

# The Architecture of Prediction

Understanding N-Gram Language Models:  
The Mathematical Ancestors of Modern AI

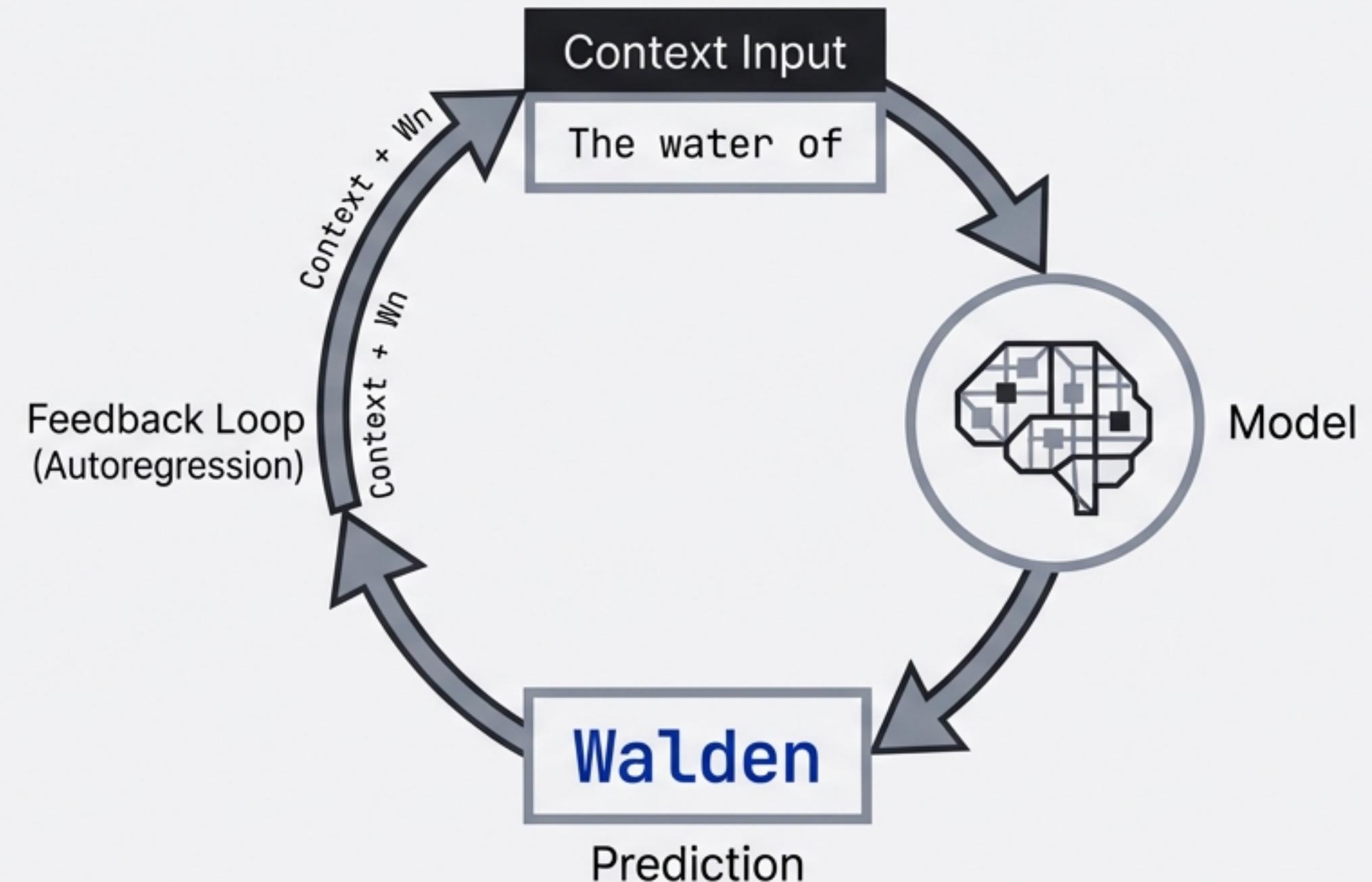
The water of Walden Pond

A scatter plot of lowercase letters and punctuation marks (., ?, !) in various sizes and gray tones, resembling a stylized word cloud or a visual representation of text frequency.

# At the Heart of Every LLM is a Simple Question: ‘What Comes Next?’

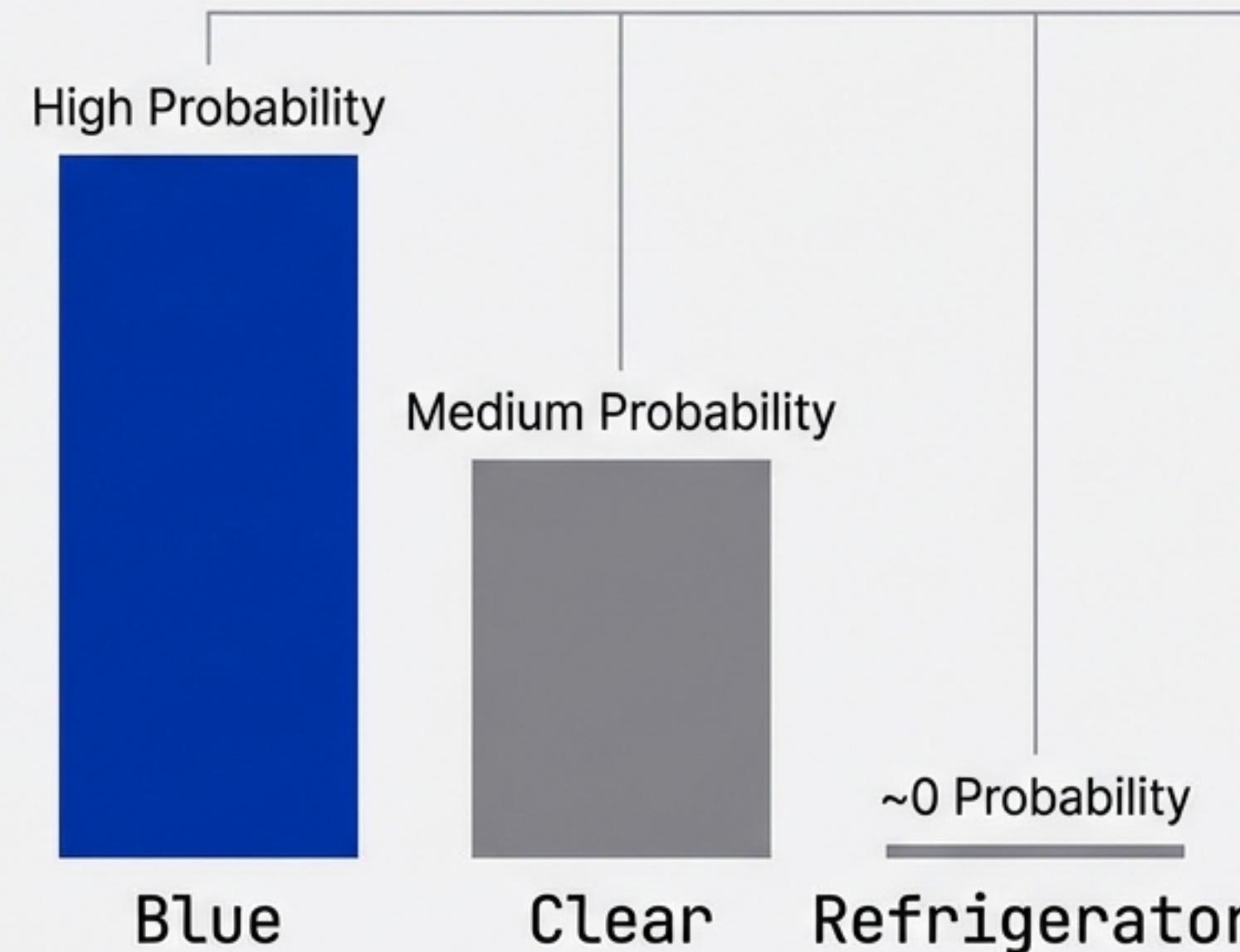
Modern Large Language Models (LLMs) are fundamentally autoregressive. Whether it is a simple N-gram model or a massive Transformer, the training objective remains consistent: predict the next word based on the context of previous words.

- **Pre-training:** Hiding a word to learn context.
- **Generation:** Predicting  $W_n$ , appending it, and repeating.



# Language is Probabilistic, Not Random

The water of Walden Pond is beautifully [\_\_\_\_\_]



A Language Model (LM) assigns a probability to every word in the vocabulary ( $V$ ). The most contextually appropriate word receives the highest mass.

# Probability Resolves Ambiguity

## Context-Sensitive Spelling

~~Their~~ is no plan. ✗ Low P(Sequence)

There is no plan. ✓ High P(Sequence)

## Speech Recognition (ASR)



Language models act as statistical referees for noisy input.

# Formalizing the Intuition

## 1. The Sequence

Let  $W$  be a sequence of words:  $W_1, W_2, \dots, W_n$

## 2. The Dual Objectives

1. Next-Word Prediction

$$P(W_4 | W_1, W_2, W_3) \xrightarrow{\text{Mathematically Connected}} P(W) = P(W_1, W_2, \dots, W_n)$$

2. Sequence Probability

# The Chain Rule of Probability

The joint probability is the product of conditional probabilities.

$$\begin{aligned} P(\text{The water of Walden}) &= \\ \rightarrow P(\text{The}) & \\ \rightarrow \times P(\text{water} \mid \text{The}) & \\ \rightarrow \times P(\text{of} \mid \text{The water}) & \\ \rightarrow \times P(\text{Walden} \mid \text{The water of}) & \end{aligned}$$

This is the “Gold Standard” calculation. **It requires no loss of information.**

# The Data Bottleneck: Sparsity

Observed in Corpus

The water of Walden Pond

20 Count

Target Sequence

The water of Walden Pond is so beautifully blue

0 Count



## Zero Probability Problem.

If the model hasn't seen this exact sequence,  
the **Chain Rule** fails to estimate it.

# The Solution: The Markov Assumption

Approximating the future by looking only at the immediate past

$$P(W_n \mid W_1, W_2, \dots, W_{n-1})$$



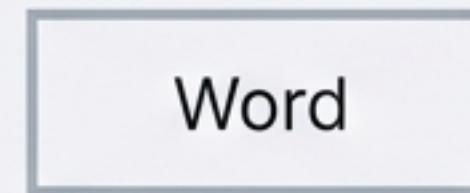
$$\approx P(W_n \mid W_{n-1})$$

We trade perfect accuracy for calculability. We assume the next word depends only on the previous few words.

# Defining the N-Gram Architecture

Unigram (N=1)

$$P(W_n)$$



Bigram (N=2)

$$P(W_n | W_{n-1})$$



Trigram (N=3)

$$P(W_n | W_{n-2}, W_{n-1})$$



General Rule: An N-gram looks at N-1 words of past context.

# Maximum Likelihood Estimation (MLE)

Converting corpus counts into probabilities.

$$\frac{\text{Count}(\text{question, paper})}{\text{Count}(\text{question})} = \frac{180}{600} = 0.3$$

P(paper | question) = 30%

Bigram Frequency

Unigram Frequency

# The Evolution of Meaning: Generating Shakespeare

## Unigram (Random)

To him swallowed  
confess here  
both which.

Chaos

## Bigram (Local Sense)

I confess all  
sorts... He is...  
Would seem...  
Last December  
through the way

Local Coherence

## Quadrigram (Near-Sentences)

It cannot be but  
so indeed the  
short and the  
long.

Structural Coherence

# Engineering Reality: The Log Space

## The Problem

Multiplying Probabilities leads to Underflow.

$$0.003 \times 0.0004 \times 0.00002 \approx 0$$

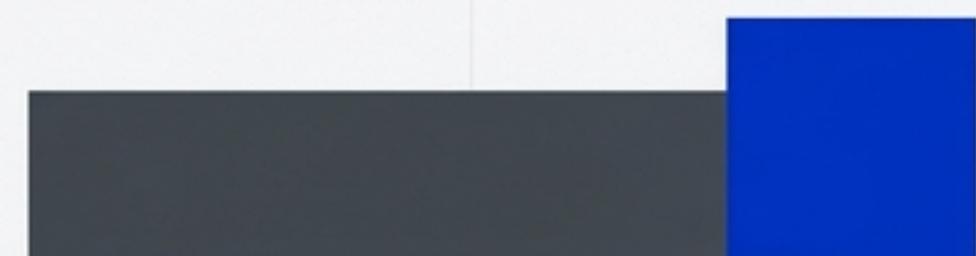


## The Solution

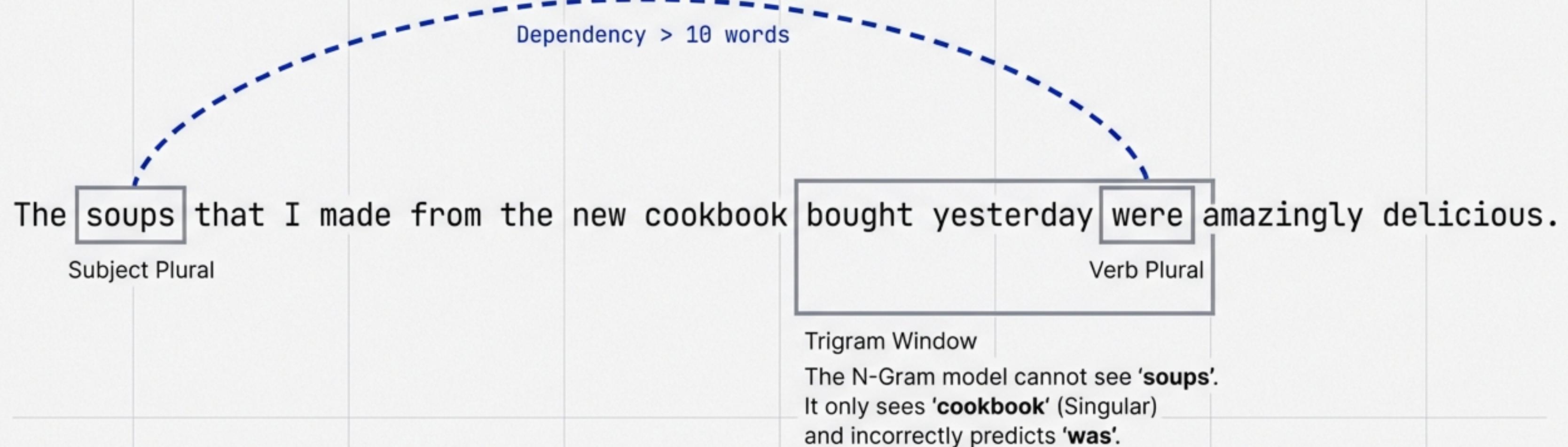
Logarithms convert multiplication to addition.

$$\log(P_1) + \log(P_2) + \log(P_3)$$

Addition is numerically stable.

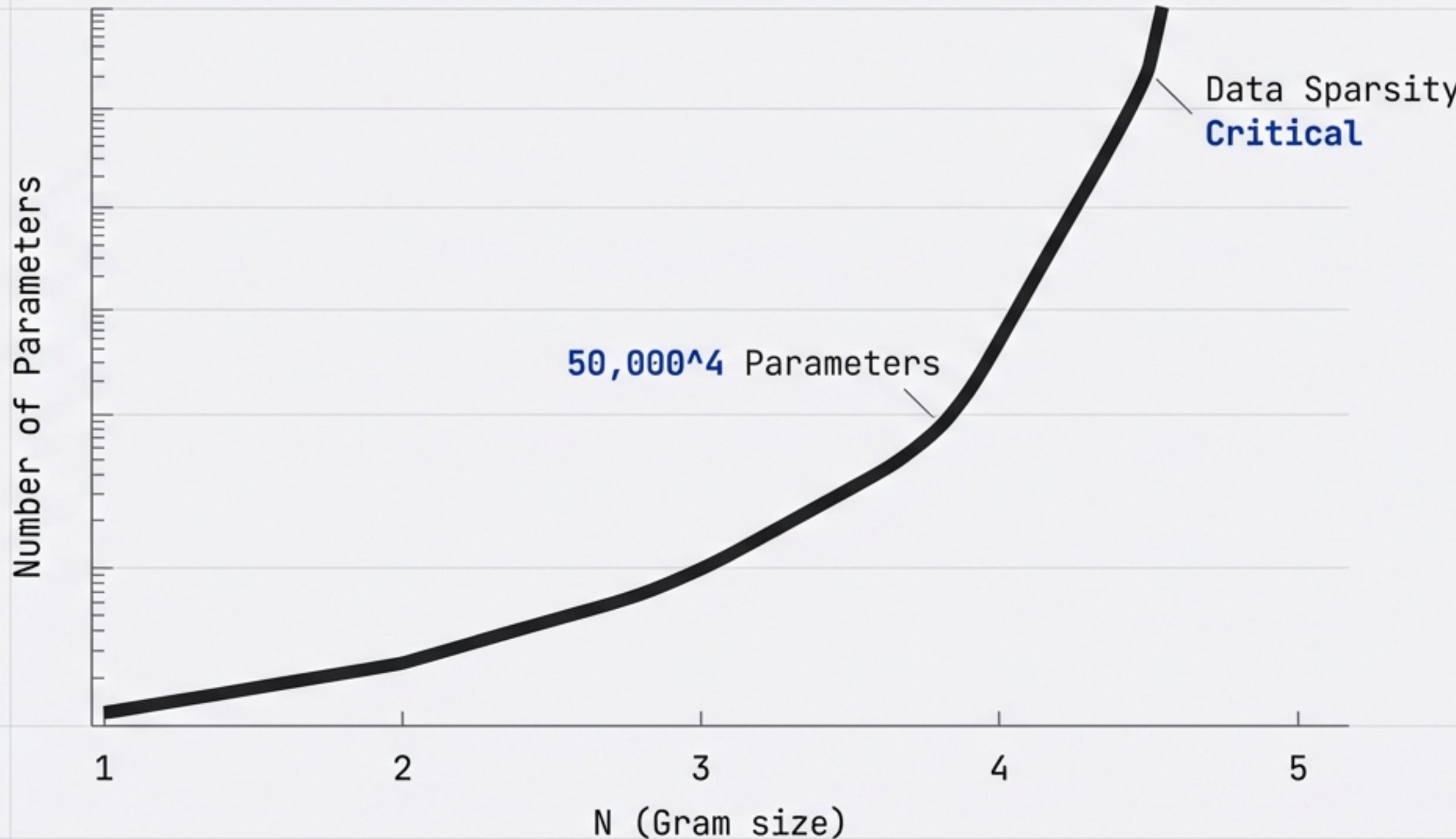


# Limitation: Long-Distance Dependencies



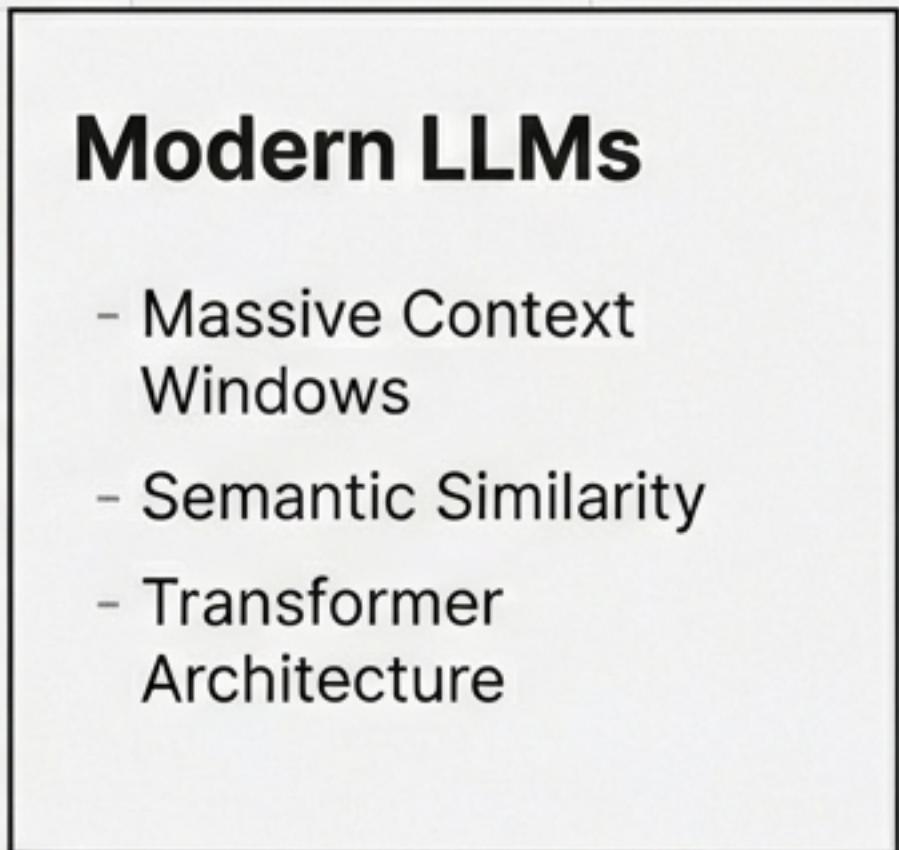
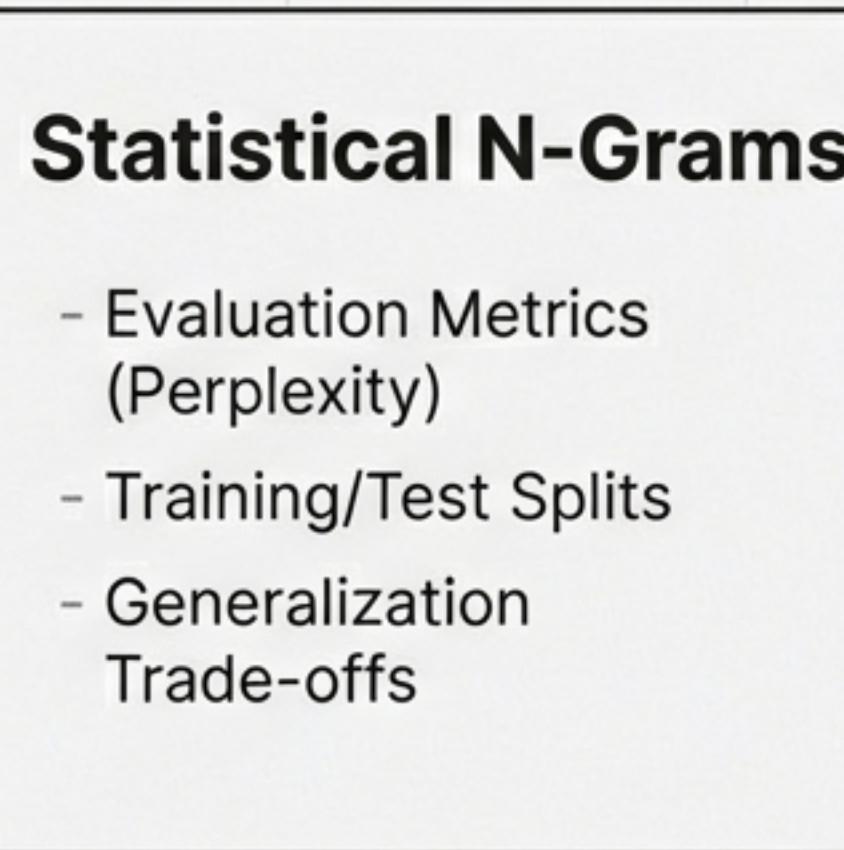
# The Trade-off: The Curse of Dimensionality

Why we can't just increase N indefinitely.



As N grows, the parameter space  $V^N$  explodes. We run out of data before we can learn the relationships.

# The Legacy of N-Grams



**N-grams taught us how to frame the problem.  
Deep Learning gave us the capacity to solve it.**