

From N-grams to LSTMs

N-grams · Word2Vec · RNNs · LSTMs · Motivation for Attention

The core problem

“The cat sat on the mat”



Representation

How do we turn words into numbers that capture meaning?

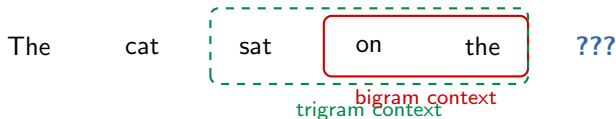
Sequence modeling

How do we capture context and word order?



N-gram language models

$$P(w_t \mid w_{t-n+1}, \dots, w_{t-1}) = \frac{\text{count}(w_{t-n+1} \cdots w_t)}{\text{count}(w_{t-n+1} \cdots w_{t-1})}$$



Bigram:

$$P(\text{mat} \mid \text{the}) = \frac{\text{count}(\text{"the mat"})}{\text{count}(\text{"the"})}$$

Trigram:

$$P(\text{mat} \mid \text{on, the}) = \frac{\text{count}(\text{"on the mat"})}{\text{count}(\text{"on the"})}$$

Idea: estimate probabilities by counting how often word sequences appear in a large corpus

N-gram limitations

12 words apart — trigram window can't see “cat”

“The **cat** that sat on the mat near the door by the window **was** _____”

1. Fixed context window

A trigram only sees 2 previous words. Long-range dependencies are invisible.

2. Curse of dimensionality

Vocabulary $V = 50k \Rightarrow$ possible 5-grams: $50k^5 = 3 \times 10^{23}$. Most never observed.

3. Data sparsity

“armadillo aardvark” has count 0 in most corpora. Smoothing helps but doesn't solve it.

4. No generalization

Knowing “dog sat on” doesn't help predict after “cat sat on.” Words are discrete symbols.

We need a way to represent words so that similar words share information.

One-hot encoding

Vocabulary

1. cat

2. dog

3. sat

4. the

⋮

50k. zoo



One-hot vectors

cat =


1	0	0	0	⋯	0
---	---	---	---	---	---

dog =

0	1	0	0	⋯	0
---	---	---	---	---	---

sat =

0	0	1	0	⋯	0
---	---	---	---	---	---


50,000 dimensions

No similarity

$$\text{cat} \cdot \text{dog} = 0$$

$$\text{cat} \cdot \text{banana} = 0$$

Every pair is equally different!

Huge & sparse

50k-dim vectors with exactly one non-zero entry

We need: dense, low-dimensional vectors
where similar words have similar representations

Word2Vec — the distributional hypothesis

“You shall know a word by the company it keeps”
— J. R. Firth, 1957

Words appearing in similar contexts have similar meanings:

“The movie was surprisingly **good** and the audience loved it”

“The movie was surprisingly **great** and the audience loved it”

“The movie was surprisingly **excellent** and the audience loved it”

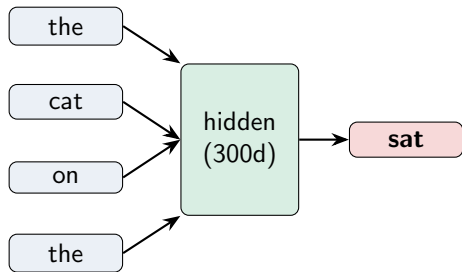


good, **great**, and **excellent** share context \Rightarrow
their vectors should be close in embedding space

Word2Vec — CBOW & Skip-gram

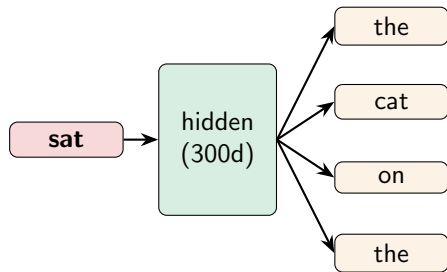
CBOW

Context → Center word



Skip-gram

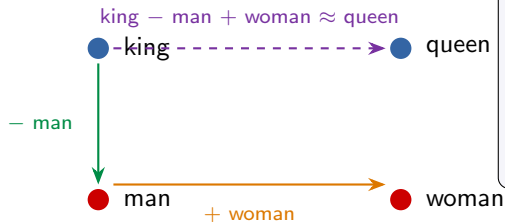
Center word → Context



Result: the hidden layer weights *are* the word embeddings — dense vectors (100–300 dims)

Word2Vec — emergent properties

Vector arithmetic in embedding space



More analogies:

Paris - France + Italy \approx Rome
bigger - big + small \approx smaller
walking - walk + swim \approx swimming

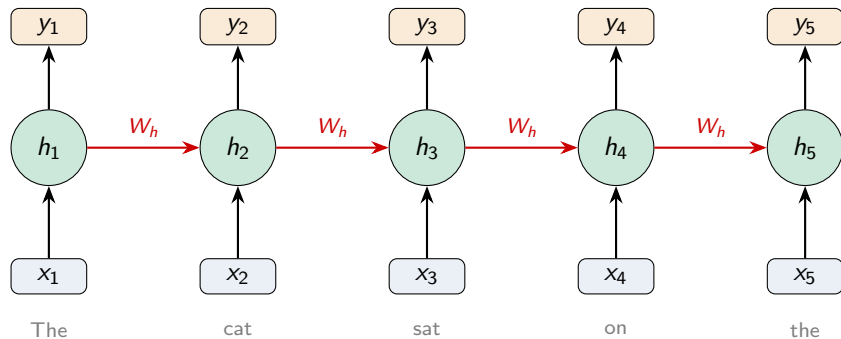
Critical limitation: one vector per word

“I went to the river **bank**” vs. “I deposited money at the **bank**”
Both map to the *same* vector — no polysemy, no context-dependence

We need representations that change based on context. Enter: recurrent neural networks.

Recurrent Neural Networks (RNNs)

Unrolled RNN across time steps

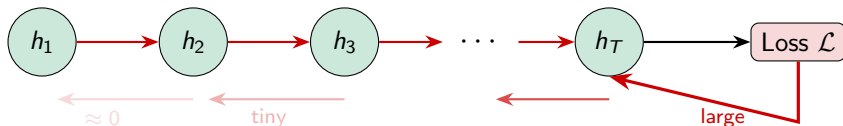


$$h_t = \tanh(W_h \cdot h_{t-1} + W_x \cdot x_t + b)$$

Key: hidden state h_t is a “memory” — no fixed window!

The vanishing gradient problem

Backpropagation through time



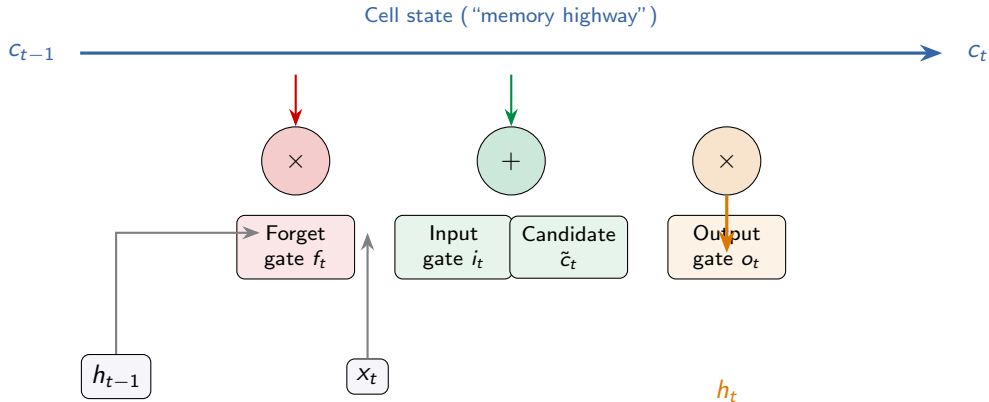
$$\frac{\partial \mathcal{L}}{\partial h_1} = \frac{\partial \mathcal{L}}{\partial h_T} \cdot \prod_{t=2}^T \frac{\partial h_t}{\partial h_{t-1}} \quad \text{Prod-}$$

uct of many values $< 1 \Rightarrow$ gradient vanishes

“**The cat** that I saw yesterday in the garden near the old oak tree **was** cute”
The gradient from “was” can’t reach “cat” — the RNN *forgets* early tokens

Solution attempt: LSTM (1997) and GRU (2014)

LSTM — Long Short-Term Memory

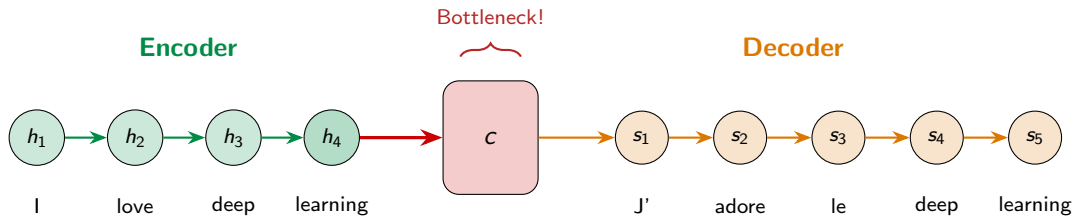


Cell state flows with minimal modification — gradient highway

Gates learn what to forget, store, and output at each step

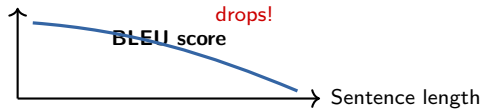
Sequence-to-sequence (Encoder-Decoder)

Machine translation with LSTMs (Sutskever et al., 2014)



The entire input sentence is compressed into a single fixed-size vector c !

Works for short sentences, but performance degrades for long ones.

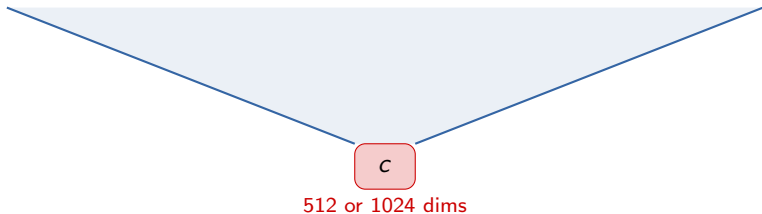


What if the decoder could **look back** at all encoder states, not just c ?

→ This is attention!

Why LSTMs are not enough — the bottleneck

“The quick brown fox jumped over the lazy dog that was sleeping by the fireplace in the old cottage”



Information lost:

Early tokens get overwritten by later ones as h_t is updated sequentially

Fixed capacity:

Whether the input is 5 tokens or 500, it must fit in the same vector

Fundamental issue: you can't losslessly compress arbitrary-length sequences into a fixed-size vector

Why LSTMs are not enough — sequential processing

LSTM: must process tokens one at a time



Each step **waits** for the previous one

What if we could process all tokens in parallel?



All processed **simultaneously** on GPU

LSTM: $O(T)$ sequential steps
Can't utilize GPU parallelism
Training is *slow*

Transformer: $O(1)$ parallel steps
Fully utilizes GPU cores
Training is *fast*

Why LSTMs are not enough — long-range dependencies

“The **trophy** didn't fit in the **suitcase** because **it** was too _____”

If **it** = **trophy**:
“it was too **big**”

If **it** = **suitcase**:
“it was too **small**”

Path length from “it” to its referent:

LSTM: $O(n)$ steps

Information must pass through
every intermediate hidden state
Signal degrades over distance

Attention: $O(1)$ steps

Direct connection be-
tween any two tokens
No degradation over distance

Attention: “let me directly look at any token I need” — regardless of distance

The stage is set

N-grams

Fixed window
No generalization
Data sparsity
Can't scale

Word2Vec

Good embeddings
But static (one
vector per word)
No sequence
modeling

RNN / LSTM

Sequential
(× parallel)
Bottleneck vector
Long-range deps
still hard



We need: **parallel processing, direct access
to all positions, contextual representations**

⇒ **Self-Attention** and the **Transformer** (Vaswani et al., 2017)

Questions?

Next: The Transformer Architecture