# Recurrent Neural Networks

## Week 3 - Teaching Networks to Remember

NLP Course 2025

September 22, 2025

From Feedforward to Recurrent: Adding Memory to Neural Networks

**Week 3: The Memory Revolution**

## Processing Sequences with State

**The Challenge**
- Sequential data everywhere
- Order matters
- Variable length inputs
- Long-term dependencies

**The Solution**
- **Recurrent connections**
- Hidden state memory
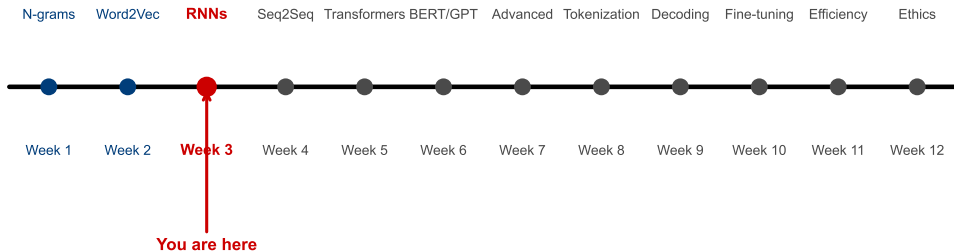- Parameter sharing
- Backprop through time

**The Evolution**
- Vanilla RNN → LSTM
- Solving gradients
- Gating mechanisms
- Modern variants (GRU)

The foundation of sequence modeling before transformers

# Course Journey: Where We Are

## NLP Course Journey

N-grams   Word2Vec   **RNNs**   Seq2Seq   Transformers   BERT/GPT   Advanced   Tokenization   Decoding   Fine-tuning   Efficiency   Ethics

Week 1   Week 2   **Week 3**   Week 4   Week 5   Week 6   Week 7   Week 8   Week 9   Week 10   Week 11   Week 12

**You are here**

**Journey So Far:**
- Week 1: Statistical language models (n-grams)
- Week 2: Word embeddings (Word2Vec, dense vectors)
- **Week 3: Sequential processing with memory**

**Coming Next:**
- Week 4: Encoder-decoder architectures
- Week 5: Attention mechanisms
- Week 6+: Transformers and beyond

# The Importance of Order

**Word Order Changes Meaning**

- "Dog bites man" $\neq$ "Man bites dog"
- "Not bad" $\neq$ "Bad, not!"
- Context flows through sequence

**Feedforward Limitations**

- Fixed input size
- No memory between inputs
- Can't model sequences naturally
- Position information lost

Sequential processing is fundamental to understanding language

**Sequential Tasks in NLP**

| Task | Type |
|------|------|
| Language Model | Many-to-many |
| Translation | Seq-to-seq |
| Sentiment | Many-to-one |
| Named Entity | Many-to-many |
| Speech Rec. | Seq-to-seq |

**Most NLP is inherently sequential**

# The Core Idea: Recurrence

**Mathematical Definition**

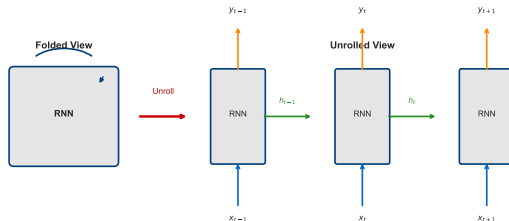$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

$$y_t = W_{hy}h_t + b_y$$

Where:

- $h_t$ = hidden state at time $t$
- $x_t$ = input at time $t$
- $y_t$ = output at time $t$
- $W_*$ = weight matrices (shared!)

**Key Insight**: Same weights at every timestep

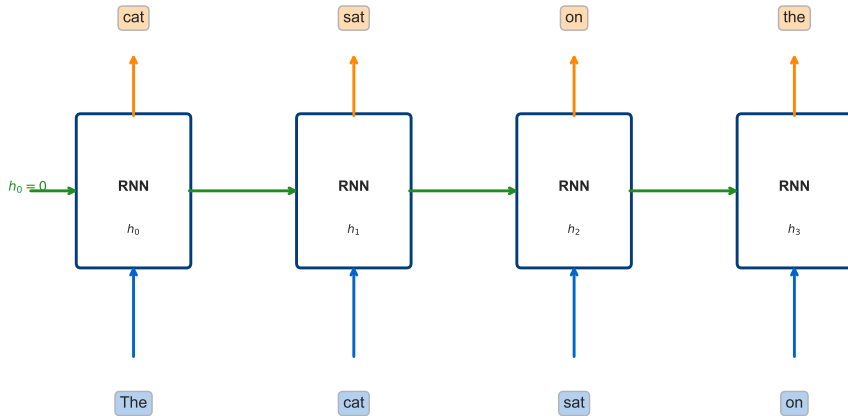One network, applied repeatedly through time



**Unrolled View Shows:**

- Information flows left-to-right
- Hidden state carries memory
- Parameters shared across time
- Can handle any sequence length

# Forward Pass: Step by Step

**RNN Forward Pass: Processing "The cat sat on"**

## Implementation: Simple RNN Cell

```python
import numpy as np

class RNNCell:
    def __init__(self, input_size, hidden_size):
        # Initialize weights
        self.Wxh = np.random.randn(input_size,
            hidden_size) * 0.01
        self.Whh = np.random.randn(hidden_size,
            hidden_size) * 0.01
        self.Why = np.random.randn(hidden_size,
            output_size) * 0.01
        self.bh = np.zeros((1, hidden_size))
        self.by = np.zeros((1, output_size))

    def step(self, x, h_prev):
        # Single timestep forward
        h = np.tanh(np.dot(x, self.Wxh) +
                    np.dot(h_prev, self.Whh) + self.bh)
        y = np.dot(h, self.Why) + self.by
        return y, h

    def forward(self, inputs):
        h = np.zeros((1, self.hidden_size))
        outputs = []

        for x in inputs:
```

**PyTorch Equivalent:**

```python
import torch.nn as nn

# Built-in RNN
rnn = nn.RNN(
    input_size=100,
    hidden_size=256,
    num_layers=1,
    batch_first=True
)

# Or use LSTM/GRU
lstm = nn.LSTM(
    input_size=100,
    hidden_size=256,
    num_layers=2,
    dropout=0.2,
    bidirectional=True
)

# Forward pass
output, (hn, cn) = lstm(input_seq)
```

**Modern frameworks handle the complexity**

# Why Simple RNNs Fail
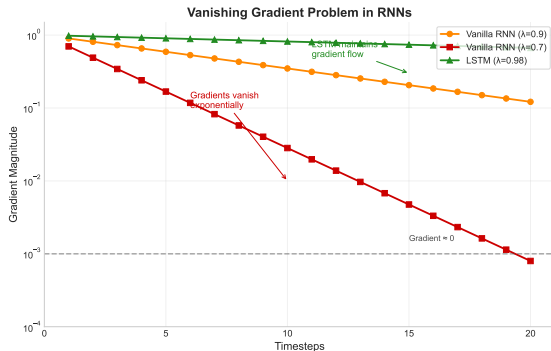
**The Problem**

Gradient through time:

$$\frac{\partial L}{\partial h_0} = \frac{\partial L}{\partial h_T} \prod_{t=1}^{T} \frac{\partial h_t}{\partial h_{t-1}}$$

Each term: $\frac{\partial h_t}{\partial h_{t-1}} = W_h^T \cdot \text{diag}(f'(h_{t-1}))$

For tanh: $|f'(x)| \leq 1$

If $||W_h|| < 1$: gradients $\rightarrow 0$ (vanish)
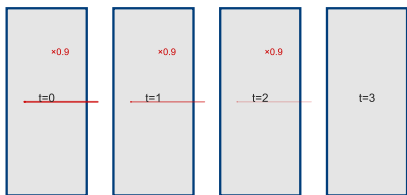
If $||W_h|| > 1$: gradients $\rightarrow$ (explode)



Vanishing Gradient Problem in RNNs

**Consequences:**

- Can't learn long dependencies
- Gradient 0 after 10-20 steps
- Network "forgets" early inputs

# Visualizing Gradient Flow

**Vanilla RNN Gradient Flow**



×0.9    ×0.9    ×0.9

t=0    t=1    t=2    t=3

*Gradient vanishes through multiplications*

**LSTM Gradient Flow**

+    +    +

t=0    t=1    t=2    t=3

*Gradient flows through cell state highway*

**Vanilla RNN**
- Exponential decay/growth
- Gradient magnitude: $O(\lambda^T)$
- Effective memory: 5-10 steps
- Cannot learn long patterns

**LSTM (Next Section)**
- Constant error flow
- Gradient highways
- Effective memory: 100+ steps
- Learns long dependencies

**Part 3: Long Short-Term Memory (LSTM)**

# Engineering Memory

**The Innovation (1997)**
Hochreiter Schmidhuber's insight:

- Add a **memory cell** $C_t$
- Control flow with **gates**
- Create gradient highways
- Selective reading/writing

**Three Gates:**

1. **Forget**: What to discard
2. **Input**: What to store
3. **Output**: What to expose

The architecture that made deep sequence modeling possible

**LSTM Equations**

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$
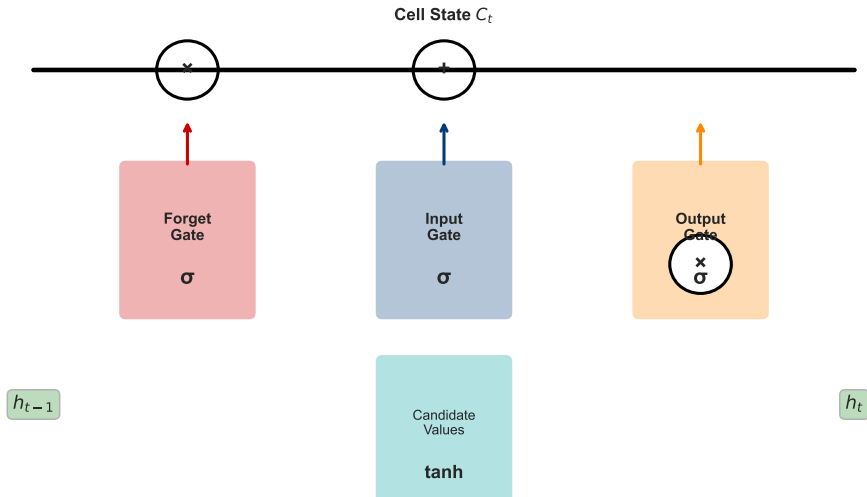$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$
$$h_t = o_t * \tanh(C_t)$$

---

**Gates use sigmoid (0-1) for control**

---

# LSTM Architecture: Gate Mechanisms



LSTM Architecture: Information Flow Through Gates

# RNN vs LSTM: Key Differences

| Aspect | Vanilla RNN | LSTM |
|---|---|---|
| Parameters | $O(h^2)$ | $O(4h^2)$ |
| Memory | Short (5-10 steps) | Long (100+ steps) |
| Gradient flow | Multiplicative | Additive |
| Training speed | Fast | Slower (4x params) |
| Gradient problem | Severe | Largely solved |
| Use cases | Short sequences | Most applications |

**When to Use RNN:**

- Very short sequences
- Real-time constraints
- Limited compute
- Simple patterns

**When to Use LSTM:**

- Long dependencies
- Complex patterns
- Production systems
- Default choice (pre-2017)

**LSTM's complexity is justified by superior performance**

# GRU: Gated Recurrent Unit

**Simplification of LSTM (2014)**

Cho et al. merged gates:

- Only **2 gates** instead of 3
- No separate cell state
- Fewer parameters (3x vs 4x)
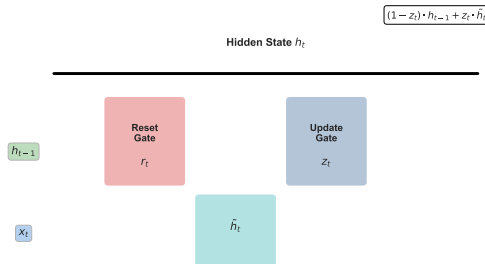- Similar performance

**GRU Equations:**

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$
$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

GRU Architecture: Simplified Gating

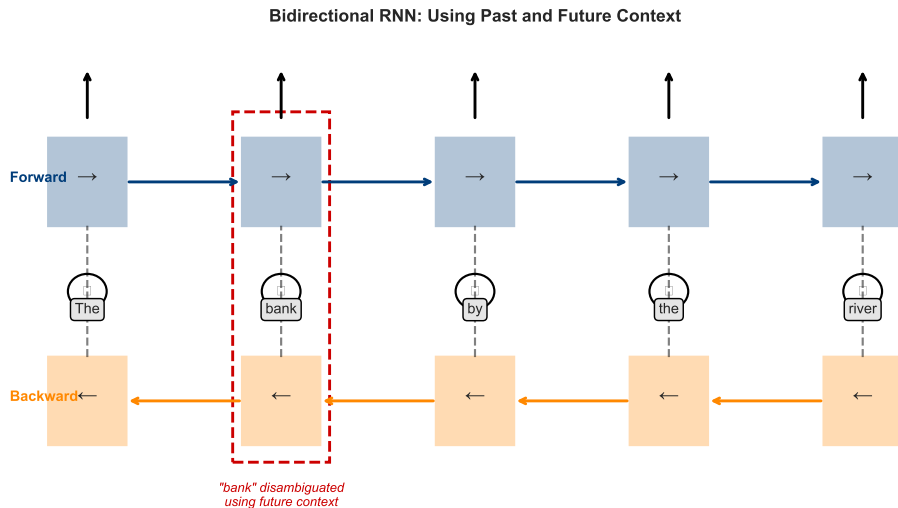$$(1 - z_t) \cdot h_{t-1} + z_t \cdot \tilde{h}_t$$

Hidden State $h_t$

$h_{t-1}$

$x_t$

Reset Gate $r_t$

Update Gate $z_t$

$\tilde{h}_t$

**Gates:**

- **Update gate** ($z_t$): How much to update
- **Reset gate** ($r_t$): How much past matters

# Bidirectional RNNs: Using Future Context



Bidirectional RNN: Using Past and Future Context

Forward →

Backward ←

The    bank    by    the    river

"bank" disambiguated
using future context

# Training RNNs: Practical Tips

## Common Issues & Solutions

1. **Gradient Explosion**
   - Solution: Gradient clipping
   - 'torch.nn.utils.clip_grad_norm_'
   - Typical value: 1.0 - 5.0

2. **Initialization**
   - Xavier/He initialization
   - Forget gate bias = 1.0 (LSTM)
   - Helps gradient flow

3. **Overfitting**
   - Dropout (between layers)
   - Recurrent dropout (careful!)
   - Weight decay

## Hyperparameters

| Parameter | Typical Range |
|---|---|
| Hidden size | 128 - 512 |
| Num layers | 1 - 3 |
| Learning rate | 1e-3 - 1e-2 |
| Batch size | 32 - 128 |
| Sequence length | 20 - 200 |
| Gradient clip | 1.0 - 5.0 |
| Dropout | 0.2 - 0.5 |

**Training Strategy:**

- Start with small sequences
- Gradually increase length
- Monitor gradient norms
- Use teacher forcing wisely

**RNN training requires careful tuning and monitoring**

# Real-World Applications
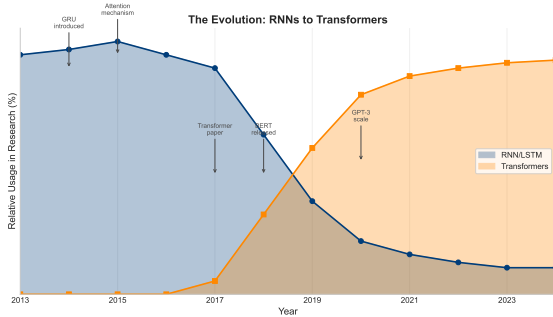
## Natural Language Processing

- Language modeling (pre-2018)
- Machine translation (pre-2017)
- Speech recognition (still used)
- Named entity recognition
- Sentiment analysis

## Time Series

- Stock price prediction
- Weather forecasting
- Anomaly detection
- Signal processing

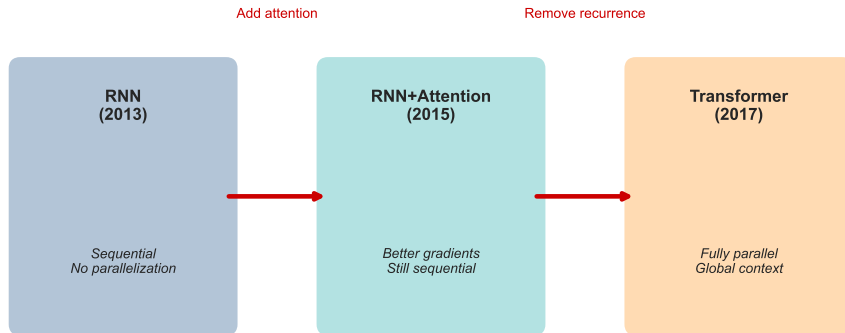## Modern Context (2024)



The Evolution: RNNs to Transformers

## Where RNNs Still Win:

- Streaming/online processing
- Edge devices (memory constraints)
- Variable-length sequences
- Time series with clear temporal patterns

RNNs remain relevant for specific use cases

# From RNNs to Transformers: Evolution

**Evolution of Sequence Models**

Add attention                    Remove recurrence

**RNN
(2013)**

*Sequential
No parallelization*

**RNN+Attention
(2015)**

*Better gradients
Still sequential*

**Transformer
(2017)**

*Fully parallel
Global context*

**Key Takeaways**

## What We Learned About RNNs

**Core Concepts**
- **Recurrence** for sequences
- Hidden state as memory
- Parameter sharing
- Backprop through time

**Challenges**
- Vanishing gradients
- Sequential bottleneck
- Training difficulty
- Limited context window

**Solutions**
- LSTM/GRU gates
- Gradient clipping
- Bidirectional processing
- Attention (next week!)

**RNNs introduced memory to neural networks - a crucial innovation**

Next Week: Sequence-to-Sequence Models
How to translate, summarize, and generate with encoder-decoder architectures

# References & Further Reading

**Foundational Papers:**

- Hochreiter & Schmidhuber (1997). "Long Short-Term Memory"
- Cho et al. (2014). "Learning Phrase Representations using RNN Encoder-Decoder" (GRU)
- Graves (2013). "Generating Sequences With RNNs"
- Karpathy (2015). "The Unreasonable Effectiveness of RNNs" (blog)

**Practical Resources:**

- PyTorch RNN Tutorial: `pytorch.org/tutorials/intermediate/char_rnn`
- Understanding LSTMs: `colah.github.io/posts/2015-08-Understanding-LSTMs/`
- Stanford CS224N Lecture 6: RNNs and Language Models

**Code Examples:**

- Week 3 Lab: 'week03_rnn_lab.ipynb'
- GitHub: Various char-RNN implementations
- Hugging Face: Modern RNN models