
Position: An Inner Interpretability Framework for AI Inspired by Lessons from Cognitive Neuroscience

Martina G. Vilas^{1,2} Federico Adolfi¹ David Poeppel^{1,3} Gemma Roig^{2,4}

Abstract

Inner Interpretability is a promising emerging field tasked with uncovering the inner mechanisms of AI systems, though how to develop these mechanistic theories is still much debated. Moreover, recent critiques raise issues that question its usefulness to advance the broader goals of AI. However, it has been overlooked that these issues resemble those that have been grappled with in another field: Cognitive Neuroscience. Here we draw the relevant connections and highlight lessons that can be transferred productively between fields. Based on these, we propose a general conceptual framework and give concrete methodological strategies for building mechanistic explanations in AI inner interpretability research. With this conceptual framework, Inner Interpretability can fend off critiques and position itself on a productive path to explain AI systems.

1. Introduction

Inner Interpretability is an emerging subfield of Artificial Intelligence (AI) that is drawing increasing attention as models get larger and better. Interpreting the internal mechanisms of these models in human-understandable terms is of great interest for scientific, engineering, and safety reasons (Räuker et al., 2023). However, despite interesting recent results, the field seems to be missing a conceptual framework that guides the development and analysis of these mechanistic explanations. Although proposals have been made on how to quantify the alignment of human-intelligible high-level theories with the internal operations of neural networks (e.g. causal abstraction approaches, see Geiger et al., 2023), a

¹Ernst Strüngmann Institute for Neuroscience, Frankfurt am Main, Germany ²Department of Computer Science, Goethe University, Frankfurt am Main, Germany ³Department of Psychology, New York University, New York, USA ⁴The Hessian Center for AI, Hessen, Germany. Correspondence to: Martina G. Vilas <martina.vilas@esi-frankfurt.de>.

more general framework for developing, discussing, analyzing, and refining these high-level mechanistic explanations is still needed. The lack of such a conceptual framework has made this subfield vulnerable to critiques that question its usefulness to advance the broader goals of AI.

In addition, responses to these criticisms have not taken notice of how similar the issues raised are those to another field that has extensively dealt with them: Cognitive Neuroscience (but see Lindsay & Bau, 2023, for more general links). In this position paper, we address these gaps. **We explain the connections between issues in AI Inner Interpretability and those in Cognitive Neuroscience, and we propose that we can take advantage of these links to derive lessons and concrete conceptual and methodological strategies to interpret AI systems mechanistically.**

Overview: The remaining sections are organized as follows. Section 2 gives an overview of the field of Inner Interpretability and current critiques directed at it. Section 3 describes the issues that give rise to these critiques, draws parallels with longstanding issues in Cognitive Neuroscience, and explains how these have been tackled. Section 4 draws on these lessons, proposes a conceptual framework for Inner Interpretability, and derives concrete methodological strategies. Finally, in Section 5 we explain how adopting this framework allows us to better integrate past and future studies and address critiques.

2. AI inner interpretability research

Recent years have seen an increase in research aiming to understand the internal structural components, operations, and representations of deep neural networks. This line of work has recently been given the name of *inner interpretability* (for a review see Räuker et al., 2023). A significant amount of work in this field is concerned with explaining how the inner mechanisms of these models give rise to their capabilities¹. This differs from other lines of interpretability research where the behavior of the model is attributed to specific properties of the input or the training dataset (Bommasani et al., 2021).

¹In this paper, the term *inner interpretability* is used to refer to this subset of work.

For example, inner interpretability studies have tried to characterize the model components (e.g. Elhage et al., 2021; Geva et al., 2021; McDougall et al., 2023; Olsson et al., 2022), functions (e.g. Merullo et al., 2023b; Todd et al., 2024), and algorithms (e.g. Zhong et al., 2023), behind a variety of emergent behaviors. Other work has developed methods to automate the discovery and analysis of activation sub-spaces (e.g. Burns et al., 2022), circuits (e.g. Conmy et al., 2023; Lepori et al., 2023), and internal representations (e.g. Belrose et al., 2023; Hernandez et al., 2023) that have a causal effect on the output of the model.

Inner interpretability work is motivated, on the one hand, by goals related to safety and transparency in AI (Casper et al., 2024). Mechanistic explanations may, for example, improve the predictability of model behavior (Bommasani et al., 2021), and enable editing out harmful or incorrect representations and steer decisions (e.g. Hernandez et al., 2023). On the other hand, the field could help improve aspects of model performance. For instance, mechanistic discoveries can lead to efficiency improvements (Zhang et al., 2023), or the development of better architectures (e.g. Akyurek et al., 2024).

Recent critiques of the field, however, paint a more pessimistic picture. It has been questioned whether current methodological strategies, which are not yet well understood, will lead to any meaningful insights. Critics believe that research practices in use are likely to lead to a false sense of understanding, which results in misleading or contradictory claims (viz. R auker et al., 2023). It has also been argued that these procedures can only achieve weak generalization to real-world problems or models (Doshi-Velez & Kim, 2017; R auker et al., 2023). More importantly, there seems to be a lack of clarity regarding the overarching questions of the field (Krishnan, 2020), and what it means to mechanistically understand a model. Consequently, there is a heightened risk of letting the choice of research questions be guided by the availability of technology and heuristics (i.e., a hardware-software lottery; Hooker, 2021). These issues, if not addressed, put the field at risk of stagnation.

The outlined issues seem to arise from a lack of consensus on how to build mechanistic explanations, and how to evaluate and compare them through a common conceptual framework. The field of Inner Interpretability is therefore in need of an analysis of its research practices and a general framework to guide them. In this work, we propose that these tools can be adapted from the Cognitive Neuroscience field. The rapid pace of current AI research often causes a disconnect from knowledge gained in adjacent fields. In this case, Cognitive Neuroscience has confronted many of the issues that give rise to the critiques of Inner Interpretability, and has developed conceptual frameworks and methodological strategies to deal with them. In the next section, we

set the groundwork to transfer these lessons to AI Inner Interpretability.

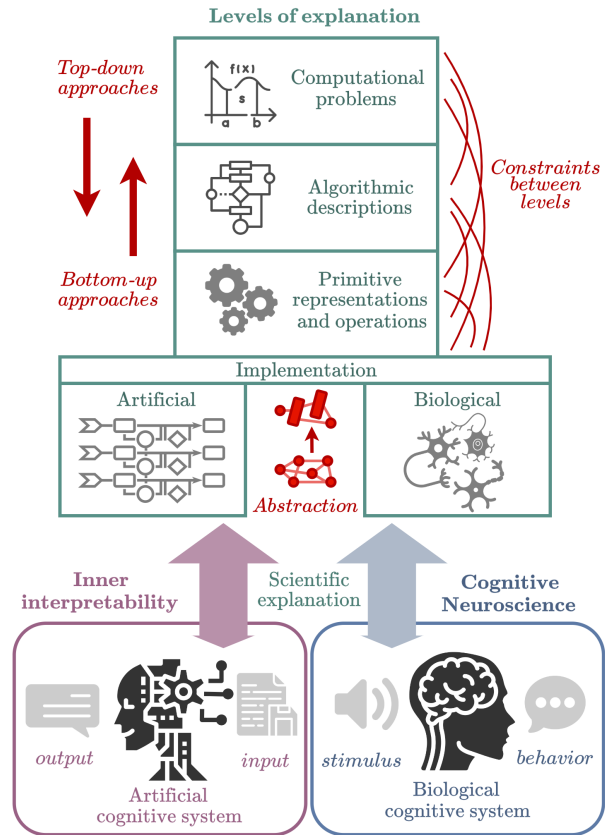


Figure 1. The fields of Inner Interpretability and Cognitive Neuroscience aim to mechanistically explain the behavior of artificial and biological systems, respectively. The multilevel explanatory framework proposed here draws out the parallels and suggests strategies that can be transferred between fields to tackle current issues in Inner Interpretability (shown in red).

3. Parallel issues between fields

Cognitive Neuroscience and AI Inner Interpretability share similar goals, concepts, and methodology. Both fields aim to uncover the mechanisms giving rise to the behavior of complex systems. In both cases, researchers want to determine what capacities these systems possess, how they are implemented, and how they can be described in a human-understandable way. The difference between the two fields, in principle, lies at the implementation level: Cognitive Neuroscience must grapple with opaque biological substrates (e.g., human brains), whereas Inner Interpretability deals with the more accessible virtual substrate of artificial learning systems (e.g., transformer networks).

Naturally, obstacles for explaining these systems, and ways to overcome them, will be similar between fields. Indeed,

here we will demonstrate that overlapping issues exist. Luckily, Cognitive Neuroscience is an older discipline than Inner Interpretability. Therefore, the community has had time to confront these obstacles by developing conceptual frameworks and methodological strategies. This allows us to look at these proposals and apply insights to Inner Interpretability, to tackle the issues undermining the field’s role in advancing the broader goals of AI.

In this section, we give some examples of these overlapping issues, namely, issues with mechanistic explanations, levels of abstraction, and bottom-up versus top-down approaches. For each issue, we (i) show how it arises in Inner Interpretability and fuels critiques of the field, and (ii) explain how it played out and was tackled in Cognitive Neuroscience using multilevel conceptual frameworks (Marr & Poggio, 1976) and other methodological strategies that guide the development of mechanistic explanations. Section 4 discusses how to apply these lessons from Cognitive Neuroscience to inner interpretability work.

3.1. Issues with mechanistic explanation

Incomplete mechanisms in Inner Interpretability. Inner Interpretability is primarily concerned with uncovering the inner mechanisms of AI systems. A classic notion of mechanistic explanation is an account of the relevant entities, activities, and organizational features (spatial and temporal relationships) that interact to have a causal effect in producing the phenomenon to be explained (Machamer et al., 2000). In practice, however, a great number of inner interpretability studies seem to equate mechanistic explanations solely with localizing model components that have a causal effect on behaviors of interest. For example, researchers develop and make use of intervention methods such as ablation or activation patching, to find model components, such as neurons, sub-modules, and circuits, that have a causal link with the behavior (e.g. Meng et al., 2022; Vig et al., 2020). Yet, mechanistic proposals in these studies contain gaps and are thus *incomplete*. For instance, a neuron with a causal effect on the behavior may be localized, but its operation is left unspecified. Or a causally relevant circuit may be mapped, but it is not decomposed into a sequence of entities performing activities (for good examples of such decomposition see Merullo et al., 2023a; Wang et al., 2022).

Proposing incomplete mechanisms as explanations can lead to a false sense of understanding (Craver, 2006), and can easily result in mischaracterizing elements in the mechanism. For example, a recent study probing the mechanisms of factual recall has shown that localizing an entity that carries factual information in a GPT model (e.g., MLP middle layers) does not mean that intervening in the operations of this entity will lead to the intuitive outcome: effective editing of this factual information (Hase et al., 2024). This

could be the result of a post-hoc misspecification of the operation performed by the MLP layers, in the context of an incomplete mechanistic proposal that did not spell out step-by-step how the output was produced.

Attempts at complete mechanisms in Cognitive Neuroscience. Similarly to the field of Inner Interpretability, a significant number of studies in Cognitive Neuroscience operate under a lenient definition of mechanistic explanation (Ross & Bassett, 2024). The term is often used to refer to findings showing a causal relationship between a neural component and a cognitive phenomenon, without providing a process description of how the outcome is produced. However, researchers have long argued that this type of work is insufficient to build mechanistic explanations (viz. Krakauer et al., 2017). For instance, previous studies have shown the difficulty of inferring cognitive function from neural recordings (Poldrack, 2006), and the insufficiency of uncovering necessary and sufficient neural circuits for building mechanistic explanations (Gomez-Marin, 2017).

To build better and more complete mechanistic explanations, the strategy of multilevel analysis was introduced (more on this in Section 4; Marr & Poggio, 1976). Multilevel explanations include comprehensive functional characterizations of the behavior (*what* is the system doing and *why*), algorithmic descriptions of how the function is computed (*how* is it doing it), and decomposition of the algorithms in a list of fine-grained and human-interpretable primitive operations and representations implemented in the neural substrate (see Fig. 1). A complete multilevel explanation results in a characterization of the relevant brain components, their activities and interactions, that implement the capacity of interest.

More recent work has pointed out that mutual constraints between levels provide an avenue to construct more complete mechanistic explanations in practice (Danks, 2013; Love, 2021). That is, results at one level can suggest what to look for at other levels. The higher levels provide a conceptual and formal structure that can guide the search for and characterization of neural mechanisms (Griffiths et al., 2010; Krakauer et al., 2017). For example, the location of circuits in the brain that are causally involved in speech processing has been known for a long time, but it was not until theoretically and empirically motivated computational steps were hypothesized that a better understanding of the functions of each of the neural structures composing the circuit was achieved (Hickok & Poeppel, 2007). In analogy to the *factual recall* example discussed previously, if a theoretically motivated sequence of primitive representations and operations had been mapped to the uncovered circuit for factual recall, perhaps the role assigned to the MLP layers would have been more accurate and interventions would have produced the predicted result.

3.2. Issues with levels of abstraction

Weak motivation for abstractions in Inner Interpretability. Mechanisms can be explained at different levels of abstraction at the implementation level (e.g., neurons, circuits, modules, representational trajectories²; see Fig. 1, bottom). That is, complete explanations can be given at levels of abstraction that ignore various details at other levels. For example, explanations of how large language models can retrieve facts have been given at the level of sub-module operations, without referring to the neurons, layers, and nonlinearities composing them (Chughtai et al., 2023; Geva et al., 2023). Importantly, some levels of abstraction may lead to explanations that are more human-intelligible and can be efficiently uncovered across model sizes.

In practice, work on inner interpretability typically chooses the level of abstraction not based on mechanistic principles but rather arbitrarily (e.g., based on previous studies), without examining their implications. Weak motivations for abstractions may block progress in building a robust understanding of the behavior. For example, some assume that human-interpretable representations are to be found at the level of neurons, based on previous empirical findings (e.g. Hernandez et al., 2021). However, recent studies suggest that relevant features may be encoded in superposition (Elhage et al., 2022), representing a switch in abstraction level to ‘directions in activation space’. In addition, it has been suggested that finding these explanations at microscopic levels of abstraction may not be computationally feasible in large models (Adolfi et al., 2024; Zou et al., 2023).

Attempts at choosing better abstractions in Cognitive Neuroscience. Neurobiological processes can also be characterized at different levels of abstraction and the choice is consequential (Barack & Krakauer, 2021). The problem of choosing an appropriate level of abstraction has been discussed in Cognitive Neuroscience in the context of a *mapping problem*: to build mechanistic explanations, the basic parts of neurobiology (e.g., synapses, neurons, brain regions) must be mapped onto the basic operations and representations of human cognition (Poehpel, 2016). For this match to succeed, the right level of abstraction for the basic parts of the brain must be discovered.

In the context of multilevel explanations (see Section 3.1; Marr & Poggio, 1976), the level of primitive operations and representations can be conceptualized as depicting the set of cognitive ‘parts’, while the implementation level encodes the set of neurobiological components (viz. Poehpel, 2016). The selection of primitive candidates is motivated by theories of cognitive functions that have been formally and

²Levels of abstraction are not to be equated with levels of physical organization (i.e. spatio-temporal scale), nor with the levels of explanation of the multilevel framework. These are orthogonal.

empirically validated. But their choice is also constrained by the type of operations that can be implemented in the brain. The reverse is also true: the decomposition into neural components at the implementation level is constrained by the list of plausible primitives. Therefore, the choice of abstraction both at the level of primitives and at the implementation level can be guided by the quality of the match it affords between neural and cognitive ‘parts’. For instance, *segmentation* operations were long hypothesized as computational primitives of speech recognition. To understand their implementation in the brain, cognitive neuroscientists found it useful to abstract away from neural circuit connectivity to focus instead on the level of oscillations of neuronal ensembles (Giraud & Poeppel, 2012; Poeppel & Assaneo, 2020). This abstraction has led to productive research programs.

Comparably, in Inner Interpretability, some recent approaches have proposed to adjust the level of abstraction of high-level theories (e.g. of the variables in causal models) to better align with empirical data (Geiger et al., 2023). More work is needed to also guide this abstraction at the implementation level.

3.3. Issues with bottom-up versus top-down approaches

Overoptimistic bottom-up approaches in Inner Interpretability. A frequent distinction in the field of Inner Interpretability is between *top-down* and *bottom-up* research methodologies. The term *top-down* is linked to work mapping pre-defined and human-interpretable representations and operations to model components (e.g. Kim et al., 2018; Meng et al., 2022; Mu & Andreas, 2020). Recently, it has been shown that some findings in this line of work can be misleading when assumptions are not well-tested (e.g. Bolukbasi et al., 2021), since alternative mechanisms may lead to the same empirical observation. As an alternative, a bottom-up approach named *mechanistic interpretability*³ was introduced to study models “without a priori theories” (Olah, 2023). They propose to decompose the network into the smallest elements possible, thoroughly investigate their functions and interactions by observing, perturbing, and describing them, and work upward to build abstractions from these foundations until their role in the behaviors of the model can be explained (Olah, 2023; Olsson et al., 2022).

However, bottom-up approaches are not free of assumptions and, if left unexamined, these can easily lead to roadblocks, false understanding, and non-generalizable claims. On the one hand, a choice is being made regarding which neural component to analyze, at which level of abstraction, what conditions to perturb, which behaviors to analyze as responses, etc. In addition, it heavily relies on the interpretations made by a human observer. Automatic interpretability

³The term *mechanistic interpretability* may cause confusion, as both top-down and bottom-up research study mechanisms.

ity methods are not free of these assumptions either. The neurons-versus-directions example discussed in the last section illustrates the assumed levels of abstraction in automatic discovery methods. On the other hand, these methods carry the risk of making inferences that do not scale up to the capacities and models of interest that originally motivated the research. Given the vast space of research strategies without a priori theories, bottom-up studies start by tackling small problems in small networks that are easily interpretable with available tools. However, the ultimate goals of understanding are related to the more complex capacities of large AI systems. Most work assumes, without guarantees, that these research strategies will eventually scale beyond toy problems and models.

Methodology-aware approaches in Cognitive Neuroscience. Whether *top-down* or *bottom-up* approaches are more effective in discovering mechanistic explanations has been discussed in Cognitive Neuroscience since its origins (see appendix Fig. 3). Top-down approaches start by defining and decomposing the behavior of interest computationally. Algorithmic candidates for these computations can then be linked to neural processes (e.g., Krakauer et al., 2017; Egan, 2018). Bottom-up approaches first thoroughly describe and manipulate neural parts and activities, and then try to infer the cognitive capacity they implement (e.g., Buzsáki, 2020; Churchland & Sejnowski, 1999). Like in Inner Interpretability, bottom-up approaches were advanced in response to fears of being misled by inaccurate theories.

Radical bottom-up approaches have highlighted the benefits of carrying out comprehensive brain mapping efforts at the finest levels of detail (i.e., building a ‘connectome’; Sporns et al., 2005). Such efforts (e.g., the Human Connectome Project) were expected to turn low-level descriptions of the brain into high-level explanations of cognitive abilities. The idea was that looking for regularities at the level of neural ‘stuff’ would eventually reveal their underlying mechanistic organization. However, these approaches are said to have “overpromised and underdelivered” (Gomez-Marin, 2021). There are at least two reasons for this failure. First, there is rarely a one-to-one mapping between neural parts (e.g. circuits) and the algorithmic or functional descriptions of the cognitive capacities they implement (viz. Gunaratne et al., 2017). Second, bottom-up approaches are not free of assumptions. These are needed to narrow down vast search spaces that cannot be explored exhaustively. The selection of neural parts and processes for investigation is necessarily based on such preconceptions. Candidate cognitive operations and representations are also selected for exploration based on explicit or implicit assumptions. Lack of explicit assumptions does not mean that assumptions are not being made about these aspects. Explicit assumptions can be examined and tested. Implicit or unexamined ones can easily result in misleading conclusions.

Cognitive Neuroscience has learned this lesson the hard way. Various radical top-down and bottom-up approaches were put forth (e.g., Buzsáki, 2020; Krakauer et al., 2017; Niv, 2021) and many of the promises have been overstated on both sides (viz. Gomez-Marin, 2021). But it is now clear that the simultaneous execution of combined approaches is necessary to discover their invariants and reach useful mechanistic explanations (Poeppele & Adolphi, 2020). This more pluralistic process takes advantage of the mutual constraints between implementation and other levels of explanation (see Fig. 1, right) to arrive at consistent mechanistic theories. In addition, to reduce the impact of incorrect assumptions and theories, a variety of methodological strategies have been proposed to rigorously test the validity of these conjectures. We will discuss how to apply these to the field of Inner Interpretability in the next section.

4. A framework for Inner Interpretability

In this section, we describe a conceptual framework where the lessons from Cognitive Neuroscience discussed in the previous sections can be translated into methodological strategies for AI Inner Interpretability. Our running example will be *factual recall*, a frequently studied topic in Inner Interpretability (Hernandez et al., 2023; Meng et al., 2022; Geva et al., 2023; Chughtai et al., 2023; Yu et al., 2023). We will illustrate, without loss of generality, how the strategies we propose can be applied to study the inner mechanisms of language transformers that implement the capacity to recall facts.

4.1. Building multilevel mechanistic explanations

Previous work has already pointed out the usefulness of applying a multilevel conceptual framework (viz. Marr & Poggio, 1976) for better analyzing and comparing the performance of machine learning models (Hamrick & Mohamed, 2020). Here, we show that a multilevel analysis of capacities also provides a useful conceptual structure to investigate their inner mechanisms. Each level offers a qualitatively different description of the mechanism under study, and as such each level employs a specialized terminology and provides a different angle of analysis (see appendix C for a discussion on their separability). A productive research program in Inner Interpretability makes use of mutual constraints across the levels to arrive at a complete mechanistic explanation. Next, we explain how to locate each level of explanation in Inner Interpretability research projects, and how to use mutual constraints between levels to converge on useful mechanistic explanations.

4.1.1. COMPUTATIONAL PROBLEMS

To build a mechanistic explanation, it is necessary to define and thoroughly characterize the phenomena to be explained

(Craver, 2006). In the framework of multilevel explanation (see Fig. 1), a systematic behavior of interest (e.g., factual recall) is selected and described at the *computational* level (Marr & Poggio, 1976). A computational description gives a functional specification of the capacity underlying the observed behavior (i.e., it describes *what* the system is doing). Here, the capacity can be characterized as an information-processing task where the system maps inputs to outputs. This level also formalizes the input and output domains, and may specify additional properties or parameters of the mapping (viz. Shagrir & Bechtel, 2017).

The difference between observed behavior and underlying capacity is important. It is intuitive to observe some complex-looking model behavior (e.g., the classification of images of different animals using an abstract category such as ‘animal’) and infer an interesting capacity of the model (e.g., the ability to build rich representations that abstract away from particular animals such as cats or dogs). However, the same behavior can be the consequence of different underlying capacities. The propensity of AI models to exploit ‘shortcuts’ means that often some of the true underlying capacities turn out to be uninteresting. For instance, cats and dogs can be distinguished from inanimate objects by building abstract representations, but this may also be achieved by exploiting contextual cues given by confounds in training datasets. Ignoring these issues regularly leads to claims that models possess interesting capacities, followed by more rigorous experimentation eliciting behaviors that evidence their absence (viz. Mitchell, 2021; Mitchell & Krakauer, 2023; Bowers et al., 2022).

Consider the example of *factual recall*, a capacity intuitively described as the recall of truthful knowledge about entities in the world. For instance, a fact might be ‘Paris is the capital of France’, and the retrieval of ‘Paris’ in response to the prompt ‘The capital of France is [—]’ is an example of factual recall behavior of the model. This intuitive verbal definition is insufficient to carry out a scientific analysis of the mechanisms that enable its emergence in a large language model. To be more precise, we can define factual recall as the capacity to recall an attribute (e.g., Paris) when prompted with a subject (e.g., France) and a relationship (e.g., capital), from a particular knowledge domain (e.g., political geography). A formal definition of the capacity could be constructed as follows (to be refined iteratively).

Definition 4.1 (Facts and fact domains). A *fact* is a 3-tuple $F = (S, R, A) \in \mathcal{D}$, where S is a subject, R is a relation, A is an attribute, and $\mathcal{D} = \{F_1, F_2, \dots, F_n\}$ represents a *fact domain*. Incompleteness of a fact tuple is denoted with \perp in the corresponding component.

Definition 4.2 (Factual recall).

Computational problem: \mathcal{D} -FACTUALRECALL

Input: An incomplete *fact* tuple (Def. 4.1), $F_I = (S, R, \perp)$

corresponding to a complete fact $F = (S, R, A) \in \mathcal{D}$, where \mathcal{D} is a constant fact domain.

Output: A completion F_C of F_I such that $F_C \in \mathcal{D}$.

A computational level explanation should also specify *why* the behavior occurs in its specific context (Shagrir & Bechtel, 2017, see also resource-rational analysis: Lieder & Griffiths 2020). That is, it should explain how environmental properties constrain and shape the function of the system. In our example, a context where the model is expected to deliver truthful information (e.g., a chatbot interacting with users) encourages the emergence of factual recall.

Computational-level descriptions (see Fig. 1, top) can be tested both formally and empirically. Descriptions at this level take the form of system capacities as computational problems (e.g., Def. 4.2). But, not all computational problems that can be written down describe capacities that are possible in practice (e.g., they are uncomputable or intractable; Wareham, 1998; or fall outside the expressive power of the AI architecture; viz. Strobl et al., 2023). Overlooking this can lead to explanations that are inconsistent between levels. For instance, one could inadvertently propose a computational-level description of a capacity that is outside the class of problems that the model architecture (e.g., transformer) can solve (viz. Strobl, 2023; Strobl et al., 2023), yielding inconsistency between computational and implementation levels. Indeed, the choice of the computational-level description can be guided by examining the functions that can be yielded by implementational properties of the model, such as its architectural design, optimization strategy, or initialization procedure. In sum, proposals at the computational level can be formally tested to determine if they are possible in the real world (e.g., using the formal tools of theoretical computer science; Garey & Johnson, 1979; Downey & Fellows, 2013; Wareham, 1998), under relevant constraints such as the model architecture.

Empirically, researchers can study the behavior of interest to characterize the reliability and flexibility of the underlying capacity, and to determine any relevant restrictions. Behaviors of foundation models, for instance, cannot be taken for granted, even for simple problems (Bommasani et al., 2021). Established benchmarks, often named after capacities (“language understanding”, “commonsense reasoning”), do not always test these fully and might not address questions of interest (viz. Mitchell & Krakauer, 2023). Customized benchmark datasets are needed to properly test computational-level proposals (e.g., Moskvichev et al., 2023).

4.1.2. ALGORITHMIC DESCRIPTIONS

Multilevel explanations contain a description of the algorithms and data structures that implement the computational theory (Marr & Poggio, 1976). At this level, a sequence

of human-understandable steps that make the capacity possible is spelled out. Algorithmic descriptions can also be provided in the form of causal models (Geiger et al., 2023). Algorithms cannot be ascertained from the study of the capacity alone because a single capacity can be realized by different algorithms (e.g., Zhong et al., 2023).

The Attribute Extraction procedure (Algorithm 1) is a candidate algorithm for the computational problem \mathcal{D} -FACTUALRECALL (Def. 4.2), inspired by Chughtai et al. (2023). It consists of a series of steps that extract attributes related to the subject and relationship entities included in a factual recall prompt, and outputs the attribute that is highly associated to both the subject and relationship. The algorithm starts from an incomplete fact tuple containing input representations, subject S and relation R . The ATTRIBUTEBOOST subroutine takes as input an element of the incomplete tuple, and outputs a vector whose components correspond to a pre-defined vocabulary, where the entries associated with attributes of the input are numerically boosted. The ATTRIBUTECOMBINE step combines the boosted attribute vectors. ATTRIBUTEMAX then outputs the maximum-value attribute of the combined vectors. This description clarifies in human-understandable terms how the capacity of factual recall may be implemented step-by-step.

Algorithm 1 Factual Recall via Attribute Extraction

Input: incomplete fact tuple $(S, R, \perp) \in \mathcal{D}$

$H_s = \text{AttributeBoost}(S)$

$H_{s,r} = \text{AttributeBoost}(S, R)$

$H_r = \text{AttributeBoost}(R)$

$H_i = \text{AttributeCombine}(H_s, H_{s,r}, H_r)$

$A = \text{AttributeMax}(H_i)$

Return attribute A

4.1.3. PRIMITIVE REPRESENTATIONS AND OPERATIONS

The third level from the top specifies how the algorithm is executed using *primitive operations and representations* (Marr & Poggio, 1976; Poeppel, 2016). Primitives are the basic building blocks of the system, without which the phenomena cannot occur (Poeppel, 2016). They must be grounded theoretically, supported by abundant empirical evidence, and be realizable at the implementation level. Primitives, which can emerge through training (i.e., *post hoc*), should not be confused with the basic components of model architectures that were programmed before learning (e.g., single neurons, activation functions). It is possible that certain *post-hoc* primitives can emerge consistently in distinct models as a result of similar inductive biases given, for example, by architectures, training datasets, or learning rules.

In the factual recall example, a primitive candidate for ATTRIBUTEBOOST is a *key-value memory pair system* (see Fig. 2). Generally, these systems are comprised of a set of

paired vectors called *key* and *value*, where the key of a pair detects a pattern in the input (e.g. subject, relation, or their combination) and the value outputs associated information (e.g., attributes). Key-value systems emerge in transformer models across domains, models, tasks, and sub-modules (e.g., Geva et al., 2021; Meng et al., 2022; Vilas et al., 2023). Moreover, their existence is also theoretically motivated (Sukhbaatar et al., 2019). This makes the *key-value memory pair system* a suitable candidate for an operational primitive.

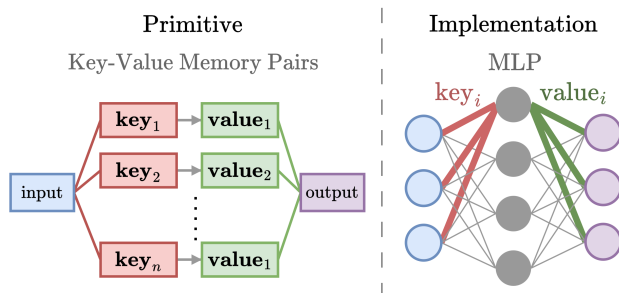


Figure 2. Examples of the *primitives* and *implementation* levels using the key-value memory pairs system in MLP layers.

4.1.4. IMPLEMENTATION

The *implementation level* characterizes how the primitives are implemented in the model. For this, the network must be decomposed at a certain level of abstraction. These choices need not match those of the architecture design before training (viz. Pylyshyn, 1978). The decomposition can be guided by the proposals at higher levels of explanation including the list of primitives. Exploratory work and the use of localization heuristics can also be used as guidance. Identifying model parts that have a causal effect on the output may help constrain the possible primitive operations, by examining the computations that can be implemented by these components (Bechtel, 2007). For example, *concept localization* efforts can be used to determine where relevant information emerges, and subsequently the function responsible for this emergence can be explored. Upper levels can then be revised to accommodate these empirical discoveries.

In the context of factual recall example, studies have demonstrated that key-value memory pair systems can both theoretically and practically be implemented by MLP or self-attention layers (viz. Dar et al., 2023; Geva et al., 2021; Vilas et al., 2023). Therefore, these sub-modules are suitable candidates to describe the implementation of key-value system primitives (see Fig. 2).

4.2. Building hypotheses and conducting severe tests

Previous work has already emphasized the need for rigorous hypothesis-testing procedures in interpretability research (Leavitt & Morcos, 2020; Rauker et al., 2023). While an

increasing number of studies are adopting these practices, there is ample room for improving how hypotheses are derived and empirically tested when evaluating mechanistic proposals. To begin with, candidate hypotheses should be tightly linked to mechanistic conjectures. Continuing with the factual recall example, one possible hypothesis could be: “extraction of subject attributes is conducted in subject tokens via key-value memory pair systems implemented by mid-layer MLP modules”. This only tackles one element of the mechanism but in a realistic experiment, hypotheses should be made about all components of the proposal.

Concrete and falsifiable empirical predictions are derived from mechanistic hypotheses. In the factual recall example, the following empirical predictions could be made: “one or more vectors in the first MLP parameter matrix will encode values that highly activate with subject tokens. The corresponding vectors of the second MLP matrix will produce distributional updates that promote attributes related to the subject.” Empirical predictions should differ from those of a reasonable baseline condition. For instance, they should minimally differ from those made about untrained models. Similarly, it is important to assess whether any other plausible mechanism could produce the same empirical predictions and to compare candidates. Competing mechanistic proposals should have true potential for explaining the capacity (Wilson & Collins, 2019; see e.g., Zhong et al., 2023).

Hypotheses can be evaluated using *severe tests* (Mayo, 2018), which have been recently introduced to the field of Cognitive Neuroscience to better evaluate how empirical practices and findings contribute to the assessment of models of human cognition (Aktunç, 2014; Bowers et al., 2023). Concretely, empirical observations and the methods used to evaluate them should have a high probability of falsifying a mechanistic proposal if it is incorrect. Severe tests can be extended to the field of AI (Bowers et al., 2023). Inner interpretability researchers should ask themselves: if the hypothesized mechanism is absent, how likely is the method to reveal its absence? Severe tests require work examining the adequacy of the methods themselves, for instance, by exploring whether they can falsify mechanistic hypotheses about a system whose design principles are known.

4.3. Designing experiments

Insights about the adequacy of mechanistic proposals can be gained by investigating input conditions that modulate the behavior of the model (Craver, 2006). Certain environmental conditions can precipitate the occurrence of a behavior, inhibit it (e.g., failure modes, see Hardcastle & Hardcastle, 2015), or modulate it. Experiments can mimic these conditions and evaluate whether hypothesized mechanisms effectively account for the modulation of the model’s capabilities.

Failure to provide adequate explanations and predictions suggests missing or erroneous elements in mechanistic proposals. As an example, models may fail to retrieve a fact because it was not available in the training set and thus not learned, or because the prompt was constructed in a way that led to failures in outputting the information (Jiang et al., 2020). Presumably, the mechanisms behind these failures are different, and mechanistic proposals should be able to differentiate them.

In addition to exploring these controlled experimental conditions, mechanisms should also be able to explain behavior under *naturalistic conditions*. For example, during model deployment, the way users construct prompts to extract factual information is more variable than those of controlled experiments. To be useful for the broader goals of AI, studies should ideally demonstrate that they can explain the manifestation of the capacity in these naturalistic settings.

4.4. Testing invariances

Multilevel explanations are often expected to generalize across conditions. Some of these conditions are explicit in the computational description. For instance, mechanisms are expected to remain invariant across input subdomains. In the factual recall case, the capacity could in principle be investigated separately for various fact subdomains (e.g., geography, politics, literature). However, the focus is to be kept on discovering the invariants, since these capacities are hypothesized to be manifestations of a more general capability of the system.

Other invariant conditions are left implicit in the computational explanation. For example, which kind of models are expected to implement the capacity is often left unsaid. Researchers may want to generalize their claims to narrower or broader classes of models (e.g. GPT-3, Large Language Models, Transformers). Typically, mechanisms invariant to initialization and hyperparameter values are sought (these often affect the inner workings of models; Zhong et al., 2023). Similarly, researchers look for mechanisms that remain valid when the models are scaled up. Whether mechanistic proposals can scale to describe larger models accurately can be supported or challenged by formal analyses (viz. Arora & Barak, 2009).

4.5. Refining mechanistic proposals and conceptual frameworks

Mechanistic proposals are to be continuously improved with new theoretical and empirical findings. When a change to one of the levels of explanation is made, all other levels need to be iteratively revised by imposing mutual constraints. Importantly, mechanistic proposals do not have to be complete to be tested or to be useful more generally. Mechanistic sketches (i.e. explanations that contain gaps), can point

to productive avenues for research (Bechtel, 2007; Craver, 2006), within and across research projects, *as long as their role in more comprehensive explanations is on the horizon*. Indeed, mechanistic proposals of factual recall in the Inner Interpretability literature have been increasingly refined and filled in, evidencing a productive line of work.

Moreover, it is essential to continuously revisit the utility of the adapted conceptual framework and methodological strategies for guiding the development of mechanistic explanations. Newer reappraisals of the multilevel framework, and how it is adapted to Inner Interpretability research, may lead to the formulation of better mechanistic explanations. For example, it has recently been suggested that separate levels should be added to the multilevel framework to describe how intelligent systems learn, and how such learning is the product of evolution (Poggio, 2012). However, adapting these levels to inner interpretability research entails challenges, since the learning constraints/goals and the conceptualization of evolution processes in biological systems cannot be easily extrapolated to that of artificial systems. Future work is needed to conceptualize this adaptation.

5. Implications for Inner Interpretability

The framework presented here can be used to examine the mechanisms of AI models regardless of the complexity of their architecture, training data, or task performed. Studies probing more complex systems and capacities may especially benefit from this framework, as it promotes simplified and human-understandable explanations that abstract from irrelevant implementational details and are generalizable to a variety of controlled and naturalistic contexts. More broadly, as detailed in the following subsections, the framework helps elucidate the state of the Inner Interpretability field and offers insights on how to move forward. In addition, adopting this framework helps tackle common criticisms of the field.

Situating previous studies. Studies in Inner Interpretability that appear hard to reconcile given their rationale and methods now have a coherent relationship in the multilevel framework proposed here. For example, some work can be understood as providing algorithmic-level descriptions of model capacities, as well as how they are realized at the implementation level (e.g., Chughtai et al., 2023; Merullo et al., 2023a; Wang et al., 2022; Zhong et al., 2023). Similarly, approaches like causal abstraction and structural equation modeling have been developed to evaluate how algorithmic descriptions fit the empirical data acquired at the implementation level (Beckers & Halpern, 2019; Chalupka et al., 2017; Geiger et al., 2023; Rubenstein et al., 2017). Other work can be viewed as focusing on uncovering primitive operations and representations across domains, architectures, and sub-modules (e.g., Geva et al., 2021; Olsson et al.,

2022; Vilas et al., 2023). Still, another body of work can be framed as developing heuristics for localization of causally relevant components at the implementation level, via automated circuit- and feature-finding procedures (e.g., Burns et al., 2022; Conmy et al., 2023; Gurnee et al., 2023).

Identifying research gaps. The framework also sheds light on potentially productive research avenues. Future directions derive directly from the aspects of the framework that are absent in current and past studies. Overall, there is a need for work formalizing capacities and theoretically analyzing their computational viability. A better understanding of the assumptions of common inner interpretability methods is also needed. More generally, studies proposing mechanistic sketches that span all four levels of the explanation are essential for progress.

The framework can also illustrate how a new research domain could be studied, or make it easy to locate gaps in some less-studied lines of research. For example, it can be used to identify aspects for improvement in studies probing the formation of abstract representations in neural networks. As explained in section 4.1.1, instead of providing descriptions of behaviors, future work could focus on formalizing capacities and later determining with theoretical and empirical work the behaviors that would be adequate reflections of them in a variety of scenarios. At the algorithmic level, Inner Interpretability studies have traditionally analyzed at which stages of the model hierarchy abstract representations emerge (e.g. Ilin et al., 2017). However, work investigating the sequence of algorithmic steps that build these abstract representations is still needed. Similarly, no work has been carried out to understand if the primitives put in use to form abstract representations differ from those of more concrete concepts, as suggested by research in Cognitive Neuroscience. At the implementation level, no studies have probed the right abstraction level or distributive nature of the model components supporting abstract representations.

Addressing criticisms. This framework offers guidance on how to address critiques arguing that the field lacks consensus and clarity on what a mechanistic understanding of a system is and how to build it. Regarding the criticism that inner interpretability methods have limited applicability to practical problems or realistic models, this work shifts the focus to studying well-defined capacities that link to real-world applications, and provides concrete strategies to do so. In addition, this framework helps avoid building a false sense of understanding and making misleading claims, by (i) encouraging a more comprehensive characterization of model capacities via multilevel analysis, inter-level constraints, and converging lines of evidence, (ii) emphasizing the use of methods involving severe tests to yield robust findings with assumption-aware interpretations.

Impact statement

The framework can be used to research the safety of the internal mechanisms of AI models. For example, the computational level can facilitate the fine-grained formalization and evaluation of the set of capacities (not behaviors) that the system should, or should not, possess according to the safety standards. The algorithmic level encourages spelling out and investigating if the algorithms that the system deploys are acceptable and guarantee safe behaviors of the model. In turn, the primitives and implementation levels call for research on how fragile the uncovered mechanisms are, and how they can be edited for safety reasons.

Overall, effective inner interpretability techniques should ultimately be useful to make systems safer to interact with, and more energy-efficient. However, other lines of work are more urgent than Inner Interpretability to address these issues. Tackling training data curation problems (viz. [Birhane et al., 2024](#)) and fostering responsible use with respect to carbon emissions (viz. [Luccioni & Hernandez-Garcia, 2023](#); [Luccioni et al., 2023](#)) are concrete avenues to deal with well-documented issues, whereas the benefits of inner interpretability in this context are currently more speculative.

Acknowledgment

This project was partly funded by the Ernst Strüngmann Foundation and the German Research Foundation (DFG) - DFG Research Unit FOR 5368. We are grateful for access to the computing facilities of the Center for Scientific Computing at Goethe University, and of the Ernst Strüngmann Institute for Neuroscience.

References

- Adolfi, F., Vilas, M., and Wareham, T. Complexity-theoretic limits on the promises of artificial neural network reverse-engineering. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, 2024.
- Aktunç, M. E. Severe tests in neuroimaging: what we can learn and how we can learn it. *Philosophy of Science*, 81(5):961–973, 2014.
- Akyürek, E., Wang, B., Kim, Y., and Andreas, J. In-context language learning: Architectures and algorithms, 2024.
- Arora, S. and Barak, B. *Computational complexity: a modern approach*. Cambridge University Press, Cambridge ; New York, 2009. ISBN 978-0-521-42426-4.
- Barack, D. L. and Krakauer, J. W. Two views on the cognitive brain. *Nature Reviews Neuroscience*, 22(6):359–371, 2021.
- Bechtel, W. *Mental mechanisms: Philosophical perspectives on cognitive neuroscience*. Psychology Press, 2007.
- Beckers, S. and Halpern, J. Y. Abstracting causal models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 2678–2685, 2019.
- Belrose, N., Furman, Z., Smith, L., Halawi, D., Ostrovsky, I., McKinney, L., Biderman, S., and Steinhardt, J. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112*, 2023.
- Birhane, A., Han, S., Boddeti, V., Luccioni, S., et al. Into the LAION’s den: Investigating hate in multimodal datasets. *Advances in Neural Information Processing Systems*, 36, 2024.
- Bolukbasi, T., Pearce, A., Yuan, A., Coenen, A., Reif, E., Viégas, F., and Wattenberg, M. An interpretability illusion for BERT. *arXiv preprint arXiv:2104.07143*, 2021.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Bowers, J. S., Malhotra, G., Dujmović, M., Montero, M. L., Tsvetkov, C., Biscione, V., Puebla, G., Adolfi, F., Hummel, J. E., Heaton, R. F., Evans, B. D., Mitchell, J., and Blything, R. Deep Problems with Neural Network Models of Human Vision. *Behavioral and Brain Sciences*, pp. 1–74, December 2022.
- Bowers, J. S., Malhotra, G., Adolfi, F., Dujmović, M., Montero, M. L., Biscione, V., Puebla, G., Hummel, J. H., and Heaton, R. F. On the importance of severely testing deep learning models of cognition. *Cognitive Systems Research*, 82:101158, 2023.
- Burns, C., Ye, H., Klein, D., and Steinhardt, J. Discovering latent knowledge in language models without supervision. *arXiv preprint arXiv:2212.03827*, 2022.
- Buzsáki, G. The Brain–Cognitive Behavior Problem: A Retrospective. *eNeuro*, 7(4), July 2020.
- Calhoun, A. and Hady, A. E. What is behavior? no seriously, what is it? *BioRxiv*, pp. 2021–07, 2021.
- Casper, S., Ezell, C., Siegmann, C., Kolt, N., Curtis, T. L., Bucknall, B., Haupt, A., Wei, K., Scheurer, J., Hobbhahn, M., Sharkey, L., Krishna, S., Hagen, M. V., Alberti, S., Chan, A., Sun, Q., Gerovitch, M., Bau, D., Tegmark, M., Krueger, D., and Hadfield-Menell, D. Black-box access is insufficient for rigorous ai audits. *arXiv preprint arXiv:2401.14446*, 2024.

- Chalupka, K., Eberhardt, F., and Perona, P. Causal feature learning: an overview. *Behaviormetrika*, 44:137–164, 2017.
- Chughtai, B., Cooney, A., and Nanda, N. Summing up the facts: Additive mechanisms behind factual recall in LLMs. 2023.
- Churchland, P. S. and Sejnowski, T. *The Computational Brain*. Computational Neuroscience. MIT Press, Cambridge, Mass., 1999. ISBN 978-0-262-53120-7 978-0-262-03188-2.
- Conmy, A., Mavor-Parker, A. N., Lynch, A., Heimersheim, S., and Garriga-Alonso, A. Towards automated circuit discovery for mechanistic interpretability. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Craver, C. F. When mechanistic models explain. *Synthese*, 153(3):355–376, 2006.
- Danks, D. Moving from levels & reduction to dimensions & constraints. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 35, 2013.
- Dar, G., Geva, M., Gupta, A., and Berant, J. Analyzing transformers in embedding space. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023.
- Doshi-Velez, F. and Kim, B. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.
- Downey, R. G. and Fellows, M. R. *Fundamentals of parameterized complexity*. Texts in computer science. Springer, London, 2013. ISBN 978-1-4471-5559-1.
- Egan, F. Function-theoretic explanation and the search for neural mechanisms. In *Explanation and Integration in Mind and Brain Science*, pp. 145–163. Oxford University Press, 2018.
- Elhage, N., Nanda, N., Olsson, C., Henighan, T., Joseph, N., Mann, B., Askell, A., Bai, Y., Chen, A., Conerly, T., et al. A mathematical framework for transformer circuits, 2021. URL <https://transformer-circuits.pub/2021/framework/index.html>.
- Elhage, N., Hume, T., Olsson, C., Schiefer, N., Henighan, T., Kravec, S., Hatfield-Dodds, Z., Lasenby, R., Drain, D., Chen, C., et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.
- Garey, M. R. and Johnson, D. S. *Computers and intractability*. W.H. Freeman, 1979.
- Geiger, A., Potts, C., and Icard, T. Causal abstraction for faithful model interpretation. *arXiv preprint arXiv:2301.04709*, 2023.
- Geva, M., Schuster, R., Berant, J., and Levy, O. Transformer feed-forward layers are key-value memories. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2021.
- Geva, M., Bastings, J., Filippova, K., and Globerson, A. Dissecting recall of factual associations in auto-regressive language models. *arXiv preprint arXiv:2304.14767*, 2023.
- Giraud, A.-L. and Poeppel, D. Cortical oscillations and speech processing: emerging computational principles and operations. *Nature Neuroscience*, 15(4):511–517, April 2012.
- Gomez-Marin, A. Causal circuit explanations of behavior: Are necessity and sufficiency necessary and sufficient? *Decoding Neural Circuit Structure and Function: Cellular Dissection Using Genetic Model Organisms*, pp. 283–306, 2017.
- Gomez-Marin, A. Promisomics and the Short-Circuiting of Mind. *eNeuro*, 8(2), March 2021.
- Griffiths, T. L., Chater, N., Kemp, C., Perfors, A., and Tenenbaum, J. B. Probabilistic models of cognition: Exploring representations and inductive biases. *Trends in cognitive sciences*, 14(8):357–364, 2010.
- Gunaratne, C. A., Sakurai, A., and Katz, P. S. Variations on a theme: species differences in synaptic connectivity do not predict central pattern generator activity. *Journal of Neurophysiology*, 118(2):1123–1132, 2017.
- Gurnee, W., Nanda, N., Pauly, M., Harvey, K., Troitskii, D., and Bertsimas, D. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*, 2023.
- Hamrick, J. and Mohamed, S. Levels of analysis for machine learning. *arXiv preprint arXiv:2004.05107*, 2020.
- Hardcastle, V. G. and Hardcastle, K. Marr’s levels revisited: understanding how brains break. *Topics in Cognitive Science*, 7(2):259–273, 2015.
- Hase, P., Bansal, M., Kim, B., and Ghandeharioun, A. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Hernandez, E., Schwettmann, S., Bau, D., Bagashvili, T., Torralba, A., and Andreas, J. Natural language descriptions of deep visual features. In *International Conference on Learning Representations*, 2021.

- Hernandez, E., Li, B. Z., and Andreas, J. Inspecting and editing knowledge representations in language models. In *Arxiv*, 2023. URL <https://arxiv.org/abs/2304.00740>.
- Hickok, G. and Poeppel, D. The cortical organization of speech processing. *Nature reviews neuroscience*, 8(5): 393–402, 2007.
- Hooker, S. The hardware lottery. *Communications of the ACM*, 64(12):58–65, 2021.
- Ilin, R., Watson, T., and Kozma, R. Abstraction hierarchy in deep learning neural networks. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pp. 768–774. IEEE, 2017.
- Jiang, Z., Xu, F. F., Araki, J., and Neubig, G. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438, 2020.
- Kaplan, D. M. and Craver, C. F. The explanatory force of dynamical and mathematical models in neuroscience: A mechanistic perspective. *Philosophy of science*, 78(4): 601–627, 2011.
- Kim, B., Wattenberg, M., Gilmer, J., Cai, C., Wexler, J., Viegas, F., et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, pp. 2668–2677. PMLR, 2018.
- Krakauer, J. W., Ghazanfar, A. A., Gomez-Marin, A., MacIver, M. A., and Poeppel, D. Neuroscience needs behavior: correcting a reductionist bias. *Neuron*, 93(3): 480–490, 2017.
- Krishnan, M. Against interpretability: a critical examination of the interpretability problem in machine learning. *Philosophy & Technology*, 33(3):487–502, 2020.
- Leavitt, M. L. and Morcos, A. Towards falsifiable interpretability research. *arXiv preprint arXiv:2010.12016*, 2020.
- Lepori, M. A., Serre, T., and Pavlick, E. Uncovering intermediate variables in transformers using circuit probing. *arXiv e-prints*, pp. arXiv–2311, 2023.
- Lieder, F. and Griffiths, T. L. Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behavioral and brain sciences*, 43:e1, 2020.
- Lindsay, G. W. and Bau, D. Testing methods of neural systems understanding. *Cognitive Systems Research*, 82: 101156, 2023.
- Love, B. C. Levels of biological plausibility. *Philosophical Transactions of the Royal Society B*, 376(1815):20190632, 2021.
- Luccioni, A. S. and Hernandez-Garcia, A. Counting Carbon: A Survey of Factors Influencing the Emissions of Machine Learning, February 2023. URL <http://arxiv.org/abs/2302.08476>. arXiv:2302.08476 [cs].
- Luccioni, A. S., Jernite, Y., and Strubell, E. Power Hungry Processing: Watts Driving the Cost of AI Deployment?, November 2023. URL <http://arxiv.org/abs/2311.16863>. arXiv:2311.16863 [cs].
- Machamer, P., Darden, L., and Craver, C. F. Thinking about mechanisms. *Philosophy of science*, 67(1):1–25, 2000.
- Marr, D. and Poggio, T. From understanding computation to understanding neural circuitry. *Neuroscience Research Program Bulletin*, 15(3):470–488, 1976.
- Mayo, D. G. *Statistical inference as severe testing: How to get beyond the statistics wars*. Cambridge University Press, 2018.
- McDougall, C., Conmy, A., Rushing, C., McGrath, T., and Nanda, N. Copy suppression: Comprehensively understanding an attention head. *arXiv preprint arXiv:2310.04625*, 2023.
- Meng, K., Bau, D., Andonian, A., and Belinkov, Y. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- Merullo, J., Eickhoff, C., and Pavlick, E. Circuit component reuse across tasks in transformer language models. *arXiv preprint arXiv:2310.08744*, 2023a.
- Merullo, J., Eickhoff, C., and Pavlick, E. Language models implement simple word2vec-style vector arithmetic. *arXiv preprint arXiv:2305.16130*, 2023b.
- Mitchell, M. Why AI is harder than we think. In *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO '21*, pp. 3, New York, NY, USA, June 2021. Association for Computing Machinery.
- Mitchell, M. and Krakauer, D. C. The debate over understanding in AI’s large language models. *Proceedings of the National Academy of Sciences*, 120(13):e2215907120, March 2023.
- Moskvichev, A., Odouard, V. V., and Mitchell, M. The ConceptARC Benchmark: Evaluating Understanding and Generalization in the ARC Domain, May 2023. URL <http://arxiv.org/abs/2305.07141>. arXiv:2305.07141 [cs].

- Mu, J. and Andreas, J. Compositional explanations of neurons. *Advances in Neural Information Processing Systems*, 33:17153–17163, 2020.
- Niv, Y. The primacy of behavioral research for understanding the brain. *Behavioral Neuroscience*, 135(5):601–609, October 2021.
- Olah, C. Interpretability dreams, 2023. URL <https://transformer-circuits.pub/2023/interpretability-dreams/index.html>.
- Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A., et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- Poepfel, D. The maps problem and the mapping problem: two challenges for a cognitive neuroscience of speech and language. *Understanding Cognitive Development*, pp. 34–55, 2016.
- Poepfel, D. and Adolfs, F. Against the epistemological primacy of the hardware: The brain from inside out, turned upside down. *Eneuro*, 7(4), 2020.
- Poepfel, D. and Assaneo, M. F. Speech rhythms and their neural foundations. *Nature Reviews Neuroscience*, 21(6): 322–334, June 2020.
- Poggio, T. The levels of understanding framework, revised. *Perception*, 41(9):1017–1023, 2012.
- Poldrack, R. A. Can cognitive processes be inferred from neuroimaging data? *Trends in cognitive sciences*, 10(2): 59–63, 2006.
- Pylyshyn, Z. W. Computational models and empirical constraints. *Behavioral and brain sciences*, 1(1):91–99, 1978.
- Räuker, T., Ho, A., Casper, S., and Hadfield-Menell, D. Toward transparent ai: A survey on interpreting the inner structures of deep neural networks. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pp. 464–483. IEEE, 2023.
- Richards, B. A., Lillicrap, T. P., Beaudoin, P., Bengio, Y., Bogacz, R., Christensen, A., Clopath, C., Costa, R. P., de Berker, A., Ganguli, S., et al. A deep learning framework for neuroscience. *Nature neuroscience*, 22(11): 1761–1770, 2019.
- Ross, L. N. and Bassett, D. S. Causation in neuroscience: keeping mechanism meaningful. *Nature Reviews Neuroscience*, 2024.
- Rubenstein, P. K., Weichwald, S., Bongers, S., Mooij, J. M., Janzing, D., Grosse-Wentrup, M., and Schölkopf, B. Causal consistency of structural equation models. *arXiv preprint arXiv:1707.00819*, 2017.
- Shagrir, O. and Bechtel, W. Marr’s Computational Level and Delineating Phenomena. In Kaplan, D. M. (ed.), *Explanation and Integration in Mind and Brain Science*, pp. 0. Oxford University Press, 2017.
- Sporns, O., Tononi, G., and Kötter, R. The Human Connectome: A Structural Description of the Human Brain. *PLOS Computational Biology*, 1(4):e42, 2005.
- Storrs, K. R. and Kriegeskorte, N. Deep learning for cognitive neuroscience. *arXiv preprint arXiv:1903.01458*, 2019.
- Strobl, L. Average-Hard Attention Transformers are Constant-Depth Uniform Threshold Circuits, August 2023. URL <http://arxiv.org/abs/2308.03212>. arXiv:2308.03212 [cs].
- Strobl, L., Merrill, W., Weiss, G., Chiang, D., and Angluin, D. Transformers as Recognizers of Formal Languages: A Survey on Expressivity, October 2023. URL <http://arxiv.org/abs/2311.00208>. arXiv:2311.00208 [cs].
- Sucholutsky, I., Muttenthaler, L., Weller, A., Peng, A., Bobu, A., Kim, B., Love, B. C., Grant, E., Achterberg, J., Tenenbaum, J. B., et al. Getting aligned on representational alignment. *arXiv preprint arXiv:2310.13018*, 2023.
- Sukhbaatar, S., Grave, E., Lample, G., Jegou, H., and Joulin, A. Augmenting self-attention with persistent memory. *arXiv preprint arXiv:1907.01470*, 2019.
- Todd, E., Li, M. L., Sharma, A. S., Mueller, A., Wallace, B. C., and Bau, D. Function vectors in large language models. In *Proceedings of the 2024 International Conference on Learning Representations*, 2024.
- Vig, J., Gehrmann, S., Belinkov, Y., Qian, S., Nevo, D., Sakenis, S., Huang, J., Singer, Y., and Shieber, S. Causal mediation analysis for interpreting neural NLP: The case of gender bias. *arXiv preprint arXiv:2004.12265*, 2020.
- Vilas, M. G., Schaumlöffel, T., and Roig, G. Analyzing vision transformers for image classification in class embedding space. In *Advances in Neural Information Processing Systems*, volume 36, pp. 40030–40041, 2023.
- Wang, K., Variengien, A., Conmy, A., Shlegeris, B., and Steinhardt, J. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv preprint arXiv:2211.00593*, 2022.

- Wareham, H. T. *Systematic Parameterized Complexity Analysis in Computational Phonology*. PhD thesis, University of Victoria, Canada, 1998. URL <http://roa.rutgers.edu/files/318-0599/roa-318-wareham-2.pdf>.
- Wilson, R. C. and Collins, A. G. Ten simple rules for the computational modeling of behavioral data. *eLife*, 8: e49547, 2019.
- Yu, Q., Merullo, J., and Pavlick, E. Characterizing mechanisms for factual recall in language models. *arXiv preprint arXiv:2310.15910*, 2023.
- Zhang, E., Lepori, M. A., and Pavlick, E. Instilling inductive biases with subnetworks. *arXiv preprint arXiv:2310.10899*, 2023.
- Zhong, Z., Liu, Z., Tegmark, M., and Andreas, J. The clock and the pizza: Two stories in mechanistic explanation of neural networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Zou, A., Phan, L., Chen, S., Campbell, J., Guo, P., Ren, R., Pan, A., Yin, X., Mazeika, M., Dombrowski, A.-K., et al. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.

A. Critical concepts and references from Cognitive Neuroscience

Below we provide a discussion of concepts used in the Cognitive Neuroscience field that are critical for the proposed AI Inner Interpretability framework:

- **Mechanistic explanations** in neuroscience are described as detailing the entities (or parts), activities (or operations), properties, organizational features (both temporal and spatial), and causal relationships among these components, that produce the target phenomena (Bechtel, 2007; Craver, 2006; Kaplan & Craver, 2011; Machamer et al., 2000). The entities or parts “[...] are the things that engage in activities” (Machamer et al., 2000), and “[...] are the structural components of the mechanism” (Bechtel, 2007). The activities or operations “[...] are the producers of change” (Machamer et al., 2000), and “[...] refer to processes or changes involving the parts” (Bechtel, 2007).
- **Mechanistic sketches** are incomplete models of a mechanism, and signal that further work is needed (Craver, 2006; Machamer et al., 2000). These sketches contain gaps where certain components of the mechanistic explanation, such as entities or activities, are missing. Gaps are sometimes masked by filler terms. For example, terms like cause, encode, produce, and represent “[...] are often used to indicate a kind of activity in a mechanism without providing any detail about how that activity is carried out” (Craver, 2006).
- The target of explanation in Cognitive Neuroscience is the neural implementation of a cognitive capacity. A **capacity** is an underlying ability of a system or organism to transform certain input to output states (Egan, 2018). It can be fully defined at the computational level of analysis by specifying the input domain and the function that maps inputs to outputs. A **behavior** can be understood as the concrete and observable manifestation of a capacity. In the literature, the term behavior is often used to refer to a behavioral phenomenon, a set of concrete actions by a system or organism that may reflect the performance of an underlying capacity. Defining what constitutes a behavior is challenging (Calhoun & Hady, 2021). Behaviors of interest are chosen based on observations and theoretical arguments that they represent a manifestation of the target cognitive capacity.
- It has been argued that cognition can be mechanistically explained as a sequence of **computational operations** performed over **representations** (Bechtel, 2007). Although the definition and properties of the term ‘representation’ continue to be a matter of debate in the field, they are frequently conceptualized as states of the neural system that carry information about external objects or events relevant to the capacity being explained (Bechtel, 2007). A list of the elementary mental representations and operations has been called the *human cognome* and corresponds to the **primitive** units on analysis of the cognitive sciences “[...] without which they could not account for the elementary phenomena of their field” (Poehpel, 2016).
- Different **levels of explanation** can be used to analyze the mechanisms of intelligent systems. Marr & Poggio (1976) propose that four “nearly independent” levels of description can be used to study machines that solve an information processing problem: “(1) that at which the nature of a computation is expressed; (2) that at which the algorithms that implement a computation are characterized; (3) that at which an algorithm is committed to particular mechanisms; and (4) that at which the mechanisms are realized in hardware”. Marr & Poggio (1976) put special emphasis on the importance of the computational level, which is often neglected in Cognitive Neuroscience studies. Later work removed the third level from the multilevel framework, while in this work we re-conceptualize it as the encoding of cognitive parts that need to be mapped to other levels. Furthermore, recent reappraisals of this framework have highlighted the importance of utilizing mutual constraints among the levels to achieve better explanations of the cognitive system, a point we support in this position paper.
- Neuroscience studies can choose to provide a mechanistic explanation at different **levels of neural organization**. For example, some studies seek to explain cognition “[...] as the result of operations on signals performed at nodes in a network and passed between them that are implemented by specific neurons and their connections in circuits in the brain” (Barack & Krakauer, 2021), while others explain cognition “[...] as the result of transformations between or movement within representational spaces that are implemented by neural populations” (Barack & Krakauer, 2021). The levels of neural organization are orthogonal to the levels of explanation. Although only at the implementation level the neural components are explicitly examined, their choice implies commitments that constrain the possible explanations at other levels.
- In Cognitive Neuroscience, it has been greatly debated whether neural mechanisms should be studied before or after having decomposed and analyzed the behavior of interest at the computational and algorithmic levels (i.e. **bottom-up**

or **top-down** approaches, see Fig. 3 and section 3.3). For example, in favor of a top-down approach, Krakauer et al. (2017) argue that “higher-level concepts are needed to understand neuronal results” (higher-level concepts are those derived from behavioral work) and provide a variety of examples of how “behaviorally driven neuroscience yields more complete insights”. In contrast, bottom-up proponents like Buzsáki (2020) argue that “[...] most of our behavior-related terms emerged before and independent of neuroscience, and there is little guarantee that these terms correspond to circumscribed brain mechanisms”. In their view, the field should instead “[...] start with the brain (independent variable) and define descriptors of behavior (dependent variables) that are free from philosophical connotations” (Buzsáki, 2020). In short, we can “[...] recast the inherent tension between these epistemic procedures as that between *What is a mechanism for X?* versus *What is Y a mechanism for?*” (Poepfel & Adolfi, 2020). Recently, it has been argued that more pluralistic approaches that combine methods from both research positions may lead to more robust mechanistic theories: “[...] we might view the process of doing research in the field of cognitive neuroscience as the iterative abduction of certain kinds of mechanistic theories about human capacities” (Poepfel & Adolfi, 2020), where abduction “[...] jointly captures the process by which a set of candidate explanations is generated from observations and background knowledge (sometimes termed abduction proper), and how the choice among them is justified” (Poepfel & Adolfi, 2020).

B. Key ideas and models of bottom-up and top-down approaches in Cognitive Neuroscience

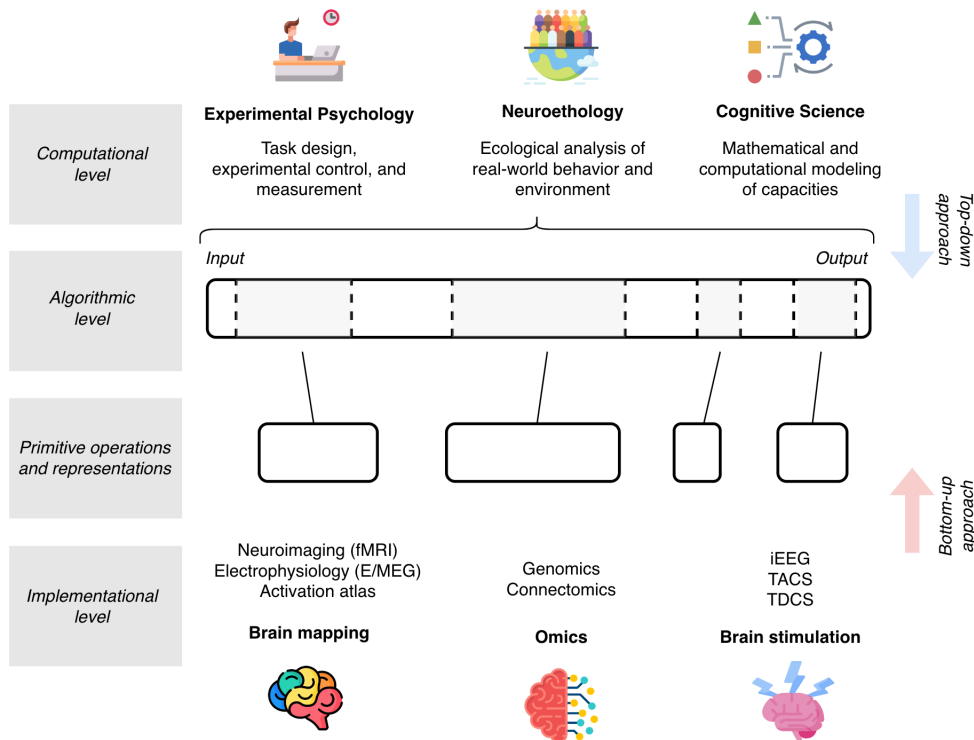


Figure 3. Schematic of the key techniques deployed by different disciplines in Cognitive Neuroscience organized according to whether they promote top-down or bottom-up approaches to the discovery of mechanistic explanations. Note that only radical (top-down/bottom-up) approaches propose to reach inner levels through a one-directional use of these techniques. In practice, the discovery of mechanistic explanations involves a healthy combination of techniques from top-down and bottom-up approaches.

C. Separability of the levels of explanation

Certain levels of the multi-level explanation framework are more easily distinguishable from neighboring levels than others. The computational and algorithmic levels provide descriptions that are formally distinguished in Computer Science, and as such cannot be easily confounded. In contrast, the algorithmic and primitive levels use a similar vocabulary. Determining

whether a particular sub-computation is a primitive relies on the theoretical and empirical evidence supporting its role as a building block of the system. Finally, the implementation level is also easily distinguishable from other levels, as it provides descriptions using the terminology used when designing the model (e.g. MLP layers, self-attention heads, activation functions, etc.).

D. Using AI to understand biological cognitive functions

Previous work has also leveraged the similarities between AI research and Cognitive Neuroscience by employing models and techniques from the AI field to better understand biological cognitive functions. Among other proposals, it has been suggested that deep learning models can serve as tools for testing cognitive theories (Storrs & Kriegeskorte, 2019). For example, they can be used to probe the learning rules, goals, and anatomical properties of the brain (Richards et al., 2019). Moreover, they can be employed to test the structure and content of cognitive representations in the human brain (Sucholutsky et al., 2023). Regarding AI techniques, it has recently been proposed that tools to interpret neural networks can be used to test the methods in neuroscience (Lindsay & Bau, 2023).