# Naturalistic Computational Cognitive Science

## Towards generalizable models and theories
## that capture the full range of natural behavior

**Wilka Carvalho**[*]
Kempner Institute
Harvard University
wcarvalho@g.harvard.edu
[*]equal contribution

**Andrew Lampinen**[*]
Google DeepMind

lampinen@google.com

### Abstract

Artificial Intelligence increasingly pursues large, complex models that perform many tasks within increasingly realistic domains. How, if at all, should these developments in AI influence cognitive science? We argue that progress in AI offers timely opportunities for cognitive science to embrace experiments with increasingly naturalistic stimuli, tasks, and behaviors; and computational models that can accommodate these changes. We first review a growing body of research spanning neuroscience, cognitive science, and AI that suggests that incorporating a broader range of naturalistic experimental paradigms (and models that accommodate them) may be necessary to resolve some aspects of natural intelligence and ensure that our theories generalize. We then suggest that integrating recent progress in AI and cognitive science will enable us to engage with more naturalistic phenomena without giving up experimental control or the pursuit of theoretically grounded understanding. We offer practical guidance on how methodological practices can contribute to cumulative progress in naturalistic computational cognitive science, and illustrate a path towards building computational models that solve the real problems of natural cognition—together with a reductive understanding of the processes and principles by which they do so.

## 1  Introduction

Cognitive scientists build models to make our theories concrete — which offers testable predictions, eliminates ambiguities (Guest and Martin, 2021), and can reveal unexpected properties of cognition and behavior (McClelland, 2009). Cognitive models can range from simplified models of a high-level process (e.g. Frank and Goodman, 2012; Cohen et al.,

1

1990), to task-performing models of a particular domain (Newell, 1973), or beyond (Newell, 1994). Many cognitive paradigms—connectionism (McClelland et al., 1986), bayesian inference (Tenenbaum and Griffiths, 2001), and cognitive architectures (e.g. Ritter et al., 2019; Laird, 2019)—have sought generalizable models that can explain a broad range of behavior from a simpler set of principles. Ultimately, these approaches could aspire to build towards "unified theories of cognition" that can explain and predict natural behavior across many tasks and domains (Newell, 1994). However, we are not there yet. In recent years there have been substantial debates about the generalizability of our theories (Yarkoni, 2022; Eckstein et al., 2021), whether our models and theories are adequately constrained (e.g. Jones and Love, 2011; Rahnev and Denison, 2018), and calls to accommodate a broader scope of naturalistic phenomena (Wise et al., 2023; Cisek and Green, 2024). Our goal in this paper is to motivate and outline a path that we believe will begin to address many of these challenges.
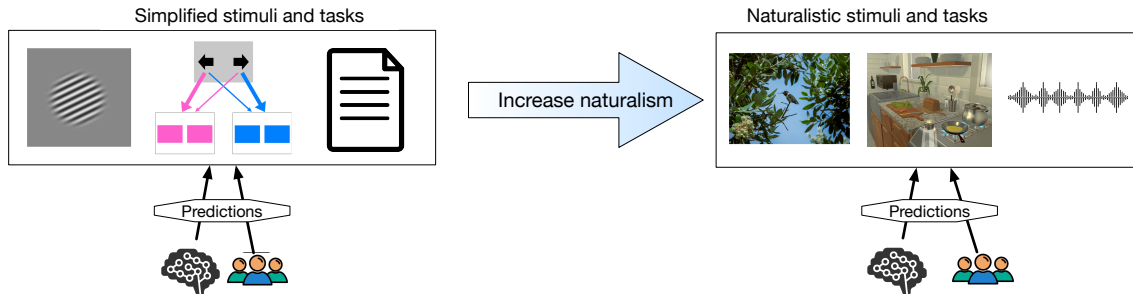
The path we propose is motivated by recent practical and conceptual developments in AI. In the last ten years AI has progressed substantially in building models that can perform a wide range of tasks in various domains. We now have vision models that can classify and segment real-world images in objects without external supervision (He et al., 2020; Caron et al., 2021), we have reinforcement learning models that can learn new tasks in a domain with human-like efficiency (Team et al., 2023), and we have language models which can produce human-like language and solve many language-specified tasks (Radford et al., 2019; Wei et al., 2022). Surprisingly, such models produce internal representations that capture many features of brain representations (Yamins et al., 2014; Schrimpf et al., 2021). These findings have led to discussion of how such complex models can be explanatory (Cao and Yamins, 2024a), and how they should alter our theories (Hasson et al., 2020; Perconti and Plebe, 2020; Piantadosi, 2023). Yet they have also led to debate over how we test our models and interpret the results (Bowers et al., 2023; DiCarlo et al., 2023; Dentella et al., 2024). Thus, it remains an open question if and how the progress in AI can contribute to addressing the challenges of cognitive science.

In this work, we attempt to weave a thread of arguments that synthesizes these literatures, spanning from cognitive science's motivation and theory development to the practicalities of model engineering and reproducibility. We articulate a research strategy that we believe will be helpful in enabling theory-driven cognitive science to achieve deeper understanding of the full range of natural intelligence. We present an overview of this framework and our arguments in Figure 1.

Our paper proceeds as follows. First, in §2, we unpack what we mean by "naturalistic" computational cognitive science, by providing a pragmatic outline of the features that we believe deserve more emphasis in our experimental paradigms and models. In particular, we argue that researchers should expand their task paradigms to to more closely approximate the breadth of the settings in which their theoretical constructs would be expected to generalize by both adding relevant variables that may interact with those in question, and by incorporating more variability within existing parameters. Likewise, researchers should expand their models
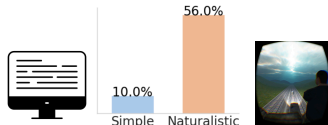
**What?**

**Section 2:** Develop **learning-based** models of intelligence that can predict human behavior on both simplified and naturalistic stimuli and tasks
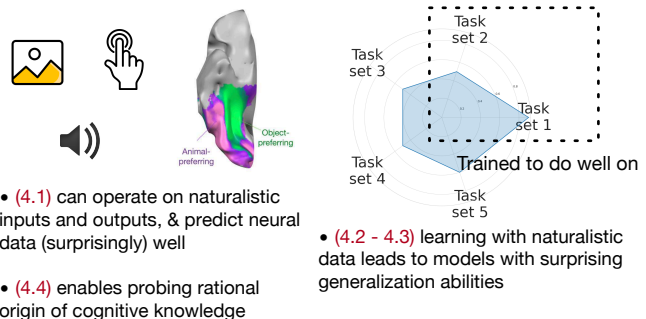
Simplified stimuli and tasks

Increase naturalism

Naturalistic stimuli and tasks

Predictions

Predictions

**Why?**

**Section 3:** Why increase the naturalism of our experimental paradigms?
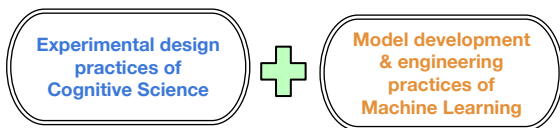
56.0%

10.0%

Simple Data

Naturalistic Data

- (3.1) Learning and behavior can be different

- (3.2) Brain systems can be engaged differently

- (3.3) Simpler settings can lose ecological validity with the phenomena they target

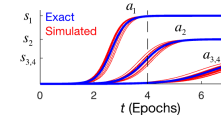**Section 4:** Why *learn* models of intelligence?

Animal-preferring

Object-preferring

Task set 1

Task set 2

Task set 3

Task set 4

Task set 5

Trained to do well on

- (4.1) can operate on naturalistic inputs and outputs, & predict neural data (surprisingly) well

- (4.4) enables probing rational origin of cognitive knowledge

- (4.2 - 4.3) learning with naturalistic data leads to models with surprising generalization abilities

**How?**

**Section 5:** How do we build generalizable models that scale to naturalistic settings?

**Experimental design practices of Cognitive Science**

+

**Model development & engineering practices of Machine Learning**

- (5.1) frictionless reproducibility

- (5.2) generalization-focused model development

- (5.3) Embrace phenomena-oriented benchmarks to catalyze progress

**Section 6:** How do we build and test theories with naturalistic tasks and models?

$s_1$
$s_2$
$s_{3,4}$

Exact
Simulated

$a_1$
$a_2$
$a_{3,4}$

0   2   4   6
$t$ (Epochs)

- (6.1 - 6.2) Controlled experiments by parametrically generating stimuli and tasks while maintaining naturalistic variation to enhance generalizability

- (6.3 - 6.4) Theories that combine predictive models with reductive explanations

Figure 1: **Naturalistic computational cognitive science**: the what, why, and the how. The first section (§2) provides an overview of "naturalistic computational cognitive science". Afterwards (§3-§4), we motivate the utility of naturalistic experimental paradigms and learning-based approached for developing a more complete understanding of cognition. The rest of the paper (§5-§6) focuses on **how** to achieve these goals; how to develop models for naturalistic settings, and how to use naturalistic experiments and models as part of explanatory cognitive theories. (Figures reproduced from Saxe et al., 2019; Geirhos et al., 2018; Doshi and Konkle, 2023.)

to more faithfully approximate the computational task that the natural system is solving.

In the next two sections, we motivate these goals, by arguing that increasing naturalism is necessary for developing generalizable understanding. First, in §3, we review strands of evidence spanning neuroscience, cognitive science, and AI, showing that naturalistic experimental paradigms can lead to different behavior, engage brain systems differently, and expose expose computational challenges that engage mechanisms differently. As an example, Francis et al. (2017) found that moral judgments on a trolley problem were dramatically different in a virtual reality presentation compared to text vignettes. Other works that we review have found that processes ranging from visual processing (Amme et al., 2024), learning (Rosenberg et al., 2021; Collins, 2024), memory (Helbing et al., 2020) and more are changed when naturalistic elements are introduced. The works that we review in this section thus suggest that experiments that span a broader range of naturalistic settings are essential for drawing correct inferences about the system as a whole. Next, in §4, we highlight the benefits of models that learn from naturalistic data—which can perform a wide range of tasks, and yield qualitatively different patterns of generalization than models trained in simplified settings. In turn, these differences can change the inferences we make from our models, and yield new insights about the origins of cognitive and neural phenomena. Together, these two sections of our paper argue that pursuing naturalistic experiments, and models trained on them, will be necessary for achieving complete understanding of cognition.

In the remainder of the paper, we turn to the practice of naturalistic computational cognitive science, and its contributions to building generalizable cognitive theories. First, in §5, we review key lessons computational cognitive scientists can adopt from AI in developing generalizable models that can operate in increasingly naturalistic data conditions: the benefits of frictionlessly reproducing prior research artifacts, the importance of generalizability-focused model development, and the value of benchmarking for concentrating research effort on important questions. Finally, in §6, we outline how one can leverage these potentially complicated and opaque models to contribute to generalizable theories of cognition—uniting task-performing predictive models with explanatory reductive understanding, using established techniques like rational analysis (Anderson, 1990) and new approaches to interpreting complex models (Geiger et al., 2021). We close in §7 by discussing the broader context and implications of naturalistic computational cognitive science.

## 2 What is "naturalistic" computational cognitive science?

Naturalistic computational cognitive science is a research strategy for theory-driven cognitive science that aims to predict human behavior across increasingly naturalistic stimuli and tasks. This approach gradually increases the ecological validity of experimental designs while ensuring models can accommodate both simplified and more natural conditions. As demonstrated in machine learning (§5), this strategy has enabled models that can work with real-world data, albeit imperfectly.
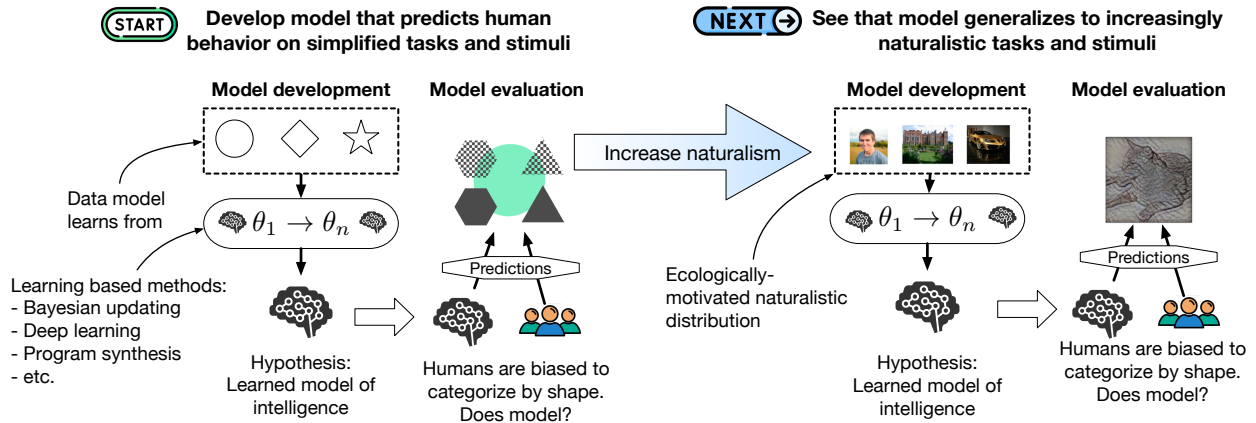
Figure 2: Overview of naturalistic computational cognitive science. Given a novel computational model, researchers generate predictions to test their model in a simplified setting. If this model successfully predicts human behavior, researchers see if their model continues to generalize human predictions in increasingly naturalisic conditions.

As in other approaches within computational cognitive science, this approach starts with formulating models and designing controlled experiments to isolate core phenomena. Like prior work, we place an emphasis on *model generalization*, i.e. predicting human behavior on data the model hasn't been exposed to—including data from new experimental conditions or tasks (Busemeyer and Wang, 2000). Therefore,

1. Model parameters are estimated with a "training" distribution before being evaluated.

2. Models should predict human behavior on held-out examples that neither the model nor humans have encountered.

One key facet of naturalistic computational cognitive science is that we place emphasis on having models—in particular, the same model—predict human behavior across increasingly naturalistic conditions that cover a broader space of settings in which the theoretical construct would be expected to generalize. We see two main paths toward increasing naturalism in experimental paradigms:

1. Changing existing experiment parameters to better cover the scope of natural distributions. For example, in an object classification task, we might expand the category set to include the broad range of objects people encounter in daily their life.

2. Adding ecologically motivated parameters that might interact with the theoretical constructs in question. For example, in a linguistic judgment task, we might test judgments across both spoken and written utterances.

We provides examples in Figure 3.

Our goal is to encourage cognitive scientists to embrace experimental paradigms and models

**2D object recognition**

**3D object recognition**

**Object search in AR/VR**

· simple 2D shapes    ✛ add spatial dimension    → objects can be stateful/have dynamics (e.g. open vs closed)
✛ can select where to look for object (action)
✛ object is in semantically meaningful location

**2D Minecraft environment**

**Legend**
→ parameter made more naturalistic
✛ New parameter added

**2D gridworld**

Small agent in a big world

Agent's view

· Small world
· Small number of objects
· Object acquisition task
· World is fully observable

→ Much larger world
→ Larger number of objects
→ Tasks are hierarchical (e.g. have subtasks)
→ World is *partially* observable

✛ There are other agents
✛ Some tasks require coordination and role assignment

**Judgments with context**

**Judgments across modalities**

**Grammaticality judgments**

"Who did he make the petition?"

"The opposition contested several previous filings. Who did the lawyer most recently make the petition?"

Grammatical / Ungrammatical    Grammatical / Ungrammatical    Grammatical / Ungrammatical

· Short sentences.
· Few other features.

✛ Additional context
→ Longer sentence
→ More structural variety

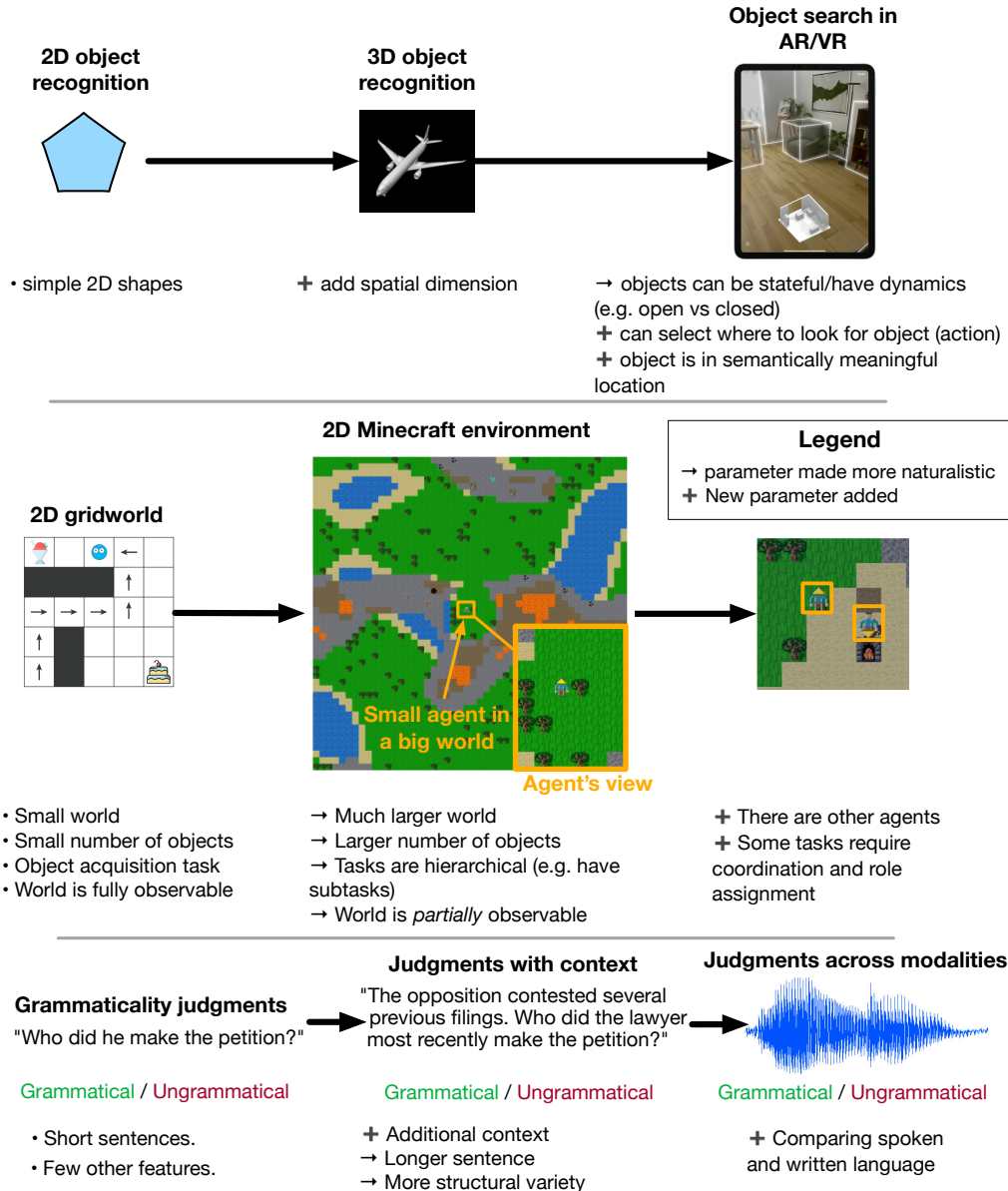✛ Comparing spoken and written language

Figure 3: **Examples of increasingly naturalistic settings that we can now study in a theory-driven manner.** All of these are settings where tasks and stimuli can now be parametrically generated—i.e., thanks to "generative AI", we can now *automate* the generation of photorealistic synthetic data and virtual worlds; thanks to virtual and augmented reality, we can now scan and parametrically manipulate real environments. Thanks to progress in AI, we can now build task-performing models that can operate in these stimuli and task settings.

that capture a broader spectrum of the variability in inputs, interactions, and tasks present in the natural environment of humans or animals—including, for humans, the rich cultural

constructs that form so much of our modern experience. That is, we identify a paradigm as more naturalistic if it incorporates a broader scope of the range of natural variability over which the theoretical constructs would be expected to generalize. Correspondingly, we identify a model as more naturalistic if it is capable of generalizing to make predictions across a broader range of contexts. Note that increasing naturalism does not always mean faithfully approximating all the joint statistics of the natural distribution—increasing the scope of variation along naturalistic dimensions can also enable moving outside the natural data distribution for controlled experiments (see §6.1).

Generally, identifying the parameters for increasing naturalism in a particular experiment is not well defined or obvious—it is task- and theory-specific. To illustrate "naturalistic" more explicitly, we will therefore give some examples of dimensions along which naturalism can be increased.

**Task paradigm**

- Incorporating multiple paradigms across which a hypothesis would be expected to hold, rather than testing on a single task; for example testing multi-step RL tasks as well as bandit tasks.

- Incorporating a broader set of stimuli, for example augmenting synthetic images with real images, or images with videos.

- Having task stimuli generated by many varying latent factors, not just the ostensible variables of theoretical interest.

**Environment**

- Large state spaces that reflect the complexity of real-world environments (Wise et al., 2023).

- Including the presence of other social agents that continue to learn with the agent (Hu et al., 1998) and that the organism must interact with (Vélez and Gweon, 2021).

- A continually changing environment (Abel et al., 2024)

**Model architecture**

- Architectures that can operate over sensory inputs like natural images (Yamins et al., 2014), speech (Kell et al., 2018), or natural language (Schrimpf et al., 2021) rather than simplified stimuli (e.g. low-dimensional or discrete inputs that contain only task-relevant features).

- Architectures with naturalistic action spaces—for example, modeling low-level motor control and embodiment rather than abstract, symbolic actions (Merel et al., 2019).

**Learning algorithm**

- Unsupervised (Higgins et al., 2021), self-supervised (Konkle and Alvarez, 2022), or intrinsic (Chentanez et al., 2004; Oudeyer et al., 2013) learning algorithms.

- Social (Henrich, 2016) or cultural (Cook et al., 2024) learning algorithms.

- Prospective learning algorithms that aim to model how the tasks evolve in continually changing environments (Seligman et al., 2013; De Silva et al., 2023).

Of course, it is easy to say that increasing the naturalism of our experiments and models would be useful, but how can we actually achieve it? A key goal for the later sections of our paper (§5-6) will be to give practical perspectives on how to achieve these goals, and how to develop theories that generalize across a broader range of naturalistic behavior. We briefly sketch these perspectives here.

**How do we scale up to naturalistic experiments?** We believe that recent developments in computer science and engineering make increasing naturalism in our experimental paradigms much more accessible to researchers, as we can use natural language to programmatically generate stimuli such as 3D scenes (Zhou et al., 2024) and videos of objects (Villegas et al., 2022), and leverage haptic feedback in virtual reality to simulate textures (Lu et al., 2022b). While the specific parameters will vary by domain, the core principle remains: gradually increasing ecological validity while maintaining experimental control.

**How do we glean understanding from potentially opaque learning based models? (§6)** We believe that task performing models offer powerful new tools for building cognitive theories. In a sense, these tools amount to treating the cognitive model as a proxy object for study—but one that admits much more opportunity for experimentation and understanding. For example, we can use "neuroscience"-like causal methods to probe a model's computations by intervening on them (e.g. Geiger et al., 2021). We can also interpret its behavior as an adaptation to the naturalistic properties of its training data (Chan et al., 2022; Prystawski et al., 2024, e.g.)—as a kind of *rational analysis* (Anderson, 1990; Lewis et al., 2014) of behavior and computations being a rational response to an environment. When applying these approaches to computational models, we have the ability to explicitly test our assumptions in ways we cannot with human subjects (e.g., how would the system behave if it were trained on only ungrammatical sentences). We can also attempt to distill our insights further, into analytic abstractions that illustrate the minimal instantiation of a phenomenon. By combining these routes, we can build rigorous links from an opaque task-performing model to conceptual understanding of cognition.

**In summary** naturalistic computational cognitive science seeks to (1) explain real-world intelligent behavior by developing models that operate over the scope of naturalistic inputs, produce naturalistic outputs, and are optimized according to the actual constraints and affordances of an person's environment, and (2) develop cognitive theories that unite these models with reductive understanding. In the coming sections, we will first motivate more deeply how naturalistic experimental settings can help us develop more complete accounts

of intelligence, before returning in more detail to the questions of how these goals can be achieved.
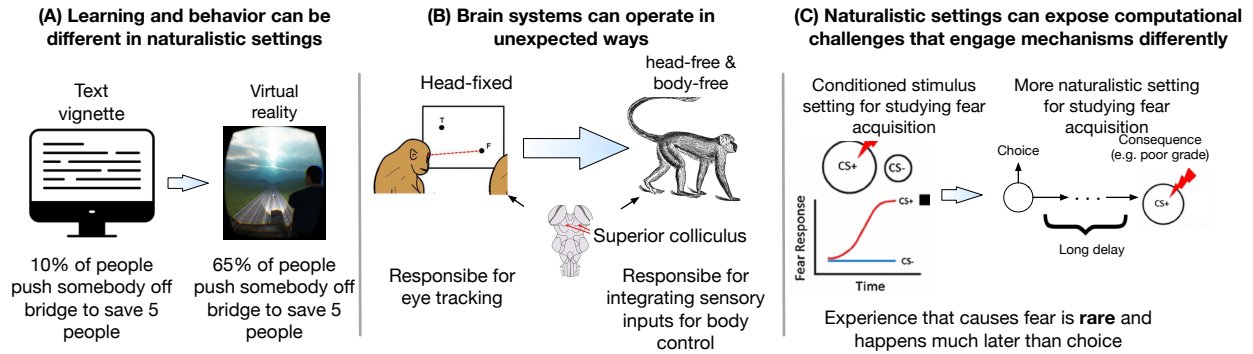
# 3  The benefits of naturalistic experimental settings



**(A) Learning and behavior can be different in naturalistic settings**

Text vignette → Virtual reality

10% of people push somebody off bridge to save 5 people

65% of people push somebody off bridge to save 5 people

**(B) Brain systems can operate in unexpected ways**

Head-fixed → head-free & body-free

Superior colliculus

Responsible for eye tracking

Responsible for integrating sensory inputs for body control

**(C) Naturalistic settings can expose computational challenges that engage mechanisms differently**

Conditioned stimulus setting for studying fear acquisition

More naturalistic setting for studying fear acquisition

CS+  CS-

Fear Response

CS+

CS-

Time

Choice

Consequence (e.g. poor grade)

Long delay

CS+

Experience that causes fear is **rare** and happens much later than choice

Figure 4: Overview of the benefits of increasingly naturalisic experimental conditions.

In this section, we outline the benefits of naturalistic experimental settings for cognitive science. Specifically, we argue that naturalistic experimental settings can engage mechanisms that are qualitatively different than those engaged in simpler settings — and thereby change the scientific inferences we draw. We review examples where naturalistic experimental settings have changed our understanding of a system, and have helped to disentangle competing models that simpler settings could not.

As an orienting conceptual example, imagine that we are doctors studying heart function, and we chose to only study it while the participants are lying down and resting. We could certainly learn interesting things thereby; the resting setting offers high test-retest reliability, and a highly-controlled environment to study the processes like how respiratory cycles affect heart rate. However, it offers this control and reliability precisely by removing the many factors of natural variation that interact with heart rate, and indeed for which heart function evolved — such as rapid adaptation to movement, stress, etc. Studying heart function across a broader range of naturalistic settings would be necessary to understand why the system is the way it is, and how all the physiological (and psychological) processes involved interact. For example, if we experimentally manipulated breathing, but did so across a range of different naturalistic activities (sitting, speaking, running, etc.) we would more effectively determine how breathing affects heart function, and how that effect depends on the state of the overall system. That is, by studying phenomena across a range of naturalistic settings we elaborate our understanding of the system, and can build more complete theories of its function.

## 3.1 Behavior can be different in naturalistic experiments

The starting point for many cognitive analyses is behaviors on a task designed to isolate some cognitive capability. However, there are many cases in which task-relevant behavior can be altered by seemingly-orthogonal features of the task—thus raising the risk that our theories will be overfit to a restricted task setting unless we also consider a broading range of variation. In this section, we review some examples illustrating changes in behavior or capabilities in naturalistic settings.

**Example 1: moral behavior can change when moving from textual vignettes to virtual reality**. Recent research has shown that behavioral choices can change when moving from text to virtual reality as it can mark a shift from moral *judgment* (i.e. deciding whether things are moral) to moral *action* (Francis et al., 2017). When presented with the classic trolley dilemma in text form, only 40% of participants say they would pull a switch to divert a trolley that would kill one person to save five others. However, when the same scenario is presented in virtual reality where participants must simulate the action, 55% choose to pull the switch. This difference becomes even more pronounced in the "footbridge" variant of the dilemma. In the text version, only 10% of participants say they would push someone off a bridge to save five others. Yet when able to simulate this action in virtual reality, 65% of participants choose to push the person. These findings demonstrate that people's stated moral preferences in hypothetical scenarios can differ substantially from their actions in more naturalistic settings.

**Example 2: memory performance differs between naturalistic search and explicit memorization**. When participants are explicitly instructed to memorize objects in a 3D home environment, their subsequent recall accuracy was significantly lower than when they incidentally encounter the same objects during visual search tasks (Helbing et al., 2020). This effect extends to spatial memory as well. Not only do participants better remember the identity of objects encountered during search, they also show more accurate memory for object locations—placing objects closer to their original positions when they were initially seen during search versus explicit memorization. These findings demonstrate that tasks focused explicitly on memory may not capture all aspects of how memory naturally operates during everyday behavior. Additionally, it suggests that computational models of memory may need to be revised to account for these differences in naturalistic settings. This is supported by research showing that people display detailed memory of object and spatial information in real-world scenes (Bainbridge et al., 2019).

**Example 3: mice can learn** $1000\times$ **faster during natural behavior, compared to two-alternative forced choice tasks**. Recent research has shown that learning rates can dramatically differ between simplified laboratory tasks and more naturalistic settings (Rosenberg et al., 2021). When mice are tested in traditional two-alternative forced choice (2AFC) tasks, they typically require 10,000 trials over 3-6 weeks of training to reach asymptotic performance of only about 67% accuracy. However, when allowed to *freely explore* a labyrinth

environment (similar to a natural burrow systems (Small, 1901)), mice demonstrate remarkably accelerated learning (Rosenberg et al., 2021). In this setting, mice learned to navigate a complex maze with 63 T-junctions in just a few hours, discovering and optimizing their path to a water reward after only about 10 reward experiences. The learning rate in the naturalistic setting was approximately 1000× times faster than in traditional 2AFC paradigms. This dramatic difference in learning speed demonstrates how traditional simplified experimental paradigms may fundamentally underestimate an organism's learning capabilities, and suggests that natural learning mechanisms may operate differently when they are engaged through more naturalistic behavior.

## 3.2    Neural systems can operate differently under naturalistic conditions

In addition to behavioral differences like those reviewed above, naturalistic conditions can result in changes to neural coding and computation. Thus, in order to arrive at a complete understanding of neural function, we need to consider the responses of the system across a range of simpler and more naturalistic settings. In this section, we highlight some examples of different computations or computational roles that arise when moving from standard paradigms to more naturalistic ones.

**Example 1: the model of superior colliculus shifted from a model of eye movement to one of integrating sensory inputs for body control**. Cisek (1999) suggests that the focus on computational behavior as an input-output mapping neglects the fundamental fact that natural intelligence evolved not to produce single responses, but for closed-loop control in an environment. This setting yields a rather different interpretation of the system's representations and processes. More generally, many researchers have argued that cognition cannot be understood completely outside its embodiment and the environment in which cognitive processes are instantiated (e.g., Newen et al., 2018). For example, Cisek and Green (2024) argues that as we increased naturalism from head-fixed to head-free to body-free settings when studying monkeys, we expanded our theory of superior colliculus from controlling saccadic eye movement to generally integrating multimodal cues to guide bodily action-selection (Cisek and Green, 2024). That is, as we increased the naturalism of the experimental conditions that we used to study monkeys, we arrived at a more complete model of superior colliculus.

**Example 2: visual processing systems operate differently during natural viewing compared to passive viewing paradigms**. Traditional research has shown that one of the earliest visual responses in the brain (known as P100/M100) occurs about 100ms after a person's eyes land and fixate on a new location (Vinje and Gallant, 2000). However, this finding comes primarily from experiments where participants passively view images while keeping their eyes still. In contrast, natural vision involves actively moving our eyes multiple times per second to sample information from different parts of a scene. When the researchers examined brain responses during this more natural active viewing, they found that the P100/M100 response actually begins when the eye movement (saccade) starts, not when it ends at the new fixation

point (Amme et al., 2024). Furthermore, they discovered that the neural patterns during active viewing were opposite (anti-correlated) to those seen during passive viewing, suggesting fundamentally different processing mechanisms. These findings reveal that simplified experimental paradigms, while valuable, may not capture all aspects of how visual processing actually operates during natural behavior.

**Example 3: model-based learning systems are used for learning moral judgments, but model-free learning are used for learning to avoid harming others**. Despite moral reasoning being strongly associated with model-based learning systems in the prefrontal cortex, neural evidence shows a surprising shift when people actually learn to avoid harming others (Lockwood et al., 2020). Moral reasoning commonly engages model-based systems in the lateral prefrontal cortex (LFPC) (Spitzer et al., 2007; Carlson and Crockett, 2018) with LPFC disruption leading to reduced norm compliance and enforcement (Knoch et al., 2006; Ruff et al., 2013). However, when participants engage in tasks requiring them to learn moral behavior (rather than just reason about it), researchers instead observed neural activity consistent with model-free learning. Similar to our moral judgment vs. moral action example from §3.1, this shows that when an experiment becomes more naturalistic (not just reasoning about morality but learning moral behavior), neural systems can operate in unexpected ways.

## 3.3 Naturalistic experimental paradigms can expose computational challenges that engage mechanisms differently

Finally, naturalistic paradigms can likewise introduce challenges that are important to understanding functioning of the system as a whole. In this section, we highlight examples where the challenges posed by increasing naturalism can engage computational mechanisms differently—which can be useful for disentangling underlying mechanisms, and is important for achieving generalizable understanding.

**Example 1: Different learning mechanisms only yield distinct predictions under working memory load**. Overly-simplified tasks can yield *aliasing* of different solutions, where many different computational approaches yield the same behavior. Collins (2024) presents an example, by showing that what seems superficially to be implemented as a standard reinforcement learning algorithm, on closer examination may instead be implemented through a combination of working memory and simpler value-free associative learning.[1] However, the difference between these two implementations cannot be identified in standard task settings, as Collins notes: "Even simple tasks designed to elicit a target process (such as bandit tasks for RL) recruit multiple other processes, but those processes may be unidentifiable in such tasks. Disentangling multiple processes requires considering more complex tasks to elicit differentiable behavior." The more complex tasks in question simply increase the stimulus set size within learning blocks to a handful of objects, rather than restricting to two or three—i.e.,

---

[1]Whether this alternative counts as an alternative implementation of RL or not is orthogonal to our argument; the important point is that the two models yield distinct predictions in some regimes but not others.

more broadly covering the range of set sizes more clearly disentangles the learning processes. Thus, by restricting to minimal paradigms, we enable a system to use many solutions and prevent ourselves from discriminating between them. By contrast, if we impose a broad range of evaluations on the system (cf. Nau et al., 2024)—and especially naturalistic evaluations that increase demands along different axes of task difficulty—the increased constraints actually make it easier to map the model on to natural intelligence, as is highlighted in the *contravariance principle* of Cao and Yamins (2024b).

**Example 2: fear conditioning**. Traditional fear conditioning research has provided fundamental insights into how humans learn associations between stimuli and aversive outcomes (Sehlmeyer et al., 2009; Maren, 2001; Britton et al., 2014). Through carefully controlled laboratory studies involving paired stimuli and outcomes, we have developed a rich mechanistic understanding of basic fear learning processes. These paradigms have proven especially valuable for studying certain types of naturally occurring fear learning, such as immediate aversive responses to foods (Riley and Tuck, 1985) or traumatic events like combat exposure (Engelhard et al., 2008). However, real-world fear learning often involves additional temporal dynamics that complement these well-studied processes. For example, a student may develop fear of exams after receiving a single poor grade weeks after taking the test, without requiring repeated negative experiences (Mobbs et al., 2021). This type of naturalistic learning involves delayed *delayed* and *sparse* feedback for choices, i.e. it relies on a computational mechanism for *credit assignment* (Sutton and Barto, 2018). Credit assignment describes how an agent credits choices (e.g. studying, socializing, diet, etc.) as responsible for future events (e.g. a poor grade). This is not to say humans have optimal credit assignment but that they have *some mechanism* for it and we are empowered to study it with appropriately naturalistic experimental settings. Clinical researchers have additionally pointed out that increased ecological validity can increase the applicability of our experimental results as the gap between lab-based fear conditioning and real-world fear has been a barrier to successful treatment of pathological fears (Beckers et al., 2013; Scheveneels et al., 2016; Krypotos et al., 2018).

**Example 3: seemingly-similar tasks that engage brain mechanisms differently can be reconciled by naturalistic models**. There has been a long debate about whether perceptual processes like "oddity" tasks (identifying an object that is different from the rest of the objects in a set) recruit the Medial Temporal Lobe (MTL). Specifically, some researchers (e.g Murray et al., 2007) have suggested that MTL is recruited to discriminate among "complex" stimulus sets that visual cortex cannot discriminate on its own (and thus MTL lesions impair performance specifically on those stimuli); other researchers have found no such interaction and suggested it stems from methodological issues (Buffalo et al., 1998; Suzuki, 2009). Recently, stimulus-computable models have shed new light on this debate. Specifically, Bonnen et al. (2021); Bonnen and Eldridge (2023) used pretrained convolutional vision models to directly process *all* the *raw* stimuli used in prior experiments, and evaluated whether the "oddity" stimulus could be distinguished from the rest of the set in the representation

space of these models. By doing so, Bonnen et al. showed that some of the purportedly "complex" stimuli from some of the prior experiments could nevertheless be distinguished with the vision model alone, while the complex stimuli from other experiments could not. Furthermore, the purportedly "complex" stimuli that in fact could be discriminated by the pretrained vision network alone did not require the MTL, while the complex stimuli that could not be distinguished by the vision model were precisely the ones that required MTL. Thus, the way that prior works operationalized "complexity" was underspecified; differences in operationalization that got abstracted away in conceptual descriptions of the experiments likely drove the seemingly-conflicting findings. This example illustrates the value of building models that directly perform the same naturalistic tasks as the subjects, and testing these models on benchmarks that incorporate stimuli from many experimental paradigms. This naturalistic approach can elucidate important differences among experimental paradigms, and clarify conceptual understanding.

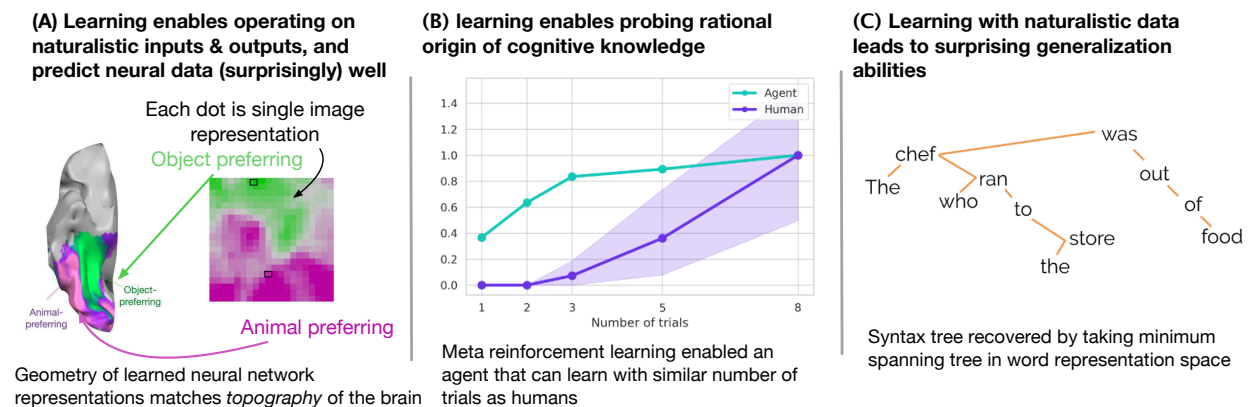# 4   The (surprising) benefits of learning with naturalistic data

**(A) Learning enables operating on naturalistic inputs & outputs, and predict neural data (surprisingly) well**

Each dot is single image representation
Object preferring

Animal preferring

Geometry of learned neural network representations matches *topography* of the brain

**(B) learning enables probing rational origin of cognitive knowledge**

Meta reinforcement learning enabled an agent that can learn with similar number of trials as humans

**(C) Learning with naturalistic data leads to surprising generalization abilities**

Syntax tree recovered by taking minimum spanning tree in word representation space

Figure 5: Overview of the benefits of learning with naturalistic data. (Figures reproduced from Doshi and Konkle, 2023; Manning et al., 2020; Team et al., 2023.)

The standard lesson in science is that we must simplify an experimental setting so that we can better arrive at a causal conclusion. However, we have already seen examples illustrating that more is different (Anderson, 1972)—certain properties of intelligence may only emerge in more "complicated" naturalistic settings. In this section, we detail examples where naturalistic data *itself* seems to play an importance role in producing the empirical phenomena of intelligence that we wish to study.

Many of the findings we will discuss here will be modern results from AI. One common trend is that there is a positive interaction between learning-based systems and naturalistic data—learning-based systems can accommodate naturalistic data, and reciprocally, learning from naturalistic data results in qualitatively better results than learning in simpler settings. While

we do not yet know the implications for natural intelligence, they open interesting questions and challenge prior assumptions. In particular, we argue that cognitive science should consider the limitations placed on our models when we fail to consider naturalistic data. (Note that while we draw many case studies from particular areas of deep learning, we believe that some of these factors could likewise contribute to generalization under other learning paradigms, such as program synthesis, or even evolution—and thus should be of interest even to those who do not work with deep learning models.)

## 4.1   Learning enables models that can operate on natural data

Perhaps the first surprising benefit of (deep) learning is that it enables a system to operate over natural data. Until AlexNet (Krizhevsky et al., 2012), the prevailing wisdom was that one had to design useful representations of natural data for models. However, AlexNet and subsequent work showed that a model could learn useful representations for naturalistic data simply by being trained to make predictions about this data. Surprisingly, these models learn representations that have notable correspondence to the internal representations in visual cortex (Yamins et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014). Here, we review some of this work and how learning-based methods can enable more flexible cognitive models that operate on natural data.

**Learning-based vision models can operate on natural data and predict neural data (surprisingly) well.** There is a rich history of learning-based methods revealing how neural networks can develop internal representations from pixel-based images that strongly correspond to neural data. A seminal example comes from Olshausen and Field (1996), who showed that V1-like Gabor filters could emerge from a convolutional neural network optimized to transmit input while minimizing overall activation (essentially, a sparse autoencoder). Prior to this, V1 activations were modeled using analytically defined Gabor filters. This breakthrough demonstrated that V1-like receptive fields could arise naturally from optimization to environmental statistics of edges at varying frequencies and orientations. Building on this foundation, Lee et al. (2007) demonstrated that V2-like convolutional filters could emerge from unsupervised optimization of sparse deep belief networks. A significant advance came when Yamins and DiCarlo (2016) showed that task-optimized deep convolutional neural networks trained for object classification developed intermediate layers that matched neural representations in V4 and IT better than any previous models. Recent research has revealed an even more striking finding: these task-optimized networks represent images in a latent space whose geometry has remarkable correspondence to brain *topography* (Doshi and Konkle, 2023). When fitting a 2D sheet to image representations such that neighboring points correspond to nearby points in the original high-dimensional space, the resulting sheet shows striking similarities to actual brain topography. Thus, learning-based approaches, particularly neural networks, not only enable effective visual processing but also provide surprisingly accurate predictions of brain organization and function.

**Integrating learning-based models with program synthesis enables flexible Bayesian models of human concept learning.** Recent work has shown how large language models can be combined with program synthesis to create more generalizable Bayesian models of human reasoning (Ellis, 2024). Hand-designed symbolic hypothesis spaces don't always transfer across domains. Ellis (2024) showed that one can circumvent these challenges by leveraging large language models to enable use of natural language as a flexible hypothesis space with a data-driven prior tuned to human judgments. By translating natural language hypotheses to executable programs for likelihood computation, such models can capture human learning dynamics in online concept acquisition tasks while continuing to predict human behavior robustly in new domains. This overcomes traditional limitations of pure symbolic program synthesis approaches which can fail to transfer due to incomplete symbolic vocabularies—a manifestation of the frame problem (Shanahan, 2004), where when deciding the relevant variables for a problem, we might accidentally choose a framing that doesn't generalize to other settings of interest. The integration of learning-based components with program synthesis thus enables Bayesian models that can flexibly adapt their hypothesis space while maintaining interpretable symbolic computation—a goal of many cognitive modellers.

## 4.2 Learning from naturalistic data can improve generalization

Studying how systems generalize is a key method in cognitive science and AI. Conceptual arguments often rely on empirical studies of generalization; however, these generalization properties are often studied in simple settings, with minimal models learning from minimal features (e.g. Marcus, 1998). A fundamental assumption underlying this approach is that the simplifications do not alter the generalization problem. However, here we review studies from AI and neuroscience showing that the variability of experience according to seemingly-orthogonal variables can alter how systems learn and generalize—thus implying that studying models within more restricted settings may be misleading about the naturalistic computational problem.

**Example 1: Compositional generalization of image classifiers vs. agents**. Naturalistic tasks can fundamentally change what models learn and how they generalize. Hill et al. (2020) trained two models to do a vision-language grounding task, and tested their compositional generalization to held-out language instructions. Both models were trained on the same language examples, and tested on the same held-out examples. However, one model was trained as an agent that interacts with a simulated environment, while the other was a simple image-language classifier trained on screenshots from that environment. The authors found that the interacting agent exhibited perfect compositional generalization to the novel language utterances, while the image classifier was substantially worse in novel settings. The authors also explored a range of other settings, including generalization benefits of more naturalistic 3D (rather than 2D) environments, or egocentric embodiments. Their results illustrate how richer, more naturalistic settings can enhance the generalizability of the solutions a system learns. Thus, if we are interested in understanding how a system generalizes, we may need

to build models that learn from appropriately naturalistic data.

**Example 2: Generalization to novel syntactic structures is enhanced by variability in other structures' fillers** Generalization is of deep interest in linguistics. Recently, Misra and Mahowald (2024) performed controlled experiments in which challenging linguistic constructions are systematically held out from the language model training data—and showed that models trained on naturalistic language can generalize to held out constructions by composing pieces of simpler constructions. Critically, the authors also found that generalization depended on the variability of the semantic fillers observed in the structures in training—thus showing how incorporating naturalistic variation in model training data can impact our theoretical inferences about generalization.

**Example 3: Rats raised in enriched environments**. Analogously, animals raised in more complex environments can be more skilled than those raised in simpler environments. This phenomenon has been well-documented in neuroscience, where rats raised in enriched environments — with social interaction and/or more space or toys in their home cage — show greater exploration, and improved learning and memory on novel tasks (Simpson and Kelly, 2011). This illustrates how, for natural as well as artificial intelligence, increasing the naturalistic variation in experience can alter the system's learning in ways that may impact our experimental and theoretical inferences.

## 4.3 Learning with naturalistic data can yield good performance across a range of seemingly disparate tasks

One interesting finding from the deep learning literature is that when deep learning architectures with many parameters are trained with a lot of naturalistic data and an appropriate "basic" learning objective, these architectures can develop mechanisms that go beyond the learning objective they were trained on. Crucially, these learning paradigms can allow models to transfer well (through initial performance or accelerated learning) on novel tasks beyond the training distribution.

Human learning may likewise benefit from transfer among the disparate tasks we learn. We argue that these AI findings motivate computational cognitive science research studying whether "down-stream human behavior" on a set of tasks can be recapitulated by a model which is trained on a large set of naturalistic tasks or stimuli representative of some of aspect of human experience—as in meta-learning approaches (cf. Wang, 2021). Below, we provide examples of this kind of transfer spanning computer vision, reinforcement learning, and natural language processing.

**Example 1: Computer Vision**. One of the earliest successes from deep learning came in computer vision. Soon after AlexNet (Krizhevsky et al., 2012) achieved strong results on the ImageNet dataset (Deng et al., 2009), researchers showed that the features discovered by AlexNet could be repurposed to novel tasks ranging from scene recognition to medical

diagnosis (Donahue et al., 2014; Sharif Razavian et al., 2014; Litjens et al., 2017). This was striking because these features were trained (a) via supervision on a fixed set of objects but (b) enabled transfer to novel objects or even different types of tasks. For years this pattern continued, and researchers even showed that the representations learned by these models could (by 2022, unsurprisingly) transfer from Imagenet to both 3D dexterious manipulation tasks and 3D household navigation tasks, sometimes achieving *better* performance than using "ground-truth" state information (Parisi et al., 2022; Yuan et al., 2022). In addition to generally useful representations, deep learning vision models trained with naturalistic data also develop mechanisms they were not trained for. For example, scene-oriented CNNs develop mechanisms for object detection (Zhou et al., 2014) and vision transformers develop mechanisms for segmenting objects (Caron et al., 2021)—both without an explicit training signal.

**Example 2: Reinforcement learning**. Likewise in reinforcement learning, researchers have found that reinforcement learning algorithms trained on many (e.g. billions) of tasks exhibit strong "out-of-distribution" generalization on unknown tasks (Team et al., 2021), with the ability to generalize to novel tasks in the same number of samples as humans (Team et al., 2023). Some algorithms can even generalize to collaborating with humans without any human data (Strouse et al., 2021). Recently, reinforcement learning algorithms for learning a modern variant of the successor representation (Gershman, 2018; Carvalho et al., 2024) have shown that they develop mechanism for exploration and behavioral skills without explicit training signals (Liu et al., 2024).

**Example 3: Large language models**. Perhaps the most striking example of this general phenomenon comes from large language models. Large language models are only trained to predict the next "token" (or word) to appear following a sequence of tokens. However, the distribution of internet language effectively includes a broad mixture of many tasks (Radford et al., 2019). When trained on vast amounts of this naturalistic data, large language models develop mechanisms for disparate tasks such as modular arithithmic, solving word problems, etc. (Wei et al., 2022)—and even for adapting to new tasks from examples presented in context (Brown et al., 2020). Moreover, the representations learned by these models can transfer to many superficially dissimilar downstream tasks (Lu et al., 2022a), even ones as dissimilar as playing video games (Reid et al., 2022). These examples illustrate how learning from a broad naturalistic distribution can induce many abilities—and can provide important transfer to downstream tasks.

## 4.4 Learning from naturalistic data allows us to ask new questions about the origins of knowledge

The points outlined above have an important consequences; using models that learn from naturalistic data can change our theoretical conclusions in cognitive science. For example, if a model fails to generalize when trained on simple data, we cannot be sure if the model

or the data are inadequate. By contrast, using models that learn from naturalistic data can enable us to ask more precise questions about which features—of the models or the data—are necessary to reproduce the theoretically-relevant features of cognitive and neural processes.

**Example 1: language models and the learnability of language.** Recent progress in language modeling from naturalistic data has suggested challenges to prior assumptions about language learnability and innateness. We briefly sketch these developments; Piantadosi (2023) and Futrell and Mahowald (2025) give much more thorough accounts. Classical approaches to language processing focused on simple models and ideas (e.g. Chomsky, 2014), but in doing so imposed strong assumptions. Under certain strong assumptions, Gold (1967) proved that even relatively simple languages cannot be learned from input alone. This proof contributed to arguments that there must be some innate universals such as recursion, and separation of syntax from semantics (Chomsky, 1957, 1965, e.g.)—ideas that were subsequently influential in arguments about the nature of cognition more broadly (Fodor et al., 1975; Fodor and Pylyshyn, 1988). These theories were based on studying particular language features in great detail[2]. Because of this focus, the theories did not attempt to model language processing in all its naturalistic detail.

However, recent progress in deep learning based language models (e.g. Brown et al., 2020) shows that many of these assumptions may need revision. By learning from large amounts of naturalistic data, these models acquire both syntax (Manning et al., 2020) and semantics (Li et al., 2021a)—and account for a broad range of behavioral and neuroscientific phenomena (Schrimpf et al., 2021). Of course, these models often learn from inhuman quantities of language (Wilcox et al., 2024). However, even language models trained on a human-like amount of language data can learn complex syntactic features that have been previously considered evidence for innate gramatical knowledge (e.g. Wilcox et al., 2023). Likewise, while some "impossible" languages that (adult) humans struggle to learn had been considered evidence for universal grammar, Kallini et al. (2024) show that language models can more easily learn real languages than impossible ones. Finally, controlled experiments—like those from Misra and Mahowald (2024), reviewed above—allow us to begin to understand which features of experience can contribute to generalizing to truly-novel structures, and how naturalistic variability contributes to this generalization. Thus, incorporating naturalistic data into has helped to reshape our understanding of the necessary and sufficient features for acquiring linguistic capabilities.

**Example 2: naturalistic data can lead to brain-like functional specialization.** Given the importance of perceiving faces in our social lives, it might seem likely that face recognition is innately encoded in the brain. Indeed, the discovery of the Fusiform Face Area (FFA)—a visual region specialized for face perception (Kanwisher and Yovel, 2006)—would seem to lend further support to this hypothesis. However, recent computational work shows that

---

[2]Or sometimes even studying the intuitions of linguists, with which the average language user would not necessarily even agree (Spencer, 1973).

this kind of specialization can emerge from a domain-general computer-vision model purely through training on a naturalistic distribution of object and face recognition tasks (Dobs et al., 2022). As the authors state: "It may be that the only inductive bias humans need to develop their face system is the already well-established early preference of infants to look at faces." Of course, this does *not* necessarily imply that the FFA is not innately specified. However, it does illustrate how training computational models over naturalistic data can help us to understand the natural constraints that affect neural organization—whether those effects occur over developmental or evolutionary timescales.

# 5    Building generalizable models

ML research has excelled in developing generalizable models. While we focus on their utility for working on naturalistic stimuli and tasks, we believe cognitive science would broadly benefit from improved practices for developing generalizable models. In this section, we highlight three key practices of empirical ML research that we believe computational cognitive science can adopt. The first is frictionless reproducibility: the practice of developing research artifacts which can be reused, and repurposed with minimal effort (§5.1). The second is a focus on developing generalizable models that can learn successfully across many datasets (§5.2). We argue that these two offer the foundation for having groups of research collectively work on important research problems through the utilization of benchmarks (§5.3).

Before we continue, we acknowledge that modern empirical ML is a young field which, like other fields, has faced challenges in rigor and reproducibility (Melis et al., 2017; Lucic et al., 2018; Riquelme et al., 2018; Henderson et al., 2018; Recht et al., 2019; Agarwal et al., 2021). Despite these challenges, we are optimistic that some of the practical knowledge ML has developed is valuable for improving rigor and reproducibility, as we highlight below.



(A) leverage frictionless reproducibility practices to easily reproduce and improve each other's models

(B) develop models with an eye towards to increase chance of external validity

(C) Embrace *dynamic* **phenomena-oriented benchmarks** to catalyze community progress
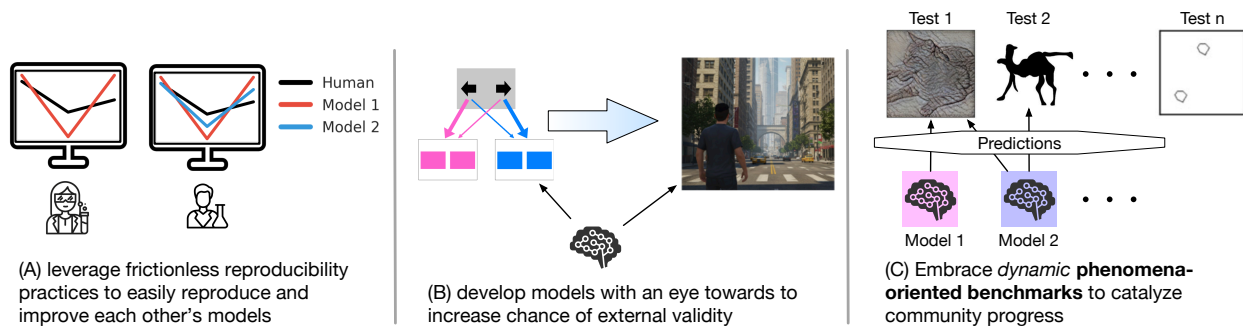
Figure 6: Overview of key strategies cognitive science can adopt to develop generalizable models. (Figures reproduced from (Geirhos et al., 2018; Bowers et al., 2023))

**Section summary.** We envision a future where cognitive scientists identify a cognitive phenomenon of interest (e.g., theory of mind, working memory, or object recognition) and create

a phenomena-oriented benchmark **(R8)** containing experiments specifically designed to test different theoretical aspects of that phenomenon. Here, experiments are implemented with standardized evaluation protocols and interfaces **(R1)** that enable direct comparison between competing theoretical frameworks. All experimental data and stimuli are made available through one-line download access **(R2)**. Models are developed and released following frictionless reproducibility standards **(R1-R3)** for easy adoption and modification. Researchers identify gaps in theoretical coverage and contribute new experiments to the benchmark, facilitating maintaining a separation between model and task design **(R6)**. Thanks to the existence of a benchmark and frictionless reproducibility standards, models are evaluated comprehensively across all experimental paradigms with minimal additional human labor. This helps mitigate the challenge of disconnected, incommensurate findings that resist integration across studies (Almaatouq et al., 2024). Thanks to the leaderboard principle **(R7)**, this also helps ensure that progress is defined by substantial rather than incremental advances. This approach encourages the development of models that work across both simplified and naturalistic contexts **(R4)** and that generalize across multiple experimental paradigms **(R5)**. Importantly, in addition to a model that predicts behavior across many tasks, we also have a model that can perform many tasks that humans can—**enabling the building of human-like models that scale and generalize**.

## 5.1 Frictionless reproducibility: a superpower

Within cognitive science and psychology, there is currently substantial friction when using data or code of prior work (Nosek et al., 2015; Hardwicke et al., 2018). Researchers may not share data in a standard or easy to use format, often don't fully report their analysis methods, or provide code that contains errors (Poldrack et al., 2017). In a large-scale replication analysis of open code and data, researchers found that 38% of code was not usable and *only* 31% of workflows had reproducible results (Hardwicke et al., 2018). These factors may have contributed to the replication crisis (Lilienfeld, 2017; Shrout and Rodgers, 2018).

In recent years, data-centric sciences are experiencing a profound transition to what David Donoho calls "*frictionless reproducibility*" (Donoho, 2024; cf. Recht, 2024). Donoho identifies three key pillars to frictionless reproducibility: data sharing, code sharing, and competitive challenges. We will focus on the first two in this section, and competitive challenges in §5.3. Critically, with frictionless reproducibility, the research artifacts that are produced allow future researchers to exactly re-execute the same *complete* workflow with minimal effort. With this, good ideas are spread, adopted, and improved with little friction. Donoho argues that ML is ostensibly the *most successful* at frictionless reproducibility. Below we review some supporting strategies that may aid cognitive science.

## Recommendations for cognitive science

**R1. Use a standardized evaluation protocol where models can be compared via a common interface**. Machine learning's foundation of standardized evaluation protocols traces back to 1959, when Highleyman created the first alphanumeric dataset and the first ever "train-test" split for comparing pattern recognition (Highleyman and Kamentsky, 1959). This strategy continued with the famous MNIST alphanumeric dataset (LeCun et al., 1998). MNIST simplified data-sharing and evaluation by pre-processing data, fixing the train and test distributions, and sharing the data online. This standardization and usage of a common interface enabled fair comparisons between diverse approaches, from boosted stumps (Kégl and Busa-Fekete, 2009) to support vector machines (Cristianini and Scholkopf, 2002) to nearest neighbor methods (Khan, 2017).

Standardized evaluation protocols serve a critical function: **mitigating personal bias** in model assessment. The evolution from MNIST through Caltech 101 (Fei-Fei et al., 2004) to ImageNet (Deng et al., 2009) demonstrates this principle, as increasingly naturalistic datasets enabled fair comparisons between diverse model families. This enabled AlexNet (Krizhevsky et al., 2012) to validate deep learning in an era when many researchers doubted its potential.

Thus, we recommend that cognitive scientists establish clear *hold-out* data for either model development or data analysis. This is especially critical for confirming exploratory results (Poldrack et al., 2017). Combined with standardized evaluation protocols, this will enable fair comparison between diverse theoretical frameworks ranging from program synthesis to deep neural networks to reinforcement learning and active inference.

**R2. Make data accessible with a single line of code**. *Fully processed* research data should be *programmatically* accessible with a single line of code, requiring no permissions. While uploading data to platforms like the Open Science Framework (OSF) is valuable, we additionally recommend uploading data to platforms like Hugging Face that offer version control, free hosting, authentication-free access, the ability to download data with one line of code, and to upload data with a few lines of code. This removes friction from reusing and updating datasets.

**R3. Adopt the "single file" philosophy**. Ideally, key workflows (e.g. model or evaluation definititions) should be defined in a single file. This has emerged as a powerful force for reproducibility in ML research (see Table 1 for examples). While this approach contradicts traditional software engineering principles of modularity and abstraction, it offers distinct advantages for research code: (1) complete understanding of an important component (e.g. how a model works) requires reading only one file, (2) minimal abstraction layers and dependencies, and (3) easy adaptation of either complete implementations or specific components. One can simply *copy* this single file into another codebase. Removing friction from future reuse also benefits us—as we can see from Table 1, these methods get a large number of citations.

| Model | Citations | Age | GitHub Stars | Single file examples |
|---|---|---|---|---|
| ResNet<br>He et al. (2016) | 247,921 | 8 years | 2,300 | model |
| Contrastive Learning<br>He et al. (2020) | 13,480 | 4 years | 4,800 | model, training |
| Object Discovery<br>Locatello et al. (2020) | 790 | 4 years | 394 | model |
| RL Library<br>Huang et al. (2022) | 252 | 2 years | 5,600 | model 1, model 2 |
| RLHF Method<br>Rafailov et al. (2024) | 1,730 | 1 year | 2,200 | train |
| Multi-agent RL<br>Rutherford et al. (2023) | 24 | 0 years | 433 | model |

Table 1: Comparison of ML libraries that adopted the single file philosophy.

## 5.2 Developing models with an eye towards generalizability

Task-performing computational models have been a goal of cognitive science since Newell (1973). This sentiment still exists today (Kriegeskorte and Douglas, 2018; Almaatouq et al., 2024) and is exemplified in conferences that explicitly seek to integrate cognitive science and neuroscience with AI (Naselaris et al., 2018). Cognitive scientists are well aware that they can't rely on a single paradigm for the studying an effect if they want to achieve generalizable conclusions (Holzmeister et al., 2024; Yarkoni, 2022). Ideally, the models we develop generalize beyond the setting (e.g. stimulus set or task-specification) they were designed for. This requires two things. First, that the model can *operate* over new stimulus sets and generalize its behavior to new tasks. Second, that the behavioral (and potentially neural) predictions made by our models capture human behavior on new tasks and new stimulus sets. Clearly, the first is a requirement for the second. Here, we argue that machine learning, with its emphasis on frictionless reproducibility, has made significant strides in developing models that can generalize. We detail recommendations for cognitive science below.

**Recommendations for cognitive science**

R4. **Develop models with simplified stimluli but ensure they work with increasingly naturalistic stimuli**. The frame problem (Shanahan, 2004) and the "simulation is doomed to succeed" critique (Grim et al., 2013) highlight a fundamental challenge in model development: the evaluation tasks we design often fail to capture the full complexity of the phenomenon we care about. Empirical machine learning addresses this by developing models that must succeed in both simplified and increasingly complex settings. Below we detail two examples.

*Example 1: Generative Models.* One of the most successful models of machine learning

has been the variational autoencoder (Kingma, 2013), a generative model that learns to infer latent variables which describe an observation. When this model was released they showed results inferring latent variables correctly on both "simplified" MNIST digits and more complex images of real-world faces. Importantly, in this more complicated setting where the latent dimensions did not have obvious correspondence, they performed a qualitative analysis showing their model had learned a face manifold capturing semantically meaningful dimensions of faces such as sentiment and facial direction.

*Example 2: Reinforcement Learning.* The history of RL is one of progress towards increasing naturalism: initial environments were grid-worlds (Sutton and Barto, 2018), later 2D video games like atari (Bellemare et al., 2013), then virtual 3D home environments (Szot et al., 2021; Li et al., 2021b), and now an emphasis on open-world environments (Fan et al., 2022; Matthews et al., 2024; Hughes et al., 2024). Despite the increased emphasis on naturalistic settings, there is still emphasis on the importance of toy settings for studying specific aspects of your model (Osband et al., 2019; Obando Ceron et al., 2024). Indeed, the strategy is almost always to *first* develop your model for simplified settings, but one should always then see that your model continues to "work" with more naturalistic settings. Collectively, this has allowed RL researchers to both advance a theoretical understanding of how these models work (Dadashi et al., 2019; Grimm et al., 2020; Lyle et al., 2022), an empirical understanding of how they work (Obando Ceron et al., 2024; Farebrother et al., 2024), but also develop performant models for a wide range of "complex" settings including beating the world's best players in games like Go (Silver et al., 2017) and controlling robots in the real world (Levine et al., 2016; Cheng et al., 2024).

**R5. Develop models that work with many stimulus sets**. The ultimate goal of ML has always been generalization (Highleyman and Kamentsky, 1959; Hardt and Recht, 2022; Recht, 2024). Initially, this meant models that generalize to new data. Now, it means both models and learning procedures that generalize to new settings. Thus, the practice of evaluating models across diverse datasets has become foundational in machine learning. This has manifested with research papers where ML models are evaluated on numerous datasets to showcase their generality. Early examples include "dropout" (Srivastava et al., 2014), which demonstrated generalization improvements across 6 distinct image and speech datasets, and MoCo (He et al., 2020), which validated its contrastive learning approach on 9 computer vision datasets. Another notable example is the seminal "Deep Q-Network" paper (Mnih et al., 2015), which showed that a single neural network architecture with fixed hyperparameters[3] could master 50 different Atari games. The field has since evolved to group environments by the cognitive capabilities they test (e.g., exploration, generalization, manipulation) (Patterson et al., 2023), with researchers developing models on one set of tasks and validating on entirely different ones to ensure robust generalization (Machado et al., 2018).

---

[3]Hyperparameters are parameters that are set for a model and algorithm before training, such as the number of layers in a neural network or the learning rate of the optimizer

We recommend cognitive scientists also test the models that support their theories more broadly. Importantly, we should develop our models with some stimulus sets or tasks, and then evaluate them (e.g. learn the parameters) on a different set of stimulus sets or tasks (Machado et al., 2018). This further helps prevent our computational model (and thereby theory) from "overfitting" to a particular paradigm or set of stimuli.

**R6. Iterate between between (1) model-design (2) task design where the other is fixed.** One crucial lesson from ML research is the importance of separating model development from task design. Patterson et al. (2023) strongly caution against simultaneously developing both the problem setting (datasets or tasks) and the solution method, as researchers may inadvertently design evaluation environments that favor their proposed solution. Historical examples illustrate this risk. When developing Q-learning algorithms, researchers modified the classic pendulum control problem by adding episode cutoffs and random state resets (Machado et al., 2018). While these modifications made Q-learning more tractable by implicitly improving exploration, they obscured fundamental limitations in Q-learning's exploration capabilities compared to policy gradient methods. Maintaining a strict separation between task design and model development enforces a clear delineation between hypothesis formation and testing, which is **important for reducing experimenter bias**. This is especially important because researchers tend to find "positive" results (Scheel et al., 2021; Sarafoglou et al., 2022), a bias exacerbated by an increasingly competitive academic landscape (Lee, 2014; Reithmeier et al., 2019; Woolston, 2021). Much of ML's progress can be attributed to researchers evaluating their models on independently developed datasets and tasks—a practice that not only simplifies the research process but also promotes scientific rigor.

Thus, we recommend that cognitive science also develop models on tasks/stimuli that have been developed by others. This may seem to contradict science, where researchers develop a hypothesis and create an experiment to test their hypothesis. The data from prior work was for a different hypothesis! We recommend that the field begin to explore phenomena-oriented benchmarks that aggregate stimuli and tasks designed to test specific theoretical predictions for broad phenomena. This approach would allow researchers to leverage existing experimental paradigms while maintaining scientific rigor through independent validation. We detail this more in the next section.

## 5.3 Benchmarks: catalysts for progress and innovation

Another challenge that cognitive science faces is fragmentation. This has been true since Newell's seminal critique (Newell, 1973), where his observation that piece-wise investigation of brain components would not yield a cohesive understanding remains relevant forty years later (Schrimpf et al., 2020; DiCarlo et al., 2023). Current experimental paradigms often produce disconnected, incommensurate findings that resist integration across studies, even within the same domain (Almaatouq et al., 2024). This fragmentation manifests not only in methodology but also in how researchers distribute their attention across different problems,

potentially diluting collective progress.

The history of machine learning offers an instructive parallel. In its early days, researchers worked on disparate, self-defined problems, leading to diffuse progress (Recht, 2024). The introduction of shared benchmarks marked a crucial shift: researchers began focusing on standardized problems that others had identified as important. We argue that this led sub-populations of researchers to converge on small subsets of problems. Putting aside whether this is good or bad, it's been effective. Indeed, this is the third pillar of frictionless reproducibility which (Donoho, 2024) argues has been a major catalyst for progress in ML: competitive challenges. We argue that cognitive science can similarly stand to benefit from benchmarks.

**Recommendations for cognitive science**

**R7. Embrace the leaderboard principle**. A common concern with benchmarks is that researchers will "overfit" to test sets, which theoretically should underestimate true test error (Hastie et al., 2009). Yet despite widespread iteration on benchmarks in machine learning—a practice that ostensibly violates statistical principles—the field has produced models with strong performance on held-out data. This apparent paradox is partially explained in a recent machine learning textbook (Hardt and Recht, 2022) by a phenomenon known as the "leaderboard principle". Researchers typically only publish and build upon models that demonstrate substantial improvements over prior state-of-the-art results, rather than minor variations. This selective pressure towards meaningful advances effectively constrains the degree of adaptation to test sets, as researchers focus on substantial improvements rather than exhaustively exploring the test set's peculiarities. Still, a risk of overfitting to data remains. We argue that the remedy to this for cognitive science is the creation of phenomena-oriented benchmarks, where researchers can contribute new challenges that identify a missing component of a phenomenon of interest. We detail this below.

**R8. Develop dynamic benchmarks to iteratively capture important phenomena**. ML benchmarks are not collected to test specific hypotheses (Hardt and Recht, 2022) and so may seem unscientific. However, we argue that cognitive science can embrace "scientific" benchmarks oriented toward testing specific hypotheses. We use the story of "Brain-Score" as an illustrative example of *initial* steps in this direction.

Brain-Score (Schrimpf et al., 2020) began as a benchmark with a collection of behavioral and neural data from humans and monkeys on object recognition tasks. One common critique of Brain-Score is that the behavioral analysis essentially boils down to whether models and humans found the same stimuli easy and challenging to identify. However, it did not test *hypotheses* about how humans identify objects, nor show that models identified objects in this same way. This was perhaps best described by Bowers et al. (2023), who noted that matching predictions of humans was not enough to show that this was a good model of human object recognition. In a scathing critique, they reference 9 experiments from the human psychophysics literature which presented stimuli generated to test specific hypothesis and

showed that deep neural networks failed to capture human behavior across these. Examples include stimuli sets that tested the hypothesis that object recognition was based on local shape features (Baker et al., 2020), stimuli sets that tested the hypothesis that human object recognition is defined by shape more than by texture (Geirhos et al., 2018), and stimuli sets that tested the hypothesis that pairs of objects were judged to be similar based on shape contours (Puebla and Bowers, 2022). In response, DiCarlo et al. (2023) added many of these stimulus sets to Brain-Score so that when future models are tested, they are tested on whether they capture human behavior across all of these stimulus sets.

We believe this story illustrates what cumulative, dynamic naturalistic computational cognitive science research could look like: where researchers can develop models and experiments that easily integrate with prior findings on a phenomenon thanks to phenomena-oriented benchmarks, frictionless reproducibility, and the development of generalization-oriented models. Of course, we must still continue to analyze and study our models to build comprehensive theories, as we discuss in the next section.

# 6    Building from naturalistic experiments to cognitive theories



(A) Combining naturalistic variation with parametric manipulation for experimental control together with generalizability.

(B) Tracing complex model behaviors to mechanisms and data properties (rational analysis).

(C) Theories that combine predictive models with reductive explanations.

Figure 7: Overview of how we can develop theories with potentially opaque models and unnatural manipulations of natural data. (Panel C bottom figure is reproduced from Saxe et al. 2019.)

A primary goal of cognitive science is understanding. Experiments and computational modeling are key constraints on theory building. By instantiating our theories in a model, we are forced to make our theories more precise by concretizing the ambiguous details of the mapping from a verbal hypothesis to its implementation (Guest and Martin, 2021).

From this perspective, building computational models that perform the task in as naturalistic a way as possible, across as wide a variety of settings as possible, imposes much stronger constraints on our theories than an abstracted model at a higher-level (cf. Cao and Yamins,

2024b). These added constraints are important; otherwise, high-level theories that do not directly engage with the details of the implementation can be so unconstrained as to lack explanatory value (cf. Erev and Greiner, 2015; Jones and Love, 2011; Rahnev and Denison, 2018; Andrews, 2021), can be intractable for real problems (Van Rooij, 2008), or can elide important details, as we have outlined above.

Yet moving towards naturalistic settings, and the complex models they require, alters the way that we need to approach experimentation and theory building. In this section we discuss the consequences for experiment design (§6.1-6.2), the role and interpretation of models (6.3), and theory building (§6.4).

## 6.1 Performing controlled experiments by parametrically manipulating naturalistic data (in unnatural ways)

Developing and testing models on truly naturalistic data introduces new issues that are not present in simple artificial settings. Naturalistic data distributions may be too "easy"—for example, they may only rarely include edge-cases that test key capabilities (Zhang et al., 2023); thus, testing average performance over a naturalistic distribution may not identify important failures. Relatedly, in naturalistic data, features may be confounded, which can prevent readily disentangling how a model is solving a task (Rust and Movshon, 2005). For example, models trained on ImageNet (Deng et al., 2009) rely on features like texture more than humans do (Geirhos et al., 2018), but still perform similarly to humans on ImageNet. These issues limit our ability to achieve deep understanding from experimenting with only the distribution of naturalistic data (Rust and Movshon, 2005).

However, we believe that these challenges are surmountable by parametrically manipulating rich, naturalistic data (Fig. 8)—preserving or enhancing its richness, while more carefully accounting for theoretically motivated constructs. For example, to effectively test processing of challenging syntactic structures, Futrell et al. (2021) edited natural stories to increase the density of these rare constructions, while preserving the richness and naturalism of the original content. Similarly, the work of Geirhos et al. (2018) engineered datasets that combine natural shapes and textures in unnatural ways. To do so, the authors exploited the contemporary ML technique of iterative style transfer Gatys et al. (2016)—combining the features of different images at different spatial scales—thus illustrating how progress in AI models unlocks new experimental approaches. These conflicting shape-texture datasets helped to inspire similar experiments on humans (Jagadeesh and Gardner, 2022b,a); surprisingly, those studies found that human visual cortex likewise appears to use textural rather than shape-focused representations, and shape-driven behavior is thus mediated by downstream readout processes. Thus, manipulating features to test models can lead to new insights into brain function.

A growing range of studies have taken similar approaches to systematically manipulate or ablate features of a natural stimulus, while preserving as much of the natural variety as
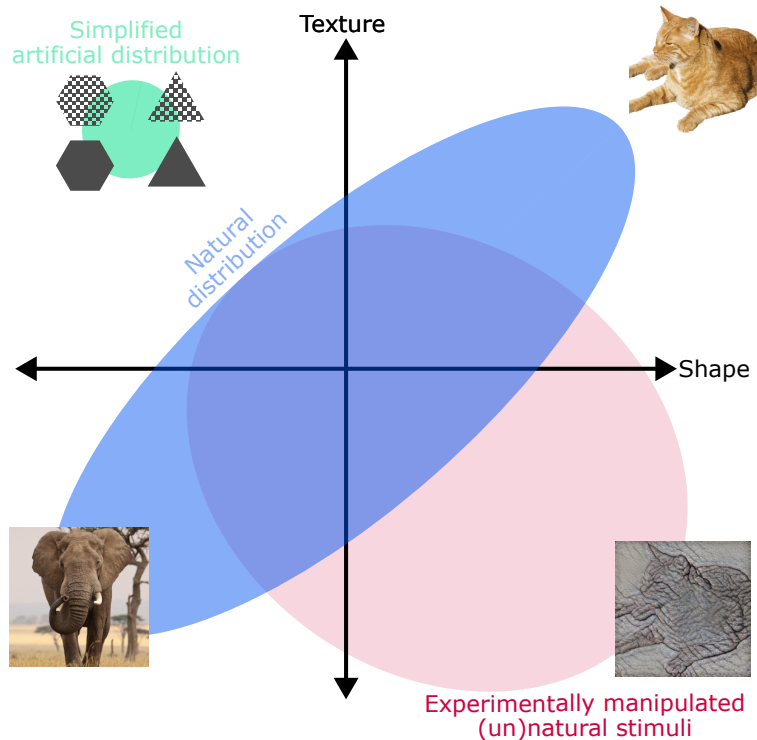
Figure 8: Testing hypotheses by parametrically manipulating naturalistic data. In the natural data distribution (blue), features like shape and texture are highly correlated—cats usually have cat fur, and elephants usually have elephant skin. This confounding can make it difficult to discern how these variables independently contribute to processing in neural systems or models. One approach to deconfounding these variables is to use a narrow distribution of highly-simplified artificial stimuli (light green) that isolate the variables of interest. However, we advocate for (additionally) creating broader distributions of stimuli that preserve more of the breadth of the natural data distribution, while deconfounding or manipulating features parametrically in order to test hypotheses—as in the unnatural combinations of natural shapes and textures created by Geirhos et al. (2018).

possible—whether introducing new parameters like modality manipulation of a stimulus (e.g. Deniz et al., 2019) or systematically selecting stimuli from a broad distribution to drive particular responses while preserving variety (e.g. Tuckute et al., 2024; Hosseini et al., 2023). We likewise advocate for testing on tasks that augment natural data distributions with unnatural systematic manipulations in accordance with the hypothesis to be tested (cf. Jain et al., 2024). Experimenting with rich settings does not mean giving up theory-driven experimental manipulation; it simply means situating those controlled experiments in more complex, naturalistic settings. For example, there has been a growing research direction using games as richer settings for the experimental study of cognition (Allen et al., 2024). Moreover, naturalistic experiments *enhance* our ability to develop generalizable theories, as we argue below.

## 6.2  Incorporating naturalistic variation to enhance generalizability

By testing hypotheses in settings where other aspects of the task and data generating process are sampled broadly from the naturalistic distribution, we address some of the challenges of generalizability that arise from testing on a much narrower distribution than the written hypothesis implies (cf. Yarkoni, 2022). For example, if the hypothesis is about visual processing in general, but the task paradigm relies solely on abstract shape stimuli, the conclusions may not generalize in the way the written results would suggest. Thus, testing our hypotheses across varied settings—"design[ing] with variability in mind" (Yarkoni, 2022)—increases the likelihood of identifying robust effects that will generalize (cf. Baribault et al., 2018). In the language of regression, varying as many other data generating factors as possible brings us closer to estimating the "main effect" of an experimentally-manipulated construct, rather than estimating the "simple effect" of the manipulation in the single setting we have tested.

For example, it is a long-standing principle of experimental design to test a hypothesis with multiple stimuli, and statistically quantify the stimulus effect to estimate the generalizability (e.g. Clark, 1973). However, many other experiment features—even exact task formulations such as multi-arm bandit tasks—are often preserved within a paper, or even across an entire literature. Just like testing a hypothesis with a single stimulus, this lack of variation poses challenges for generalizability (cf. Eckstein et al., 2022)—especially because the way that researchers operationalize a theory can substantially affect their conclusions (e.g. Holzmeister et al., 2024; Schweinsberg et al., 2021; Strickland and Suben, 2012). While naturalistic variation cannot resolve this issue in full, we believe that incorporating a broader range of the natural variation of settings, stimuli, and tasks within which we expect those theories to hold, will help us to develop more generalizable theories.

Incorporating variation into our experiments can also help us to revise the constructs underlying our theories. For example, Eisenberg et al. (2019) study the putative construct of self-regulation by simultaneously testing a large variety of existing self-regulation measures. The authors find that these measures do not form a unitary construct. Rather, task- and survey-based measures each tap into distinct, multi-factorial constructs, with individual tasks loading primarily onto a subset of the factors in their domain. That is, two different studies that each use only a single measure for self-regulation may be tapping into entirely different constructs. By contrast, by substantially varying the tasks we use to measure a theoretical component like self-regulation, we can more fully disentangle the underlying constructs—and thereby improve the extent to which our theories generalize.

## 6.3  Interpreting complex models to yield explanations

If we perform experiments in naturalistic settings, and develop complex models that can predict human behavior across them, how does this support explaining the cognitive phenomena of interest? In this section we will review some paths to deriving explanations of phenomena

from complex models.[4]

**Data properties as explanation & rational analysis.** One route to explaining model behavior is as a consequence of the properties of the data from which it learns. This idea stretches back through behaviorism as well as cognitive science; for example, connectionist research has often focused on how naturalistic data could drive cognitive phenomena (McClelland and Jenkins, 1991; Elman, 1996; Rogers and McClelland, 2004). Analogously, in AI, various researchers have taken inspiration from behaviors that emerge in large pretrained models, and similarly attempted to identify minimal properties of naturalistic data that give rise to those behaviors. For example, Chan et al. (2022) explore how the bursty, long-tailed nature of natural data can give rise to in-context learning, and Prystawski et al. (2024) examine how local dependencies can yield certain kinds of sequential reasoning. These works illustrate how rich properties of natural data can be distilled down to the minimal elements that give rise to a behavior—and can thus likewise offer a candidate explanation for the origin of those capabilities. Moreover, explanations of behavior in terms of the property of the data from which a system learns offer a route toward *normative* explanations of that behavior as a rational solution to constraints—precisely the perspective taken by *rational analysis* (Anderson, 1991).

**Mechanistic explanation.** Understanding that a behavior arises from properties of data does not explain how that behavior is implemented. Fortunately, the internal workings of task-performing computational models can generally be inspected and intervened upon. For example, in contemporary AI research, a variety of works have analyzed the mechanisms that implement model behaviors, either in simplified settings (e.g. Nanda et al., 2023; Zhong et al., 2024), or even in large models trained on natural data (Geiger et al., 2021; Wu et al., 2024). A particularly promising approach for cognitive science is illustrated by Distributed Alignment Search (Geiger et al., 2024), which attempts to map postulated abstract causal models onto low-level features of the model that play corresponding causal roles—thus enabling the link between abstract models or hypotheses about the computational solution, and evidence about concrete implementational mechanisms. Morever, mechanistic interventions can even be combined with data studies, to examine the causal role of different mechanisms over the course of learning (Singh et al., 2024). Appropriate regularization of models can also produce representations that yield greater interpretability—e.g. Miller et al. (2024) train networks to predict animal behavior while restricting their representations, and used these to identify non-optimal features influencing the decisions—thus allowing bottom-up discovery of possible implementations. Thus, various methods for mechanistic study of models can offer powerful routes to experimentally determining sufficient implementations of cognitive processes.

**Formal theories.** From the way that behavior emerges from data or mechanisms, we can

---

[4]N.B., not all complex systems admit simple, abstract explanations. For example, where natural intelligence incorporates chaotic (Freeman, 1995) or critical (O'Byrne and Jerbi, 2022) dynamics, they may alter the kind of explanations we can seek (cf. Kellert, 1993). A benefit of seeking predictive models as a route to explanation is that the models may be useful even if simple explanations do not follow.

progress to formal theories. Saxe et al. (2019) offers an illustrative example. Prior works observed empirically that pseudo-naturalistic data give rise to various aspects of human-like semantic development and representation in neural networks (Rogers and McClelland, 2004), such as progressive differentiation. However, these works had not explained *why* these parallels emerge. Saxe et al. (2019) bridged this gap by formally deriving how the data properties combine with gradient-based learning to yield these cognitive phenomena. This illustrates how preliminary explanations in terms of data can be developed into more rigorous analytic theories. Thus, building naturalistic models can be a stepping stone toward more rigorous, normative theories of intelligence as a rational solution to a particular problem setting—much like prior work in cognitive science (e.g. Anderson, 1991; Oaksford and Chater, 2009).

## 6.4 Theories that combine task-performing models with reductive explanation.

Above, we have argued for the value of embracing the richness of naturalistic tasks, even if it leads us to build more complex task-performing models that are hard to interpret, but that perhaps make more accurate predictions of behavior. In the previous section, we have outlined some of the possibilities (and challenges) of deriving understanding from these models. So what, ultimately, should we seek?

We argue that naturalistic computational cognitive science should seek theories of cognitive phenomena that consist of two components:

1. Task-performing (predictive) models that reproduce the phenomena across the same range of naturalistic stimuli and paradigms as the human/animal subjects.

2. Reductions of these task-performing models to simpler mechanisms, properties, and theories of *why* these models reproduce the phenomena.

These two components serve different purposes. The first helps to demonstrate that the models are real candidate models of the phenomena in question, rather than models that simply could not scale to the real problem or that are overfit to a particular instantiation of the task. We believe this component is necessary because otherwise the problem space is under-constrained, as we have argued above. Moreover, it provides a test-bed for performing experimentation that would be impossible in reality (e.g. beyond what would be feasible or ethical in the lab), which is useful for both scientific and practical reasons. The second component—where feasible—allows linking these predictions to the explanatory understanding that cognitive scientists have usually sought, and that may help us to generalize our insights beyond the current paradigms. This perspective aligns with past arguments on how deep learning can contribute to understanding in cognitive neuroscience (Saxe et al., 2021; Kanwisher et al., 2023; Cao and Yamins, 2024b; Doerig et al., 2023). However, we emphasize that combining these two components makes explicit links between the actual problem the

system solves and the theorized constructs underlying that solution. It also ensures that our theories make testable predictions beyond the narrow settings of a particular task paradigm.

# 7 Discussion

In this paper, we have tried to outline a direction of research that we call "naturalistic computational cognitive science" that aims to build generalizable models of cognition that scale to naturalistic tasks, while still offering routes to explanatory and theoretical understanding. This perspective synthesizes a growing literature on the importance of naturalistic experiments and generalizable models, and grows out of a long-standing focus in cognitive science on explaining a broad range of phenomena. We outline these connections here.

**The quest for building general models of natural intelligence.** A long history of cognitive science research has sought to develop frameworks with the generality to explain a wide scope of cognitive phenomena. Researchers in connectionist (e.g. McClelland et al., 1986, 2003), Bayesian (e.g. Tenenbaum and Griffiths, 2001), and cognitive architectures (e.g. Ritter et al., 2019; Laird, 2019) paradigms have tried to identify underlying principles or mechanisms that explain many phenomena. However, while the modeling frameworks themselves are general, typical *instantiations* of these frameworks have built specialized models focused on individual tasks.

Other recent studies have built upon foundation models from AI that can perform many tasks. Some works have applied cognitive paradigms to study the behaviors of these models (e.g. Binz and Schulz, 2023; Lampinen et al., 2024; Buschoff et al., 2025). Other recent works have finetuned these models using cognitive data, to create generalizable cognitive models that can predict behavior on new experimental paradigms in areas like vision (e.g. Muttenthaler et al., 2024b,a; Fu et al., 2024), or the broader space of cognitive tasks that can be presented in language (Binz and Schulz, 2024; Binz et al., 2024).

We see the arguments that we have laid out in this paper as being broadly compatible with many of these approaches and frameworks, but offering a stronger emphasis on the virtuous cycle of expanding the scope and naturalism of our experimental designs, and expanding the generalizability of our models and theories—and anchoring these in the practicalities of achieving frictionless reproducibility.

**Towards experiments that involve naturalistic behavior** In emphasizing more naturalistic experiments, our perspective aligns with a variety of works that have likewise highlighted the importance of considering the scope of complex naturalistic behaviors. This perspective has been advocated frequently in neuroscience (e.g. Mobbs et al., 2021; Cisek and Green, 2024), where even prominent visual processing signals can be fundamentally different in active, naturalistic paradigms than in simpler classic ones (e.g. Amme et al., 2024). Likewise, Wise et al. (2023) advocate for the need to explore more naturalistic environments and behaviors in reinforcement learning, to grapple with the complexity of behavior in environments closer to

the real world. We concur with these works, and emphasize that advances in technology both enable richer experimental paradigms, and the models that can perform them. We differ from many of these works in emphasizing the role of models that generalize across task paradigms.

**Linking between learning models and the brain** A large recent literature has explored the surprising alignment between the representations learned by task-optimized deep learning models and those in the brain (Khaligh-Razavi and Kriegeskorte, 2014; Yamins et al., 2014; Yamins and DiCarlo, 2016; Schrimpf et al., 2021; Sucholutsky et al., 2023). The observations of surprising alignment have led to substantial interest and debate. The reactions range from critical argument that failures to capture particular phenomena doom these models (e.g. Bowers et al., 2023), to suggestions that deep learning models upend the theoretical assumptions of broad areas of the field (Perconti and Plebe, 2020; Hasson et al., 2020). Other authors have written frameworks that attempt to integrate deep learning models within the existing scientific paradigms (Doerig et al., 2023). In this vein, several papers have proposed frameworks for understanding the role of deep learning in cognitive neuroscience (e.g Richards et al., 2019; Storrs and Kriegeskorte, 2019; Cichy and Kaiser, 2019)—and how these models can contribute to deriving explanatory and theoretical understanding (Cao and Yamins, 2024a,b; Saxe et al., 2021). In this context, Feather et al. (2025) argue for the importance of benchmarks that place stronger constraints on models, by testing alignment of both behavior and representations. As above, our perspective aligns with many of these works, but we place greater emphasis on the complementary role of increasing naturalism of experimental paradigms, and learning from the engineering paradigms of AI, rather than simply adopting its models as artifacts.

**Conclusions** We have outlined a perspective that advocates for studying cognition through a broad spectrum of naturalistic experimental paradigms, building generalizable models of cognition that can perform naturalistic tasks, and deriving explanatory theories from these models and their simplifications. We have supported this perspective with examples illustrating the value of naturalistic settings and generalizable models, and concrete guidance on the practicalities of building generalizable models and using them to shape our theories. Many individual aspects of our argument align with prior works, but we hope that there is value in synthesizing these perspectives and highlighting their synergy. We hope that this perspective will inspire computational cognitive scientists to continue embracing richer naturalistic experimental paradigms and models.

## Acknowledgements

# References

Abel, D., Barreto, A., Van Roy, B., Precup, D., van Hasselt, H. P., and Singh, S. (2024). A definition of continual reinforcement learning. *Advances in Neural Information Processing Systems*, 36.

Agarwal, R., Schwarzer, M., Castro, P. S., Courville, A. C., and Bellemare, M. (2021). Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320.

Allen, K., Brändle, F., Botvinick, M., Fan, J. E., Gershman, S. J., Gopnik, A., Griffiths, T. L., Hartshorne, J. K., Hauser, T. U., Ho, M. K., et al. (2024). Using games to understand the mind. *Nature Human Behaviour*, pages 1–9.

Almaatouq, A., Griffiths, T. L., Suchow, J. W., Whiting, M. E., Evans, J., and Watts, D. J. (2024). Beyond playing 20 questions with nature: Integrative experiment design in the social and behavioral sciences. *Behavioral and Brain Sciences*, 47:e33.

Amme, C., Sulewski, P., Spaak, E., Hebart, M. N., König, P., and Kietzmann, T. C. (2024). Saccade onset, not fixation onset, best explains early responses across the human visual cortex during naturalistic vision. *bioRxiv*, pages 2024–10.

Anderson, J. R. (1990). The adaptive character of thought.

Anderson, J. R. (1991). Is human cognition adaptive? *Behavioral and brain sciences*, 14(3):471–485.

Anderson, P. W. (1972). More is different: Broken symmetry and the nature of the hierarchical structure of science. *Science*, 177(4047):393–396.

Andrews, M. (2021). The math is not the territory: navigating the free energy principle. *Biology & Philosophy*, 36(3):30.

Bainbridge, W. A., Hall, E. H., and Baker, C. I. (2019). Drawings of real-world scenes during free recall reveal detailed object and spatial information in memory. *Nature communications*, 10(1):5.

Baker, N., Lu, H., Erlikhman, G., and Kellman, P. J. (2020). Local features and global shape information in object classification by deep convolutional neural networks. *Vision research*, 172:46–61.

Baribault, B., Donkin, C., Little, D. R., Trueblood, J. S., Oravecz, Z., Van Ravenzwaaij, D., White, C. N., De Boeck, P., and Vandekerckhove, J. (2018). Metastudies for robust tests of theory. *Proceedings of the National Academy of Sciences*, 115(11):2607–2612.

Beckers, T., Krypotos, A.-M., Boddez, Y., Effting, M., and Kindt, M. (2013). What's wrong with fear conditioning? *Biological psychology*, 92(1):90–96.

Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. (2013). The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279.

Binz, M., Akata, E., Bethge, M., Brändle, F., Callaway, F., Coda-Forno, J., Dayan, P., Demircan, C., Eckstein, M. K., Éltető, N., Griffiths, T. L., Haridi, S., Jagadish, A. K., Ji-An, L., Kipnis, A., Kumar, S., Ludwig, T., Mathony, M., Mattar, M., Modirshanechi, A., Nath, S. S., Peterson, J. C., Rmus, M., Russek, E. M., Saanum, T., Scharfenberg, N., Schubert, J. A., Buschoff, L. M. S., Singhi, N., Sui, X., Thalmann, M., Theis, F., Truong, V., Udandarao, V., Voudouris, K., Wilson, R., Witte, K., Wu, S., Wulff, D., Xiong, H., and Schulz, E. (2024). Centaur: a foundation model of human cognition.

Binz, M. and Schulz, E. (2023). Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120.

Binz, M. and Schulz, E. (2024). Turning large language models into cognitive models. In *Twelfth International Conference on Learning Representations (ICLR)*.

Bonnen, T. and Eldridge, M. A. (2023). Inconsistencies between human and macaque lesion data can be resolved with a stimulus-computable model of the ventral visual stream. *Elife*, 12:e84357.

Bonnen, T., Yamins, D. L., and Wagner, A. D. (2021). When the ventral visual stream is not enough: A deep learning account of medial temporal lobe involvement in perception. *Neuron*, 109(17):2755–2766.

Bowers, J. S., Malhotra, G., Dujmović, M., Montero, M. L., Tsvetkov, C., Biscione, V., Puebla, G., Adolfi, F., Hummel, J. E., Heaton, R. F., et al. (2023). Deep problems with neural network models of human vision. *Behavioral and Brain Sciences*, 46:e385.

Britton, J. C., Evans, T. C., and Hernandez, M. V. (2014). Looking beyond fear and extinction learning: considering novel treatment targets for anxiety. *Current behavioral neuroscience reports*, 1:134–143.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Buffalo, E. A., Reber, P. J., and Squire, L. R. (1998). The human perirhinal cortex and recognition memory. *Hippocampus*, 8(4):330–339.

Buschoff, L. M. S., Akata, E., Bethge, M., and Schulz, E. (2025). Visual cognition in multimodal large language models. *Nature Machine Intelligence*.

Busemeyer, J. R. and Wang, Y.-M. (2000). Model comparisons and model selections based on generalization criterion methodology. *Journal of Mathematical Psychology*, 44(1):171–189.

Cao, R. and Yamins, D. (2024a). Explanatory models in neuroscience, part 1: Taking mechanistic abstraction seriously. *Cognitive Systems Research*, page 101244.

Cao, R. and Yamins, D. (2024b). Explanatory models in neuroscience, part 2: Functional intelligibility and the contravariance principle. *Cognitive Systems Research*, 85:101200.

Carlson, R. W. and Crockett, M. J. (2018). The lateral prefrontal cortex and moral goal pursuit. *Current opinion in psychology*, 24:77–82.

Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., and Joulin, A. (2021). Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660.

Carvalho, W., Tomov, M. S., de Cothi, W., Barry, C., and Gershman, S. J. (2024). Predictive representations: building blocks of intelligence. *arXiv preprint arXiv:2402.06590*.

Chan, S., Santoro, A., Lampinen, A., Wang, J., Singh, A., Richemond, P., McClelland, J., and Hill, F. (2022). Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891.

Cheng, X., Shi, K., Agarwal, A., and Pathak, D. (2024). Extreme parkour with legged robots. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11443–11450. IEEE.

Chentanez, N., Barto, A., and Singh, S. (2004). Intrinsically motivated reinforcement learning. *Advances in neural information processing systems*, 17.

Chomsky, N. (1957). *Syntactic structures*. Mouton de Gruyter.

Chomsky, N. (1965). *Aspects of the Theory of Syntax*. MIT press.

Chomsky, N. (2014). *The minimalist program*. MIT press.

Cichy, R. M. and Kaiser, D. (2019). Deep neural networks as scientific models. *Trends in cognitive sciences*, 23(4):305–317.

Cisek, P. (1999). Beyond the computer metaphor: Behaviour as interaction. *Journal of Consciousness Studies*, 6(11-12):125–142.

Cisek, P. and Green, A. M. (2024). Toward a neuroscience of natural behavior. *Current Opinion in Neurobiology*, 86:102859.

Clark, H. H. (1973). The language-as-fixed-effect fallacy: A critique of language statistics in psychological research. *Journal of verbal learning and verbal behavior*, 12(4):335–359.

Cohen, J. D., Dunbar, K., and McClelland, J. L. (1990). On the control of automatic processes: a parallel distributed processing account of the stroop effect. *Psychological review*, 97(3):332.

Collins, A. G. (2024). RL or not RL? Parsing the processes that support human reward-based learning.

Cook, J., Lu, C., Hughes, E., Leibo, J. Z., and Foerster, J. (2024). Artificial generational intelligence: Cultural accumulation in reinforcement learning. *arXiv preprint arXiv:2406.00392*.

Cristianini, N. and Scholkopf, B. (2002). Support vector machines and kernel methods: the new generation of learning machines. *Ai Magazine*, 23(3):31–31.

Dadashi, R., Taiga, A. A., Le Roux, N., Schuurmans, D., and Bellemare, M. G. (2019). The value function polytope in reinforcement learning. In *International Conference on Machine Learning*, pages 1486–1495. PMLR.

De Silva, A., Ramesh, R., Ungar, L., Shuler, M. H., Cowan, N. J., Platt, M., Li, C., Isik, L., Roh, S.-E., Charles, A., et al. (2023). Prospective learning: Principled extrapolation to the future. In *Conference on Lifelong Learning Agents*, pages 347–357. PMLR.

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee.

Deniz, F., Nunez-Elizalde, A. O., Huth, A. G., and Gallant, J. L. (2019). The representation of semantic information across human cerebral cortex during listening versus reading is invariant to stimulus modality. *Journal of Neuroscience*, 39(39):7722–7736.

Dentella, V., Günther, F., Murphy, E., Marcus, G., and Leivada, E. (2024). Testing ai on language comprehension tasks reveals insensitivity to underlying meaning. *Scientific Reports*, 14(1):28083.

DiCarlo, J. J., Yamins, D. L., Ferguson, M. E., Fedorenko, E., Bethge, M., Bonnen, T., and Schrimpf, M. (2023). Let's move forward: Image-computable models and a common model evaluation scheme are prerequisites for a scientific understanding of human vision. *Behavioral and Brain Sciences*, 46:e390.

Dobs, K., Martinez, J., Kell, A. J., and Kanwisher, N. (2022). Brain-like functional specialization emerges spontaneously in deep neural networks. *Science advances*, 8(11):eabl8913.

Doerig, A., Sommers, R. P., Seeliger, K., Richards, B., Ismael, J., Lindsay, G. W., Kording, K. P., Konkle, T., Van Gerven, M. A., Kriegeskorte, N., et al. (2023). The neuroconnectionist research programme. *Nature Reviews Neuroscience*, 24(7):431–450.

Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., and Darrell, T. (2014). Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, pages 647–655. PMLR.

Donoho, D. (2024). Data science at the singularity. *Harvard Data Science Review*, 6(1).

Doshi, F. R. and Konkle, T. (2023). Cortical topographic motifs emerge in a self-organized map of object space. *Science Advances*, 9(25):eade8187.

Eckstein, M. K., Master, S. L., Xia, L., Dahl, R. E., Wilbrecht, L., and Collins, A. G. (2022). The interpretation of computational model parameters depends on the context. *Elife*, 11:e75474.

Eckstein, M. K., Wilbrecht, L., and Collins, A. G. (2021). What do reinforcement learning models measure? interpreting model parameters in cognition and neuroscience. *Current opinion in behavioral sciences*, 41:128–137.

Eisenberg, I. W., Bissett, P. G., Zeynep Enkavi, A., Li, J., MacKinnon, D. P., Marsch, L. A., and Poldrack, R. A. (2019). Uncovering the structure of self-regulation through data-driven ontology discovery. *Nature communications*, 10(1):2319.

Ellis, K. (2024). Human-like few-shot learning via bayesian reasoning over natural language. *Advances in Neural Information Processing Systems*, 36.

Elman, J. L. (1996). *Rethinking innateness: A connectionist perspective on development*, volume 10. MIT press.

Engelhard, I. M., van den Hout, M. A., and McNally, R. J. (2008). Memory consistency for traumatic events in dutch soldiers deployed to iraq. *Memory*, 16(1):3–9.

Erev, I. and Greiner, B. (2015). The 1-800 critique, counter-examples, and the future of behavioral economics. *Handbook of experimental economic methodology*, pages 151–165.

Fan, L., Wang, G., Jiang, Y., Mandlekar, A., Yang, Y., Zhu, H., Tang, A., Huang, D.-A., Zhu, Y., and Anandkumar, A. (2022). Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362.

Farebrother, J., Orbay, J., Vuong, Q., Taïga, A. A., Chebotar, Y., Xiao, T., Irpan, A., Levine, S., Castro, P. S., Faust, A., et al. (2024). Stop regressing: Training value functions via classification for scalable deep rl. *arXiv preprint arXiv:2403.03950*.

Feather, J., Khosla, M., Murty, N. A. R., and Nayebi, A. (2025). Brain-model evaluations need the neuroai turing test.

Fei-Fei, L., Fergus, R., and Perona, P. (2004). Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *2004 conference on computer vision and pattern recognition workshop*, pages 178–178. IEEE.

Fodor, J. A. et al. (1975). *The language of thought*, volume 5. Harvard university press Cambridge, MA.

Fodor, J. A. and Pylyshyn, Z. W. (1988). Connectionism and cognitive architecture: A critical analysis. *Cognition*, 28(1-2):3–71.

Francis, K. B., Terbeck, S., Briazu, R. A., Haines, A., Gummerum, M., Ganis, G., and Howard, I. S. (2017). Simulating moral actions: An investigation of personal force in virtual moral dilemmas. *Scientific Reports*, 7(1):13954.

Frank, M. C. and Goodman, N. D. (2012). Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998.

Freeman, W. J. (1995). Chaos in the brain: Possible roles in biological intelligence. *International Journal of Intelligent Systems*, 10(1):71–88.

Fu, S., Tamir, N., Sundaram, S., Chai, L., Zhang, R., Dekel, T., and Isola, P. (2024). Dreamsim: Learning new dimensions of human visual similarity using synthetic data. *Advances in Neural Information Processing Systems*, 36.

Futrell, R., Gibson, E., Tily, H. J., Blank, I., Vishnevetsky, A., Piantadosi, S. T., and Fedorenko, E. (2021). The natural stories corpus: a reading-time corpus of english texts containing rare syntactic constructions. *Language Resources and Evaluation*, 55:63–77.

Futrell, R. and Mahowald, K. (2025). How linguistics learned to stop worrying and love the language models.

Gatys, L. A., Ecker, A. S., and Bethge, M. (2016). Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423.

Geiger, A., Lu, H., Icard, T., and Potts, C. (2021). Causal abstractions of neural networks. *Advances in Neural Information Processing Systems*, 34:9574–9586.

Geiger, A., Wu, Z., Potts, C., Icard, T., and Goodman, N. (2024). Finding alignments between interpretable causal variables and distributed neural representations. In *Causal Learning and Reasoning*, pages 160–187. PMLR.

Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., and Brendel, W. (2018). Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*.

Gershman, S. J. (2018). The successor representation: its computational logic and neural substrates. *Journal of Neuroscience*, 38(33):7193–7200.

Gold, E. M. (1967). Language identification in the limit. *Information and control*, 10(5):447–474.

Grim, P., Rosenberger, R., Rosenfeld, A., Anderson, B., and Eason, R. E. (2013). How simulations fail. *Synthese*, 190:2367–2390.

Grimm, C., Barreto, A., Singh, S., and Silver, D. (2020). The value equivalence principle for model-based reinforcement learning. *Advances in neural information processing systems*, 33:5541–5552.

Guest, O. and Martin, A. E. (2021). How computational modeling can force theory building in psychological science. *Perspectives on Psychological Science*, 16(4):789–802.

Hardt, M. and Recht, B. (2022). *Patterns, predictions, and actions: Foundations of machine learning*. Princeton University Press.

Hardwicke, T. E., Mathur, M. B., MacDonald, K., Nilsonne, G., Banks, G. C., Kidwell, M. C., Hofelich Mohr, A., Clayton, E., Yoon, E. J., Henry Tessler, M., et al. (2018). Data availability, reusability, and analytic reproducibility: Evaluating the impact of a mandatory open data policy at the journal cognition. *Royal Society open science*, 5(8):180448.

Hasson, U., Nastase, S. A., and Goldstein, A. (2020). Direct fit to nature: an evolutionary perspective on biological and artificial neural networks. *Neuron*, 105(3):416–434.

Hastie, T., Tibshirani, R., Friedman, J. H., and Friedman, J. H. (2009). *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer.

He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. (2020). Momentum contrast for unsu-pervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738.

He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.

Helbing, J., Draschkow, D., and Vo, M. L.-H. (2020). Search superiority: Goal-directed attentional allocation creates more reliable incidental identity and location memory than explicit encoding in naturalistic virtual environments. *Cognition*, 196:104147.

Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., and Meger, D. (2018). Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.

Henrich, J. (2016). *The secret of our success: How culture is driving human evolution, domesticating our species, and making us smarter*. princeton University press.

Higgins, I., Chang, L., Langston, V., Hassabis, D., Summerfield, C., Tsao, D., and Botvinick, M. (2021). Unsupervised deep learning identifies semantic disentanglement in single infer-otemporal face patch neurons. *Nature communications*, 12(1):6456.

Highleyman, W. H. and Kamentsky, L. A. (1959). A generalized scanner for pattern-and character-recognition studies. In *Papers presented at the the March 3-5, 1959, western joint computer conference*, pages 291–294.

Hill, F., Lampinen, A., Schneider, R., Clark, S., Botvinick, M., McClelland, J. L., and Santoro, A. (2020). Environmental drivers of systematicity and generalization in a situated agent. In *International Conference on Learning Representations*.

Holzmeister, F., Johannesson, M., Böhm, R., Dreber, A., Huber, J., and Kirchler, M. (2024). Heterogeneity in effect size estimates. *Proceedings of the National Academy of Sciences*, 121(32):e2403490121.

Hosseini, E. A., Zaslavsky, N., Casto, C., and Fedorenko, E. (2023). Teasing apart the representational spaces of ann language models to discover key axes of model-to-brain alignment. In *Conference on Cognitive Computational Neuroscience (CCN 2023), Oxford, UK, Aug*, pages 24–27.

Hu, J., Wellman, M. P., et al. (1998). Multiagent reinforcement learning: theoretical frame-work and an algorithm. In *ICML*, volume 98, pages 242–250.

Huang, S., Dossa, R. F. J., Ye, C., Braga, J., Chakraborty, D., Mehta, K., and AraÃšjo, J. G. (2022). Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18.

Hughes, E., Dennis, M., Parker-Holder, J., Behbahani, F., Mavalankar, A., Shi, Y., Schaul, T., and Rocktaschel, T. (2024). Open-endedness is essential for artificial superhuman intelligence. *arXiv preprint arXiv:2406.04268*.

Jagadeesh, A. V. and Gardner, J. L. (2022a). Human visual cortex as a texture basis set for object perception. *Journal of Vision*, 22(14):3694–3694.

Jagadeesh, A. V. and Gardner, J. L. (2022b). Texture-like representation of objects in human visual cortex. *Proceedings of the National Academy of Sciences*, 119(17):e2115302119.

Jain, S., Vo, V. A., Wehbe, L., and Huth, A. G. (2024). Computational language modeling and the promise of in silico experimentation. *Neurobiology of Language*, 5(1):80–106.

Jones, M. and Love, B. C. (2011). Bayesian fundamentalism or enlightenment? on the explanatory status and theoretical contributions of bayesian models of cognition. *Behavioral and brain sciences*, 34(4):169–188.

Kallini, J., Papadimitriou, I., Futrell, R., Mahowald, K., and Potts, C. (2024). Mission: Impossible language models. *arXiv preprint arXiv:2401.06416*.

Kanwisher, N., Khosla, M., and Dobs, K. (2023). Using artificial neural networks to ask 'why'questions of minds and brains. *Trends in Neurosciences*, 46(3):240–254.

Kanwisher, N. and Yovel, G. (2006). The fusiform face area: a cortical region specialized for the perception of faces. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 361(1476):2109–2128.

Kégl, B. and Busa-Fekete, R. (2009). Boosting products of base classifiers. In *Proceedings of the 26th annual international conference on machine learning*, pages 497–504.

Kell, A. J., Yamins, D. L., Shook, E. N., Norman-Haignere, S. V., and McDermott, J. H. (2018). A task-optimized neural network replicates human auditory behavior, predicts brain responses, and reveals a cortical processing hierarchy. *Neuron*, 98(3):630–644.

Kellert, S. H. (1993). *In the wake of chaos: Unpredictable order in dynamical systems*. University of Chicago press.

Khaligh-Razavi, S.-M. and Kriegeskorte, N. (2014). Deep supervised, but not unsupervised, models may explain it cortical representation. *PLoS computational biology*, 10(11):e1003915.

Khan, H. A. (2017). Mcs hog features and svm based handwritten digit recognition system. *Journal of Intelligent Learning Systems and Applications*, 9(02):21–33.

Kingma, D. P. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

Knoch, D., Pascual-Leone, A., Meyer, K., Treyer, V., and Fehr, E. (2006). Diminishing reciprocal fairness by disrupting the right prefrontal cortex. *science*, 314(5800):829–832.

Konkle, T. and Alvarez, G. A. (2022). A self-supervised domain-general learning framework for human ventral stream representation. *Nature communications*, 13(1):491.

Kriegeskorte, N. and Douglas, P. K. (2018). Cognitive computational neuroscience. *Nature neuroscience*, 21(9):1148–1160.

Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25.

Krypotos, A.-M., Vervliet, B., and Engelhard, I. M. (2018). The validity of human avoidance paradigms. *Behaviour Research and Therapy*, 111:99–105.

Laird, J. E. (2019). *The Soar cognitive architecture*. MIT press.

Lampinen, A. K., Dasgupta, I., Chan, S. C., Sheahan, H. R., Creswell, A., Kumaran, D., McClelland, J. L., and Hill, F. (2024). Language models, like humans, show content effects on reasoning tasks. *PNAS nexus*, 3(7).

LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.

Lee, H., Ekanadham, C., and Ng, A. (2007). Sparse deep belief net model for visual area v2. *Advances in neural information processing systems*, 20.

Lee, I. (2014). Publish or perish: The myth and reality of academic publishing. *Language teaching*, 47(2):250–261.

Levine, S., Finn, C., Darrell, T., and Abbeel, P. (2016). End-to-end training of deep visuo-motor policies. *Journal of Machine Learning Research*, 17(39):1–40.

Lewis, R. L., Howes, A., and Singh, S. (2014). Computational rationality: Linking mechanism and behavior through bounded utility maximization. *Topics in cognitive science*, 6(2):279–311.

Li, B. Z., Nye, M., and Andreas, J. (2021a). Implicit representations of meaning in neural language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1813–1827.

Li, C., Xia, F., Martín-Martín, R., Lingelbach, M., Srivastava, S., Shen, B., Vainio, K., Gokmen, C., Dharan, G., Jain, T., et al. (2021b). igibson 2.0: Object-centric simulation for robot learning of everyday household tasks. *arXiv preprint arXiv:2108.03272*.

Lilienfeld, S. O. (2017). Psychology's replication crisis and the grant culture: Righting the ship. *Perspectives on psychological science*, 12(4):660–664.

Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., Van Der Laak, J. A., Van Ginneken, B., and Sánchez, C. I. (2017). A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88.

Liu, G., Tang, M., and Eysenbach, B. (2024). A single goal is all you need: Skills and exploration emerge from contrastive rl without rewards, demonstrations, or subgoals. *arXiv preprint arXiv:2408.05804*.

Locatello, F., Weissenborn, D., Unterthiner, T., Mahendran, A., Heigold, G., Uszkoreit, J., Dosovitskiy, A., and Kipf, T. (2020). Object-centric learning with slot attention. *Advances in neural information processing systems*, 33:11525–11538.

Lockwood, P. L., Klein-Flügge, M. C., Abdurahman, A., and Crockett, M. J. (2020). Model-free decision making is prioritized when learning to avoid harming others. *Proceedings of the National Academy of Sciences*, 117(44):27719–27730.

Lu, K., Grover, A., Abbeel, P., and Mordatch, I. (2022a). Frozen pretrained transformers as universal computation engines. In *Proceedings of the AAAI conference on artificial intelligence*, number 7, pages 7628–7636.

Lu, S., Zheng, M., Fontaine, M. C., Nikolaidis, S., and Culbertson, H. (2022b). Preference-driven texture modeling through interactive generation and search. *IEEE transactions on haptics*, 15(3):508–520.

Lucic, M., Kurach, K., Michalski, M., Gelly, S., and Bousquet, O. (2018). Are gans created equal? a large-scale study. *Advances in neural information processing systems*, 31.

Lyle, C., Rowland, M., and Dabney, W. (2022). Understanding and preventing capacity loss in reinforcement learning. *arXiv preprint arXiv:2204.09560*.

Machado, M. C., Bellemare, M. G., Talvitie, E., Veness, J., Hausknecht, M., and Bowling, M. (2018). Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562.

Manning, C. D., Clark, K., Hewitt, J., Khandelwal, U., and Levy, O. (2020). Emergent linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the National Academy of Sciences*, 117(48):30046–30054.

Marcus, G. F. (1998). Rethinking eliminative connectionism. *Cognitive psychology*, 37(3):243–282.

Maren, S. (2001). Neurobiology of pavlovian fear conditioning. *Annual review of neuroscience*, 24(1):897–931.

Matthews, M., Beukman, M., Ellis, B., Samvelyan, M., Jackson, M., Coward, S., and Foerster, J. (2024). Craftax: A lightning-fast benchmark for open-ended reinforcement learning. *arXiv preprint arXiv:2402.16801*.

McClelland, J., Rumelhart, D., and Hinton, G. (1986). The appeal of parallel distributed processing. In *Parallel distributed processing: explorations in the microstructure of cognition, vol. 1: foundations*, pages 3–44.

McClelland, J. L. (2009). The place of modeling in cognitive science. *Topics in Cognitive Science*, 1(1):11–38.

McClelland, J. L. and Jenkins, E. (1991). Nature, nurture, and connections: Implications of connectionist models for cognitive development. In *Architectures for intelligence*, pages 41–73. Psychology Press.

McClelland, J. L., Plaut, D. C., Gotts, S. J., and Maia, T. V. (2003). Developing a domain-general framework for cognition: What is the best approach? *Behavioral and Brain Sciences*, 26(5):611–614.

Melis, G., Dyer, C., and Blunsom, P. (2017). On the state of the art of evaluation in neural language models. *arXiv preprint arXiv:1707.05589*.

Merel, J., Botvinick, M., and Wayne, G. (2019). Hierarchical motor control in mammals and machines. *Nature communications*, 10(1):1–12.

Miller, K., Eckstein, M., Botvinick, M., and Kurth-Nelson, Z. (2024). Cognitive model discovery via disentangled rnns. *Advances in Neural Information Processing Systems*, 36.

Misra, K. and Mahowald, K. (2024). Language models learn rare phenomena from less rare phenomena: The case of the missing aanns. *arXiv preprint arXiv:2403.19827*.

Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. (2015). Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533.

Mobbs, D., Wise, T., Suthana, N., Guzmán, N., Kriegeskorte, N., and Leibo, J. Z. (2021). Promises and challenges of human computational ethology. *Neuron*, 109(14):2224–2238.

Murray, E. A., Bussey, T. J., and Saksida, L. M. (2007). Visual perception and memory: a new view of medial temporal lobe function in primates and rodents. *Annu. Rev. Neurosci.*, 30(1):99–122.

Muttenthaler, L., Greff, K., Born, F., Spitzer, B., Kornblith, S., Mozer, M. C., Müller, K.-R., Unterthiner, T., and Lampinen, A. K. (2024a). Aligning machine and human visual representations across abstraction levels. *arXiv preprint arXiv:2409.06509*.

Muttenthaler, L., Linhardt, L., Dippel, J., Vandermeulen, R. A., Hermann, K., Lampinen, A., and Kornblith, S. (2024b). Improving neural network representations using human similarity judgments. *Advances in Neural Information Processing Systems*, 36.

Nanda, N., Chan, L., Lieberum, T., Smith, J., and Steinhardt, J. (2023). Progress measures for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning Representations*.

Naselaris, T., Bassett, D. S., Fletcher, A. K., Kording, K., Kriegeskorte, N., Nienborg, H., Poldrack, R. A., Shohamy, D., and Kay, K. (2018). Cognitive computational neuroscience: A new conference for an emerging discipline. *Trends in cognitive sciences*, 22(5):365–367.

Nau, M., Schmid, A. C., Kaplan, S. M., Baker, C. I., and Kravitz, D. J. (2024). Centering cognitive neuroscience on task demands and generalization. *Nature Neuroscience*.

Newell, A. (1973). You can't play 20 questions with nature and win: Projective comments on the papers of this symposium.

Newell, A. (1994). *Unified theories of cognition*. Harvard University Press.

Newen, A., De Bruin, L., and Gallagher, S. (2018). *The Oxford handbook of 4E cognition*. Oxford University Press.

Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G., et al. (2015). Promoting an open research culture. *Science*, 348(6242):1422–1425.

Oaksford, M. and Chater, N. (2009). Précis of bayesian rationality: The probabilistic approach to human reasoning. *Behavioral and Brain Sciences*, 32(1):69–84.

Obando Ceron, J., Bellemare, M., and Castro, P. S. (2024). Small batch deep reinforcement learning. *Advances in Neural Information Processing Systems*, 36.

Olshausen, B. A. and Field, D. J. (1996). Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609.

Osband, I., Doron, Y., Hessel, M., Aslanides, J., Sezener, E., Saraiva, A., McKinney, K., Lattimore, T., Szepesvari, C., Singh, S., et al. (2019). Behaviour suite for reinforcement learning. *arXiv preprint arXiv:1908.03568*.

Oudeyer, P.-Y., Baranes, A., and Kaplan, F. (2013). Intrinsically motivated learning of real-world sensorimotor skills with developmental constraints. *Intrinsically motivated learning in natural and artificial systems*, pages 303–365.

O'Byrne, J. and Jerbi, K. (2022). How critical is brain criticality? *Trends in Neurosciences*, 45(11):820–837.

Parisi, S., Rajeswaran, A., Purushwalkam, S., and Gupta, A. (2022). The unsurprising effectiveness of pre-trained vision models for control. In *international conference on machine learning*, pages 17359–17371. PMLR.

Patterson, A., Neumann, S., White, M., and White, A. (2023). Empirical design in reinforcement learning. *arXiv preprint arXiv:2304.01315*.

Perconti, P. and Plebe, A. (2020). Deep learning and cognitive science. *Cognition*, 203:104365.

Piantadosi, S. (2023). Modern language models refute chomsky's approach to language. *Lingbuzz Preprint, lingbuzz*, 7180.

Poldrack, R. A., Baker, C. I., Durnez, J., Gorgolewski, K. J., Matthews, P. M., Munafò, M. R., Nichols, T. E., Poline, J.-B., Vul, E., and Yarkoni, T. (2017). Scanning the horizon: towards transparent and reproducible neuroimaging research. *Nature reviews neuroscience*, 18(2):115–126.

Prystawski, B., Li, M., and Goodman, N. (2024). Why think step by step? reasoning emerges from the locality of experience. *Advances in Neural Information Processing Systems*, 36.

Puebla, G. and Bowers, J. S. (2022). Can deep convolutional neural networks support relational reasoning in the same-different task? *Journal of Vision*, 22(10):11–11.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. (2019). Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. (2024). Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

Rahnev, D. and Denison, R. N. (2018). Suboptimality in perceptual decision making. *Behavioral and brain sciences*, 41:e223.

Recht, B. (2024). The mechanics of frictionless reproducibility.

Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. (2019). Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, pages 5389–5400. PMLR.

Reid, M., Yamada, Y., and Gu, S. S. (2022). Can wikipedia help offline reinforcement learning? *arXiv preprint arXiv:2201.12122*.

Reithmeier, R., O'Leary, L., Zhu, X., Dales, C., Abdulkarim, A., Aquil, A., Brouillard, L., Chang, S., Miller, S., Shi, W., et al. (2019). The 10,000 phds project at the university of toronto: Using employment outcome data to inform graduate education. *PLoS One*, 14(1):e0209898.

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. (2019). A deep learning framework for neuroscience. *Nature neuroscience*, 22(11):1761–1770.

Riley, A. L. and Tuck, D. L. (1985). Conditioned food aversions: A bibliography. *Annals of the New York Academy of Sciences*, 443(1):381–437.

Riquelme, C., Tucker, G., and Snoek, J. (2018). Deep bayesian bandits showdown: An empirical comparison of bayesian deep networks for thompson sampling. *arXiv preprint arXiv:1802.09127*.

Ritter, F. E., Tehranchi, F., and Oury, J. D. (2019). Act-r: A cognitive architecture for modeling cognition. *Wiley Interdisciplinary Reviews: Cognitive Science*, 10(3):e1488.

Rogers, T. T. and McClelland, J. L. (2004). *Semantic cognition: A parallel distributed processing approach*. MIT press.

Rosenberg, M., Zhang, T., Perona, P., and Meister, M. (2021). Mice in a labyrinth show rapid learning, sudden insight, and efficient exploration. *Elife*, 10:e66175.

Ruff, C. C., Ugazio, G., and Fehr, E. (2013). Changing social norm compliance with noninvasive brain stimulation. *Science*, 342(6157):482–484.

Rust, N. C. and Movshon, J. A. (2005). In praise of artifice. *Nature neuroscience*, 8(12):1647–1650.

Rutherford, A., Ellis, B., Gallici, M., Cook, J., Lupu, A., Ingvarsson, G., Willi, T., Khan, A., de Witt, C. S., Souly, A., et al. (2023). Jaxmarl: Multi-agent rl environments in jax. *arXiv preprint arXiv:2311.10090*.

Sarafoglou, A., Kovacs, M., Bakos, B., Wagenmakers, E.-J., and Aczel, B. (2022). A survey on how preregistration affects the research workflow: Better science but more work. *Royal Society Open Science*, 9(7):211997.

Saxe, A., Nelli, S., and Summerfield, C. (2021). If deep learning is the answer, what is the question? *Nature Reviews Neuroscience*, 22(1):55–67.

Saxe, A. M., McClelland, J. L., and Ganguli, S. (2019). A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546.

Scheel, A. M., Schijen, M. R., and Lakens, D. (2021). An excess of positive results: Comparing the standard psychology literature with registered reports. *Advances in Methods and Practices in Psychological Science*, 4(2):25152459211007467.

Scheveneels, S., Boddez, Y., Vervliet, B., and Hermans, D. (2016). The validity of laboratory-based treatment research: Bridging the gap between fear extinction and exposure treatment. *Behaviour research and therapy*, 86:87–94.

Schrimpf, M., Blank, I. A., Tuckute, G., Kauf, C., Hosseini, E. A., Kanwisher, N., Tenenbaum, J. B., and Fedorenko, E. (2021). The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the National Academy of Sciences*, 118(45):e2105646118.

Schrimpf, M., Kubilius, J., Lee, M. J., Murty, N. A. R., Ajemian, R., and DiCarlo, J. J. (2020). Integrative benchmarking to advance neurally mechanistic models of human intelligence. *Neuron*, 108(3):413–423.

Schweinsberg, M., Feldman, M., Staub, N., van den Akker, O. R., van Aert, R. C., Van Assen, M. A., Liu, Y., Althoff, T., Heer, J., Kale, A., et al. (2021). Same data, different conclusions: Radical dispersion in empirical results when independent analysts operationalize and test the same hypothesis. *Organizational Behavior and Human Decision Processes*, 165:228–249.

Sehlmeyer, C., Schöning, S., Zwitserlood, P., Pfleiderer, B., Kircher, T., Arolt, V., and Konrad, C. (2009). Human fear conditioning and extinction in neuroimaging: a systematic review. *PloS one*, 4(6):e5865.

Seligman, M. E., Railton, P., Baumeister, R. F., and Sripada, C. (2013). Navigating into the future or driven by the past. *Perspectives on psychological science*, 8(2):119–141.

Shanahan, M. (2004). The frame problem.

Sharif Razavian, A., Azizpour, H., Sullivan, J., and Carlsson, S. (2014). Cnn features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 806–813.

Shrout, P. E. and Rodgers, J. L. (2018). Psychology, science, and knowledge construction: Broadening perspectives from the replication crisis. *Annual review of psychology*, 69(1):487–510.

Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. (2017). Mastering the game of go without human knowledge. *nature*, 550(7676):354–359.

Simpson, J. and Kelly, J. P. (2011). The impact of environmental enrichment in laboratory rats—behavioural and neurochemical aspects. *Behavioural brain research*, 222(1):246–264.

Singh, A. K., Moskovitz, T., Hill, F., Chan, S. C., and Saxe, A. M. (2024). What needs to go right for an induction head? a mechanistic study of in-context learning circuits and their formation. In *Forty-first International Conference on Machine Learning*.

Small, W. S. (1901). Experimental study of the mental processes of the rat. ii. *The American Journal of Psychology*, pages 206–239.

Spencer, N. J. (1973). Differences between linguists and nonlinguists in intuitions of grammaticality-acceptability. *Journal of psycholinguistic research*, 2(2):83–98.

Spitzer, M., Fischbacher, U., Herrnberger, B., Grön, G., and Fehr, E. (2007). The neural signature of social norm compliance. *Neuron*, 56(1):185–196.

Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958.

Storrs, K. R. and Kriegeskorte, N. (2019). Deep learning for cognitive neuroscience. *arXiv preprint arXiv:1903.01458*.

Strickland, B. and Suben, A. (2012). Experimenter philosophy: The problem of experimenter bias in experimental philosophy. *Review of Philosophy and Psychology*, 3:457–467.

Strouse, D., McKee, K., Botvinick, M., Hughes, E., and Everett, R. (2021). Collaborating with humans without human data. *Advances in Neural Information Processing Systems*, 34:14502–14515.

Sucholutsky, I., Muttenthaler, L., Weller, A., Peng, A., Bobu, A., Kim, B., Love, B. C., Grant, E., Groen, I., Achterberg, J., et al. (2023). Getting aligned on representational alignment. *arXiv preprint arXiv:2310.13018*.

Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.

Suzuki, W. A. (2009). Perception and the medial temporal lobe: evaluating the current evidence. *Neuron*, 61(5):657–666.

Szot, A., Clegg, A., Undersander, E., Wijmans, E., Zhao, Y., Turner, J., Maestre, N., Mukadam, M., Chaplot, D. S., Maksymets, O., et al. (2021). Habitat 2.0: Training home assistants to rearrange their habitat. *Advances in neural information processing systems*, 34:251–266.

Team, A. A., Bauer, J., Baumli, K., Baveja, S., Behbahani, F., Bhoopchand, A., Bradley-Schmieg, N., Chang, M., Clay, N., Collister, A., et al. (2023). Human-timescale adaptation in an open-ended task space. *arXiv preprint arXiv:2301.07608*.

Team, O. E. L., Stooke, A., Mahajan, A., Barros, C., Deck, C., Bauer, J., Sygnowski, J., Trebacz, M., Jaderberg, M., Mathieu, M., et al. (2021). Open-ended learning leads to generally capable agents. *arXiv preprint arXiv:2107.12808*.

Tenenbaum, J. B. and Griffiths, T. L. (2001). Generalization, similarity, and bayesian inference. *Behavioral and brain sciences*, 24(4):629–640.

Tuckute, G., Sathe, A., Srikant, S., Taliaferro, M., Wang, M., Schrimpf, M., Kay, K., and Fedorenko, E. (2024). Driving and suppressing the human language network using large language models. *Nature Human Behaviour*, 8(3):544–561.

Van Rooij, I. (2008). The tractable cognition thesis. *Cognitive science*, 32(6):939–984.

Vélez, N. and Gweon, H. (2021). Learning from other minds: An optimistic critique of reinforcement learning models of social learning. *Current opinion in behavioral sciences*, 38:110–115.

Villegas, R., Babaeizadeh, M., Kindermans, P.-J., Moraldo, H., Zhang, H., Saffar, M. T., Castro, S., Kunze, J., and Erhan, D. (2022). Phenaki: Variable length video generation from open domain textual descriptions. In *International Conference on Learning Representations*.

Vinje, W. E. and Gallant, J. L. (2000). Sparse coding and decorrelation in primary visual cortex during natural vision. *Science*, 287(5456):1273–1276.

Wang, J. X. (2021). Meta-learning in natural and artificial intelligence. *Current Opinion in Behavioral Sciences*, 38:90–95.

Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., et al. (2022). Emergent abilities of large language models. *Transactions on Machine Learning Research*.

Wilcox, E. G., Futrell, R., and Levy, R. (2023). Using computational models to test syntactic learnability. *Linguistic Inquiry*, pages 1–44.

Wilcox, E. G., Hu, M., Mueller, A., Linzen, T., Warstadt, A., Choshen, L., Zhuang, C., Cotterell, R., and Williams, A. (2024). Bigger is not always better: The importance of human-scale language modeling for psycholinguistics.

Wise, T., Emery, K., and Radulescu, A. (2023). Naturalistic reinforcement learning. *Trends in Cognitive Sciences*.

Woolston, C. (2021). Researchers' career insecurity needs attention and reform now, says international coalition. *Nature*.

Wu, Z., Geiger, A., Icard, T., Potts, C., and Goodman, N. (2024). Interpretability at scale: Identifying causal mechanisms in alpaca. *Advances in Neural Information Processing Systems*, 36.

Yamins, D. L. and DiCarlo, J. J. (2016). Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience*, 19(3):356–365.

Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., and DiCarlo, J. J. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the national academy of sciences*, 111(23):8619–8624.

Yarkoni, T. (2022). The generalizability crisis. *Behavioral and Brain Sciences*, 45:e1.

Yuan, Z., Xue, Z., Yuan, B., Wang, X., Wu, Y., Gao, Y., and Xu, H. (2022). Pre-trained image encoder for generalizable visual reinforcement learning. *Advances in Neural Information Processing Systems*, 35:13022–13037.

Zhang, Y., Kang, B., Hooi, B., Yan, S., and Feng, J. (2023). Deep long-tailed learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9):10795–10816.

Zhong, Z., Liu, Z., Tegmark, M., and Andreas, J. (2024). The clock and the pizza: Two stories in mechanistic explanation of neural networks. *Advances in Neural Information Processing Systems*, 36.

Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., and Torralba, A. (2014). Object detectors emerge in deep scene cnns. *arXiv preprint arXiv:1412.6856*.

Zhou, M., Hou, J., Luo, C., Wang, Y., Zhang, Z., and Peng, J. (2024). Scenex: Procedural controllable large-scale scene generation via large-language models. *arXiv preprint arXiv:2403.15698*.