

# Natural Language Processing Course

## Week 4: Sequence-to-Sequence Models

Joerg R. Osterrieder  
[www.joergosterrieder.com](http://www.joergosterrieder.com)

## Week 4

# Sequence-to-Sequence Models

Breaking the Fixed-Length Barrier

# Why Google Translate Struggled for Years

## The Challenge:

English: "I love you" (3 words)

French: "Je t'aime" (2 words)

German: "Ich liebe dich" (3 words)

Japanese: "aishiteru" (1 word)

Different languages express the same idea with different lengths!

**The Problem:** RNNs produce one output per input

**The Solution:** Separate encoding from decoding

*This breakthrough improved Google Translate accuracy by 60% in 2016<sup>1</sup>*

---

<sup>1</sup>Wu et al. (2016). "Google's Neural Machine Translation System", arXiv

# Seq2Seq Powers Your Daily Interactions

## Translation (2024):

- Google: 1B+ translations daily<sup>2</sup>
- DeepL: Seq2seq + attention
- Real-time conversation mode

## Chatbots:

- Customer service (80% first-line)<sup>3</sup>
- Variable-length responses
- Context-aware replies

## Text Summarization:

- News article → headline
- Email → one-liner
- Document → abstract

## Code Generation:

- Comment → code
- Natural language → SQL
- Bug description → fix

Key Innovation: Input and output can have different lengths!

---

<sup>1</sup>Google Translate statistics 2024

<sup>2</sup>Gartner report on AI customer service

## Week 4: What You'll Master

**By the end of this week, you will:**

- **Understand** why variable-length I/O is crucial
- **Build** intuition for encoder-decoder architecture
- **Implement** a complete seq2seq translator
- **Discover** why attention changes everything
- **Create** a chatbot that handles any conversation length

**Core Insight:** Compress entire input to a "thought", then expand to output

# The Fixed-Length Prison

## What RNNs can do:

- Input: "The cat sat" → Output: "on the mat" (same length)
- Each input produces exactly one output

## What we actually need:

- "Hello" → "Bonjour" (different lengths)
- "How are you?" → "¿Cómo estás?" (different structure)
- Long paragraph → Short summary (compression)

## Failed Approaches:

- 1 Pad everything to max length (wastes computation)
- 2 Truncate to fixed length (loses information)
- 3 Multiple passes (complicated and slow)

We need to decouple input processing from output generation!

# The Brilliant Solution: Encoder-Decoder

## How humans translate:

- 1 Read entire English sentence
- 2 Understand the complete meaning
- 3 Express that meaning in French

## Seq2seq does exactly this:<sup>4</sup>

- 1 **Encoder:** Process entire input into a "thought vector"
- 2 **Context Vector:** Fixed-size representation of meaning
- 3 **Decoder:** Generate output from the thought vector

## The Magic:

- Encoder handles any input length
- Decoder produces any output length
- Middle "thought" is always same size!

Context vector = Compressed understanding of entire input

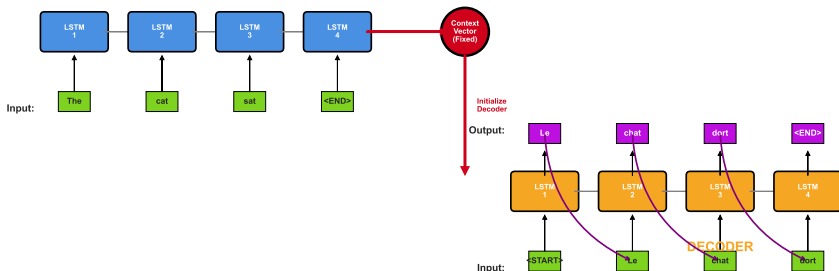
---

<sup>4</sup>Sutskever, Vinyals & Le (2014). "Sequence to Sequence Learning with Neural Networks", NeurIPS

# Visualizing Seq2Seq Architecture

## Sequence-to-Sequence Architecture with Encoder-Decoder

### ENCODER



### Key insights:

- Encoder reads left-to-right (or bidirectional)
- Final hidden state contains entire input meaning
- Decoder generates one word at a time
- Special tokens: START begins generation, END stops it



## Seq2Seq Mathematics: Two RNNs Working Together

**Encoder (processes input):**

$$h_t^{enc} = \text{LSTM}_{enc}(x_t, h_{t-1}^{enc})$$
$$c = h_{final}^{enc} \text{ (context vector)}$$

**Decoder (generates output):**

$$h_0^{dec} = c \text{ (initialize with context)}$$
$$h_t^{dec} = \text{LSTM}_{dec}(y_{t-1}, h_{t-1}^{dec})$$
$$y_t = \text{softmax}(W_y h_t^{dec})$$

**In plain English:**

- Encoder: Compress input into fixed-size thought
- Decoder: Expand thought into variable-length output

# Building Seq2Seq: Complete Implementation

```
1 import torch
2 import torch.nn as nn
3
4 class Encoder(nn.Module):
5     def __init__(self, input_size, hidden_size):
6         """Compress input sequence to fixed-size vector"""
7         super().__init__()
8         self.hidden_size = hidden_size
9         self.embedding = nn.Embedding(input_size, hidden_size)
10        self.lstm = nn.LSTM(hidden_size, hidden_size)
11
12    def forward(self, input_seq):
13        """Process entire input sequence"""
14        embedded = self.embedding(input_seq)
15        outputs, (hidden, cell) = self.lstm(embedded)
16        # Return final hidden state as context
17        return hidden, cell
18
19 class Decoder(nn.Module):
20     def __init__(self, output_size, hidden_size):
21         """Generate output sequence from context"""
22         super().__init__()
23         self.hidden_size = hidden_size
24         self.embedding = nn.Embedding(output_size, hidden_size)
25         self.lstm = nn.LSTM(hidden_size, hidden_size)
26         self.out = nn.Linear(hidden_size, output_size)
27
28    def forward(self, input_token, hidden, cell):
29        """Generate one output token at a time"""
30        embedded = self.embedding(input_token.unsqueeze(0))
31        output, (hidden, cell) = self.lstm(embedded, (hidden, cell))
32        prediction = self.out(output.squeeze(0))
33        return prediction, hidden, cell
```

## Design Choices:

- Hidden size: 256-512 typical<sup>5</sup>
- Bidirectional encoder often better
- Teacher forcing during training

## Usage Pattern: # Encode

context = encoder(input\_seq)

# Decode

output = decoder(context)

---

Britz et al. (2017) extensive comparison

# Complete Seq2Seq Model

```
1 class Seq2Seq(nn.Module):
2     def __init__(self, encoder, decoder, device):
3         """Complete translation model"""
4         super().__init__()
5         self.encoder = encoder
6         self.decoder = decoder
7         self.device = device
8
9     def forward(self, src, trg, teacher_forcing_ratio=0.5):
10        """Translate source to target"""
11        batch_size = src.shape[1]
12        trg_len = trg.shape[0]
13        trg_vocab_size = self.decoder.out.out_features
14
15        # Store decoder outputs
16        outputs = torch.zeros(trg_len, batch_size, trg_vocab_size).
17            to(self.device)
18
19        # Encode entire source sequence
20        hidden, cell = self.encoder(src)
21
22        # First decoder input is <START> token
23        decoder_input = trg[0,:]
24
25        for t in range(1, trg_len):
26            # Decode one token
27            output, hidden, cell = self.decoder(decoder_input,
28                hidden, cell)
29            outputs[t] = output
30
31            # Teacher forcing: use true target or predicted
32            teacher_force = random.random() < teacher_forcing_ratio
33            top1 = output.argmax(1)
34            decoder_input = trg[t] if teacher_force else top1
```

## Training Tricks:

- Teacher forcing prevents drift<sup>6</sup>
- Gradually reduce forcing ratio
- Beam search for better decoding

## Performance (2014):

- BLEU: 34.8 on EN-FR<sup>7</sup>
- Beats phrase-based SMT
- 1000x faster than RNN search

---

<sup>a</sup>Williams & Zipser (1989)

<sup>b</sup>Original seq2seq paper results

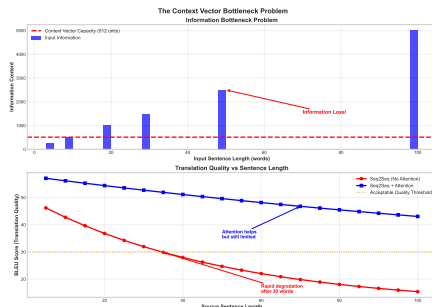
# The Information Bottleneck Problem

Imagine compressing a book into one sentence...

What happens with long sequences:

- 10 words: Context vector remembers well
- 20 words: Starting to forget early words
- 50+ words: Information bottleneck!<sup>8</sup>

The Problem Visualized:



Fixed-size context vector must capture EVERYTHING!

<sup>8</sup>Cho et al. (2014) showed performance degrades after 30 tokens

# The Attention Revolution (2015)

## How humans actually translate:

"The black cat sat on the mat" → "Le chat noir..."

- When translating "chat" (cat), we look back at "cat"
- When translating "noir" (black), we look back at "black"
- We don't compress everything into one thought!

## Attention mechanism:<sup>9</sup>

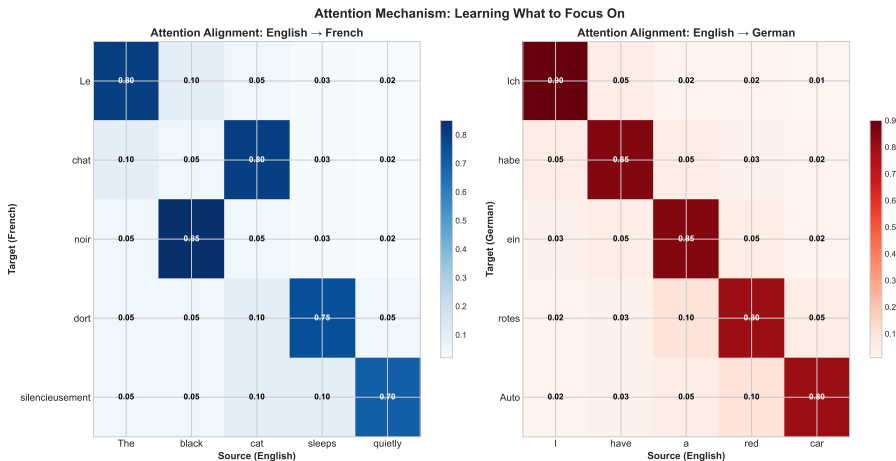
- 1 Keep ALL encoder hidden states
- 2 When decoding, look back at relevant parts
- 3 Weight importance of each input word
- 4 Focus on what matters right now

Attention = Let the decoder peek back at the input whenever needed

---

<sup>9</sup>Bahdanau, Cho & Bengio (2015). "Neural Machine Translation by Jointly Learning to Align and Translate"

# Visualizing Attention: What the Model Looks At



## Key observations:

- Diagonal pattern shows alignment
- Some words need multiple source words
- Reordering handled naturally

# Implementing Attention

```
1 class Attention(nn.Module):
2     def __init__(self, hidden_size):
3         """Calculate attention weights"""
4         super().__init__()
5         self.attn = nn.Linear(hidden_size * 2, hidden_size)
6         self.v = nn.Linear(hidden_size, 1, bias=False)
7
8     def forward(self, hidden, encoder_outputs):
9         """Compute attention over encoder outputs"""
10        batch_size = encoder_outputs.shape[1]
11        src_len = encoder_outputs.shape[0]
12
13        # Repeat decoder hidden state
14        hidden = hidden.unsqueeze(1).repeat(1, src_len, 1)
15        encoder_outputs = encoder_outputs.permute(1, 0, 2)
16
17        # Calculate attention scores
18        energy = torch.tanh(self.attn(
19            torch.cat((hidden, encoder_outputs), dim=2)
20        ))
21        attention = self.v(energy).squeeze(2)
22
23        # Convert to probabilities
24        return F.softmax(attention, dim=1)
25
26 class AttentionDecoder(nn.Module):
27     def __init__(self, output_size, hidden_size, attention):
28         """Decoder with attention mechanism"""
29         super().__init__()
30         self.attention = attention
31         self.lstm = nn.LSTM(hidden_size * 2, hidden_size)
32         self.out = nn.Linear(hidden_size, output_size)
```

## How Attention Works:

- Score each encoder state
- Softmax creates distribution
- Weighted sum of encoder states
- Concatenate with decoder state

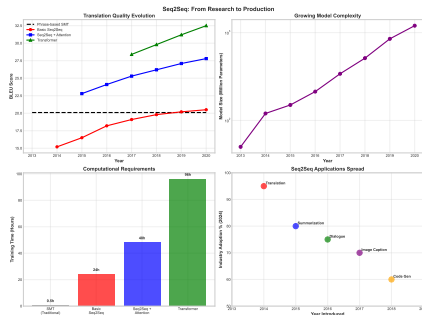
## Performance Boost:

- BLEU: 34.8  $\rightarrow$  41.8<sup>10</sup>
- Handles 2x longer sequences
- Interpretable alignments

---

15% improvement with attention

# Seq2Seq Impact: Translation Quality Leap



## Key Insights

- Phrase-based SMT dominated for decades
- Seq2seq matched SMT in 2014
- Attention surpassed all previous methods
- Google deployed in production 2016
- Now foundation for all translation systems



# Beyond Translation: Seq2Seq Everywhere (2024)

## Text Generation:

- Email auto-complete<sup>11</sup>
- Code documentation
- Story continuation
- Dialogue systems

## Summarization:

- News → headlines
- Papers → abstracts
- Meetings → notes
- Videos → descriptions

## Data Transformation:

- Text → SQL queries<sup>12</sup>
- Natural language → code
- Speech → text
- Image → caption

## Still Competitive:

- Smaller than transformers
- Streaming applications
- Low-latency requirements
- Resource-constrained devices

Seq2seq principle: Any input → any output transformation

---

<sup>1</sup>Gmail Smart Compose uses seq2seq variants

<sup>2</sup>Text-to-SQL still uses specialized seq2seq models

## Week 4 Exercise: Build a Smart Chatbot

**Your Mission:** Create a chatbot that gives variable-length responses

**Example Behavior:**

- "Hi" → "Hello! How can I help you today?"
- "What's the weather?" → "I'd need to check current conditions. What's your location?"
- "Tell me a joke" → (Generates different length jokes)

**Implementation Steps:**

- 1 Build encoder-decoder architecture
- 2 Train on conversation pairs
- 3 Implement attention mechanism
- 4 Compare responses with/without attention
- 5 Visualize what words the model attends to

**Bonus Challenges:**

- Implement beam search (top-k responses)
- Add personality tokens (formal/casual)
- Handle multi-turn conversations
- Measure response diversity

**You'll discover:** Why chatbots sometimes give generic responses!

## Key Takeaways: Breaking Free from Fixed Length

### What we learned:

- Language tasks need variable-length input/output
- Encoder-decoder separates understanding from generation
- Context vector creates information bottleneck
- Attention lets decoder access all input information
- Seq2seq enables any-to-any sequence transformation

### The evolution:

Fixed I/O (RNN) → Variable I/O (Seq2seq) → Full visibility (Attention)

**Critical Innovation:** Attention shows us what the model is "thinking" - interpretability!

### Next week: The Transformer

What if we used ONLY attention, no RNNs at all?

## References and Further Reading

### Foundational Papers:

- Sutskever et al. (2014). "Sequence to Sequence Learning with Neural Networks", NeurIPS
- Bahdanau et al. (2015). "Neural Machine Translation by Jointly Learning to Align and Translate"
- Cho et al. (2014). "Learning Phrase Representations using RNN Encoder-Decoder"

### Applications:

- Wu et al. (2016). "Google's Neural Machine Translation System"
- See et al. (2017). "Get To The Point: Summarization with Pointer-Generator Networks"
- Vinyals & Le (2015). "A Neural Conversational Model"

### Recommended Resources:

- TensorFlow Seq2seq Tutorial (with attention visualization)
- Visualizing and Understanding Neural Machine Translation
- PyTorch Chatbot Tutorial