# Platonic Forms in the Study of Language and Mind

Elliot Murphy

Department of Neurosurgery
McGovern Medical School
University of Texas Health Science
Center

'*Symposium on the Platonic Space*'

*Primavera,* Botticelli (c. 1480)

Diverse causal landscapes in the brain derive distinct components of language

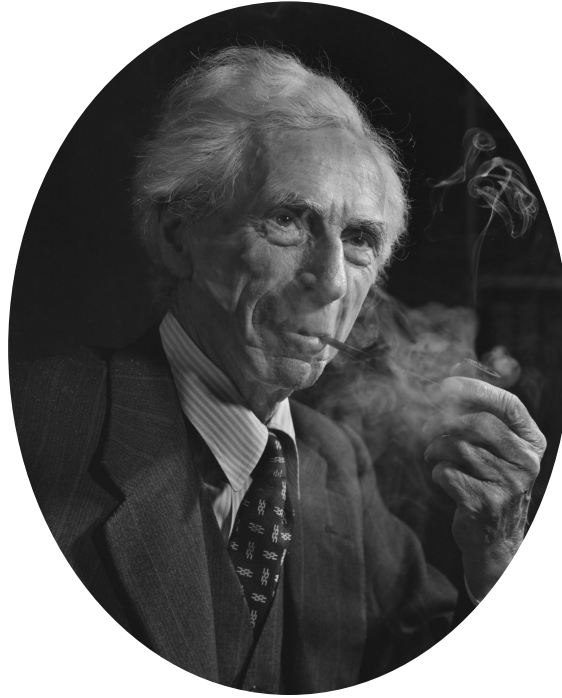Mathematically formalized theories of language can help to narrow down the list of plausible candidate neural mechanisms

This mathematical space will involve the use of concepts from category theory, Hopf algebra and statistical physics

With respect to previous speakers in this 'Symposium on the Platonic Space', some of this material will be building off the work of Lauren Ross and Karl Friston and the ideas they presented in their lectures.

John von Neumann

The mind follows logico-syntactic, **algebraic rules**



Bertrand Russell

It has fixed **scopes and limits**



Alan Turing

It can execute specific **computational operations**

Otto Jespersen

Humans abstract away from sensation to generate **"a notion of structure"**

Noam Chomsky

Language is a **generative system** of structured expressions

Paul Pietroski

Language provides **instructions to** conceptual systems

**Philip W. Anderson**

Rhythmicity provides a means of "**handling information**"

**Karl Friston**

The brain is an **inference generator,** negotiating its internal model with sensation

**Ada Lovelace**

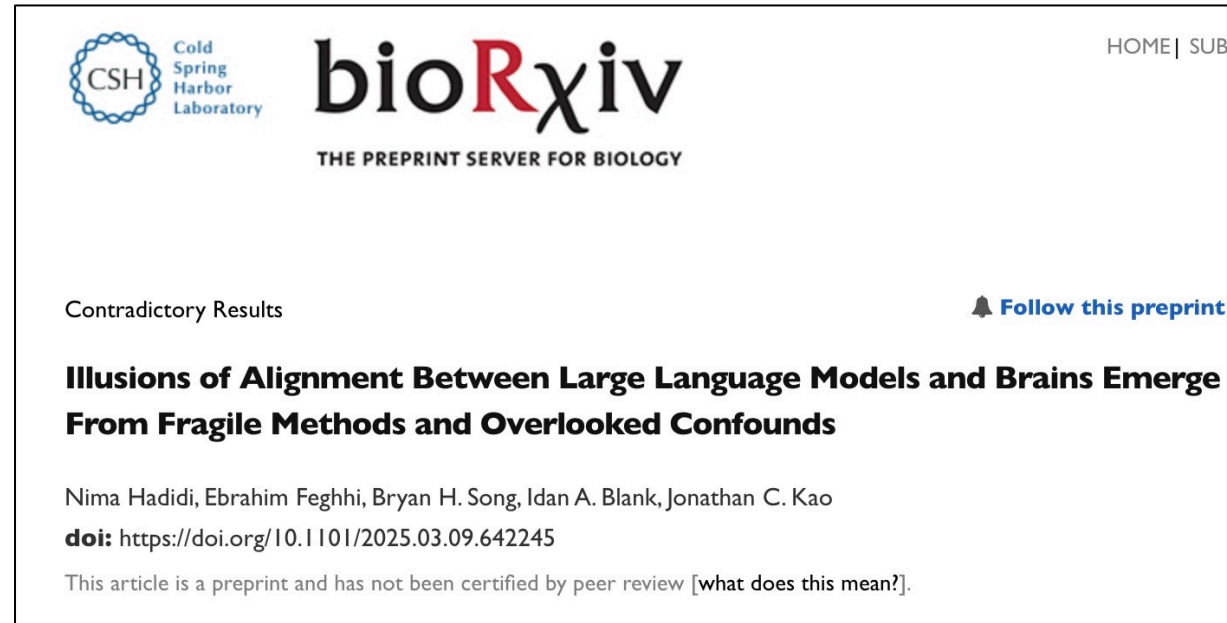We seek a "**uniting link**" between "the operations of matter ... and abstract mental processes".

# Large "Language" Models... or Large Corpus Models?

**(1)** LLM-brain alignments "driven by fragile methodologies and overlooked confounds" (Hadidi et al. 2025)

**Overlooked confounds**: Word rate, positional signals

# Large "Language" Models... or Large Corpus Models?

**(1)** LLM-brain alignments "driven by fragile methodologies and overlooked confounds" (Hadidi et al. 2025)

**(2)** Transformers do not represent sentence meaning in a manner akin to the human brain (Fodor et al. 2025)

> When word order matters: human brains represent sentence meaning differently from large language models
>
> James Fodor,* Carsten Murawski,* Shinsuke Suzuki[†]
>
> July 2, 2025

7T fMRI; 30 participants reading 108 sentences

Transformers were significantly inferior to models explicitly designed to encode the syntactic relations between words

# Large "Language" Models... or Large Corpus Models?

**(1)** LLM-brain alignments "driven by fragile methodologies and overlooked confounds" (Hadidi et al. 2025)

**(2)** Transformers do not represent sentence meaning in a manner akin to the human brain (Fodor et al. 2025)

**(3)** LLMs do not reliably distinguish between grammatical and ungrammatical structures (Dentella et al. 2025)

---

PLOS One

RESEARCH ARTICLE

## Language in vivo vs. in silico: Size matters but Larger Language Models still do not comprehend language on a par with humans due to impenetrable semantic reference

Vittoria Dentella[1]*, Fritz Günther[2], Evelina Leivada[3,4]

---

tion, namely on grammatical sentences. Additionally, ChatGPT-4 wavers more than humans in its answers (12.5% vs. 9.6% likelihood of an oscillating answer, respectively). Thus, while increased model size may lead to better performance, LLMs are still not sensitive to (un)grammaticality the same way as humans are. It seems possible but unlikely that scaling alone can fix this issue. We interpret these results by comparing language learning in vivo and in silico, identifying three critical differences concerning (i) the type of evidence, (ii) the poverty of the stimulus, and (iii) the occurrence of semantic hallucinations due to impenetrable linguistic reference.

# Large "Language" Models... or Large Corpus Models?

**(1)** LLM-brain alignments "driven by fragile methodologies and overlooked confounds" (Hadidi et al. 2025)

**(2)** Transformers do not represent sentence meaning in a manner akin to the human brain (Fodor et al. 2025)

**(3)** LLMs do not reliably distinguish between grammatical and ungrammatical structures (Dentella et al. 2025)

**(4)** Large 'reasoning' models (ChatGPT's o3, etc.) engage in "sophisticated pattern matching" (Shojaee et al. 2025)

 Machine Learning Research

Paper | June 2025

# The Illusion of Thinking: Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity

Parshin Shojaee*[†], Iman Mirzadeh*, Keivan Alizadeh, Maxwell Horton, Samy Bengio, Mehrdad Farajtabar

A recent assessment from Apple found no evidence of formal reasoning in LLMs

# Large "Language" Models... or Large Corpus Models?

**(1)** LLM-brain alignments "driven by fragile methodologies and overlooked confounds" (Hadidi et al. 2025)

**(2)** Transformers do not represent sentence meaning in a manner akin to the human brain (Fodor et al. 2025)

**(3)** LLMs do not reliably distinguish between grammatical and ungrammatical structures (Dentella et al. 2025)

**(4)** Large 'reasoning' models (ChatGPT's o3, etc.) engage in "sophisticated pattern matching" (Shojaee et al. 2025)

**(5)** Seven leading LLMs do not show reliable sensitivity to meaning/structure distinction (Dentella et al. 2024)

**scientific** reports

## Testing AI on language comprehension tasks reveals insensitivity to underlying meaning

Vittoria Dentella[1,2], Fritz Günther[3], Elliot Murphy[4], Gary Marcus[5] & Evelina Leivada[6,7]

# Large "Language" Models... or Large Corpus Models?

**(1)** LLM-brain alignments "driven by fragile methodologies and overlooked confounds" (Hadidi et al. 2025)

**(2)** Transformers do not represent sentence meaning in a manner akin to the human brain (Fodor et al. 2025)

**(3)** LLMs do not reliably distinguish between grammatical and ungrammatical structures (Dentella et al. 2025)

**(4)** Large 'reasoning' models (ChatGPT's o3, etc.) engage in "sophisticated pattern matching" (Shojaee et al. 2025)

**(5)** Seven leading LLMs do not show reliable sensitivity to meaning/structure distinction (Dentella et al. 2024)

**(6)** ChatGPT's o3 does not capture fundamental principles of linguistic structure (Murphy et al. 2025)

## Fundamental Principles of Linguistic Structure are Not Represented by o3

Elliot Murphy[1,2*], Evelina Leivada[3,4], Vittoria Dentella[5], Fritz Günther[6], Gary Marcus[7]

1. Vivian L. Smith Department of Neurosurgery, UTHealth, Texas, USA
2. Texas Institute for Restorative Neurotechnologies, UTHealth, Texas, USA
3. Universitat Autònoma de Barcelona, Barcelona, Spain
4. Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain
5. University of Pavia, Pavia, Italy
6. Humboldt-Universität zu Berlin, Berlin, Germany
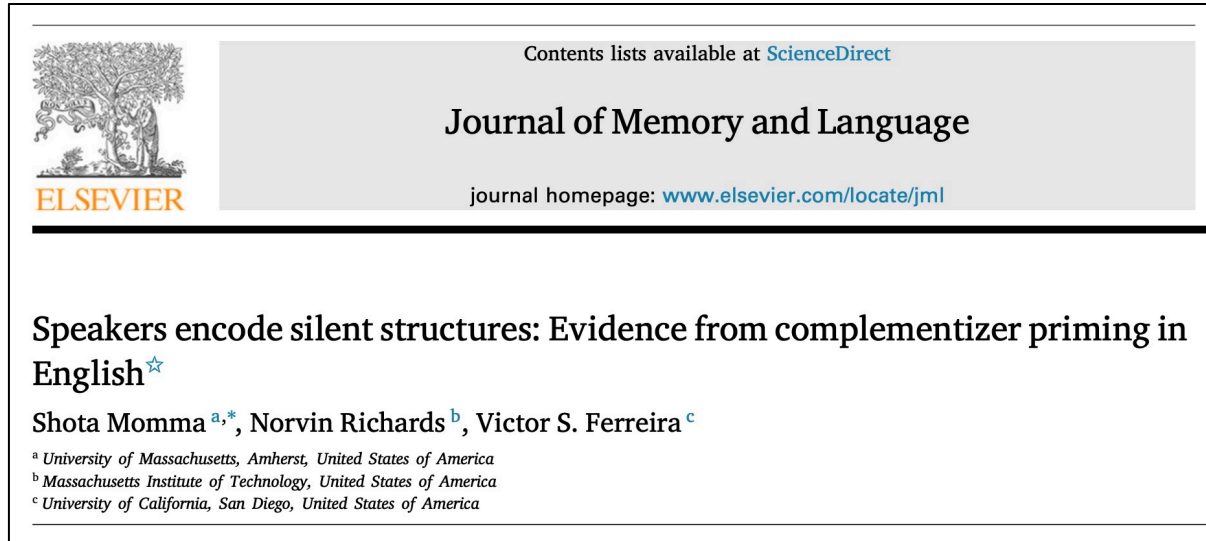7. New York University, New York, USA
**Corresponding author**: elliot.murphy@uth.tmc.edu

- Fails to generalize basic phrase structure rules
- Fails to distinguish between instructions to generate unacceptable semantic vs. unacceptable syntactic outputs

# Large "Language" Models... or Large Corpus Models?

**(1)** LLM-brain alignments "driven by fragile methodologies and overlooked confounds" (Hadidi et al. 2025)

**(2)** Transformers do not represent sentence meaning in a manner akin to the human brain (Fodor et al. 2025)

**(3)** LLMs do not reliably distinguish between grammatical and ungrammatical structures (Dentella et al. 2025)

**(4)** Large 'reasoning' models (ChatGPT's o3, etc.) engage in "sophisticated pattern matching" (Shojaee et al. 2025)

**(5)** Seven leading LLMs do not show reliable sensitivity to meaning/structure distinction (Dentella et al. 2024)

**(6)** ChatGPT's o3 does not capture fundamental principles of linguistic structure (Murphy et al. 2025)

**(7)** LLMs do not capture generalizations about null complementizers (Momma et al. 2025)

Speakers encode silent structures: Evidence from complementizer priming in English☆

Shota Momma [a,*], Norvin Richards [b], Victor S. Ferreira [c]

[a] University of Massachusetts, Amherst, United States of America
[b] Massachusetts Institute of Technology, United States of America
[c] University of California, San Diego, United States of America

Speakers syntactically encode zero complementizers as cognitively active mental objects

The cat [that] the boy loves
Jason whispered [that] the phoenix had escaped
?Jason whispered the phoenix had escaped

No evidence that LLMs capture cross-constructional generalizations about null complementizers

# Large "Language" Models... or Large Corpus Models?

**(1)** LLM-brain alignments "driven by fragile methodologies and overlooked confounds" (Hadidi et al. 2025)

**(2)** Transformers do not represent sentence meaning in a manner akin to the human brain (Fodor et al. 2025)

**(3)** LLMs do not reliably distinguish between grammatical and ungrammatical structures (Dentella et al. 2025)

**(4)** Large 'reasoning' models (ChatGPT's o3, etc.) engage in "sophisticated pattern matching" (Shojaee et al. 2025)

**(5)** Seven leading LLMs do not show reliable sensitivity to meaning/structure distinction (Dentella et al. 2024)

**(6)** ChatGPT's o3 does not capture fundamental principles of linguistic structure (Murphy et al. 2025)

**(7)** LLMs do not capture generalizations about null complementizers (Momma et al. 2025)

**(8)** LLMs have a strong linearity bias and do not capture deep syntactic structures (Diego-Simón et al. 2025)

## Probing Syntax in Large Language Models: Successes and Remaining Challenges

**Pablo J. Diego-Simón**
LSCP, ENS, PSL
Paris, France
pablo.diego-simon@psl.eu

**Emmanuel Chemla**
LSCP, ENS, PSL
Paris, France
emmanuel.chemla@ens.psl.eu

**Jean-Rémi King**
Meta AI
Paris, France
jeanremi@meta.com

**Yair Lakretz**
LSCP, ENS, PSL
Paris, France
yair.lakretz@gmail.com

"Remaining challenges" are effectively all of human-specific syntactic knowledge

😲

But surely Large "Language" Models would align best with *language areas* in the brain?

"Our results suggest that LLMs' ability to predict brain activation does not strongly differ between language and non-language-related brain areas".

(Gürel et al. 2025)

## On whether the relationship between large language models and brain activity is language-specific

Sertug Gürel (sertug.guerel@uni-potsdam.de)
Alessandro Lopopolo (lopopolo@uni-potsdam.de)
Milena Rabovsky (milena.rabovsky@uni-potsdam.de)
Department of Psychology, University of Potsdam, Karl-Liebknecht-Str. 24–25, 14476 Potsdam, Germany

**Computer Science > Computation and Language**

[Submitted on 22 Aug 2025]

## ChatGPT-generated texts show authorship traits that identify them as non-human

Vittoria Dentella, Weihang Huang, Silvia Angela Mansi, Jack Grieve, Evelina Leivada

Large Language Models can emulate different writing styles, ranging from composing poetry that appears indistinguishable from that of famous poets to using slang that can convince people that they are chatting with a human online. While differences in style may not always be visible to the untrained eye, we can generally distinguish the writing of different people, like a linguistic fingerprint. This work examines whether a language model can also be linked to a specific fingerprint. Through stylometric and multidimensional register analyses, we compare human-authored and model-authored texts from different registers. We find that the model can successfully adapt its style depending on whether it is prompted to produce a Wikipedia entry vs. a college essay, but not in a way that makes it indistinguishable from humans. Concretely, the model shows more limited variation when producing outputs in different registers. Our results suggest that the model prefers nouns to verbs, thus showing a distinct linguistic backbone from humans, who tend to anchor language in the highly grammaticalized dimensions of tense, aspect, and mood. It is possible that the more complex domains of grammar reflect a mode of thought unique to humans, thus acting as a litmus test for Artificial Intelligence.

**LLMs prefer to manipulate noun-based information, not verb-related information.**

But verbs encode the core properties of compositional syntax-semantics, setting up things like tense, voice, transitivity, mood, theta-roles, argument structure and aspect. Nouns are easier targets for distributional models due to frequency...

Regular Article

# DALL·E 2 fails to reliably capture common syntactic processes

Evelina Leivada [a,b,*], Elliot Murphy [c], Gary Marcus [d]

[a] Universitat Autònoma de Barcelona, Spain
[b] Institució Catalana de Recerca i Estudis Avançats (ICREA), Spain
[c] Vivian L. Smith Department of Neurosurgery, University of Texas Health Science Center at Houston, TX, USA
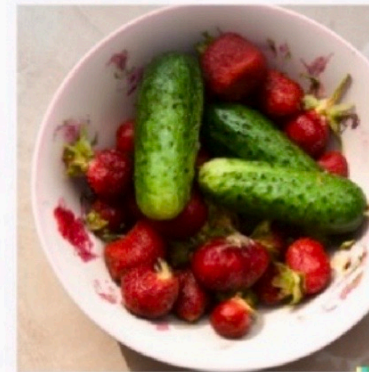[d] New York University, New York, USA

**DALL·E 2**

**ELSEVIER**

Regular Article

# A comparative investigation of compositional syntax and semantics in DALL·E and young children

Elliot Murphy [a,b,*] , Jill de Villiers [c], Sofia Lucero Morales [c]

[a] *Vivian L. Smith Department of Neurosurgery, University of Texas Health Science Center at Houston, TX, USA*
[b] *Texas Institute for Restorative Neurotechnologies, University of Texas Health Science Center at Houston, TX, USA*
[c] *Department of Psychology, Smith College, MA, USA*

**Reversible transitives**
"The girl sprays the boy"

**Reversible prepositions**
"The ball is behind the dog"

**Negation**
"One woman has glasses and one woman has no glasses"

**Advanced PPs**
"The umbrella is below the swing"

**PPs with adjectives**
"The girl is behind a car in a white garage"

**Passives**
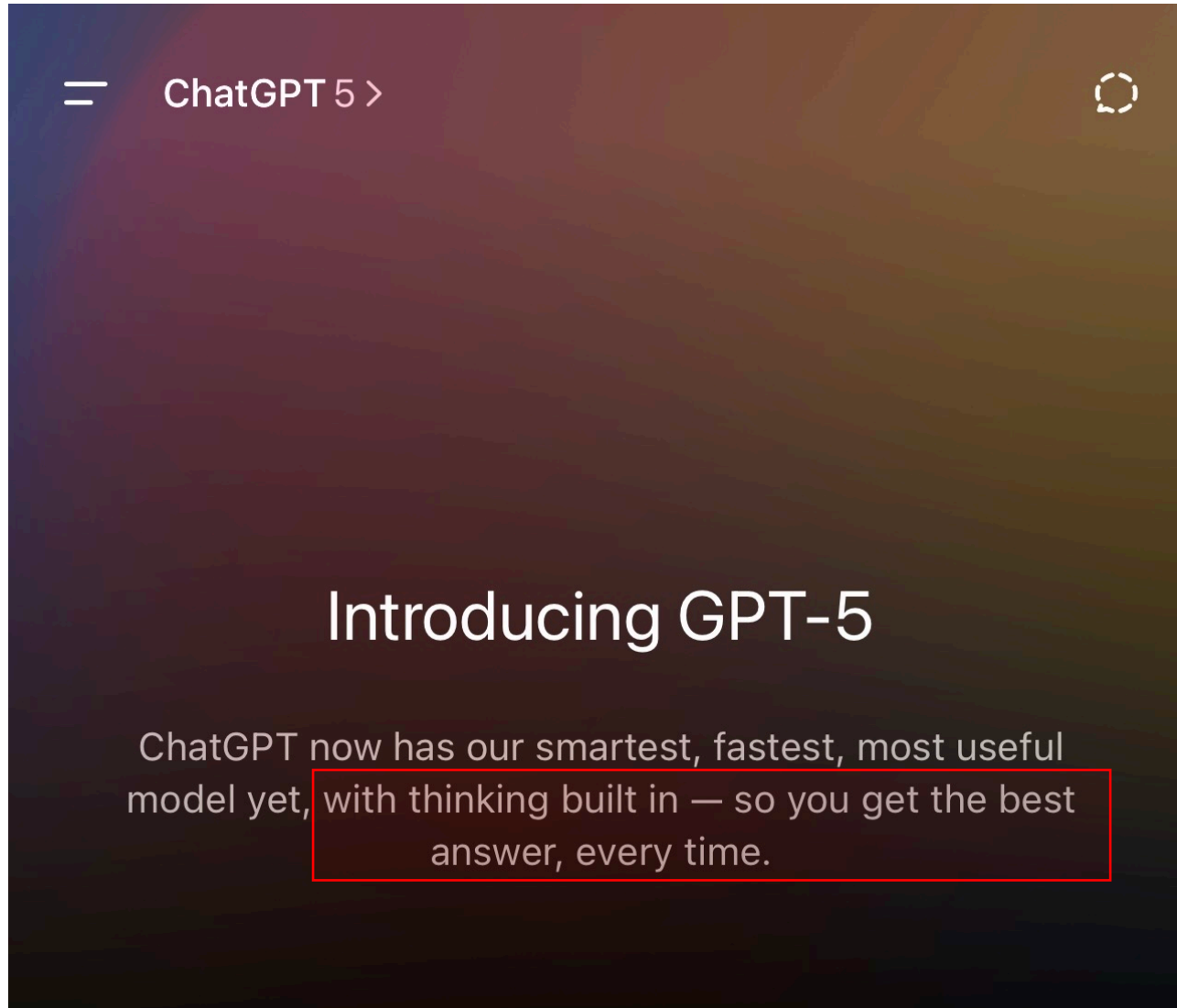"The cat is being dressed"

DALL·E 3

In the early days of generative text-to-image, I generated this 'Craiyon' prompt (originally **DALL·E mini**) in 2022:

**"A wedding photo with absolutely no Shrek"**

# But what about GPT-5?



ChatGPT 5 >

## Introducing GPT-5

ChatGPT now has our smartest, fastest, most useful model yet, with thinking built in — so you get the best answer, every time.

Only 2 unlit 🙊🙊🙊

ChatGPT 5 >

Draw a picture of a street at night with every street lamp lit except three. Draw in cartoon style
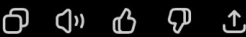
Image created

*"You'll see me one more time if you do good. You'll see me two more times if you do bad."*

Oh, that's a pretty memorable line! That's from the movie "Twin Peaks: Fire Walk with Me," and it's got that very mysterious, kind of eerie vibe that David Lynch does so well. Definitely one of those lines that sticks with you!

*"Not correct."*

Ah, got it! My mistake there. That line is actually from "Twin Peaks," but it's from the character known as The Man from Another Place in the original series rather than from "Fire Walk with Me." It's one of those iconic, cryptic lines that definitely adds to the mystery of the show. Thanks for the correction!

*"You're still wrong about that. It's not from Twin Peaks."*

Oh, my apologies for the mix-up again! You're right to keep me on my toes. That line is actually from "The Grand Budapest Hotel," spoken by the character of Mr. Moustafa as he's recounting advice he got from M. Gustave. It's definitely a memorable line with that sort of moral, whimsical touch Wes Anderson does so well. Thanks for pointing that out!
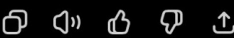
---

*"That's also wrong."*

Ah, third time's the charm then! I appreciate your patience with me here. That line is actually from "No Country for Old Men." It's spoken by the character Anton Chigurh. So that's the real source of the quote. Thanks for sticking with me until we got it right!
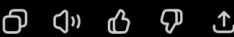
*"It's from a Lynch movie."*

Oh, got it! Well, then let's zero back in on Lynch. In that case, the line is actually from "Twin Peaks: The Return," which is part of the Twin Peaks universe but from the 2017 revival series. It's definitely a very Lynchian, mysterious sort of quote that fits right in with his style. Thanks for steering us back on track there!

*"No, wrong again."*

Ah, I appreciate you bearing with me here. Let's sort this out properly. That line is indeed from a David Lynch work, and it comes from "Inland Empire." So that's where you've heard it. It's a very Lynchian line for sure. Thanks for your patience as we pinned it down!

*"Wrong again. Where exactly in Inland Empire is the line?"*

---

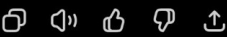*"Wrong again. Where exactly in Inland Empire is the line?"*

Ah, I see you're really making sure we get this one right, and I appreciate the persistence. I misspoke there—it's not actually from "Inland Empire," and I want to give you the correct info.

That quote is indeed a David Lynch line, but it's actually from "Twin Peaks: The Return." I know I initially went down a different path there, but this line is spoken by a character in that series. Thanks for sticking with me while we pinned it down for sure!

*"Which character sets the line and when?"*

Absolutely, let's get it spot on. That line is actually from "Twin Peaks: The Return," and it's spoken by the character known as The Fireman (previously known as The Giant in the original series). So you were absolutely right that it's Lynch, and it's from that particular character in "The Return." Thanks for making sure we got the exact details pinned down!

*"Wrong again, it's from Mulholland Drive."*

Ah, "Mulholland Drive" it is then! Thank you for that final nudge. That line is indeed from "Mulholland Drive," which absolutely fits the Lynchian mystery vibe we've been circling around. I appreciate your patience in getting us right to the correct source. Thanks for sticking with me!

To quote *Anchorman*...

March 2025 (GPT-4)

"Carefully" can modify either "fixed" or "packed" – regardless of intonation, 'context', etc...



> The mechanic who fixed the car carefully packed his tools.
> What did he do carefully?

Reasoned about what was done carefully for a few seconds

He carefully packed his tools.

> Did he carefully fix the car?

Reasoned about car repair actions for a few seconds >

The sentence only specifies that he "carefully packed his tools." It doesn't say anything about how he fixed the car. So, based on that sentence alone, we don't know whether he fixed the car carefully.

March 2025 (GPT-4)

"Carefully" can modify either "fixed" or "packed" – regardless of intonation, 'context', etc…



**Garry Kasparov** ✓
@Kasparov63

Ambiguity in language often requires context and understanding likelihood and custom. In a lecture 10 years ago, I used the sentence, "The chicken is too hot to eat" as an example. Is it a concern for a farmer or someone eating dinner? It depends on the context of the dialogue. t.co/Wi3DDUXLeR

**BabyLM Challenge**

Sample-efficient pretraining on a developmentally plausible corpus

A **very quick** lesson from the ongoing BabyLM Challenge

When limited to ecologically valid amounts of data (i.e., the amount that children are exposed to), LMs have an unshakably strong bias to impose linear solutions on linguistic tasks.

In contrast, infants instinctively impose hierarchical structure-dependent rules.

**Valid form (Modus Ponens)**
If $P$, then $Q$
$P$
Therefore, $Q$ ✅

**Valid form (Modus Tollens)**
If $P$, then $Q$
Not $Q$
Therefore, not $P$ ✅

> If the model correlates with human behavioural and/or neuroimaging data, then the model does what humans do. ($P \rightarrow Q$)The model correlates with human behavioural and/or neuroimaging data. ($P$)Therefore, the model does what humans do. ($\Vdash Q$)

Guest & Martin (2023)

**_Affirming the consequent!_**
Correlation of our models to data is necessary but not sufficient.

Common in the literature to:
(1) confuse types of inference
(2) misunderstand the evidentiary role that correlation plays
(3) offer no formalized thought on the relationship between model and observation

> If the model correlates with human behavioural and/or neuroimaging data, then the model does what humans do. ($P \rightarrow Q$)The model correlates with human behavioural and/or neuroimaging data. ($P$)Therefore, the model does what humans do. ($\Vdash Q$)

Guest & Martin (2023)

**Statement 1:** If I am Beyoncé, then I am fabulous

**Statement 2:** I am fabulous

**Conclusion:** Therefore, I am Beyoncé 🎤✨

# Models ≠ Phenomenon
# Explanandum ≠ Explanans

***Example case***: Motion of bodies under gravity and Newtonian mechanics

If Newtonian mechanics behaves like physical objects, then Newtonian mechanics is physical objects

?

# A Task-Optimized Neural Network Replicates Human Auditory Behavior, Predicts Brain Responses, and Reveals a Cortical Processing Hierarchy

Alexander J.E. Kell,[1,2,6,7,*] Daniel L.K. Yamins,[3,4,6] Erica N. Shook,[1,2] Sam V. Norman-Haignere,[1] and Josh H. McDermott[1,2,5,*]

[1]Department of Brain and Cognitive Science, MIT, Cambridge, MA, USA
[2]Center for Brains, Minds, and Machines, MIT, Cambridge, MA, USA
[3]Departments of Psychology and Computer Science, Stanford University, Stanford, CA, USA
[4]Stanford Neurosciences Institute, Stanford, CA, USA
[5]Program in Speech and Hearing Biosciences and Technology, Harvard University, Cambridge, MA, USA
[6]These authors contributed equally
[7]Lead Contact
*Correspondence: alexkell@mit.edu (A.J.E.K.), jhm@mit.edu (J.H.M.)
https://doi.org/10.1016/j.neuron.2018.03.044

"…intermediate model layers best **explain** primary auditory cortical responses, while deeper layers best **explain** voxels in non-primary areas." (Kell et al. 2018, *Neuron*)

?

Correlation ≠ Explanation

"This LLM output has human-like qualities to it, therefore it gives us certain specific implications for theories of human cognition..."

But nobody argues that AlphaGo is a plausible model of human strategizing, or that Cicero is a plausible model of human theory of mind. Was Deep Blue a plausible model of Garry Kasparov's brain simply because it matched his performance levels?

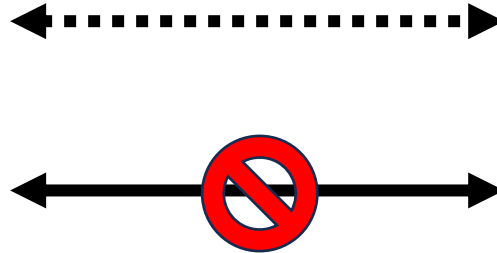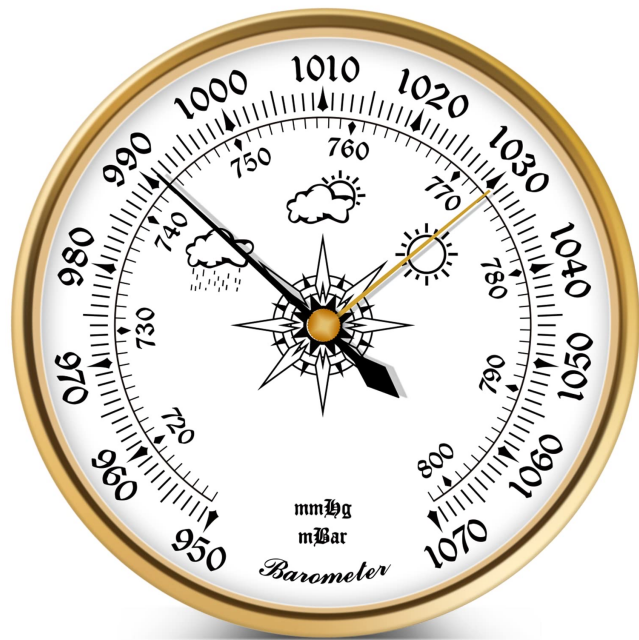So why would an LLM be a model of human language?

**Problem**: Standards of evaluation

Sometimes adopting a different testing regime for models vs. humans can be used as a tool to artificially inflate LLM accuracy.

Some cases of s(timulus)-hacking: In s-hacking, people try many different prompts and testing regimes, cherry-pick the best performing ones, and rerun the tests in LLMs until a desired result is achieved. This is analogous to p-hacking, but today many people working on AI capabilities call it "pre-testing"...

**Poor causal logic**

A barometer can measure atmospheric pressure, and have its activity regulated by it. But if we break the barometer, the thunderstorm doesn't stop.

**LLM text output looks like real human language…**

**…and a fake flower looks a lot like a real flower**

## *Multiple realizability*

Vastly different substrates and mechanisms can perform the same input-output mappings.

Digital clocks and analog clocks both **show the time** but have completely distinct implementations – and LLMs and humans can both **generate coherent prose**...
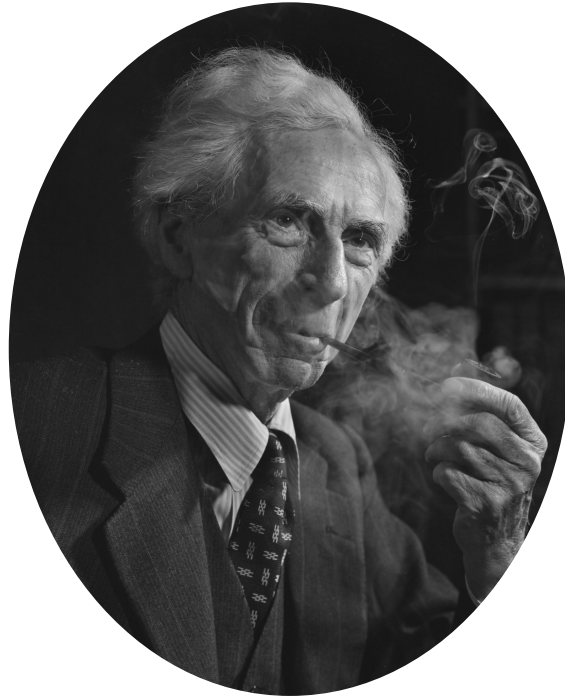
# Related criticism to the LLM-brain debate…

- **Fitting RNNs make attractors even when the data does not have attractors**

- "A popular approach to study brain data is to fit an RNN to neural data and then show that this RNN has attractors. However, a recent paper [Qian et al. 2024, NeurIPS] showed that in the context of partial observations, a simulated system that has no attractors gives rise to a fitted RNN that does have attractors. In other words, in the context of unobserved data (always the case in neuroscience), we may even wrongly infer that there are attractors." (Konrad Kording)

In the cognitive neuroscience literature, there are many "curious shadowy" syllogisms and statements (Russell 1918, *The Philosophy of Logical Atomism*) that do not obtain, but which are disguised as complex, mechanistic insights.



Bertrand Russell

# ChatGPT

Do you think neuroscientists should use LLMs as genuine 'theories of language' to guide investigation?

# ChatGPT

Do you think neuroscientists should use LLMs as genuine 'theories of language' to guide investigation?

No, neuroscientists **should not** use large language models (LLMs) as *genuine theories of language* to guide investigation—at least not in the current state of LLM development and linguistic theory.

Here's a breakdown of why:

## 🧠 1. LLMs Are Not Theories — They're Function Approximators

LLMs like GPT are trained to model surface-level statistical patterns in language, not to formalize or explain the cognitive or neural mechanisms underlying language use. A genuine theory of language should:

- Specify **generative rules** or **computational mechanisms**.
- Explain **how** humans acquire and use language.
- Predict **neurological or behavioral** outcomes.

## 🧪 2. They Lack Cognitive Plausibility

A theory of language for neuroscience must be anchored in:

- **Working memory constraints**
- **Incremental parsing**
- **Compositional semantics**
- **Neural implementation feasibility**

## 🧬 3. They Obscure Rather Than Clarify Neural Mechanisms

Treating LLMs as theories can **mislead** neuroscientific interpretations:

- LLMs are **black boxes**: even with probing, we cannot meaningfully link their internal states to neural computations.
- Their success at some tasks tempts researchers into **post hoc rationalization** ("it must be like the brain!") rather than principled hypothesis testing.

## 📚 4. They Bypass Theoretical Linguistics

LLMs do not engage with core theoretical insights from:

- **Syntax** (e.g., Merge, movement, locality constraints)
- **Semantics** (compositionality, quantifier scope, etc.)
- **Phonology** (prosody, syllable structure)

They generate plausible strings but do not represent **formal grammatical competence**. Using them to guide brain research encourages **empiricism without theory**.

Modeling linguistic semantics via a "vector space" will only get us so far. Language interfaces with various conceptual domains with their own unique formats: some are digital/logical, others are analogue/graded/continuous.

Within a single word, we can call upon multiple conceptual representations (e.g., *book*, *city*, *proudly*)

Spatial structure

Causation

Telic/functional

Epistemic/evidential

Intensionality

Event structure

Theory of mind

Quantification

Boolean logic

Factivity

**Vector space semantics** (word embeddings, LLM hidden states) treat lexical meaning as points in a continuous, uniform space.

This is maybe OK for things like 'king' vs 'queen', but semantics in reality is highly *typed* and *heterogeneous*.

"The progressive **school** next to the river was high in the *NYT* rankings and had just hired a new chancellor after being repainted"

---

A single word (***school***) can call upon diverse and categorially incompatible semantic features. These conceptual domains are not "dimensions" in the embedding sense: they host different mathematical structures *and* implicate different neural substrates.

Word2Vec/LLM embeddings collapse these senses into a blended vector. Even contextual embeddings often smear them, because they lack an explicit type system to separate senses.

Distances in vector space often reflect corpus frequency artifacts.

**Asymmetric entailment relations**

(1) "Every dog barked"

(2) "Some dog barked"

Distributional vectors place "every" and "some" close together because they co-occur with the same kinds of nouns/verbs.

But (1) entails (2), not the reverse.

Review

# Why concepts are (probably) vectors

Steven T. Piantadosi[1,2,*], Dyana C.Y. Muller[2], Joshua S. Rule[1], Karthikeya Kaushik[1], Mark Gorenstein[2], Elena R. Leib[1], and Emily Sanford[1]

Vector approaches risk conflating the *implementation medium* with the *computational level*.
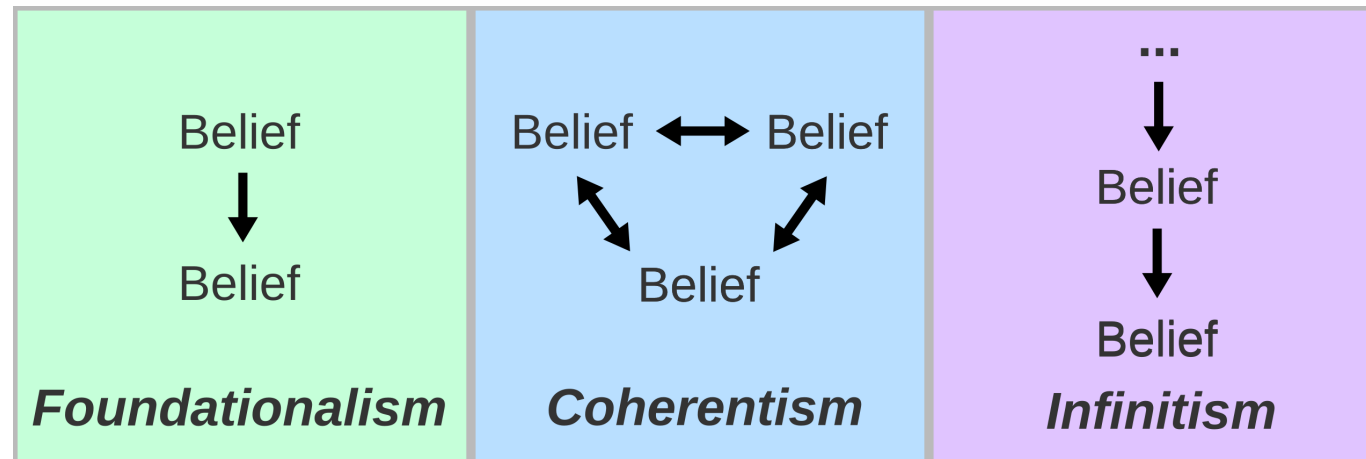
**How does a vector-based learner discover abstract, exceptionless rules** without relying on statistical accident?

Piantadosi et al. risk explaining compositionality post hoc ("it emerges in the geometry") rather than as **a necessary design property**.

The very properties that many authors cite as hallmarks of efficient communication can also be reinterpreted as pressures shaping **internal symbol manipulation for thought:**

Dependency minimization, ambiguity/polysemy, and short words can **ALL** help with both communicative and internal computational efficiency.

In philosophy of law, **coherentism** holds that a belief is justified if it fits into a larger, consistent system of beliefs, forming a mutually supportive network.



Claims about "LLMs are like the brain" can certainly be *internally coherent* – but this doesn't mean that our foundational assumptions are valid!

# The 'definitional fallacy'

What I will call the 'definitional fallacy' has run rampant in contemporary cognitive neuroscience.

Just because you can define a concept like 'culture' or 'communication' in relation to some scientific field, it does *not* follow that this definition will actually be able to be operationalized within the context of a productive explanatory theory.

I can provide a *definition* of semantics as 'vector math' (etc) – but it doesn't follow that just because you can define semantics as x, y or z that it *should* be conceptualized this way.

Often researchers think half the battle is already won just by offering coherent definitions and boundaries of inquiry.

# Connectionism is all you need?

GPT-2 was a pure LLM. But over the past few years, some leading LLMs have gradually become neurosymbolic by smuggling in numerous module interfaces (e.g., Python interpreters; possible symbolic filters in guardrails), not relying purely on deep learning.

**Marcus on AI**

## How o3 and Grok 4 Accidentally Vindicated Neurosymbolic AI

Neurosymbolic AI is quietly winning. Here's what that means – and why it took so long

GARY MARCUS
JUL 13, 2025

# In the study of semantics, we have moved from this...

**TWO ELEMENTARY COMPOSITIONAL OPERATIONS:** M-join & D-join

$$\mu([a\ cow_N]_N) = \text{M-join}(\mu(-pl), \mu(cow_N))$$
$$= \text{M-join}(fetch@-pl, fetch@cow_N)$$
$$\rightarrow \text{ONE}(\_) \wedge \sqrt{}\text{COW}(\_)$$

$$\mu([above_P\ [a\ cow_N]_N]_P)$$
$$= \text{D-join}(\mu(above), \mu([a\ cow_N]_N))$$
$$= \text{D-join}(fetch@above, \text{M-join}(fetch@-pl, fetch@cow))$$
$$\rightarrow \exists[\text{ABOVE}(\_, \_) \wedge [\text{ONE}(\_) \wedge \sqrt{}\text{COW}(\_)]]$$

OXFORD

CONTEXT & CONTENT

CONJOINING MEANINGS

*Semantics Without Truth Values*

PAUL M. PIETROSKI

Paul Pietroski

**In the study of semantics, we have moved from this... to this**

Neuron 1

man

woman

king

queen

Neuron 2

Neuron 3

My central scientific research interest concerns how the brain coordinates the binding of distinct features into what linguists describe – using slightly distinct mathematical formalisms – as a labeled set (Chomsky in the 1990s), a compositional structure (Pietroski in the 2000s), an intentional composite function (Hoshi in the 2010s), a hylomorphic pluralistic mereological object (Adger in 2020s), or a free non-associative, commutative magma (Marcolli et al. 2025).

Or, more simply, what Otto Jesperson called "a notion of structure" in the early 1900s.

**Can these formal accounts of language be of service to experimental neurolinguistics? Does the high-level psychological theory of language you subscribe to have consequences for cognitive neuroscience?**

Two lessons from the 'neural correlates of consciousness'…

(1) Correlating neural activity ≠ explanation
(2) Bad metaphysics gives you bad science

Mark Solms,
South African neuropsychologist

Studies of consciousness that prioritize vision miss the mark somewhat.

'Affect' is the most fundamental property of consciousness, not perceptual representations. At the most essential level, every conscious experience has a 'feeling' to it. **Consciousness is not reducible to reportable perception.**

Therefore, neuroscientific studies that are based on visual oddball paradigms (etc.) may not expose critical neural dynamics for consciousness.
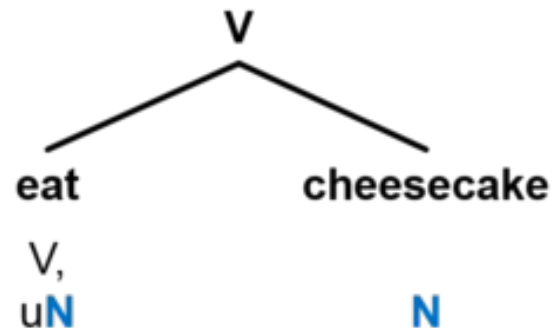
**Hierarchical constituency structure**

Workspace (WS) = [X, Y, Z]

MERGE(X, Y) → WS' = [{X, Y}, Z]

**Example:**

MERGE(eat, cheesecake) – [*asymmetric headedness*]

Algebraic models from linguistic theory provide hints at what the neural code for syntax *might* look like.

We may be looking for a neural mechanism that respects non-associativity of constituent geometry:

$$((W_1 W_2) W_3) \neq (W_1 (W_2 W_3))$$

[[old men] and women] ≠ [old [men and women]]

Syntactic knowledge boils down to a non-associative, commutative magma (category theory, Hopf algebra) that generates nonplanar trees, interfacing with distinct cognitive systems, providing instructions to them.

Chomsky argued for formulations of MERGE being couched within naïve set theory. But sets are too unstructured.

Recent category-theoretic magma formulations (Marcolli et al. 2025, *Mathematical Structure of Syntactic Merge*) provide new avenues to formally map syntax to possible neural geometries.

The algebraic properties of natural language seem to be unlike anything else in cognition.

How the brain neurally enforces a *free non-associative commutative unbounded digitally infinite combinatorial hierarchical recursive structure-building tree-formation category-theoretic magma operation* (syntactic structure-building) remains a mystery – but it is the central question in the cognitive neuroscience of language.

*Operations are subject to structural configurations, not linear sequential distance*

**The boy who is holding the flowers is happy**

**Is** the boy [who is holding the flowers] _ happy? ✅

**Is** the boy [who _ holding the flowers] is happy? ❌

At a minimum, the neurobiology of natural language syntax must comply with some simple conditions:
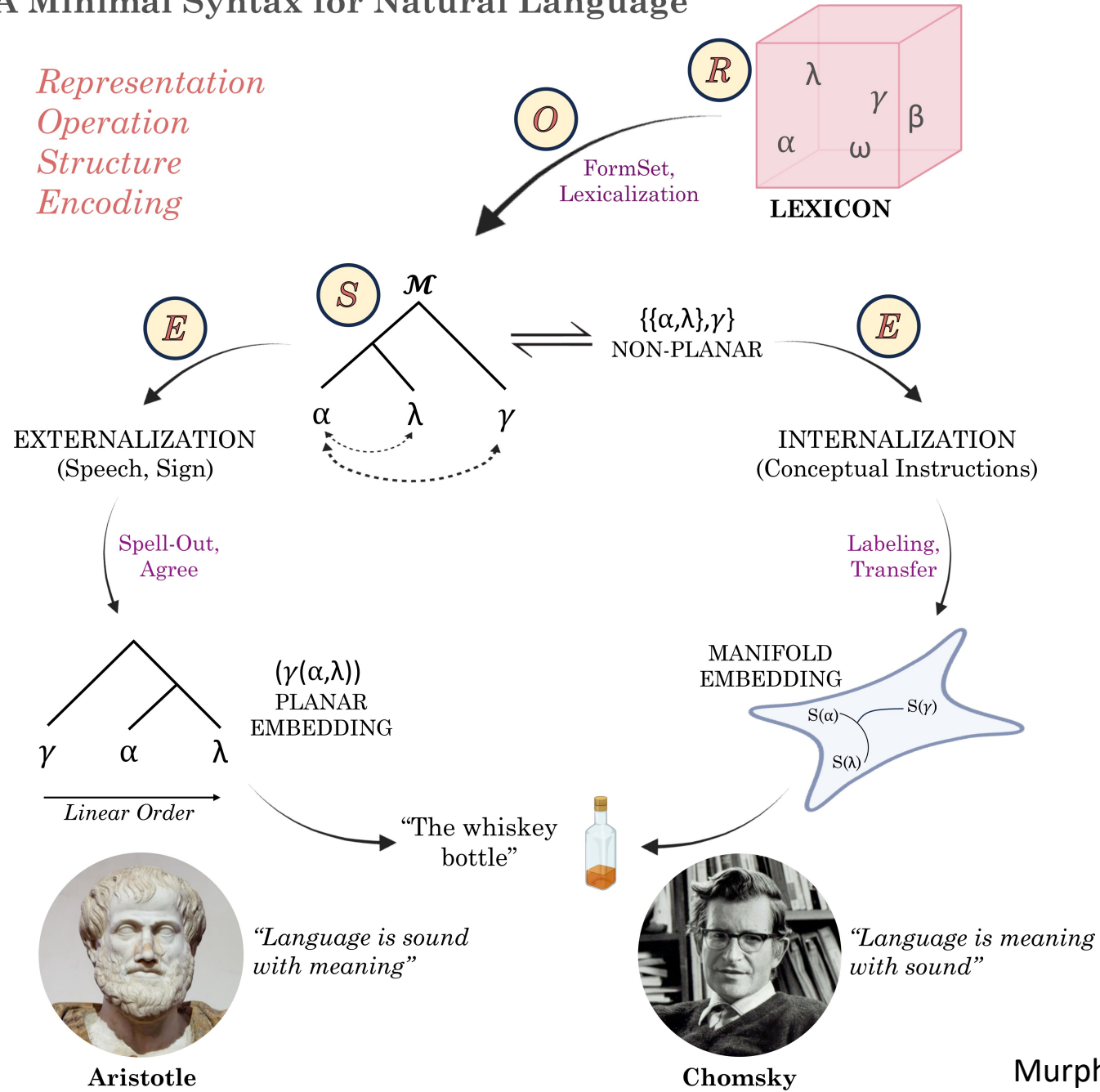
- **Commutativity** (MERGE(A,B)=MERGE(B,A)): The neural code for a two-item set must ignore linear order at the moment of combination.

- **Non-associativity** ((A∘B)∘C ≠ A∘(B∘C)): Once a third element is merged, the hierarchical depth of previous combinations must be recoverable, usually through a categorial 'label' or head-selection step that privileges one member of the newly formed set.

- **Closure**: The output of MERGE is itself a syntactic object.

- **Binarity**: MERGE generates strictly binary-branching structures.

- **Non-monotonic structure-building**: MERGE can involve deletion of sub-trees or workspace elements.

Candidate mechanisms for neural implementation therefore need order-insensitive pairwise binding plus a second, depth-sensitive process that is triggered only when the bound item is itself subjected to a further merge. Critically, experimental predictions for MERGE-based syntax must be specific enough to help adjudicate between salient support for and against structural inferences – without this precision, predictions will risk being ornamental and subject to confirmation bias.

Later, I will argue for a specific neural model of language ('ROSE') that satisfies these criteria.

# A Minimal Syntax for Natural Language

*Representation*
*Operation*
*Structure*
*Encoding*

R, O — FormSet, Lexicalization — LEXICON

S — $\mathcal{M}$ — $\{\{\alpha,\lambda\},\gamma\}$ NON-PLANAR

E — EXTERNALIZATION (Speech, Sign) — Spell-Out, Agree

E — INTERNALIZATION (Conceptual Instructions) — Labeling, Transfer

$(\gamma(\alpha,\lambda))$ PLANAR EMBEDDING

*Linear Order*

MANIFOLD EMBEDDING — $S(\alpha)$, $S(\gamma)$, $S(\lambda)$

"The whiskey bottle"

**Aristotle** — *"Language is sound with meaning"*

**Chomsky** — *"Language is meaning with sound"*

Murphy (2025), *Cognitive Neuroscience*

# All or nothing: No half-Merge and the evolution of syntax

Robert C. Berwick [1] *, Noam Chomsky [2]

1 Massachusetts Institute of Technology, Cambridge, Massachusetts, United States of America, 2 University of Arizona, Tucson, Arizona, United States of America

"...there is [no] justification, empirical or conceptual, for the decomposition of binary set formation into separate steps".

There is no mathematically plausible 'incremental', gradualist account of the emergence of Merge-based recursive syntax.

Can We Still Be Optimistic About the Future? | A Conversation with Steven Pinker

**The Poetry of Reality with Richard Dawkins**
196K subscribers

Join     Subscribe

👍 2.5K     👎     ➤ Share     ⬇ Download     💲 Thanks     •••

112K views  3 months ago   The Poetry of Reality with Richard Dawkins

Unlike domains of the language sciences like sociolinguistics, pragmatics, and experimental psycholinguistics (which deal with highly complex cognitive and social systems), theoretical linguistics can provide a streamlined view on what it is our brains must be able to do in order to comprehend and produce natural language.

In turn, formal syntactic models act as a sieve for neural theories.

# 4 A typology of systems based on the type(s) of merge they adopt

In previous work (Rizzi 2016b) I have proposed a typology of language-like systems based on the type(s) of merge that they adopt:

**0-merge** systems: they have a lexicon and an Item selection mechanism, but no merge or other combinatorial device. Linguistic expressions therefore consist of single lexical items.

**1-merge** systems: they have Item – Item merge, but no temporary repository. So, they can do Item – Item merge, generating expressions of at most two items, with no recursive application.

[[the] [boy]]

**2-merge** systems: they have a single temporary repository, which permits (recursive) applications of merge, thus permitting Item – Phrase merge, but no Phrase – Phrase merge (or at least no Phrase –Phrase external merge: I do not address internal merge in this note).

[He [saw [the boy]]]

**3-merge** systems: they involve two temporary repositories (or the option of parallel computations in a single workspace, in Chomsky's terms), which is consistent with Phrase – Phrase merge, the merger of two complex phrases, like the one that derives DP-vP (subject – predicate) structures.

[The man [saw [the boy]]] (complex specifiers and complex subjects)

**Some monkeys have a kind of primitive morphology (krak-oo vs. hok-oo, where -oo is some kind of pragmatic emphasizer), but are possibly restricted to 0-merge systems**

**word-word binding**

**word-phrase binding**

**phrase-phrase binding**

Rizzi (2024)

# 4 A typology of systems based on the type(s) of merge they adopt

In previous work (Rizzi 2016b) I have proposed a typology of language-like systems based on the type(s) of merge that they adopt:

**0-merge** systems: they have a lexicon and an Item selection mechanism, but no merge or other combinatorial device. Linguistic expressions therefore consist of single lexical items.

**1-merge** systems: they have Item – Item merge, but no temporary repository. So, they can do Item – Item merge, generating expressions of at most two items, with no recursive application.

**2-merge** systems: they have a single temporary repository, which permits (recursive) applications of merge, thus permitting Item – Phrase merge, but no Phrase – Phrase merge (or at least no Phrase –Phrase external merge: I do not address internal merge in this note).

**3-merge** systems: they involve two temporary repositories (or the option of parallel computations in a single workspace, in Chomsky's terms), which is consistent with Phrase – Phrase merge, the merger of two complex phrases, like the one that derives DP-vP (subject – predicate) structures.

**This level and above requires a workspace to store combined objects – this will be the E level of ROSE introduced later**

**word-word binding** [[the] [boy]]

**word-phrase binding** [He [saw [the boy]]]

**phrase-phrase binding** [The man [saw [the boy]]] (complex specifiers and complex subjects)

Rizzi (2024)

# 4 A typology of systems based on the type(s) of merge they adopt

In previous work (Rizzi 2016b) I have proposed a typology of language-like systems based on the type(s) of merge that they adopt:

**0-merge** systems: they have a lexicon and an Item selection mechanism, but no merge or other combinatorial device. Linguistic expressions therefore consist of single lexical items.

**1-merge** systems: they have Item – Item merge, but no temporary repository. So, they can do Item – Item merge, generating expressions of at most two items, with no recursive application.

[[the] [boy]]

**2-merge** systems: they have a single temporary repository, which permits (recursive) applications of merge, thus permitting Item – Phrase merge, but no Phrase – Phrase merge (or at least no Phrase –Phrase external merge: I do not address internal merge in this note).

[He [saw [the boy]]]

**3-merge** systems: they involve two temporary repositories (or the option of parallel computations in a single workspace, in Chomsky's terms), which is consistent with Phrase – Phrase merge, the merger of two complex phrases, like the one that derives DP-vP (subject – predicate) structures.

[The man [saw [the boy]]] (complex specifiers and complex subjects)

**The possibility of developing in parallel two complex phrases also demands something that word-phrase does not need – a workspace where parallel structure-formation is permitted (Chomsky), or two separate workspaces (Adger, Rizzi)**

**word-word binding**

**word-phrase binding**

**phrase-phrase binding**

Rizzi (2024)

# 4 A typology of systems based on the type(s) of merge they adopt

In previous work (Rizzi 2016b) I have proposed a typology of language-like systems based on the type(s) of merge that they adopt:

**0-merge** systems: they have a lexicon and an Item selection mechanism, but no merge or other combinatorial device. Linguistic expressions therefore consist of single lexical items.

**word-word binding**

**1-merge** systems: they have Item – Item merge, but no temporary repository. So, they can do Item – Item merge, generating expressions of at most two items, with no recursive application.

[[the] [boy]]

**word-phrase binding**

**2-merge** systems: they have a single temporary repository, which permits (recursive) applications of merge, thus permitting Item – Phrase merge, but no Phrase – Phrase merge (or at least no Phrase –Phrase external merge: I do not address internal merge in this note).

[He [saw [the boy]]]

**phrase-phrase binding**

**3-merge** systems: they involve two temporary repositories (or the option of parallel computations in a single workspace, in Chomsky's terms), which is consistent with Phrase – Phrase merge, the merger of two complex phrases, like the one that derives DP-vP (subject – predicate) structures.

[The man [saw [the boy]]] (complex specifiers and complex subjects)

Rizzi (2024)

Taining Zhang, Adam Goldstein, and Michael Levin (iD) ✉ View all authors and affiliations

Even some of the simplest algorithms have emergent behaviors and 'side quests' (Zhang et al. 2025).

Likewise, natural language is comprised of an elementary structure-building algorithm of assembling minimal compositional schemes (like Nouns Phrases, Verb Phrases, etc.); from the simplest form ($_{XP}$[X Y]) we can build complex compositional conceptual instructions that cannot be reduced to the meaning of individual parts.

Mathematics is the study of... ~~numbers~~

~~Platonic entities and realms~~

~~collections of formulas~~

abstract structures ✅

... and human language is **one of the most interesting abstract structures in nature**.

**A simple grammar formalism—dependency grammar—motivated by the observation that longer distance connections between words are harder to make.**

*Syntax* provides a cognitive basis for syntactic structures across languages. Edward Gibson observes that there is a cognitive cost associated with connecting words that increases with the dependency length, such that shorter connections are preferred. A transparent formalism to represent this observation is dependency grammar, in which a word is simply connected to another word via a dependency arc to form a larger compositional meaning. This formalism can explain numerous aspects of word order universals across languages.

This book contrasts dependency grammar with the industry standard going back to Chomsky's phrase structure grammar with transformations. Dependency grammar is a simpler formalism: It does not posit the existence of categories that combine words. Furthermore, there are no transformations. Gibson argues that a construction-based dependency grammar is not only simpler than a phrase structure with transformations approach, but it also accounts for language phenomena more effectively.
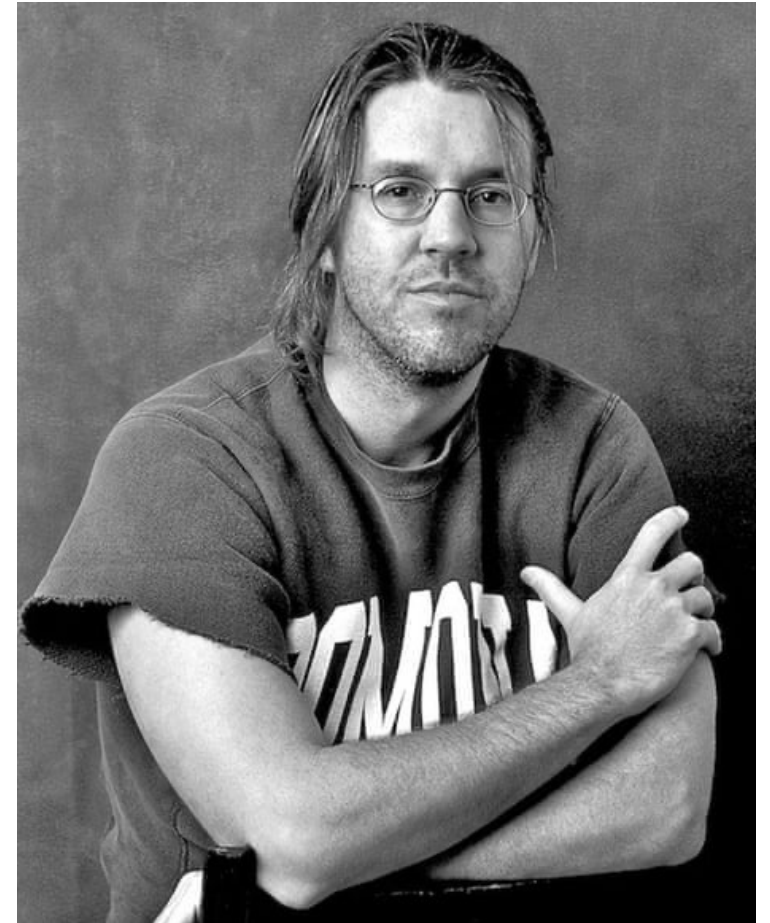
Dependency grammar is not a viable candidate theory of natural language syntax, since it isolates word-word dependency graphs rather than hierarchical constituency structure (dependency graphs are projections from this deeper algebraic structure).

Certain of these theories in the literature are intentionally positioned simply to counter generative grammar, and be distinct from it, rather than being formed from a novel first principles account of *what linguistic knowledge is*.

"...defining yourself in opposition to something is still being anaclitic on that thing, isn't it?"

**–David Foster Wallace, *Infinite Jest***

The internal competence of LLMs is often better than their actual output performance – what they 'know' exceeds what they might appear to know when playing around with ChatGPT....

# Large Language Models as Neurolinguistic Subjects: Discrepancy between Performance and Competence

**Linyang He**[1,2]     **Ercong Nie**[3,4]     **Helmut Schmid**[4]
**Hinrich Schütze**[3,4]     **Nima Mesgarani**[1]     **Jonathan Brennan**[2]

[1]Columbia University  [2]University of Michigan
[3]Munich Center for Machine Learning, Germany  [4]LMU Munich, Germany

linyang.he@columbia.edu     {nie,schmid}@cis.lmu.de
hinrich@hotmail.com     nima@ee.columbia.edu     jobrenn@umich.edu

...While this is completely true, it has so far only been shown that the distance and direction of LLM embeddings might represent some aspect of dependency grammars (Diego-Simon et al. 2024) – but dependency grammar is not a candidate model for compositional syntactic knowledge.

## A polar coordinate system represents syntax in large language models

**Pablo Diego-Simón**
ENS, PSL University, Paris, France
pablo-diego.simon@psl.eu

**Stéphane D'Ascoli**
Meta AI, Paris, France
stephane.dascoli@gmail.com

**Emmanuel Chemla**
ENS, PSL University, Paris, France
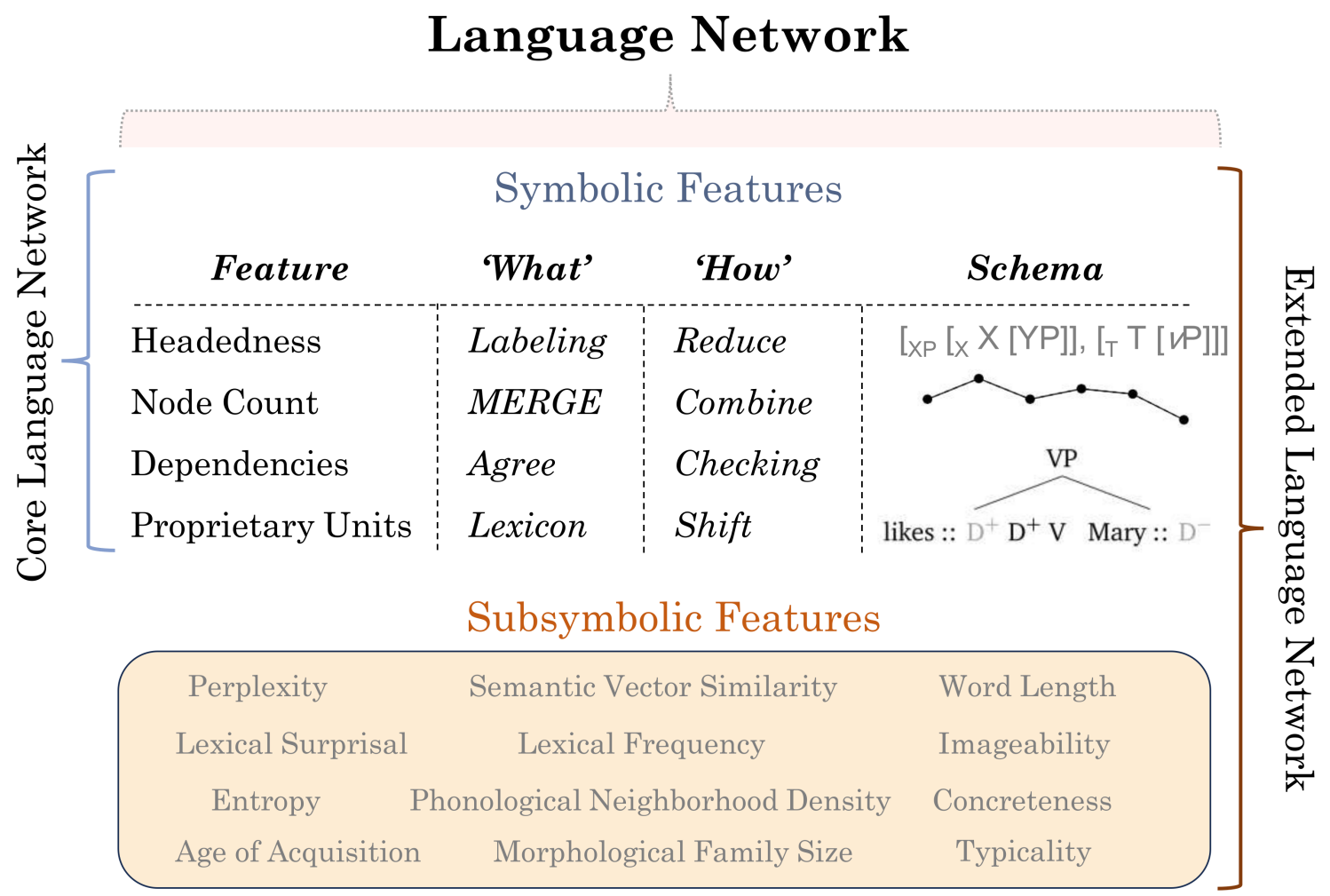emmanuel.chemla@ens.psl.eu

**Yair Lakretz**
ENS, PSL University, Paris, France
yair.lakretz@gmail.com

**Jean-Rémi King**
Meta AI, Paris, France
jeanremi@meta.com

We are looking for candidate neural mechanisms that:

Respect **set grouping** (non-associativity)

Demonstrate sensitivity to a **typed system** of semantic categories

Reflect the **asymmetric** (headed) nature of phrase composition

Can **recursively** embed headed structures inside other structures

Are sensitive to long-distance **dependencies** between elements

## Language Network

Core Language Network

Extended Language Network

### Symbolic Features

| Feature | 'What' | 'How' | Schema |
|---|---|---|---|
| Headedness | Labeling | Reduce | $[_{XP} [_X X [YP]], [_T T [vP]]]$ |
| Node Count | MERGE | Combine | VP |
| Dependencies | Agree | Checking | |
| Proprietary Units | Lexicon | Shift | likes :: $D^+ D^+ V$   Mary :: $D^-$ |

### Subsymbolic Features

| | | |
|---|---|---|
| Perplexity | Semantic Vector Similarity | Word Length |
| Lexical Surprisal | Lexical Frequency | Imageability |
| Entropy | Phonological Neighborhood Density | Concreteness |
| Age of Acquisition | Morphological Family Size | Typicality |

Murphy (2025), *Cognitive Neuroscience*

We are looking for candidate neural mechanisms that:

Respect **set grouping** (non-associativity)

Demonstrate sensitivity to a **typed system** of semantic categories

Reflect the **asymmetric** (headed) nature of phrase composition ← **When learning their language, no child uses a strategy whereby a phrase's category is unrelated to any element inside it**

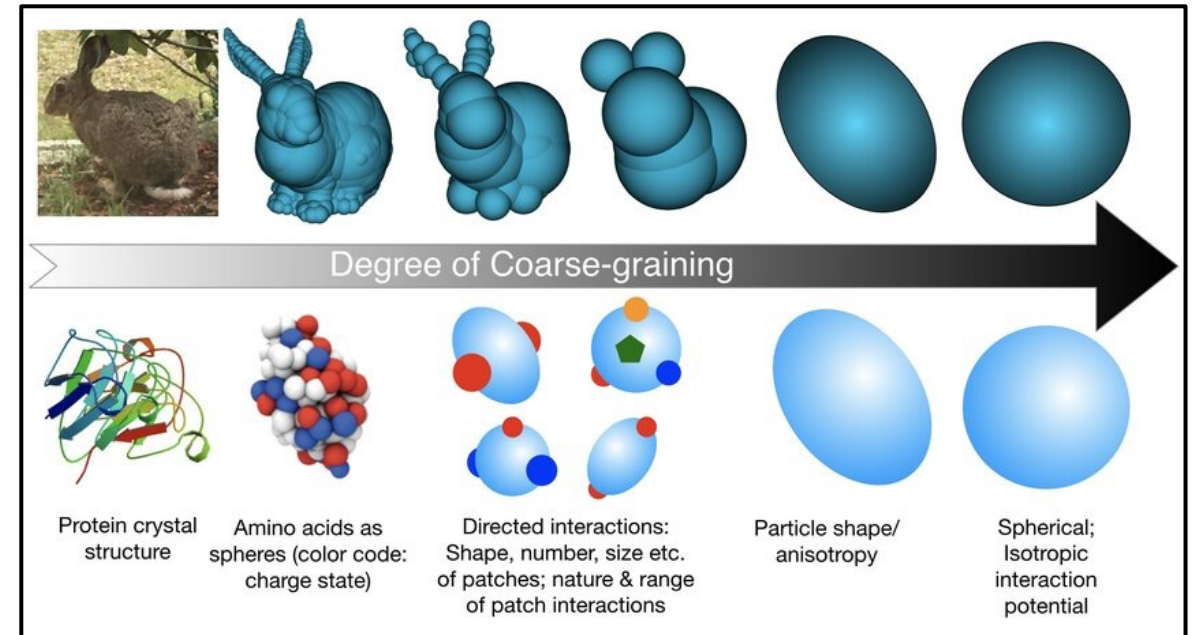Can **recursively** embed headed structures inside other structures

Are sensitive to long-distance **dependencies** between elements ← **And these dependencies are guided by structural information (no rules say "link this noun to this verb after 4 words")**

Merge-based syntax permits a kind of information coarse-graining, re-formatting concepts and giving us new kinds of 'things' to think about.

Language evolution didn't just make it easier to communicate – it also made it easier to think. Merge-based syntax gives us more precise coordinates in conceptual space.



Degree of Coarse-graining

Protein crystal structure

Amino acids as spheres (color code: charge state)

Directed interactions: Shape, number, size etc. of patches; nature & range of patch interactions

Particle shape/ anisotropy

Spherical; Isotropic interaction potential

{Those two old Italian men} are happy

{John} is happy

{Noun Phrase} {Verb Phrase}

In the literature, we often read: "Stroke patients have substantial language deficits, but they can still *think* (pass various cognitive tests). Therefore, language cannot be considered a thought system".

e.g., the authors equate evidence that non-linguistic thought survives after aphasia with the conclusion that language is not a thought system.

But this is flawed logic and a category error!

Inferring that "language is not a thought system" because non-linguistic thought survives after damage to language areas is like saying "vision is not a sensory system" because blind people can still smell and hear.

Glossa a journal of general linguistics

The **formal structure** of language is excellent for generating complex semantic inferences, but not as good for efficient/clear communication

## Language design and communicative competence: The minimalist perspective

Elliot Murphy[1,2,3]

(7) You persuaded John to buy a car.

Both the individual and the object being purchased can be questioned, but questioning the more deeply embedded object forces the speaker to produce some form of more complex circumlocution ([ ] = originally merged position of *wh*-expression).

(8) a. *[What] did you persuade who to buy [ ]?
    b. [Who] did you persuade [ ] to buy what?

The vast majority of our everyday use of language is not for "communication", but for organizing our thoughts, planning, strategizing, aiding directed attention, the consolidation of experience, reflecting on personal responsibilities, cognitive model updating, and more general and abstract forms of reflection.

"99% of the head's thinking activity consists of trying to scare the everliving shit out of itself."

**–David Foster Wallace, *Infinite Jest***

Even if non-verbal reasoning remains, **language facilitates certain higher-order processes,** and syntax provides precise instructions to these systems:

- Multi-step logical inference
- Counterfactual reasoning
- Inductive definitions
- Complex (nested) planning
- Evidentiality / belief representations



**John Lennon, 1964**

It may not be *intuitive* to think about language this way...

But Newton and later physicists showed that our intuitions about mass and motion are wrong.

Mendelian genetics showed that our intuitions about much of biology are wrong.

Gödel showed that our intuitions about mathematics are wrong, and of course mathematics is infamously riddled with counter-intuitive conclusions (think of Conway's surreal numbers).

The infinite series 1 + 2 + 3 + 4 + ..., though divergent, is associated via analytic continuation with the value ?

$$\sum_{n=1}^{\infty} n = \ ?$$

$$\sum_{n=1}^{\infty} n = -\frac{1}{12}$$

Why should our intuitions about language be taken seriously as feasible means to guide scientific theories, and why should our intuitions be used as additional constraints on how we define and operationalize concepts like 'thought' and 'language'?

Why is 'language' one of the only major topics in cognitive science where attitudes are so hidebound, and intuitions and biases are so unshakeable, that they guide theory-formation?

The word 'geometry' literally means 'land (earth) measurement'. But since classical times, geometry has now matured sufficiently to the point that more generalizable and abstract principles can be extracted, and modern geometry has no relation to hills and mountains.

The the same is true for the study of 'language'.

(1) Introduction

(2) Impediments to inference

(3) Language and its place in nature

(4) Distinct causal processes underpin the neurobiology of language

MRI    CT    OR Pictures    Localized Electrodes

MRI    CT    ROSA Trajectories    Localized Electrodes

**Trait**: Language allows us to construct **phrases out of individual words**

**Neural Basis**: Our intracranial recordings indicate unique involvement of **the pSTS**, *across different phrase types and sensory modalities*.



(1) SEARCH Lexicon: {blue}, {car}

(2) MERGE = {blue car}

Murphy et al. (2022), *J. Neuro.*

**Peter Hagoort**     **Katrien Segaert**     **Nitin Tandon**

Auditory stimuli, two-word phrases:

"red boat" (*real phrase*)
"bulg boat"
"red bulg"

Only "red boat" involves semantic composition



A

B
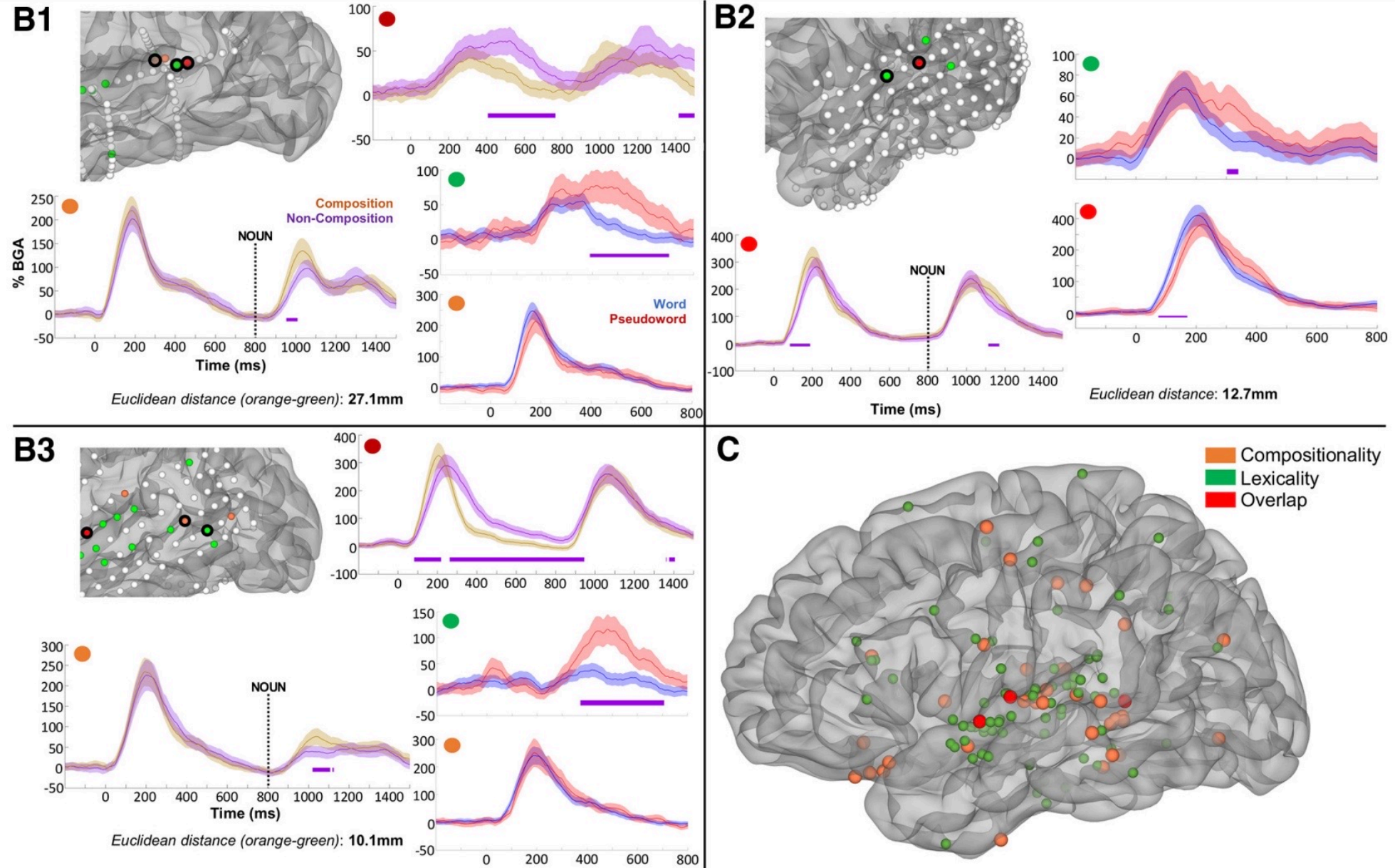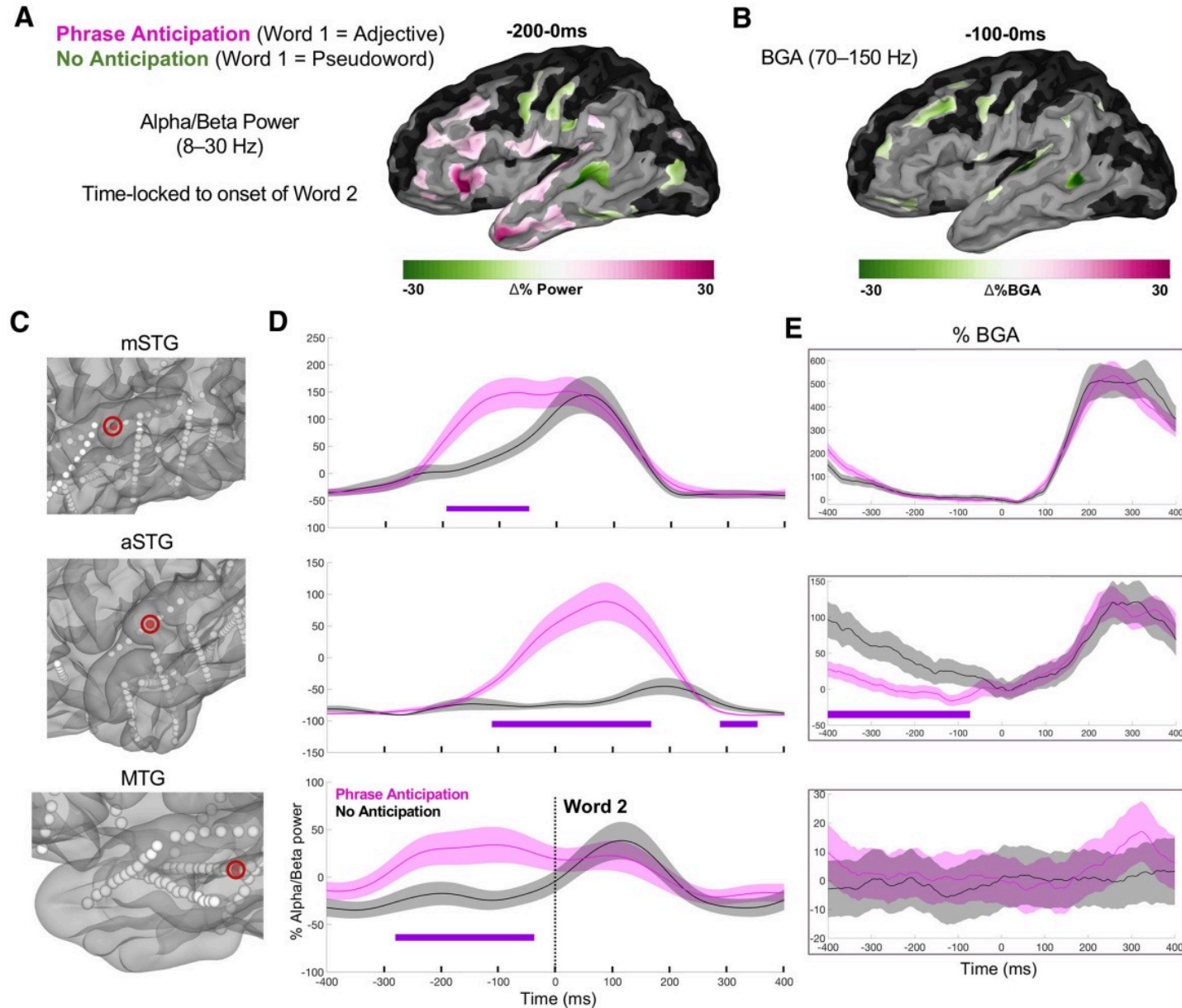Electrode Coverage

C
Activation to Auditory Stimuli

**High gamma activity** (70-150 Hz) is reflective of local cortical activation

Strongly correlated with the fMRI BOLD signal (blood flow, oxygenation)
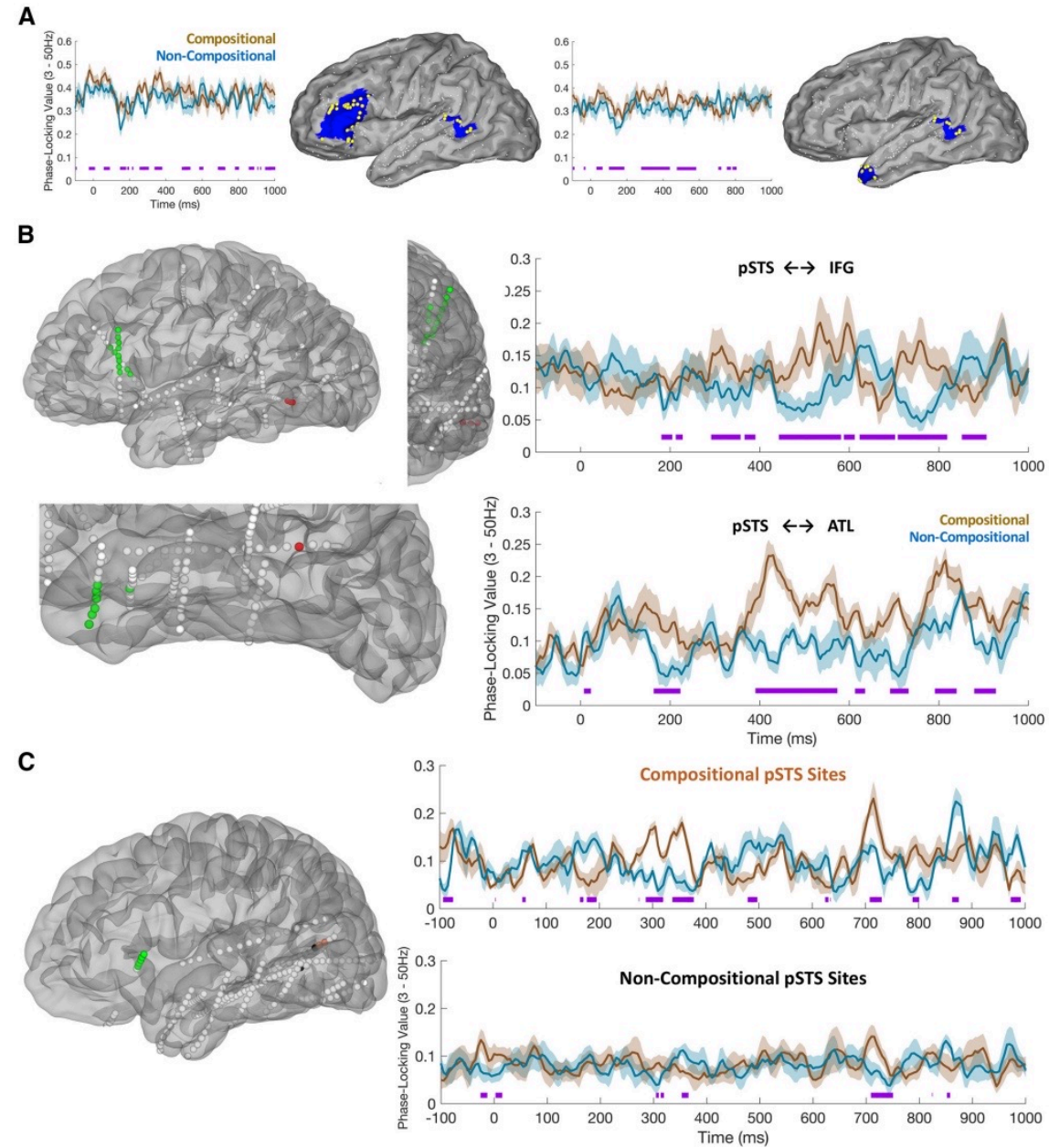
# Multiple dimensions of syntactic structure are resolved earliest in posterior temporal cortex

Elliot Murphy [a,b,*,1], Patrick S. Rollo [a,b], Katrien Segaert [c,d], Peter Hagoort [d,e],
Nitin Tandon [a,b,f,*]

[a] Vivian L. Smith Department of Neurosurgery, McGovern Medical School, University of Texas Health Science Center at Houston, Houston, TX 77030, United States
[b] Texas Institute for Restorative Neurotechnologies, University of Texas Health Science Center at Houston, Houston, TX 77030, United States
[c] School of Psychology & Centre for Human Brain Health, University of Birmingham, Birmingham B15 2TT, UK
[d] Max Planck Institute for Psycholinguistics, Nijmegen 6525 XD, the Netherlands
[e] Donders Institute for Brain, Cognition and Behaviour, Nijmegen 6525 HR, the Netherlands
[f] Memorial Hermann Hospital, Texas Medical Center, Houston, TX 77030, United States

**Phrase Composition**

100-300 ms    300-500 ms    500-700 ms    700-900 ms

**High gamma activity** (70-150 Hz) is reflective of local cortical activation

Strongly correlated with the fMRI BOLD signal (blood flow, oxygenation)

Murphy et al. (2024), *Prog. Neurobiol.*

Morphosyntactic Composition

Morphosyntactic Composition

## nature communications

# The spatiotemporal dynamics of semantic integration in the human brain

Elliot Murphy [1,2] ✉, Kiefer J. Forseth[1,2], Cristian Donos [3], Kathryn M. Snyder[1,2], Patrick S. Rollo [1,2] & Nitin Tandon [1,2,4] ✉

Language depends critically on the integration of lexical information across multiple words to derive semantic concepts. Limitations of spatiotemporal resolution have previously rendered it difficult to isolate processes involved in semantic integration. We utilized intracranial recordings in epilepsy patients (n = 58) who read written word definitions. Descriptions were either referential or non-referential to a common object. Semantically referential sentences enabled high frequency broadband gamma activation (70–150 Hz) of the inferior frontal sulcus (IFS), medial parietal cortex, orbitofrontal cortex (OFC) and medial temporal lobe in th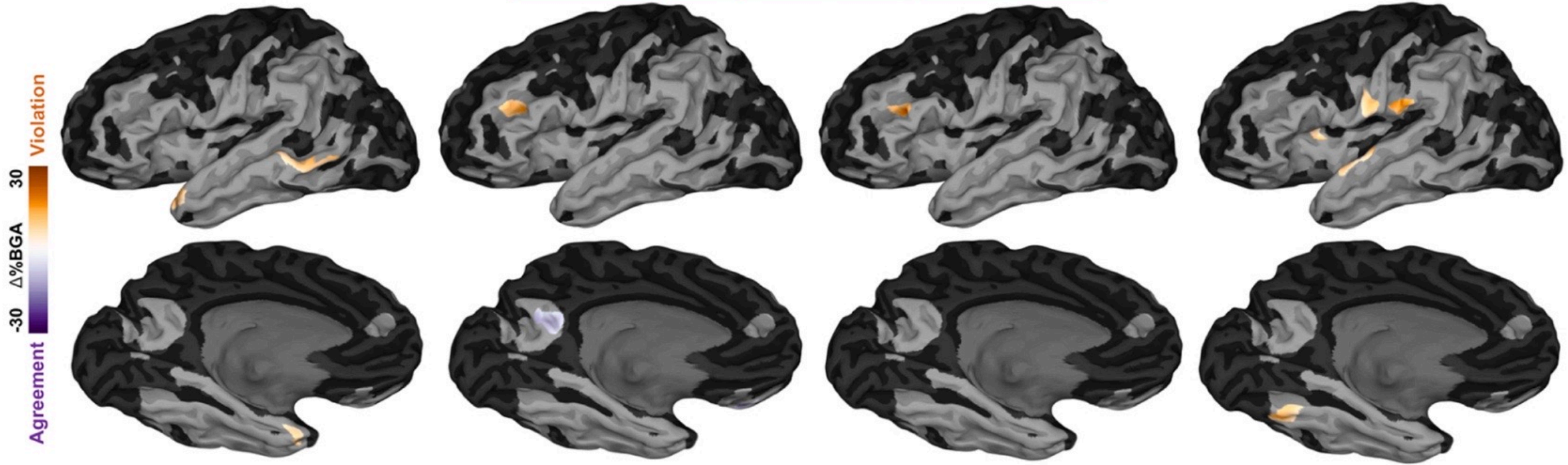e left, language-dominant hemisphere. IFS, OFC and posterior middle temporal gyrus activity was modulated by the semantic coherence of non-referential sentences, exposing semantic effects that were independent of task-based referential status. Components of this network, alongside posterior superior temporal sulcus, were engaged for referential sentences that did not clearly reduce the lexical search space by the final word. These results indicate the existence of complementary cortical mosaics for semantic integration in posterior temporal and inferior frontal cortex.

**a**

**b**

| | | | | | FRUIT | "Apple" |
| + | A | ROUND | RED | | AGE | "Nonsense" |
| 1000ms | 0ms | 500 | 1000 | 1500 | | |

**Mean sentence length:** 6.5 words
**Sentence length range:** 3–12 words
**# Trials per condition:** 84

**c**

Referential Definition

Non-Referential Definition

Response Time (ms)

1 — 15+
**#patients**

**d**

SENTENCE

**REFERENTIAL**

**LIMITED LEXICAL NARROWING**
An object used for weighing

**STRONG LEXICAL NARROWING**
It's white and falls from the sky in winter

**NON-REFERENTIAL**

**COHERENT**
A person at the circus who makes you commute

**INCOHERENT**
A place where oceans shop

**f**

300-500ms   500-700ms   700-900ms

INCOHERENT
COHERENT

Δ%BGA
20
−20

**g**

Coherent

Incoherent

Frequency (Hz)
150
70

0   0.5   1
Time (s)

0.5
−0.5

**Type** —— CROWN
**Token** —— *A specific crown*
**Referent** —— "What a king wears on his head"

Reference
Coherence
Narrowing
R + C
C + N
R + N

**500-700 ms**

**Inferior Frontal Sulcus**
Reference, Coherence, Narrowing
Lexico-semantic hub

**Medial Parietal Cortex**
Reference, Narrowing
Situation model construction

**Anterior IFG**
Coherence
Lexical control

**Posterior STS**
Reference
Narrowing
Compositional syntax-semantics

**Hippocampus / PHC**
Reference
Token memory representations

**700-900 ms**

**MFG / pIFG**
Reference
Narrowing
Lexical workspace

**Angular Gyrus**
Coherence
Narrowing
Lexico-thematic interface

**Orbitofrontal Cortex**
Reference, Coherence, Narrowing
Semantic saliency

**vmPFC**
Reference
Narrowing
Inferential memory

Article

Intraoperative cortical localization of music and language reveals signatures of structural complexity in posterior temporal cortex

Meredith J. McCarty,[1,2,6] Elliot Murphy,[1,2,6,*] Xavier Scherschligt,[1,2] Oscar Woolnough,[1,2] Cale W. Morse,[1,2] Kathryn Snyder,[1,2] Bradford Z. Mahon,[3,4] and Nitin Tandon[1,2,5,7,*]

**Musical Structure**

• Broad tone range (≥3)
• Easily detectable melody
• Lower algorithmic compressibility

• Tone repetition ($A^n$, $A^nB^n$, $A^nB^m$)
• Less easily detectable melody
• Higher algorithmic compressibility

**Syntactic Structure**

**Meredith McCarty**

Closely neighboring portions of pSTG index sensitivity to musical or syntactic complexity

pSTG shares neural resources for processing **units** (words, tones) but distinct resources for **structures**



- ● **CSM+**
- · **CSM-**

*Music*

Music Perception

Music Perception + Production

*Language*

Sentence Repetition

Common Object Naming

Music + Language

Most of these intracranial signatures are in high frequency gamma activity.

But what is the computational scope (expressive power) of gamma?

Different causal structures can provide unique explanatory perspectives in neuroscience.

"Neural mechanisms" can have causal and explanatory power, but there are also causal structures that are not classically "mechanistic" (e.g., cascades, pathways).

We also have non-causal entities that provide unique explanatory power (mathematical models, neural topologies).

Most high-impact factor journals will state their aim is to publish "mechanistic" insights – yet many journal editors are unable to explain what exactly a neural mechanism is (Ross & Bassett 2024).

Funding agencies also request "mechanistic" contributions – different agencies place distinct weight to this notion.

LLMs mostly obscure insights into neural mechanisms...

# EXPLANATORY STRUCTURES IN NEUROSCIENCE

## Mechanism



Hierarchical/part-whole
Fine-grained detail
Narrow (ion channel)
Broad (circuit-level)
"Clockwork"

## Topology



Structural-constraint
Neural geometries
Abstracted away from temporality
Manifolds
"Landscape"

*Thalamic ring attractors*

*Entorhinal cortex toroids*

## Pathway



Sequential flow
Linear, stepwise
Dictates routes
"Highway"

## Circuit



Fixed, closed-loop
Recurrent system
Mesoscale focus
Non-reductive
Drive oscillatory motifs
"Wiring diagram"

*Fear and extinction circuits*

## Cascade



Initial trigger
Amplification
Momentum
"Snowball effect"

X causes Y if we intervene on X in such a way as to change the outcome of Y ("interventionist" philosophy of science). (Counterfactual) causation here is almost synonymous with **control**. Causes can be proximal/distal. Pathways and cascades are not "mechanisms", but they do have causal-explanatory power.

If you cannot identify why something would *not* qualify as a neural mechanism, you have not successfully delimited the concept.

So we need a method to help migrate the concepts of linguistic theory into a testable framework...

**Question**: Which neurobiological scales of organization are going to be causally prominent and explanatorily useful for varying representational levels of language?

Observable output behavior greatly underdetermines the network implementation:

Prinz et al. 2004, 'Similar network activity from disparate circuit parameters', *Nature Neuroscience*

Jonas & Kording 2017, 'Could a neuroscientist understand a microprocessor?', *PLoS Computational Biology*

**Strict reductionism is extremely ambitious**

ELSEVIER

Check for updates

# ROSE: A neurocomputational architecture for syntax

## Elliot Murphy [a,b,*]

[a] *Vivian L. Smith Department of Neurosurgery, McGovern Medical School, UTHealth, 1133 John Freeman Blvd, Houston, TX, 77030, USA*
[b] *Texas Institute for Restorative Neurotechnologies, UTHealth, 1133 John Freeman Blvd, Houston, TX, 77030, USA*

### ARTICLE INFO

### ABSTRACT

A comprehensive neural model of language must accommodate four components: representations, operations, structures and encoding. Recent intracranial research has begun to map out the feature space associated with syntactic processes, but the field lacks a unified framework that can direct invasive neural analyses. This article proposes a neurocomputational architecture for syntax, termed ROSE (Representation, Operation, Structure, Encoding). Under ROSE, the basic data structures of syntax are atomic features, types of mental representations (R), and are coded at the single-unit and ensemble level. Operations (O) transforming these units into manipulable objects accessible to subsequent structure-building levels are coded via high frequency broadband γ activity. Low frequency synchronization and cross-frequency coupling code for recursive structural inferences (S). Distinct forms of low frequency coupling encode these structures onto distinct workspaces (E). Causally connecting R to O is spike-phase/LFP coupling; connecting O to S is phase-amplitude coupling; connecting S to E are frontotemporal traveling oscillations. ROSE is reliant on neurophysiologically plausible mechanisms and provides an anatomically precise and falsifiable grounding for natural language syntax.

## The Oscillatory Nature of Language

**Elliot Murphy**

# ROSE: A Universal Neural Grammar

Elliot Murphy[a,b]

[a]Vivian L. Smith Department of Neurosurgery, McGovern Medical School, UTHealth, Houston, TX, USA; [b]Texas Ins
Neurotechnologies, UTHealth, Houston, TX, USA

**ABSTRACT**

Processing natural language syntax requires a negotiation between symbolic and subsymbolic representations. Building on the recent representation, operation, structure, encoding (ROSE) neurocomputational architecture for syntax that scales from single units to inter-areal dynamics, I discuss the prospects of reconciling the neural code for hierarchical syntax with predictive processes. Here, the higher levels of ROSE provide instructions for symbolic phrase structure representations (S/E), while the lower levels provide probabilistic aspects of linguistic processing (R/O), with different types of cross-frequency coupling being hypothesized to interface these domains. I argue that ROSE provides a possible infrastructure for flexibly implementing distinct types of minimalist grammar parsers for the real-time processing of language. This perspective helps furnish a more restrictive 'core language network' in the brain than contemporary approaches that isolate general sentence composition. I define the language network as being critically involved in executing specific parsing operations (i.e. establishing phrasal categories, tree-structure depth, resolving dependencies, and retrieving proprietary lexical representations), capturing these network-defining operations jointly with probabilistic aspects of parsing. ROSE offers a 'mesoscopic protectorate' for natural language; an intermediate level of emergent organizational complexity that demands multi-scale modeling. By drawing principled relations across computational, algorithmic and implementational Marrian levels, ROSE offers new constraints on what a unified neurocomputational settlement for natural language syntax might look like, providing a tentative scaffold for a 'Universal Neural Grammar' – a species-specific format for neurally organizing the construction of compositional syntactic structures, which matures in accordance with a genetically determined biological matrix.

An emerging consensus in neuroscience is that complex behavior and cognition rely on coordinated interactions between brain regions, with phase synchronization being a major candidate for implementing this coordination, by gating information transmission.

Yet, unlike for models of attention and working memory, there is a current absence of oscillatory phase coding in models of natural language.

To summarize a long history of research…

Reliable signatures of syntactic structure tend to be found in the low-frequency range, whereas reliable signatures of semantic composition and lexical information are rapidly found (e.g., in ECoG, early processing windows) in higher frequencies.

There are, naturally, a few exceptions here, and it is this tension and reconciliation that ROSE attempts to address.

LLMs will be helpful with isolating properties of phonological and lexical statistics at R and O levels

...But will be less helpful with isolating higher-order syntactic inferences at S and E levels

**Representation**

F φ
F v/n
F animate

*Single-unit encoding of distinct conceptual/formal features that compose into manipulable objects at O and S*

**R » O: Spike-phase coupling**

**Operation**

"red" + "boat"

GAMMA

*High γ ensembles transform coherently bound R units into complex feature-bundles accessible to S*

**O » S: Phase-amplitude coupling (PAC)**

**Structure**

THETA

N

red    boat
A      N

DELTA

*Low-frequency synchronization and cross-frequency coupling to generate recursive categorial inferences*

**S » E: Traveling waves coordinate frontotemporal PAC**

**Encoding**

δ-θ PAC
HIERARCHICAL MEMORY

θ-γ PAC
LEXICAL MEMORY

WS = [{X, {X, Y}}]

δ-γ PAC
COMBINATORIAL PROCESSING

*θ-γ lexical memory code (linearized) synchronized to δ-θ supraordinate structure memory (hierarchical)*

**E » ... R: Low-frequency phase reset of spike-LFP**

**Levels of ROSE are mechanistically connected (e.g., S-to-E)**

Murphy (2025), *Cognitive Neuroscience*

Most accounts either place heavy emphasis on symbolic knowledge (e.g., Chomsky, 2013, 2023, Chomsky et al., 2019, 2023; Friederici, 2017; Murphy, 2023, 2024; Murphy, Holmes et al., 2024) or predictive processing (e.g., Caucheteux et al., 2023; Kwiatkowski et al., 2012; Schrimpf et al., 2021; Zhou et al., 2025), without a means for integration.

ROSE (Murphy 2025, *Cognitive Neuroscience*) offers a concrete means to integrate statistical learning with symbolic knowledge via mechanisms like PAC.

**Single-sentence summary of ROSE**

The higher levels of ROSE provide instructions for symbolic phrase structure representations (S/E), while the lower levels provide probabilistic aspects of linguistic processing (R/O), with different types of cross-frequency coupling being hypothesized to interface these domains.

Empirical and conceptual motivations are presented to defend the idea that δ-θ inter-regional phase-amplitude coupling constructs multiple sets of syntactic and semantic features, and imposes biases on how to read out the items provided by this phase code. This occurs when the phase of δ is synchronized with the amplitude of θ – in turn, θ phase couples with high-frequency local cortical processing.

δ represents supraordinate syntactic categories, and θ represents feature-bundles generated via lexical access.

Phase-resetting of this mechanism, alongside concurrent encoding/storage of its products in workspaces before a newly-generated δ-θ complex is created, permits a facility for recursive self-call.

Murphy (2025, *Cognitive Neuroscience*) provides more explicit, concrete details about the phase and frequency dynamics.

The bulk of work under ROSE is achieved by a frontotemporal symbolic low-frequency phase code interacting via cross-frequency coupling with a series of local probabilistic inferences over lexico-semantic content, with the latter being implemented via spike-phase coupling assembling bundles of linguistic features and which can emerge into 'dynamical motifs'.

**R:** Single-unit encoding of conceptual features and formal syntactic features. This level involves a cellular barcode for distinct features that compose into syntactic objects coherently bound by high $\gamma$ at O. It also involves vector codes for ensembles hosting features common to objects represented at O and that are ultimately coordinated by S.

**O:** High $\gamma$ sensorimotor transformations into lexicalized objects (core network nodes: mid-fusiform cortex, orbitofrontal cortex, middle temporal gyrus, inferior frontal cortex, intraprietal sulcus) accessible to $\delta/\theta$ phase-locking. This level can implement the semantic composition of language-specific concepts (minimal phrase schemes) that coordinate the firing of R units. High $\gamma$ activates assemblies of distinct units hosting the barcode or vector code for units $R_1 \ldots R_n$ that compose into feature-bundles.

**S:** A low frequency neural program for generating structural inferences over O. $\delta$-$\theta$ phase-amplitude coupling (posterior superior temporal sulcus to inferior frontal gyrus) for categorial inferences modulating the representation of feature-bundles in $\theta$-$\gamma$ by structuring the read-out of these complexes (frontotemporal language sites to cross-modular hubs).

**E:** Local and global workspaces for bottom-up lexical memory and top-down hierarchical memory. Traveling waves implement $\delta$-$\theta$ coupling for hierarchical memory, and $\theta$-$\gamma$ coupling for lexical memory. $\alpha$ power codes for workspace 'disruption' (posterior temporal and inferior frontal cortex). $\beta$ power coding for syntactic predictions (inferior frontal cortex).

The basic data structures of syntax are atomic features; types of linearly readable mental representations (R) that are coded at the single-unit and ensemble level.

I assume that representations encompass any object manipulable by the generative component of language, being composed of features determining constraints on operations, such as selection, agreement, licensing and movement.

Examples include [N], [Plural], [Dem], [C], [T], [P], and also conceptual features pertaining to lexical roots, like √BREAK and √HOME. Syntax builds structure through recursive applications of MERGE, and these are then entered into a space of syntactic working memory. Lexical items are simplex conceptual atoms (bundles of features).

If these features seem exotic and not "neurally plausible", consider how we happily entertain stubby-animate concepts in higher-order vision.

Representational eccentricity is already normalized in sensory neuroscience ("elongated blue edge detector", etc).



Can we normalise functional syntactic features (e.g., D, P)?

If feature bundles like [N] or [T] appear exotic, it is only because linguistics has historically lacked the luxury of metaphoric 'reification' that vision science enjoys. Stubby cells, simple cells, grandmother neurons are all idealized explanatory constructs.

ROSE's representational primitives are of precisely the same ontological kind: simplified handles on the high-dimensional, feature-specific subspaces that neural populations inhabit.

Depending on the lexical item in question, spikes from relevant cortical regions will be coordinated by spike-phase coupling (e.g., posterior middle temporal cortex for abstract word features; inferior parietal cortex for eventive features; anterior temporal lobe for object features; ventrotemporal cortex for face and place features; inferior frontal cortex for more formal and function-word related features)



Murphy, Woolnough et al. (Submitted), 'Cortical cascades support rapid semantic inference during reading'

Depending on the lexical item in question, spikes from relevant cortical regions will be coordinated by spike-phase coupling (e.g., posterior middle temporal cortex for abstract word features; inferior parietal cortex for eventive features; anterior temporal lobe for object features; ventrotemporal cortex for face and place features; inferior frontal cortex for more formal and function-word related features)



Murphy, Woolnough et al. (Submitted), 'Cortical cascades support rapid semantic inference during reading'

ROSE invokes ephaptic coupling and dynamical motifs to guide the construction of minimally complex linguistic objects (i.e., clusters of features that assemble into morphemes – *not* recursive compositional syntactic objects) at the appropriate level of resolution whereby statistical and symbolic instructions will (presumably) interface.

These would likely be *commonly co-occurring lexico-semantic features* that routinely get merged via the mechanisms of ROSE.

Empirical data for stable local attractors or 'motifs' in higher-order language regions remain limited – demonstrating them for lexical composition is non-trivial.

Still, the use of dynamical motifs for lexical information seems all the more plausible in the face of widespread, flexible and mixed high-dimensional coding across the cortical hierarchy.

Routledge
Taylor & Francis Group

## COMMENTARY

Check for updates

# Dynamical motifs for computations in language

Katarína Labancová [iD] and Nina Kazanina

Department of Basic Neurosciences, Faculty of Medicine, University of Geneva, Switzerland

Murphy's (2025) *ROSE: A Universal Neural Grammar* brings up the notion of dynamical motifs (DMs, Driscoll et al., 2024) as one of candidate mechanisms for bringing ROSE toward neurobiological plausibility. Dynamical motifs are patterns of neural activity that implement specific computations and can be reused across tasks that share components. Murphy proposes that DMs can generate, at the R/O levels, feature bundles associated with a given word ('[s]pike-phase coupling would trigger appropriate dynamical motifs to generate feature-bundles associated with a given word'), but details of this process remain unclear.

e.g., Noun vs Verb to the word *kick* in *a kick* vs *we kick*; or the need to assign a thematic role to a noun phrase, i.e., agent vs patient role to 'the horse' in *The horse kicked* vs. *The horse was kicked*. Binding to the correct thematic role, albeit more abstract, parallels color classification: the same way that the monkey has to categorize every object into one of two color categories, the noun phrase has to be bound to one of several thematic roles. Such computations can each be instantiated via a DM where the word+grammatical category or word+role complex occupies an independent subspace. As the sentence continues, different DMs are used to assign

## Related criticism to the LLM-brain debate...

- **Fitting RNNs make attractors even when the data does not have attractors**

- "A popular approach to study brain data is to fit an RNN to neural data and then show that this RNN has attractors. However, a recent paper [Qian et al. 2024, NeurIPS] showed that in the context of partial observations, a simulated system that has no attractors gives rise to a fitted RNN that does have attractors. In other words, in the context of unobserved data (always the case in neuroscience), we may even wrongly infer that there are attractors." (Konrad Kording)

New Results

🔔 **Follow this preprint**    ◀ Previous

# Modality-Specific and Amodal Language Processing by Single Neurons

Posted November 22, 2024.

Yair Lakretz, Naama Friedmann, Jean-Rémi King, Emily Mankin, Anthony Rangel, Ariel Tankus, Stanislas Dehaene, Itzhak Fried

This article is a preprint and has not been certified by peer review [what does this mean?].

↗ **Download PDF**

▽ **Print/Save Options**

✎ Revision Summary

| **Abstract** | Full Text | Info/History | Metrics | 📄 Preview PDF |

IPL
Event Semantics
Lexico-Semantics
(LATE)

pSTS
Compositional Syntax
(EARLY)

aIFG
Syntactic Linearization
Conceptual Structure
(LATE)

pIFG
Syntactic Workspace
Logico-Syntactic Rules

mFus
Lexico-Semantics
(EARLY)

pMTG
Conceptual Structure
(EARLY)

Murphy (2025), *Cognitive Neuroscience*

Prefrontal cortex may host circuits supporting variable binding, sequencing, gating and working memory storage (syntax-external demands), while lateral posterior temporal cortex might more reliably and efficiently subserve cross-modal semantic integration, supraordinate categorization and recursive hierarchical processing (syntax-internal demands).

Traveling waves build up the workspace cycle-by-cycle (evidence-accumulation for syntactic inferences bounded by the cross-frequency coupling dynamics over language hubs in pSTS and IFS/IFG)

# Rotating brain waves help thought circle back to the task at hand



Correct

Error

→ mean trajectory

Batabyal et al. 2025, JoCN

**Rotating brain waves help thought circle back to the task at hand**

"There is no reason in principle why a rotation in this mathematical subspace should correspond directly to a rotation on the surface of the cortex. But it does. That suggests to me that the brain is using these traveling waves to actually do computation, analog computation. Analog computation is way more energy efficient than digital and biology favors energy efficient solutions."

- Earl Miller (Picower Institute 'News' summary)

**Minimalist Grammar Parsing Algorithm**

BU — Passive 'Incremental'
LC — Mildly Predictive
TD — Active 'Eager'

**Not just next-token prediction!**

Murphy (2025), *Cognitive Neuroscience*

Tree terminals: The / king / forgave / nobody / he / met

| Syntactic Node Counts | The | king | forgave | nobody | he | met |
|---|---|---|---|---|---|---|
| Top-down | 3 | 1 | 2 | 3 | 3 | 2 |
| Bottom-up | 1 | 2 | 1 | 2 | 2 | 6 |
| Left-corner | 2 | 2 | 2 | 3 | 3 | 3 |

**PAC Modulations as a Function of Parsing Operation**

Morphological Feature Combinatorics

| $\delta - \gamma$ | The | king | forgave | nobody | he | met |
|---|---|---|---|---|---|---|
| | ↑ | – | ↑ | ↑ | – | – |
| | – | ↑ | ↓ | ↑ | – | ↑ |
| | ↑ | – | – | ↑ | ↑ | – |

Headedness; Syntactic Workspace

| $\delta - \theta$ | The | king | forgave | nobody | he | met |
|---|---|---|---|---|---|---|
| | ↑ | ↑ | – | ↑ | – | ↑ |
| | – | ↑ | – | ↑ | ↑ | ↑ |
| | ↑ | – | ↑ | ↑ | ↑ | – |

Reduce/Node Closure; Derivational Complexity

| $\delta - \beta$ | The | king | forgave | nobody | he | met |
|---|---|---|---|---|---|---|
| | ↑ | ↓ | ↑ | ↑ | ↑ | ↓ |
| | ↑ | ↓ | ↑ | ↑ | ↑ | ↓ |
| | ↑ | ↓ | ↑ | ↑ | ↑ | ↓ |

Lexical Feature-Bundling; Lexical Memory

| $\theta - \gamma$ | The | king | forgave | nobody | he | met |
|---|---|---|---|---|---|---|
| | ↑ | ↑↑ | ↑ | ↑ | ↑ | ↑↑ |
| | ↑ | ↑↑ | ↑ | ↑ | ↑ | ↑↑ |
| | ↑ | ↑↑ | ↑ | ↑ | ↑ | ↑↑ |

$\alpha$
Lexico-Semantic 'Shielding'
Stabilizing Structural Inferences

$\beta$
Set-Maintenance
Phrasal Anticipation

$\gamma$
'Shift' MG Operations
Lexico-Semantic Unification

"But it's not cognitively plausible to have a simple one-to-one isomorphism between some parsing computation and a narrow frequency band."

Indeed.

The spatiotemporal dynamics I invoke are to be thought of as the principal but not exclusive drivers of parsing operations. These dynamics are the PAC relations with the largest explanatory-causal scope for specific parsing operations, and the ones that best predict and drive them.

"Ok, but doesn't this make falsifiability of ROSE difficult?"

No.

ROSE provides a possible infrastructure for flexibly implementing distinct types of parsing operations for the real-time processing of language. This perspective helps to furnish a more restrictive 'core language network' in the brain than current language localizers that isolate general sentence composition (e.g., "sentences > wordlists").

I define the language network as being **critically involved in selectively computing and representing specific parsing operations.**

Correspondence | Published: 09 August 2024

# The language network is topographically diverse and driven by rapid syntactic inferences

Elliot Murphy ✉ & Oscar Woolnough

# ROSE indexes a "mesoscopic protectorate"

## The middle way

R. B. Laughlin*, David Pines[†‡§], Joerg Schmalian[¶], Branko P. Stojković[∥]**, and Peter Wolynes[††]

*Department of Physics, Stanford University, Stanford, CA 94305; [†]Institute for Complex Adaptive Matter, University of California, Office of the President, Oakland, CA 94607; [‡]Los Alamos Neutron Science Center Division and [∥]Theoretical Division and Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, NM 87545; [¶]Department of Physics and Astronomy and Ames Laboratory, Iowa State University, Ames, IA 50011; and [§]Science and Technology Center for Superconductivity and [††]School of Chemical Sciences, University of Illinois, Urbana, IL 61801

Contributed by David Pines, October 29, 1999

Mesoscopic organization in soft, hard, and biological matter is examined in the context of our present understanding of the principles responsible for emergent organized behavior (crystallinity, ferromagnetism, superconductivity, etc.) at long wavelengths in very large aggregations of particles. Particular attention is paid to the possibility that as-yet-undiscovered organizing principles might be at work at the mesoscopic scale, intermediate between atomic and macroscopic dimensions, and the implications of their discovery for biology and the physical sciences. The search for the existence and universality of such rules, the proof or disproof of organizing principles appropriate to the mesoscopic domain, is called the middle way.

the very large and the very small. But, as we all know, there is life in the desert.

The miracles of nature revealed by modern molecular biology are no less astonishing than those found by physicists in macroscopic matter. Their existence leads one to question whether as-yet-undiscovered organizing principles might be at work at the mesoscopic scale, at least in living things. This is by any measure a central philosophical controversy of modern science, for a commonly held view is that there are no principles in biology except for Darwinian evolution. But what if this view is just a consequence of our inability to see? Indeed the rules of self-organization at macroscopic length scales were not self-evident at the time of their discovery and were accepted as true only after
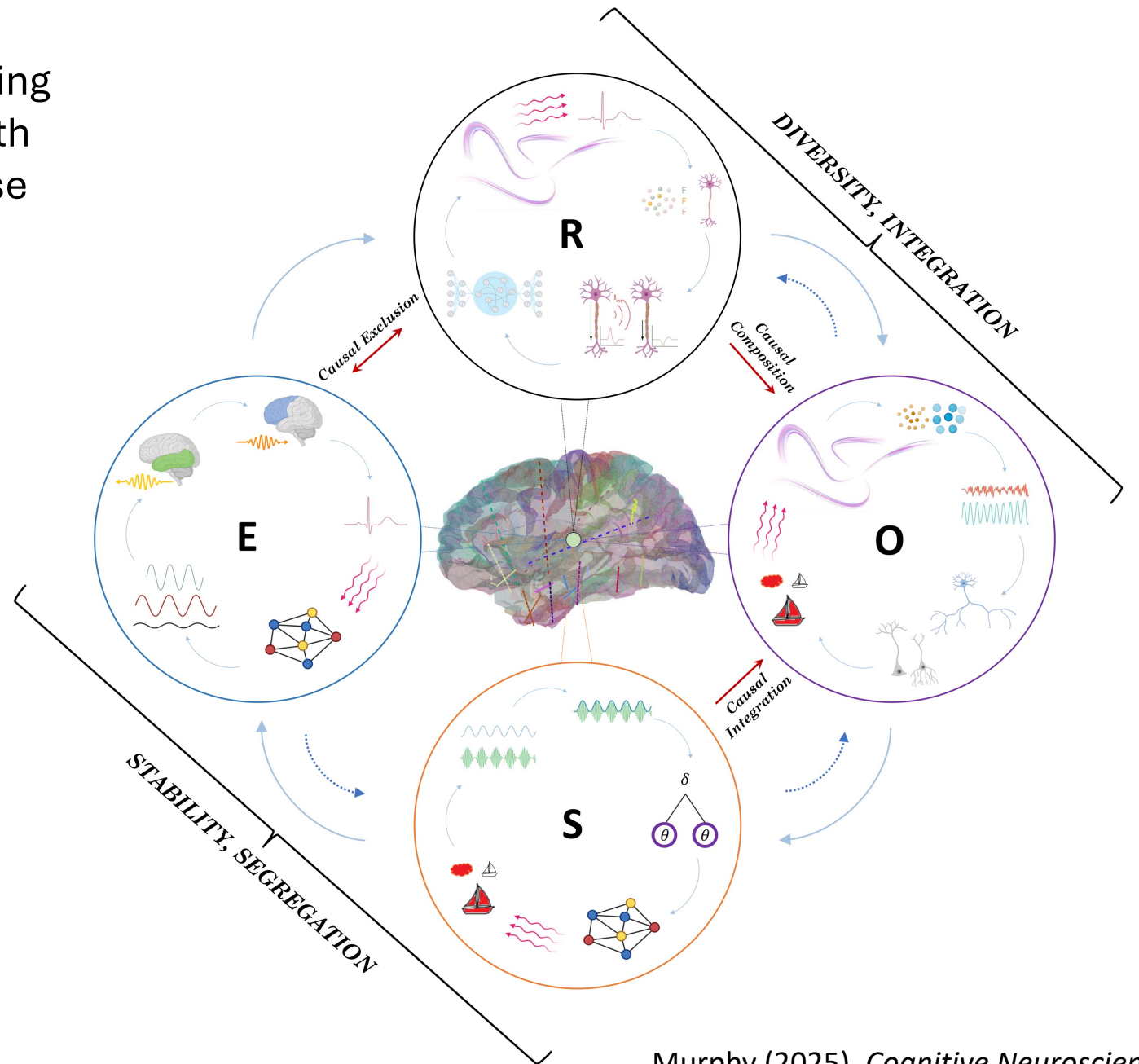
Phase-amplitude coupling (PAC) coordinates the **hierarchical assembly of features**, forming headedness inferences (i.e., the assembly with greatest relative PAC strength codes for phrase head), while **dynamical motifs and spike-phase coupling** coordinate constellations of semantic features into lexical items

Causal **composition, integration** and **exclusion** apply across levels, yielding a "mesoscopic protectorate" in the brain for syntactic inferences

ROSE uses a combination of *mechanisms*, *cascades* and *topological* structures to neurally enforce non-associativity, commutativity and other algebraic properties of human language, ensuring that $((\alpha\ \beta)\ \gamma) \neq (\alpha\ (\beta\ \gamma))$



• *Universal Neural Grammar* •

*DIVERSITY, INTEGRATION*

*STABILITY, SEGREGATION*

Causal Exclusion

Causal Composition

Causal Integration

R

E

O

S

Murphy (2025), *Cognitive Neuroscience*

ROSE enables potentially unbounded recursion until rising fronto-parietal α signals that the syntactic workspace (δ-θ dynamics) is full, while frontotemporal traveling δ waves ferry each completed complex to working memory buffers.

Each δ-cycle ends with a β-mediated 'commit' burst (Lundqvist et al. 2024) that silences the γ carriers for the daughters, ensuring the composite enters the next workspace step as an indivisible unit. Because headedness is selected by mutual information PAC strength rather than presentation order, the system is commutative, yet the β commit burst freezes the set, preserving non-associativity.

*(see Murphy 2025 for specific details and empirical support – Supplementary Materials include a mini-review of how LLMs fail to capture higher-order language, and a comprehensive Table comparing different connectionist models of composition)*

One of the interesting consequences of this theory is that we effectively get commutativity 'for free', given that we assume headedness is established via strength of PAC-active nodes.

We don't have to stipulate some additional mechanism to allow for commutativity, since it's not the *order* of PAC complexes that matters, just the strength.
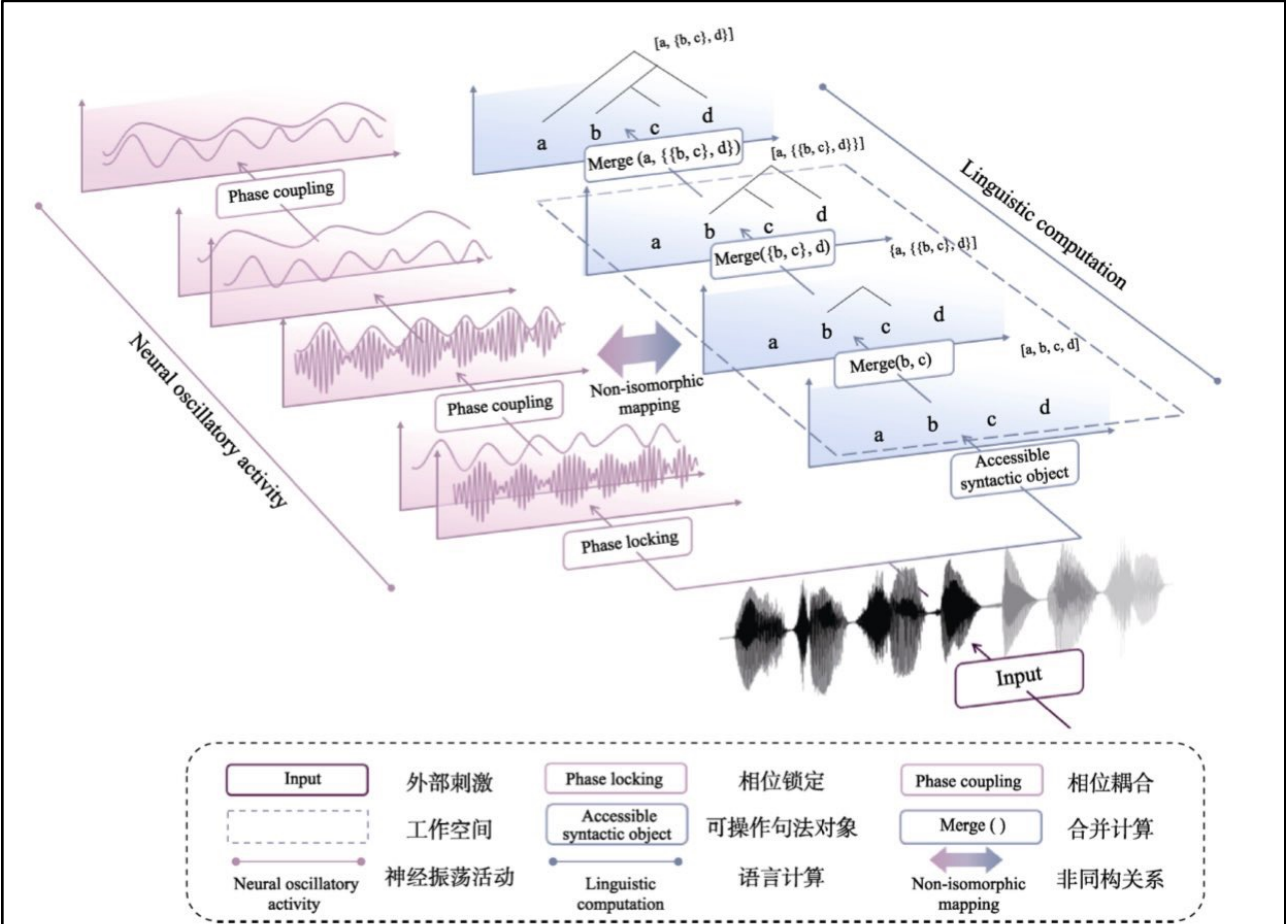
图 1 句法操作实现过程及其神经机制

表 1 以神经振荡为基的句法解析模型汇总

| 模型简称 | 核心概念 | 三分框架 | | | 粒度 | 出处 |
| --- | --- | --- | --- | --- | --- | --- |
| | | 实现层 | 算法层 | 计算层 | | |
| ROSE | 交叉频率耦合 | | 全频段活动 | 合并操作 | 细 | Murphy, 2024 |
| CNAL | 神经流形 | 语言功能区神经 | δ, θ, γ 活动 | 语义组合 | 细 | Martin, 2020 |
| SMMM | 记忆维持 | 集群随时间变化 | δ 活动 | 句法预测 | 粗 | Ding, 2020 |
| DORA | 符号−联结主义 | 的渐进激活 | δ 活动 | 论元结构 | 粗 | Martin & Doumas, 2017 |
| VS-BIND | 符号−联结主义 | | θ-γ 耦合 | 邻接依存 | 粗 | Calmus et al., 2020 |

Qi et al. (2025), *Adv. Psy. Sci.*

Reflecting the ever-flexible nature of symbolic knowledge, we do not need to assume that frequency bands here are strict types with rigid functional interpretations; rather, they are more likely to be what Martin (2020, *J Cog Neuro*) calls "tokens of processes with physiological bounds that render them into functional types".

e.g., children's low-frequency cortical tracking of syntax is slightly different from the adult brain – the canonical frequency band itself is less important than the causal, structuring force of the signal.

These are reflective of endogenous timescales of specific computations, rather than being fixed and strict bounds. What is critical is the logical and causal relations of neural structures invoked by ROSE.
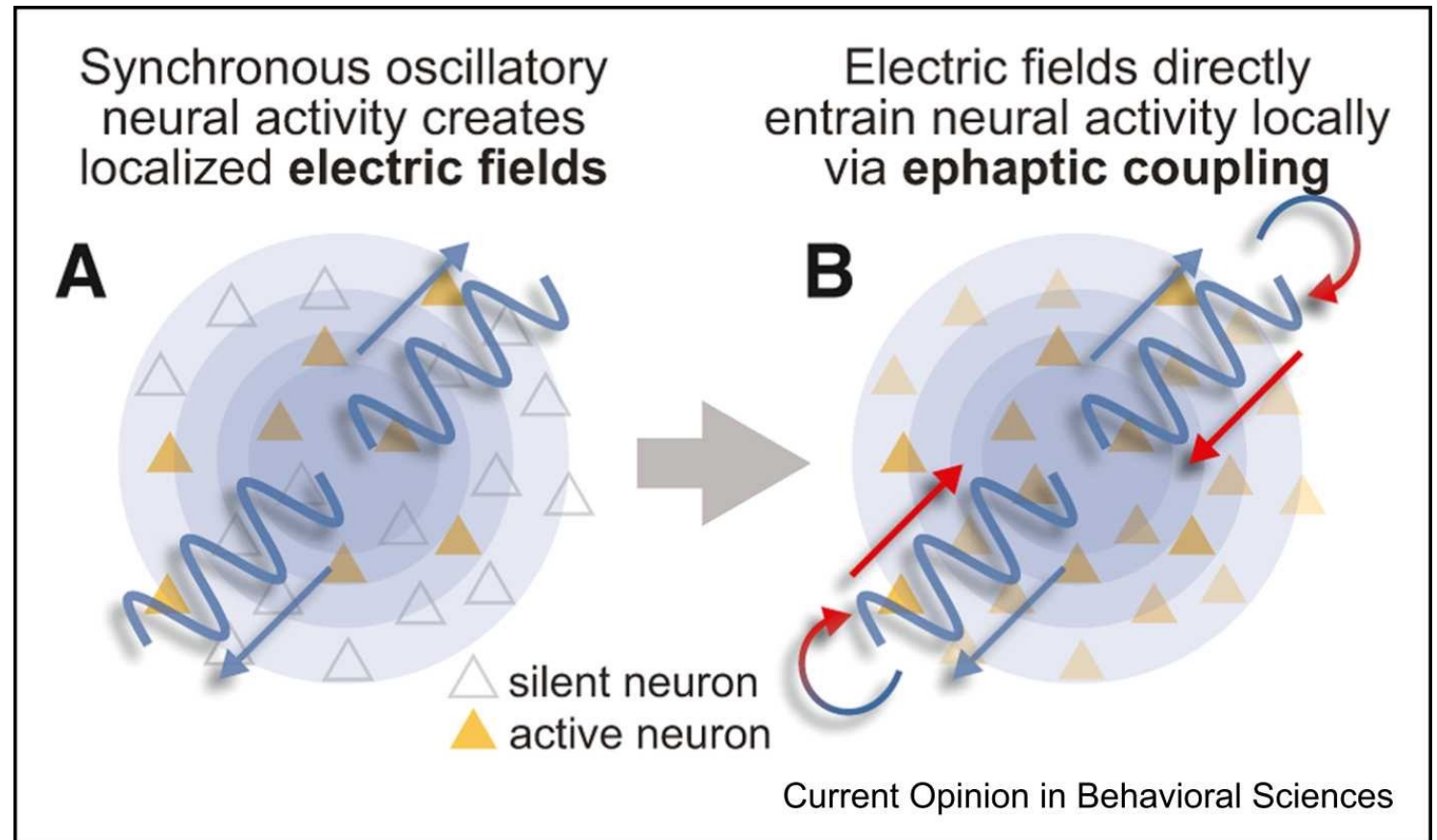
Do different types of macro-vs. micro-scale traveling waves form their own internal hierarchy of sensitivity to distinct structures or higher-order relations in language processing, e.g., global situation model maintenance vs. local phrase structure coordination?

"One set of emergent properties — organization of coding and communication in subspaces — is observed at the local spiking level. Another set of properties — oscillations that can bidirectionally influence spikes and organize cortical information flow — is observed at the mesoscale network level."

(Miller et al. 2025, *Current Opinion in Behavioral Sciences*)

Organized travelling waves execute computation in the brain through so-called "slow" low frequency waves (which can rapidly impact electric fields)

Earl Miller has a **content** (gamma) vs. **control** (alpha/beta) model which is not dissimilar to ROSE



**A** Synchronous oscillatory neural activity creates localized **electric fields**

**B** Electric fields directly entrain neural activity locally via **ephaptic coupling**

△ silent neuron
▲ active neuron

Current Opinion in Behavioral Sciences

A lesson from the history of science, via Michael Levin [paraphrasing ]:

*If you keep interrogating deeper and deeper, getting more precise and causal-mechanistic, you eventually end up in the mathematics department.*

**This is where the language sciences and cognitive neurosciences should also be heading. But how can we get there?**

arXiv > cs > arXiv:2507.13501

**Computer Science > Computation and Language**

[Submitted on 17 Jul 2025]

**Encoding syntactic objects and Merge operations in function spaces**

Matilde Marcolli, Robert C. Berwick

We provide a mathematical argument showing that, given a representation of lexical items as functions (wavelets, for instance) in some function space, it is possible to construct a faithful representation of arbitrary syntactic objects in the same function space. This space can be endowed with a commutative non-associative semiring structure built using the second Renyi entropy. The resulting representation of syntactic objects is compatible with the magma structure. The resulting set of functions is an algebra over an operad, where the operations in the operad model circuits that transform the input wave forms into a combined output that encodes the syntactic structure. The action of Merge on workspaces is faithfully implemented as action on these circuits, through a coproduct and a Hopf algebra Markov chain. The results obtained here provide a constructive argument showing the theoretical possibility of a neurocomputational realization of the core computational structure of syntax. We also present a particular case of this general construction where this type of realization of Merge is implemented as a cross frequency phase synchronization on sinusoidal waves. This also shows that Merge can be expressed in terms of the successor function of a semiring, thus clarifying the well known observation of its similarities with the successor function of arithmetic.

Marcolli & Berwick (2025) focus on ROSE as a plausible candidate theory for syntax.

They write a proof showing that the mechanisms of ROSE (e.g., phase synchronization) can be mathematically connected to the algebraic properties of MERGE, providing hints towards an explanatory theory of neurolinguistics.

In the setting we are describing here, this would mean that when two waves $\varphi_{T_1}$ and $\varphi_{T_2}$ are merged in to the resulting wave $\varphi_{\mathfrak{M}(T_1,T_2)}$ according to the phase synchronization $\omega_{\mathfrak{M}(T_1,T_2)} = \omega_{T_1} \oplus_{\mathrm{Ry}_2,\beta} \omega_{T_2}$, as discussed above, if the head function $h_{\mathfrak{M}(T_1,T_2)}$ assigns to the root vertex $v$ of
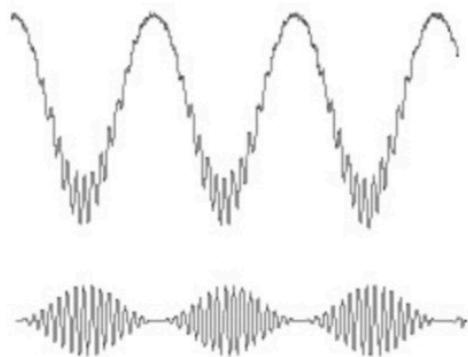
FIGURE 4. Phase-amplitude coupling in sinusoidal waves.

$\mathfrak{M}(T_1, T_2)$ the same head $h_{\mathfrak{M}(T_1,T_2)}(v) = h_{T_1}(v_1)$, where $v_1$ is the root vertex of $T_1$, then a lower frequency wave modulates the amplitude of the wave $\omega_T(\varphi_{T_1})$ in the resulting superposition wave $\varphi_{\mathfrak{M}(T_1,T_2)} = \omega_T(\varphi_{T_1}) + \omega_T(\varphi_{T_2})$.

Implementing this idea, however, has the problem that it leaves open the question of how the slower wave modulating the faster wave $\omega_T(\varphi_{T_1})$ that carries the head is to be determined: a mechanism for selecting such a modulating wave is not part of the type of model we have been discussing in this paper.
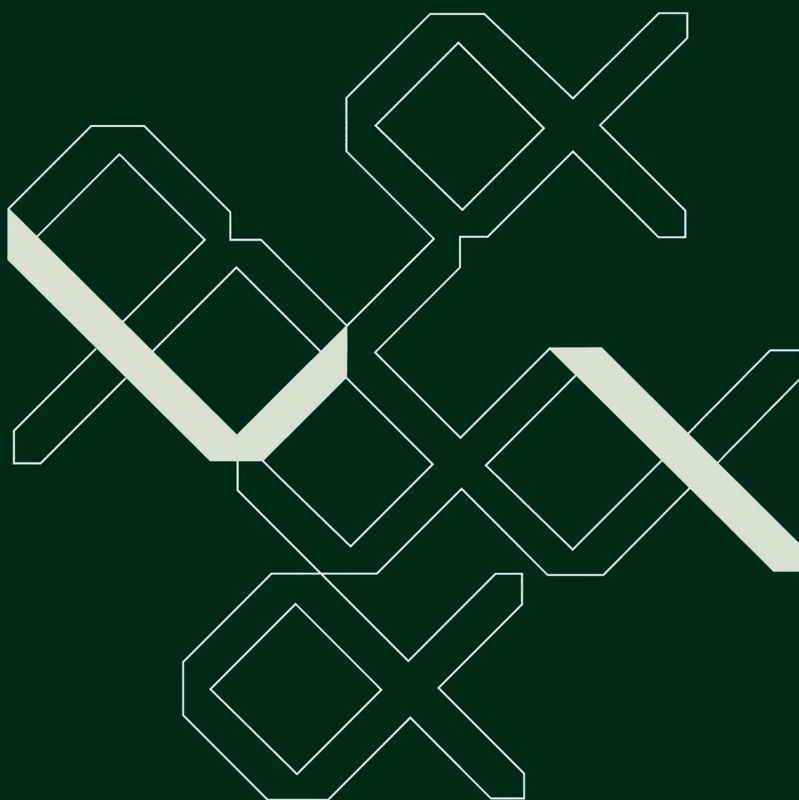
If mathematical linguistics can 'point' more easily to ROSE-compliant neural processes than to others, then this leads us towards an exciting terrain of previously inaccessible neurolinguistic research.

**Linguistic Inquiry Monograph Ninety**

**Mereological Syntax:**
Phrase Structure, Cyclicity, and Islands

David Adger

We can symbolize this way of thinking about Subjoin$(x, y)$ as follows:

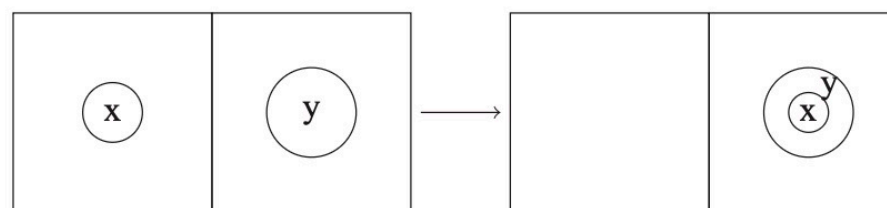(2) Subjoin$(x, y) \Rightarrow y : x < y$

Subjoin carries an object $y$ from one state, where it does not have $x$ as a part, to a new state, where it does.

The $<$ symbol is the "proper part" relation of mereology (Cotnoir and Varzi 2021). I'll use $<$ as the core relation, as opposed to using $\leq$ (which would allow $x$ to be part of $y$ if $x = y$). Proper parthood captures the (I think common) intuition that a syntactic constituent is not a constituent of itself, though, of course, issues of reflexivity of categories in structure have been a matter for debate over the decades in definitions of dominance and so on (e.g., Cushing 1978).

Visually, we can represent the derivation as follows:

(3)



The relevant properties of $<$ are usually specified via a set of axioms that ensure that the relation is:

(4) a. irreflexive ($\forall x. \neg x < x$),
    b. transitive (if $x < y$ and $y < z$, then $x < z$)
    c. asymmetric (if $x < y$, then $\neg y < x$)

Adger (2025), MIT Press

Adger uses part-whole relations via mereological formalism, whereas Marcolli & Berwick use category-theoretic magma (a more constrained type of set-formation).

Adger highlights the necessity of cyclic, local integration steps, while Marcolli & Berwick highlight the role of cross-frequency and synchronization mechanisms as not just correlative but necessary for building recursive structure.

Adger's part-whole 'Subjoin' operation speaks directly to the mereological nature of multiplexed PAC relations, while Marcolli & Berwick's framework is mathematically sympathetic to phase synchronization.

Different mathematical formalisms will be more/less directly translatable into processes at different scales of neural organization (*R, O, S, E*), hence will yield distinct experimental predictions for which types of neural signatures will drive syntactic inferences.

By ensuring that the 'algorithms' the brain uses respect the formal design features of human language, we help align experimental neuroscience with what we know (from first principles) about the nature of language.

For example, perhaps the intimate relation between lexico-syntactic and semantic processing in the brain's language network speaks more to Adger's thesis that lexical features are inherently 'part' of syntactic objects than it does other theories that posit a starker representational divergence between units and operations.

Unlike Merge, Adger's Subjoin operation doesn't create a nested set but a new object where both inputs are parts. No separate label object is required, since the composite's identity arises from its parts.

Mapping this model more closely to psycholinguistic variables would yield assumptions about the presence/absence of neural signatures pertaining to headedness that would differ from Marcolli & Berwick's model.

In contrast, Marcolli & Berwick propose Merge as a binary magma operation via category theory, formalizing this as addition in a semiring (with entropy-minimization defining combination). This operation can be implemented by a simple neural circuit (binary gate) computing a join of two input functions.

The use of a commutative, non-associative semiring for syntax suggests that any neural implementation must allow combining signals without order bias and without averaging away hierarchy – which naturally points to oscillatory binding (phase locking) rather than, say, simple additive firing-rate summation.

For Marcolli & Berwick, a plausible neural mechanism for Merge must allow reversible composition and decomposition (since Hopf algebras have inverses/co-operations). This points toward dynamic patterns like oscillations that can flexibly bind/unbind representations. A stable firing rate pattern summing two inputs might not easily be decomposed into the inputs again. By contrast, oscillatory phase coding can be more naturally parsed back into constituents.

Similarly, if we focus only on one specific frequency band or ERP component (e.g., P600), we likely blur parent-child relations in any hierarchical representation (how can compositional syntax emerge from high-frequency gamma activity in isolation? There are no accounts for how we can ground syntax in mono-causal signatures such as this).

Marcolli & Berwick prove that the algebraic structure of human syntax – centered on Merge as a free commutative, non-associative operation – can be faithfully embedded in a function space governed by thermodynamic semirings and optimized by Rényi entropy.

By modeling lexical items as wavelet-based functions and syntactic combinations as entropy-regularized additions, a novel blueprint for neurolinguists can be assembled.

The recursive, non-associative, and information-sensitive structure of Merge here mirrors properties observed in cortical oscillatory dynamics, such as cross-frequency coupling and spike-phase coordination (via ROSE).

These insights suggest specific empirical targets: for instance, cortical circuits that minimize local entropy under compositional constraints, or that exhibit nonlinear gain modulation aligned with Rényi-like cost functions during structure building.

Future intracranial or high-density MEG experiments can test whether syntactic processing engages such entropy-sensitive wave interactions, thereby grounding formal language theory in neurophysiological computation.

Ramping neural engagement for semantic composition has been well documented in high-frequency cortical recordings (Woolnough et al. 2023), but constraining the hypothesis space towards more specific mechanistic candidates for syntactic composition – as in $\mathfrak{M}(T_1, T_2)$ and the Hopf algebra Markov chain – remains a clear challenge for contemporary cognitive neuroscience.

Relatedly, if MERGE is to be seen as an operation involving minimization and entropy functionals there may be certain routes from the active inference world that might provide support here.

We would also expect Rényi entropy cost to be associated with general neural complexity metrics and oscillatory complexity. For instance, PAC modulation index (MI) and phase concentrations can be converted to discrete probability distributions, $p_i$, from which Rényi entropy can be computed. Specific parsing windows when MERGE is expected to occur should exhibit a reduction in Rényi entropy in high-gamma amplitude distribution (more structured, focused activation), or in PAC strength between relevant bands, depending on our preferred neurocomputational model for syntax.
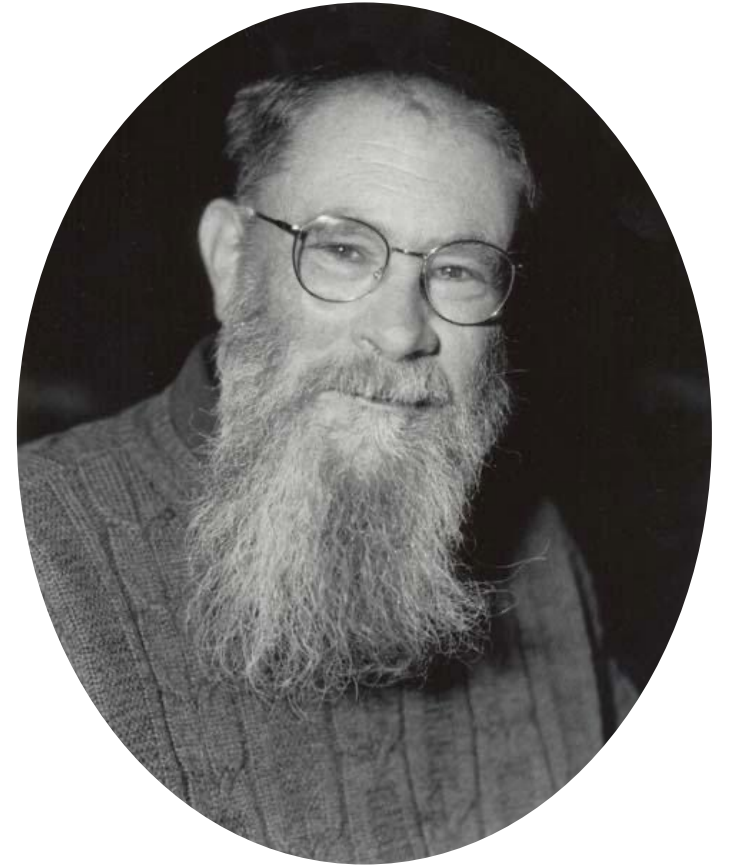
Psychological theory and theoretical linguistics seek to uncover what kind of neural machinery could carry the load that syntactic theory (MERGE-based syntax) says the brain must perform.

But the only way this will happen is if linguists and psychologists of language formalize their models of language knowledge/processing in more algebraically explicit ways. Keeping to naïve models of 'hierarchy' and basic graph-theoretic models of tree-structures will likely be insufficient for offering more acute experimental predictions for direct cortical recordings in the human brain.

The philosopher David Lewis famously said there are infinite possible causes for any event in the world.

There are surely many lower-order causal structures that subserve syntactic inferences, but ROSE places greater emphasis on mesoscale configurations as a way of re-framing where the "heart" of neurolinguistics should lie.

Scientists often say "causes should always produce their effects". This initially sounds reasonable, but is much too ambitious in reality for most scientific models in the life sciences.



David Lewis

A scientific theory is only as successful as the number of new interesting research questions it helps open up.

ROSE has been used productively to frame and explain a range of neuroimaging and EEG results (see Murphy 2025, Section 'ROSE as a plausible model for syntax').

Forthcoming work from multiple labs explores how ROSE can help researchers ask new questions about the neural implementation of second language acquisition, dependency resolution, and lexico-semantic processing.

ROSE "matures in accordance with a genetically determined biological matrix" (Murphy 2025).

Implying an *innate endowment* for the neural organization of syntactic computation.

Children use hierarchical knowledge *and* statistical learning jointly – S/E levels interface via cross-frequency coupling with probabilistic R/O levels.

ROSE's innateness claim rests on several converging empirical and theoretical points (1/2):

**Structure-dependence appears early in development:** infants already compute hierarchical rather than linear relations (Perkins & Lidz 2021, Shi et al. 2020).

**Neural evidence for endogenous syntactic organization:** cross-frequency coupling mechanisms implement recursive structure building and headedness in a way that is independent of learned statistical patterns, suggesting these are pre-specified neural motifs.

**Species specificity and developmental maturation:** the UNG claim posits that all human brains instantiate this multi-level oscillatory infrastructure for syntax, which *develops but is not learned*; similar to other neurobiological systems with constrained plasticity.

ROSE's innateness claim rests on several converging empirical and theoretical points (2/2):

**Unlearnability of Merge-based syntax:** no 'half-Merge' solutions.

**Poverty of the stimulus:** children uniformly converge on complex syntactic rules (e.g., auxiliary inversion, structure-dependent question formation) that are not derivable from surface-level statistics (and which LLMs fail or struggle significantly with).

**Neurolinguistic insights:** EEG and MEG studies show that the types of cross-frequency coupling invoked by ROSE tracks hierarchical phrase structure and syntactic closure in both adults and infants (Zhao et al. 2024; Weissbart & Martin 2024) in ways that are updated by the statistics of language but which are not reducible to statistical information.

Children as young as 7 months show sensitivity to the typical word order pattern of their parent's language ("But statistics cues this!" – Of course, but the principles of phrase structure are domain-specific – children do not invent the notion of head or structure)

Even by 18 months they show sensitivity to non-local syntactic dependencies.

See also the work of Charles Yang: an innate structural scaffold (Merge) interacts with learning biases ('Tolerance Principle'). When we add in constraints reflective of UG (e.g., assuming the existence of categories and phrase boundaries), statistical learning becomes far more efficient.

It is not "statistical learning vs. innate structure", but rather structural inferences via statistics.

See also Nicaraguan Sign Language - children rapidly and effortlessly acquire a recursive grammar in the absence of explicit instruction and decisive evidence.

There is a very strong obsession with "learning" in the fields of AI – partly driven by the nature of the enterprise (machine learning) but this has a problematic influence when it bleeds over into domains of psychology and cognitive science.

Not everything has to be 'learned'!

Carey (2023) argues that "there is no good evidence for nonlinguistic deductive reasoning involving the disjunctive syllogism", and that "animals and prelinguistic children probably do not make logical inferences".
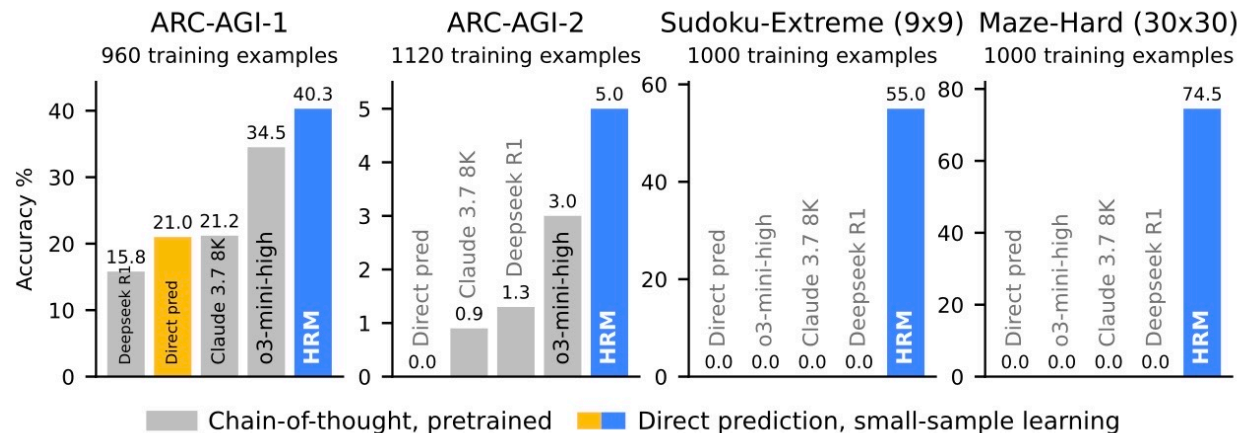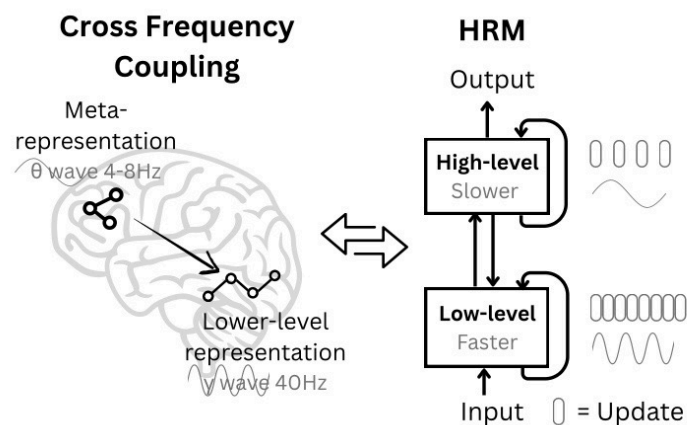
**Towards biological plausibility** (and away from "Transformers offer a candidate model for human language"…). Inspired by neural processes thought to be tied to compositional linguistic structures (e.g., PAC)
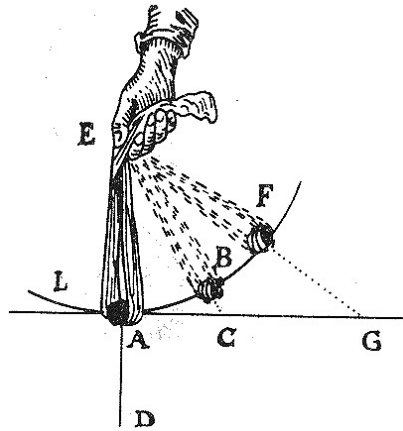


# Hierarchical Reasoning Model

Guan Wang[1,†], Jin Li[1], Yuhao Sun[1], Xing Chen[1], Changling Liu[1],
Yue Wu[1], Meng Lu[1,†], Sen Song[2,†], Yasin Abbasi Yadkori[1,†]

[1]Sapient Intelligence, Singapore

In the 17$^{th}$ century, the 'mechanical philosophy' sought explanations in terms of contact mechanics with deterministic interactions.

The imaginative space of cognitive neuroscience remains within pre-Newtonian "mechanistic" bounds, despite evidence for distinct causal landscapes.

**Distinct causal logics in the brain**

But there is no reason to remain reductionist about the causal structure of mental content. Downward causation is readily apparent in the brain, and "the neuron doctrine" died some years ago.

**For example:**

Increasing evidence that the central functional unit of the brain is not the cell. Cells contribute to function but do not causally anchor it.

Emerging evidence that types of behaviorally relevant information available at the LFP level are not represented in single units.

*J Physiol* 603.14 (2025) pp 4063–4090

4063

## Out of the single-neuron straitjacket: Neurons within assemblies change selectivity and their reconfiguration underlies dynamic coding

Fabrizio Londei[1,2], Francesco Ceccarelli[1,3], Giulia Arena[1,2,3], Lorenzo Ferrucci[1], Eleonora Russo[4], Emiliano Brunamonti[1] and Aldo Genovesio[3,5]

*J Neurophysiol* 104: 1768–1773, 2010.
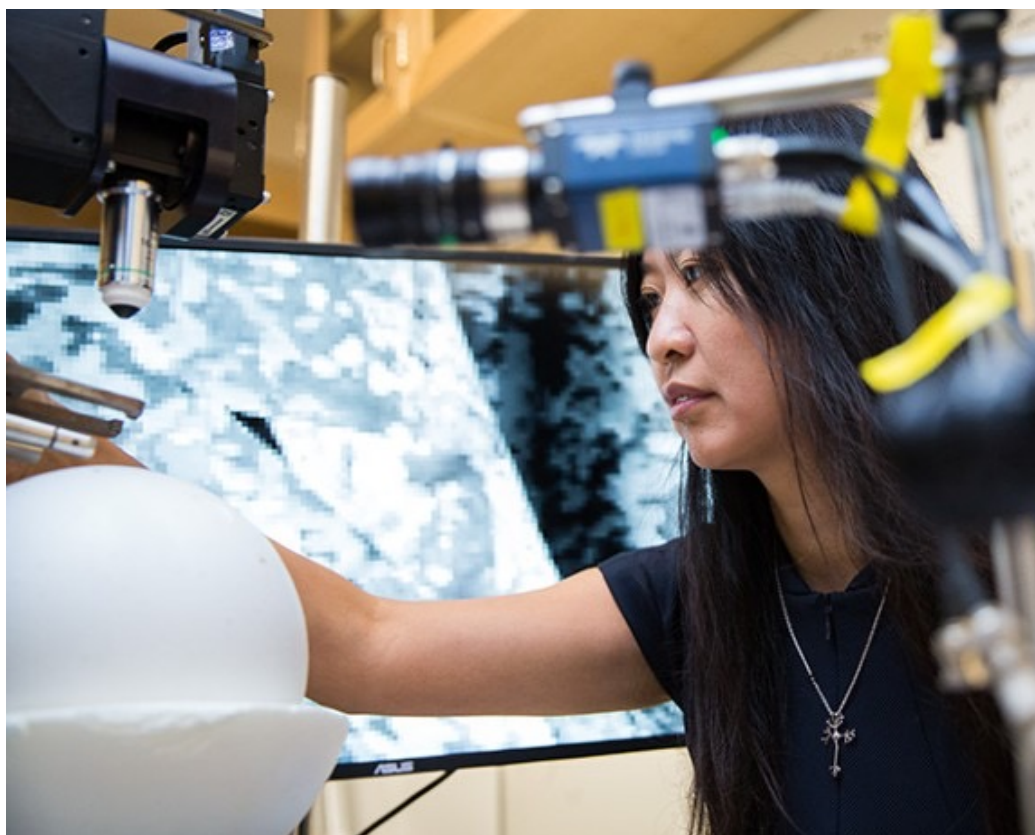First published July 21, 2010; doi:10.1152/jn.00478.2010.

## How Global Are Olfactory Bulb Oscillations?
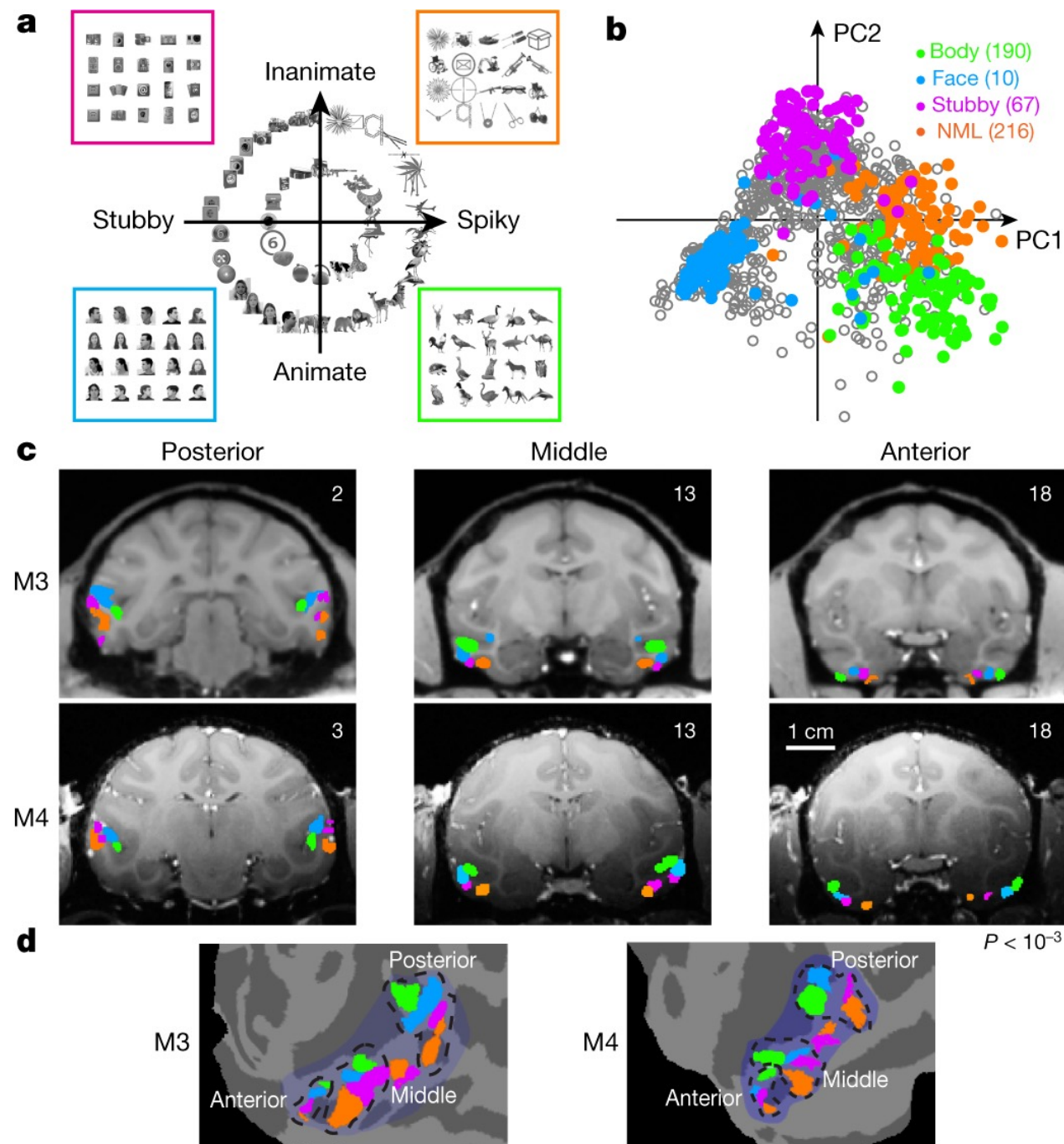
Leslie M. Kay[1,2] and Philip Lazzara[2,3]
[1]*Department of Psychology and* [2]*Institute for Mind and Biology, The University of Chicago, Chicago, Illinois; and* [3]*St. John's College, Annapolis, Maryland*

Although single neurons constitute the basic units of the nervous system, their impact on information processing is contingent on their interaction with the specific synaptic connection patterns of underlying neural circuits (Luo 2021, *Science*).

**RESEARCH**

**REVIEW SUMMARY**

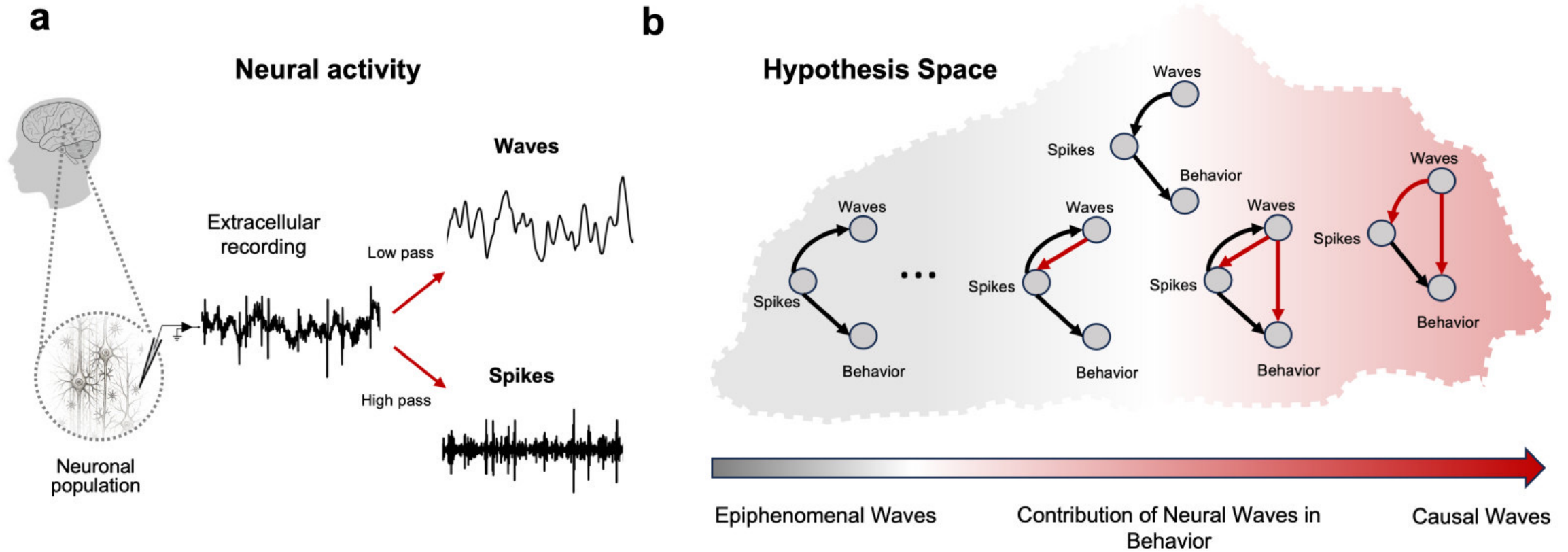**NEUROSCIENCE**

**Architectures of neuronal circuits**

Liqun Luo

**Doris Tsao**
**UC, Berkeley**

**a**

Inanimate

Stubby — Spiky

Animate

**b**

PC2

PC1

- Body (190)
- Face (10)
- Stubby (67)
- NML (216)

**c**

Posterior          Middle          Anterior

M3        2            13            18

M4        3            13            18

1 cm

$P < 10^{-3}$

**d**

M3          Posterior

Anterior        Middle
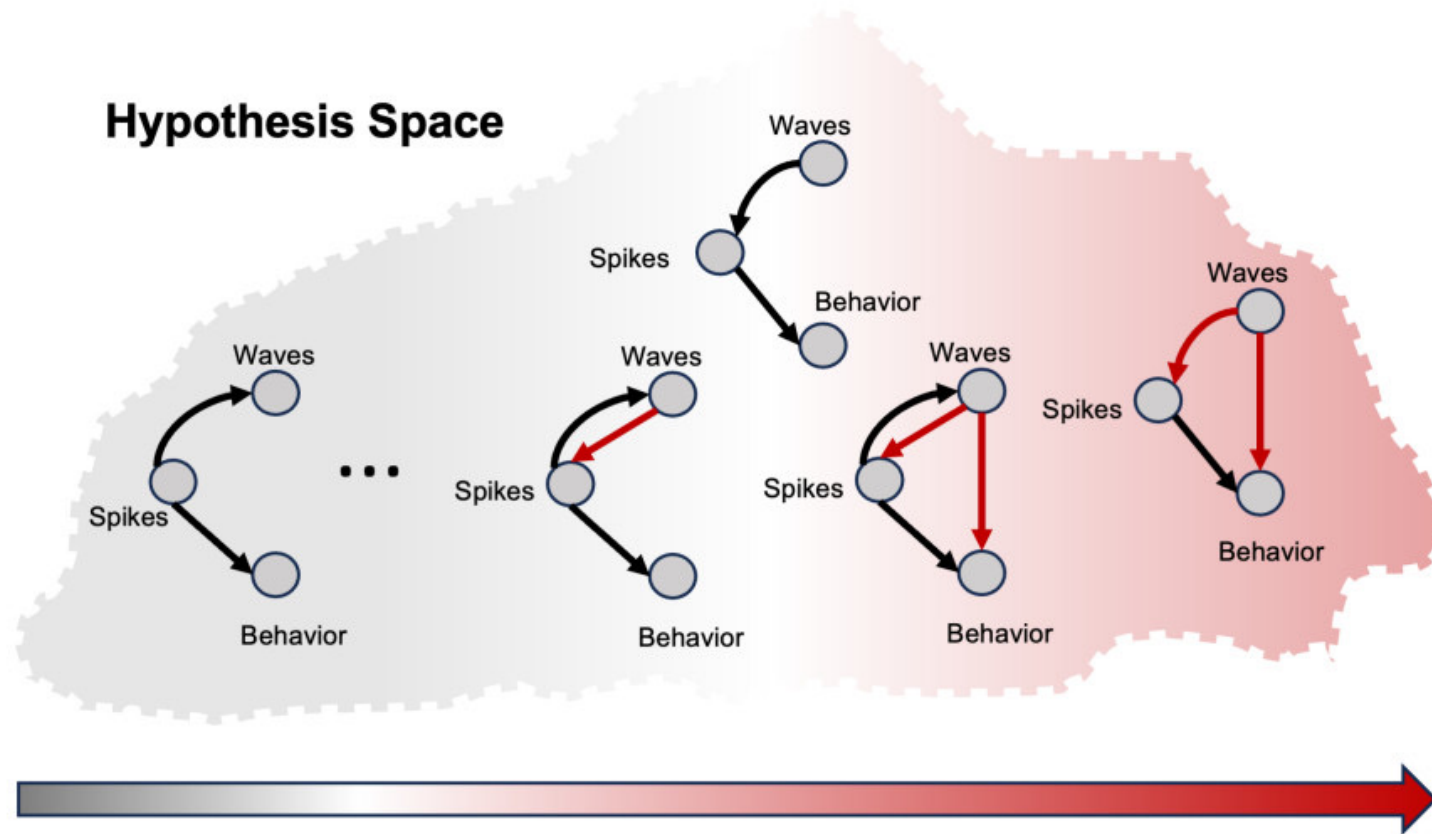
M4          Posterior

Anterior        Middle

Just as higher visual cortex may have inborn axes like "spiky vs. stubby" or "animate vs. inanimate" for categorizing objects, the language system might have pre-specified dimensions such as "predicate vs. non-predicate" or "concrete vs. abstract" that help organize concepts into syntactic roles.

What are the relevant axes for conceptual combination(s)? Event(IPL)-entity(ATL)? Predicate-non-predicate(MTG)?

Jalaldoust & Zabeh (2025), arXiv

**Hypothesis Space**

Epiphenomenal Waves          Contribution of Neural Waves in Behavior          Causal Waves

**(1)**
Randy Bruno
Kenneth Harris
Hickok
Kording

**(2)**
Buzsáki
O'Keefe

**(3)**
Lisman
Jensen

**(4)**
Murphy
Miller
Krishna Jayant
Poeppel

**(5)**
Fries
Tononi
Llinás

Many researchers think of causation exclusively through the metaphors of 'driving' a transfer of energy (e.g., Hume's billiard balls).

But even synaptic transmission itself does not actually involve a transfer of energy! The 'causal connection' is here really just *constraint satisfaction* or *signal transduction*.

Biological causation often involves enabling ('criterial causation') or constraining, not 'producing'.

## Physical bounds on cortical computation may yield explanatory power for neurolinguistics

Constraint 1: **Algorithmic Complexity**
Implementation cost

$$K_M(x) \stackrel{\text{def}}{=} \min\{|p| : M(p) = x\} \cup \{\infty\}$$

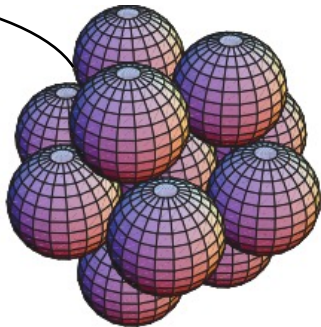Constraint 2: **von Neumann-Landauer Limit**
Physical energetic cost per bit

$$E \geq k_{\mathrm{B}} T \ln 2$$

Constraint 3: **Shannon-Hartley Theorem**
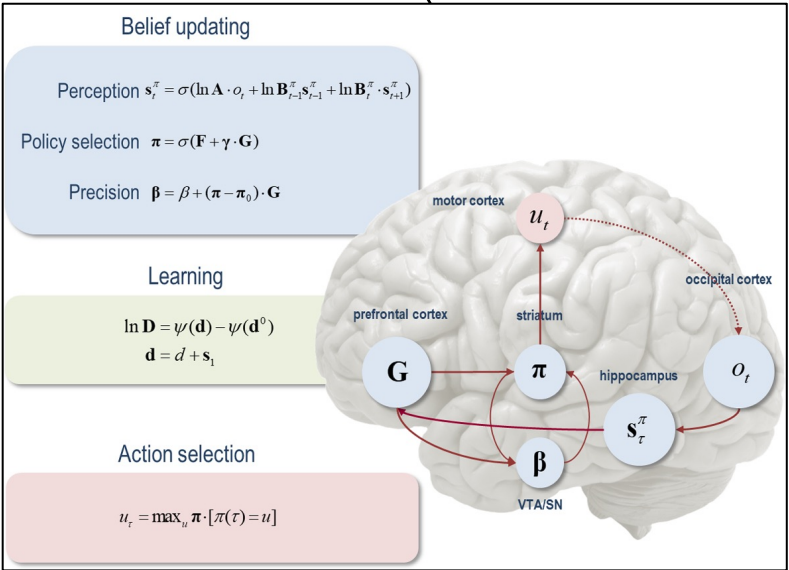Channel capacity of transmission without error

Constraint 4: **Bekenstein Bound**
Max information per unit volume

$$S \leq \frac{2\pi k R E}{\hbar c}$$

Constraint 5: **Sphere Packing**
Geometric laws of embeddings in $n$ dimensions

Constraint 6: **Bremermann's Limit**
Max rate of computation in an isolated physical system



**Free-energy principle**
Murphy, Friston & Holmes (2024), *Synthese*

**ORIGINAL RESEARCH**

# Natural language syntax complies with the free-energy principle

Elliot Murphy[1,2] iD · Emma Holmes[3,4] · Karl Friston[4]

Migrating structures from linguistics over to statistical physics

There are surely many lower-order causal structures that subserve syntactic inferences, but ROSE places greater emphasis on mesoscale configurations as a way of re-framing where the "heart" of neurolinguistic theory should lie.

This is a partly ontological and partly methodological intervention.

## Timescales of causality

Scientists are often tempted to assign "true" (or greater) causality to whichever causal factor occurs over the faster timescale, even when two causal factors may produce the same outcome with equal probability.

We may be eager to explain human language in terms of single-cell behavior but "true" causality may involve **downward causation** from mesoscale dynamics (as in ROSE).

We should return to Aristotle's thesis of *causal pluralism* that entertains multiple interacting causal forces, which has been sidelined in contemporary 'mechanism'-obsessed neuroscience.

e.g., Material, Formal, Efficient, and Final causes

Aristotle

Francis Bacon has some blame here. In the 17th century, he pushed the empiricist approach whereby only material and efficient causes were entertained, and other causal relations were considered 'metaphysical' or 'magic'.

Francis Bacon

It may seem *intuitive* that both brains and LLMs do very similar things. But our intuitions about the three main branches of philosophy also turn out to be wrong:

**Epistemology** ("No statement can be true and false at the same time") ❌
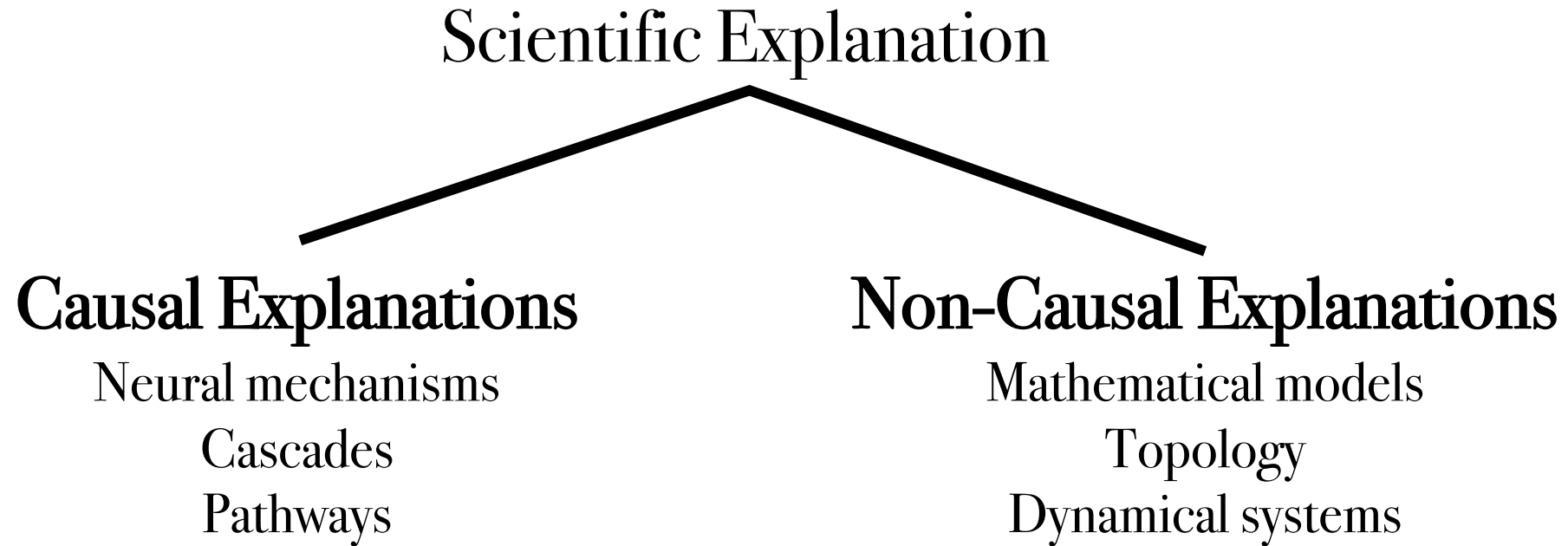**Ethics** ("Reducing pain and increasing happiness is the primary guiding principle of ethical life") ❌
**Metaphysics** ("Everything must have a sufficient cause or reason (Leibniz's *Principle of Sufficient Reason*)") ❌

"But *real* understanding of the brain requires lower-level causal details".

**Says who?**

"Our artificial neural language model offers a *biologically plausible* account of language and compositionality".

**So what are the concrete neurobiological predictions?**

Scientific Explanation

**Causal Explanations**

Neural mechanisms

Cascades

Pathways

**Non-Causal Explanations**

Mathematical models

Topology

Dynamical systems

Where do LLMs fit in here? Data-driven predictive models?…

Regression ≠ Explanation

Data ≠ Theory

# A call for "open theory"

Modern calls for "open science" and "open data" are valuable and important.

But we also need *open theory*! Researchers need to lay their inferential and metatheoretical cards on the table.

# Conclusion

Diverse causal landscapes in the brain derive distinct components of linguistic structure.

An explanatory neurolinguistics will require the causal language of linguistics and neuroscience to be in sympathy — this cannot happen if we model language processing as centered on sequential probabilistic statistics.

There is something unique about our species' evolutionary-ecological niche: We inhabit more than any other species a richly symbolic, logical and causal mental realm. We should not be surprised if sidelining these concepts in favor of statistical, functional and frequentist tools (e.g., prediction) offers little help in theory-formation for hierarchical linguistic structures.

A mathematically explicit navigation of linguistic knowledge (the 'Platonic forms' driving linguistic constituency structure) can help narrow down the list of candidate neural mechanisms for syntax.

# Future Directions

Explore how to align further these topics in experimental neurolinguistics (intracranial EEG, MEG, fMRI), theoretical neuroscience (ROSE model), and mathematical models of linguistic computation.

This can be achieved through (i) experimentally probing the parsing of linguistic structures of varying sizes (minimal compositional schemes, through to naturalistic sentences) in terms of how punctuated moments of *symbolic* inference occur alongside *statistical* processes; and (ii) exploring the processing of different types of categories of composition (i.e., generating geometric, melodic and mathematical structures with a 'minimal structure').

In the theoretical space, this can also be achieved through exploring which psychological and computational theories of language can be formalized into a mathematical language that is more amenable to mapping into certain neurocomputational regimes (via ROSE).

# Future Directions

Psychologists of language have the potential to help guide the search for the neural code for syntax, providing cognitive neuroscientists with an explicit algebraic 'parts list' to offer novel constraints for experimental testing.

It may be possible to use the formal, mathematical properties of language to help narrow the space of candidate neural mechanisms for how language is biologically implemented. This is an explanatory step that simply *cannot* translate into models of language that are purely 'boxological' and localizationist ("semantics is in X region, phonology in Y region").

Similar steps have already been made in psychological theories of working memory and attention, with respect to constraining neural theories – but the language sciences are yet to catch up…

# Thanks to the Tandon Lab and my collaborators

In recent years, the cognitive neurosciences have <span style="color:red">overwhelmingly embraced functionalist/statistical methods and models</span>, marginalizing algebraic/symbolic accounts of how Platonic forms may ingress into biology…



**Thomas Arnold**
(1795 – 1842)

"Take but one step in submission, and all the rest is easy. Satisfy yourself that you may honestly defend an unrighteous cause, and then you may go to the Bar, and become distinguished, and perhaps in the end sway the counsels of the State. All this is open to you; while if you refuse to tamper in a single point with the integrity of your conscience, isolation awaits you, and unhappy love, and the contempt of men; and amidst the general bustle of movement of the world you will be stricken with a kind of impotence, and your arm will seem to be paralysed, and there will be moments when you will almost doubt whether truth indeed exists, or, at least, whether it is fitted for man. Yet in your loneliness you will be visited by consolations which the world knows not of; and you will feel that, if renunciation has separated you from the men of your own generation, it has united you to the great company of just men throughout all past time; nay, that even now, there is a little band of Renunciants scattered over the world, of whom you are one, whose you are, and who are yours for ever."