So you want to model cognition?

A Guided Tour with MINERVA2

Sydelle de Souza

CCS Guest Lecture | 10th November 2025







Every modeling choice is a theoretical commitment

3

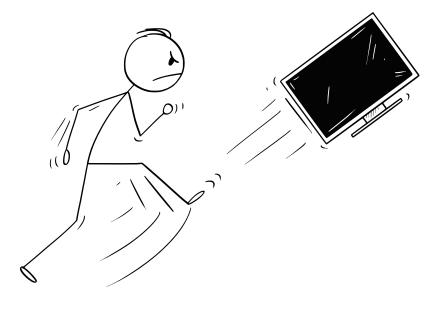
Cognitive modeling starts with a research question, not a dataset.

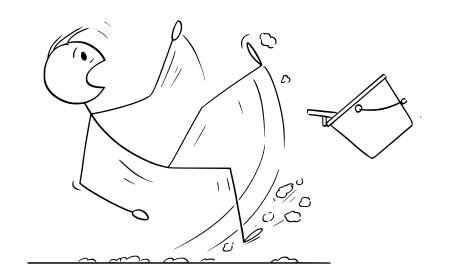
BACKGROUND & RQS

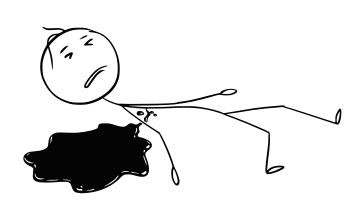
Human language is fundamentally compositional.



utterances





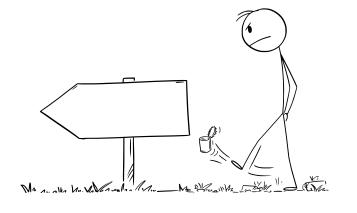




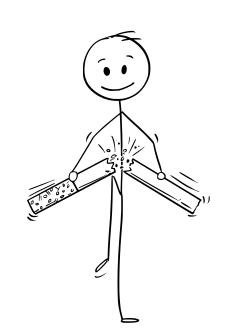


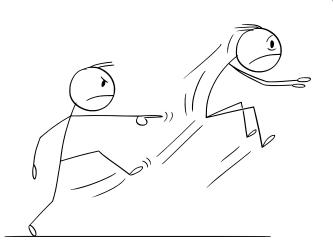


(v. transitive)





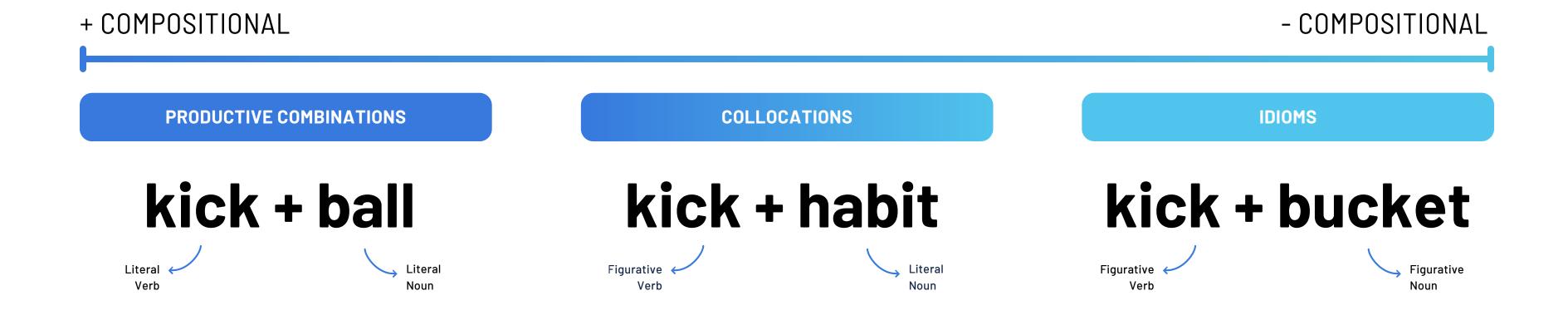






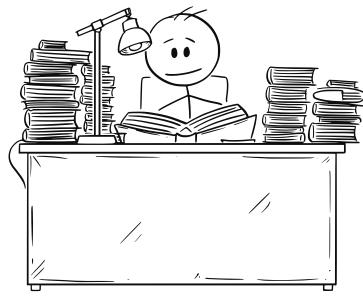
SEMANTIC COMPOSITIONALITY

The degree to which the meanings of the constituent words contribute to the overall meaning of the expression.





How do humans acquire and process semi-compositional language?







Abstraction



Generalization

freeze $\{ \bullet (, \bullet, \bullet) \rightarrow \{ \bullet (, \bullet, \bullet) \}$ becomes cold enough to become solid

freeze (v.)
to become cold enough to become solid
[+ reversible]

freeze
$$\{ , \exists \} \rightarrow \{ , \exists \}$$

becomes cold enough to become solid





Abstraction



Generalization

freeze $\{ \bullet \bullet, \bullet, \bullet \} \rightarrow \{ \bullet \bullet, \bullet, \bullet \}$ becomes cold enough to become solid

freeze(v.)

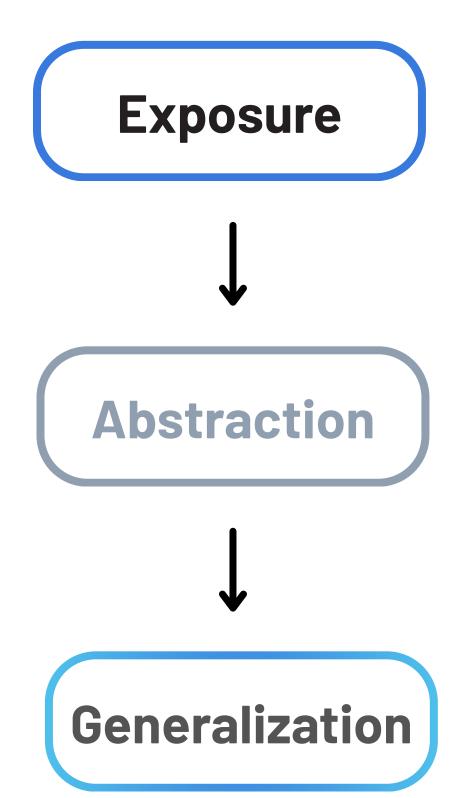
to become cold enough to become solid

liquid → solid

[+ reversible]

freeze
$$\{ , \bigcup_{\substack{N \text{ Nitrogen} \\ 14.00674}}^{7}, \bigcup_{\substack{N \text{ Nitrogen} \\ 14.00674}}^{7}}, \bigcup_{\substack{N \text{ Nitrogen} \\ 14.00674}}^$$

becomes cold enough to become solid



freeze
$$\{ \bullet (, \bullet, \bullet) \rightarrow \{ \bullet (, \bullet, \bullet) \}$$

becomes cold enough to become solid

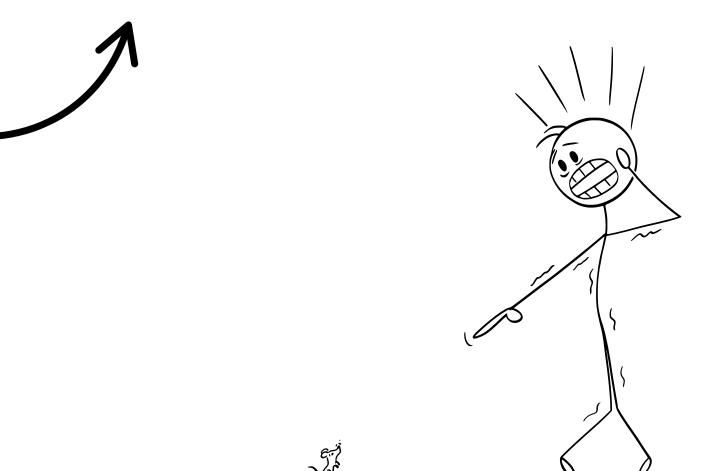
freeze {account, time} → {account, time}

LEARNING

freeze LITERAL

- to become cold enough to become solid
- liquid → solid
- fluid → rigid
- not permanent; reversible

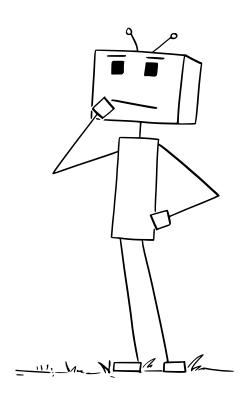




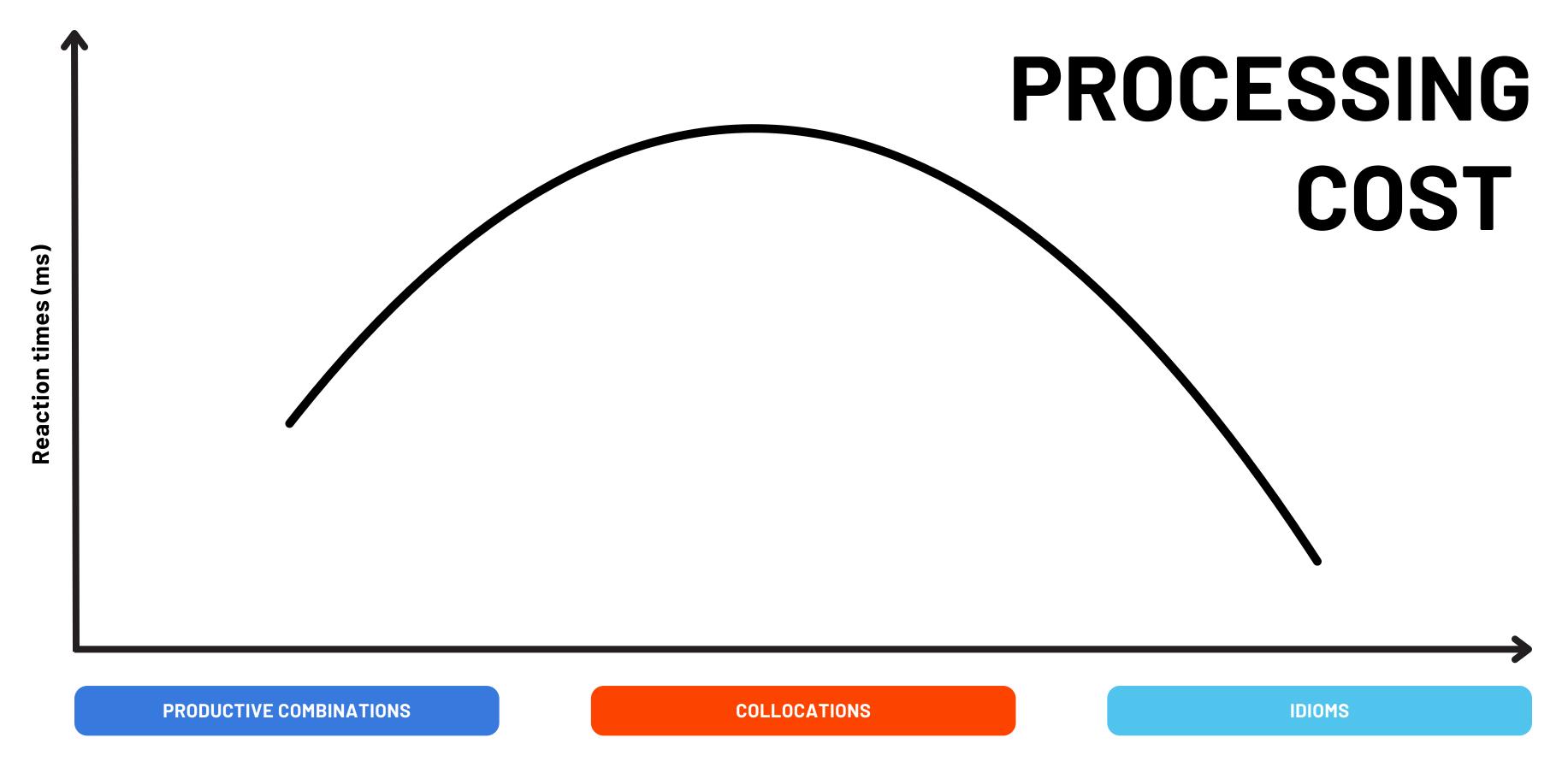
COLLOCATIONS POSE A LEARNING CHALLENGE



L2 speakers



Machines



FREQUENCY

- Roughly half of human language is comprised of compositional units
- Collocations are the largest subset of figurative language

+ COMPOSITIONAL - COMPOSITIONAL

PRODUCTIVE COMBINATIONS

COLLOCATIONS

IDIOMS

chase dreams

strong coffee

heavy rain

break promises

swing voters



flood airwaves

read minds

spill secrets

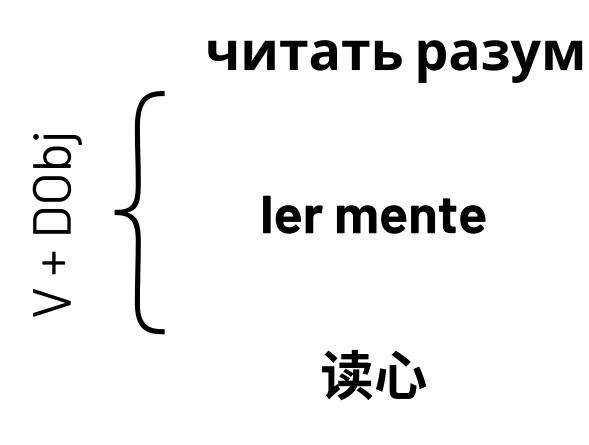
freeze accounts

harbour grudges

fight poverty

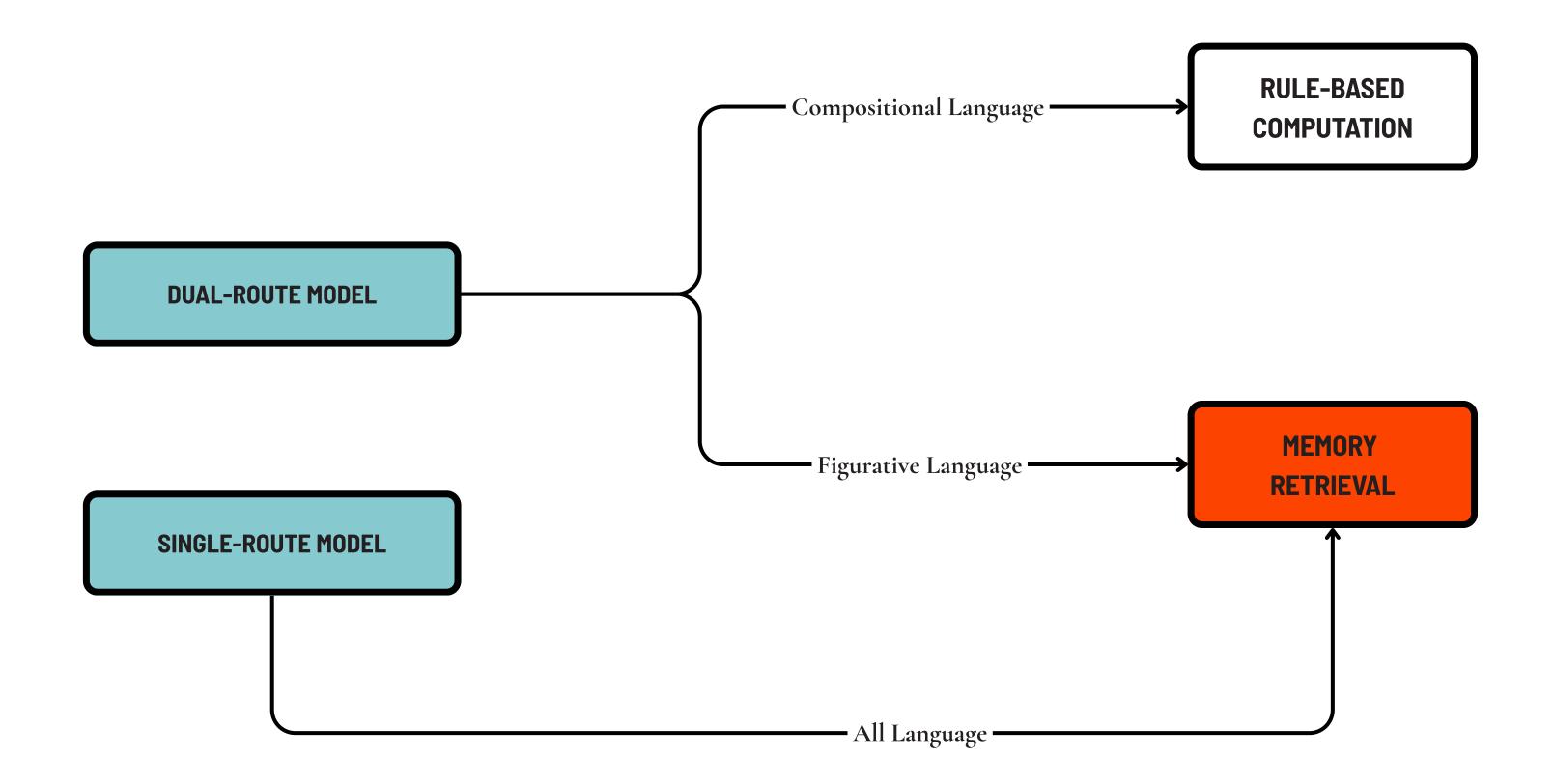
throw parties

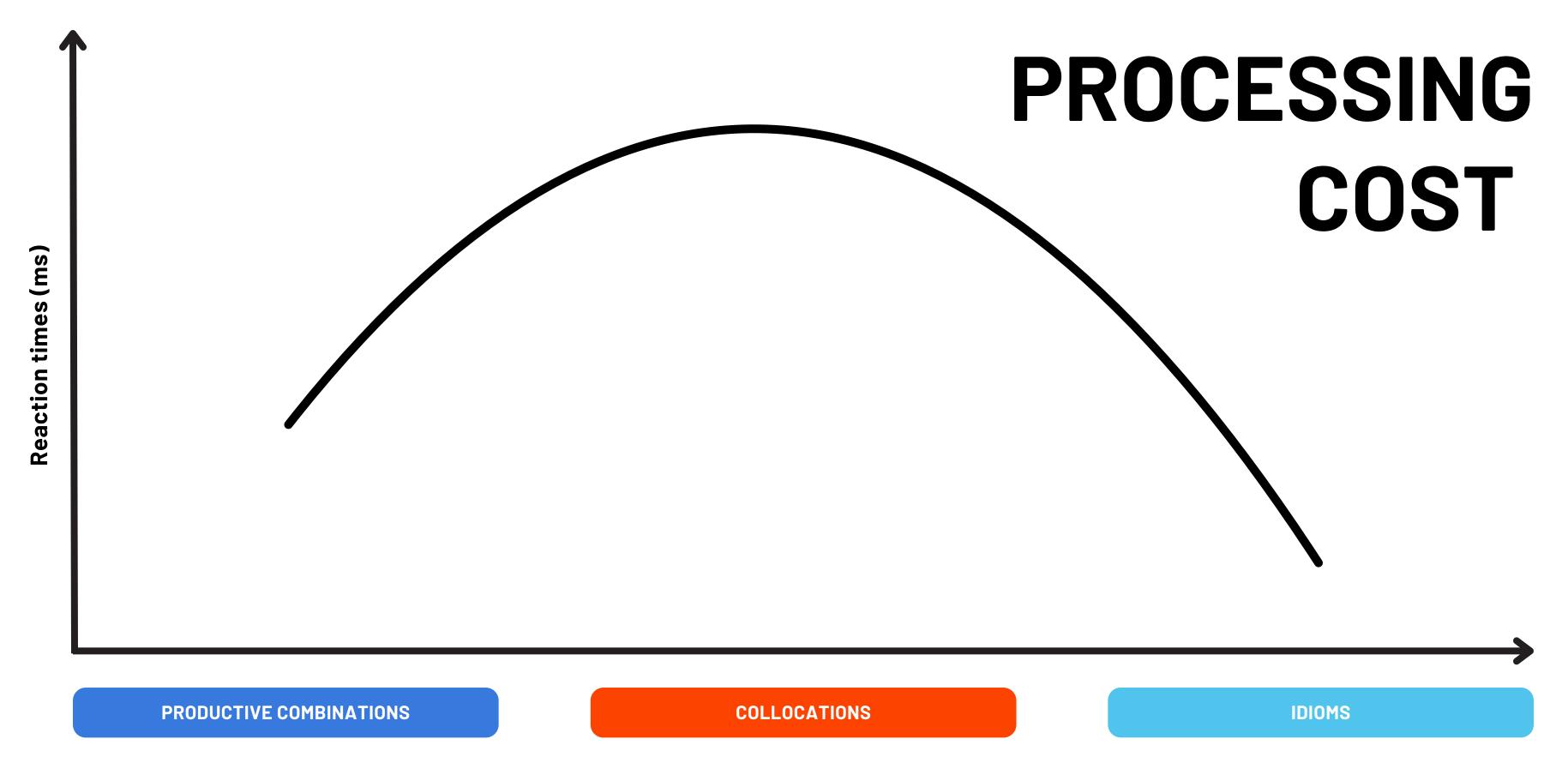
read + mind



മനസു വായിക്കുക DObj + V Why would linguistic structures that are hard to process and hard to learn emerge so ubiquitously in human language?

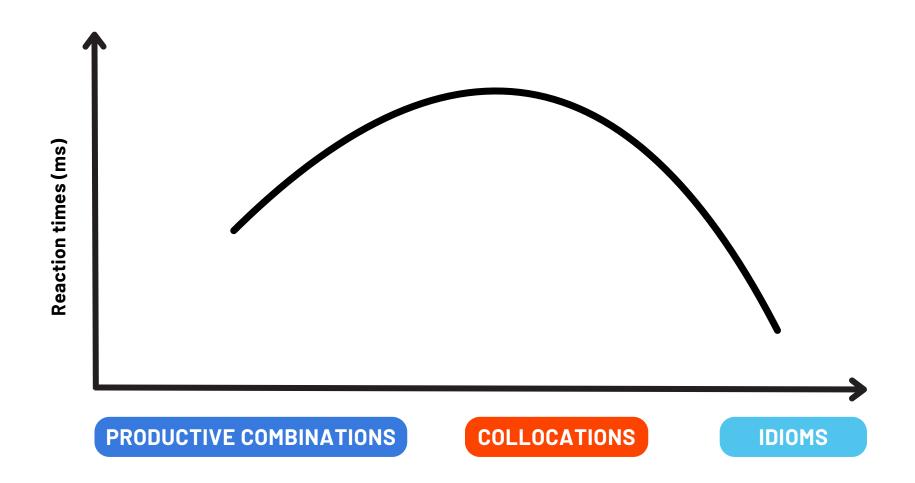
DUAL-ROUTE MODEL LANGUAGE **PROCESSING SINGLE-ROUTE MODEL**





AIM

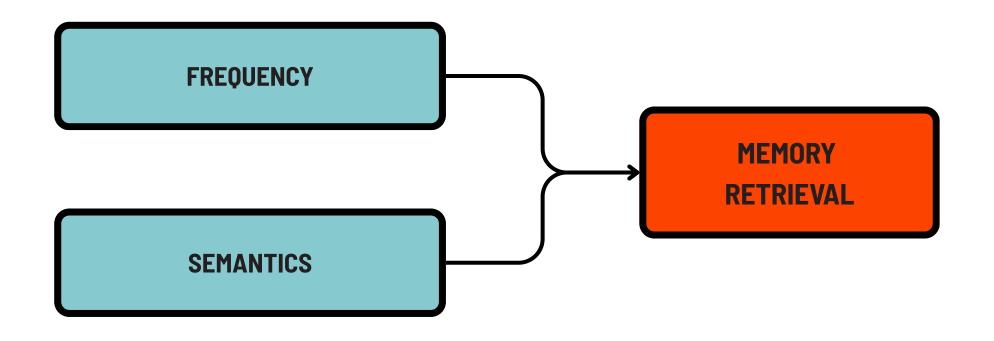
To test whether pure memory retrieval can account for the processing differences observed in humans

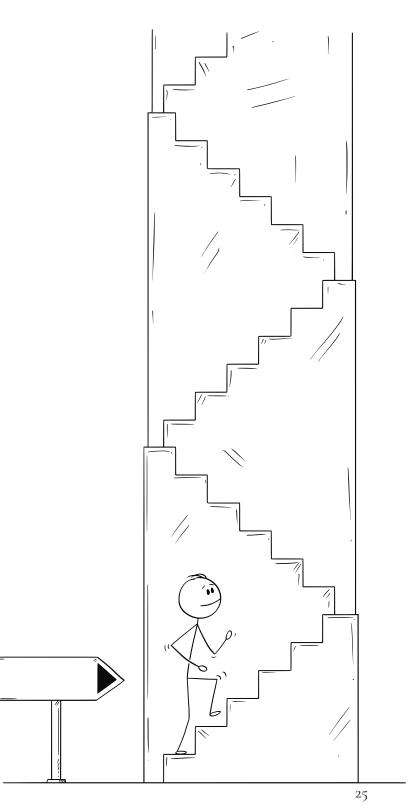


SPECIFIC RQS

- Do collocations really take longer to process?
- If yes, to what extent can pure memory retrieval account for the processing differences?

EMERGENT VARIABLES





BEHAVIORAL EXPERIMENT



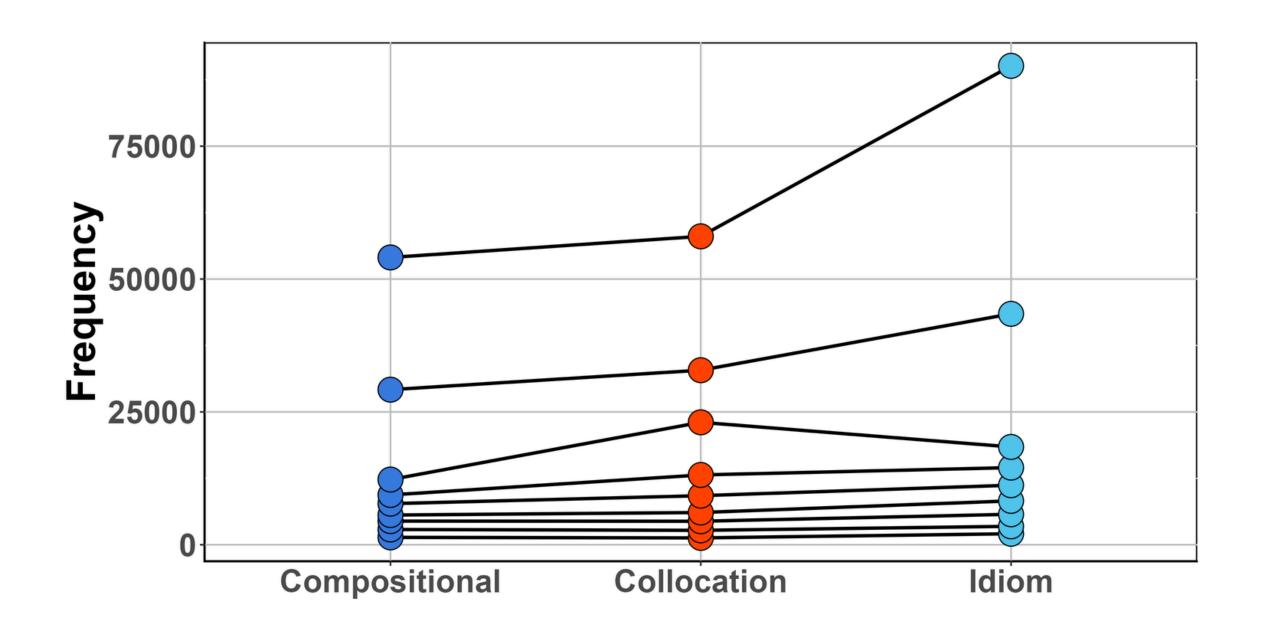
A good task elicits behavior that reflects the cognitive process you want to model.

STIMULI

240 Verb + Noun Items

- 80 Productive
- 80 Collocations
- 80 Idioms





We tried (unsuccessfully) to frequency match items!

STIMULI

240 Verb + Noun Items

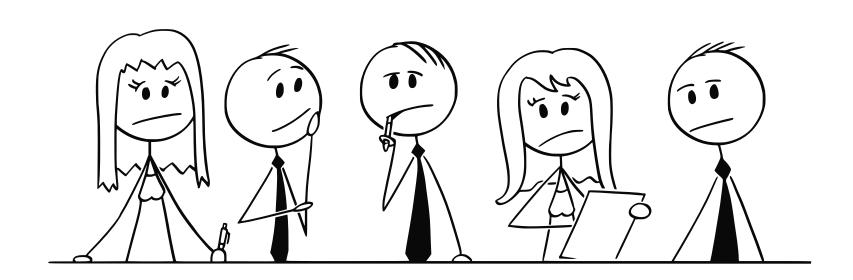
- 80 Productive
- 80 Collocations
- 80 Idioms



PARTICIPANTS

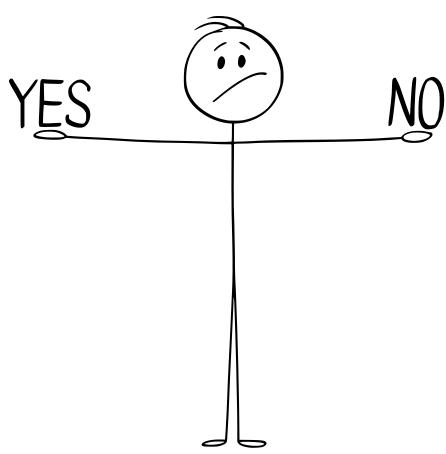
N = 186 humans

- L1 English speakers
- Mean Age: 38.57 (SD=10.81)
- 117 F, 71 M, 3 NB
- Prolific



ACCEPTABILITY JUDGEMENT TASK

- Implicit memory task
- Widely used in previous studies
- Each participant saw 160 items
- Individual randomized order
- 15 second break
- Training set with feedback



Are the word combinations that appear on the screen acceptable in English?

eat cakes

ACCEPTABLE

Productive Combination

eat cakes



read one's mind

ACCEPTABLE

Collocation

read one's mind





kick the bucket

ACCEPTABLE

Idiom

kick the bucket



(UNACCEPTABLE)

crack cakes

ACCEPTABLE

UNACCEPTABLE

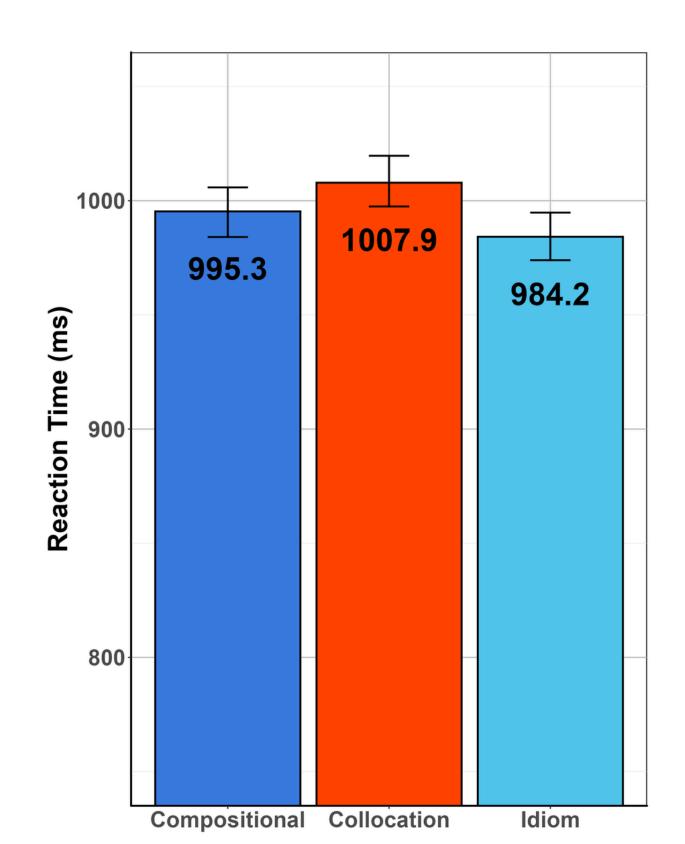
Filler

crack cakes





RESULTS



```
RT Condition + Phrasal Frequency
+ (1 | ID) + (1 | Verb)
```

Results:

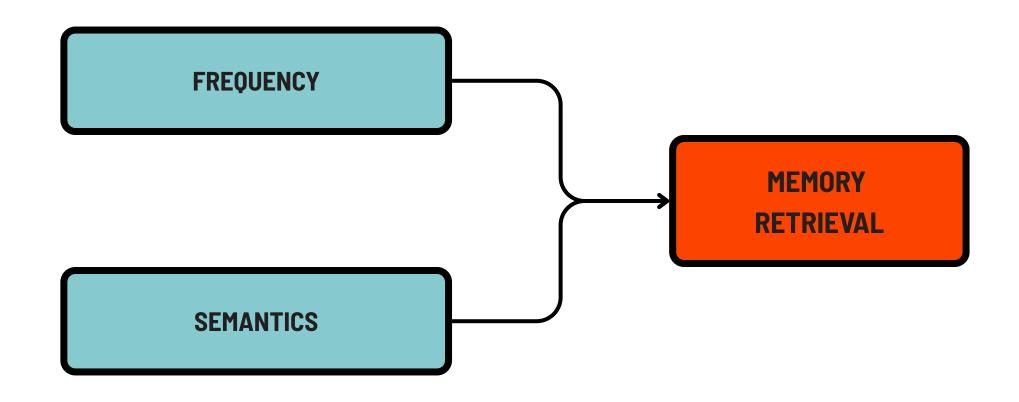
- Productive-Collocation ***
- Idiom-Collocation ***
- Idiom-Productive *

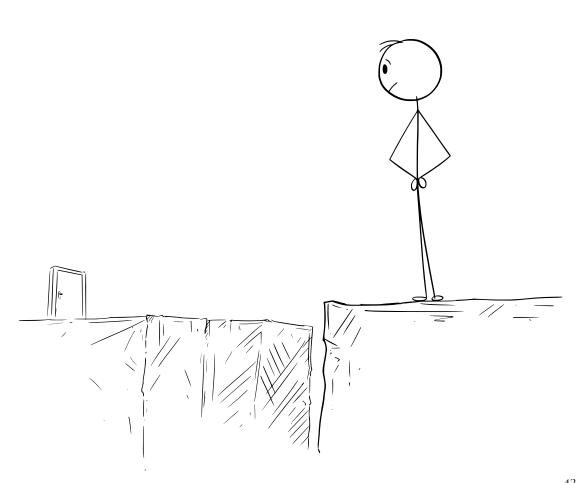
COMPUTATIONAL MODELLING



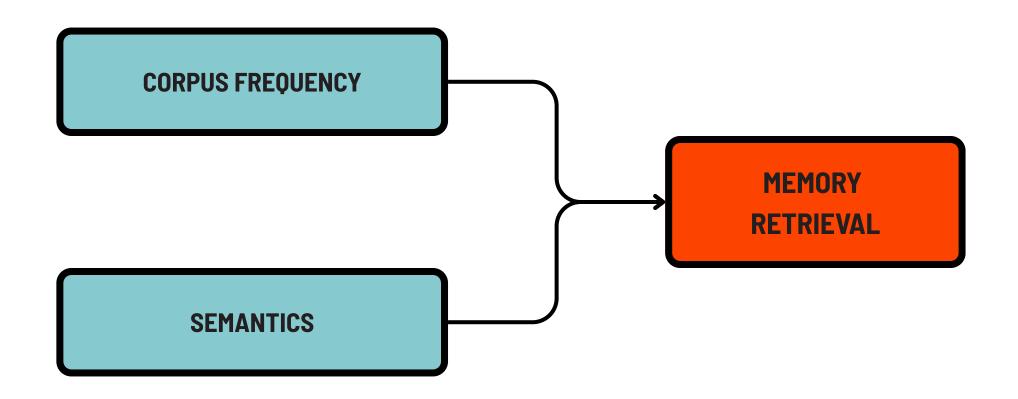
Modeling means formalizing what participants might be doing cognitively.

DESIDERATA



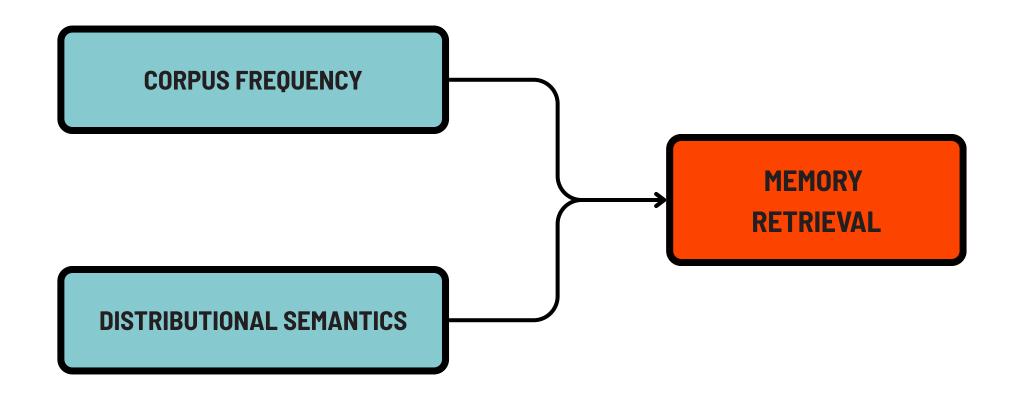


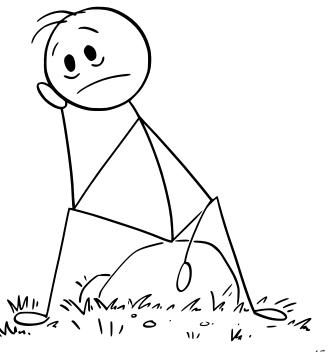
MODELLING FREQUENCY





MODELLING SEMANTICS

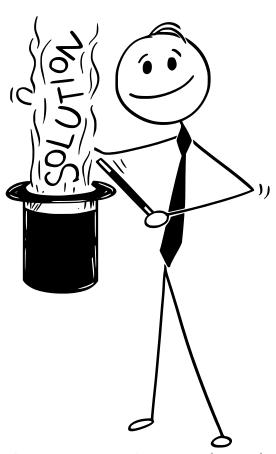




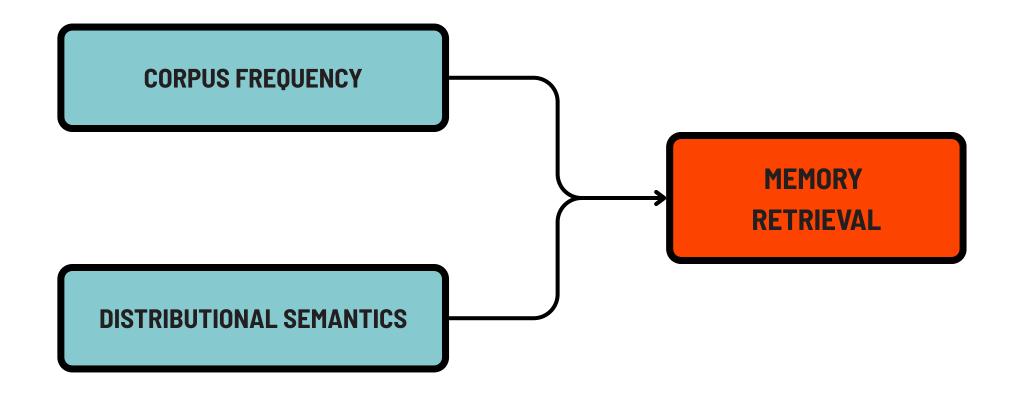
EMBEDDINGS

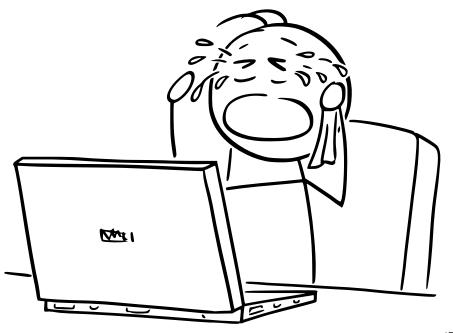
Contextual embeddings

- Collect 100 context sentences containing word combination from the corpus
- Embed sentences using <u>SentenceBERT</u>
- Pick out embeddings for verb and noun from each sentence
- Average verb and noun embeddings across sentences, separately
- Item embedding = concatenated embeddings of verb and noun

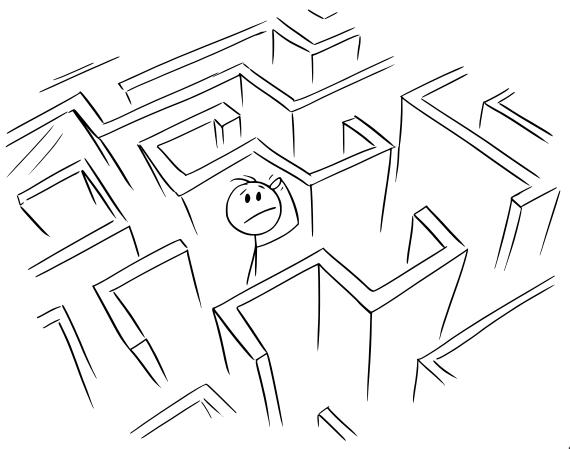


MODELLING MEMORY

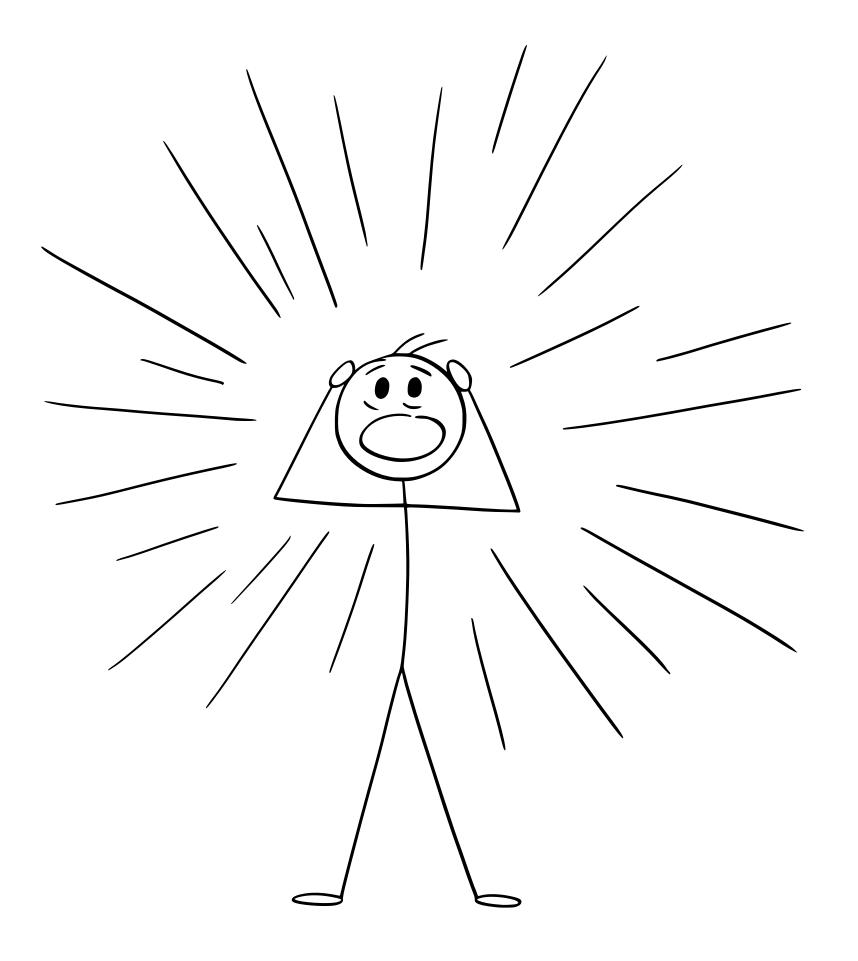




SEEMS LIKE EVERYONE AND THEIR MOM HAS BUILT A MEMORY MODEL!



Category	Models	Retrieval Mechanism	Major Limitations
Trace-Based / Global Matching	MINERVA 2, REM, SAM	Global similarity; sampling across stored traces	Ad hoc similarity metrics and scaling parameters
Context-Based / Temporal	TCM, CMR, eCMR	Reinstating prior context or temporal cues	Too many free parameters!
Connectionist	Hopfield, HRR, CLS	Pattern completion in distributed representations	Hard to interpret opaque retrieval dynamics
Probabilistic / Bayesian	BART, Bayesian Memory	Probabilistic inference under uncertainty	Only explains "optimal" behavior Computationally expensive; analytically intractable
Neural / AI-Inspired	Modern Hopfield, MANNs, RAG, NTM	Attention or key-value lookup in differentiable memory	Not cognitively plausible!
Hybrid / Cognitive Architectures	ACT-R, Soar, LEABRA, CHREST	Combine symbolic control with subsymbolic retrieval mechanisms	Complex and parameter-heavy Integrate multiple mechanisms without clear theoretical boundaries

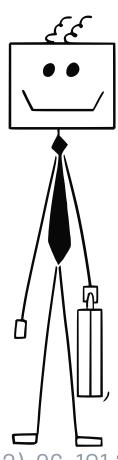


The goal isn't a perfect fit, but rather a transparent test of your theoretical assumptions.

MINERVA2

- A parsimonious, tractable model of episodic memory
- 3 core assumptions:
 - each experience leaves a trace of features
 - similar inputs → similar traces
 - retrieval driven by global similarity matching
- Integrates episodic and semantic memory through parallel, item-specific retrieval

TLDR: offers predictive precision while remaining computationally efficient



MINERVA2

- Psychologically grounded
- Makes testable predictions
- Successfully applied to diverse cognitive domains, including language:
 - Frequency judgments (Hintzman, 1988)
 - False memory (Arndt & Hirshman, 1998)
 - Artificial grammar learning (Jamieson & Mewhort, 2009)
 - Metaphor recognition (Reid & Jamieson, 2023)

TLDR: is empirically consistent and falsifiable



MINERVA2

$$s = sim(p, \mathbf{M})$$

$$a_{\tau} = s^{\tau} \operatorname{sign}(s)$$

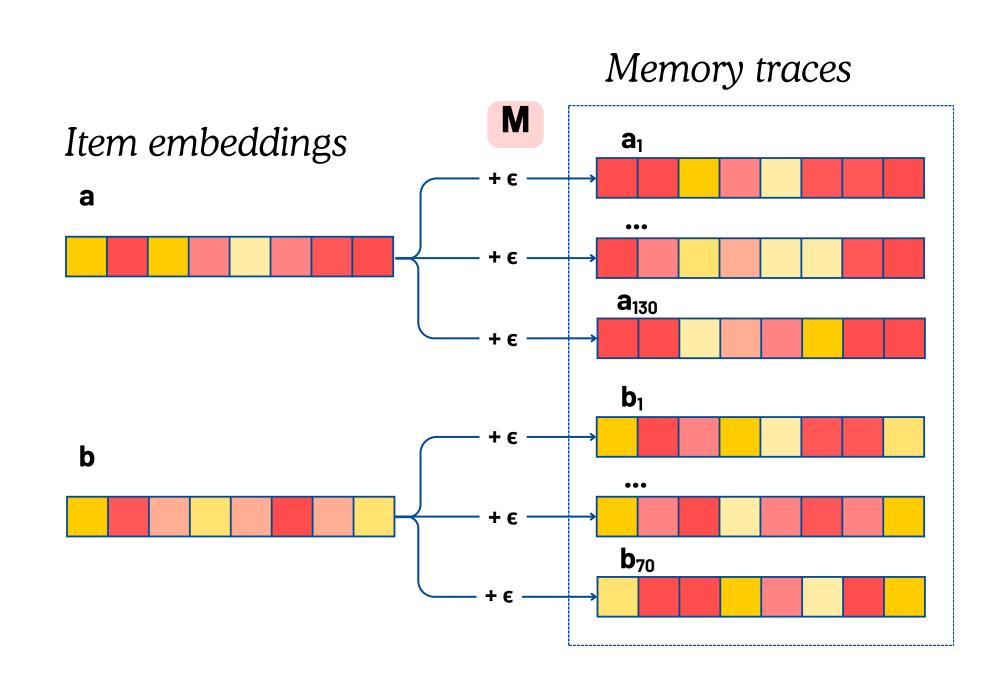
$$e_{\tau} = a_{\tau} M$$

$$f_{\tau} = sim(e_{\tau}, p)$$

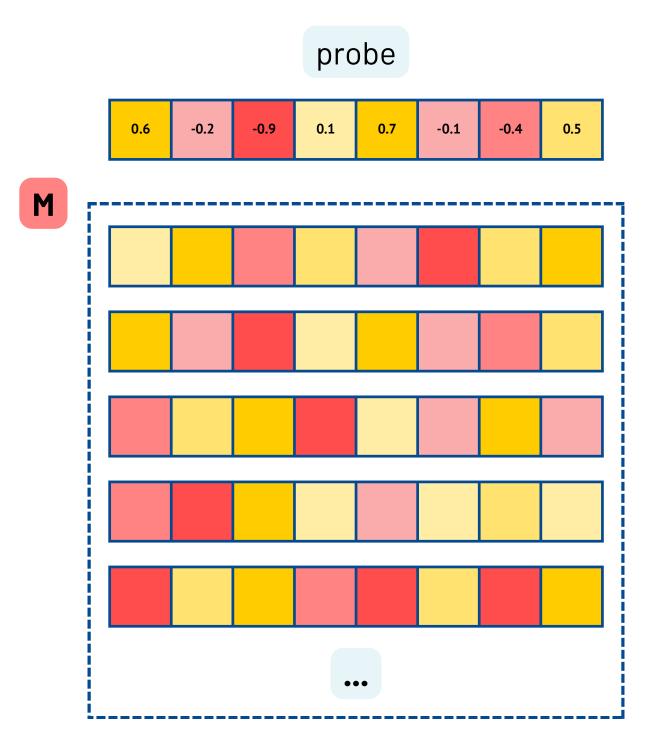
EXPOSURE + STORAGE

Sample item embeddings into MINERVA memory

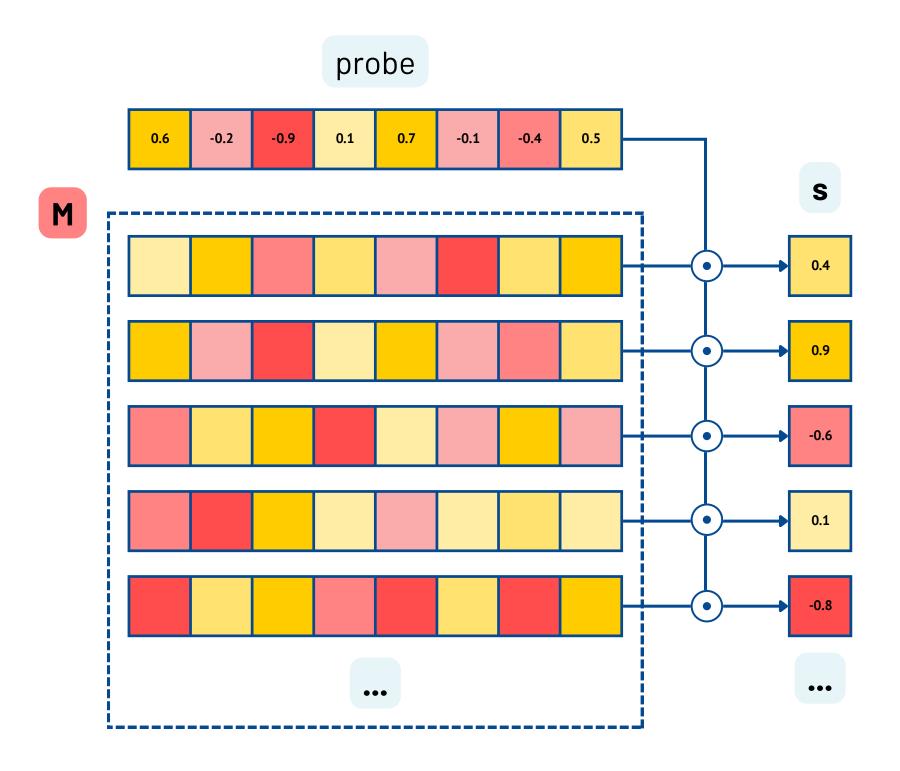
- according to corpus frequency (n)
- Each row is a d-dimensional memory trace.
- ullet $M\in\mathbb{R}^{n imes d}$
- Randomly noise each dimension of each embedding
 - Simulates forgetting



$s = sim(p, \mathbf{M})$



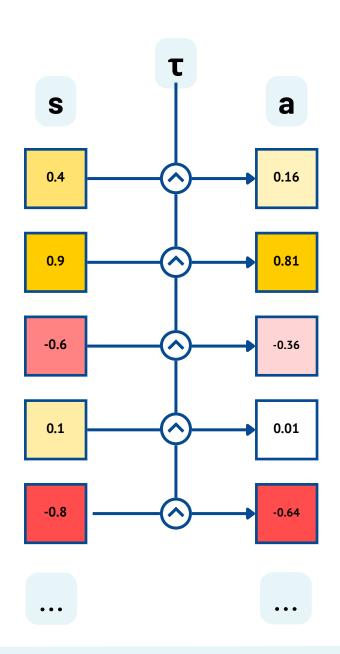
$s = sim(p, \mathbf{M})$



Compute similarity between the probe and each item in memory → activation

- ullet Probe: $p \in \mathbb{R}^d$
- represents current input

$a_{\tau} = s^{\tau} \operatorname{sign}(s)$

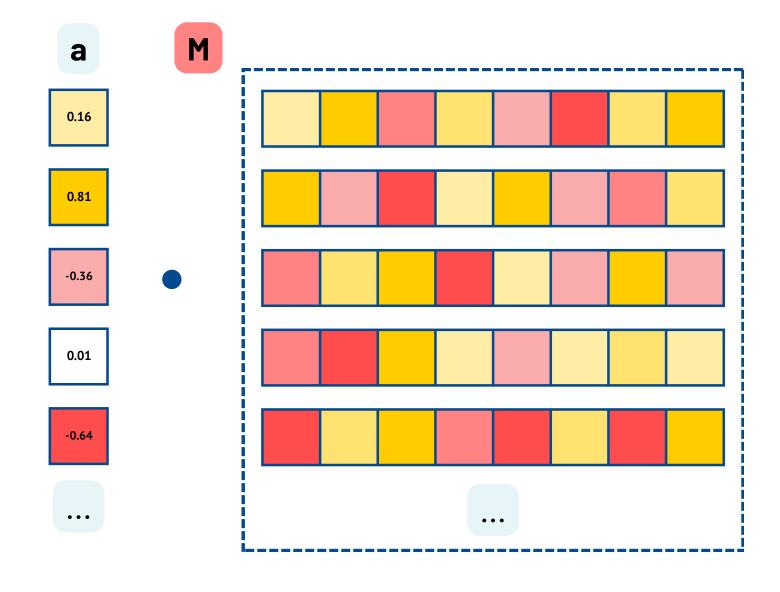


τ-accentuated similarities

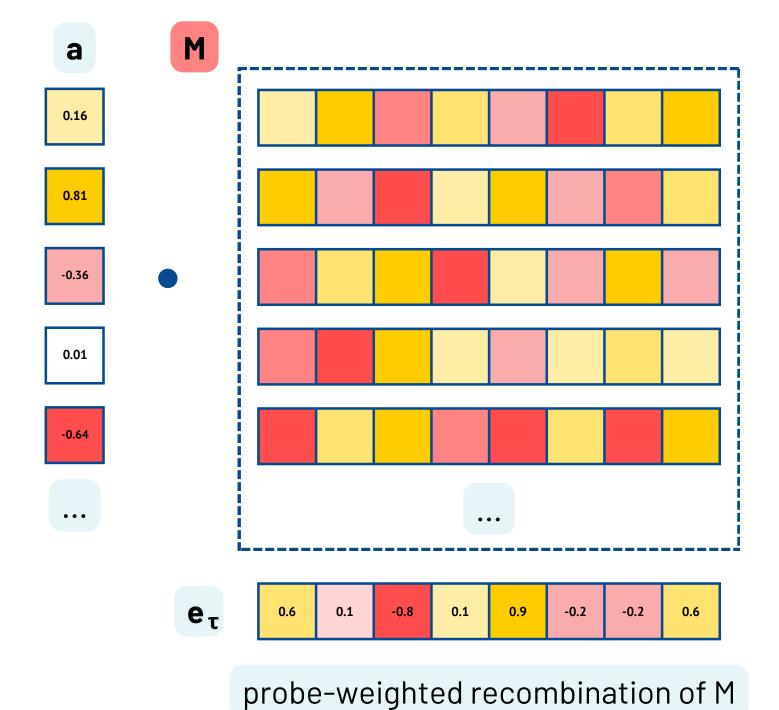
- Sharpen activations by exponentiating them by Tau
- High magnitudes stay high, low magnitudes get decayed
- Temporal index, i.e., no. of iterations required to reach a threshold

$$e_{\tau} = a_{\tau} \mathbf{M}$$

Average all items in memory according to the sharpened activation of each



$$e_{\tau} = a_{\tau} \mathbf{M}$$

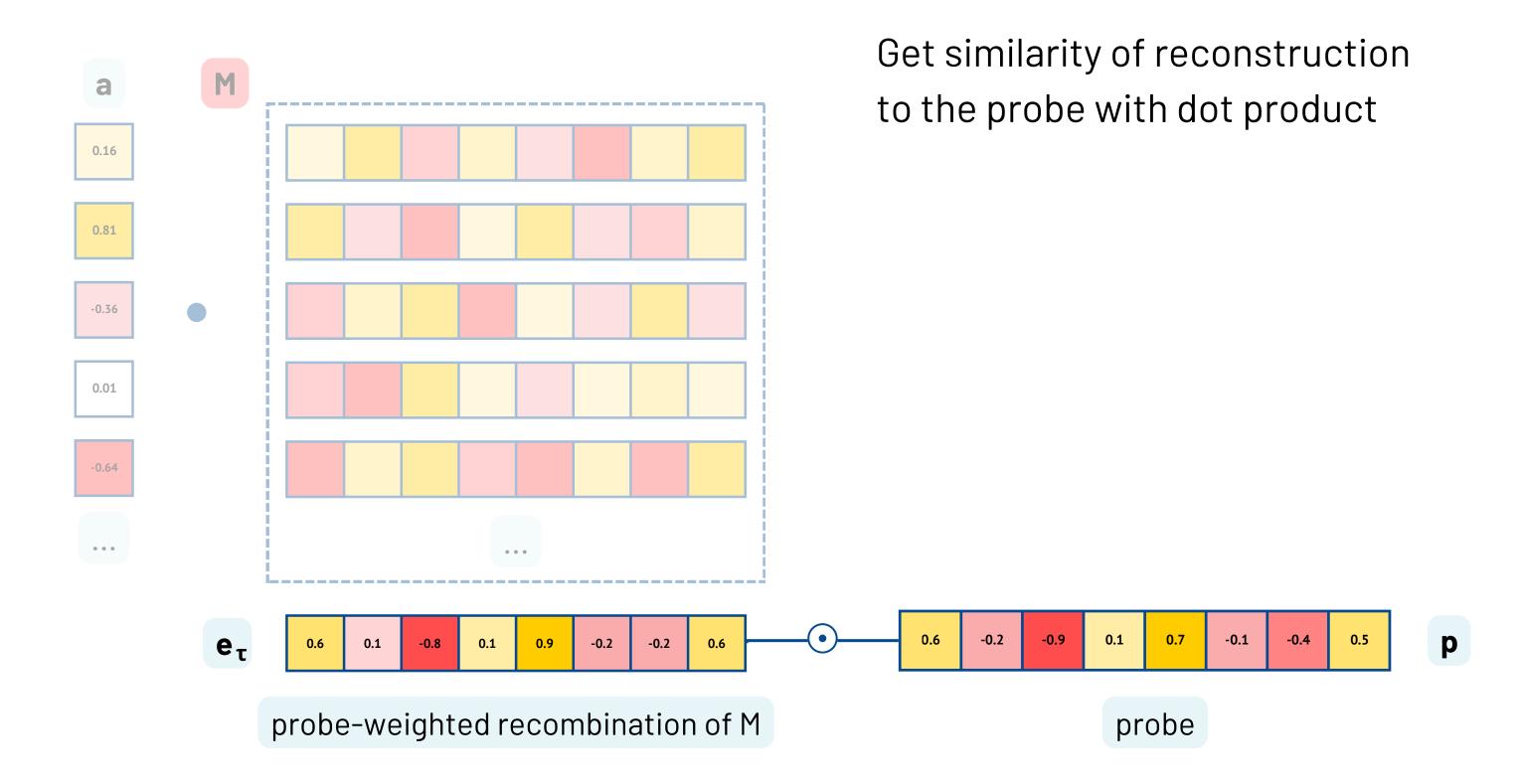


Average all items in memory according to the sharpened activation of each

- Uses the memory to reconstruct the probe
- i.e., reconstruct probe from past traces

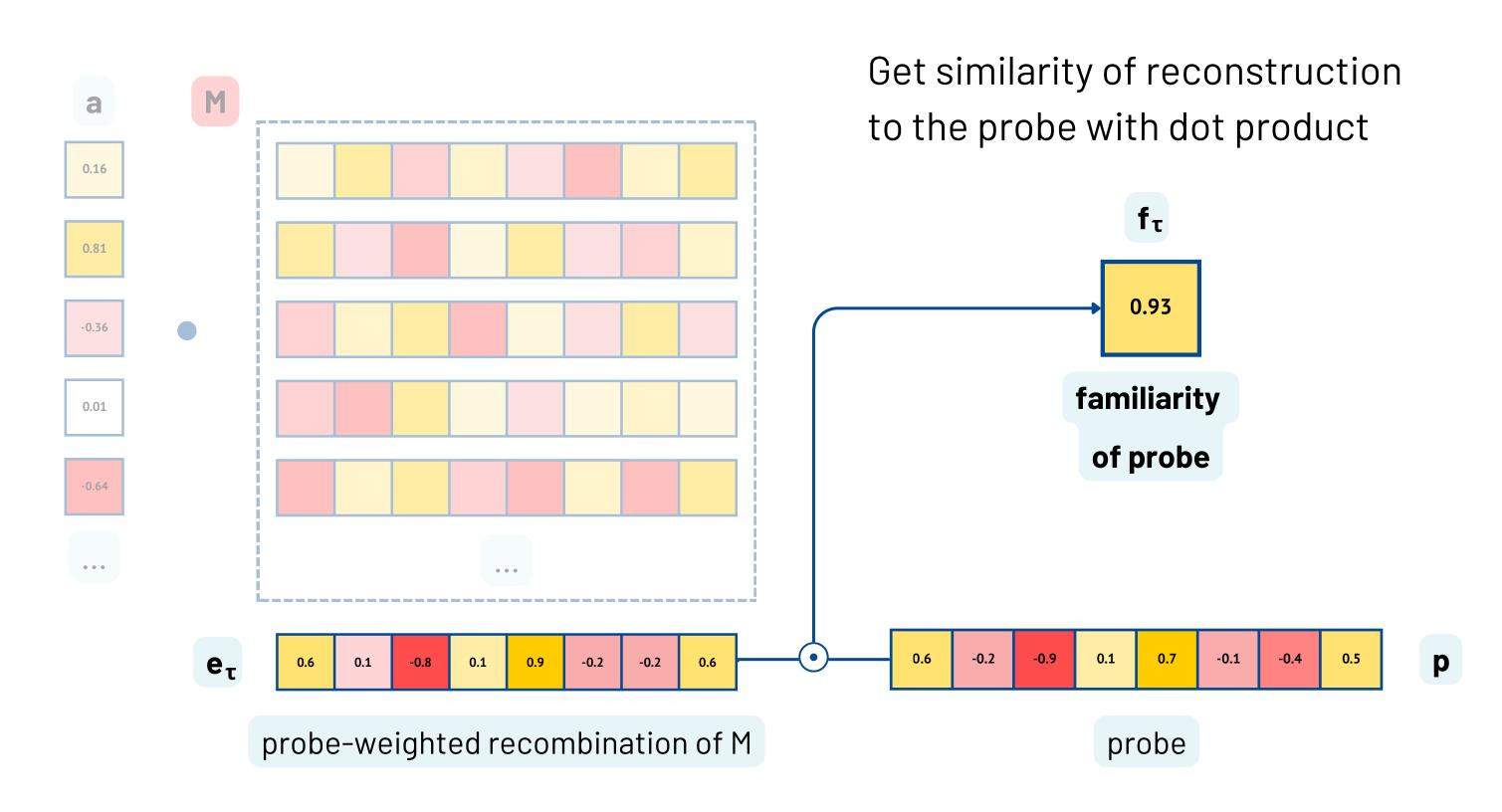
$$e_{\tau} = a_{\tau} \mathbf{M}$$

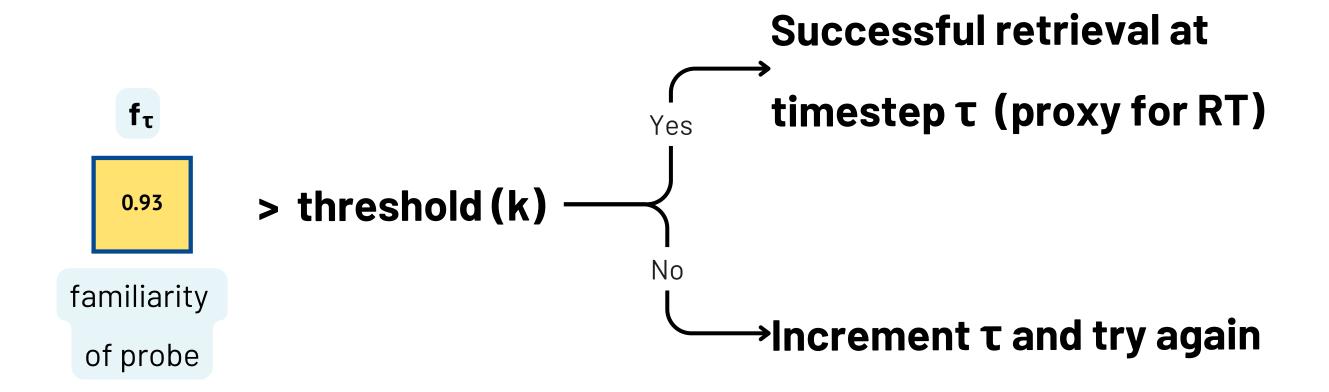
$$f_{\tau} = sim(e_{\tau}, p)$$



$$e_{\tau} = a_{\tau} \mathbf{M}$$

$$f_{\tau} = sim(e_{\tau}, p)$$





SIMULATIONS

• 2 Hyperparameters:

- Retrieval threshold k (cosine sim > k = successful retrieval)
- Forgetting probability (noise probability)

• Attempt to retrieve each item

- Keep track of tau for successful retrieval
- Timeout (max tau exceeded) = unsuccessful retrieval

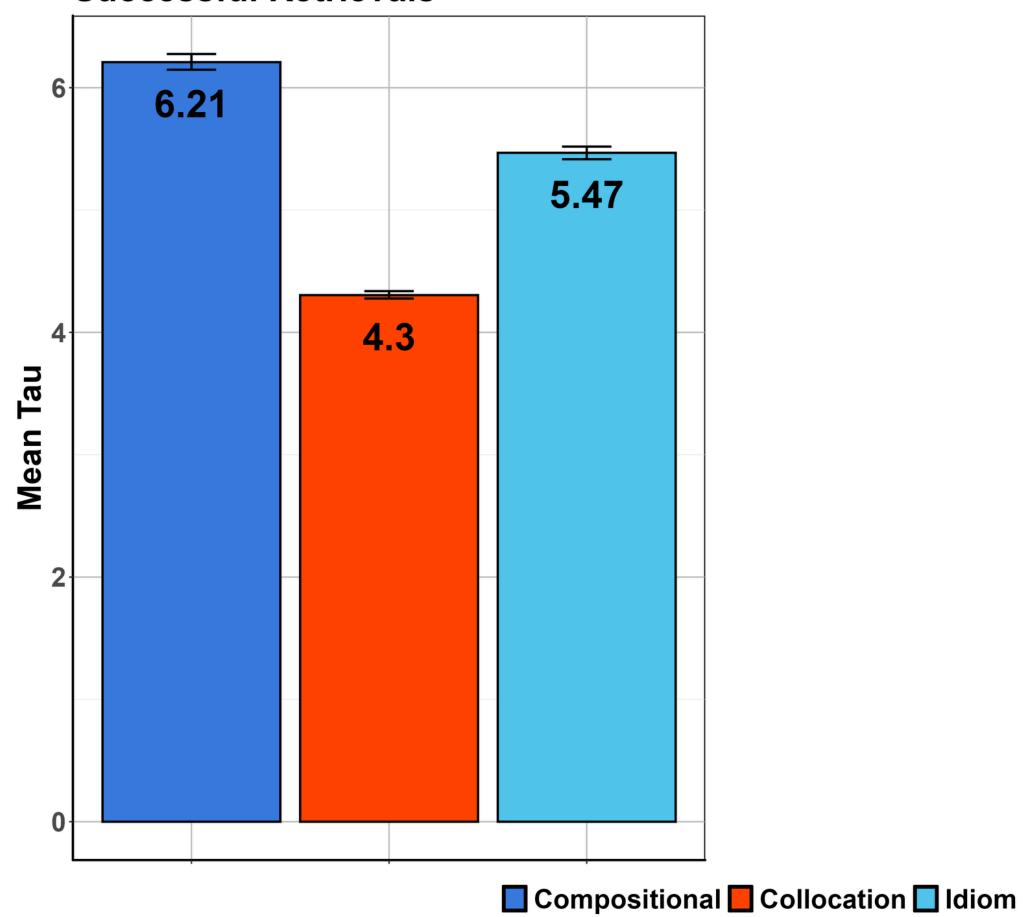
EXPERIMENTAL MANIPULATIONS

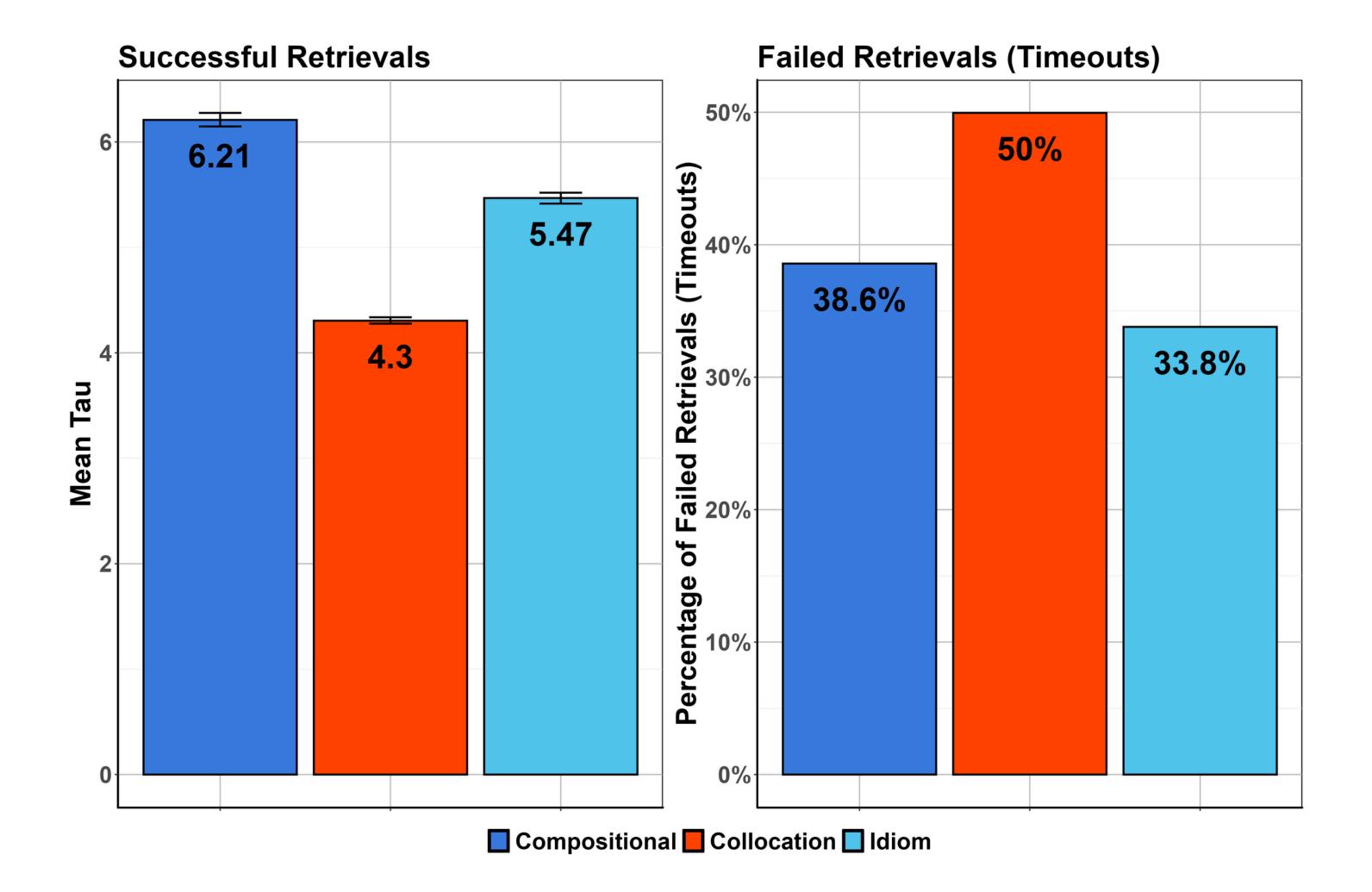
- Main Experiment: Frequency & Semantics
- Sanity Checks:
 - Frequency-only
 - Embedding vectors for each item replaced with unique random noise
 - Semantics-only
 - Each item is sampled into MINERVA memory in equal proportions
 - Null model
 - Equal frequency AND noise embeddings

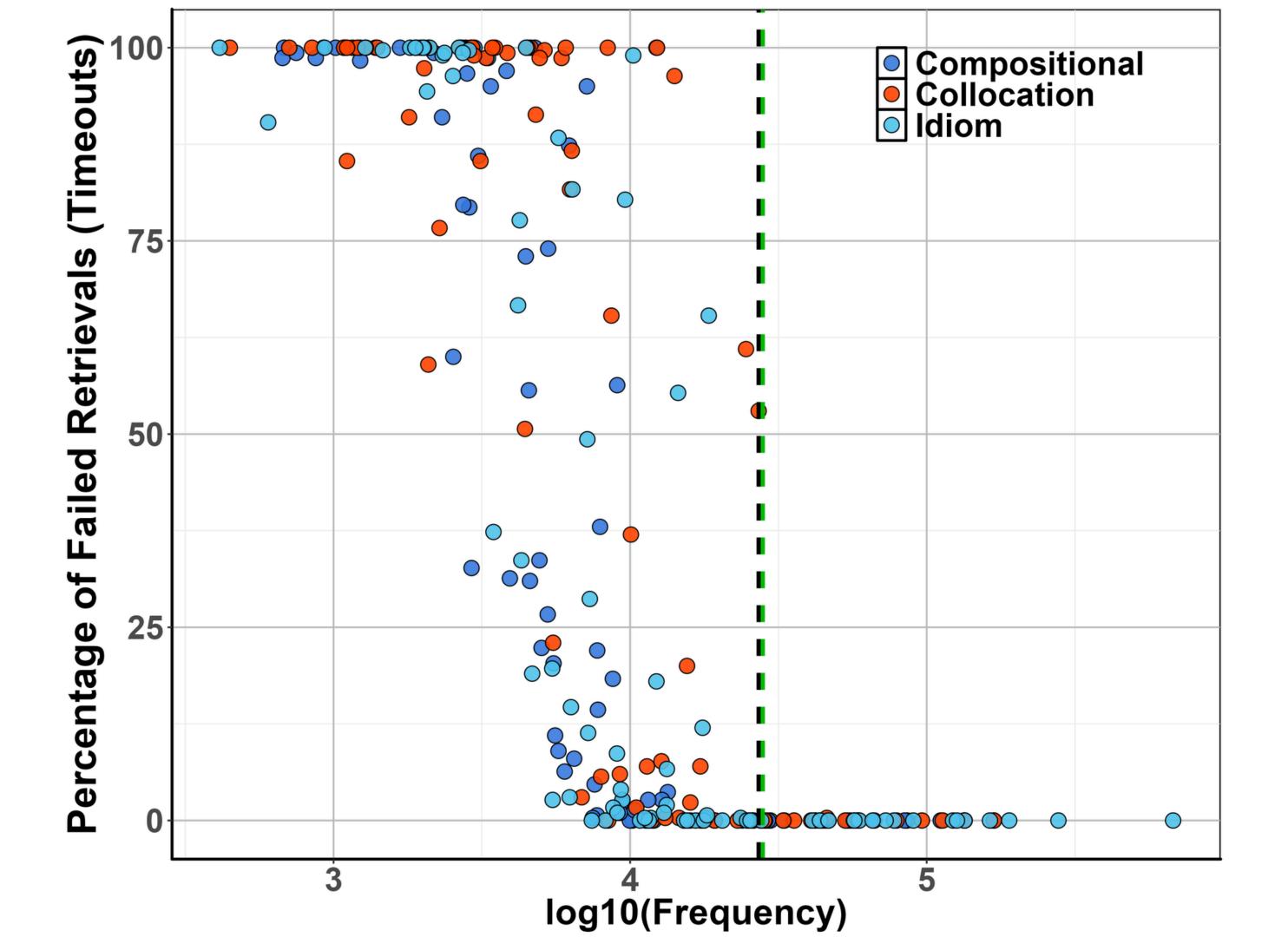
MODELLING RESULTS



Successful Retrievals







SUMMARY

1. **Human results** show:

- a. Collocations are processed the slowest
- b. Marginal difference between idioms and productive items

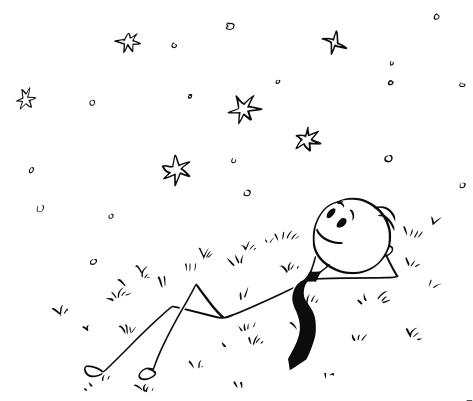
2. In **MINERVA2:**

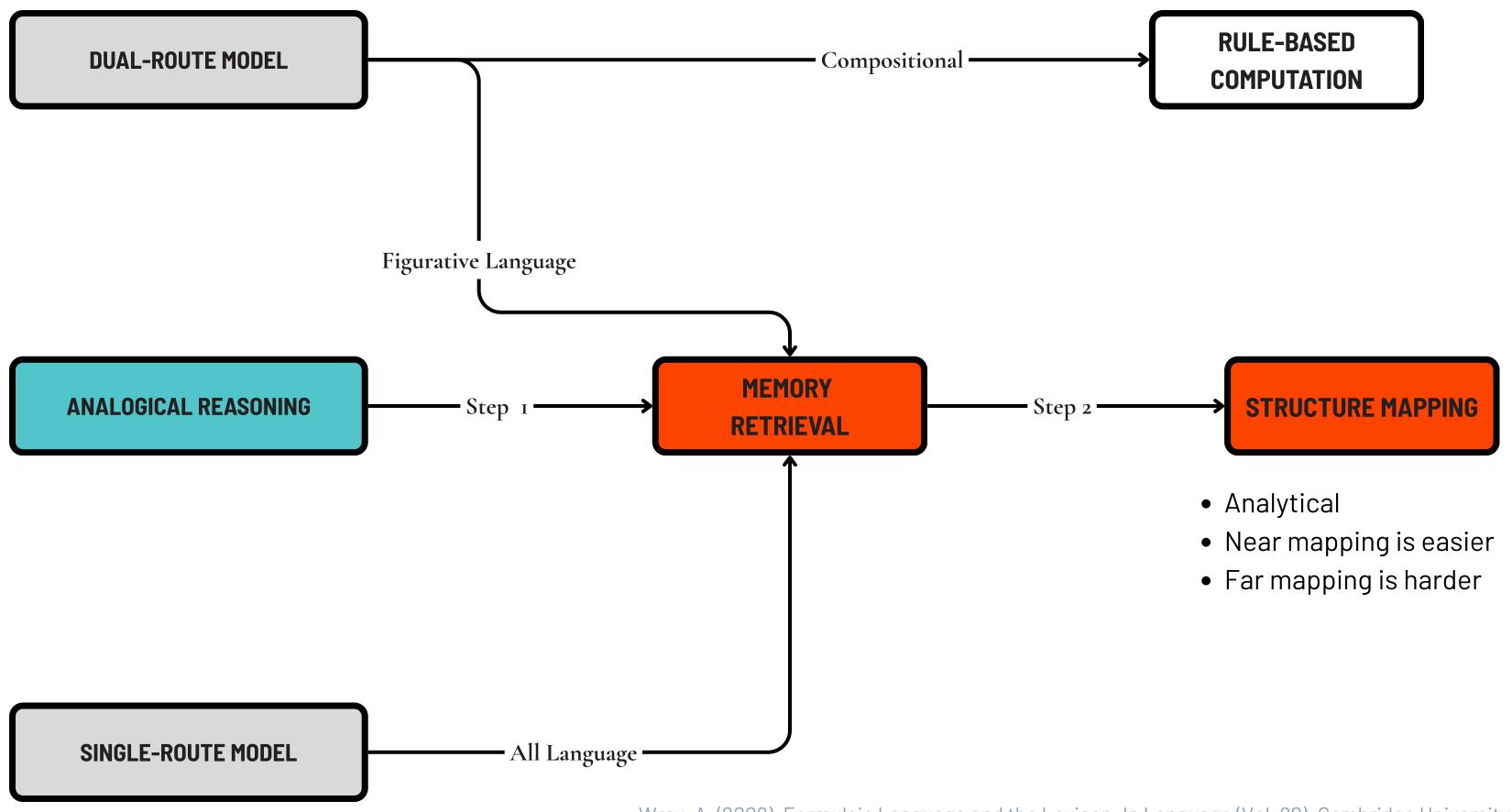
- a. frequency threshold similar to humans
- b. above threshold, MINERVA mirrors human results
- c. below threshold, failures to retrieve mirror human results

CONCLUSION

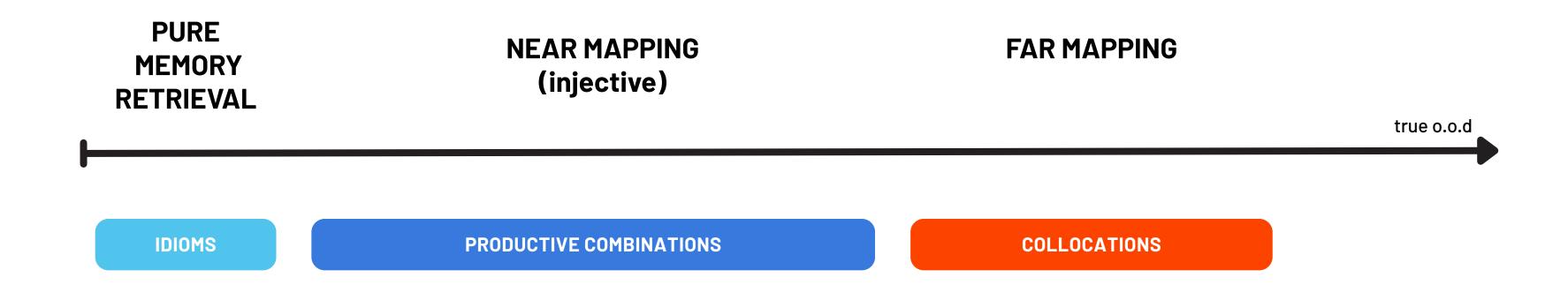
Memory retrieval is insufficient to capture human processing trends observed across the semantic compositionality continuum.

SO, WHAT ELSE IS THERE?



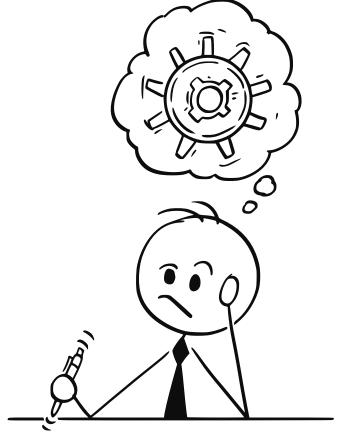


Wray, A. (2002). Formulaic Language and the Lexicon. In Language (Vol. 80). Cambridge University Press.⁷³
https://doi.org/10.1353/lan.2004.0209



This could account for the observed processing costs...

CAN WE MODEL THIS?



Back to #1: Every modeling choice is a theoretical commitment...

Slide Graveyard



LIMITATIONS

Task (for humans) and embeddings (for MINERVA) do not fully capture figurative vs compositional meanings for idioms:

- Can we give humans a memory task which elicits figurative readings of idioms?
- Can we get embeddings that better capture figurativeness?

