_dim, hidden_dim, n_layers, dropout=dropout)
8         self.dropout = nn.Dropout(dropout)
9
10    def forward(self, src):
11        # src = [src length, batch size]
12        embedded = self.dropout(self.embedding(src))
13        # embedded = [src length, batch size, embedding dim]
14        outputs, (hidden, cell) = self.rnn(embedded)
15        # outputs = [src length, batch size, hidden dim * n directions]
16        # hidden = [n layers * n directions, batch size, hidden dim]
17        # cell = [n layers * n directions, batch size, hidden dim]
18        # outputs are always from the top hidden layer
19        return hidden, cell
```

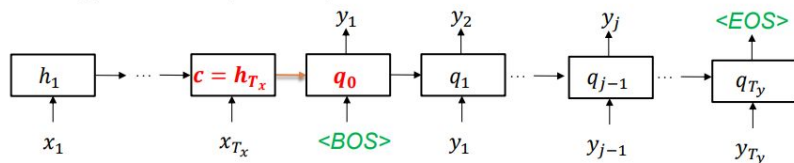
Step 2.4: Decoder



```
1 class Decoder(nn.Module):
2     def __init__(self, output_dim, embedding_dim, hidden_dim, n_layers, dropout):
3         super().__init__()
4         self.output_dim = output_dim
5         self.hidden_dim = hidden_dim
6         self.n_layers = n_layers
7         self.embedding = nn.Embedding(output_dim, embedding_dim)
8         self.rnn = nn.LSTM(embedding_dim, hidden_dim, n_layers, dropout=dropout)
9         self.fc_out = nn.Linear(hidden_dim, output_dim)
10        self.dropout = nn.Dropout(dropout)
11
12    def forward(self, input, hidden, cell):
13        # input = [batch size]
14        # hidden = [n layers * n directions, batch size, hidden dim]
15        # cell = [n layers * n directions, batch size, hidden dim]
16        # n directions in the decoder will both always be 1, therefore:
17        # hidden = [n layers, batch size, hidden dim]
18        # context = [n layers, batch size, hidden dim]
19        input = input.unsqueeze(0)
20        # input = [1, batch size]
21        embedded = self.dropout(self.embedding(input))
22        # embedded = [1, batch size, embedding dim]
23        output, (hidden, cell) = self.rnn(embedded, (hidden, cell))
24        # output = [seq length, batch size, hidden dim * n directions]
25        # hidden = [n layers * n directions, batch size, hidden dim]
26        # cell = [n layers * n directions, batch size, hidden dim]
27        # seq length and n directions will always be 1 in this decoder, therefore:
28        # output = [1, batch size, hidden dim]
29        # hidden = [n layers, batch size, hidden dim]
30        # cell = [n layers, batch size, hidden dim]
31        prediction = self.fc_out(output.squeeze(0))
32        # prediction = [batch size, output dim]
33        return prediction, hidden, cell
```


Step 3: Training

Training of seq2seq



- We need to maximize the log-likelihood:

$$\max_{\theta} J(\theta) = \sum_{(x,y) \in \mathcal{D}} \log P(y|x, \theta)$$

where $\theta = [\theta_e, \theta_d]$ and θ_e, θ_d are encoding and decoding parameters respectively.

- Product rule:

$$P(y|x, \theta) = P(y_{1:T_y} | x_{1:T_x}, \theta) = P(y_{1:T_y} | c, \theta)$$

$$= P(y_1 | c, \theta) P(y_2 | y_1, c, \theta) \dots P(y_j | y_{1:j-1}, c, \theta) \dots P(y_{T_y} | y_{1:T_y-1}, c, \theta) = \prod_{j=1}^{T_y} P(y_j | y_{1:j-1}, c, \theta)$$

$$\log P(y|x, \theta) = \log P(y|c, \theta) = \sum_{j=1}^{T_y} \log P(y_j | y_{1:j-1}, c, \theta) = \sum_{j=1}^{T_y} \log P(y_j | q_{j-1}, c, \theta)$$

- We can compute $P(y_j | q_{j-1}, c) = g(y_j, q_{j-1}, c)$ where g is a **nonlinear**, potentially **multi-layered NN** that outputs the probability of y_j .

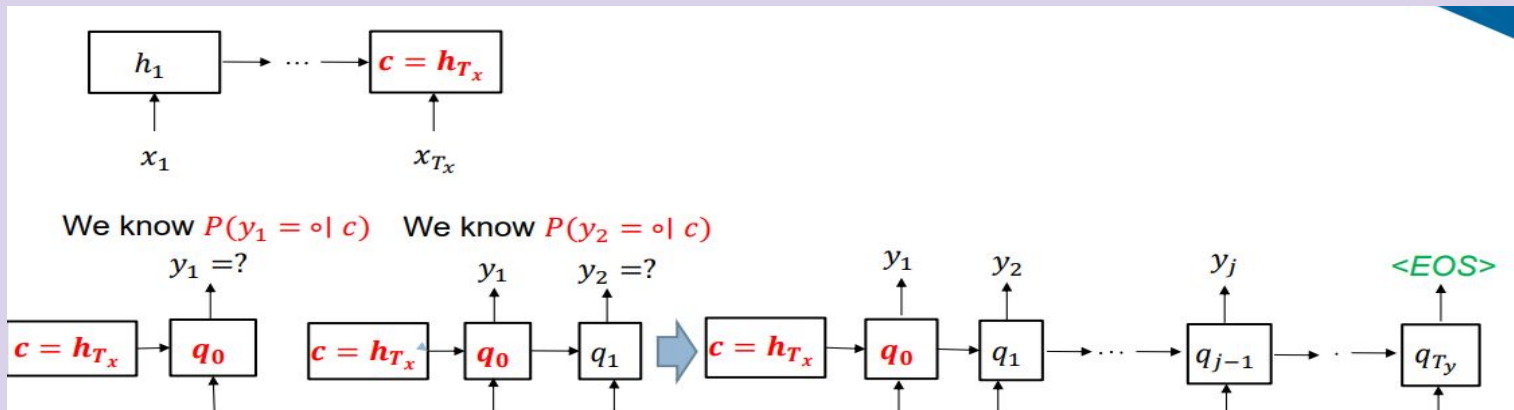
- Pay attention on how c is used in every step during decoding

```

1 input_dim = len(de_vocab)
2 output_dim = len(en_vocab)
3 encoder_embedding_dim = 256
4 decoder_embedding_dim = 256
5 hidden_dim = 512
6 n_layers = 2
7 encoder_dropout = 0.5
8 decoder_dropout = 0.5
9 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
10
11 encoder = Encoder(
12     input_dim,
13     encoder_embedding_dim,
14     hidden_dim,
15     n_layers,
16     encoder_dropout,
17 )
18
19 decoder = Decoder(
20     output_dim,
21     decoder_embedding_dim,
22     hidden_dim,
23     n_layers,
24     decoder_dropout,
25 )
26
27 model = Seq2Seq(encoder, decoder, device).to(device)

```

Step 4: Inference



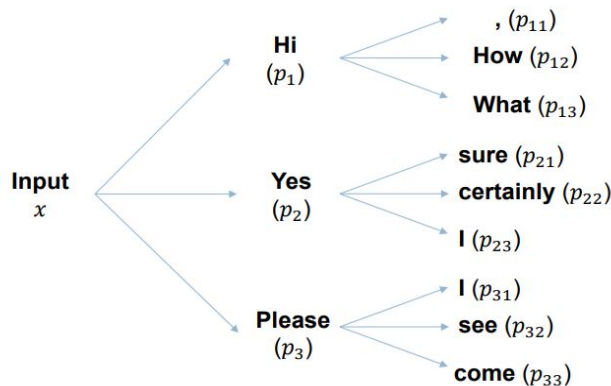
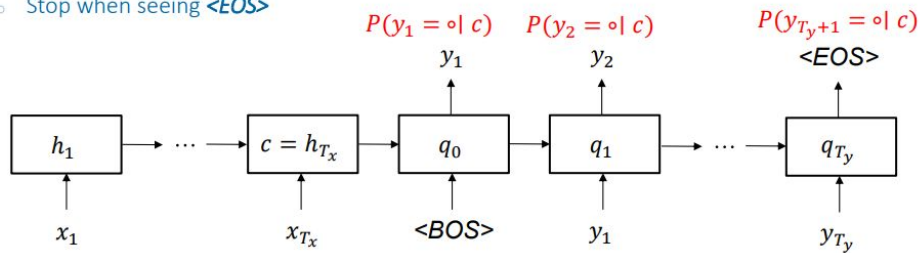
```
1 model.load_state_dict(torch.load("tut1-model.pt"))
2
3 test_loss = evaluate_fn(model, test_data_loader, criterion, device)
4
5 print(f"| Test Loss: {test_loss:.3f} | Test PPL: {np.exp(test_loss):7.3f} |")
```

Test Loss: 3.780 | Test PPL: 43.833 |

Step 4.1: Inference; Greedy Vs Beam

Greedy Decoding

- Given \mathbf{x} , find word \mathbf{y}_1 with highest probability
- Given \mathbf{y}_1 and \mathbf{x} , find word \mathbf{y}_2 with highest probability
- ...
- Stop when seeing $\langle \text{EOS} \rangle$



Beam width = 3

- We **always** choose three sentences with highest joint probabilities

Drawback of fixed context

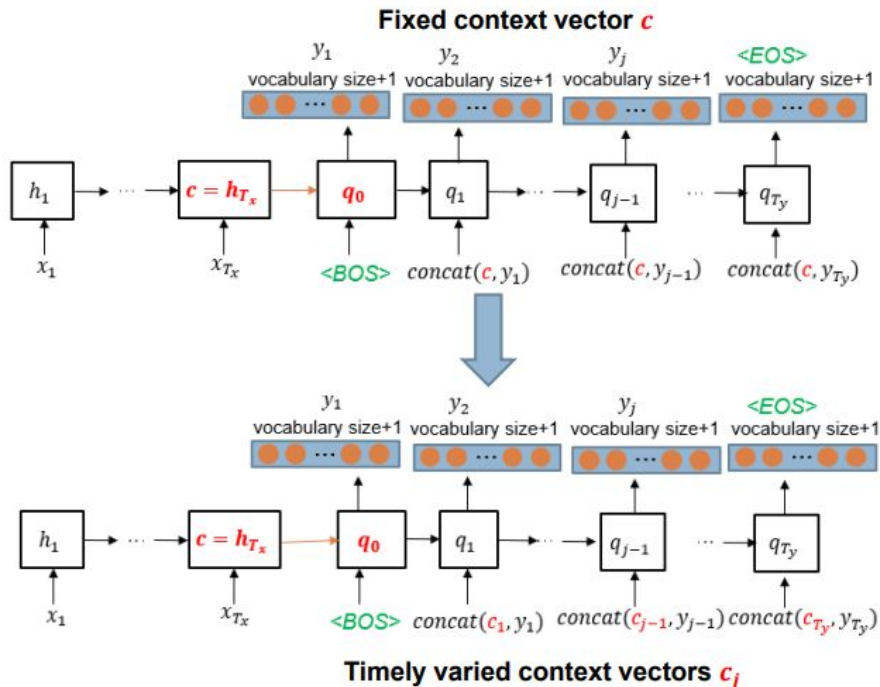
- Fixed context vector c is easily overwhelmed by long inputs or long outputs.
- At a specific timestep j , some words or items in the input sequence might possibly contribute more to the generation of next item or word in the output sequence.
 - I want to see you every day \rightarrow Je veux te voir chaque jour
 - I want to see **you** every day \rightarrow Je veux te ? (voir)
- How to timely adapt the context vector c_j ?
 - $c_j = \alpha(h_1, \dots, h_{T_x}, q_{j-1})$
 - Computed using attention mechanism

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau
Jacobs University Bremen, Germany

KyungHyun Cho Yoshua Bengio*
Université de Montréal

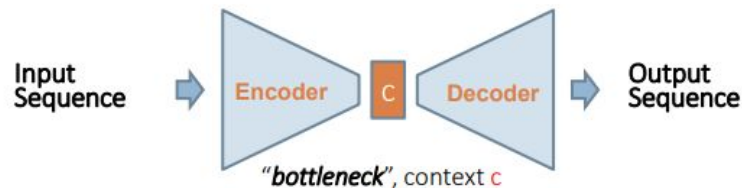
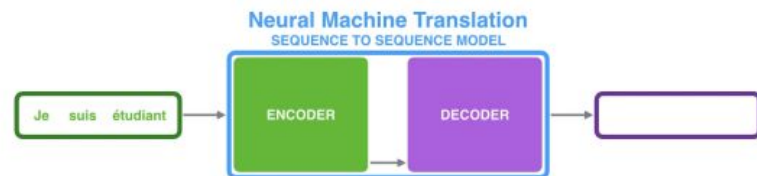
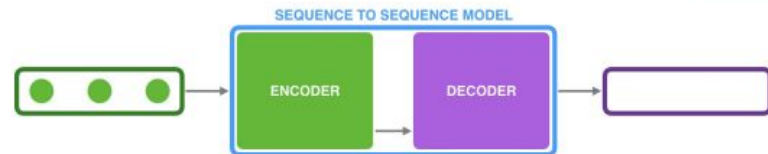
Paper: [paper link](#)



Attention is all you need

Attention mechanism

- So far, the input sequence is **summarised** by a **single** context **c** !
- Fixed-length context** could be problematic as it is **easily overwhelmed** by long inputs or long outputs
 - The **fixed-length** context **c** might **not be powerful enough** to capture long input sequences
- Some **specific items** in an input sequence might be **more relevant** and **contributing** in **generating** a given item in output sequence.



- Gratefully acknowledge the excellent visualizations used from Jay Alammar blog at:
 - <https://jalammar.github.io/>
 - <https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>

Attention mechanism

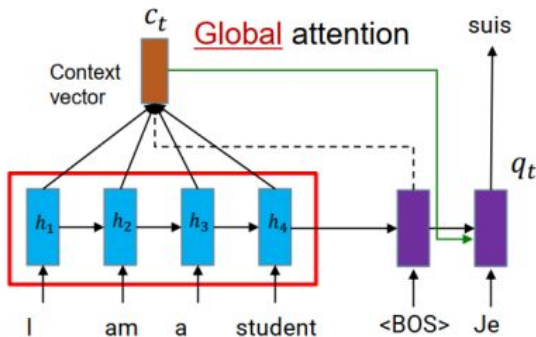
Published as a conference paper at ICLR 2015

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dmitry Bahdanau
Jacobs University Bremen, Germany

KyungHyun Cho Yoshua Bengio*
Université de Montréal

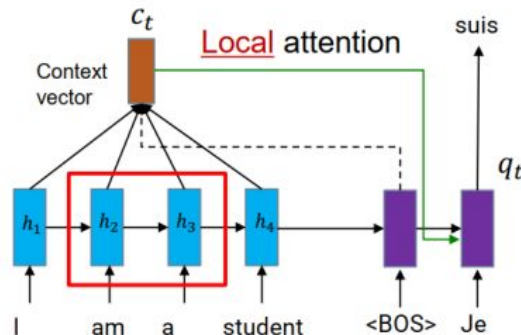
Bahdanau, Cho, Bengio, [Neural Machine Translation by Jointly Learning to Align and Translate](#), ICLR 2015



Effective Approaches to Attention-based Neural Machine Translation

Minh-Thang Luong Hieu Pham Christopher D. Manning
Computer Science Department, Stanford University, Stanford, CA 94305
{lmthang, hyhieu, manning}@stanford.edu

Luong, Pham, Manning, [Effective Approach Attention-based Neural Machine Translation](#), EMNLP, 2015



● Attention mechanism allows the decoding network to refer to the input.

- Global attention

- Use **all input hidden states** of the encoder when deriving the context c_t .

- Local attention

- Use a **selective window of input hidden states** of the encoder when deriving the context c_t .