

# Stop Thinking, Start Remembering

CALCULATING PROBABILITY DISTRIBUTION OF DEFINITE ARTICLES BASED ON 17 TRILLION PARAMETERS...

**Standard LLM** 

**The**

Inconsolata

Left: The 'Standard LLM' approach: Brute force computation.

It says "The". Cool.

Memory Lookup

**Engram** 

**The**

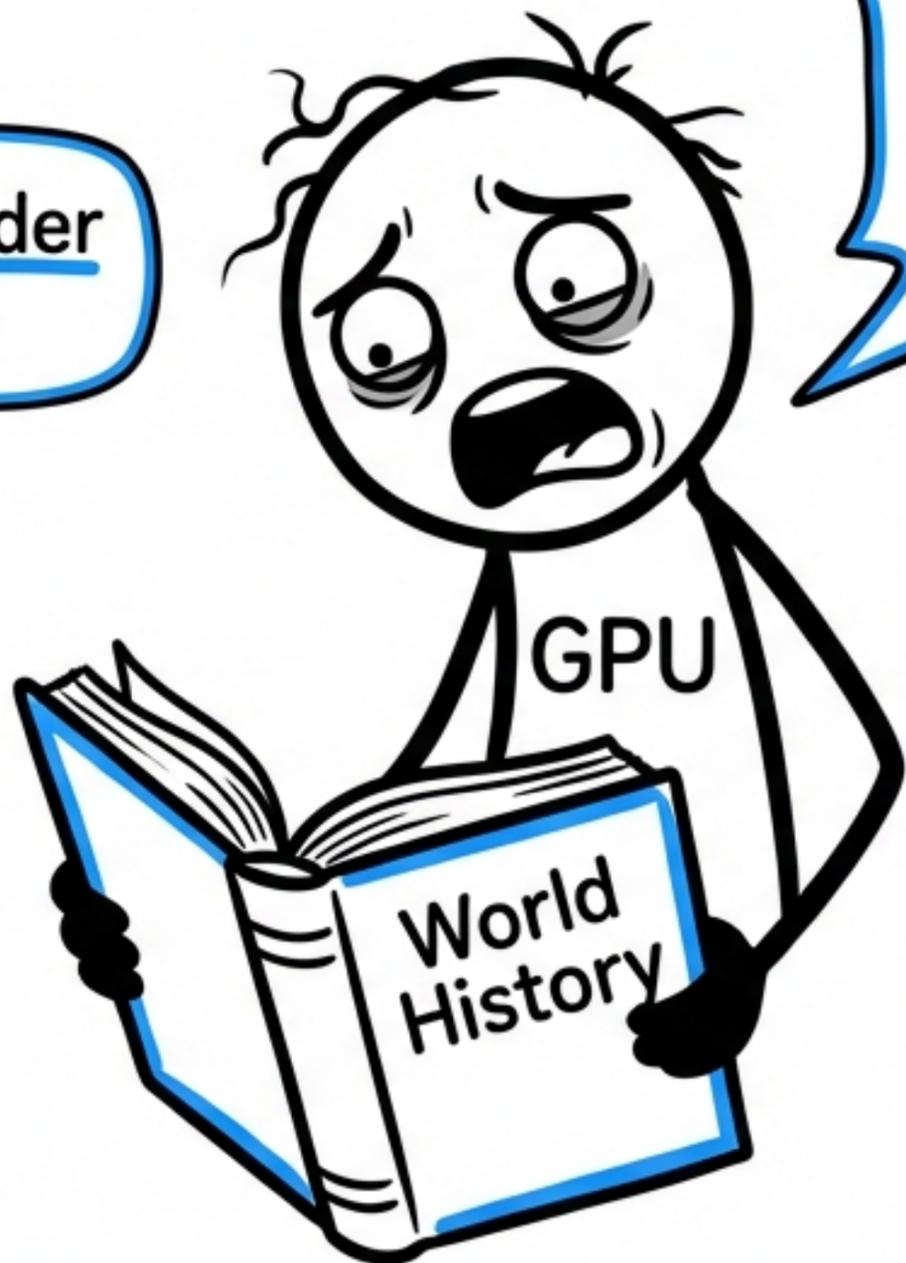
Inconsolata

Right: The 'Engram' approach: Instant retrieval.

## The Engram Architecture: A Stick Figure's Guide to Cheating at Intelligence

# The Status Quo: Simulation via Computation

Who is Alexander the Great?



WAIT! I must derive his existence from the fundamental laws of physics and the statistical probability of the token 'Great' following 'Alexander'!

## // THE PROBLEM:

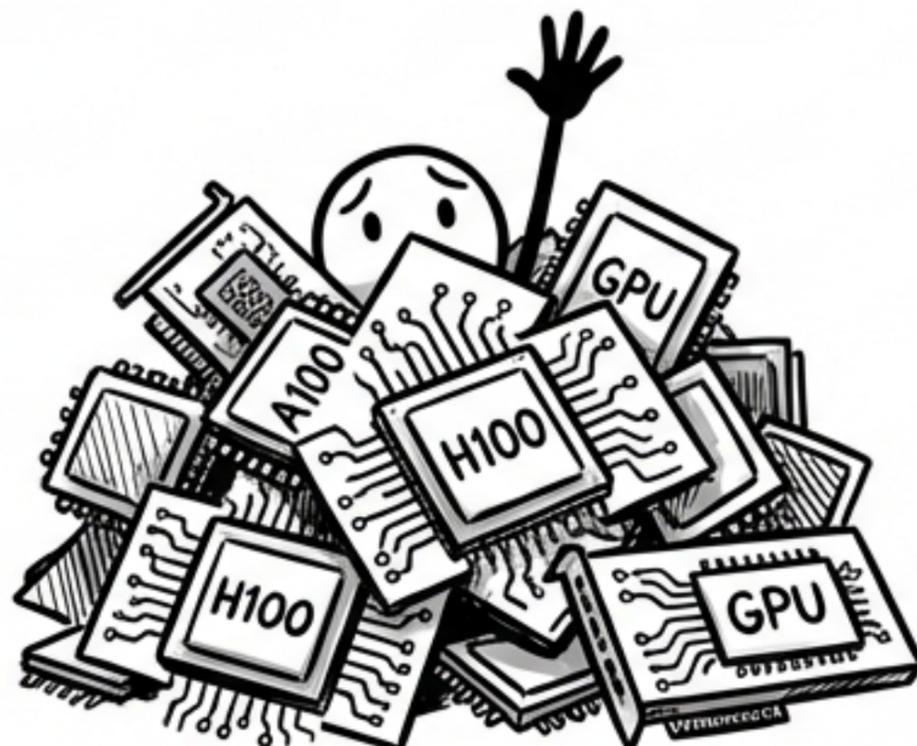
1. Transformers lack a native knowledge lookup primitive.
2. To resolve a static entity like 'Diana, Princess of Wales', the model consumes multiple layers of Attention and FFNs.
3. Result: Expensive runtime reconstruction of a static lookup table.

# The 'New' Idea: Stapling a 1990s Search Engine to a H100

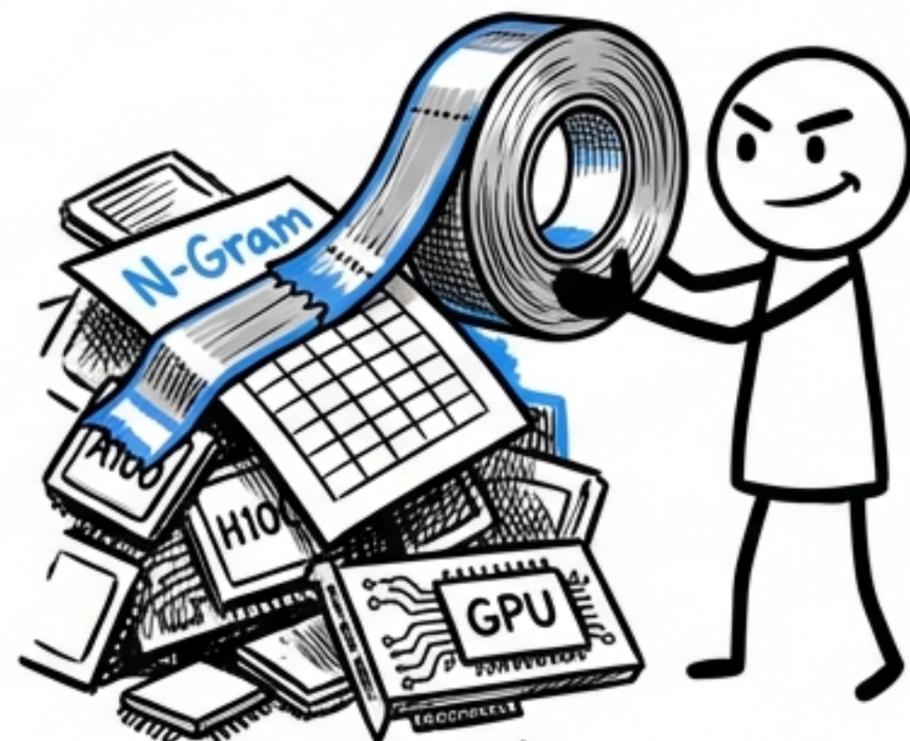
1950s



2017



2025 (Engram)



Inconsolata

$O(1)$  Lookup. Fast. Dumb.

Inconsolata

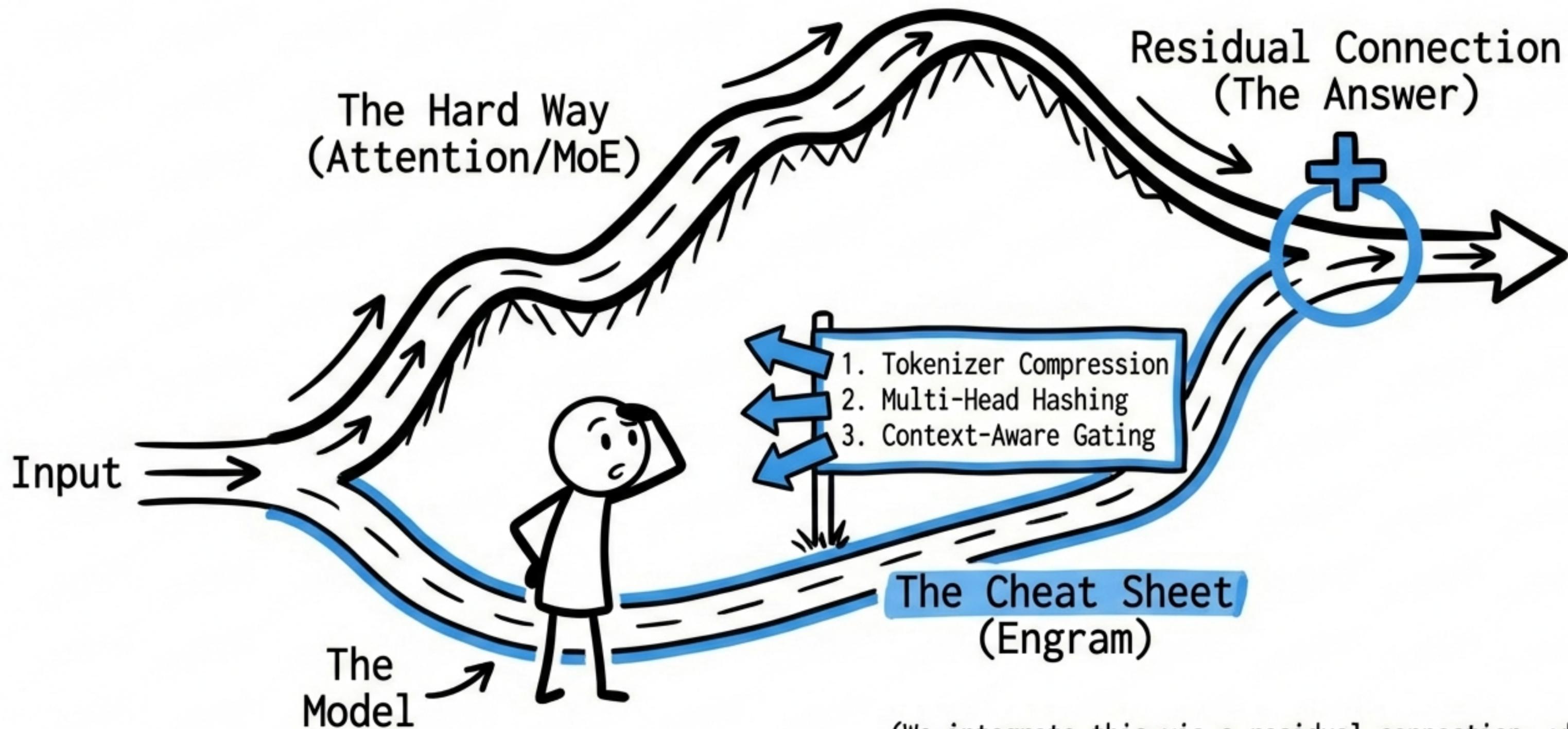
Transformers. Slow. Smart.

Inconsolata

The Best of Both Worlds.

The Insight: Use 'Conditional Memory' as a new axis of sparsity.  
Why compute what you can just find in a filing cabinet?

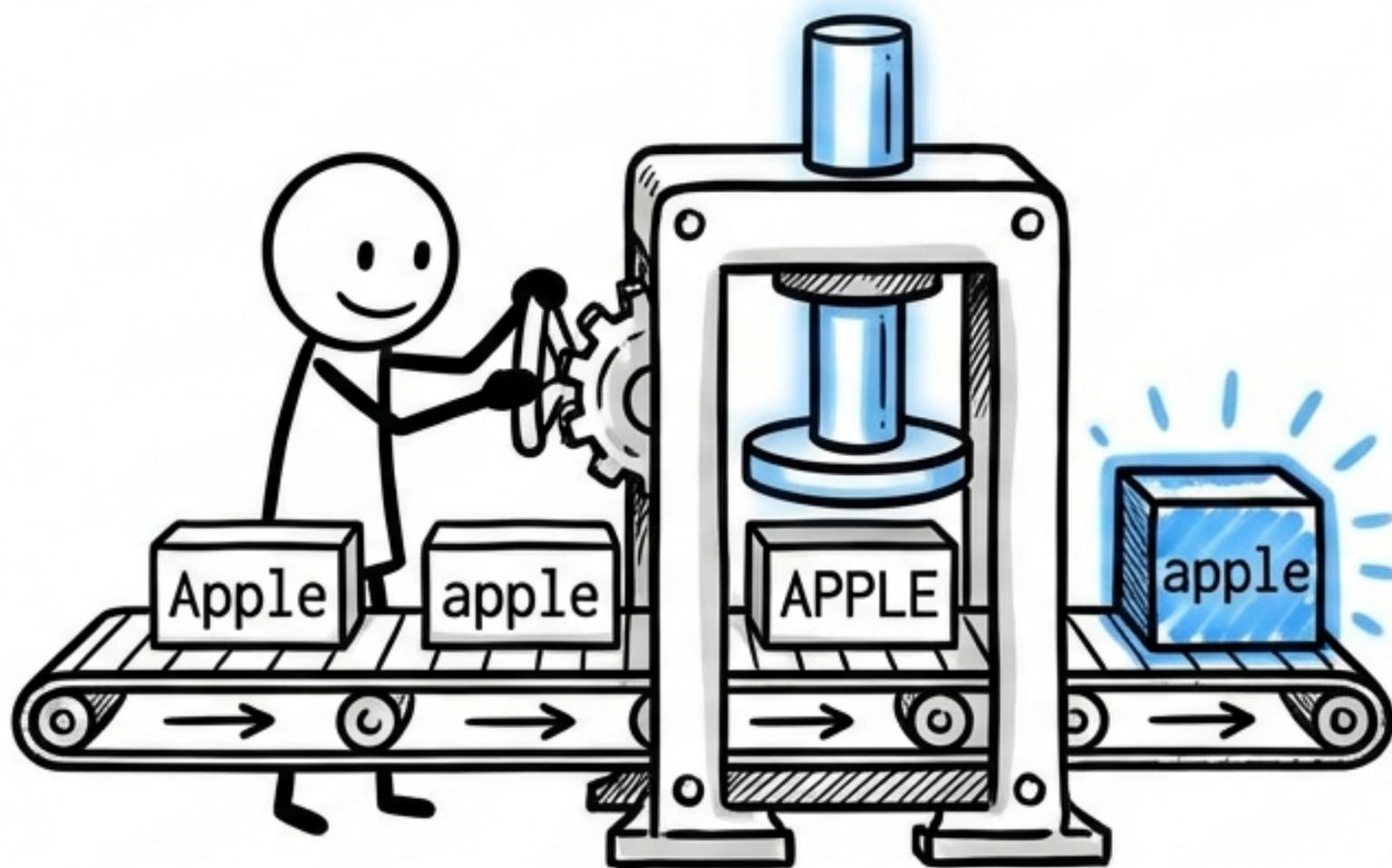
# The Architecture: The Model's 'Cheat Sheet'



(We integrate this via a residual connection, which is fancy talk for 'adding the answer at the end'.)

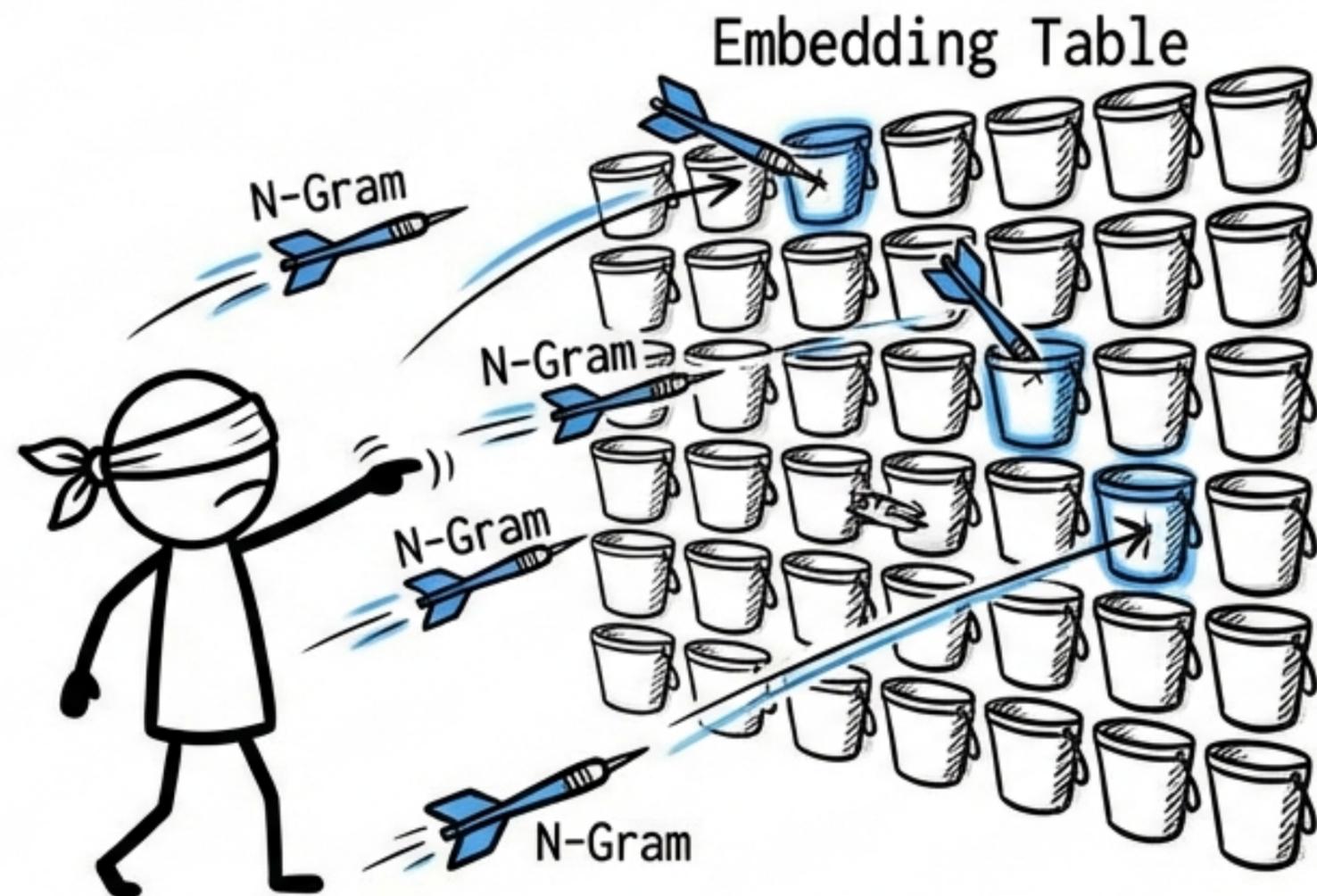
# Architects Daughter: Step 1: Squashing and Hashing

## A. Compression



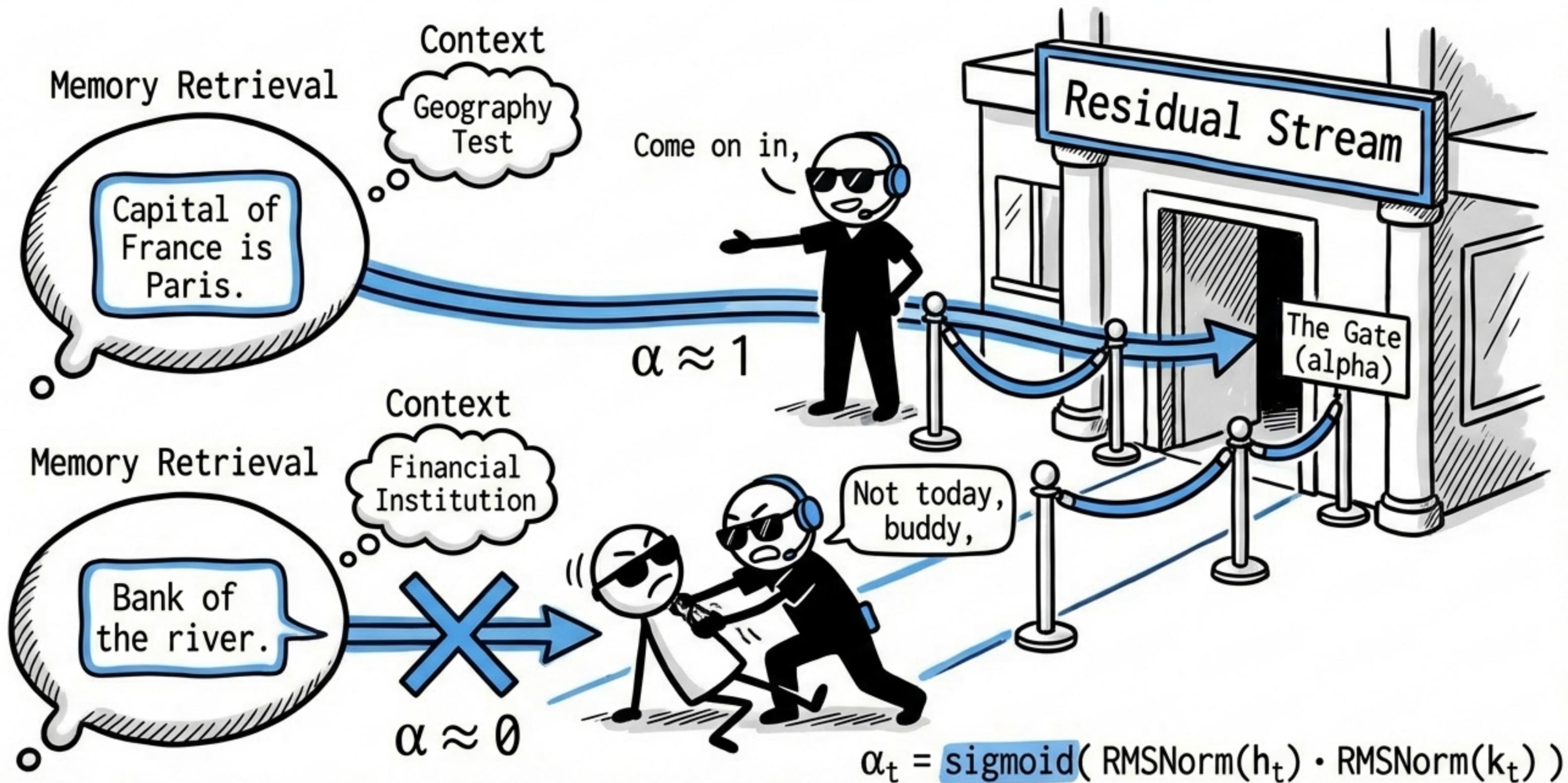
Tokenizer Compression: Surjective mapping reduces effective vocab size by ~23%.

## B. Multi-Head Hashing

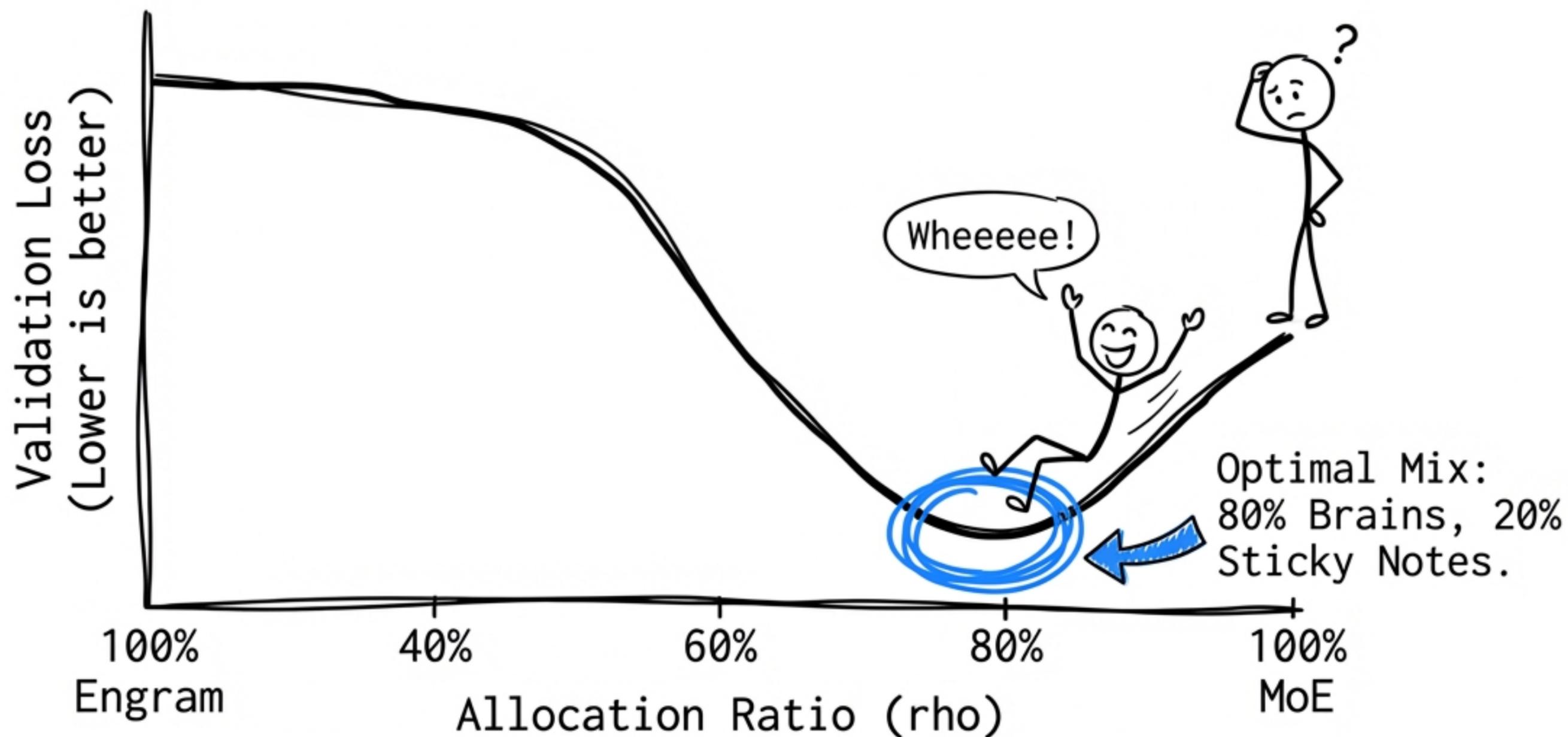


Multi-Head Hashing: We use  $K$  distinct hash heads. We rely on the magic of collisions not mattering as much as you think.

# Architects Daughter: Step 2: (Context-Aware Gating)

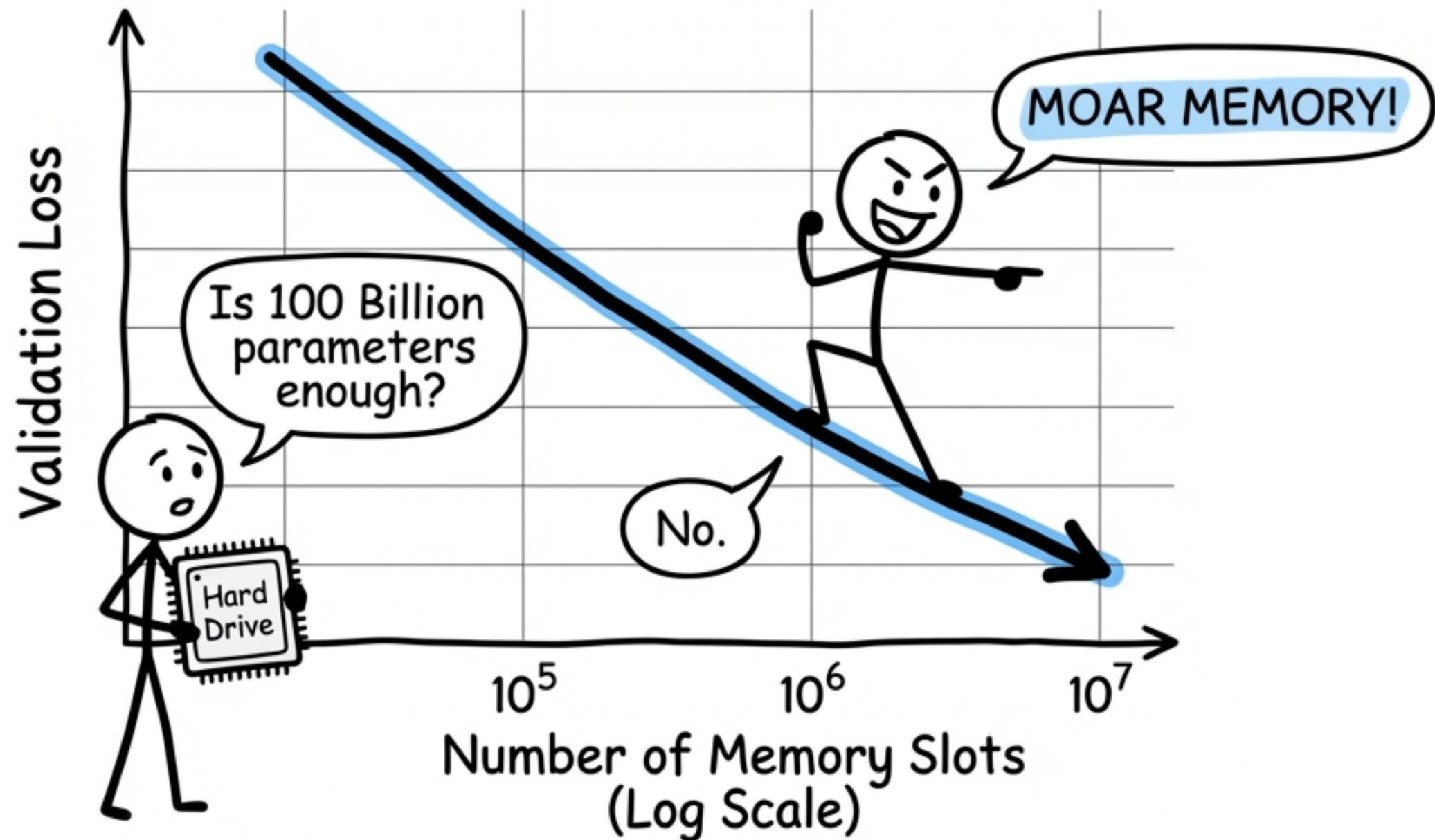


# The Sparsity Allocation Problem



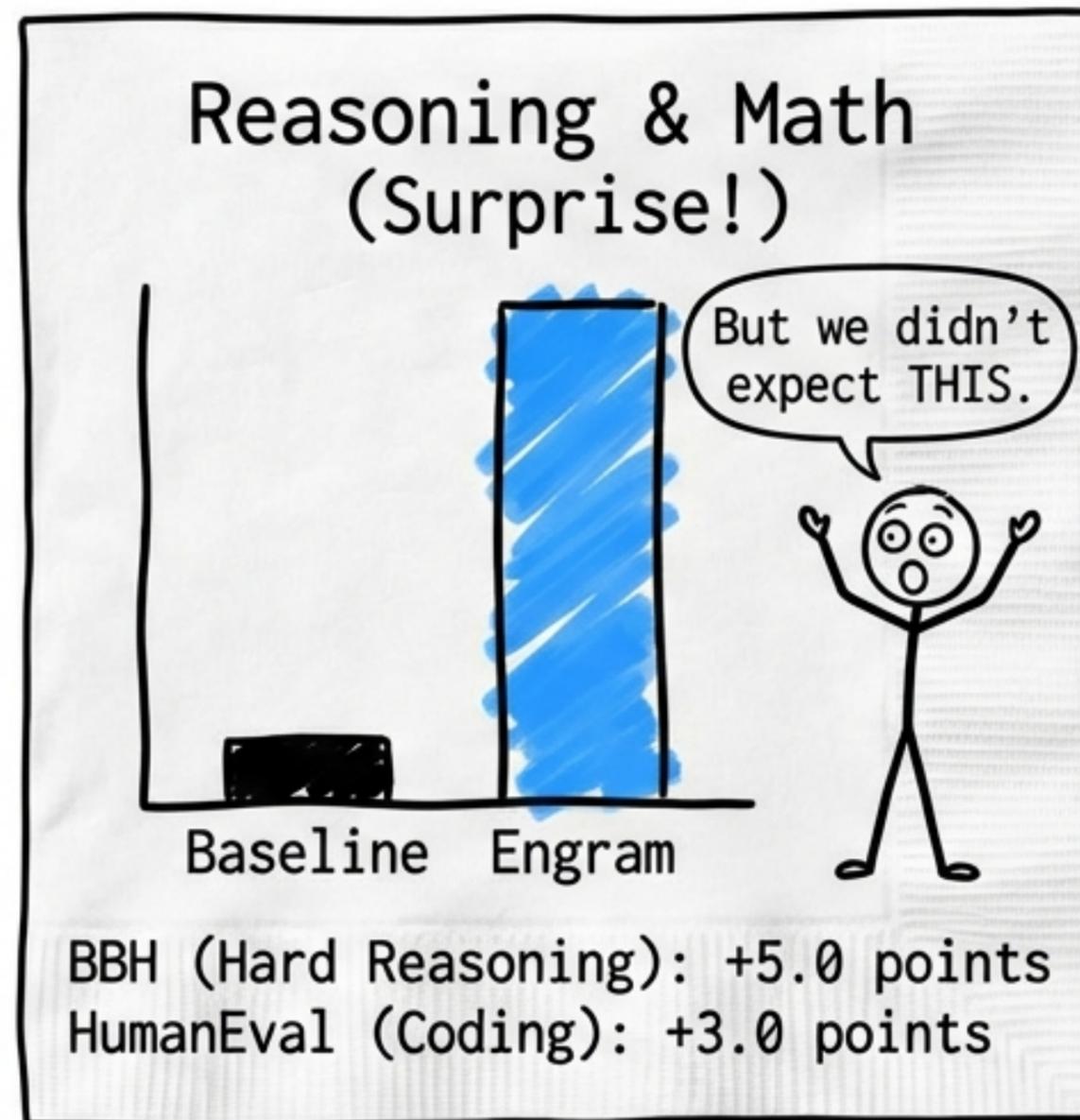
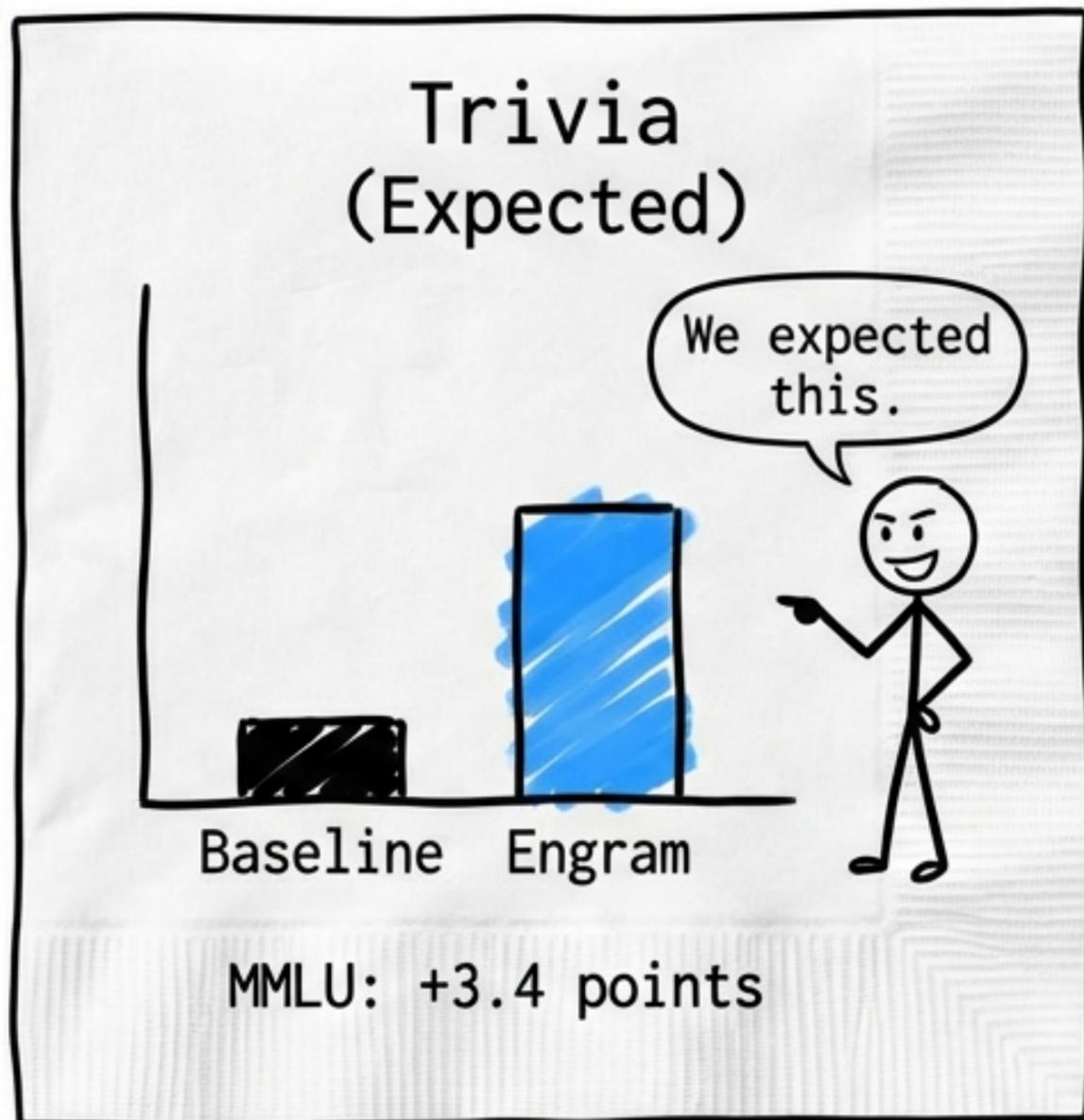
Pure MoE ( $\rho=100\%$ ) is suboptimal. A mix of computation and memory yields better performance.

# The “Infinite Memory” Regime



Fixed MoE Backbone (3B). Memory adds capacity without increasing compute cost.

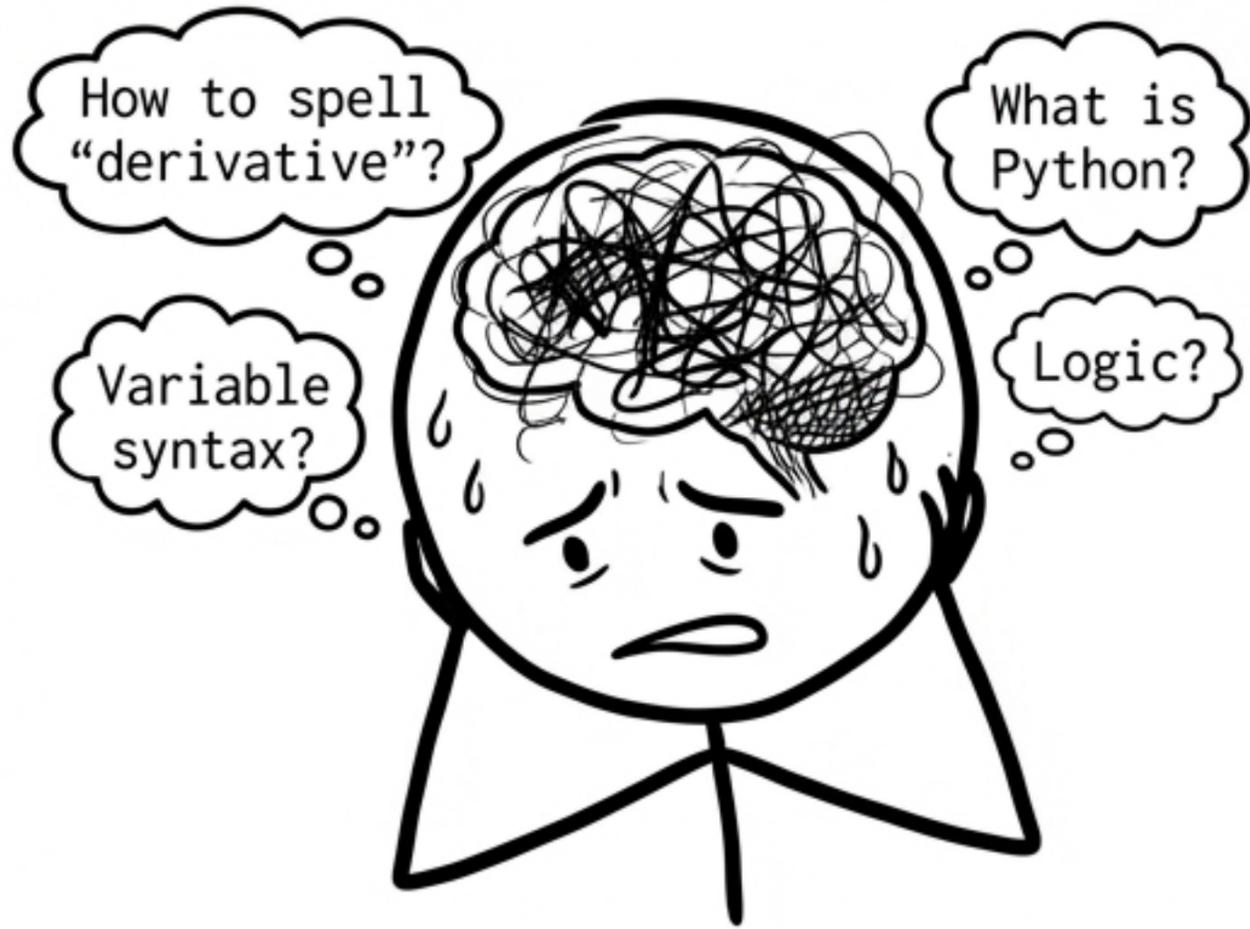
# The Plot Twist: It's Not Just for Trivia



We built a cheat sheet for history class,  
but it helped us pass Calculus. Weird.

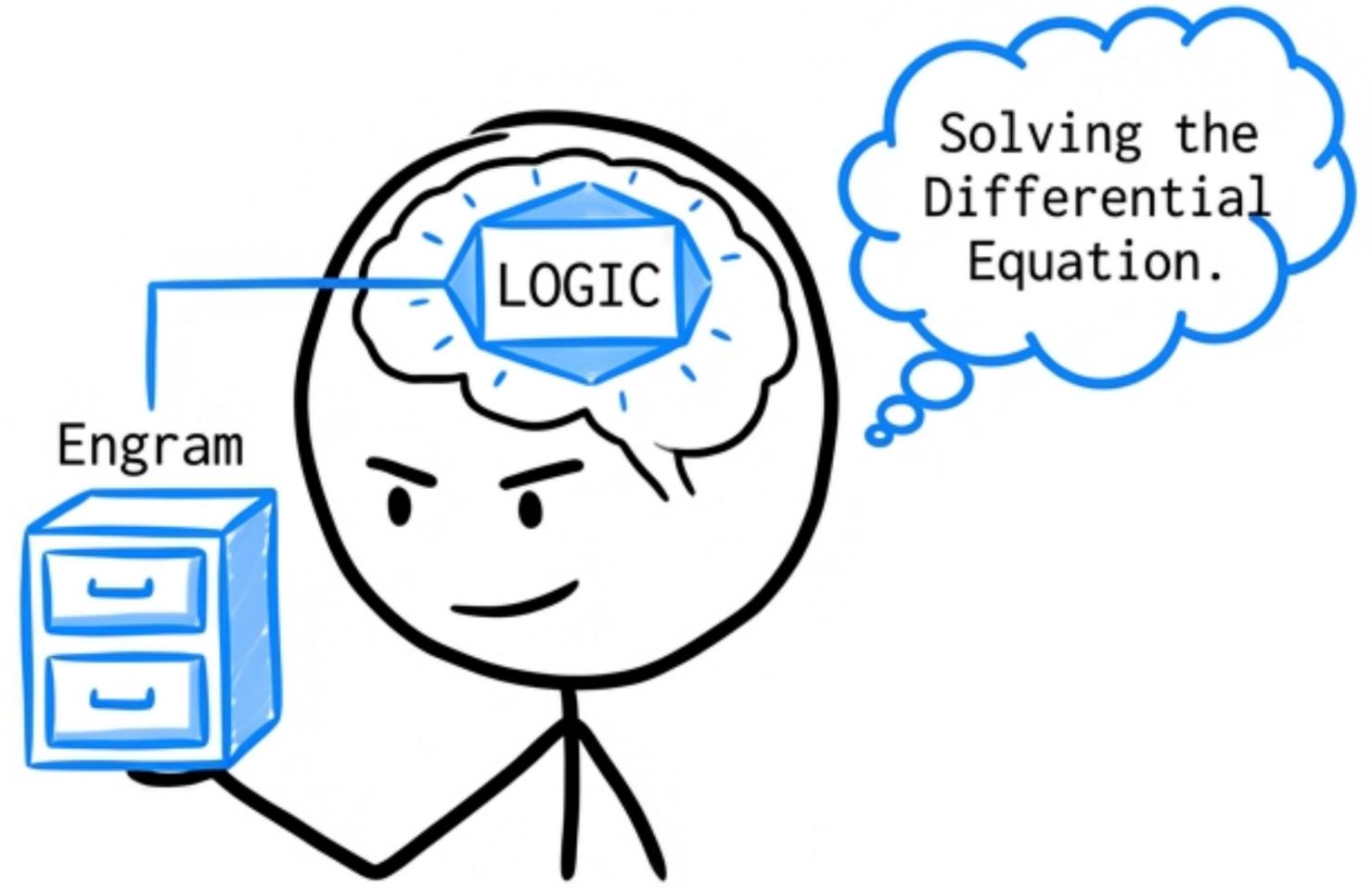
# Why does Memory help Reasoning?

## Standard Model



Brain clogged with static patterns.

## Engram Model

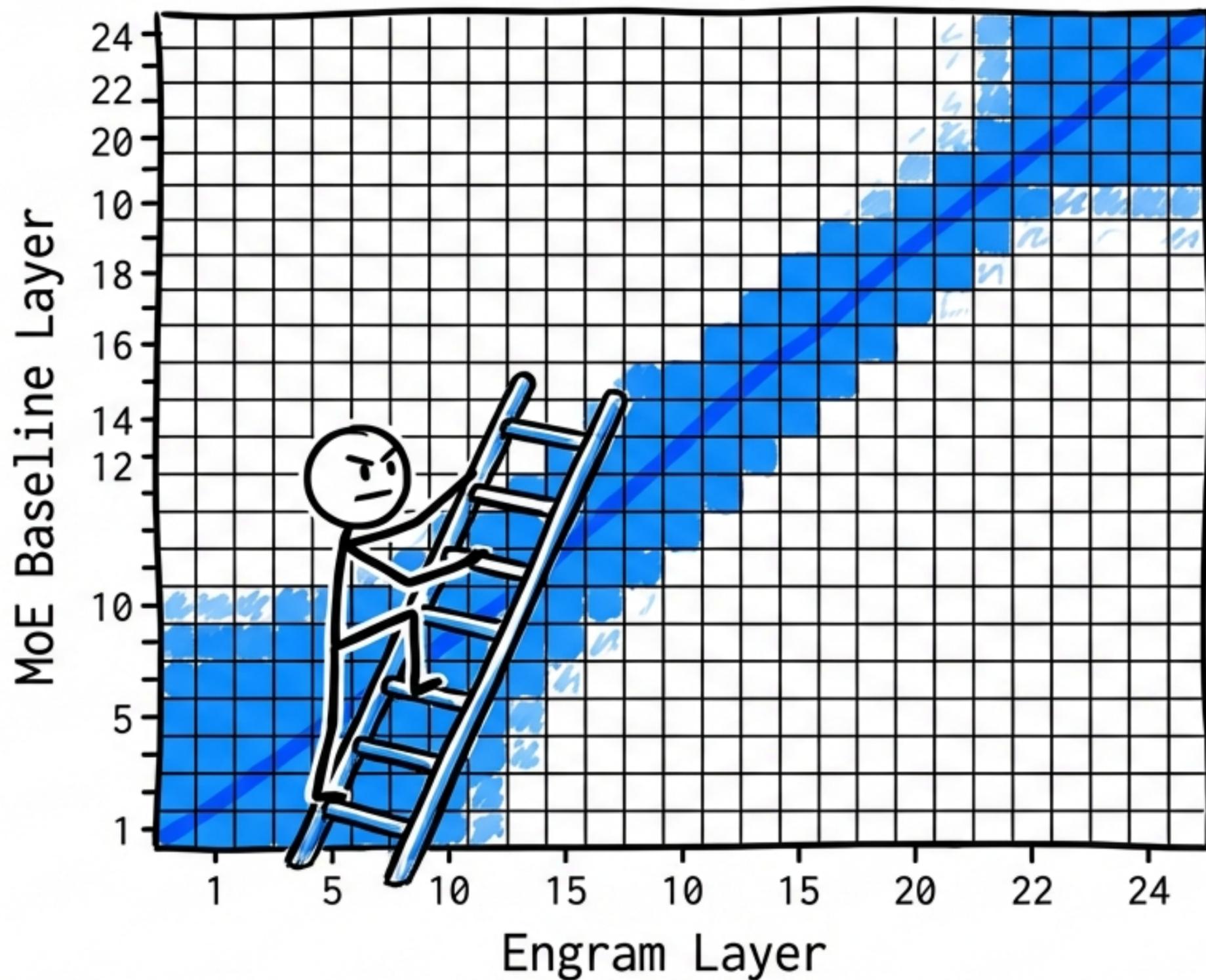


Static patterns offloaded to lookup. Brain free to think.

“Don’t sweat the small stuff. (Lookup the small stuff).”

# Proof: The Shortcut Mechanism

CKA Similarity heatmap

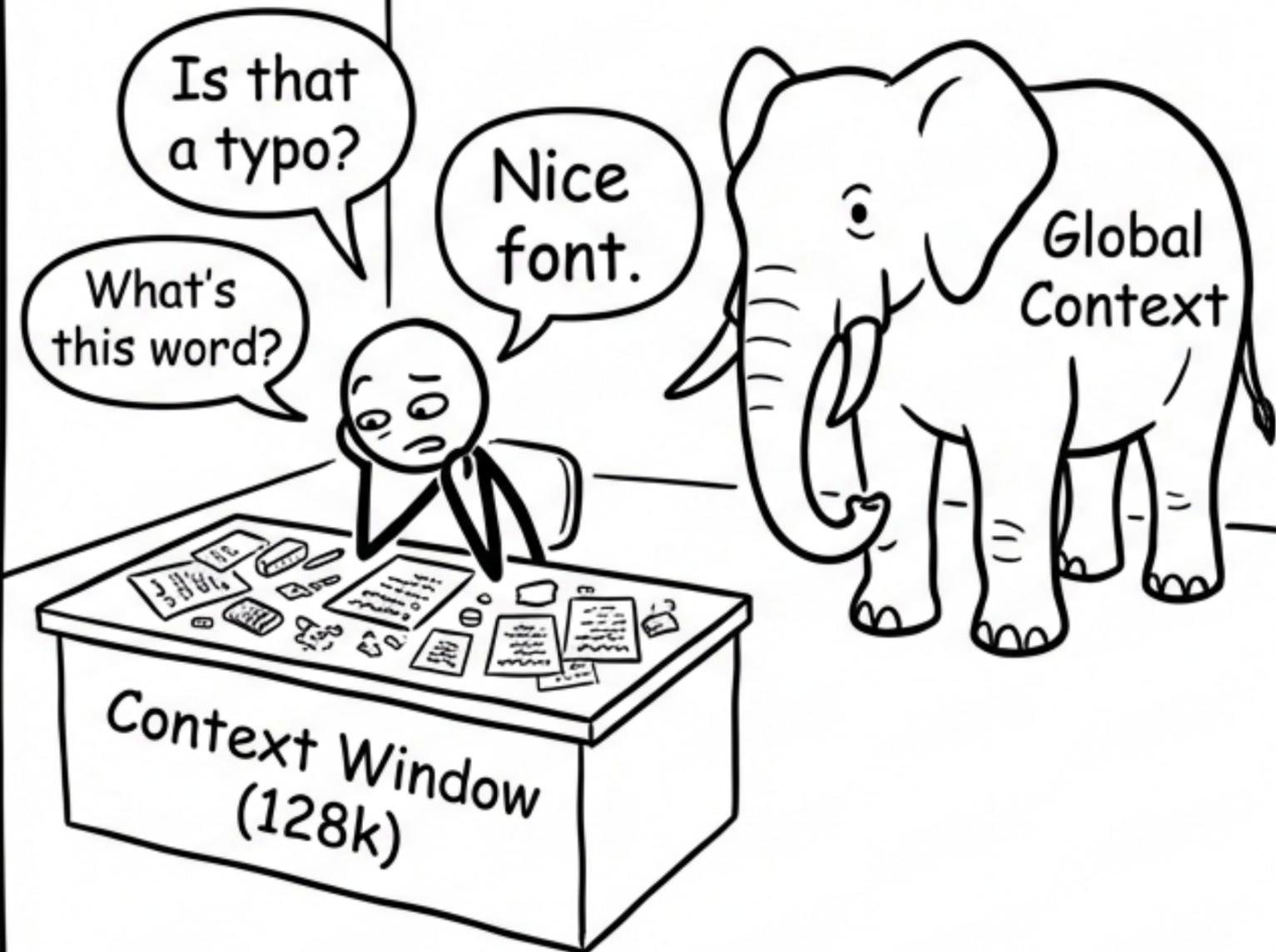


## CKA Analysis & LogitLens: Architects Daughter

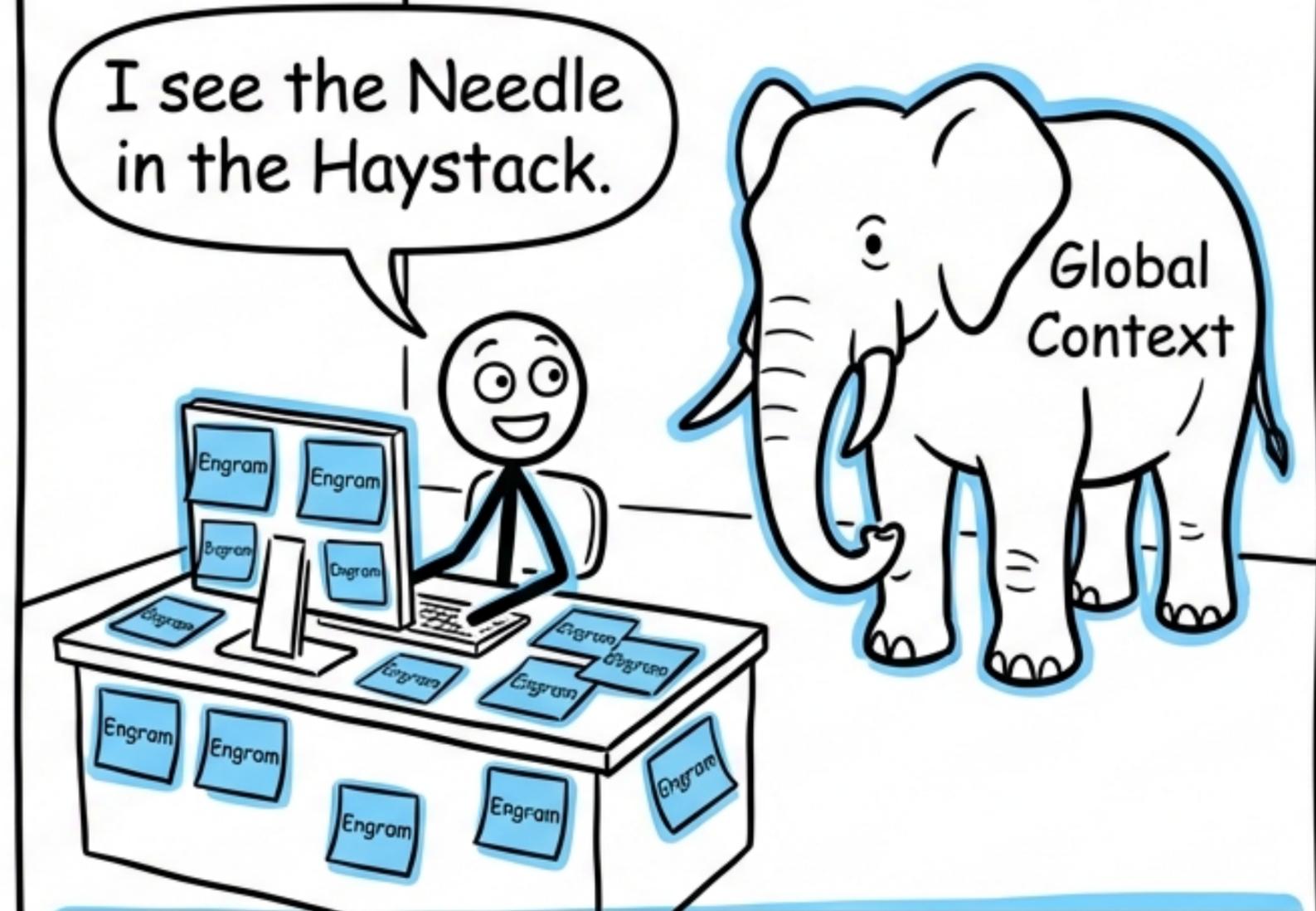
1. Engram's shallow layers are functionally equivalent to the the MoE baseline's deeper layers.
2. The model skips the "Who am I?" phase and gets straight to work.
3. Effective Depth > Actual Depth.

# Attention is Precious (Don't Waste It)

## Without Engram

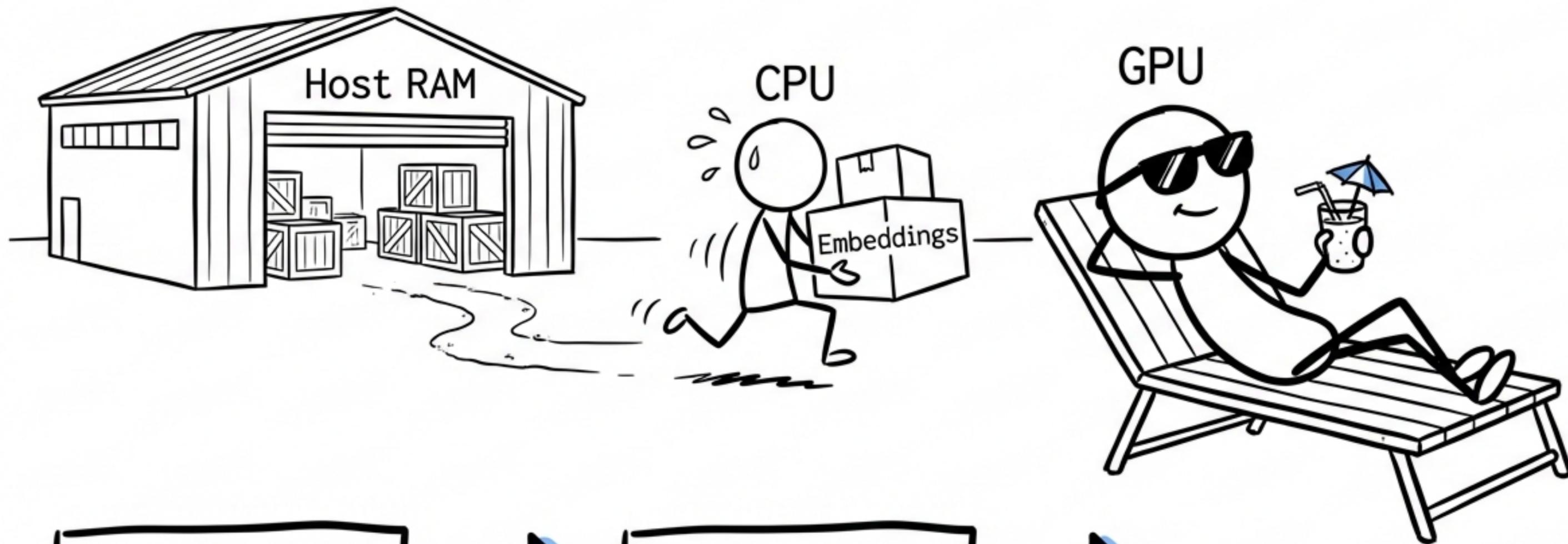


## With Engram



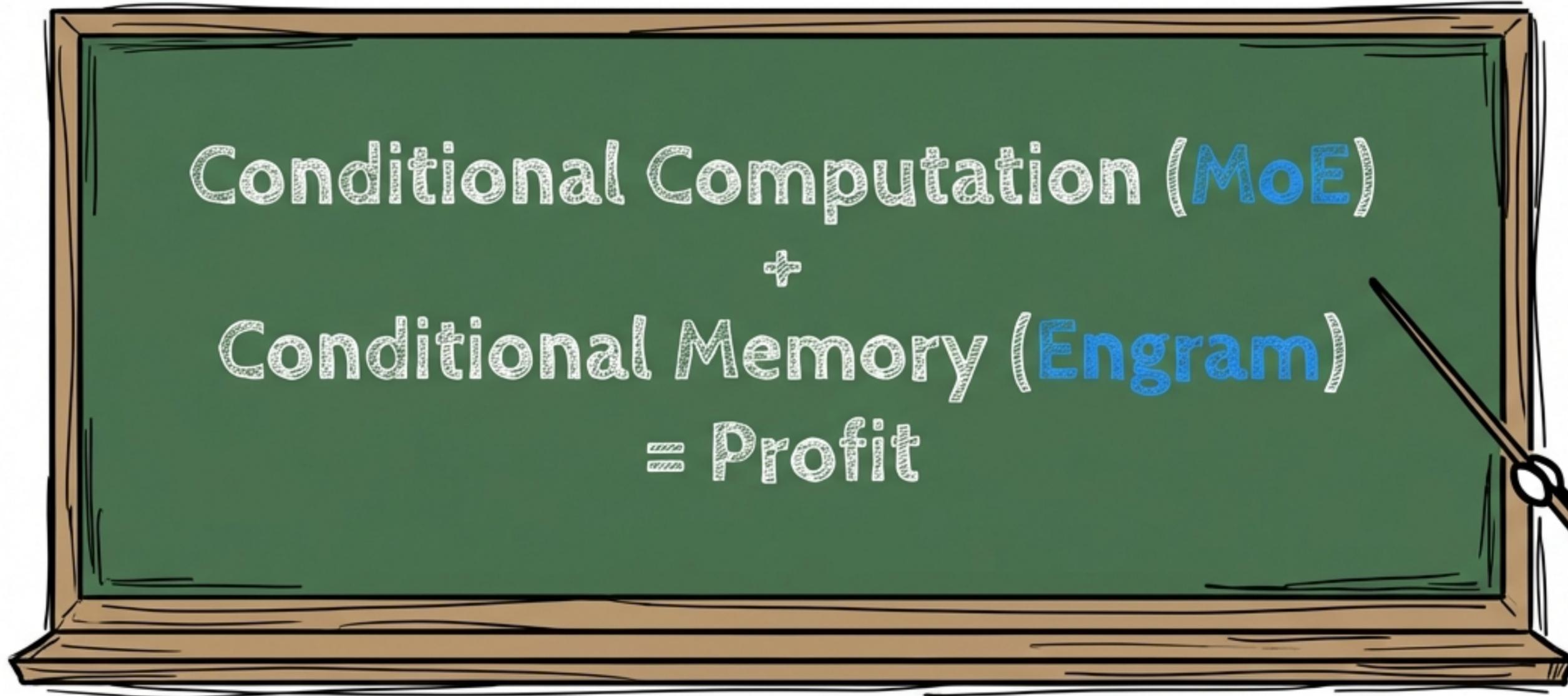
Multi-Query NIAH: 97.0 (Engram) vs 84.2 (MoE).  
Variable Tracking: 89.0 (Engram) vs 77.0 (MoE).

# The Lazy GPU Strategy (System Efficiency)



Overhead < 3% even with 100B-parameter table offloaded to CPU RAM. Because addressing is deterministic, we know what the GPU needs before it asks for it.

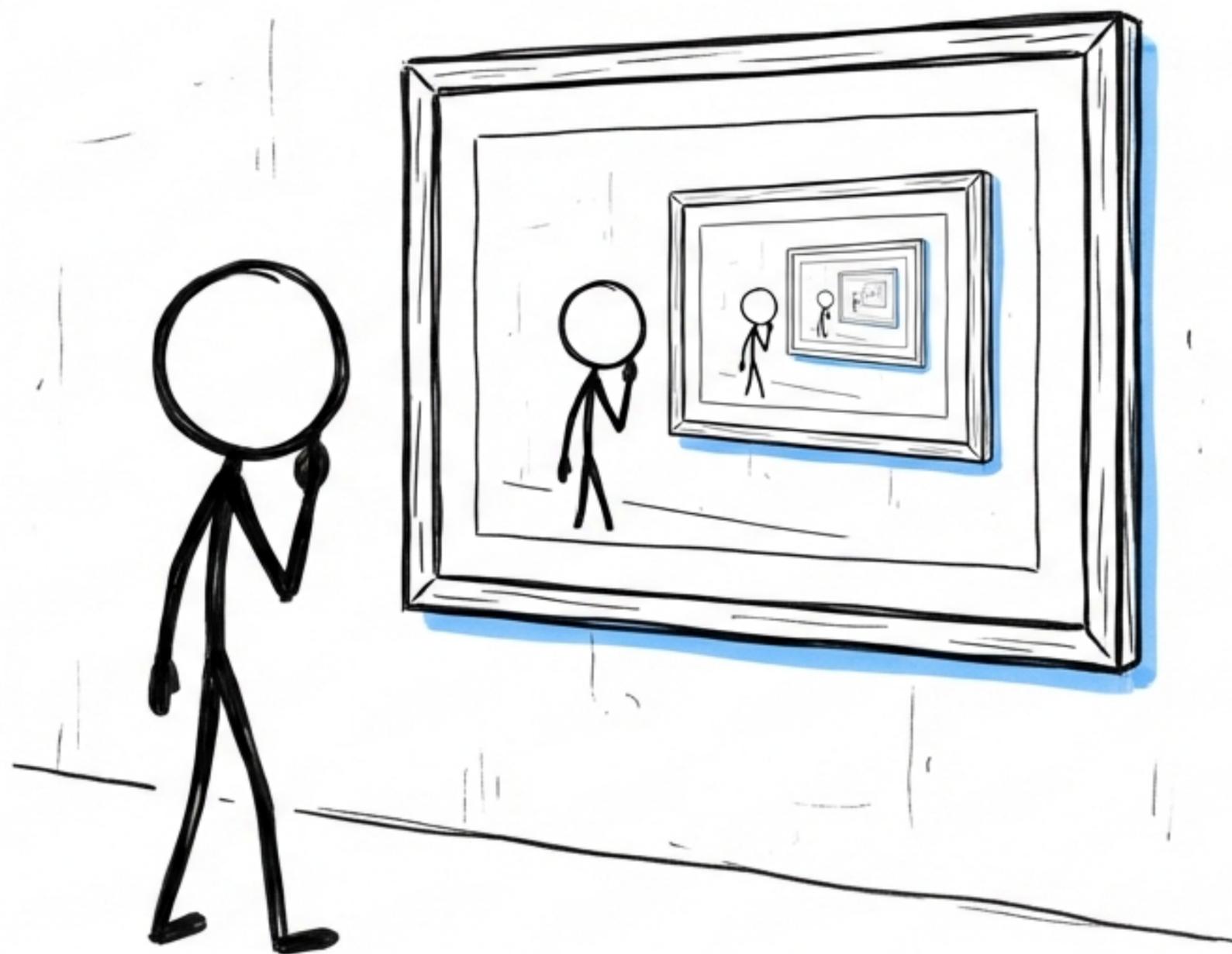
# The Takeaway



Conditional Computation (MoE)  
+  
Conditional Memory (Engram)  
= Profit

- Sparsity is Dual: Compute less, Remember more.
- The U-Curve: Optimizing the trade-off yields superior performance.
- Generalization: Helps Reasoning and Long Context, not just Trivia.
- Efficiency: Decouples storage from compute.

# Future Work



We envision **conditional memory** as an **indispensable primitive** for **next-gen models**.

We also envision eventually forgetting how to multiply because the **Engram module** does it for us.

Code available at: <https://github.com/deepseek-ai/Engram>