# WHAT DOES IT MEAN TO UNDERSTAND LANGUAGE?

**Colton Casto**
Kempner Institute for the
Study of Natural & Artificial Intelligence
Harvard University
Cambridge, MA 02138
ccasto@mit.edu

**Anna Ivanova**
School of Psychology
Georgia Institute of Technology
Atlanta, GA 30332
a.ivanova@gatech.edu

**Evelina Fedorenko**
Department of Brain and Cognitive Sciences
Massachusetts Institute of Technology
Cambridge, MA 02139
evelina9@mit.edu

**Nancy Kanwisher**
Department of Brain and Cognitive Sciences
Massachusetts Institute of Technology
Cambridge, MA 02139
ngk@mit.edu

November 26, 2025

## ABSTRACT

Language understanding entails not just extracting the surface-level meaning of the linguistic input, but constructing rich mental models of the situation it describes. Here we propose that because processing within the brain's core language system is fundamentally limited, deeply understanding language requires exporting information from the language system to other brain regions that compute perceptual and motor representations, construct mental models, and store our world knowledge and autobiographical memories. We review the existing evidence for this hypothesis, and argue that recent progress in cognitive neuroscience provides both the conceptual foundation and the methods to directly test it, thus opening up a new strategy to reveal what it means, cognitively and neurally, to understand language.

# Introduction: A Cognitive Neuroscience Perspective on Language Understanding

What does it mean to understand language? Let's consider a few sentences:

- *"Killer whales exfoliate each other with home-made scrubbers."*
- *"Democracies facing authoritarian threats push back with mass demonstrations."*
- *"The man transfers the wet glue from one piece of wood to another."*
- *"More people have been to Russia than I have."*

Did you visually imagine the killer whales? (What were those scrubbers made of?)
Did you relate the authoritarian threats to a particular democracy (perhaps one close to home)?
Did you wonder exactly how the glue got transferred? (by scraping the glue off one piece of wood with a knife then smearing it on the other, or by pressing the two pieces together?)
Did you realize that the last sentence doesn't make any sense at all [1]?

What then does it mean to understand language, cognitively, computationally, and neurally? How rich and structured is your mental representation of the sentence you are reading right now, and what cognitive and neural systems are holding this representation? And are you now linking the meaning of this sentence to the previous one, to the unfolding structure of the argument in this article, and to thoughts you have had previously about language processing? In this opinion piece, we argue that a **deep understanding** of language, which entails building mental models and rich representations of meaning that connect to our broader knowledge of the world, requires the **exportation** of information from the brain's core language system to other cognitive and neural systems that can build models of what we are hearing or reading.

The idea that language understanding entails building mental models of events and situations described in language is not new. Decades of theorizing and behavioral experimentation have supported the idea that we construct **situation models** (see *Glossary*) when processing language that diagram in our minds the characters described, their goals and spatial locations, and the temporal and causal relations among the events taking place [2-9]. However, the set of representations and computations that might be recruited during language understanding is large and heterogeneous, from reasoning about the contents of someone else's mind, to retrieving our knowledge of typical event sequences, to activating perceptual/motor representations from vivid descriptions. Although this long tradition of behavioral work has provided many clues about the representations extracted during language processing, the precise nature of these representations, and the conditions under which they are constructed during language understanding, has remained elusive.

Cognitive neuroscience has much to offer to these questions, by providing both a rich taxonomy of mental processes that selectively engage specific brain regions, and the methods to directly measure that recruitment during language understanding. In particular, neuroscientists have leveraged fMRI and related methods to identify and functionally characterize both the **core language system** itself, and many other cognitive and neural systems that it may export information to for further processing (**Figure 1A**). A key finding from this research is that the core language system is specifically engaged in understanding (and producing) language, and not in other nonlinguistic cognitive processes (see [10] for review). Further, many other cortical regions show remarkable specificity for other mental functions, including regions in the ventral visual pathway selectively engaged in perception of faces, bodies, and scenes [11-13], regions of auditory cortex selective for the perception of speech sounds and music [14-15], the right temporo-parietal junction (rTPJ), which is selectively engaged in thinking about what others are thinking [16], and the other regions comprising the theory of mind network [17-18], and multiple regions engaged in visually-guided navigation [19] and in reaching and grasping [20]. Ongoing research is investigating the precise function of many other brain regions, including those engaged in intuitive physical reasoning [21], the perception of social agents and their actions [22], amodal semantic processing [23], episodic retrieval and self-projection [24,17], as the domain-general processing characteristic of the 'Multiple-Demand' system [25] that is engaged by a broad range of cognitively demanding tasks.
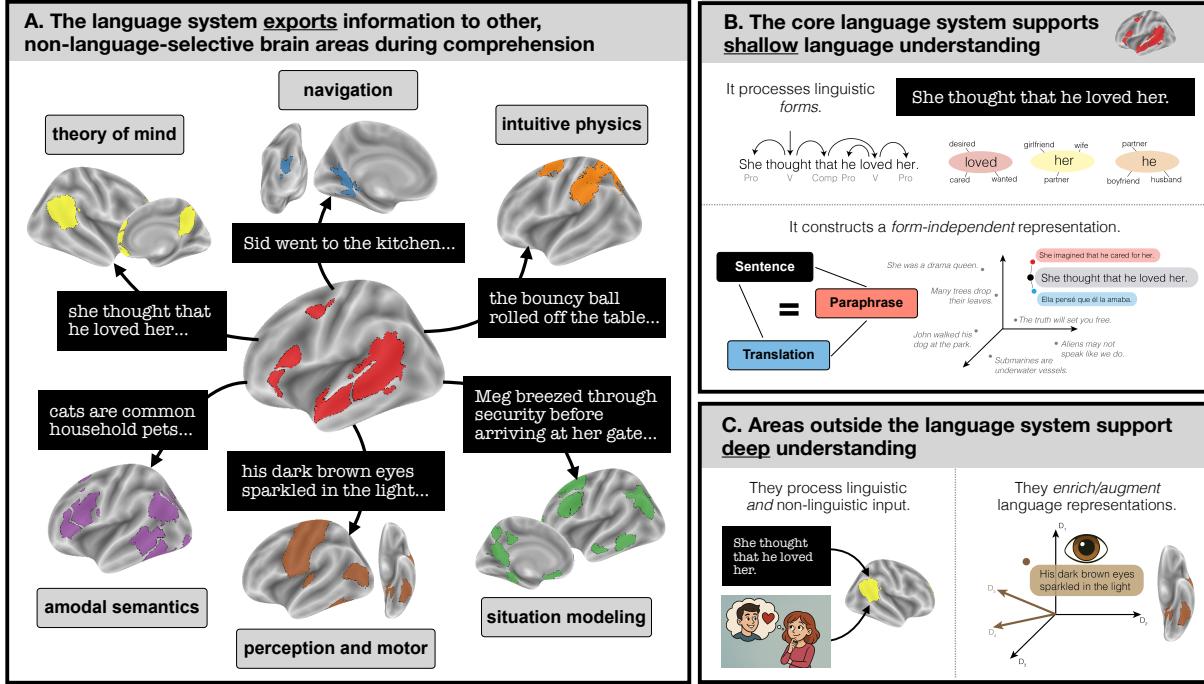
**A. The language system <u>exports</u> information to other, non-language-selective brain areas during comprehension**

navigation

theory of mind

intuitive physics

Sid went to the kitchen...

she thought that he loved her...

the bouncy ball rolled off the table...

cats are common household pets...

Meg breezed through security before arriving at her gate...

his dark brown eyes sparkled in the light...

amodal semantics

perception and motor

situation modeling

**B. The core language system supports <u>shallow</u> language understanding**

It processes linguistic *forms*.

She thought that he loved her.

She thought that he loved her.
Pro V Comp Pro V Pro

desired
girlfriend wife partner
loved her he
cared wanted partner boyfriend husband

It constructs a *form-independent* representation.

Sentence
=
Paraphrase
Translation

She was a drama queen.
Many trees drop their leaves.
John walked his dog at the park.

She imagined that he cared for her.
She thought that he loved her.
Ella pensé que él la amaba.
The truth will set you free.
Aliens may not speak like we do.
Submarines are underwater vessels.

**C. Areas outside the language system support <u>deep</u> understanding**

They process linguistic *and* non-linguistic input.

She thought that he loved her.

They *enrich/augment* language representations.

His dark brown eyes sparkled in the light

Figure 1: **Hypothesis: Exportation of Information from the Core Language System to Other Brain Regions Is Necessary for Deep Understanding. A)** The language system exports information to extra-linguistic systems that support deep understanding of linguistic content. Candidate exportation destinations include regions specialized for thinking about the contents of someone else's or one's own mind [16], regions specialized for scene understanding and navigation [19], regions that construct mental models of the physical world [21], regions that represent and process amodal semantic information [23], perceptual and motor regions (e.g., [11-13]), and regions that integrate information into a coherent representation of the situation that is being described [95-96]. These exportation destination regions are visualized with broad masks delineating where the relevant functional regions are found; although these masks may overlap, the functional regions within individuals are distinct. **B)** The language system supports shallow language understanding. Shallow language understanding entails recognizing individual words and how they go together within the structure of the sentence (i.e., parsing the surface form of the sentence), and then constructing a compositional representation of the sentence's meaning, perhaps represented in a high-dimensional vector space. Importantly, these shallow representations are only defined with respect to other linguistic elements and patterns of their use that we have encountered during our lifetime of experience with language, making no contact with other aspects of our lived experience. **C)** Extra-linguistic systems support deep language understanding. They can be engaged both by linguistic and nonlinguistic inputs, which is why we do not consider them part of the language system. During language understanding these regions may augment the shallow representations constructed by the core language system with perceptual or motor features, translate them into a format compatible with probabilistic programing engines or symbolic problem solvers, or transform them in some other way.

This emerging picture of the functional architecture of the human brain [26] empowers us to ask three key questions. First, how deeply is information processed within the language system itself (Section 1)? Second, what information is exported from our core language system, and where does it go (Section 2)? And third, what factors influence whether and when information is exported (Section 3)?

Before tackling these questions, we make two important clarifications. First, our hypothesis is distinct from the claim that the whole brain—not a single, specialized network—is responsible for language processing (**Box 1**). Brain regions that, we hypothesize, language information is exported to are not language-selective: they can be strongly engaged by nonlinguistic inputs (**Figure 1C**). These systems are not part of the language system simply because they can be accessed through language, any more than they are part of the visual system because they can be engaged through vision. The core language system, however, is selective for language only. Second, exportation obviously has to happen when linguistic inputs call for an overt response or behavior, such as when we are asked a question or given an instruction (e.g., *What is the appropriate tip for an eighty-dollar restaurant bill? What is the most efficient route from*

*work to home? Imagine your mother's face.*). Here, we argue that exportation can also happen when language simply provides a description during passive reading or when listening to a story.

---

**Box 1: A unifying account of language processing in the brain.**

Many past studies have reported neural responses in cortical regions outside the classic language areas when people process words, sentences, or entire stories (e.g., [68,72-73,79-80,86,94,23,96,101,104]). Some have interpreted these findings as evidence that no single part of the brain is specialized for language processing but rather "the entire brain is at work" (e.g., [123-124]). We argue, in contrast, for two fundamental differences between the contributions of the core language system and the many other brain regions that can be engaged during language understanding. First, all linguistic input (spoken, written, or signed) is necessarily processed by the core language system initially. Further processing by downstream systems occurs only after an initial representation has been constructed by the core language system. Second, no areas outside the core language system are language-specific: although they can be engaged by linguistic inputs, they can also be engaged by visual and other non-verbal inputs (**Figure 1C**), which is why we do not consider these regions part of the language system. This perspective reconciles the sometimes-broad activation of many brain regions in response to language with decades of neuroimaging studies demonstrating the functional specificity of many of these same regions [10-26]).

To illustrate the power of our framework in explaining diverse findings, consider a recent study by Singh, Antonello et al. [73]. These authors used a large language model (LLM) to rate a language input on several dimensions intended to capture qualitative scientific theories about the function of various brain regions (e.g., Does the input involve spatial reasoning? Does the input describe a person or social interaction?). They then used these ratings to predict responses across the brain while participants listened to naturalistic stories. Although the study did not localize functional regions within individual participants (see Antonello, Singh et al. [72] for a similar study from the same authors that did), at a coarse level, their results appear to be largely consistent with the known functional organization of the human brain. For example, the authors found that brain responses near the parahippocampal place area (PPA) and retrosplenial cortex (RSC) were most strongly predicted by ratings of whether the input mentioned a specific location, and that responses near the right temporo-parietal junction (rTPJ) were most strongly predicted by whether the input described a social interaction [73]. These and other similar findings (e.g., [72]) support our hypothesis that language understanding often entails exporting information to brain regions outside the core language system that are functionally specific for the content described in the linguistic input.

---

## 1. How much of language understanding occurs within the core language system?

fMRI research over the last 15 years has identified and characterized the brain's core language system with new precision (see [10] for review). Consisting of temporal and frontal left-hemisphere regions, this system responds when we listen to or read language [27-28] or when we speak [29], but not when we engage in apparently similar but nonlinguistic processes, such as arithmetic, music perception, understanding gestures, or reasoning [30-32] (for reviews, see [10,33]). Damage to these regions leads to linguistic deficits [34]. All regions of the language network show sensitivity to structure at the sublexical, word, and sentence level [35-38], in line with usage-based accounts of language processing (e.g., [39]), which draw no sharp distinctions by unit size (sounds, words, sentences).

During comprehension, the language system's goal is to extract information from word sequences, which requires recognizing familiar words and larger constructions and figuring out how they go together within the sentence (**Figure 1B**). The end product is a representation that captures both structural and co-occurrence-based information encoded in a sentence. Several lines of evidence suggest that this representation is quite *abstract*, no longer tied to the particular words. For example, representations appear to be similar across sentence paraphrases and across languages [40-41] (**Figure 1B**). Furthermore, the representations constructed within the language system share some features with representations extracted from meaningful visual inputs: pictures and silent videos elicit some response in the language areas, although substantially lower than verbal stimuli [42], and representations from neural network vision models can be used to predict responses to sentences in these areas [43-44]. Presumably, this similarity between language- and vision-derived representations stems from the fact that both language and visual experiences capture the structure of the world [45-46].

Critically, although abstract, we argue that these meaning-approximating representations are shallow: they are derived solely from our knowledge of language statistics, akin to representations in early versions of text-only large language models (**Box 2**). Even severe damage to these brain areas can leave conceptual processing intact, suggesting that whatever representations and computations they support cannot be critical to world knowledge and reasoning [42,47]. Furthermore, the representations in the language system lack sensitivity to real-world plausibility, as evidenced by similarly strong responses to plausible sentences such as *"Sally went to the park to walk her dog"* and well-formed but nonsensical sentences like the famous example *"Colorless green ideas sleep furiously"* [48] (**Figure 1B**). These nonsensical sentences abide by the rules of English syntax, but do not express a plausible meaning; the fact that the language areas respond to them just as strongly as to plausible sentences suggests that the interpretations constructed within the language system are not constrained by our knowledge of the world. This feature is of course "by design": although language statistics capture many aspects of our world knowledge (e.g., [49]), we can use language to tell each other anything, including things that are false, nonsensical, and surprising. (*Did you know that armadillos are actually bulletproof, or that owls don't have eyeballs?*). How would we learn new things about the world through language, if the system were not flexible enough to build representations for concepts and ideas that we have not encountered before?

Taken in tandem, the available evidence suggests that the core language system is only capable of **shallow understanding**: it doesn't really "understand" things, it only knows how people talk about them. This is not to say that the representations constructed by the language system are not useful or powerful. As noted above, language statistics reflect many aspects of the structure of the world, and—as we have learned from LLMs—training on language statistics alone is sufficient to produce grammatically well-formed and semantically coherent discourse, answer questions, and carry out simple conversations [50-51]. However, when we read or listen to language, our minds often do more than extract these shallow approximations of meaning. We construct mental models of what the language is describing and rich representations of meaning that connect to our broader knowledge of the world, which suggests that the language system is often not the end point of linguistic processing (see [10] and [52] for related discussions), much like the visual system is not the end point of visual processing.

---

**Box 2: Shallow vs. deep understanding in large language models (LLMs).**

The question of whether LLMs trained on text alone really "understand" language has been raging in the field of AI [125-126]. In a recent piece [52], we introduced a distinction between formal and functional linguistic competence, where formal competence refers to the system's ability to use linguistic rules and patterns and corresponds to computations in the brain's language system, and human-like functional competence refers to the ability to use language to do things in the world, which we argued requires additional, extra-linguistic systems. This formal-functional distinction mirrors the distinction between shallow and deep language understanding that we have focused on here, except we here focus on comprehension for the sake of comprehension, not to achieve some external goal, such as answering a question or solving a problem, which we have emphasized in our discussions of functional competence.

Early models like GPT-2 [50] primarily show signatures of shallow understanding, with limited functional competence [52], just like the brain's core language system. The more powerful the model, the better it becomes at diverse functional competence tasks [51,127], and the better it predicts responses across the brain to stories [128-129]. Importantly, however, the ability of these models to predict brain responses in the core language system saturates as soon as formal competence plateaus [130-131]. Thus, relatively small LLMs by today's standards (capable of primarily shallow language understanding) can capture the majority of variance in the core language system's responses [131], consistent with our view that most processing in the language network is shallow.

Do any AI systems understand language deeply? If so, then we might expect them to show signatures of exportation. Several recent studies have attempted to explicitly build exportation into AI systems by augmenting LLMs with diverse extralinguistic systems: vision models, physics engines, formal logic provers, theory of mind engines, memory augmentations, and many more (e.g., [55-57,132]). Other studies have explored architectures that allow for functional specializations to be learned over the course of training, rather than having to be defined *a priori* through explicit augmentations (e.g., [133]). There is also some evidence that individual units within a standard

LLM—without any additional augmentations or architectural changes—can become functionally specialized [134-135]; however, the causal importance of these functionally-specific units remains debated (e.g., [136]), leaving open the question of whether classic LLMs can develop brain-like exportation. AI systems which do have exportation (whether neuro-symbolic or neural-net-only) can serve as useful model organisms (e.g., [116]) for human-like language understanding and may shed light on some of the **Outstanding questions**, such as which computations underlie information transfer during exportation.

## 2. To what extent does language understanding engage regions outside the core language system?

If the core language system does not fully process the meaning of a linguistic input, what other systems do? In this section, we discuss some of the brain regions that may support deep language understanding (**Figure 1C**). Under our framework, these brain systems can, for example, augment the shallow representations constructed by the core language system with perceptual or motor features [53-54], translate them into a format compatible with probabilistic programing engines [55-56] or symbolic problem solvers [57], or transform them in some other way. Whether these downstream representations are continuous [58] or symbolic and propositional [59] is outside the scope of this article, although we highlight that all of our ideas are compatible with a distributional framework [58]. Below we summarize available evidence of exportation with an emphasis on studies that satisfy our methodological desiderata (**Box 3**); for many brain areas, however, direct evidence of exportation during passive language understanding is limited, so we focus on predictions of the exportation hypothesis.

---

**Box 3: What would count as evidence of extra-linguistic engagement?**

What would strong evidence for the exportation hypothesis look like? In addition to observing an above-baseline response to language in a non-language system, we advocate three criteria:

1. ***Interpreting responses is easiest when the brain region in question is functionally specific.*** Inferences of the engagement of specific mental processes from activation in particular anatomical locations (or "reverse inference" [137]) are strongest for brain areas that are functionally specific. For example, engagement of the FFA during language understanding would provide strong evidence that a face percept has been constructed, whereas engagement of anterior cingulate cortex would be consistent with numerous disparate cognitive processes.

2. ***The brain region should be functionally identified within individual participants.*** Because the anatomical location of functional regions varies across participants, standard group analyses blur the responses of each region with its anatomical neighbors. Thus, a precise measure of the engagement of a given functionally-specific cortical region requires first identifying that region within individual brains with a standard 'localizer' task. For example, if we want to test the role of the PPA in understanding navigation-relevant language content, we should identify it in each participant with a contrast of viewing scenes versus viewing objects. For many systems—such as the language network—one can also identify them using functional correlations derived from resting state (e.g., [138]) or task-based (e.g., [139]) timeseries.

3. ***The language input should not instruct the comprehender to do a task.*** To address the question of whether and when information is exported from the language system in the course of natural, everyday language understanding, we advocate using passive comprehension and not paradigms that instruct the participant to engage in imagery (e.g., *imagine your mother's face*) or answer questions (e.g., *how red is a carrot?* or *how do you grasp a pencil?*), as these trivially require exportation.

---

### a. Brain Regions that Build Mental Models of Minds, Objects, and Places

**Theory of mind.** There is at least one case where the evidence is clear that information is exported from the language system to be further processed elsewhere in the course of everyday language understanding: theory of mind (ToM), the

ability to construct a mental model of someone else's mind. Extensive studies by Saxe and colleagues have demonstrated the highly specific fMRI response of an area in the right temporo-parietal junction (rTPJ) when participants think about the content of other people's thoughts [16,60-61]. This region is functionally correlated with a set of other regions (e.g., [60-61]) that together constitute the ToM network [16-18].

In many studies of the ToM network, participants are asked to read a short passage and then choose the correct completion of the story, which requires inferring the beliefs of a character in the passage. The ToM network responds strongly during this task (and the response is locked to the part of the passage that talks about mental states), but not during a logically similar task in which the passage requires an inference about a physical representation [16]. In contrast, the mental and physical inference passages engage the language system to a similar degree (when controlling for linguistic confounds [62]). The ToM network responds even if no question is asked and the participant need only read and understand the text (e.g., [16,63]), and it is also engaged when viewing silent movies that require thinking about others' mental states [64-65]. Finally, language and ToM areas are partially synchronized during story comprehension (e.g., [65]), which provides indirect evidence of the exportation process as well. Theory of mind is therefore a case where a full understanding of information delivered through language entails exporting representations from the brain's language system to another functionally-specific region, the rTPJ, along with the broader ToM network.

**Intuitive physics.**    Just as we build mental models of each other's minds, we also build mental models of the physical world to understand what we are looking at, predict what will happen next, and plan actions [66]. Growing evidence implicates regions of the parietal and frontal cortex in this intuitive physical reasoning (aka the "Physics Network" [21,67]). Although these regions have been primarily studied using nonverbal visual scenarios, they have also been shown to respond more strongly when solving verbal problems about physical (compared to social) scenarios [68-69], suggesting that information from the language system can be exported to the Physics Network for further processing.

**Navigation and scene understanding.**    Another kind of mental model concerns the map of space we use to navigate around in the world. Extensive work in human cognitive neuroscience has identified key brain regions underlying our navigational ability (see [19] for review), including the parahippocampal place area (PPA, [12]), the occipital place area (OPA, [70]), and retrosplenial cortex (RSC, or "medial place area", MPA [12,71]). Consistent with the exportation hypothesis, preliminary evidence indicates that these regions can be engaged by understanding paragraphs that contain information about places and navigation [72-73] (see also **Box 1**). We would more specifically predict that during language understanding the PPA would become engaged when reading a vivid description of the shape of a room [12], the OPA for a description of an open doorway [74-75], and the RSC when listening to someone's experience of regaining their bearings after emerging from a subway [76-77].

### b. Brain Regions that House Perceptual, Motor, and Emotional Representations

**Perception.**    An engaging novel or story can conjure up vivid mental imagery. According to the exportation hypothesis, this is possible due to exporting information to corresponding brain areas. Consider visual processing. Some of the strongest evidence for functional specificity in the cortex comes from studies of mid-level and high-level vision, where cortical regions respond selectively to mid-level visual features like motion and to high-level categories like faces, places, and bodies. It is well established that these regions can be activated when participants are explicitly instructed to mentally imagine faces or places [71] or motion [78]. But it is not known whether the functionally-specific regions of the ventral pathway are engaged when simply understanding verbal descriptions of visual scenes. Would your FFA fire up for you to really understand a passage like this: *"He was not over thirty. His eyes were very dark brown and there was a hint of brown pigment in his eyeballs. His cheek bones were high and wide, and strong deep lines cut down his cheeks, in curves beside his mouth. His upper lip was long, and since his teeth protruded, the lips stretched to cover them, for this man kept his lips closed."* (Grapes of Wrath, Steinbeck).

**Motor.**    Early work on "embodied" representations argued that motor representations are necessarily engaged for understanding the meanings of action words. For example, Hauk et al. [54] famously reported that passive reading of words describing actions of the face ("lick"), arm ("pick") or leg ("kick") activated regions of motor cortex overlapping with or adjacent to those activated by corresponding actual movements (see also [79]). Although these early studies faced several methodological limitations (see **Box 3** for how these questions can be tested more rigorously), a more recent intracranial study found cross-decoding of neural responses from action words to videos depicting those actions [80], more strongly supporting the exportation of action-related meanings. Although the general idea behind embodied cognition [53-54,81] aligns with the view we propose here, we broaden it by proposing that virtually any brain region, including those that represent abstract content about other minds and the physical world, may become engaged when understanding language content about those domains.

**Emotion.**    Reading a novel can make us laugh or cry or make our hearts beat with fear. Multiple studies have shown activation of brain regions implicated in emotional processing by single words [82-84] or sentences/passages [63-64,85-86]. Jacoby et al. [64] found that reading passages describing a protagonist's emotional pain produces activations that substantially overlap with the ToM network. Hsu et al. [86] reported that the emotional valence of the words in short passages from Harry Potter modulated responses in the amygdala. These findings provide suggestive evidence for the exportation hypothesis although i) current evidence suggests that few if any brain regions are selectively engaged by emotional processing (e.g., [87]), making reverse inferences from those regions challenging (**Box 3**), and ii) many prior studies—especially those looking at emotional processing of single words—include an explicit task (e.g., [83]; **Box 3**).

### c. Brain Regions that Represent Episodic and Semantic Knowledge

A rich understanding of language often requires having access to information that may not be represented within the language system itself but must be retrieved from memory. Memory is classically thought to have two main forms: episodic and semantic [88]. Episodic memory represents personal memories of past events and has been linked to a distributed network of regions in the frontal, parietal, and temporal lobes ("Default Network A", DN-A [24,17]; see [89] for a review). The exportation hypothesis predicts that this same network will be engaged when a friend recounts a shared experience, or when we read a novel and use our own lived experience to better understand what the protagonist is going through. Semantic memory instead consists of our general conceptual knowledge of objects, people, places, and facts about the world [90-92]. Acquired through a combination of lived experience and learning through language, semantic knowledge is not tied to a particular place or time, nor a particular input modality (e.g., language or vision). Although where and how semantic knowledge is represented in the brain remains debated [23,90-91,93-94], everyday language routinely taps into our rich conceptual knowledge, so it seems plausible that these amodal semantic regions would get engaged when simply understanding language [23].

### d. Brain Regions that Construct Situation Models

Many argue that the goal of language understanding is the construction of mental "situation models" that integrate all the above information—domain-specific mental models, memories and world knowledge, and aspects of perceptual and motor representations—into a single coherent representation of what we are hearing or reading [2-4,8]. Where and how does this integration happen? Some studies have implicated the "default network" (DN) [95-96]. Unlike the language system, the DN has the computational capacity to integrate information over long contexts for both verbal and visual inputs [97-101], and some studies have argued that areas of this network track specific aspects of the story (e.g., changes in the spatial location or the characters [101-104]).

However, most of these studies were carried out before clear evidence emerged against this network being monolithic, and instead consisting of the Theory of Mind network and Default Network A (DN-A) [24,17], both of which we have introduced in earlier sections. These networks have similar overall anatomy but are clearly dissociable within individual

participants. Moreover, although the DN-A was originally implicated in episodic projection (remembering the past and imagining the future [118]), an alternative account argues for a domain-specific role for DN-A in representing and navigating in space [18]. Some of the early findings implicating the DN in situation model construction should therefore be revisited to more critically assess the nature of its contributions to language understanding, including testing the possibility that the situation models we construct from language arise from a dynamic—and possibly hippocampal-dependent [105-107]—coordination of representations constructed by multiple domain-specific systems, rather than a single brain network integrating information across domains.


## 3. When does exportation happen, and how might it be implemented?

In addition to discovering the full ontology of extra-linguistic mechanisms that may get engaged during language comprehension, another important research goal is to understand what factors affect the likelihood of information exportation. We briefly highlight three possible factors. First, because the language system is memory-limited and only processes short segments of language [100,108-109], input that extends beyond a sentence or so *must* be passed to other systems that can construct temporally extended contexts (such as the DN-A and the Theory of Mind networks [100]). By exporting constructed meanings to other systems, the language system frees up space for new incoming input [110]. Traits, states, and goals of the comprehender may matter as well, including how much they know about the topic (e.g., [111]), how alert they are (e.g., [112]), how able they are to engage in perceptual imagery (e.g., [113]), and how interested they are in the content (e.g., [114]). Imagine, for example, listening to a doctor explain treatments for a particular condition before vs. after you've been diagnosed with that condition: the likelihood of exportation in the latter case is much higher! Indeed, previous behavioral work has shown that people employ different types of mental simulation (e.g., perceptual, motor, mentalizing) to differing degrees during natural reading ([115]), which may be attributable, in part, to individuals having different goals when reading the same text. Third, exportation critically depends on the language system being able to extract from the surface form a representation that is abstract enough to interface with downstream systems. Consequently, the difficulty of extracting such representations due to the complexity of the form (e.g., infrequent words or complex constructions) or insufficient proficiency of the comprehender (e.g., a non-native speaker) may make exportation less likely.

When exportation does take place, how is this process implemented? One possibility is that the language system preferentially sends information to the most likely target destination (i.e., "**routing**"). A different possibility is that representations constructed by our language system are sent to all possible downstream systems all the time (i.e., "**broadcasting**"). However, because the methods currently available in cognitive neuroscience make it challenging to adjudicate between these possibilities, making progress towards a mechanistic understanding of the exportation process will likely require *in silico* experimentation in artificial systems (e.g., [116], see **Box 3**). We also emphasize that the flow of information is likely bidirectional. During production, the language system must receive information from brain regions that construct the meaning to be communicated. And even during comprehension, substantial behavioral evidence indicates that language processing is affected by extra-linguistic information sources, including the visual environment (e.g., [117]), non-verbal signals (e.g., gestures [118]), prior linguistic context (e.g., [119]), and knowledge of what information is accessible to our interlocutor (e.g., [120]).

**Outstanding questions**

1. How much of everyday language understanding takes place within the core language system? Is exportation the exception or the rule?
2. How are language representations "translated" into a format that other brain regions can process?
3. What specifically is represented in each of the exportation destinations and how is this representation formatted?
4. What stimulus-, person-, and neural-factors does exportation depend on?
5. Does the language system route information to the appropriate extra-linguistic system, or broadcast it broadly across the brain?
6. How is exportation from the language system coordinated with information flow in the opposite direction (from extra-linguistic brain areas to the language system)?

## 4. Concluding remarks

In this opinion piece, we have proposed that deep language understanding requires exporting information from the brain's core language system to other systems that build mental models of the minds, objects, and places described in the linguistic input, that store our memories and word knowledge, and that contain our perceptual and motor representations. Some evidence for this idea already exists, but our goal here is to outline a *conceptual framework* for thinking about language comprehension in the brain beyond the parsing of sentences and construction of shallow approximations of linguistic meaning carried out in the language network, and to outline many *testable hypotheses*, as we tried to do above.

Several caveats should be noted. First, strong empirical evidence for the exportation hypothesis is so far restricted to very few systems, and some studies find no evidence for exportation (e.g., [121-122]). Although suggestive evidence for the exportation hypothesis abounds (e.g., [68,72-73,79-80,86,94,23,96,101,104]) (see **Box 1**), methodological limitations (e.g., reliance on group-averaging approaches) restrict the conclusions we can draw from much past work. Second, even clear fMRI evidence for activation of an extra-linguistic brain system during language processing does not imply that deep understanding was achieved. For example, think of trying, but failing, to construct a mental model from a linguistic description (e.g., because the description is not clear enough, or because you are distracted while trying to do it). Finally, there is currently no consensus on how to behaviorally evaluate the depth of language understanding, and such measures would be important for attempts to falsify the idea that deep understanding requires exportation. The field of education research, which is currently grappling with how to probe students' understanding in the age of AI, may eventually provide some solutions.

By speculating about the brain regions that may work in concert with the language system, the circumstances under which they may be recruited, and the methods available to test these ideas, we have proposed a broad research program that harnesses recent progress in human cognitive neuroscience to advance a richer theory of what language understanding entails cognitively, neurally, and computationally.

**Glossary**

- **Exportation:** The process by which information delivered through language is transferred out of our core language system to one or more other brain regions for further processing.

- **Shallow understanding:** Building a representation of linguistic input based solely on knowledge of patterns of language use. This process requires recognizing familiar elements (e.g., morphemes, words, larger constructions), figuring out how those elements relate to one another, and ultimately combining them into a form-independent abstract representation (i.e. one that generalizes across paraphrases and languages) that can interface with downstream systems. Importantly, the representations constructed during shallow understanding are grounded only to the language statistics, and not tied to our lived experience or broader knowledge of the world.

- **Deep understanding:** Building rich mental models of the people, objects, places, and situations described in a language input. These models may incorporate our world knowledge, autobiographical memories, and/or perceptual and motor representations.

- **Situation model:** A representation of the 'situation' that is being described by a language input, including the characters and objects described, their environments and spatial locations, as well the goals of the individuals, and the temporal and causal relations among events.

- **Core language system:** A network of frontal and temporal brain regions that are highly specialized for language processing, strongly functionally and anatomically interconnected, and causality important in comprehending and producing language. The language system builds shallow representations of word- and sentence-level meanings.

- **Routing:** One hypothesis for how information is exported from the language system to downstream systems, whereby information is selectively directed (e.g., by the language system or a mediating brain region) to the most appropriate downstream system.

- **Broadcasting:** Another hypothesis for how information is exported to downstream systems, whereby information is sent from the core language system to all possible downstream systems, which then evaluate whether to process the information further.

## Acknowledgments

# References

1. Zhang, Y., Kauf, C., Levy, R. P., Gibson, E. (2025). Comparative illusions are evidence of rational inference in language comprehension. Journal of Experimental Psychology: General, 154(10), 2752-2771.

2. Bransford, J.D., Barclay, J.R., Franks, J.J. (1972). Sentence memory: A constructive versus interpretive approach. Cognitive Psychology, 3, 193-209.

3. Johnson-Laird, P.N. (1983). Mental models: Towards a cognitive science of language, inference, and consciousness (Cambridge, MA: Harvard University Press).

4. van Dijk, T.A., Kintsch, W. (1983). Strategies in discourse comprehension (New York: Academic Press).

5. Bower, G.H., Morrow, D.G. (1990). Mental models in narrative comprehension. Science, 247(4938), 44-8.

6. Gernsbacher, M.A. (1990). Language comprehension as structure building (Hillsdale, NJ: Erlbaum).

7. Graesser, A.C., Singer, M., Trabasso, T. (1994). Constructing inferences during narrative text comprehension. Psychological Review, 101(3), 371-95.

8. Zwaan, R.A., Radvansky, G.A. (1998). Situation models in language comprehension and memory. Psychological Bulletin, 123(2), 162-85.

9. Jackendoff, R.S. Erk, K.E. (2025). Toward a Deeper Lexical Semantics. Topics in Cognitive Science.

10. Fedorenko, E., Ivanova, A.A., Regev, T.I. (2024). The language network as a natural kind within the broader landscape of the human brain. Nature Reviews Neuroscience, 25(5), 289-312.

11. Kanwisher, N., McDermott, J., Chun, M. (1997). The Fusiform Face Area: A Module in Human Extrastriate Cortex Specialized for the Perception of Faces. Journal of Neuroscience, 17, 4302-4311.

12. Epstein, R. Kanwisher, N. (1998). A Cortical Representation of the Local Visual Environment. Nature, 392, 598-601.

13. Downing, P., Jiang, Y., Shuman, M., Kanwisher, N. (2001). A Cortical Area Selective for Visual Processing of the Human Body. Science, 293, 2470-2473.

14. Norman-Haignere, S., Kanwisher, N., McDermott, J. (2015). Distinct Cortical Pathways for Music and Speech Revealed by Hypothesis-Free Voxel Decomposition. Neuron, 88(6), 1281-96.

15. Overath, T., McDermott, J.H., Zarate, J.M., Poeppel, D. (2015) The cortical analysis of speech-specific temporal structure revealed by responses to sound quilts. Nature Neuroscience, 18(6), 903-11.

16. Saxe, R. Kanwisher, N. (2003). People thinking about thinking people: The role of the temporo-parietal junction in theory of mind. NeuroImage, 19, 1835-1842.

17. DiNicola, L.M., Braga, R.M., Buckner, R.L. (2020). Parallel distributed networks dissociate episodic and social functions within the individual. Journal of Neurophysiology, 123(3), 1144-1179.

18. Deen, B. Freiwald, W. (2022). Parallel systems for social and spatial reasoning within the cortical apex. bioRxiv.

19. Epstein, R.A., Baker, C.I. (2019). Scene Perception in the Human Brain. Annual Review of Vision Science, 5, 373-397.

20. Culham, J.C., Valyear, K.F. (2006). Human parietal cortex in action. Current Opinion in Neurobiology, 16(2):205-12.

21. Fischer, J., Tenenbaum, J. Kanwisher, N. (2016). The functional neuroanatomy of intuitive physical inference. PNAS, 113(34), E5072-81.

22. Isik, L., Koldewyn, K., Beeler, D., Kanwisher, N. (2017). Perceiving Social Interactions in the Posterior Superior Temporal Sulcus. PNAS, 115(1), E113-E114.

23. Ivanova, A.A. (2022). The role of language in broader human cognition: evidence from neuroscience. [Doctoral dissertation, Massachusetts Institute of Technology].

24. Braga, R.M., Buckner, R.L. (2017). Parallel Interdigitated Distributed Networks within the Individual Estimated by Intrinsic Functional Connectivity. Neuron, 95(2), 457-471.e5.

25. Duncan, J. (2010). The multiple-demand (MD) system of the primate brain: mental programs for intelligent behavior. Trends in Cognitive Sciences, 14(4), 172-9.

26. Kanwisher, N. (2025). Animal models of the human brain: Successes, limitations, and alternatives. Current Opinion in Neurobiology, 90, 102969.

27. Regev, M., Honey, C.J., Simony, E., Hasson, U. (2013). Selective and invariant neural responses to spoken and written narratives. Journal of Neuroscience, 33(40), 15978-88.

28. Malik-Moraleda, S., Ayyash, D., Gallée, J., Affourtit, J., Hoffmann, M., Mineroff, Z., Jouravlev, O., Fedorenko, E. (2022). An investigation across 45 languages and 12 language families reveals a universal language network. Nature Neuroscience, 25(8), 1014-1019.

29. Menenti, L., Gierhan, S.M., Segaert, K., Hagoort, P. (2011). Shared language: overlap and segregation of the neuronal infrastructure for speaking and listening revealed by functional MRI. Psychological Science, 22(9), 1173-82.

30. Fedorenko, E., Behr, M.K., Kanwisher, N. (2011). Functional specificity for high-level linguistic processing in the human brain. PNAS, 108(39), 16428-33.

31. Deen, B., Koldewyn, K., Kanwisher, N., Saxe, R. (2015). Functional Organization of Social Perception and Cognition in the Superior Temporal Sulcus. Cerebral Cortex, 25(11), 4596-609.

32. Jouravlev, O., Zheng, D., Balewski, Z., Le Arnz Pongos, A., Levan, Z., Goldin-Meadow, S., Fedorenko, E. (2019). Speech-accompanying gestures are not processed by the language-processing mechanisms. Neuropsychologia, 132, 107-132.

33. Fedorenko, E., Piantadosi, S.T., Gibson, E.A.F. (2024). Language is primarily a tool for communication rather than thought. Nature, 630(8017), 575-586.

34. Fridriksson, J., den Ouden, D.B., Hillis, A.E., Hickok, G., Rorden, C., Basilakos, A., Yourganov, G., Bonilha, L. (2018). Anatomy of aphasia revisited. Brain, 141(3), 848-862.

35. Bozic, M., Tyler, L.K., Ives, D.T., Randall, B. Marslen-Wilson, W.D. (2010). Bihemispheric foundations for human speech comprehension. PNAS, 107, 17439–17444.

36. Regev, T.I., Kim, H.S., Chen, X., Affourtit, J., Schipper, A.E., Bergen, L., Mahowald, K., Fedorenko, E. (2024). High-level language brain regions process sublexical regularities. Cerebral Cortex, 34(3), bhae077.

37. Pallier, C., Devauchelle, A. D. Dehaene, S. (2011). Cortical representation of the constituent structure of sentences. PNAS, 108, 2522–2527.

38. Shain, C., Kean, H., Casto, C., Lipkin, B., Affourtit, J., Siegelman, M., Mollica, F., Fedorenko, E. (2024). Distributed Sensitivity to Syntax and Semantics throughout the Language Network. Journal of Cognitive Neuroscience, 36(7), 1427-1471.

39. Bybee, J. (2010). Language, Usage and Cognition (Cambridge University Press).

40. Kauf*, C., Tuckute*, G., Levy, R., Andreas, J., Fedorenko, E. (2024). Lexical semantic content, not syntactic structure, is the main contributor to ANN-brain similarity of fMRI responses in the language network. Neurobiology of Language, 5(1), 7–42.

41. de Varda, A.G., Malik-Moraleda, S., Tuckute, G., Fedorenko, E. (2025). Multilingual computational models reveal shared brain responses to 21 languages. bioRxiv.

42. Ivanova, A.A., Mineroff, Z., Zimmerer, V., Kanwisher, N., Varley, R., Fedorenko, E. (2021). The Language Network Is Recruited but Not Required for Nonverbal Event Semantics. Neurobiology of Language, 2(2), 176-201.

43. Saha, S., Li, S. Tuckute, G., Li, Y., Zhang, R., Wehbe, L., Fedorenko, E., Khosla, M. (2025). Modeling the language cortex with form-independent and enriched representations of sentence meaning reveals remarkable semantic abstractness. bioRxiv.

44. Small, H., Lee Masson, H., Wodka, E., Mostofsky, S.H., Isik, L. (2025). Ubiquitous cortical sensitivity to visual information during naturalistic audiovisual movie viewing. PsyArXiv.

45. Huh, M., Cheung, B., Wang, T., Isola, P. (2024). The Platonic Representation Hypothesis. ICML.

46. de Varda, A.G., Petilli, M., Marelli, M. (2025). A distributional model of concepts grounded in the spatial organization of objects. Journal of Memory and Language, 142, 104624.

47. Dickey, M.W. Warren, T. (2015). The influence of event-related knowledge on verb-argument processing in aphasia. Neuropsychologia, 67, 63-81.

48. Kauf, C., Kim, H.S., Lee, E.J., Jhingan, N., She, J.S., Taliaferro, M., Gibson, E., Fedorenko, E. (2024). Linguistic inputs must be syntactically parsable to fully engage the language network. bioRxiv.

49. Grand, G., Blank, I.A., Pereira, F., Fedorenko, E.. Semantic projection recovers rich human knowledge of multiple object features from word embeddings. Nature Human Behavior, 6(7), 975-987.

50. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I. (2019). Language models are Unsupervised Multitask Learners. arXiv.

51. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., et al. (2020). Language Models Are Few-Shot Learners. 34th Conference on Neural Information Processing Systems (NeurIPS), 1877–1901.

52. Mahowald*, K., Ivanova*, A.A., Blank, I.A., Kanwisher, N., Tenenbaum, J.B., Fedorenko, E. (2024). Dissociating language and thought in large language models. Trends in Cognitive Science, 28(6), 517-540.

53. Barsalou, L.W. (2008). Grounded cognition. Annual Review of Psychology, 59, 617-45.

54. Hauk, O., Johnsrude, I., Pulvermüller, F. (2004). Somatotopic representation of action words in human motor and premotor cortex. Neuron, 41(2), 301-7.

55. Zhang, C.E., Wong, L., Grand, G., Tenenbaum, J.B. (2023). Grounded physical language understanding with probabilistic programs and simulated worlds. Proceedings of the Annual Meeting of the Cognitive Science Society.

56. Wong*, L, Grand*, G., Lew, A.K., Goodman, N., Mansinghka, V.K., Adreas, J., Tenenbaum, J.T. (2023). From Word Models to World Models: Translating from Natural Language to the Probabilistic Language of Thought. arXiv.

57. Olausson, T., Gu, A., Lipkin, B., Zhang, C., Solar-Lezama, A., Tenenbaum, J.T. (2023). LINC: A Neurosymbolic Approach for Logical Reasoning by Combining Language Models with First-Order Logic Provers. Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP), 5153-5176.

58. Piantadosi, S.T., Muller, D.C.Y., Rule, J.S., Kaushik, K., Gorenstein, M., Leib, E.R., Sanford, E. (2024). Why concepts are (probably) vectors. Trends in Cognitive Science, 28(9), 844-856.

59. Piccinini, G. (2025). Neural Hardware for the Language of Thought: New Rules for an Old Game. arXiv.

60. Saxe, R., Wexler, A. (2005). Making sense of another mind: the role of the right temporo-parietal junction. Neuropsychologia, 43(10), 1391-9.

61. Saxe, R., Powell, L.J. (2006). It's the thought that counts: specific brain regions for one component of theory of mind. Psychological Science, 17(8), 692-9.

62. Shain*, C., Paunov*, A., Chen*, X., Lipkin, B., Fedorenko, E. (2024). No evidence of theory of mind reasoning in the human language network. Cerebral Cortex, 33(10), 6299-6319.

63. Bruneau, E.G., Pluta, A., Saxe, R. (2012). Distinct roles of the 'shared pain' and 'theory of mind' networks in processing others' emotional suffering. Neuropsychologia, 50(2), 219-31.

64. Jacoby, N., Bruneau, E., Koster-Hale, J., Saxe, R. (2016). Localizing Pain Matrix and Theory of Mind networks with both verbal and non-verbal stimuli. Neuroimage, 126, 39-48.

65. Paunov, A.M., Blank, I.A., Jouravlev, O., Mineroff, Z., Gallée, J., Fedorenko, E. (2022). Differential Tracking of Linguistic vs. Mental State Content in Naturalistic Stimuli by Language and Theory of Mind (ToM) Brain Networks. Neurobiology of Language, 3(3), 413-440.

66. Battaglia, P.W., Hamrick, J.B., Tenenbaum, J.B. (2013). Simulation as an engine of physical scene understanding. PNAS, 110(45), 18327-32.

67. Pramod, R.T., Mieczkowski, E., Fang, C.X., Tenenbaum, J.B., Kanwisher, N. (2025). Decoding predicted future states from the brain's "physics engine". Science Advances, 11(22), eadr7429.

68. Jack, A.I., Dawson, A.J., Begany, K.L., Leckie, R.L., Barry, K.P., Ciccia, A.H., Snyder, A.Z. (2013). fMRI reveals reciprocal inhibition between social and physical cognitive domains. Neuroimage, 66, 385-401.

69. Pramod, R.T., Chomik-Morales, J., Schulz, L., Kanwisher, N. (2024). A region in human left prefrontal cortex selectively engaged in causal reasoning. Proceedings of the Annual Meeting of the Cognitive Science Society.

70. Nakamura, K., Kawashima, R., Sato, N., Nakamura, A., Sugiura, M., Kato, T., Hatano, K., Ito, K., Fukuda, H., Schormann, T., Zilles, K. (2000). Functional delineation of the human occipito-temporal areas related to face and scene processing—a PET study. Brain, 123, 1903–12.

71. O'Craven, K. Kanwisher, N. (2000). Mental imagery of faces and places activates corresponding stimulus-specific brain regions. Journal of Cognitive Neuroscience, 12, 1013-1023

72. Antonello*, R.J., Singh*, C., Jain, S., Hsu, A., Guo, S., Gao, J., Yu, B., Huth, A.G. (2024). Generative causal testing to bridge data-driven models and scientific theories in language neuroscience. arXiv.

73. Singh*, C., Antonello*, R.J., Guo, S., Mischler, G., Gao, J., Mesgarani, N., Huth, A.G. (2025). Evaluating scientific theories as predictive models in language neuroscience. bioRxiv.

74. Kamps, F.S., Julian, J.B., Kubilius, J., Kanwisher, N., Dilks, D.D. (2016). The occipital place area represents the local elements of scenes. NeuroImage, 132, 417–24.

75. Bonner, M.F., Epstein, R.A. (2018). Computational mechanisms underlying cortical responses to the affordance properties of visual scenes. PLOS Computational Biology, 14, e1006111

76. Baumann, O., Mattingley, J.B. (2010). Medial parietal cortex encodes perceived heading direction in humans. Journal of Neuroscience, 30, 12897–901.

77. Marchette, S.A., Vass, L.K., Ryan, J., Epstein, R.A. (2014). Anchoring the neural compass: coding of local spatial reference frames in human medial parietal lobe. Nature Neuroscience, 17, 1598–606.

78. Goebel, R., Khorram-Sefat, D., Muckli, L., Hacker, H., Singer, W. (1998). The constructive nature of vision: direct evidence from functional magnetic resonance imaging studies of apparent motion and motion imagery. European Journal of Neuroscience, 10(5), 1563-73.

79. Pulvermüller, F. (2005). Brain mechanisms linking language and action. Nature Reviews Neuroscience, 6(7), 576-82.

80. Aflalo, T., Zhang, C.Y., Rosario, E.R., Pouratian, N., Orban, G.A., Andersen, R.A. (2020). A shared neural substrate for action verbs and observed actions in human posterior parietal cortex. Science Advances, 6(43), eabb3984.

81. Bechtold, L., Cosper, S.H., Malyshevskaya, A., Montefinese, M., Morucci, P., Niccolai, V., Repetto, C., Zappa, A., Shtyrov, Y. (2023). Brain Signatures of Embodied Semantics and Language: A Consensus Paper. Journal of Cognition, 6(1), 61.

82. Isenberg, N., Silbersweig, D., Engelien, A., Emmerich, S., Malavade, K., Beattie, B., Leon, A.C., Stern, E. (1999). Linguistic threat activates the human amygdala. PNAS, 96(18), 10456-9.

83. Kensinger, E.A., Schacter, D.L. (2006). Processing emotional pictures and words: effects of valence and arousal. Cognitive, Affective, and Behavioral Neuroscience, 6(2), 110-26.

84. Ghio, M., Cassone, B., Tettamanti, M. (2025). Unaware processing of words activates experience-derived information in conceptual-semantic brain networks. Imaging Neuroscience, 3, imaga00484.

85. Ferstl, E. C., Rinck, M., von Cramon, D. Y. (2005). Emotional and Temporal Aspects of Situation Model Processing during Text Comprehension: An Event-Related fMRI Study. Journal of Cognitive Neuroscience, 17(5), 724–739.

86. Hsu, C.T., Jacobs, A.M., Citron, F.M., Conrad, M. (2015). The emotion potential of words and passages in reading Harry Potter—an fMRI study. Brain and Language,142, 96-114.

87. Landau-Wells, M., Lydic, K.O., Kennedy, J., Mittman, B.G., Thompson, T.W., Gupta, A., Saxe, R. (in press). Inclusionary and Exclusionary Preferences: A Test of Three Cognitive Mechanisms. Political Behavior.

88. Tulving, E. (1972). Episodic and semantic memory. Organization of memory, 1(1), 381-403.

89. Buckner, R.L., DiNicola, L.M. (2019). The brain's default network: updated anatomy, physiology and evolving insights. Nature Reviews Neuroscience, 20(10), 593-608.

90. Patterson, K., Nestor, P.J., Rogers, T.T. (2007). Where do you know what you know? The representation of semantic knowledge in the human brain. Nature Reviews Neuroscience, 8(12), 976-87.

91. Lambon Ralph, M.A., Jefferies, E., Patterson, K., Rogers, T.T. (2017). The neural and computational bases of semantic cognition. Nature Reviews Neuroscience, 18(1), 42-55.

92. Reilly, J., Shain, C., Borghesani, V., Kuhnke, P., Vigliocco, G., Peelle, J.E., Mahon, B.Z., Buxbaum, L.J., Majid, A., Brysbaert, M., Borghi, A.M., De Deyne, S., Dove, G., Papeo, L., Pexman, P.M., Poeppel, D., Lupyan, G., Boggio, P., Hickok, G., Gwilliams, L., Fernandino, L., Mirman, D., Chrysikou, E.G., Sandberg, C.W., Crutch, S.J., Pylkkänen, L., Yee, E., Jackson, R.L., Rodd, J.M., Bedny, M., Connell, L., Kiefer, M., Kemmerer, D., de Zubicaray, G., Jefferies, E., Lynott, D., Siew, C.S.Q., Desai, R.H., McRae, K., Diaz, M.T., Bolognesi, M., Fedorenko, E., Kiran, S., Montefinese, M., Binder, J.R., Yap, M.J., Hartwigsen, G., Cantlon, J., Bi, Y., Hoffman, P., Garcea, F.E., Vinson, D. (2025). What we mean when we say semantic: Toward a multidisciplinary semantic glossary. Psychonomic Bulletin and Review, 32(1), 243-280.

93. Wurm, M.F., Caramazza, A. (2019). Distinct roles of temporal and frontoparietal cortex in representing actions across vision and language. Nature Communications, 10(1), 289.

94. Popham, S.F., Huth, A.G., Bilenko, N.Y., Deniz, F., Gao, J.S., Nunez-Elizalde, A.O., Gallant, J.L. (2021). Visual and linguistic semantic representations are aligned at the border of human visual cortex. Nature Neuroscience, 24(11), 1628-1636.

95. Radvansky, G.A., Zacks, J.M. (2010). Event perception. Wiley Interdisciplinary Review Cognitive Science, 2(6), 608-620.

96. Baldassano, C., Chen, J., Zadbood, A., Pillow, J.W., Hasson, U., Norman, K.A. (2017). Discovering Event Structure in Continuous Narrative Perception and Memory. Neuron, 95(3), 709-721.e5.

97. Xu, J., Kemeny, S., Park, G., Frattali, C., Braun, A. (2005). Language in context: emergent features of word, sentence, and narrative comprehension. Neuroimage, 25(3), 1002-15.

98. Yarkoni, T., Speer, N.K., Zacks, J.M. (2008). Neural substrates of narrative comprehension and memory. Neuroimage, 41(4), 1408-25.

99. Tylén, K., Christensen, P., Roepstorff, A., Lund, T., Østergaard, S., Donald, M. (2015). Brains striving for coherence: Long-term cumulative plot formation in the default mode network. Neuroimage, 121, 106-14.

100. Lerner, Y., Honey, C. J., Silbert, L. J. Hasson, U. (2011). Topographic mapping of a hierarchy of temporal receptive windows using a narrated story. Journal of Neuroscience 31, 2906–2915.

101. Simony, E., Honey, C.J., Chen, J., Lositsky, O., Yeshurun, Y., Wiesel, A., Hasson, U. (2016). Dynamic reconfiguration of the default mode network during narrative comprehension. Nature Communications, 7, 12141.

102. Ferstl, E.C., von Cramon, D.Y. (2007). Time, space and emotion: fMRI reveals content-specific activation during text comprehension. Neuroscience Letters, 427(3), 159-64.

103. Whitney, C., Huber, W., Klann, J., Weis, S., Krach, S., Kircher, T. (2009). Neural correlates of narrative shifts during auditory story comprehension. Neuroimage, 47(1), 360-6.

104. Speer, N.K., Reynolds, J.R., Swallow, K.M., Zacks, J.M. (2009). Reading stories activates neural representations of visual and motor experiences. Psychological Science, 20(8), 989-99.

105. Duff, M.C., Brown-Schmidt, S. (2012). The hippocampus and the flexible use and processing of language. Frontiers in Human Neuroscience, 6, 69.

106. Milivojevic, B., Varadinov, M., Vicente Grabovetsky, A., Collin, S.H., Doeller, C.F. (2016). Coding of Event Nodes and Narrative Context in the Hippocampus. Journal of Neuroscience, 36(49), 12412-12424.

107. Dijksterhuis, D.E., Self, M.W., Possel, J.K., Peters, J.C., van Straaten, E.C.W., Idema, S., Baaijen, J.C., van der Salm, S.M.A., Aarnoutse, E.J., van Klink, N.C.E., van Eijsden, P., Hanslmayr, S., Chelvarajah, R., Roux, F., Kolibius, L.D., Sawlani, V., Rollings, D.T., Dehaene, S., Roelfsema, P.R. (2024). Pronouns reactivate conceptual representations in human hippocampal neurons. Science, 85(6716), 1478-1484.

108. Blank, I. A. Fedorenko, E (2020). No evidence for differences among language regions in their temporal receptive windows. Neuroimage, 219, 116925.

109. Jacoby, N. Fedorenko, E (2020). Discourse-level comprehension engages medial frontal Theory of Mind brain regions even for expository texts. Language, Cognition, and Neuroscience, 35, 780–796.

110. Christiansen, M.H., Chater, N. (2016). The Now-or-Never bottleneck: A fundamental constraint on language. Behavioral and Brain Sciences, 39, e62.

111. Amalric, M., Dehaene, S. (2016). Origins of the brain networks for advanced mathematics in expert mathematicians. PNAS, 113(18), 4909-17.

112. Cohen, L., Salondy, P., Pallier, C., Dehaene, S. (2021). How does inattention affect written and spoken language processing? Cortex, 138, 212-227.

113. Zeman, A. (2024). Aphantasia and hyperphantasia: exploring imagery vividness extremes. Trends in Cognitive Scienc, 28(5), 467-480.

114. Olson, H.A., Johnson, K.T., Nishith, S., Frosch, I.R., Gabrieli, J.D.E., D'Mello, A.M. (2024). Personalized Neuroimaging Reveals the Impact of Children's Interests on Language Processing in the Brain. Imaging Neuroscience, 2, 1-14.

115. Mak, M., Willems, R.M. (2019). Mental simulation during literary reading: Individual differences revealed with eye-tracking. Language, Cognition and Neuroscience, 34(4), 511-535.

116. Jain, S., Vo, V.A., Wehbe, L., Huth, A.G. (2024). Computational Language Modeling and the Promise of In Silico Experimentation. Neurobiology of Language, 5(1):80-106.

117. Tanenhaus, M.K., Spivey-Knowlton, M.J., Eberhard, K.M., Sedivy, J.C. (1995). Integration of visual and linguistic information in spoken language comprehension. Science, 268(5217), 1632-4.

118. Holle, H., Gunter, T.C. (2007). The role of iconic gestures in speech disambiguation: ERP evidence. Journal of Cognitive Neuroscience, 19(7), 1175-92.

119. Nieuwland, M.S., Van Berkum, J.J. (2006). When peanuts fall in love: N400 evidence for the power of discourse. Journal of Cognitive Neuroscience, 18(7), 1098-111.

120. Hanna, J.E., Tanenhaus, M.K., Trueswell, J.C. (2003). The effects of common ground and perspective on domains of referential interpretation. Journal of Memory and Language, 49(1), 43–61.

121. Dravida, S., Saxe, R., Bedny, M. (2013). People can understand descriptions of motion without activating visual motion brain regions. Frontiers in Psychology, 4, 537.

122. Bedny, M., Caramazza, A., Grossman, E., Pascual-Leone, A., Saxe, R. (2008). Concepts are more than percepts: the case of action verbs. Journal of Neuroscience, 28(44), 11347-53.

123. Drijvers, L., Small, S. L., Skipper, J. I. (2025). Language is widely distributed throughout the brain. Nature Reviews Neuroscience, 26(3), 189.

124. Aliko, S., Franch, M., Kewenig, V., Wang, B., Cooper, G., Glotfelty, A., Hayden, B., Small, S.L., Skipper, J.I. (2025). The entire brain, mor or less, is at work: 'Language regions' are artefacts of averaging. bioRxiv.

125. Piantadosi, S., Hill, D. (2022). Meaning without reference in large language models. 36th Conference on Neural Information Processing Systems (NeurIPS).

126. Pavlick, E. (2023). Symbols and grounding in large language models. Philosophical Transactions of the Royal Society A, 381(2251), 20220041.

127. Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y.T., Li, Y., Lundberg, S., Nori, H., Palangi, H., Tulio Ribeiro, M., Zhang, Y. (2023). Sparks of Artificial General Intelligence: Early experiments with GPT-4. arXiv.

128. Antonello, R.J., Vaidya, A.R., Huth, A.G. (2023). Scaling laws for language encoding models in fMRI. 37th Conference on Neural Information Processing Systems (NeurIPS), 21895-21907.

129. Binhuraib, T., Gao, R., Ivanova, A.A. (2025). LITcoder: A General-Purpose Library for Building and Comparing Encoding Models. arXiv.

130. AlKhamissi, B., Tuckute, G., Tang, Y., Binhuraib, T., Bosselut*, A., Schrimpf*, M. (2025). From Language to Cognition: How LLMs Outgrow the Human Language Network. Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing (EMNLP).

131. Hosseini, E.A., Schrimpf, M., Zhang, Y., Bowman, S., Zaslavsky, N., Fedorenko, E. (2024). Artificial Neural Network Language Models Predict Human Brain Responses to Language Even After a Developmentally Realistic Amount of Training. Neurobiology of Language, 5(1), 43-63.

132. Dong, C.V., Lu, Q., Norman, K.A., Michelmann, S. (2025). Towards large language models with human-like episodic memory. Trends in Cognitive Science, 29(10), 928-941.

133. AlKhamissi, B., Nicoló De Sabbata, C., Chen, Z., Schrimpf*, M., Bosselut*, A. (2025). Mixture of Cognitive Reasoners: Modular Reasoning with Brain-Like Specialization. arXiv.

134. AlKhamissi, B., Tuckute, G., Bosselut*, A., Schrimpf*, M. (2024). The LLM Language Network: A Neuroscientific Approach for Identifying Causally Task-Relevant Units. Proceedings of the 2025 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL),10887-10911.

135. Hanna, M., Belinkov, Y., Pezzelle, S. (2025). Are Formal and Functional Linguistic Mechanisms Dissociated in Language Models? Computational Linguistics, 1-40.

136. Jamaa, Y., Badr, Alkhamissi, Gosh, S., Schrimpf, M. (2025). Evaluating Contrast Localizer for Identifying Causal Units in Social Mathematical Tasks in Language Models. 2025 Conference on Language Modeling (COLM), Interplay of Model Behavior and Model Internals Workshop.

137. Poldrack, R.A. (2011). Inferring mental states from neuroimaging data: from reverse inference to large-scale decoding. Neuron, 72(5): 692-7.

138. Braga, R.M., DiNicola, L.M., Becker, H.C., Buckner, R.L. (2020). Situating the left-lateralized language network in the broader organization of multiple specialized large-scale distributed networks. Journal of Neurophysiology, 124(5), 1415-1448.

139. Shain, C., Fedorenko, E. (2025). A language network in the individualized functional connectomes of over 1,000 human brains doing arbitrary tasks. bioRxiv.