

# From Word Models to World Models: Translating from Natural Language to the Probabilistic Language of Thought

Lionel Wong<sup>1\*</sup>, Gabriel Grand<sup>1\*</sup>, Alexander K. Lew<sup>1</sup>, Noah D. Goodman<sup>2</sup>, Vikash K. Mansinghka<sup>1</sup>, Jacob Andreas<sup>1</sup>, Joshua B. Tenenbaum<sup>1</sup>

*\*Equal contribution.*

<sup>1</sup>MIT, <sup>2</sup>Stanford

## Abstract

How does language inform our downstream thinking? In particular, how do humans make meaning from language—and how can we leverage a theory of linguistic meaning to build machines that think in more human-like ways? In this paper, we propose *rational meaning construction*, a computational framework for language-informed thinking that combines neural models of language with probabilistic models for rational inference. We frame linguistic meaning as a context-sensitive mapping from natural language into a *probabilistic language of thought* (PLoT)—a general-purpose symbolic substrate for probabilistic, generative world modeling. Our architecture integrates two powerful computational tools that have not previously come together: we model thinking with *probabilistic programs*, an expressive representation for flexible commonsense reasoning; and we model meaning construction with *large language models* (LLMs), which support broad-coverage translation from natural language utterances to code expressions in a probabilistic programming language. We illustrate our framework in action through examples covering four core domains from cognitive science: probabilistic reasoning, logical and relational reasoning, visual and physical reasoning, and social reasoning about agents and their plans. In each, we show that LLMs can generate context-sensitive translations that capture pragmatically-appropriate linguistic meanings, while Bayesian inference with the generated programs supports coherent and robust commonsense reasoning. We extend our framework to integrate cognitively-motivated symbolic modules (physics simulators, graphics engines, and goal-directed planning algorithms) to provide a unified commonsense thinking interface from language. Finally, we explore how language can drive the construction of world models themselves. We hope this work will help to situate contemporary developments in LLMs within a broader cognitive picture of human language and intelligence, providing a roadmap towards AI systems that synthesize the insights of both modern and classical computational perspectives.

## 1 Introduction

Language expresses the vast internal landscape of our thoughts. We use language to convey what we believe, what we are uncertain about, and what we do not know. We talk about what we see in the world around us, and what we imagine in real or wholly hypothetical futures. We discuss what we want and what we plan to do, and dissect what others want and what we think they will do. We build and pass on new bodies of knowledge in language—we ask questions and offer explanations, give commands and instructions, and propose and refute theories. Some of these ideas can be expressed in part through other means. But language stands apart for its flexibility and breadth, and its seeming proximity to our thoughts.

What *is* language? How does language get its meaning, and when should we say that a person or machine knows, understands, and can use it? What is the relationship between language and the rest of general cognition—what allows language to inform and support so much of thought? This paper focuses on these questions as they relate to *human* language and thought, in computational terms. What integrated cognitive theory can model how language relates to the other core systems of human cognition? If we seek to build AI systems that emulate how humans talk and think, what architecture can integrate language robustly into systems that support the full scope of our thought?

---

Code for the examples in this paper is available at: [github.com/gabegrand/world-models](https://github.com/gabegrand/world-models).

Correspondence: co-primary authors ([zyzzyva@mit.edu](mailto:zyzzyva@mit.edu), [gg@mit.edu](mailto:gg@mit.edu)); co-supervisors ([jda@mit.edu](mailto:jda@mit.edu), [jbt@mit.edu](mailto:jbt@mit.edu)).

Theories of cognition have long considered human language and thinking to be deeply related, but fundamentally distinct. *Thinking*, in many traditional cognitive theories, revolves around goal-directed world modeling, inference, and decision making—constructing mental models of the world that reflect prior beliefs, can be updated from new observations, and support rational prediction and decision making toward’s one’s goals (Craig, 1967; Gentner & Stevens, 2014; Johnson-Laird, 1980, 1989; Lake, Ullman, Tenenbaum, & Gershman, 2017; Morgan, 1999; Nersessian et al., 2010). *Language*, in contrast, centers around communicating these thoughts to others, and receiving their thoughts in turn. In most linguistic theories, human languages are mappings between the internal representations of thought and an externalizable symbol system, which might be phonemes, signs, or glyphs (Frege, 1892; Heim & Kratzer, 1998; Lewis, 1976). To produce language is to map thoughts into these external symbols, and to understand language is to transduce from these external symbols back into the representations of thought.

The theoretical distinction between language and thought rests on multiple intersecting lines of evidence. Prior to learning language, infants are born equipped with a powerful toolkit for modeling and thinking about the world, including an understanding of physical objects and events, and the goals and actions of agents (Spelke, 2022; Spelke & Kinzler, 2007), and general abilities for learning statistics and structure (Saffran, Senghas, & Trueswell, 2001; Xu et al., 2021). Building on these foundations, children acquire language from relatively sparse input data, rapidly generalizing beyond the utterances they hear to produce and understand entirely new ones (Bloom, 2002; L. Gleitman, 1990; L. R. Gleitman, Cassidy, Nappa, Papafragou, & Trueswell, 2005; Landauer & Dumais, 1997; Pinker, 1998; L. Smith & Yu, 2008); they then use language to acquire new concepts they would not get merely from direct experience (Carey, 2009; Gopnik, 1996; Wellman & Gelman, 1992). Language and thought also appear to operate in distinct but interacting brain systems: neuroimaging and neurological studies reveal a “language network” specialized for processing sentences, functionally and anatomically separate from but closely connected to brain networks supporting other aspects of general cognition (Fedorenko & Varley, 2016; Mahowald et al., 2023).

“These empirical findings have shaped decades of computational models in cognitive science and AI. To model the expressiveness of human cognition, an influential computational paradigm suggests that humans compose and execute mental programs in an internal *language of thought* (Fodor, 1975), a structured symbolic substrate for representing conceptual knowledge that provides a general interface to algorithms for problem solving and reasoning. These symbolic systems are not merely logic engines; they support our probabilistic inferences, and rich intuitive simulations (Goodman, Tenenbaum, & Gerstenberg, 2014; Oaksford & Chater, 2007; Russell & Norvig, 2021). This paradigm underlies many of the success stories in cognitive science and related applications in AI. It has influenced models that capture how people draw causal and explanatory inferences about facts and observations (Pearl, 1988; Pearl et al., 2000), learn and generalize concepts from few examples (Lake et al., 2017); plan actions over long time horizons and under complex conditions (Kaelbling & Lozano-Pérez, 2013; Russell & Norvig, 2021); imagine and predict the physical world (Battaglia, Hamrick, & Tenenbaum, 2013; Ullman, Spelke, Battaglia, & Tenenbaum, 2017); and reason about other agents with their own beliefs and goals (C. Baker, Saxe, & Tenenbaum, 2011). Within linguistics and natural language processing, in turn, this paradigm underlies *semantic parsing* systems designed to map from human language into symbolic computational representations. It has yielded AI systems that could follow instructions (Tellex et al., 2011) and answer natural language queries with respect to structured knowledge representations (Klein & Manning, 2003; Liang, 2016; Steedman, 2011; Y. W. Wong & Mooney, 2007); as well as cognitive models that capture how human children learn the grammar and meaning of expressions in their native language (Abend, Kwiatkowski, Smith, Goldwater, & Steedman, 2017; Chater & Manning, 2006; Frank, Goodman, & Tenenbaum, 2009; Gauthier, Levy, & Tenenbaum, 2018; Goldwater, Griffiths, & Johnson, 2009; Perfors, Tenenbaum, & Regier, 2011; Piantadosi, Tenenbaum, & Goodman, 2012).

Despite this progress, however, modular and symbolic models of language and thought have been dogged by persistent critiques of their scalability and scope. Cognitive and AI researchers over the years have carved off specific domains of world knowledge, constructing bespoke representations to model them without a general account of whether they would generalize to all of human knowledge, or how they could be scalably learned. Semantic parsing systems inherited these critiques, and faced additional challenges in implementing the mapping from sentences into symbolic representations. These mapping functions were either hand-engineered or learned from strong supervision on specific domains of language, limiting them to brittle, imperfect models of the breadth and complexity of real human discourse.

In just the last few years, a serious challenge has emerged to the traditional view of language and thought as distinct but interacting components of the mind, each modeled using structured representations. *Large language models* (LLMs) use a new generation of attention-based deep neural networks to learn the probabilistic distributions of words from vast datasets of human language, generally training on orders of magnitude more data than a human encounters in their lifetime (Bommasani et al., 2021; T. B. Brown et al., 2020; OpenAI, 2023c; Rae et al., 2021; Vaswani et al., 2017). The underlying computational objective that drives these models is not itself new. LLMs follow in the tradition of distributional approaches to discovering structure in language (Firth, 1957; Harris, 1954; Osgood, 1952), which seek to extract representations of meaning from statistical patterns in how words are used in context (Dumais et al., 2004; Griffiths, Steyvers, & Tenenbaum, 2007; Mikolov, Sutskever, Chen, Corrado, & Dean, 2013; Sahlgren, 2008). What is new, however, is the scale and scope of today’s distributional vision, which has expanded in stages. A first generation of LLMs, trained specifically to predict words in context, produced such fluent language that they challenged traditional symbolic approaches to modeling language (Devlin, Chang, Lee, & Toutanova, 2018; Peters et al., 1802; Radford et al., 2019). Their qualitative success, as well as internal representational probes, suggested that linguistic structures sufficient for grammatically coherent language could be learned entirely from modeling the statistics of words (Piantadosi, 2023; Tenney, Das, & Pavlick, 2019). By scaling to even larger datasets and neural networks, LLMs appeared to learn not only the structure of language, but capacities for some kinds of thinking; they could learn new words in context, and extract patterns in language from a few examples that they could generalize locally to similar cases (T. B. Brown et al., 2020). The most recent LLMs have been trained not only to model the statistics of language but explicitly to reason, with targeted supervision on instruction following, writing code, and other forms of human dialog and feedback in conversational contexts (Chen et al., 2021; OpenAI, 2023a, 2023c; Ouyang et al., 2022). They produce such fluent language on a wide variety of tasks that many have begun to ask whether merely more training of this sort, with increasing scale, could learn representations sufficient for general intelligence (Bubeck et al., 2023). Proponents of the most extreme “scaling hypothesis” have argued that because language is used to express so much of human thought, a sufficiently large and performant predictive language model would effectively *have* to construct an internal model of all of cognition (Branwen, 2022).

This theoretical vision has sparked both excitement and controversy, but proponents and critics agree that it raises its own questions about its long-term scalability—most significantly, what will be required to close the outstanding gaps between today’s LLMs and general cognitive models that reason systematically and consistently about the language they receive or produce. Current LLMs can produce impressive results on a set of linguistic inputs and then fail completely on others that make trivial alterations to the same underlying domain (Ullman, 2023); they mix confident answers to complex questions with equally confident, hallucinated language that does not reflect a consistent, calibrated notion of truth or belief (Bubeck et al., 2023; OpenAI, 2023c). These issues make it difficult to evaluate whether LLMs have acquired cognitive capacities such as social reasoning and theory of mind (Ullman, 2023), or to compare different kinds of world modeling and planning tasks (Valmeekam, Sreedharan, Marquez, Olmo, & Kambhampati, 2023). One approach to solving these problems is through additional data. Perhaps fully robust, systematic reasoning will finally emerge if models are trained on still more language, or supervised more explicitly on data from complex reasoning tasks. This scaling route raises practical questions about whether it will be possible to acquire enough data to train such a model, as well as theoretical questions whether more data and more parameters alone *will* in fact yield robust systems for thought. Another strategy in recent work seeks to build more robust cognitive capacities by augmenting LLMs with various external tools for structured representation and symbolic reasoning, such as calculators (Cobbe et al., 2021), logic engines (Weir & Van Durme, 2022), databases (Alon et al., 2022; Borgeaud et al., 2022; Izacard et al., 2022; Thoppilan et al., 2022), physics simulators (R. Liu et al., 2022), planners (B. Liu et al., 2023), and APIs for executing arbitrary code (Karpas et al., 2022; OpenAI, 2023c; Schick et al., 2023). But these new hybrid approaches resurrect many of the same long-term scalability challenges that confronted earlier semantic parsing and knowledge representation systems, by designing a menagerie of bespoke representations and tools without a broader account of how they will scale towards general models of language and thought.

In this paper, we consider a different approach to integrating the strengths of modern language models and classic symbolic architectures, one that draws on but also runs counter to recent trends in AI, in a sense flipping these scaling questions on their head. Instead of trying to turn models trained to predict language into models that might genuinely think—filling each gap in reasoning we discover through yet more

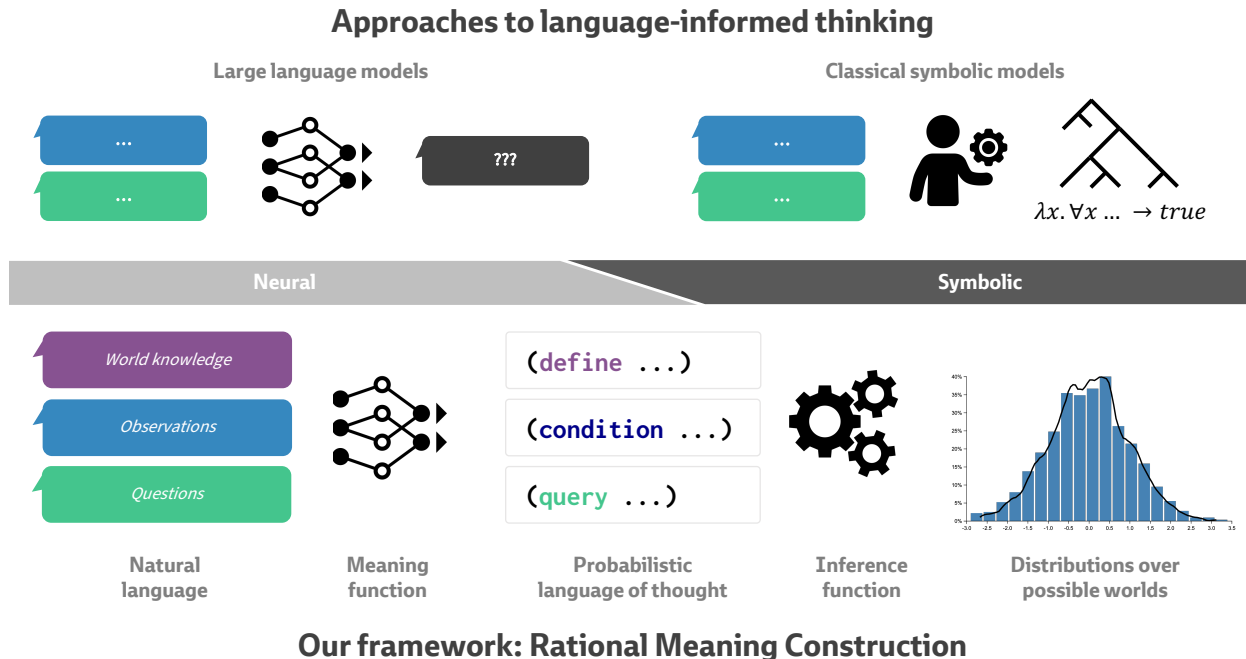


Figure 1: Human language understanding supports flexible inferences in a process we term *language-informed thinking*. Computational approaches to language-informed thinking sit on a neurosymbolic continuum: On one side, classical symbolic models (top right) yield systematic, structured inferences, but are typically limited to narrow linguistic domains and often require hand-engineering. On the other side, large language models (LLMs; top left) achieve remarkable facility with open-domain natural language, but struggle to ground reasoning in a consistent world state that supports coherent inferences, predictions and plans. Our *rational meaning construction* framework decomposes language-informed thinking into two modules: (1) A *meaning function* translates natural language into probabilistic programming language (PPL) statements that represent linguistic meaning with respect to a symbolic world model. (2) An *inference function* computes probabilities over the space of possible worlds consistent with and conditioned on information in the language. In the rest of this paper, we illustrate how our framework can combine the strengths of LLMs and PPLs, affording both broad coverage of natural language and a principled treatment of reasoning about uncertain events, outcomes, and scenarios.

data, new kinds of language training or linguistic prompting tricks, or by plugging in yet another external tool—we ask: what are the prospects for a unifying computational framework guided by the study of thought and language in the human mind and brain, as well as what we have learned from multiple eras of AI? Can we build intelligent architectures that use, learn and understand language as people do, informed by neuroscience constraints and developmental trajectories? That is, can we build models in which language is learned efficiently within one relatively small, modular computational system, which interfaces generally with other systems dedicated to robust world modeling and reasoning? What architecture lets language build on pre-existing capacities for symbolic world modeling and inference, while also allowing linguistic meanings and world knowledge to scaffold and bootstrap each other, as a learner’s experiences and competences grow?



This paper attempts to show what such a model might look like—and how it can build theoretically and practically on the insights from both classical paradigms for language and thought, *and* the recent successes of statistical learning made by large language models. We propose a framework for intelligent computational architectures that reason about and learn from language, but we begin with a proposal for what it means to *think*. As in the traditional cognitive view, thinking at its core is constructing general-purpose representations for modeling the entities and events in the world, sufficient to support rational, coherent inferences under

uncertainty and planning actions that achieve our goals. We then consider how language relates to this architecture to support *language-informed thinking*—how language sets up world modeling and inference, to guide, constrain, and drive our downstream thought, and grow new thinking capacities.

Our proposal, which we call **rational meaning construction**, rests on the integration of two computational components, each which we suggest can be instantiated using modern computational tools—a *probabilistic language of thought* for constructing structured models of arbitrary situations, which supports coherent belief updating and inferences over them; and a general mechanism for taking natural language and *constructing meaning* from it, represented as distributions over expressions in this language of thought (Fig. 1). We propose that *probabilistic programs* can formally instantiate the first component. They offer a structured representation for expressing novel situations and arbitrary problems with respect to a meaningful model over possible world states, a coherent notion of conditional belief updating, and a systematic framework for inferences with respect to queries and goals. We propose, in turn, that *meaning construction* can be modeled as *translation* from utterances in language to expressions in a general probabilistic programming language. Theoretical and empirical results have long suggested that human languages implement locally compositional, efficiently learnable mappings between symbolic representations of thought and external symbol systems. We therefore propose that code-trained large language models can be viewed as in-principle implementations of broad, context-sensitive, and *resource-rational meaning functions*, in that they can be used to efficiently infer distributions between language and programs from stored, prior patterns in the background distribution of language and code. By integrating these two components, we propose that this paradigm suggests a general framework by which language can meaningfully relate to many fundamental aspects of cognition, modeling how we might condition on language to systematically update our beliefs, pose new questions and goals in language, and convey structured background information or even define new relevant concepts about a situation or about the world.

In Section 2, we give an overview of this framework, describing the overall structure and more detailed rationale behind the computational components we build on in the remainder of this paper. We then describe a concrete but minimal implementation of this framework using contemporary probabilistic programming and language modeling tools, intended to demonstrate the basic computational components of this approach and elucidate the scope and scalability of the broader proposal.

Given this general paradigm, we first illustrate the potential breadth of this approach for integrating **meaning construction and reasoning**, showing how it might address a core set of computational and cognitive domains that we communicate about in language (Fig. 2). Each of these examples uses minimal pedagogical examples intended to suggest how this approach integrates language with important bodies of work from computational cognitive science and artificial intelligence. We first show how this framework can condition on language in order to describe and reason about *uncertain situations* with respect to an ongoing discourse (Section 2.2), then show how this approach can be extended to reason about *relational systems* (Section 3.1), *physical and perceptual scenes* (Section 3.2), and social situations involving agents with *goals and plans* (Section 3.3).

We then turn to how this approach might begin to address core scalability challenges that confront traditional approaches to modeling thinking as symbol processing, whether logical or probabilistic. In Section 4, we show how language can support growing knowledge autonomously, without hand engineering, by using the rational meaning construction framework to construct a broad range of **new concepts** in existing models and even whole **new world models**, which in turn support coherent downstream reasoning.

Ultimately, this paper is a prospective one, and the examples presented here are intended to convey a sufficiently concrete proposal to suggest avenues for future work. In Section 5, we outline what we see as some of the most significant open questions and future directions raised by this framework. These include theoretical questions that relate our approach to classical models of language, open cognitive directions for extending this approach to model language acquisition and production, and important engineering directions necessary for scaling inference, robust translation, and learning under this general paradigm. Finally, in Section 6, we conclude by looking ahead to the longer-term implications of this proposal for modeling intelligent systems that use, understand, and think about language as we do.

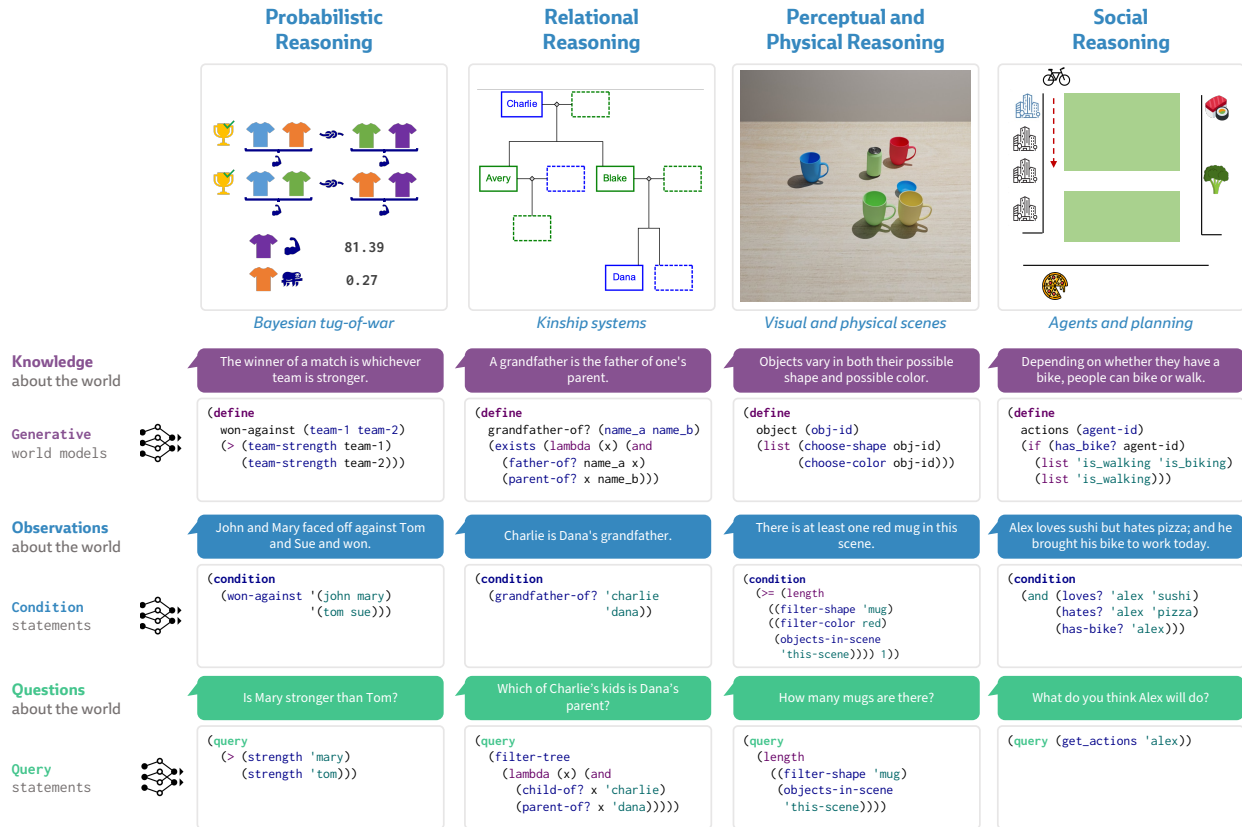


Figure 2: Understanding language in four domains of reasoning that form the core of this paper. **Probabilistic reasoning** requires integrating sparse evidence to predict the outcomes of uncertain events, like the winners of tug-of-war matches. **Relational reasoning** involves maintaining and updating coherent beliefs about structured domains, like family trees, based on relational information. **Perceptual and physical reasoning** links language to our sensory and intuitive physical knowledge of objects in the external world, such as kitchen items on a tabletop. **Social reasoning** involves reasoning about the minds of other intelligent agents, such as how their goals, preferences, and circumstances shape their actions as they navigate in the world. Across all the domains, we present a unified framework that translates language into code in a probabilistic programming language to facilitate human-like reasoning.

## 2 Overview of the key ideas

The central goal of this paper is to propose a new computational framework, *rational meaning construction*, which relates language to thought. This framework licenses a concrete class of computational architectures for building intelligent systems that use language, which we propose can be implemented using modern AI tools. In this section, we briefly overview the key ideas that form the basis of this proposal. We draw on three observations from a rational, probabilistic perspective on biological intelligence and human language:

**A rational perspective on intelligent thought.** Biological intelligence encompasses many computational capacities. The foundational notion of thought we focus on here centers on *rational inference and decision making* in service of one’s goals (Anderson, 1990; Chater & Oaksford, 1999). Under this perspective, thought comprises systems for modeling the world. These internal world models allow us to infer the particulars of a situation from whatever information is at hand, evaluate alternative world states and imagine possible future ones, and decide on actions that might bring one towards valuable future states in the world. Following extensive work in computational cognitive science, we view the world models that support biological intelligence as structured and probabilistic (Goodman et al., 2014; Griffiths, Chater, Kemp, Perfors, & Tenenbaum, 2010; Lake et al., 2017), designed to integrate the noisy evidence an agent receives into causal, explanatory models that allow them to maintain coherent beliefs about the world and generalizably infer consistent, useful predictions and plans. This basic, underlying view of intelligent thought draws on empirical evidence from essentially every species with a brain, from bees (Biernaskie, Walker, & Gegeer, 2009; R. F. Wang & Spelke, 2002), to zebrafish Bolton et al. (2019); R. E. Johnson et al. (2020), mice (English, Nejad, Sommerfelt, Yanik, & von der Behrens, 2023), birds (Isomura, Parr, & Friston, 2019), and primates (Khalvati, Kiani, & Rao, 2021). Informally, a rational view of thought can be summarized as the ability to *solve useful problems given our internal models of the world*, ranging from navigation and foraging to physical prediction and social reasoning. Against this overarching picture of thought, human intelligence further stands out for its flexibility and expressiveness. We *invent* our own problems along with new approaches to solving them, rather than sticking to a limited set of largely innate goals and strategies (Tomasello, 2022). A few other species, non-human primates, dolphins, and some birds, are creative problem-solvers and problem-creators, but none come close to the range of goals humans can adopt (Chu & Schulz, 2023). Uniquely in the natural world, humans think about and come to understand problems far beyond the narrow range necessary for our immediate survival, considering goals and questions that draw on abstract, culturally constructed, and even entirely hypothetical systems for modeling and conceptualizing the world (Dennett, 2017).

**A rational perspective on language.** As with thought, language also encompasses many systems and capacities. This paper focuses on the class of problems we refer to as *language-informed thinking*, the general means by which language informs the inferences and decisions of an intelligent agent. We take a broadly rational perspective on language—we consider language to be a system of *goal-directed actions for externalizing and communicating thoughts* to other intelligent beings (Chater & Manning, 2006; Gibson, 2014; Goodman & Frank, 2016). In this context, we frame the problem of deriving *meaning* as inferring the mappings between a language’s system of external communicative signals into the representations of rational thought. It is worth highlighting that thought does not require language and is distinct from language in the human brain (Mahowald et al., 2023). Non-human species, and pre-verbal infants (Spelke, 2022), are clearly capable of modeling the world towards their inferences and goals without language. But for humans, language clearly plays a profound role in determining the problems we think about, and how we think about them. Our natural languages allow us to communicate an extraordinarily broad range of our thoughts about the problems we pose and solve, including our abstract and general world knowledge, our specific beliefs about a situation, the particular questions or goals we have or want to pose to others, and our approaches to reasoning about them.

**A resource-rational perspective on language and thought.** Finally, our integrated computational approach to language and thought builds on extensive evidence that humans are *resource-rational* thinkers—under finite constraints of time and memory, we rationally allocate computational resources in order to make useful inferences and plans (S. J. Gershman, Horvitz, & Tenenbaum, 2015; Lieder & Griffiths, 2019). Resource rational agents *amortize* computational effort across prior experience and problems, storing and

reusing prior computation towards similar new problems that we encounter in the future (S. Gershman & Goodman, 2014; Le, Baydin, & Wood, 2017). Certain domains of inferences share more structure than others, and evidence suggests that we therefore heavily amortize them. Prior work, for instance, suggests that computations involved in basic perceptual activities (Brooke-Wilson, 2023; Dasgupta & Gershman, 2021; Fodor, 1983), such as object recognition under common lighting conditions, are highly amortizable from reusable patterns in computation that are learnable and shared across a background distribution of perceptual instances. This view suggests why fast, bottom-up pattern recognition models have made great advances in modeling perception in recent years, while it has proved much more challenging to amortize the wide range of flexible inferences required for arbitrary problem solving.

We propose an analogous resource-rational perspective on the kinds of computation implicated in language-informed thought. Under almost every theoretical and empirical account of linguistic structure and semantics, the mappings between language and meanings should be highly amortizable across the background distribution of language—there are structured, systematic, and learnable patterns in how units of language map onto units of thought. The idea that meaning construction should be highly amortizable follows from our view on language itself as an efficient communicative system. Extensive empirical evidence suggests that communicative pressures shape how language maps onto meanings at every level of linguistic structure, from individual morphemes (Bybee, 1985) to patterns in how common syntactic frames communicate meaning (L. Gleitman, 1990; Grimshaw, 1981), and even reusable pragmatic implications present across common discourse situations (White, Mu, & Goodman, 2020). But while we take the view that a resource-rational agent should intelligently learn and reuse prior computation when possible, we do not view language-informed thinking, or thinking in general, as solely a matter of learning and interpolating over statistical patterns from prior experience. When we think, including when we think *about* the meanings we recover from language—to update our beliefs, to follow instructions, or to answer questions posed in language—we must be able to flexibly model arbitrary situations and support capacities for general problem solving, including inference, planning, and simulation, under a wide range of new and unencountered circumstances.

The efficient learnability of human language also highlights that, in many senses, the computational relationship between language and thought in humans is almost the inverse of that in today’s LLMs. For humans, language could be characterized as an emergent property of *thinking*. Infants can model the world and draw inferences well before they know language (Gopnik, 1996; Spelke, 2022), and reliably acquire complete linguistic capabilities from exposure to relatively tiny amounts of language (R. Brown, 1973). Congenitally-Deaf humans born with *no* language input spontaneously develop languages to communicate their thoughts, with the same basic hallmarks of mature natural languages (Goldin-Meadow, 2012; Pyers, Shusterman, Senghas, Spelke, & Emmorey, 2010; Senghas, Kita, & Ozyurek, 2004). This paper seeks to understand and model the cognitive and computational structures underlying *this* human scaling route to intelligence and language use — one that begins with robust capacities for thought, and scaffolds language efficiently on top of them to then offer a powerful tool for driving and constructing new thought.

## 2.1 Our proposal: A framework for modeling rational meaning construction

The perspective we offer above draws from theoretical and empirical work that precedes this paper. Our core contribution in this paper is to propose a new computational framework in light of these observations, that seeks to unify prior symbolic, probabilistic inference and statistical learning traditions and to take advantage of the clear computational advances made by modern LLMs as learned statistical models of language. We describe a framework for **rational meaning construction** in which linguistic meaning is formalized as a context-sensitive mapping from natural language to a distribution over expressions in a probabilistic language of thought (PLoT) for rational world modeling and inference. Under this framework, we then propose that large language models trained on language and code can be *used to implement meaning functions in a resource-rational architecture* – they can implement learned, broad-coverage mappings between language and code; and they can be understood as part of a human-like, resource-rational system that efficiently infers these mappings using stored patterns amortized from the prior joint distribution over language and code. This motivates the concrete architecture we propose and illustrate throughout the remainder of this paper, and its two main components for modeling thinking and modeling language relative to thinking—or how language informs thinking.



### 2.1.1 Modeling thinking

We propose implementing **thinking using probabilistic programs** as a general representational substrate for building world models and specifying rational inferences over them. This proposal builds on prior work in cognitive science and AI formalizing how a broad class of problems can be expressed as probabilistic programs (Chater & Manning, 2006; Goodman et al., 2014), following a generic *inference query* motif (Goodman, Mansinghka, Roy, Bonawitz, & Tenenbaum, 2008) — a probabilistic program that combines a *generative world model* that models abstract, causal beliefs about probable world states; specific evidence that an agent *conditions* on; and a particular query being posed as the question or goal for thinking. Inference to solve a problem consists of formally computing or sampling from a probability distribution over answers to this question, specified by the world model and conditions. This computational proposal forms the backbone of the *probabilistic language of thought* model of general human cognition (Goodman et al., 2014), and has been used empirically to model a wide range of human inferences, including those that draw on visual perception (V. K. Mansinghka, Kulkarni, Perov, & Tenenbaum, 2013), physical simulation (Battaglia et al., 2013), and social reasoning (C. Baker et al., 2011). It is designed explicitly to formalize a central property of human thought — the capacity to expressively and flexibly pose problems involving entirely novel situations and goals, and to solve them relative to a computable representation of the world and internal belief.

### 2.1.2 Modeling language relative to thought

Given this model for thought, we propose formalizing **rational meaning construction** as a broad-coverage, contextual translation function that maps language into a distribution over expressions in a probabilistic language of thought. This proposal builds most closely on and draws inspiration from efforts to articulate a *probable world semantics* for natural language in prior work (Goodman & Lassiter, 2015), in order to express how language could compactly convey uncertain propositions and vague meanings with respect to a formal probabilistic generative model. It also builds on the longer history of symbolic semantic theories we overview in the introduction, including formal semantics theories that model language as mapping into formal propositions over possible worlds (eg. Heim and Kratzer (1998); Lewis (1976)), and semantic parsing systems (eg. Abend et al. (2017); Klein and Manning (2003); Liang (2016); Steedman (2001); Y. W. Wong and Mooney (2007)) that map language into formally executable program expressions.

Our goal, however, is to broaden and generalize these framings to suggest a general framework for modeling how language can interface with and inform such a broad swath of human cognition. By positing that meaning is a general mapping between sentences and expressions in a probabilistic language of thought, we believe that a rational meaning construction approach can elaborate on and concretely model core desiderata of a coherent theory of linguistic meaning – modeling how meanings drive inferences about what is true and probable; formalizing how language can pose propositions and queries that are then evaluated with respect to an internal model over probable worlds; and relating meaning to the general computational systems for representing, thinking about, and receiving new information about the world from broader cognition.

This proposal suggests a wide class of possible architectures that map from language into probabilistic programs—in principle, any general mapping function that expresses a distribution over programs conditioned on sentences in context. Under this umbrella of possible implementations, we propose finally that **large language-to-code models** can be used to generally instantiate these meaning functions. Unlike prior semantic parsers or attempts to hand implement mappings between language and code, LLMs offer a concrete means of instantiating far more broad-coverage mappings between human sentences and meanings than have been previously possible. They are also context-sensitive, in that they can construct meanings for an utterance that condition both on the general distribution of language and thought and a local linguistic and thinking context. They can condition translation on a local discourse context, when prompted with prior utterances, and on a local problem under consideration, when prompted with existing code in a probabilistic program.

By using LLMs to map between language and code, this proposal is also closely related to the recent lines of work we review in the introduction that seek to augment and connect LLMs with various structured and symbolic reasoning tools—both domain-specific reasoning engines like planners and physics engines (eg. B. Liu et al. (2023); R. Liu et al. (2022)), and more general APIs for code execution (eg. Karpas et al. (2022); OpenAI (2023b); Schick et al. (2023)). As we demonstrate throughout this paper, however, we propose that the probabilistic language of thought can offer a cognitively-motivated, unifying symbolic substrate for

interfacing between language and many core aspects associated with general cognition. It provides a general motif for structuring and constructing generative world models, which can nest calls to other domain-specific systems (such as planners and physics engines); and an overarching framework for modeling how diverse kinds of observations can be used to update these models and answer new queries, framed as Bayesian conditioning and inference. With respect to the more general landscape of large statistical language models, this proposal finally suggests one way to situate the strengths of LLMs into a more human-like, modular framework for language-informed thinking. Rather than look to statistical patterns to capture all of the ways we think, plan, and reason about language, this resource-rational approach seeks to ground distributional aspects of language into a framework that can leverage learned prior patterns when they are useful—while also modeling how language can construct and relate to coherent world models and algorithms for explicit, novel decision making and inference.

### 2.1.3 Illustrating the architecture by example

This general architecture is best explained through concrete implemented examples, which we give in the next sections. For each of the four domains of reasoning shown in [Fig. 2](#), we work through a representative dialog between a speaker of English and our language-informed thinking computational architecture, which could stand in for how we model another human being’s understanding and thinking about the speaker’s language, or the ways we hope a human-like AI system would similarly respond.

For pedagogical reasons, we have chosen to implement these examples using one particular probabilistic programming language and one particular language-to-code model. These particular tools are not necessarily the most performant or scalable AI solutions; nor the best accounts we have of the corresponding components of human architecture. Nevertheless, they are familiar and simple, and provide the most direct route we know to illustrate our ideas in ways others can also experiment with. To elaborate on these choices:


- The probabilistic language of thought we use to express inference problems is *Church* ([Goodman et al., 2008](#)), a Turing-universal probabilistic programming language constructed on top of the functional programming language Scheme. We have used the WebChurch dialect which implements several general inference procedures, but we have chosen the simplest and most general—and least efficient—approach based on rejection sampling: Inference is based on drawing samples from the prior over world states described by the generative model, and rejecting those that fail satisfy the constraints of any observation conditions. The samples that remain constitute a posterior sample over possible worlds consistent with the observed information, sufficient to answer the queries under consideration in the language discourse. Other similarly functional PPLs such as WebPPL or Gen could have been chosen instead. In [Section 5](#), we discuss future directions for extending and scaling inference beyond these simple illustrative implementations.
- The language-to-code model we use to amortize meaning construction over programs is Codex model ([Chen et al., 2021](#)), a GPT-3-based language model fine-tuned on source code, which provides pairings between natural language and code with comments, drawn from programs on GitHub and other sources. Since the release of Codex, many other language-to-code models have been developed, and more recent versions of GPT-based language models are now routinely trained on large amounts of source code; we believe these could be used to similar effect. In [Section 5](#), we also discuss future directions for more cognitively plausible training and updating of neural models that amortize inference in joint distributions over natural language and probabilistic languages of thought.

Finally, before turning to the examples, we want to add an important note about our intentions and goals. The examples are designed to be illustrative and pedagogical—we choose them for their simplicity and clarity, and to show how prior empirical and computational work from cognitive science can be related under this general framework to language. Each example gestures at a larger domain of reasoning, but, of course, each domain is much broader than what we can implement here. Each example is also representative of a wide class of computational cognitive models that can be instantiated in a probabilistic language of thought, and which we propose can be integrated with natural language inputs and outputs under a rational meaning construction framework. In each section we therefore also discuss how this framework might be scaled, and what more work may be necessary, to scale from these examples towards a richer model of language in relation to those domains.

We also hope that these examples, and other variations that elaborate on them and on the core domains of reasoning we discuss here, will offer useful starting points for more rigorous, systematic, cognitively-oriented evaluation and interpretation of the reasoning processes emergent in large language models and other language-based AI systems. In our own preliminary evaluations of these domains, we find that current large language models show many of the properties we discuss in the introduction. In some cases they appear to approximate implicitly the representations and algorithms we seek to model explicitly. In others, particularly with more complex modifications beyond these simple examples, we find that large language models left to their own devices produce outputs that diverge from our intuitions. We seek here to model the representations with which *people* make meaning from language in relation to all of these domains, but hope that these frameworks will be useful for understanding other computational systems that use language as well, including interpreting the representations that large language models already learn or should seek to acquire.

### Graphical conventions

*Throughout the examples presented in this paper:*

 **Translations** mapping from language into probabilistic programs, produced by Codex, are indicated by a neural network icon.

 **Probabilistic inferences**, performed by Church, are indicated by a cog icon.

## 2.2 Understanding language with probabilistic reasoning

To illustrate our framework, let’s consider a concrete scenario that involves reasoning from language in the face of uncertainty. Suppose a friend is telling you about a tug-of-war tournament that took place the prior weekend in which the authors were participating:

*Right off the bat, Josh won against Lio. He then proceeded to claim victory against Alex. Even working as a team, Lio and Alex still could not beat Josh!*

In order to understand this story, it is useful to construct a little mental model: there are different players, they face each other solo or in teams, and based on his track record, Josh appears to be particularly strong. Now, suppose your friend tells you about a newcomer: *In a huge upset, Gabe managed to best Josh in the fourth round.* Maybe Gabe is even stronger than Josh! Or, perhaps Josh was simply feeling lazy in the last match, in which case, Gabe might not actually be so strong. To clarify, you might ask a question, *Who is stronger: Gabe or Josh?* Your friend’s answer, which might itself express uncertainty, will nevertheless provide further information for you to incorporate into your understanding.

In making meaning from language about a scenario like the above, you are engaging in *probabilistic reasoning*: integrating over different possibilities in order to infer likely explanations. People are remarkably proficient at making inferences from exactly this kind of sparse evidence. Sometimes, we acquire this evidence through direct experience—by watching the tournament, for instance—but often, this kind of information comes to us through language that cues us to update our beliefs accordingly. Critically, in order to reason *consistently*, we need to represent core aspects of the situation: who are the different actors, what events took place, and what inferences have we already made? To this end, it is extremely useful to have a *world model*, which we defined earlier as a probabilistic generative model that encapsulates the key mechanics of a domain and facilitates coherent, causal explanations of events. In this section, our aim is to further formalize what exactly we mean by world models and how large-scale neural models might serve as an interface between natural language and these kinds of cognitive representations.

**World models as generative programs.** The core of each example in this paper is a probabilistic generative model that defines the mechanics of a domain. For the purposes of this demonstration, and throughout [Section 3](#), we focus on reasoning from language given a pre-specified world model. Later, in [Section 4](#), we show how language can be used to grow out and construct new world models.

As a playground for this initial demonstration, we consider the “Bayesian tug-of-war,” a classic experimental domain in cognitive science that requires making inferences about the latent traits of individuals from sparse evidence. Prior work establishes that Bayesian inference in a probabilistic generative model closely captures people’s predictions about scenarios in the tug-of-war ([Gerstenberg & Goodman, 2012](#); [Goodman et al., 2014](#)), and that simple sentences can be mapped onto queries in this model ([Goodman & Lassiter, 2015](#)). Here, we build on this work to give an account for how people might turn open-ended natural language into statements in the probabilistic language-of-thought.

In tug-of-war, we start with a generative model of a tournament in which players of varying strengths compete in a series of matches, facing off either solo or as part of fluid teams ([Fig. 3A](#)). Each player has a latent strength value randomly sampled from a Gaussian distribution (with parameters arbitrarily chosen as  $\mu = 50$  and  $\sigma = 20$ ). As an observer, our goal is to infer the latent strength of each individual based on their win/loss record. However, players sometimes don’t pull at their full strength and each player has a different intrinsic “laziness” value (uniformly sampled from the interval  $[0, 1]$ ) that describes how likely they are to be lethargic in a given match. The full Church code for the tug-of-war is given in [Appendix A.1.1](#).

**Linguistic meanings as probabilistic program expressions.** While the generative model defines the generic mechanics of the domain, we want to be able to talk about *specific* people and events. In our framework, we focus on two kinds of linguistic utterances:

**Observations** provide information about people, objects, and events in the world; e.g., “Josh faced off against Lio and won.” In our framework, we translate observations into condition statements in Church, which update the state of the world model to reflect new facts. Note that **condition** statements have no return value; instead, they constrain the world model such that downstream inferences must be consistent with respect to the conditioning statement.

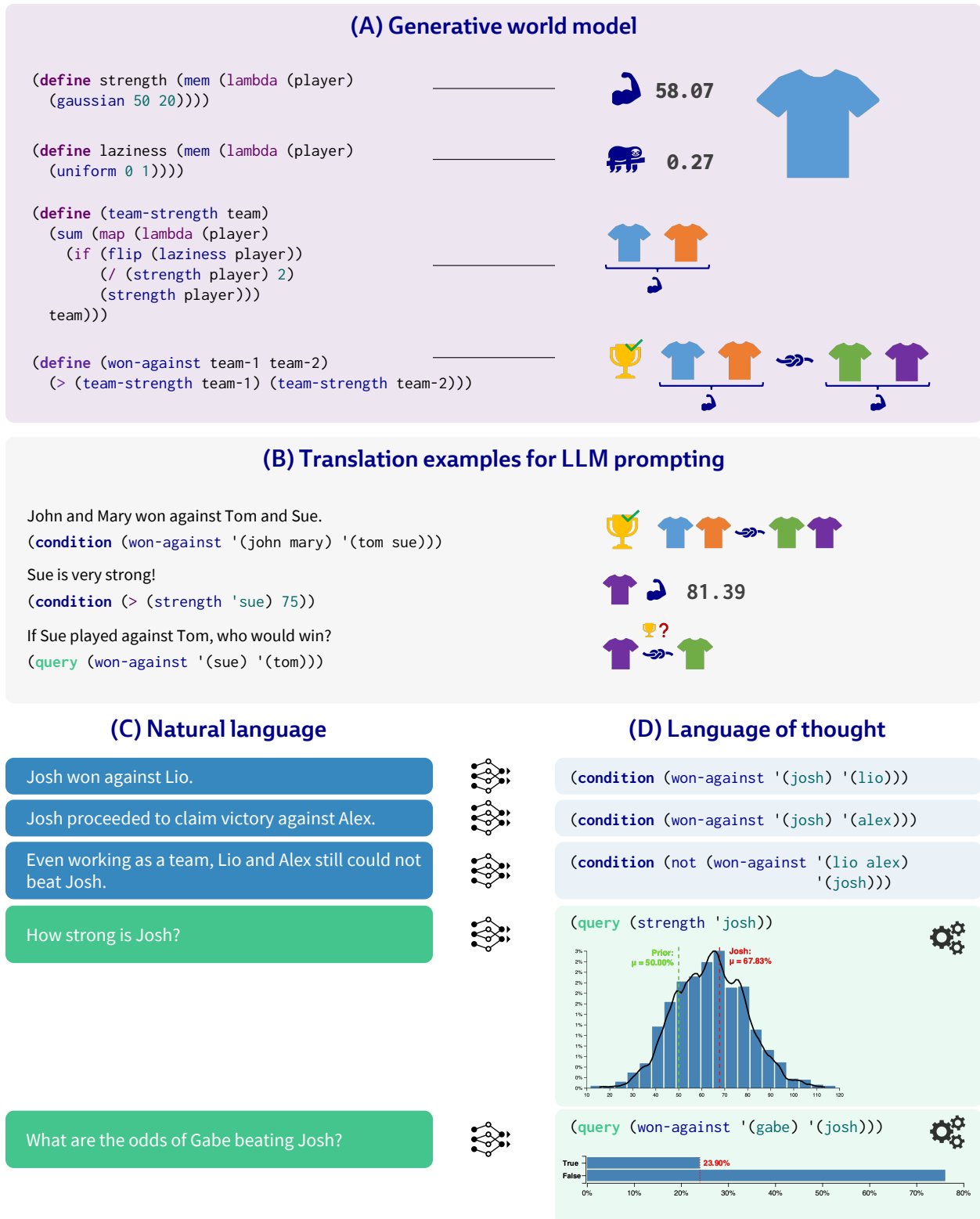


Figure 3: Illustration of probabilistic reasoning via language-to-code translation in the tug-of-war domain. (A) The generative model defines two latent traits, strength and laziness, and specifies how these interact to determine team-strength. By combining (A) and (B), we can few-shot prompt an LLM to translate open-ended natural language (C) into Church statements (D) that capture linguistic meaning with respect to the domain. The resulting probabilistic inferences transparently represent the model’s beliefs and naturally capture human-like intuitions about players’ latent traits.

**Questions** seek information in the face of uncertainty about the world; e.g., “Would Josh beat Gabe if they played again?” In our framework, we translate questions into **query** statements in Church, which evaluate the quantity of interest. Calling **query** triggers a probabilistic computation that simulates possible worlds under the model, constrained by any observations so far. The query expression is evaluated in each simulated world, yielding multiple samples that form a posterior distribution over the value of interest.

Throughout the examples in this work, we freely interleave **query** and **condition** statements, much as questions might occasionally arise between statements of fact in a natural dialogue. Implementationally, this behavior is achieved through a read-evaluate-print loop (REPL) inspired by Venture’s (V. Mansinghka, Selsam, & Perov, 2014), that evaluates queries against all **condition** statements that have appeared up to that point in the dialogue history. In our model, we assume that the user specifies whether each utterance is a **condition** or a **query**, but LLMs could likely classify unannotated utterances accurately.

**Translating from natural language to program expressions.** Inspired by the work of Goodman and Lassiter (2015), if we had some way to translate linguistic utterances into probabilistic program statements, we could perform a wide variety of probabilistic inferences from plain English. Up until recently, however, it was unclear how to construct a *meaning function* sufficiently general to translate open-ended natural language into highly structured expressions compatible with a Church model. Our core observation is that language-code LLMs have many of the properties necessary to serve as a useful meaning function: broad-coverage exposure to natural language, a robust capacity to model joint language-code text distributions, and the ability to quickly grasp domain-specific syntax and semantics from a few examples.

In this work, we leverage the few-shot prompting capabilities of one such LLM, the Codex model from OpenAI, to induce a translation model from English to Church code. As it turns out, we only need to provide a small handful of *example translations* (represented in Fig. 3B) to achieve a variety of interesting behaviors. To translate a new language utterance to Church, we simply concatenate the generative model (full text in Appendix A.1.1) and the translation examples (full text in Appendix A.1.2) into a prompt whose final line is the utterance. We then generate from Codex, which, based on the comment-code pattern in the prompt, infers that the completion should be written in Church, using the function definitions and constructs provided in the prompt.

Notice the high degree of variation in phrasing and lexical choice in Fig. 3C; none of the utterances contain “won” or “against,” yet Codex still maps these to the *won-against* function. Here, we start to see some of the advantages of using an LLM over more traditional semantic parsing techniques like CCG parsers (Artzi, Lee, & Zettlemoyer, 2015; Artzi & Zettlemoyer, 2013). Because the model is pre-trained on a vast amount of linguistic data, it fluently handles many different kinds of linguistic variation. However, by including the Church generative model in the prompt, we can effectively constrain the output space; the model infers that the generated code should use the functions defined in the generative model.

As a semantic parsing tool, this combination of pre-training and prompting manages to achieve broad invariance to spurious linguistic variation while remaining sensitive to wording choices that might affect meaning. We can see this tradeoff at work in Fig. 3C, where the translation uses a negation, closely reflecting the structure of “Lio and Alex still could not beat Josh.” Of course, there are multiple aspects of the utterance that this translation does *not* capture (e.g., “Even working as a team...” suggests that Lio and Alex’s efforts were well-coordinated; as opposed to something like, “Stepping on each other’s toes the whole match...,” which would imply the opposite). Our point is not that the LLM translation perfectly captures all aspects of the utterance meaning, but rather, that it encodes those that are relevant to and compatible with the domain model so as to facilitate downstream reasoning.

**Reasoning about scenarios with probabilistic inference.** So far, we’ve illustrated how we might **condition** a PLoT model on natural language, but what about reasoning? After hearing the information in Fig. 3C, we might assume that the player named Josh is quite strong. Exactly how strong is Josh, though? And how likely is it that he would beat another player who isn’t Lio or Alex? Just as we used Codex to translate facts into **condition** statements, we can use it to translate questions into **query** statements in Church. The Church inference engine then automatically simulates scenarios (in this case, 1000 times) that are consistent with the given **condition** statements in order to produce an approximate posterior distribution over each query.

By offloading reasoning from the LLM to the PLoT, we can obtain a much richer picture of our model’s beliefs about the world (Fig. 3D). While the LLM alone can only respond with textual statements like “Josh is very strong,” Church gives us an entire probability density over Josh’s strength (on expectation, he is a little less than one standard deviation above the average strength = 50). Likewise, we can easily obtain a distribution over the outcomes of a Gabe-Josh match (given Josh’s strong track record, our model finds Gabe’s chances slim, at 23.90%). Critically, Church is doing much of the heavy lifting of inference in the background in order to produce these posterior distributions.

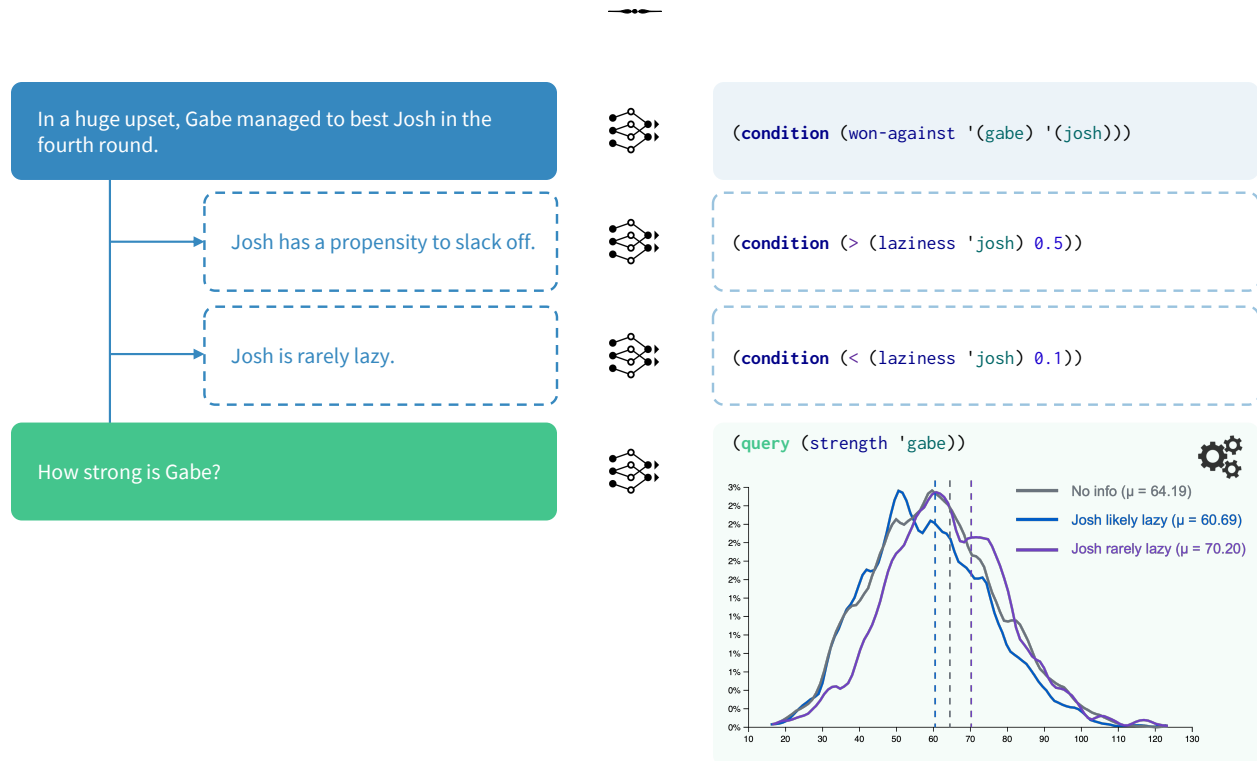


Figure 4: Reasoning about a pair of hypothetical scenarios with language-code translation. In a world where Josh is often lazy, Gabe’s win is counteracted by a high likelihood that Josh threw the match. Conversely, in a world where Josh is rarely lazy, Gabe’s win is surprising and suggests a high strength value. Rational meaning construction with an LLM appropriately resolves the linguistic meaning of these two scenarios, selecting reasonable probability parameters for the conditioning statements. Meanwhile, probabilistic inference about Gabe’s strength is finely sensitive to the implications of these competing hypotheses.

In addition to providing useful interpretability, reasoning in Church models is sensitive to each new piece of information. Much like human learners, Church models can flexibly update their beliefs when presented with low-probability or unanticipated events. Picking up our tug-of-war saga, consider the plot twist in Fig. 4:

*In a huge upset, Gabe managed to best Josh in the fourth round.*

How might this new information shape our interpretation of the match 4 outcome? If Josh is likely to be lazy, then it’s possible that Gabe simply got lucky and wasn’t so strong after all. If, on the other hand, Josh is rarely lazy, we might start to regard Gabe as particularly strong. In Fig. 4, we can observe how Church reasons about these two possibilities, shifting the probability density over Gabe’s strength left if Josh is likely lazy and right if Josh is rarely lazy.

Note how, in order to translate a phrase like “Josh has a propensity to slack off,” Codex must choose a particular probability threshold. This choice is arbitrary and, while there is no “correct” answer, we see that Codex is able to choose valid probability values between  $[0, 1]$  that feel appropriate to the wording: a

“propensity to slack off” doesn’t necessarily imply that someone slacks off *all the time*, while, in contrast, “rarely lazy” offers more certainty. Indeed, across many different contexts, we observe that Codex is able to pick reasonable parameter values that respect both the language and the parametrization of defined distributions. We consider these inferences to represent a form of “amortized pragmatics” (Goodman & Lassiter, 2015), which we will revisit in Section 5.

**Putting it together: the power of probabilistic reasoning.** We conclude this section with a final example that underscores the flexibility of our framework to model complex reasoning from language and foreshadows multiple themes that we will revisit later in the paper. Consider the dialogue in Fig. 5, in which the students and faculty team up to face one another. The interlocutor poses two questions: “Is Gabe stronger than the weakest player on the faculty team?” and “Who would win in a match between the students and the faculty?” As we saw in the prior tug-of-war examples, the answers to both of these questions are expressed as probability distributions derived from simulation of the generative tug-of-war model. Moreover, in both cases, the introduction of new information flips the model’s belief state in a way that aligns with human intuitions. In this way, the PLoT framework is natively capable of *defeasible inference*—a phenomenon of human reasoning that was of great interest to early AI pioneers of non-monotonic logics (Ginsberg, 1987; McCarthy, 1980).

A key advantage of our framework is that achieving these kinds of defeasible and flexible inferences from natural language reduces to grounding utterances into appropriate **condition** and **query** statements. While the observations and questions in Fig. 5 are semantically more complex than those that appeared in the prior examples, and though there are many degrees of freedom involved in the translation problem, we confirm that an appropriately-prompted LLM can produce translations that intuitively capture the meaning of each utterance with respect to the tug-of-war domain. Moreover, as we saw in Fig. 4, Codex is able to *amortize* certain pragmatic inferences in resolving “pretty strong” to a threshold of strength  $> 60$ , “real slackers” to a threshold of laziness  $> 0.9$ , and “several of the faculty” to count  $\geq 3$ . How far can we go with these kinds of amortizations? Throughout Section 3 and Section 4, we will see examples of context-driven amortizations across different domains; and in Section 5, we will regroup to discuss how these different examples of amortization might inform our theories of language understanding and pragmatics.

In this dialogue, we also give a preview of **define**, a powerful construct in our framework that is discussed in depth in Section 4. Just as people come up with terms like “20th-century pragmatists” or “Meatless Monday” to pick out entire hierarchies of people, things, and events, a core feature of the probabilistic LoT is the ability to **define** new concepts that can later be referenced symbolically. In the Fig. 5 dialogue, language about team structures defines two new concepts, `faculty-team` and `student-team`, that facilitate concise translation of language like, “Is Gabe stronger than the weakest player on the faculty team?” Moreover, while `faculty-team` is a static list, other defined concepts can ground out in *functions* that take arguments. In fact, `stronger-than?`, which is defined in the prompt (Appendix A.1.2), is one such example, illustrating how programming languages are well-suited to capture the infinite productivity of language that arises through structured composition. Through this lens, we can start to imagine how our tug-of-war world model might be expanded to ground many new kinds of language:

- The tug-of-war tournament is organized into three leagues for novices, amateurs, and professionals. In order to be considered a professional, a player must win 20 one-on-one matches against other professionals.
- Players often get increasingly tired over the course of a tournament, though some players have more stamina than others.
- The tournament has an entry fee of \$20 per contestant and a grand prize of \$10,000 for the winning team.

How can we grow our world models to incorporate new language, or even construct new world models entirely from scratch? In Section 4, we revisit the tug-of-war domain with an eye to precisely these questions.

**Conclusions.** As an introductory example, the tug-of-war domain serves as a minimal illustration of the kind of reasoning from language that our framework is concerned with. Our goal here was to build intuition



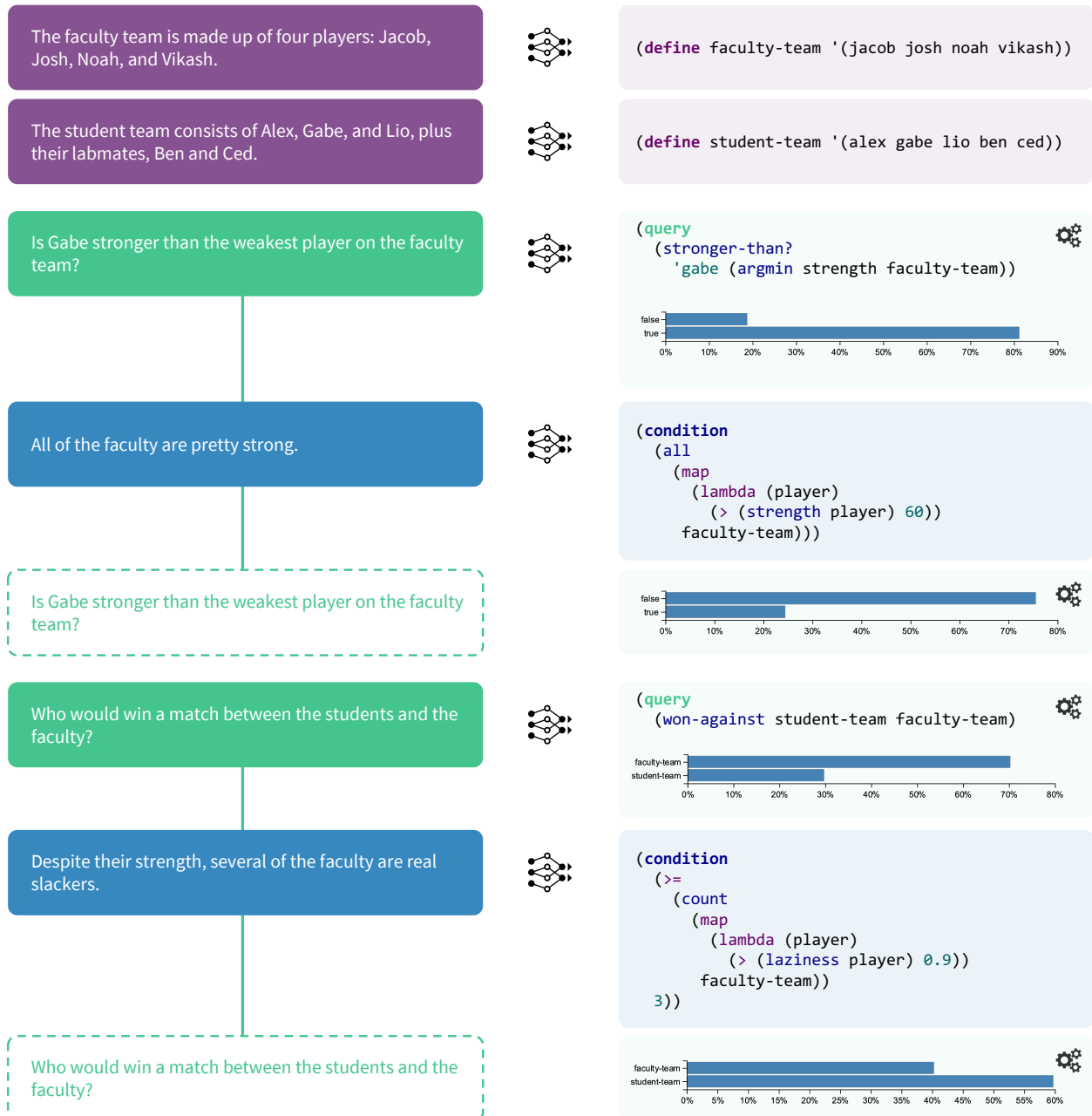


Figure 5: In this final tug-of-war dialogue, natural language plays three interlaced roles in interfacing with the language-of-thought. Definitions (purple) introduce new concepts, such as specific player-teams, that can later be referenced symbolically. Observations (blue) translate into condition statements that probabilistically constrain the world state, sometimes amortizing the resolution of linguistic ambiguity (e.g., “pretty strong” or “real slackers”). Finally, questions (green) translate into queries that trigger inference by probabilistic simulation over possible worlds that is both sensitive to and consistent with prior definitions and observations.

for our general approach: by translating natural language into **condition** and **query** statements as inputs to a probabilistic inference engine, we can achieve forms of reasoning from language that are consistent with respect to a mental model of the world. Nonetheless, in scaling this approach beyond the toy domain of tug-of-war, many questions arise. How does probabilistic inference relate to models of relational and deductive reasoning of the sort that classical AI approaches excel at? How do we ground linguistic meaning in the visual and physical world? And how does language understanding inform our actions and interactions with other agents through goal-directed planning? In **Section 3**, we will progressively expand our scope to touch on each of these questions and show that, in each case, new kinds of language understanding and reasoning can be naturally incorporated into our framework.

### 3 Understanding and reasoning about language with world models

In this section, we illustrate how the general framework we propose in [Section 2](#) can be applied and extended to integrate natural language with core domains of human-like thought. In each, we build on the idea that language that conveys observations and questions about uncertain situations, constructing meanings from a generative world modeling program that supports *probabilistic reasoning*. In [Section 3.1](#), we show how this approach can be integrated to understand language that conveys structured, logical lexical relations. In [Section 3.2](#), we show how generative programs that support perceptual and physical simulation can be used to ground language about scenes into visual world. Finally, in [Section 3.3](#), we consider language about agents with preferences and goals, and show how we can make meaning from sentences with respect to a generative program that supports planning.

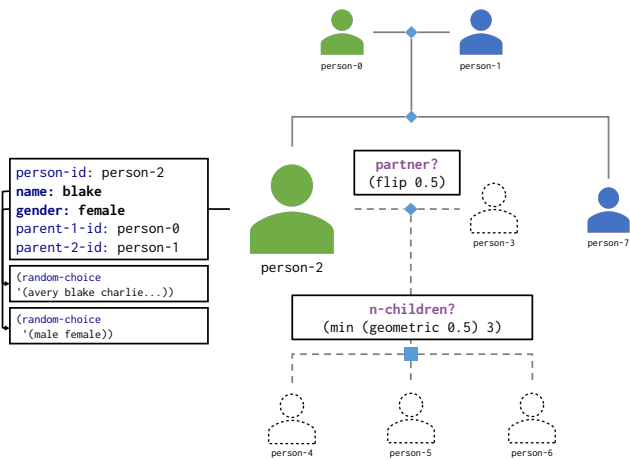
#### 3.1 Language for logical and relational reasoning

In the previous section, we examined how translation from natural language into the probabilistic language of thought naturally captures a certain form of reasoning in which uncertainty plays a key role. How does this framework relate to earlier computational theories of reasoning, such as classical AI approaches to logical and relational reasoning ([Russell & Norvig, 2021](#))? Historically, systems like Prolog ([Colmerauer, Kanoui, Pasero, & Roussel, 1972](#); [Philippe, 1972](#)) were designed for similar goals to ours here, to allow people to directly interact with computers via natural language (French, originally), specifying only the background knowledge and goals for computation without the algorithmic details ([Colmerauer & Roussel, 1996](#)). In this section, we demonstrate how the PLoT not only fully supports the style of deductive, logical reasoning characteristic of classical AI, but extends it to support *inductive* inferences as well. Moreover, we argue that many kinds of real-world reasoning problems that are traditionally modeled using structured logic-based approaches actually require a mix of both symbolic and probabilistic reasoning. In doing so, we aim to illustrate how our approach of translating from natural language to the PLoT fluidly integrates both kinds of reasoning in a way that comes naturally to people, but that has proved elusive for both traditional deductive programming systems and purely statistical language models.

**Language about kinship relations.** Suppose you are again with your friend from [Section 2.2](#), who is telling you about a part of their extended family. “Avery has a sister named Blake, and their father is named Charlie,” your friend says. Immediately, you start to sketch a picture in your mind of this family, which you can update on-the-fly as you get more information: “Charlie is the grandfather of Dana.” At this point, you can infer that one of Charlie’s kids is also Dana’s parent, but which one? In the absence of additional information, it’s a toss-up between Avery and Blake, with some outside chance that there could be another, unmentioned sibling who is Dana’s parent. Hearing that “Blake has two kids” might initially shift your beliefs towards Blake. However, upon learning that “Dana is an only child,” you’d have to rule Blake out entirely! This kind of relational reasoning, which freely intermixes deductive and inductive inferences, comes quite naturally to people. How do we make such rich inferences from a relatively sparse sequence of words?

In this section, our domain of interest will be **kinship**: relationships between people in a family. The kinship domain provides fertile ground for the study of logical reasoning for several reasons. First, during development, one of the first domains where we learn about logical relations is in describing families ([Elkind, 1962](#); [Piaget, 1951](#)). Language has evolved to describe family structures in highly economical terms that naturally express composition (e.g., my mother’s father is my *grandfather*) and symmetry (e.g., if Avery is my cousin, then I am Avery’s cousin; together, we are *cousins*). Nevertheless, while certain kinship references are relatively straightforward (e.g., “Blake’s mother”), others involve ambiguity (e.g., “Blake’s uncle” could refer to the brother of *either* of Blake’s parents; or even, perhaps, a close older male who is not related by blood or marriage). Finally, kinship reasoning freely intermixes deductive and inductive inferences: for instance, “Charlie has a grandson named Dana” *deductively* implies the existence of a child of Charlie who is also a parent to Dana; and it *inductively* implies that Charlie was possibly partnered at some point, such that Dana might have another grandparent in the picture. Traditional logical accounts of reasoning in this domain capture the deductive inferences but not the inductive inferences in cases like this. People, in contrast, routinely make statements such as “This is Kendall, the partner of Avery’s niece” with the expectation that others will draw roughly the same inferences they would in building a mental model of this family: Avery has

## (i) Generative domain theory of family trees



## (ii) Translations into LoT predicates

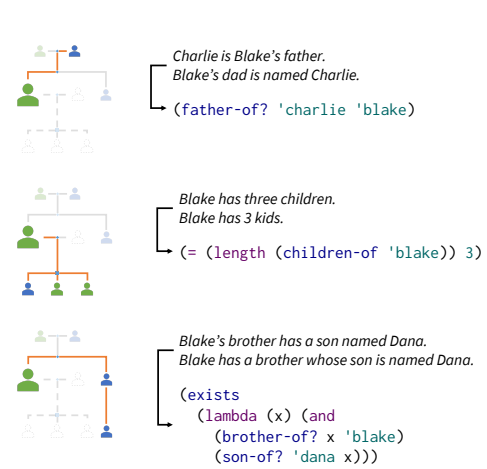


Figure 6: Illustration of a simple kinship domain theory and conceptual system implemented in Church. (i) The generative model specifies a process by which individuals form couples and have children to form family trees. Each tree represents a “possible world” in which certain relationships hold. (ii) These relationships are expressed using predicates in a conceptual system that supports quantificational logic and composition, giving rise to an expressive domain semantics that aligns well with natural language.

a brother or sister, and that sibling has a female child, and Kendall is that person’s partner. In sum, the kinship domain offers a rich set of relations and possible inferences, and comes equipped with an extensive natural language vocabulary, making it an ideal playground to explore our translation hypothesis.

**World models of kinship as probabilistic generative programs.** Owing to the richness of the domain, recent years have seen a steady interest in computational cognitive models of various aspects of kinship, ranging from development and acquisition of kinship terms across cultures (Mitchell & Jordan, 2021; Mollica & Piantadosi, 2022), tradeoffs in communicative efficiency in natural (Jones, 2010; Kemp & Regier, 2012) and artificial (K. Smith, Frank, Rolando, Kirby, & Loy, 2020) kinship systems, and probabilistic inferences about kinship relations from sparse evidence (Katz, Goodman, Kersting, Kemp, & Tenenbaum, 2008). In this work, our primary interest is in how people represent and reason about kinship relations conditioned on language. Following Katz et al. (2008), we construct an intuitive domain theory of kinship using a probabilistic generative model and a small number of rules that form a conceptual system.

As in Section 2.2, our kinship domain theory is expressed as a generative model in Church. In the Bayesian tug-of-war, the generative model consisted of random variables over continuous quantities like *strength* and *laziness*. In contrast, in this section, our generative model specifies a series of discrete random choices that describe events in a family’s genealogy: people are born, find partners, have children, and the process repeats. All of these events involve random choices that shape the makeup of the family tree.

Fig. 6 (i) shows a schematic of the kinship generative domain theory. When a person is born, they are assigned a unique *person-id*, a name<sup>1</sup> sampled from a list of gender-neutral names, and a gender sampled from {male, female}. Next, with fixed  $p = 0.5$ , the person partners with a new individual from outside the family. Finally, if partnered, the couple has  $n = \{0, 1, 2, 3\}$  children, with the number of kids drawn from a geometric distribution ( $p = 0.5$ ). This process repeats recursively until a full family tree is generated. To support efficient inference using Church’s generic sampling algorithms, we cap the trees at 3 generations and limit each couple to 3 children. Further implementation details of the generative model can be found in Appendix A.2.1.

<sup>1</sup>For simplicity, names uniquely reference individuals in the tree, so as to avoid confusing scenarios like “Avery is the mother of Avery.” Additionally, for efficiency of inference, the only names that are assigned are ones that are used in the conversational context.

As with any computational model of a social phenomenon, this toy kinship model is reductive of many important nuances of identities and relationships. For instance, while the model includes both same- and opposite-gender couples, these couples never split, so step-relations aren’t well-captured. While these kinds of compromises are designed to keep inference tractable, still others stem from limitations of the language itself. For example, many colloquial English kinship terms are gender-binary (e.g., mother, grandfather, daughter), so instantiating them as truth-conditional predicates coerces the generative model towards traditional gender assignments. Similarly, many English names carry strong gender associations, which NLP systems trained on large linguistic corpora pick up on (Caliskan, Bryson, & Narayanan, 2017; Grand, Blank, Pereira, & Fedorenko, 2022). In our examples, we intentionally select gender-neutral names (e.g., Avery, Blake, Charlie, Dana) to emphasize that these naming-based gender inferences are deliberately not part of the reasoning task.

To summarize, language both reflects and constrains our intuitive theories of complex domains like kinship (Sapir, 1929; Whorf, 1956; c.f. Gentner & Goldin-Meadow, 2003 for a review of contemporary perspectives on linguistic relativity), and these tradeoffs manifest concretely in the toy model presented in this section. Fortunately, where this initial “off-the-shelf” kinship model lacks social and cultural nuance, our framework offers opportunities to extend and modify these areas. In section Section 4.1, we look at ways of growing our kinship model to include concepts from non-English-speaking cultures and more inclusive concepts of gender.

**Relational meanings as program statements.** Given a generative model of family trees, we can define a rich conceptual system to make statements about relationships between individuals. Our conceptual system consists primarily of a dozen-odd *derived predicates* that are binary operators over pairs of names; e.g., (father-of? 'charlie 'blake) is true iff Charlie is the father of Blake in a particular tree instance.<sup>2</sup> These derived predicates build on a small number of low-level accessor functions that operate directly on nodes in the tree data structure. For instance, (children-of 'blake) returns a list of names corresponding to the children of Blake in the tree. Finally, our conceptual system includes several higher-order functions, like map-tree, filter-tree, and exists that take custom predicates as inputs and return a boolean. These functions facilitate the expression of a rich compositional semantics by allowing for compound predicates containing conjunctions and disjunctions. Fig. 6 (ii) illustrates several examples of the kinds of statements that can be made using combinations of derived predicates, low-level accessors, and higher-order functions. The full set of definitions making up the conceptual system is given in Appendix A.2.3.

**Translating from language to program expressions.** As in Section 2.2, we use a handful of paired natural language / code examples (Appendix A.2.4) to induce a meaning function via Codex. Because the prompt also includes the generative model source code and the full set of derived predicates, the LLM is able to resolve statements like “Blake has two kids” to the appropriate function (in this case, children-of) using the available definitions. Moreover, we observe zero-shot generalization to linguistic constructs that are not explicitly defined in the prompt, such as the concept of an “only child” (Fig. 7).

**Putting it together: Reasoning from language about kinship relations.** What is the purpose of all of this domain-specific machinery that we have now built up? The answer is two-fold. First, the *generative domain theory* compactly captures the key dynamics of our domain, allowing us to reason about a combinatorially vast space of possible family trees. Meanwhile, the *conceptual system* serves as a higher-level program interface, defining certain relationships that we would like to be able to talk about. Finally, the *large language model* bridges the domain model with natural language, providing a flexible and context-aware way to ground language into conditioning and query statements.

In Fig. 7, we can see how these components come together to facilitate naturalistic reasoning from language about kinship relations. Each natural language utterance translates to a **condition** statement in Church that serves as a constraint on family trees. With each successive **condition**, our uncertainty decreases and our picture of the family tree in question starts to crystallize. Samples from the conditioned domain theory model therefore serve as hypotheses about possible worlds that are consistent with the information provided through language. Furthermore, the *distribution over conditioned samples* provides a principled way to reason about

<sup>2</sup>Note that because our model includes same-gender couples, Blake may have one father, two fathers, or no fathers. Blake also may not exist in the tree in the first place! Crucially, these aspects of the generative model don’t matter to the derived predicate, which simply evaluates whether the relationship in question holds somewhere in the tree.

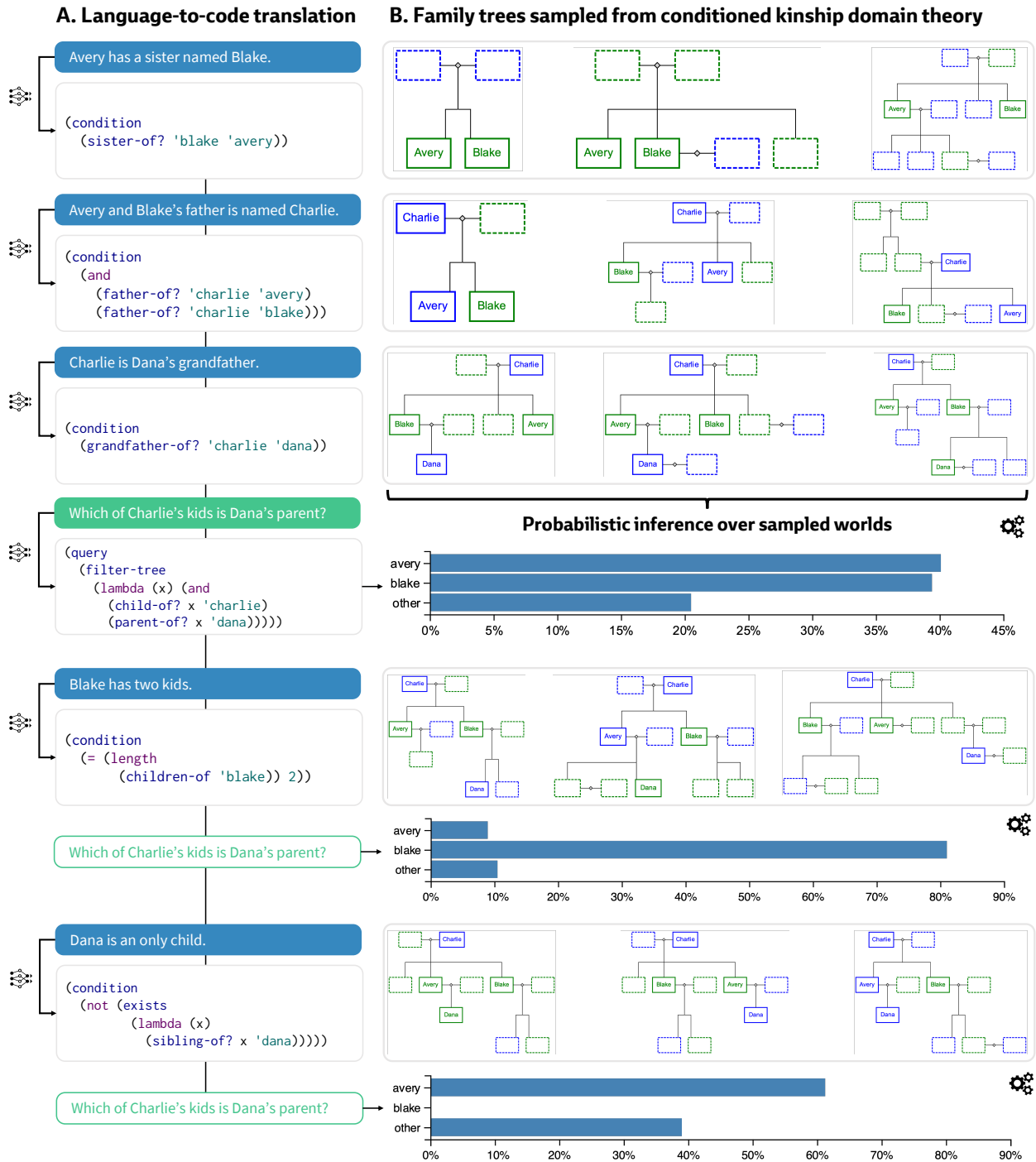


Figure 7: Kinship reasoning from natural language, backed by a domain theory model in the probabilistic language of thought. (A) Natural language utterances about a particular family are readily translated into Church conditioning statements by a LLM. (B) Samples from the conditioned generative domain model are possible family trees that adhere to the growing set of constraints (conditioning statements are cumulative). Reasoning about unknown kinship relations is accomplished through posterior inference against a translated query. With each new piece of information, the model's beliefs reflect both deductive and inductive inferences.

queries, such as *Which of Charlie’s kids is the parent of Dana?* Posterior inference (in this case, accomplished via rejection sampling) faithfully reflects various possible configurations and their relative probabilities. For instance, in Fig. 7, after conditioning on *Blake has two kids*, the model puts > 80% probability on Blake being Dana’s parent, but also factors in low-probability possible worlds where Avery or a third unnamed sibling is Dana’s parent. Yet, despite this confident answer, the model can correctly infer that this same probability drops to 0% in the face of the contradictory information that *Dana is an only child*. Note that the distributional parser plays a crucial role in this inference by providing a correct interpretation of this utterance. Meanwhile, the Church inference engine does the heavy lifting of representing possible worlds and reasoning about them in a principled manner.

**Future directions: Logical and relational reasoning with language models.** Significant recent attention has been directed towards studying reasoning in LLMs. Typical approaches involve engineering prompts so as to induce structured generations in text space that approximate “step-by-step” reasoning (Kojima, Gu, Reid, Matsuo, & Iwasawa, 2022; Nye et al., 2021; Wei et al., 2022). Nevertheless, current evaluations find that even with such methods, LLMs are prone to producing unfaithful reasoning chains in which conclusions do not follow logically from the premises (Golovneva et al., 2022; H. Liu et al., 2023; Lyu et al., 2023; Ribeiro et al., 2023). These issues of consistency have motivated several systems that connect LLMs to external symbolic inference engines that perform deductive inference using Prolog-style backwards chaining (Dalvi, Tafjord, & Clark, 2022; Pan, Albalak, Wang, & Wang, 2023; Weir & Van Durme, 2022). We see this work as closely-related in spirit to our approach, but fundamentally limited to deductive reasoning. (See Appendix A.2.5 for a technical explanation of these limitations.) Of course, we make no claim that Church or its derivatives are the *only* languages that can capture human-like relational reasoning. For instance, ProbLog (De Raedt, Kimmig, & Toivonen, 2007; Dries, Kimmig, Davis, Belle, & De Raed, 2017; Suster et al., 2021), a probabilistic extension of Prolog in which deduction rules can be annotated with probabilities, offers a compelling alternative. Indeed, interfacing ProbLog with a natural language via an LLM-backed meaning function would constitute a promising instantiation of our rational meaning construction framework. Our core assertion here, and in the rest of this paper, is that *representing probabilistic, generative models over possible worlds* is critical to reasoning coherently about a structured domains.

## 3.2 Language for visual and physical reasoning

Sensory detail and physical knowledge pervade our everyday language. We can describe and imagine highly visual objects and scenes—a *few red mugs on a tabletop*, a *tall stack of blue plates*, a *heavy box*, and objects that *move*, *bounce*, and *collide*. We flexibly make predictions about physical events (*what will happen if a kid crashes into that table stacked with plates?*), or infer the underlying physical properties of the world (*how heavy is that box that no one can lift?*), based on situations described entirely in words. As with the other domains we have considered thus far, understanding this language requires integrating over the uncertainty inherent to language, like the possible heights picked out by *tall* and motions picked out by a *bounce*, as well as the uncertainty inherent to how we imagine the physical world itself.

How can we so flexibly relate language to our more general perceptual and physical reasoning? In this section, we illustrate how our overarching framework for language understanding can be modularly extended to capture both of these capabilities. We begin with perception, extending our framework to integrate a **graphics rendering engine** to relate linguistic meanings to visual knowledge (Section 3.2.1). We then build on this approach to integrate a **physics simulation engine** to further interface between language and intuitive, probabilistic physical reasoning (Section 3.2.2). By incorporating these external engines, these sections blueprint how computational models that ground linguistic meaning in a PLoT can interface with other cognitive modules for perception and physical reasoning.

### 3.2.1 Language about visual scenes

To illustrate the structured relationship between language and visual knowledge, imagine how we might talk about a very simple domain of scenes (Fig. 8, top)—tables on which someone was placing some household objects (*mugs*, *cans*, or *bowls*) that come in different colors (*red*, *green*, *yellow*, or *blue*.)

Given descriptions of particular scenes (Fig. 8, bottom), most of us can easily picture tabletop scenes that fit these descriptions, updating what we imagine to incorporate arbitrary new information, like that *everything on the table is blue*, and also that *there are no mugs*, and *lots of bowls*. We can do this despite uncertainty in the language itself—a phrase like *lots of bowls* leaves open just how many bowls there are, though we have general intuitions that there should be more than one or even two bowls on our imagined table. You can also draw a host of fundamentally probabilistic inferences to answer many arbitrary questions about the scenes you imagine, like *how many green mugs* there might be, or whether there are *more red objects or green ones*. The set of scenes you imagine, and the way you answer these questions, is structured and compositional at the level of individual objects and their properties (*a mug*, *a green mug*, *a bunch of green mugs*), and over successive sentences (like *there are many red objects on the table*, *there are just a few green mugs*, and *there are also at least three green bowls*.) The way we talk about scenes like these suggests the level of abstraction with which we mentally represent them. We describe and reason over object categories, lexical properties, numeric quantities, and set relations, and we can easily visualize scenes from these abstract, linguistic descriptions. In contrast, recent evaluations of current multimodal models—large language models fine tuned on corpora of images (Ramesh, Dhariwal, Nichol, Chu, & Chen, 2022; Ramesh et al., 2021)—suggest that even large models struggle with just these kinds of simple but abstract relational concepts in language, such as producing images consistent with quantifiers like *more red things than green things*, or relations like *a plate on top of a cup* (Conwell & Ullman, 2022; Marcus, Davis, & Aaronson, 2022; Radford et al., 2019).

In this section, we propose that the basic motif outlined in our framework also suggests an alternate approach for relating language and visual reasoning. Our architecture draws on the traditions of viewing perception as “analysis by synthesis” or “vision as inverse graphics” from cognitive science and classic computer vision (Battaglia et al., 2013; Gothoskar et al., 2021; Kulkarni, Kohli, Tenenbaum, & Mansinghka, 2015; Lee & Mumford, 2003; J. Wu, Yildirim, Lim, Freeman, & Tenenbaum, 2015a; Yuille & Kersten, 2006). This approach frames visual imagination and visual scene understanding as two sides of the same coin, modeling visualization in a mental **graphics rendering engine** over internal scene representations and perception as probabilistic inference to **invert the renderer** and thereby recover the physical content of scenes from vision. In this section, we show how this general approach to modeling human perception can integrate cleanly into the framework we have sketched so far, augmenting the probabilistic language of thought with an interface to a rendering engine so it can serve as a general, flexible intermediary for relating language, world models, and visual scenes.



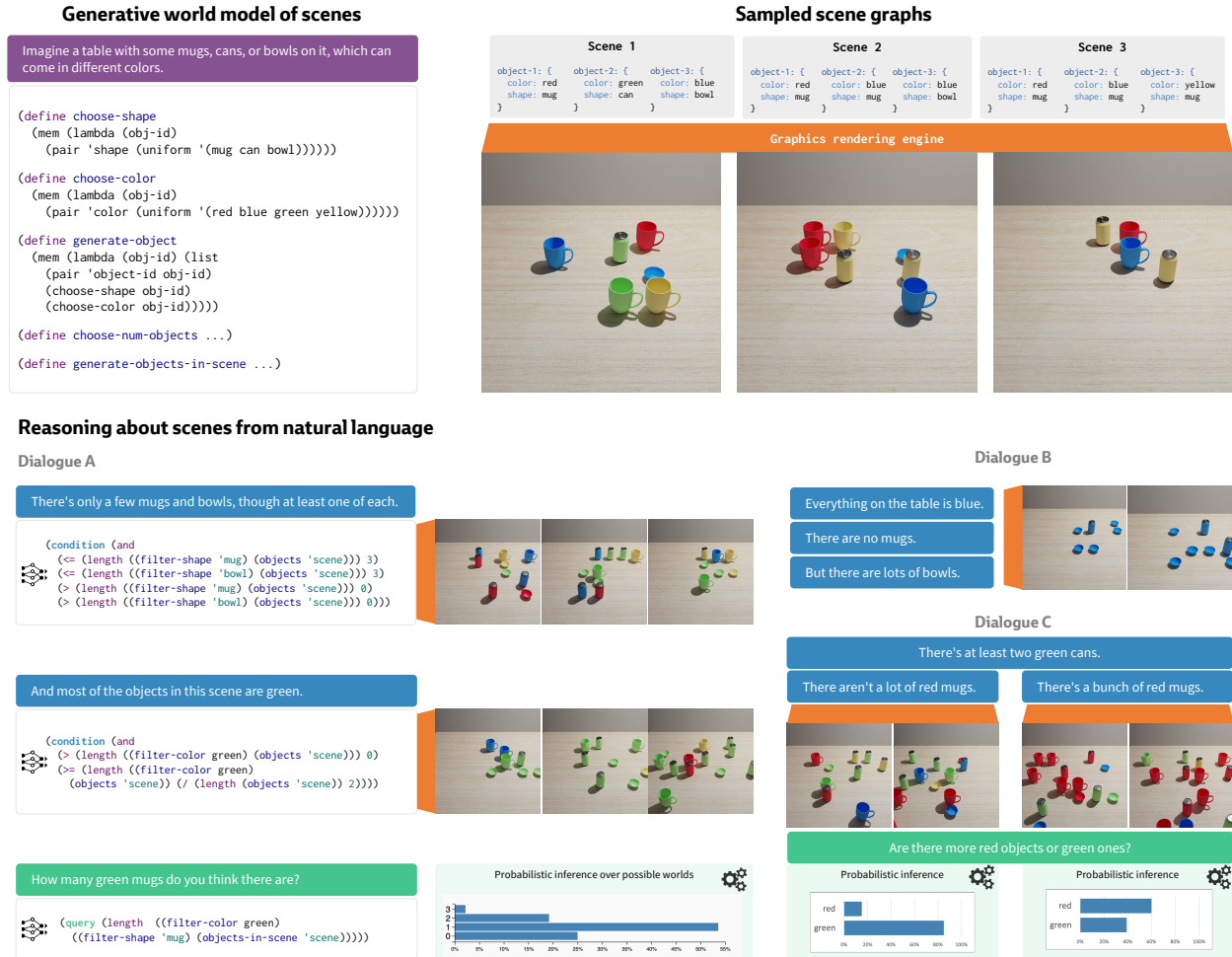


Figure 8: Human language understanding draws on our structured knowledge of the visual world. (Top) A probabilistic generative model describes a prior over tabletop scenes with varying configurations of colored mugs, cans, and bowls. Sampled world states describe a scene based on symbolic object concepts. Interfacing this world model with a *graphics rendering engine* models visual imagination of a given scene. (Bottom) Language about particular visual scenes can now be translated as before into conditions (*blue*) and queries (*green*) on the distribution over scenes, which can be rendered into visual scenes that reflect language.

**Integrating the probabilistic generative model over scenes with a rendering engine.** To model the domain of tabletop scenes, we begin with a probabilistic generative model like those in the preceding sections. The generative program excerpted at the top of Fig. 8 (*purple*) describes a prior over the number of objects in a given scene, and the shape and color of each object. This program is similar in many ways to the kinship model in Section 3.1, which generates possible family trees as a collection of entities and stochastic choices about each one. Similarly, the generative model in this domain generates a particular *scene* by making stochastic choices over the number of objects in the scene (*choose-num-objects*), then generates each individual object (*generate-object*) based on stochastic choices over its possible properties (e.g. *choose-shape* and *choose-color*). This basic structure can be augmented in many ways to model more complex scenes, with more variation over possible properties like *size* or *material*, hierarchical classes of object categories like *dishware*, *cups*, and *mugs*, or hierarchical object structures like a *stack of plates*.

Each sample from the generative model in Fig. 8 is a structured, symbolic representation of a particular scene state, represented in our particular implementation as a list of object dictionaries that map between attribute kinds (like *object-shape*) and values (like *'mug'*). These scene states are very simple instances of the many symbolic scene representations used throughout computer graphics and computational models

of human scene understanding, data structures which model the abstract and semantic contents of scenes (Armeni et al., 2019; Bar-Zeev, 2003; Clark, 1976; Gothoskar et al., 2021; J. Johnson et al., 2017, 2015; Zinberg, Cusumano-Towner, & Vikash, 2019).

We can now extend this probabilistic generative program so that it expresses not just a distribution over possible scene states, but over the visual percepts of each scene. We do so by extending our base probabilistic programming language with a new function, `render`, that takes in scene graphs as inputs and calls out to Blender, a 3D computer graphics engine.<sup>3</sup> Our `render` implementation builds on the basic capabilities of any programmable graphics engine. It defines how symbolic object entities with the properties defined in our model (shapes like *mug*) are rendered and colored into 3D CAD shapes, and can forward render any sampled scene graph into a visual scene with the requisite object types and colors, and overall structure (Fig. 8, top, *Rendered possible worlds*). Collectively, this generative model and rendering interface unites the underlying belief distribution over possible scene states with how each of these scenes might look.

More broadly, this implementation is intended as a simple, illustrative example of how our framework could be integrated to model many, complex relationships between the objects we talk about in a scene and how they look — recent work in scene understanding, for instance, models variation in lighting, viewer angle and distance from the scene, stereo depth sensing, and sources of noise in perception (such as from a viewer who only looks briefly at an image, or an imperfect, non-idealized visual sensor) (e.g., in Deng, Zhi, Lee, and Ahn (2021); Gothoskar et al. (2021); Hughes, Chang, and Carlone (2022); Kulkarni et al. (2015); V. K. Mansinghka et al. (2013); Zinberg et al. (2019)).

**Grounded meanings as program expressions.** By augmenting probabilistic generative models with a graphics rendering engine, we have now extended our framework to allow *language* that describes and asks questions about scenes to interface with visual depictions of those scenes.

In our simple tabletops scene domain, for instance, we can ground linguistic descriptions of the number, kinds, and colors of objects in a scene (Fig. 8, *blue*) like *there’s at least two green cans* or *a few mugs and bowls* into probabilistic program condition statements on scene states in the generative model. As with preceding sections, the translations shown in Fig. 8 are quite straightforward and interpretable, because the generative model we have defined expresses compositional predicates on object properties at the grain of language. Constraints on objects of specific types, like *green cans* are translated into a sequence of conditions on the relevant properties of object entities, successively filtering on the set of objects that are green (`filter-color green`) and then further filtering to the set of objects that are also cans (`filter-shape 'can'`).<sup>4</sup>

Sampling scene states from the conditioned generative model, and rendering these scenes into images with the `render` interface, then produces visual depictions that are consistent with any sequence of observations made in language. This approach disentangles reasoning, as probabilistic inference over a structured generative model, from the perceptual properties of scenes. As with before, we can translate questions like *How many green mugs do you think there are?* into probabilistic query expressions. Our approach reasons about these questions as inferences over the distribution of possible scenes, adapting beliefs about the scenes to condition systematically and coherently on sequences of new statements made in language.

**Translating from language to program expressions.** As with the previous sections, we can now translate actual descriptions and questions about scenes, by using a large language-to-code model conditioned on the generative domain model and a few example pairs of language and code (see Appendix A.3.1 for the full prompt we provide to condition the language-program model).

The translations in Fig. 8 and Fig. 9 generally showcase the local generalizability and flexibility we illustrate in the other sections—the translation is robust to conjunction and syntactic variation, differing numbers of object predicates (*yellow object, red mug*), compositions of object predicates (eg. *a few mugs and bowls*), negations over set quantity (*there aren’t any*), and comparatives over object sets (*more red mugs than green cans*).

<sup>3</sup><https://www.blender.org/>

<sup>4</sup>In our implementation, which can be found in Appendix A.3.1, we derive named color predicates like `green` over the base generative model, which samples color properties over a continuous space of RGB values. This implementation suggests a more general point—that any number of lexical concepts, such as many more arbitrary color names over the underlying color space, can be derived as symbolic predicates over a richer continuous space reflected in the generative model. A similar approach could be taken for other lexical terms that carve up continuous spaces, such as prepositions like *left*, *center*, or *near* over geometric space.

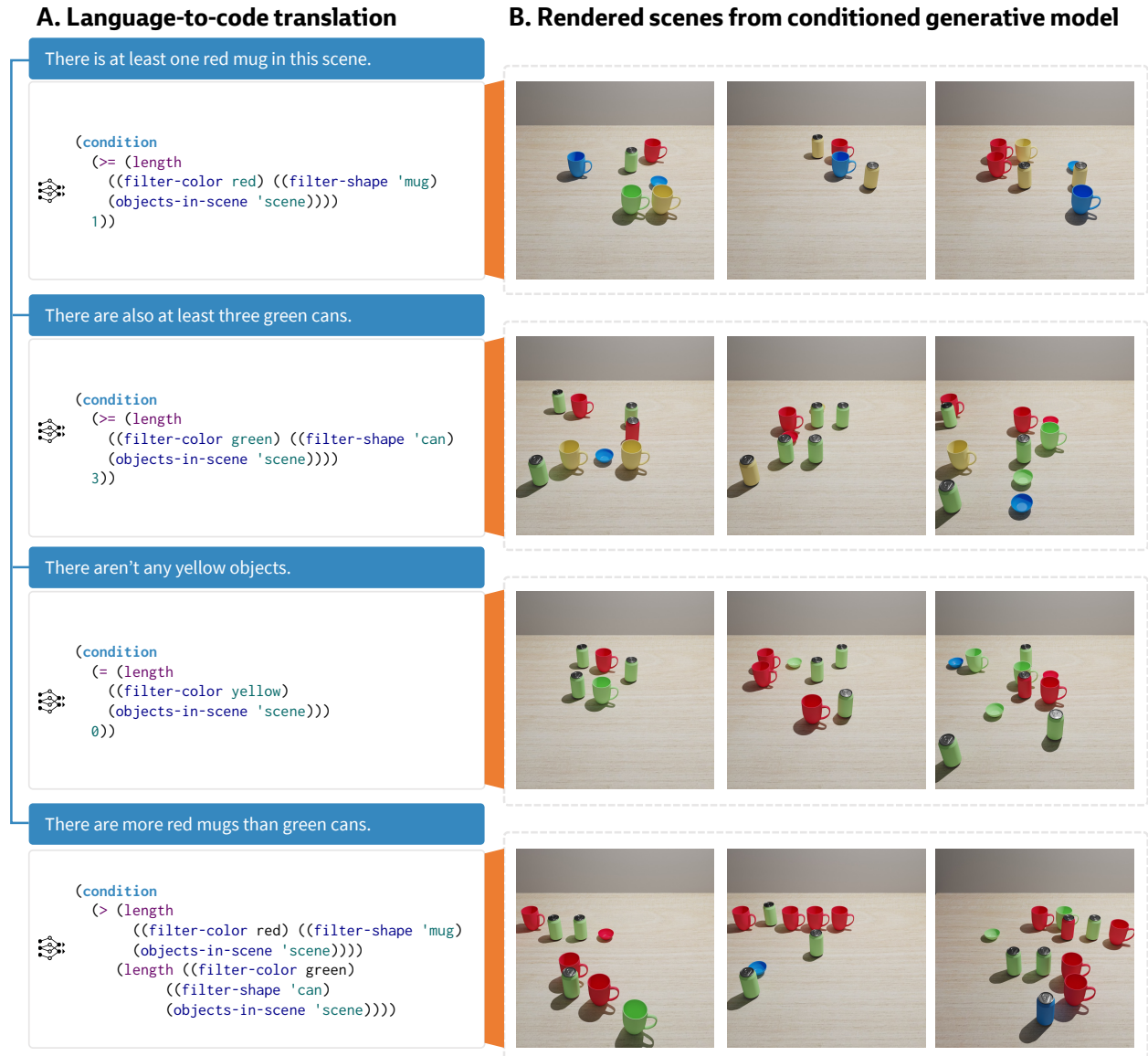


Figure 9: Each sentence in this sequence (*left*) translates into a separate, composable condition expressions that updates the underlying generative model over scene states. After each sentence, sampling symbolic scene states from the updated distribution and rendering them (*right*) yields images that reflect the prior over scenes and are consistent with the information in all successive sentences.

Even on this relatively simple domain, Fig. 8 and Fig. 9 also showcase ways in which the LLM can represent conditional inferences from language to program expressions that go beyond simple, literal semantic meanings. These examples build on what we already find in Section 2.2, in which the LLM can contextually interpret vague language like *very strong* as thresholds on continuous variables in the generative world model.

In this domain, find the LLM can translate vague quantifiers (like *few*, *most*, *aren't a lot*, *a bunch*, or *aren't many*) without explicit program predicates defining each lexical term—the model can directly translate these terms into reasonable, interpretable quantities over sets of objects (such as translating *only a few* to  $(\leq 3)$  objects). We also find that sampling from the distribution over meanings further supports the idea that the LLM represents a broader distribution over intended meanings, including acceptable semantic *variation* in the interpretation of vague lexical terms. Sampling from the distribution at higher temperatures, for instance, we find that our implementation variously translates *most* into program expressions that interpret

this as *more than half*, or *more than 80%*, or other greater fractions, of the set of overall objects in the scene. These translations draw on the language-to-code model’s background prior on language itself (we do not prompt it with examples of these particular phrases), its amortized understanding of how these phrases relate to continuous, named variables in code (like length of a set of objects), and the particular context of the generative world model itself (which defines the prior over numbers of objects that determines the context-specific scale of these graded quantifiers.)

Translations of vague quantifiers like these have been handled in classical semantics and recent accounts as explicit, pragmatic and probabilistic inferences based on context-specific priors—the acceptable quantity most people would infer for *many* mugs on a table is intuitively very different from the quantity intended by *many* grains of sand (Edgington, 1992, 1997; Graff, 2000; Lassiter & Goodman, 2017). The results we show here provide further evidence that LLMs can often amortize many of these inferences, to directly predict common interpretations from language. As we discuss in Section 5, future work might explore more fluid, joint integrations of these approaches to inferring meanings, trading off between the amortized interpretations the LLM can produce and more explicit probabilistic inference, such as conditioning on other information in language. Learning that *Sally is a wholesale porcelain supplier who owns thousands of mugs in a nearby warehouse* might lead you to infer an updated meaning of *Sally has many mugs*, but is a complex inference that we might not expect to be amortized in an LLM from the background distribution of language and commented code.

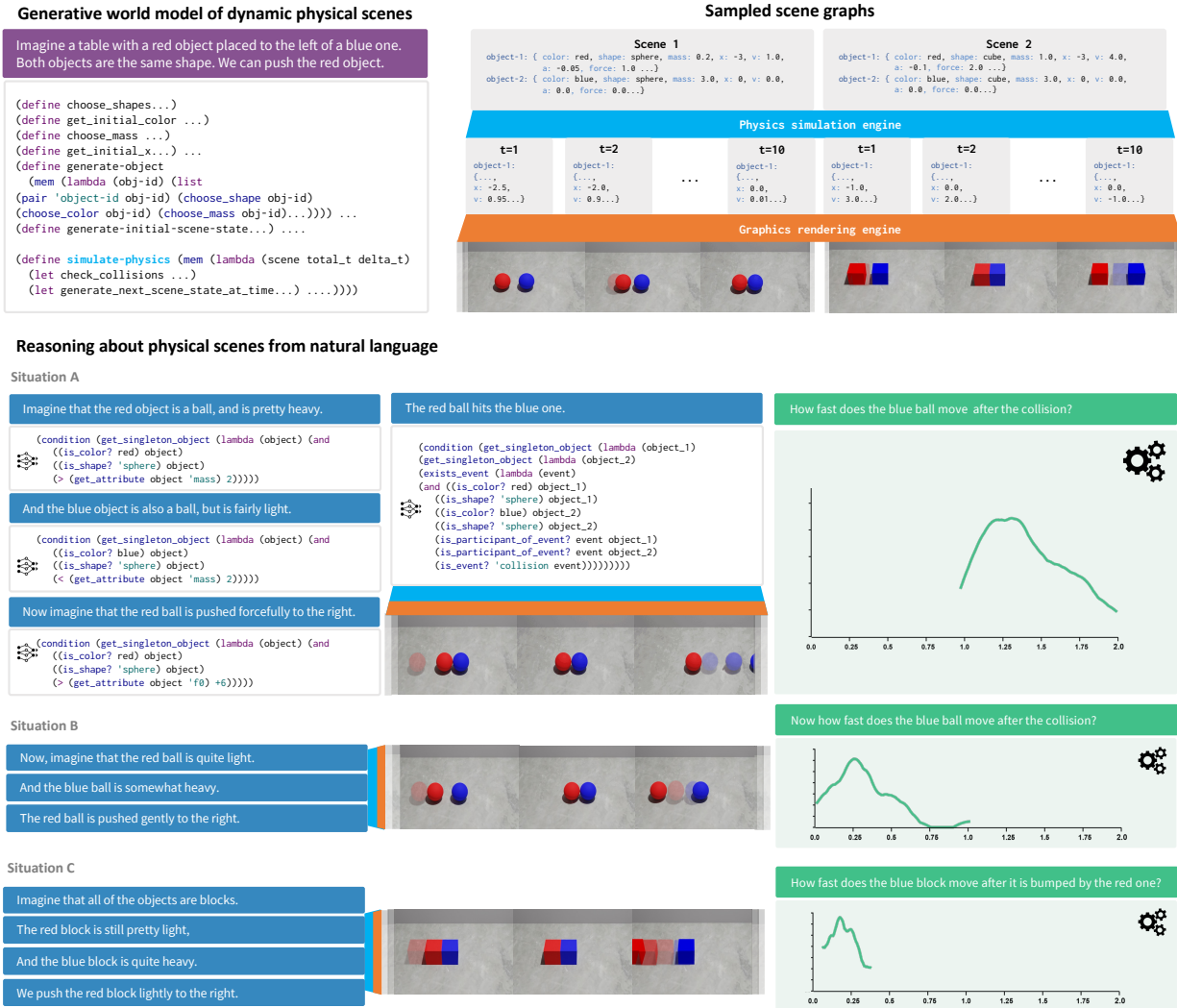
**Putting it together: Reasoning from language about visual scenes.** Taken together, the examples in Fig. 8 show how this approach naturally extends the components of this framework—the ability to describe possible worlds in language, flexibly updating a background distribution of beliefs within a conditioned generative model, and query this model to draw probabilistic inferences—to also ground out in visual scenes. The more extended example in Fig. 9 highlights the more granular, intuitive way in which the distribution over scenes changes to reflect successive new sentences, updating a flexible distribution over scenes that still remains consistent with all of the previous observations.

### 3.2.2 Language about dynamic physical scenes

When we talk about a scene, however, we describe more than just the colors and shapes of objects sitting on a table. We talk in verbs, describing events unfolding in the changing, physical world around us. Consider, for instance, descriptions of another set of tabletop scenes—ones that just involve a red object placed to the left of a blue one on a table (Fig. 10). These scenes are initially even simpler than our tables of cans and dishware, but still afford a range of dynamic and physics-specific descriptions.

You can easily imagine, for instance, what would happen if someone *pushed the red ball gently to the right*—you might say that it would *bump into the blue ball*, and you could likely imagine *how fast the blue ball would be moving* as a result. You can infer how these scenes would change if someone *pushed the red ball much harder*, as if shooting a billiard ball, or *tapped it even more gently*, nudging it forward with their finger, so that perhaps it wouldn’t collide with the blue ball at all. These inferences are sensitive to many other properties of the scene, and of these objects, that we could describe in language, like whether *the red ball is really heavy*, or *the blue ball is very light*, or at least *much lighter than the red one*. If we changed the objects in the scene, and now *placed a red block to the left of the blue one*, your intuitive understanding of how different shapes relate to different degrees of friction would again change how you might see these scenes play out in your mind, and how you might answer questions about their collision and motion.

As adults, we have a deep, general understanding of how physical objects move and behave, and extensive developmental evidence suggests that well before we acquire language, we understand many core physical principles that govern our world (Baillargeon, 2004; Hespos & Baillargeon, 2008; Rips & Hespos, 2015; Spelke, 1990; Spelke, Gutheil, & Van de Walle, 1995; Téglás et al., 2011). A productive line of computational cognitive models, in turn, has modeled human physical understanding as probabilistic inferences over a **mental physics engine**, modeled as programmable physics simulation engines (Battaglia et al., 2013; de Avila Belbute-Peres, Smith, Allen, Tenenbaum, & Kolter, 2018; Lake et al., 2017; Ullman et al., 2017; Yi et al., 2019) like those used in video games, computer animation, and robotics (Coumans & Bai, 2016; Erez, Tassa, & Todorov, 2015; Todorov, Erez, & Tassa, 2012).



everyday physical language—including with respect to uncertainty and vagueness in language about the underlying world state, or combined with inputs from visual reasoning, as discussed in the prior section.

**Integrating the probabilistic generative model over scenes with a physics engine.** To model language about the example scenes we described here—red and blue balls, or blocks, placed on a tabletop (Fig. 10)—we implement a probabilistic generative model that is similar by design to the previous visual scenes domain (a very short excerpt appears in Fig. 10, and the full model appears in Appendix A.3.3). This generative program describes a prior over the possible properties of the objects initially set on a table, modeling scenes as a collection of objects in which each individual object is again generated (`generate-object`) based on stochastic choices over its possible properties (e.g. `choose_shapes`). In this domain, however, we also model an explicit prior over the physical properties of each object, such as its mass (`choose_mass`), and the relationship between shape and friction (as a simple `get_friction_constants` function returns different constants, with a higher constant for blocks than spheres).

As with the visual scenes example, each sample from this generative model again returns a structured, symbolic representation of a possible initial scene state, as a list of object entities that represents each object as a dictionary-like mapping from attribute kinds. This dictionary also stores each object’s initial kinematic state, such as its *position*, *velocity*, *acceleration*, and any *forces* applied to it. To model the various ways we can *push* the objects around, our generative model over scene also implements a stochastic function over possible initial forces (`choose_initial_forces`) applied to an object.

To model how each possible world unfolds as a dynamic scene over time, we implement a `simulate_physics` function (Fig. 10) that integrates the basic functionality of any programmable physics engine into the probabilistic model—this function takes in a scene state that specifies the relevant physical properties of objects, and returns a sequence of scene states forward simulated in time under the laws of physics. In our implementation, this sequence is a list of scene states at each timestep, each which contains its own set of the objects and their properties with relevant kinematic properties (like *position*, *velocity*, *acceleration*) updated at each timestep. The physics model we use in our example is simple enough that we implement it fully within the body of the probabilistic program itself (see Appendix A.3.3) for illustrative purposes—our `simulate_physics` updates each object at each timestep under the basic kinematic laws of Newtonian mechanics, includes a simple implementation of static and kinetic friction under gravity, and models simple collisions as impulse exchanges in momentum.

The rendered simulations we show in Fig. 10 also showcase the interplay between these modular, API-like interfaces integrated into a probabilistic language of thought—combined with the `render` interface from the previous section, we can not only simulate underlying physical scene states, but visualize them by rendering each individual scene state in the sequence over time. Collectively, this model now captures a prior over tabletop scenes, models how any given scene in the distribution unfolds dynamically under physics, and captures how each scene appears visually over time.

**Grounding physical language in program expressions.** By extending the underlying probabilistic world model to interface with a physics engine, we can ground the semantics of language about the physical world in intuitive, human-like physical reasoning modeled by the physics simulation engine over world states.

Descriptions of the physical properties of objects, for instance, like *the blue ball is not very heavy* (Fig. 10) translate into conditions on the mass property of an object in the world state, and maintain uncertainty inherent to language—phrases like *very heavy* translate into conditions that threshold a continuous distribution of possible masses. As in the visual scene example, sampling from the conditioned generative model produces dynamic scene simulations that reflect language. Descriptions of *heavy blue balls*, or *red blocks that are relatively light*, or scenes in which a *red ball is pushed forcefully*, or in which a *red block bumps into a blue one*, all connote sets of scenes that model explicit, physical simulation. In turn, queries about distributions over physical scenes (like *how fast a heavy blue ball will move after it is bumped*) reflect probabilistic inferences that condition on all of these relevant descriptions in language, estimated by sampling and running physical simulations over the possible world states.

In this example, we highlight an approach to translating verbs and descriptions of physical events (*the red ball pushed forcefully to the right*, *the red ball hits the blue ball*) that grounds them directly over continuous variables in our world model. In Fig. 10, for example, our implementation translates *pushed forcefully to the right* into a condition expression that picks out a distribution of initial forces, over a space of continuous force

vectors with direction and magnitude, as the meaning of *push* in a physical world. Similarly, we translate *hits* with respect to collisions simulated by the physics engine between the two object entities.

In our appendix, however, we also implement and show how a discrete event semantics can also be constructed over variables in the physics engine, to highlight potential future connections between our implementation and more classical event semantics representations. Neo-Davidsonian event semantics and related approaches (D. Davidson & Rescher, 1967; Parsons, 1990), for instance, have long modeled events in language with discrete event entities and lexical event predicates (eg. *is\_hitting*) that describe particular categories of events in time. Prior work in classical event semantics has also considered how discrete event representations relate to underlying physical forces (Talmy, 1988), with particularly close connections to lexical semantics approaches (Jackendoff, 1985; Levin, 1993; Pinker, 1984; Schuler, 2005; Talmy, 1988) that realize verb meanings into cognitively-grounded physical concepts of motion and forces.

Our implementation concretely realizes these semantic events and predicates as functions derived entirely on top of a fully realized, continuous world state modeled in a physics engine—*is\_hitting*, for instance, is an event derived on top of the collision mechanics in the underlying physics engine. Other event predicates, like *is\_moving*, or *is\_resting*, can be similarly as thresholds on continuous kinematic properties (here, *velocity*) represented in the world state. Our broader goal is to show all of these can be broadly constructed over a probabilistic language of thought, which grounds out concretely with respect to states in an implementable physics engine.

**Translating from language to program expressions.** As with our visual scenes example, the translations we show in Fig. 10 are again chosen to illustrate the generalization and amortized inferences that the language-to-code LLM can make. Much like vague quantifiers, we find that the context-conditioned LLM can directly infer reasonable meanings for graded terms that pick out thresholds over a numeric distribution—translating phrases like *not very heavy*, *pretty heavy*, and *pretty light* directly into reasonable, context-specific thresholds on continuous masses, or *pushed gently* and *pushed forcefully* into thresholds on forces. Again, we see interesting future grounds for further integrating these kinds of amortized inferences with more explicit, probabilistic inference mechanisms for deriving them—such as to integrate inferences over language with new contextual observations from other modalities, such as perceptual or motor observations from seeing or actually moving these objects that might update one’s background beliefs over the distribution of masses.

**Putting it together: Probabilistic inference and physics simulation from language.** The examples in Fig. 10 show how this approach can capture the nuanced relationships between language and physical reasoning. Language that modulates any of the physical properties in our introduction to this section, from the masses of objects, their shapes and corresponding friction when moving, and the forces they receive, changes the distribution over internally simulated scenes, and is reflected in updated inferences about downstream events.

**Future directions: Perception as inverse rendering and complex physical reasoning as intuitive physics.** As with all of our other examples, its important to emphasize that our *simulate\_physics* interface is almost the simplest possible world model we might construct over physical scenes. The approach we take here is inspired by, but much simpler than, many other probabilistic generative models (Allen, Smith, & Tenenbaum, 2020; Battaglia et al., 2013; Ullman et al., 2017; J. Wu et al., 2015a; Xu et al., 2021) of more complex object configurations in more complex environments (such as *ramps*, *stacks of objects*), many other properties that we can describe about objects themselves (such as their *material*), and arbitrary forces (like *bumping* the table or *dropping* objects from above).

Our approach in these sections also suggests a rich line of future work for reasoning jointly about observations in language, and from perception. While we do not implement a perceptual module in our example, the framework we sketch here can be directly integrated with the broad body of work on *inverse graphics*, which frames scene understanding as inference from observed visual inputs to recover structured representations of a scene’s contents (D. Kersten, Mamassian, & Yuille, 2004; D. K. D. Kersten & Yuille, 1996; Lee & Mumford, 2003; J. Wu, Tenenbaum, & Kohli, 2017; J. Wu, Yildirim, Lim, Freeman, & Tenenbaum, 2015b; Yi et al., 2018; Yildirim, Belledonne, Freiwald, & Tenenbaum, n.d.; Yuille & Kersten, 2006). Our framework suggests a particularly fruitful integration between language and the growing body of work that

combines probabilistic programs and graphics rendering engines (Gothoskar et al., 2021; Kulkarni et al., 2015; V. K. Mansinghka et al., 2013; Zinberg et al., 2019). To draw inferences about visual scenes from perceptual inputs, models like these incorporate convolutional neural networks to make fast, amortized proposals about the scene state from vision, but with respect to a generative program that defines the underlying scene state and guide inferences about particular scenes, such as to reason about occlusion.

Integrated with the approach we describe here, this framework could ground linguistic queries directly into vision, allowing structured inferences for visual question answering (e.g., *counting the number of unique colors of the dishes in a scene*). Moreover, it could enable more complex, joint inferences that integrate visual observation with linguistic information about latent physical properties of objects in a scene (e.g., mass, friction) or the presence or identity of occluded objects. Such multimodal integration holds the potential to shed further light on the ways that linguistic knowledge can shape our understanding of and reasoning about physical scenes.



### 3.3 Language for reasoning about agents and plans

One of the most deeply human things we can talk about is other people. To conclude this section, we turn to language about other social beings—agents who want things, chase goals, and plan how to act in the world around them.

As an illustrative example, we consider a domain (Fig. 11) inspired by C. L. Baker, Tenenbaum, and Saxe (2007), which evaluated commonsense social inferences about agents with different preferences and goals. In our slightly modified example, we consider a set of people with varying food preferences who are making plans for lunch. Based on the map shown in Fig. 11, we’ll imagine which restaurant they might go to based on what foods they like, how far each restaurant is from their office (shown in blue), and whether restaurants happen to be open or closed. We’ll also note that students can bike or walk to any restaurant, but include the intuitive fact that biking is faster on roads, but slower than walking on the lawns.

The original experiments in C. L. Baker et al. (2007) used visual stimuli to depict agents’ paths and plans, but language is a particularly natural and nuanced way to communicate information about other agents. Consider the range of situations we can describe in this simple example. We might leverage our wide vocabulary for describing the spectrum of someone’s preferences and desires—whether they *crave pizza* or *hate vegetables*, or whether they *love sushi* rather than merely *liking it*. We might describe their more concrete, discrete goals, like *getting to a pizza place* or generally *getting to the closest restaurant to the office*. The inferences we draw from language also depend on our intuitions about agents themselves. All else being equal, we expect people to minimize the effort it takes to act, while trying to maximize the value they gain from acting. We might generally expect someone to *walk down Ames Street* if they wanted to go to the pizza place rather than taking an unnecessarily convoluted path, or to *jump on a bike* if they owned one, rather than taking a slower walk there. We also understand, of course, that people need to accommodate the world itself in their plans, and might not go to the pizza place no matter how much they love pizza, if they were told that *the pizza place was closed*.

Perhaps more subtly, but equally importantly, what we know about agents also informs the wide range of inferences we can draw from language about their *actions*. Consider, for instance, what you can infer from being told that someone had started at the office, and was now *walking across the southern lawn*. Because they’re on a direct route towards the vegetarian place, you might infer that they are more likely to prefer vegetarian food, and that they either know or at least believe that the vegetarian place is open. Because they are walking on foot, you might also suspect that they do not own a bike, which would have allowed them to get to the restaurant more quickly. All of these inferences build on a cohesive picture of agents as a whole—our expectations about agents as goal-directed, efficient actors inform how we think about any given action.

As with visual and physical reasoning, this section builds more generally on extensive work in cognitive science and artificial intelligence on *social reasoning*, and seeks to integrate this broader literature into our framework for language. Developmental evidence suggests that we have a core conceptual understanding of agents as goal-directed actors from a very young age (Csibra, 2008; Csibra, Bíró, Koós, & Gergely, 2003; R. M. Scott & Baillargeon, 2013; Spelke & Kinzler, 2007). Computational cognitive models, and the broader AI planning literature, have long approached social inferences like those we describe here under a unifying model of **planning** and **inverse planning** (C. Baker et al., 2011; C. L. Baker, Saxe, & Tenenbaum, 2009; C. L. Baker et al., 2007; M. F. Cusumano-Towner, Radul, Wingate, & Mansinghka, 2017; Jara-Ettinger, Schulz, & Tenenbaum, 2020; Seaman, van de Meent, & Wingate, 2018). This framing couples the forward planning of actions that achieve goals or maximize utilities, to the inverse problem of inferring latent variables about the agent or the world from observations of their actions.

The example below shows how we can extend the modeling motif from our previous discussion of visual and physical reasoning, which shows how our framework can relate language to other core cognitive modules via interfaces implemented in a general probabilistic language of thought. In this section, we introduce *model-based planners* as another such core computational module, which can be integrated into this framework to support a wide range of probabilistic forward and inverse inferences about agents and their actions as they are referenced in language.

**Integrating a probabilistic generative model over agents with a planner.** As a concrete example, the generative program excerpted in Fig. 11 (shown in full in Appendix A.4.1) illustrates how an integrated



model the varying preferences of any given person and the effort of taking particular actions (see [Russell and Norvig \(2021\)](#) for a review). Incorporated into a probabilistic generative model, we can express the distribution of preferences any particular agent could have, and the way these preferences interact with the stochastic mechanics of any given world. In our implementation, we model these varying preferences as a stochastic utility function associated with particular agents and restaurants (`restaurant_utility`). Our implementation shows a bimodal Gaussian distribution, in which people tend to have distinctly negative or positive preferences for any given restaurant, but any other formulation would be easily expressible. We also model how these preferences interact with other aspects of the world—we condition the value an agent derives from actually arriving at a restaurant (`utility_at_restaurant_state`) on whether or not it is open. These utilities interact with possible actions an agent can take to get to different restaurants. We model the distribution of possible actions an agent might take (our `available_actions` function conditions on whether an agent `has_bike`), and the varying effort of individual actions. Our `motion_utility` conditions on the type of action and the state in which it is used, to model the greater effort of biking on grass and the relative ease of biking on the road.

Up to this point, the generative model simply expresses a general prior over world states that includes agent preferences. Now, to model how an agent actually decides on a course of action conditioned on the world state, we can finally introduce a plan interface ([Fig. 8](#)) that calls out to a model-based planner. Our implementation, while simple, implements the basic functionality core to an AI planner—it computes a sequence of actions that achieves a goal or maximizes a value function, subject to an agent’s underlying preferences, available actions, and the conditions of the environment. As with the physics interface in our previous section, our example planner implementation is simple and generic enough that we also implement it fully within the body of the probabilistic program itself (see [Appendix A.4.1](#)) for illustrative purposes. Our implementation here uses a simple value-iteration algorithm, which computes an optimal policy of action to trade off between the value an agent derives from particular restaurants, and the cost of taking actions (walking or biking in any given direction, from any location in the map) towards them.

**Language about agents as program expressions.** By augmenting the probabilistic generative model with a planner, we can now ground many of the basic ways we talk about agents themselves in probabilistic conditions and queries to this model.

Language about what people want and prefer, like whether someone *wants*, *likes*, *loves*, *doesn’t mind*, or *hates* a given restaurant in our example domain, can construct formal conditions over underlying utility variables, that in turn drive inferences about how the agent will act. In the examples shown in [Fig. 8](#), we illustrate the semantics of these terms as conditions constructed directly over the continuous utility variables defined in this domain. We could also derive a more explicit set of predicates (like a Boolean `likes?` predicate, defined over the underlying utilities), but as in several previous sections, we show these more transparent semantics (like translating *likes* into a  $> \theta$  threshold on utilities) to illustrate how language relates to our model of agents, and to demonstrate the amortized inferences that a language-to-code model can make in directly inferring these threshold values in context, and for new preference words.

Observations about relevant aspects of the environment, like whether the *sushi place is closed* or *Alex has a bike*, are translated as in previous sections into conditions on the generative world model. In this integrated framework, these observations now support downstream inferences about how agents might change their behavior with respect to what we are told about the world.

Finally, of course, explicit observations and queries about someone’s goals, plans, and individual actions (*Gabe was biking East on Barlow Street*, or *What restaurant will Alex go to for lunch?*) can be interpreted with respect to the underlying, model-based planner, to drive inferences about forward planning agents choosing actions in the world, and inverse inferences over the many latent variables in the world that collectively explain language about someone’s actions.

**Translating language using language-program distributions.** We showcase several distinct examples ([Fig. 11](#)) of context-sensitive, pragmatic inferences derived using a language-to-code meaning function conditioned on language and this generative world model.

As in previous sections, we find that the LLM can directly ground vague, graded terms in context-specific thresholds over particular continuous variables in the probabilistic world model. Here, this approach grounds preference terms (*doesn’t mind*, *loves*) into reasonable thresholds over the utility variables in the world model

(Fig. 11). We find that the LLM can both infer reasonable utility thresholds and generalize to words not explicitly given as example translations: we prompt the model with a handful of examples pairs, such as a translation that maps the word *likes* to a  $> 0$  threshold on utilities, and the LLM successively generalizes this parse to ground other preference terms like *hate* and *love*, presumably based on the comparative valences of these preference terms in the broader distribution of language.

We also find that the LLM can directly translate quantifiers over contextual sets in this domain—like *likes all of the nearby restaurants*—into a conjunction over the set of restaurant literals in this domain, by conditioning on the generative world model during parsing. More concretely, this means the LLM identifies the relevant restaurants list (shown in the excerpted generative world model in Fig. 11), and conditions on it to directly produce the unrolled conjunction over the list contents (`(and (is_open 'sushi) (is_open 'pizza)...`) intended by *all restaurants*, amortizing the computation that would have otherwise been necessary over a more literal semantics like `(all restaurants)`. Together with the previous sections, these examples suggest how our framework might jointly support explicit inferences from language into various expressions in a language of thought, and learned patterns from the large language-to-code model that amortize some of these inferences over time, which we discuss directly as grounds for future work in Section 5.

**Putting it together: Probabilistic inference and planning over language.** The several example dialogues shown in Section 2 show how this approach captures the integrated social inferences we make about agents in language. We can now query the plans and goals of agents, deriving inferences with respect to the *forward planning* module incorporated into the underlying generative model, conditioning flexibly on arbitrary information in context, and updating expectations about where agents will go, and how they will change their plans based on new observations about the world. In turn, we can derive *inverse planning* inferences, like whether *the pizza place is open*, based on relatively tangential information about someone’s actions—knowing that an agent really likes pizza, but is seen taking a path that wouldn’t efficiently lead them there. All of these inferences fall out of the same underlying generative model, which unifies these distinct observations about people and the world in language with respect to a formal model of how agents tend to behave.

**Future directions: Scaling integrated world models for planning and inference.** The plan function in our example implements a very simple but *model-based* planner—it computes actions based on an underlying, structured model of the world. In comparison to the other domains in this paper, linguistic planning and social reasoning have received perhaps the attention in recent work, in part because complex reasoning about other agents (Ullman, 2023) and precise general planning tasks (Bubeck et al., 2023; Valmeekam et al., 2023) appear to pose outstanding challenges for even the largest current language models. Recent work has sought to interface large language models with classical planning languages and symbolic planners (eg. Collins, Wong, Feng, Wei, and Tenenbaum (2022); Ding, Zhang, Paxton, and Zhang (2023); B. Liu et al. (2023); Xie et al. (2023)), as well as general purpose programming languages used to express code-based policies (G. Wang et al., 2023). All of these approaches suggest directions for scaling the simple planning implementation we show here—our goal is to show how classical planning approaches can be nested within and integrated into probabilistic generative models to support a range of complex reasoning about other agents to infer their goals and actions from information in language.

Collectively, the broader cognitive science and AI planning literature suggests many directions for scaling up this model towards more of the nuance in human social reasoning, each which would in turn allow this paradigm to ground richer and more nuanced descriptions of agents, and the inferences we draw from this language. Some of the most important ones include planners and planning languages that designed to express explicit and discrete *goals*, like *wanting to be at the highest rated pizza place within a half-mile radius* or *trying to get a plate of sushi for under ten dollars*, rather than continuous values and utilities (G. Davidson, Gureckis, & Lake, 2022; Fikes & Nilsson, 1971; D. McDermott, 1982; D. M. McDermott, 2000; Pednault, 1989); planners that model explicit uncertainty about the world itself, like agents who *don’t know whether a restaurant is open or closed until they get there* (C. Baker et al., 2011; Kaelbling & Lozano-Pérez, 2013; Zhi-Xuan, Mann, Silver, Tenenbaum, & Mansinghka, 2020); hierarchical planners that recursively turn goals into more specific subgoals to account for plans over longer timescales, at differing levels of abstraction (Kaelbling and Lozano-Pérez (2011)); and recursive models of agents who are themselves thinking about other agents, such as models of *two people trying to meet up* at a restaurant that they think will satisfy both of

them, or where they might be most likely to find the other (C. Baker et al., 2011; Krafft, Baker, Pentland, & Tenenbaum, 2016; S. A. Wu et al., 2021). Each of these could allow this paradigm to ground richer and more nuanced descriptions of agents, and the inferences we draw from this language.

**Conclusions.** Together with the previous sections on vision and physics, our approach to grounding language about social agents highlights the more general computational account suggested by our framework. By translating language into probabilistic programs, language can construct, describe, and drive inferences over our internal world models. These may in turn incorporate many more specific computational engines—modeling how scenes are visualized, how physics unfolds in the world, or how agents plan towards their goals—as modular interfaces that can be called upon in a general probabilistic language of thought.

## 4 Growing and constructing world models from language

In [Section 3](#), we illustrated how domain theories expressed in a probabilistic language-of-thought can provide flexible and powerful scaffolding for language understanding. In each, generative world modeling programs provided a unified substrate for defining a structured domain model and representing the meanings of sentences. But where do these world models come from? If we want our PLoT account of language understanding to scale beyond the knowledge that can be hand-coded by a programmer, we need to provide some account of how such a system might acquire new concepts and domain theories.

One of the hallmarks of human communication is our ability to *teach* each other fundamentally new concepts in language. We coin new words, define interrelated conceptual systems, and describe entirely new world models, explaining the abstract underlying structure of whole domains. Because language spans so many aspects of human thought, it is perhaps a uniquely powerful tool for structuring learning. In language, we can define new concepts and domains that are integrated into our inferences, relational reasoning, understanding of the visual and physical world, and goals and plans.

How do we learn new concepts and world models from language? And how can we build computational systems that can be taught in language as we teach each other? In this section, we showcase the extensibility of the framework we have proposed as a unified model for relating language to thinking. Because world models in a PPL are expressed as programs, the same core computational components can be used to extend and construct world models themselves from language. In [Section 4.1](#), we show how we can extend an existing domain model with new lexical concepts. Then, in [Section 4.2](#), we turn to language that communicates an entire background domain model from scratch. Through these simple explorations, we aim to point towards a near-term horizon where systems might construct rich and nuanced probabilistic models to make sense of their linguistic environments and the broader world around them.

### 4.1 Growing a world model from language

How can we enrich our world models with concepts learned from language? Let’s consider, for instance, the kinship domain model used in the relational reasoning example in [Section 3.1](#). The probabilistic program used in this example described a basic generative model over family trees, and then defined a handful of primitives, such as concepts for grandparent and sibling. But most people know and can talk about many more kinship relations than those included in that simple example, such as uncles, aunts, and cousins. What happens when we use language that invokes one of these undefined concepts?

**Condition:** *Avery is Blake’s uncle.*

```
(condition
  (exists (lambda (x) (and
    (sibling-of? x 'avery)
    (parent-of? x 'blake))))))
```

The LLM proposes an initial translation that includes some of the important components in the concept of an “uncle.” However, several key details are not quite right: an uncle should be the *brother* of Avery’s parent, not just a generic sibling. Moreover, an uncle can come from outside the bloodline, in which case this definition would not fit. Much like a person learning English, the LLM has a partial notion of this concept, but could benefit from more explicit instruction from a knowledgeable teacher. In this section, we introduce a new **define** construct that does just this by prompting the LLM to generate a new definition from language.

**Define:** *An uncle is the brother of one’s father or mother, or the husband of one’s aunt.*

```
(define (uncle-of? name_a name_b)
  (or (exists (lambda (x) (and
    (brother-of? name_a x)
    (parent-of? x name_b))))))
```

```
(exists (lambda (x) (and
  (husband-of? name_a x)
  (aunt-of? x name_b)
))))
```

We’ve used `define` to fill in a bit of common knowledge that was missing from our conceptual system. But the mental frameworks we use to reason about the world are constantly undergoing conceptual change, both at an individual and a societal level (Carey, 1999; Posner, Strike, Hewson, & Gertzog, 1982). For instance, shifts in cultural notions of gender and identity have introduced new kinship terms into English. One of the hallmarks of language is the ease with which we can coin and communicate new concepts, like the following:

*“Pibling” is a gender-neutral term for “aunt” or “uncle” that refers to the sibling of one’s parent.*

Finally, as we touched on in Section 3.1, kinship systems vary widely; certain cultures have kinship concepts that are more granular than those found in English. For instance:

*In the language of the Northern Paiute, a group of peoples indigenous to the Great Basin region of the US, “pāan’i” refers specifically to the sister of one’s father.<sup>5</sup>*

From this definition, we can incorporate the concept of a *pāan’i* into our growing set of kinship concepts. Our framework elegantly captures this ability to learn new concepts in language that we can then use productively to construct new sentences and reason about coherently against the background of our existing world knowledge. Here, we walk concretely through how the basic components of our framework are combined to grow the original kinship model with new concepts.

**Linguistic meanings as program expressions.** Much as we interpreted observations as program expressions that conditioned an existing world model, and questions as program expressions that queried it, a sentence like, *The term “pāan’i” refers to the sister of one’s father*, can be modeled as a program expression that *defines* a new such primitive relation, `paani-of?`. The examples in Fig. 12 show how the semantics of this sentence, along with the other kinship concepts introduced in the introduction to this section, can be similarly understood as expressions that define new conceptual primitives. These expressions are particularly interesting because they are defined in terms of other concepts, like `sister-of?` and `father-of?`, that make up this conceptual system. In this way, our treatment of concept learning is closely linked to the idea of a *conceptual role semantics* (Block, 1998; Field, 1977; Greenberg & Harman, 2005; Harman, 1982), in which concepts (including lexical concepts) derive meaning from their interrelated roles and relationships to other concepts. In these examples, interpreting these sentences as program expressions defined over the base generative model showcases the flexible role that the generative modeling program can play, in relation to language about the domain. While our example showcases simple relational definitions over the underlying world model, it is worth noting that these are not the only kinds of functional definitions that we could learn to extend a world model from language. This general approach can be used to make meaning from sentences that grow an underlying world model in other ways, such as by defining new random variables (like phenotypic eye colors or other inherited traits) that extend the probabilistic generative model.

**Translating with a language-program distribution.** While the meanings of these sentences play a different role in our framework—they extend the world modeling program, rather than condition or query it—they are still program expressions. Therefore, with minor adjustments, we can use the same language-to-code LLM approach to ground these new concepts in our world model. To derive each of the translations shown in Fig. 12, we feed the LLM the same prompt as in Section 3, which includes the existing generative model and example translations. The final line of the prompt begins with `Define:` and contains the language describing the new concept definition. Each sentence is then translated into the new `define` statements which construct new conceptual kinship primitives. In sum, linguistic definitions are simply another kind of program expression we can translate into from language.

<sup>5</sup>At the time of writing, a Google search for the term “pāan’i” yielded zero results. The term itself was pulled from a non-searchable table in a century-old manuscript (Lowie, 1930). As far as real-world kinship terms go, it is comparatively unlikely—though not impossible—that “pāan’i” was part of Codex’s pretraining data.

**A. Existing generative world model**

```
(define (person person-id parent-1-id parent-2-id)
  (list
    (pair 'person-id person-id)
    (pair 'name person-id)
    (pair 'gender (person->gender person-id))
    (pair 'parent-1-id parent-1-id)
    (pair 'parent-2-id parent-2-id)))
...
(define (parent-of? name_a name_b)
  (member? name_a (parents-of name_b)))
(define (father-of? name_a name_b)
  (and (equal? (get-property name_a 'gender) 'male)
    (parent-of? name_a name_b)))
(define (sister-of? name_a name_b)
  (and (equal? (get-property name_a 'gender) 'female)
    (sibling-of? name_a name_b)))
...

```

**C. Extended world model**

```
(define (uncle-of? name_a name_b)
  ...
(define (pibling-of? name_a name_b)
  ...
(define (paani-of? name_a name_b)
  ...

```

**B. Defining new concepts via language-to-code translation**

An uncle is the brother of one's parent, or the husband of one's aunt.

```
(define (uncle-of? name_a name_b)
  (or (exists (lambda (x) (and
    (brother-of? name_a x)
    (parent-of? x name_b))))
    (exists (lambda (x) (and
    (husband-of? name_a x)
    (aunt-of? x name_b))))))
```

“Pibling” is a gender-neutral term for “aunt” or “uncle” that refers to the sibling of one's parent.

```
(define (pibling-of? name_a name_b)
  (or (uncle-of? name_a name_b)
    (aunt-of? name_a name_b)))
```

In Northern Paiute, “pāan'i” refers to the sister of one's father.

```
(define (paani-of? name_a name_b)
  (exists (lambda (x) (and
    (sister-of? name_a x)
    (father-of? x name_b))))))
```

**D. Grounding new language in learned concepts**

According to historical records, Chief Winnemucca was either the father or uncle of Numaga.

```
(condition
  (or (father-of? 'winnemucca 'numaga)
    (uncle-of? 'winnemucca 'numaga)))
```

Numaga's son would have called Sarah Winnemucca his “pāan'i”.

```
(condition
  (exists (lambda (x) (and
    (son-of? x 'numaga)
    (paani-of? 'sarah x))))))
```

Sarah had a sibling named Natchez.

```
(condition
  (sibling-of? 'sarah 'natchez))
```

How many piblings would Numaga's son have had?

```
(query
  (length (filter-tree
    (lambda (x)
      (exists (lambda (y)
        (and (pibling-of? x y)
          (son-of? y 'numaga))))))))
```

**Probabilistic inference**

Figure 12: Extending the kinship world model with linguistic descriptions of kinship relations drawn from contemporary English and a low-resource language (Northern Paiute). A language-to-code LLM is prompted with (A) the existing generative model code and (B) language describing novel kinship relations to produce new concept definitions in Church. The extended world model (C) now supports probabilistic reasoning from language that contains these new concepts (D).



**Growing the domain model with new program expressions.** Finally, by incorporating the meanings of sentences like *The term “pāan’i” refers to the sister of one’s father* back into the domain model itself, we have formalized a simple approach for enriching a world models with concepts learned from language. Each sentence shown in Fig. 12 is translated into a program expression that defines a new relational function which extends the set of conceptual primitives that comprise the extended kinship domain.

The more general principle here is not limited, of course, to kinship concepts. We could extend any of the domain models in each of our previous examples with new concepts learned from language. For example:

- *In tug of war, the strongest person on a team is referred to as the “anchor”.*
- *A “monochrome” scene is one in which every object is the same color.*
- *On “National Restaurant Day”, all the restaurants in town are guaranteed to be open.*

Our proposal in this section is closely related to other work which formalizes the learning of new concepts as the learning of new program components, such as program synthesis systems that bootstrap a growing library of domain-specific concepts constructed out of an initial programming language (Bowers et al., 2023; Dechter, Malmaud, Adams, & Tenenbaum, 2013; Ellis et al., 2020); work that formalizes the learning of new concepts from language as the learning of new program primitives (Shin, Brockschmidt, Allamanis, & Polozov, 2018; Summers, Hawkins, Ho, Griffiths, & Hadfield-Menell, 2022; C. Wong, Ellis, Tenenbaum, & Andreas, 2021); and semantic parsers that bootstrap lexicons of compositional word meanings, defined in a formal logical language, for interpreting new sentences (Artzi, Das, & Petrov, 2014; Cai & Yates, 2013; Kwiatkowski, Zettlemoyer, Goldwater, & Steedman, 2011).

The framing we describe here showcases the tight integration between language, meanings, and the probabilistic programs that form the formal substrate for modeling the world in our framework. Language that specifies new parts of a world model can be cleanly interpreted as program expressions, which are used to extend the generative world modeling program itself. These generative models in turn provide the basis for reasoning about new observations that build on these learned and structured bodies of conceptual knowledge. Returning to the themes of our introduction, human-like thinking, under the broader computational approach we take throughout this paper, is formalized as probabilistic programming and inference over probabilistic programs. This is how we construct models of and reason about the world. Language, then, is an especially powerful tool for constructing programs of all kinds—ones that condition and query existing world models, and ones that actually construct and extend the flexible domain models themselves that undergird linguistic meaning and thought.

## 4.2 Constructing new world models from language

So far, we have assumed that language understanding happens in the context of a particular world model appropriate for the situation at hand, containing definitions of key concepts like *sibling* for kinship reasoning, or *strength* for reasoning about playground games. We have now seen how these models can be extended with new lexical definitions on the fly, but the question remains of where these background world models come from in the first place. The full answer to this question is likely complex: people learn about the world in all sorts of ways. But in some settings, people do seem to acquire new world models largely *through* language: we read the rules of new games, are taught the workings of machines, and take classes on the causal structure of many other complex systems (the human body, the solar system, the government). In this section, we broaden our scope beyond language that conveys new concepts that extend an existing domain model to consider how language can define entire new domain models from scratch.

As a concrete example, let’s return to the scenario from Section 2.2. Suppose your friend is telling you about a tug-of-war tournament that took place the prior weekend—only this time, you’ve never heard of tug-of-war before and don’t know how it’s played. Your friend might explain the scenario to you using language—indeed, their description might sound similar to the one our paper itself uses to convey the concepts of this particular situation:

*Tug-of-war is a game played between teams of players. First, strength levels vary widely from person to person. Furthermore, each person has a percentage of the time that they are lazy. The strength of a team is the combined strength of its members, except that in any given match, each*

*player may decide to be lazy, and thus contribute only half of their strength. Whether one team beats another just depends on which team pulls stronger that match.*

Given this language, you can learn the underlying domain model necessary to reason about future observations (*Even working as a team, Lio and Alex could not beat Josh*) and answer questions (*How strong is Josh?*). In this section, we explore how the components of our framework can be used to construct an entire domain model as it is communicated in language, using the tug-of-war domain as an illustrative example.

**Linguistic concepts as program expressions.** Considering the vignette above, we might distinguish between two kinds of statement in your friend’s description of tug-of-war:

- Some statements introduce new concepts solely in terms of previously introduced concepts (e.g., *Whether one team beats another just depends on which team pulls stronger that match*).
- Other statements posit the existence of new primitive concepts, like *strength* and *laziness*, that have certain properties (e.g., *Strength levels vary widely from person to person*).

The first case is similar to the sentences we saw in [Section 4.1](#), and we can interpret them as language-of-thought definitions. The second case, however, is genuinely new: these sentences neither define new words in terms of an existing domain theory, nor encode predicates over possible worlds. Rather, they define *random variables* that we expect to have different values in each possible world.<sup>6</sup>

In Church, such variables can be defined using `mem`: for example,

```
(define strength (mem (lambda (person) (normal 100 20))))
```

declares that expressions of the form `(strength person)` are well-formed and evaluate to a number in each possible world, and that our prior distribution for a new person’s strength is a Gaussian centered at 100. (The `mem` construct *memoizes* the defined function, so that repeatedly evaluating `(strength 'lio)` in the same world will always give the same result.) It might seem strange to claim that the *meaning* of the sentence “*Players have different strength levels*” includes a specific prior over player strengths, like `(normal 100 20)`. We do not make this claim: rather, the meaning function induces a *distribution* over possible definitions of strength, each of which uses a different prior. What the different possible translations have in common is that they model strength as a continuous variable assigned on a per-player basis, with some population-level variation. See [Footnote 6](#) for further discussion of this distribution, and how it might arise from the literal meaning of the sentence being translated.

**Translating new concepts from language.** As before, because each sentence *means* some distribution over program fragments in a probabilistic language, we can use probabilistic language-to-code translation models like Codex as models of the meaning function. In [Fig. 13](#), we prompt Codex with an unrelated example world model in domain about diseases and symptoms, and then ask it to translate sentences defining the tug-of-war domain.

<sup>6</sup> An alternative perspective is that the sentences we consider in this section—both straightforward definitions, and sentences introducing new primitive concepts—*do* still encode predicates on possible worlds. According to this viewpoint, a sentence like “*The term ‘uncle’ refers to the brother of one’s parent, or the husband of one’s aunt*” is an assertion that can be true or false; maybe *uncle* means something different in another possible world. To understand this viewpoint within our framework, we need to imagine that there is a background world model that models uncertainty about the *code* of a first-order world model (which definitions exist, and how they are defined). If we had such a *model over world models*, then sentences like “*Players have different strength levels*” could be interpreted as conditioning statements, observing that *strength* exists as a variable and that its value should vary from person to person. Conditioning on this constraint, we could then sample from the posterior over world models that satisfy this property. In this posterior, there would be some uncertainty over exactly how strength is modeled: e.g., does it vary according to a Gaussian distribution, and if so, with what parameters?

We find this view appealing, and believe that making it practical would be an intriguing technical challenge, requiring new developments in the field of *Bayesian probabilistic program synthesis* ([Saad, Cusumano-Towner, Schaechtle, Rinard, & Mansinghka, 2019](#)). In this section, we take a shortcut, of assuming that the meaning distribution induced by a sentence like “*Players have different strength levels*” directly samples model fragments consistent with the statement. That is, we ask our meaning function to *amortize* inference in the hierarchical model, directly proposing code defining *strength*, rather than first translating to a conditioning statement about *strength* existing, and then using a slower inference algorithm to infer its definition.

**A. Prompt, containing unrelated example world model**

```
;; We define a probabilistic model in Church of the following scenario.
;; At any given time, about 1% of the population has lung cancer,
;; 20% have a cold, 10% have a stomach flu, and 0.5% have TB.
(define lung-cancer (mem (lambda (person) (flip 0.01))))
(define cold (mem (lambda (person) (flip 0.2))))
(define stomach-flu (mem (lambda (person) (flip 0.1))))
(define TB (mem (lambda (person) (flip 0.005))))

;; If you have a cold, there's a 50% chance you have a cough.
;; 30% of people with lung cancer have a cough, and 70% with TB.
;; There's also a small chance you have a cough even if you're otherwise healthy.
(define cough (mem (lambda (person)
  (or (and (cold person) (flip 0.5))
      (and (lung-cancer person) (flip 0.3))
      (and (TB person) (flip 0.7))
      (flip 0.01)))))

;; Whether a person coughs during a particular visit to the doctor's office
;; depends on whether they have a cough, and a bit of random chance.
;; Note that this will differ each time they go to the doctor's office, so
;; we do not use 'mem' (which memoizes the result).
(define coughs-on-particular-visit (lambda (person) (and (cough person) (flip 0.7))))
```

**B. Defining a new world model from scratch via language-to-code translation**

```
;; Now, let's define a different probabilistic model of the following scenario.
;; It is totally unrelated to the previous model and does not reference the functions above.
```

First, strength levels vary widely from person to person.

```
(define strength (mem (lambda (person) (normal 100 20))))
```

Furthermore, each person has a percentage of the time that they are lazy.

```
(define laziness (mem (lambda (person) (uniform 0 1))))
```

The strength of a team is the combined strength of its members, except that in any given match, each player may decide to be lazy, and thus contribute only half of their strength.

```
(define team-strength
  (lambda (members)
    (apply + (map (lambda (member)
      (if (flip (laziness member))
          (/ (strength member) 2)
          (strength member)))
      members))))
```

Whether one team beats another just depends on which team pulls stronger that match.

```
(define team-beats-team
  (lambda (team1 team2)
    (> (team-strength team1) (team-strength team2))))
```

Figure 13: Constructing the tug-of-war model from scratch. This can be accomplished with the same overarching language-to-code approach. (A) We provide a prompt containing one or more unrelated world models as examples. (In this case, the world model defines a medical diagnosis domain.) (B) Prompted line-by-line with language explaining the tug-of-war, Codex constructs a generative model from scratch that is semantically equivalent to the one from Section 2.2 (modulo some superficial naming and parameter choices).

**Constructing the domain model from new program expressions.** By translating each sentence in a domain description in sequence, we can—starting with *no* definitions beyond those built into Church—build a domain model just as rich as the ones we hand-coded in earlier sections. In Fig. 13, although the specific priors may vary slightly, Codex recovers all the essential structure of our hand-coded tug-of-war model. Once we have a new domain model, we can immediately begin interpreting observations and queries, like those in Section 2.2, or continue to extend the domain model with new definitions.

**Putting it together: Growing and constructing world models from language.** In this section, we’ve illustrated how the same basic building blocks used in the rest of the paper — language-to-code translation and probabilistic programs — can be used to extend and construct new world models. Hopefully, these simple sketches highlight a much deeper point: systems that have the ability to author world models in a universal programming language like Church can take advantage of the infinite expressivity of code to generalize to new kinds of language and thinking.

Nevertheless, the examples presented in Section 4 were limited to cases where there was an explicit connection between linguistic instructions and the resulting probabilistic programming expressions. In reality, this relationship is often indirect; language typically only provides us clues about how to think about a situation. In still other instances, we assemble world models in the absence of language, drawing instead on prior experience of similar situations. How can we build systems that *learn* to build world models on-the-fly? How can such systems *remember* and expand on prior world models to understand new situations? And how can they incorporate not just language, but the full spectrum of experiences in the world? In Section 5, we consider these questions as part of a discussion of the many future research directions needed to scale our framework to a general model of cognition.

## 5 Open questions and future directions

By using neural models to translate sentences into probabilistic programs, the sections above demonstrated how LLMs could extract meaning from—and inference engines could reason about—language describing uncertain situations, relational structures, embodied situations and goal-directed reasoning. However, these vignettes also leave open many questions about how to scale this framework to more complex language, and how to automate the process of building meaning representations for new domains. Together, these questions offer a roadmap for progress on central challenges in modeling language, reasoning, and their interaction, across many sub-fields of artificial intelligence and cognitive science.

### 5.1 Scaling models of rational meaning construction

We begin by describing several of the most important research directions necessary for scaling the framework we have articulated throughout this paper towards a more complete model of integrated cognition and language understanding.

#### 5.1.1 Building new world models on the fly

A key aspect of our proposed architecture is that language is interpreted relative to a probabilistic model of a domain, capturing just enough structure to represent the situation at hand. In [Section 4.2](#), we saw that LLMs could generate these programmatic world models, assuming the model was communicated via a sequence of natural language definitions. But people rarely need such elaborate scene-setting: we can understand language about the world even if no teacher has carefully drawn our attention to the relevant concepts beforehand. A key question is how to model this capability. How do minds craft bespoke world models on the fly, drawing in just enough of our knowledge about the world to answer the questions of interest? How does this process balance competing priorities, such as fidelity to what we know about the world, relevance to the problem at hand, and the efficiency and robustness of inference? These tradeoffs can sometimes seem to evolve during the course of a single chain of human thought. These questions are related to the classic *frame problem* ([McCarthy, 1980](#)) in artificial intelligence and cognitive science, and to recent proposals for addressing it in the setting of causal, probabilistic reasoning ([Icard & Goodman, 2015](#)). These approaches view the problem as one of retrieval: from a vast array of knowledge we have about the world, how can we select just the relevant parts for reasoning about a particular problem? It remains unclear, however, whether the sequences of bespoke models and approximate inferences produced by our minds can be understood as resource-rational approximations to coherent reasoning and planning in some larger, unifying world model, even in principle.

Most probabilistic programming languages were designed for inference in a single, unifying world model ([Bingham et al., 2019](#); [Carpenter et al., 2017](#); [Goodman et al., 2008](#); [Milch et al., 2007](#)) that was written by an external mechanism, not to dynamically explore a sequence of probabilistic programs that are being synthesized, learned, and/or edited on the fly. But some progress in language-level support for dynamic world modeling has already been made. Probabilistic programs in Gen ([M. F. Cusumano-Towner, Saad, Lew, & Mansinghka, 2019](#)) have been used to synthesize and edit other probabilistic programs ([Saad et al., 2019](#); [Witty, Lew, Jensen, & Mansinghka, 2019](#)), and to approximate globally coherent inferences by bridging across sequences of probabilistic programs describing translations among only partially-overlapping worlds ([M. Cusumano-Towner, Bichsel, Gehr, Vechev, & Mansinghka, 2018](#); [M. Cusumano-Towner, Lew, & Mansinghka, 2020](#); [A. K. Lew, Matheos, et al., 2023](#); [V. K. Mansinghka et al., 2018](#)). Analogous language-level support for dynamic abstraction for planning with symbolic world models has also been developed ([Zhi-Xuan, 2022](#)). It remains to be seen to what extent these new degrees of freedom can be exploited by language-to-code models targeting these newer probabilistic programming platforms.

How could the common-sense background knowledge needed for dynamic world model synthesis be represented, even in principle? Modern game engines may provide important clues. They can be reconfigured and scripted to simulate diverse imaginary worlds and narratives, featuring interactions between physical objects and goal directed agents in both realistic and physically impossible environments. They routinely combine simulations of the same environment at multiple levels of detail, making computational tradeoffs that are in some ways analogous to the tradeoffs faced by human thinking. The level of scale, coherence, realism, and computational efficiency that they achieve still vastly outstrips the best multi-modal neural

models. Although some progress is already being made by synthesizing lightweight, probabilistic game engine scripts using language-to-code models (C. E. Zhang, Wong, Grand, & Tenenbaum, 2023), many fundamental challenges remain. Game engines lack crucial affordances for robustly fitting world models to sparse data, simulating rare events, and planning under uncertainty. And despite promising progress in neurally-guided program learning (Ellis et al., 2020), showing that libraries and DSLs can be learned from sparse data, there seems to be a long way to go before we can learn game-engine like rules that are sufficient to robustly model common sense. Flexible synthesis and learning mechanisms that can hope to scale across the vast scope of human thought thus seems to require new ideas that span and integrate probabilistic programming, cognitive architecture, and hierarchical program learning.

### 5.1.2 **Scaling probabilistic inference in dynamically synthesized world models**

A central challenge not addressed by this paper is how to scale probabilistic inference to begin to approach the robustness, speed, efficiency, and flexibility of human thought. Consider that the rejection sampling algorithm used in Sections 3 and 4 requires an exponentially-growing number of proposal attempts as the scenario becomes less likely under the prior. Although many exact inference methods for probabilistic programs are much faster and more reliable, they are too restrictive to support many of the world models in this paper (Gehr, Misailovic, & Vechev, 2016; Gehr, Steffen, & Vechev, 2020; Holtzen, Van den Broeck, & Millstein, 2020; Saad, Rinard, & Mansinghka, 2021; Shan & Ramsey, 2017). And although there are many approaches to generic approximate inference in probabilistic programs, drawing on MCMC (Carpenter et al., 2017; Goodman et al., 2008; Wingate, Stuhlmüller, & Goodman, 2011), sequential Monte Carlo (V. Mansinghka et al., 2014; Tolpin, van de Meent, Yang, & Wood, 2016), and variational methods (Bingham et al., 2019; Hoffman, Blei, Wang, & Paisley, 2013; V. Mansinghka et al., 2014; Ranganath, Gerrish, & Blei, 2014), they all routinely struggle to solve simple problems that can be solved by custom algorithms.

One potential way forward is to explicitly generate models of thinking processes that augment the world models with which they are thinking, by synthesizing inference programs (M. F. Cusumano-Towner et al., 2019; V. K. Mansinghka et al., 2018) tailored to specific problems. For example, Venture’s inference meta-programming language is designed to enable concise specification of sequential inference processes that combine SMC, dynamic programming, MCMC, gradient-based optimization, and variational inference to perform inference in a sequence of world models and queries that grows dynamically. Data-driven proposals for use with these thinking strategies can also be generated in real-time, without any offline learning, using dynamic programming over blocks of highly coupled variables. This approach has recently outperformed machine learning methods on hard common-sense reasoning problems in databases with millions of records (A. Lew, Agrawal, Sontag, & Mansinghka, 2021). Scaling this approach will require not just synthesizing world models but automatically analyzing and decomposing them, analogously to how inference algorithm designers decompose large inference problems into sequences of more tractable subproblems.

Another promising approach is to train neural networks to make data-driven proposals via amortized inference, potentially using synthetic data from an open-ended simulator of world models and queries (M. Wu & Goodman, 2022). This can be seen as an alternative to inference programming, avoiding the need for explicit symbolic analysis of the process of thought. It can also be seen as a potential technique by which inference programs might eventually be synthesized, once a suitable training corpus can be generated synthetically — as well as a source of data-driven proposals that can be recombined by inference programs.

### 5.1.3 **Resource rational amortization in meaning construction and problem solving**

In some of our examples, a sentence (e.g., “Gabe is stronger than Josh”) is translated to a meaning representation that looks very much like its classical formal semantics, composing the literal meanings of each word in the sentence. But in other examples (e.g., “several of the faculty are real slackers”), the translations appear to incorporate complex contextual and pragmatic judgments, judgments that might otherwise have been arrived at via probabilistic inference in a model of speakers, listeners, and their intents (Goodman & Frank, 2016). This raises the question of where to draw the line between translation and inference. Versions of this question have been extensively studied (e.g., does a word like “some” imply “not all” as part of its meaning, or does this implicature arise via after-the-fact pragmatic reasoning (Tessler, Tenenbaum, & Goodman, 2022)?), and some past work has offered a unifying view via theories of *amortized pragmatics* (White et al., 2020), whereby RSA-style inferences are “compiled down” into new word meanings.

A key feature of our architecture is that it is largely *agnostic* to where exactly the boundary should lie, and as such could help to model and extend this process of amortized inference in language understanding. For example, as expanded on below, we could extend our symbolic world models to include aspects of the language understanding process itself (such as those described in symbolic derivations of *semantics* (Heim & Kratzer, 1998; Montague, 1970; Pollard & Sag, 1994; Steedman, 2001, 2011), and those used explicitly to compute its *pragmatic* interpretations (Fox, 2007; Goodman & Frank, 2016)). Symbolic inferences about meanings could then be used to train the language understanding module to directly generate the results of this symbolic inference process—for use either as a fully amortized pragmatic translator, or as a proposal distribution within a larger Monte Carlo algorithm that could score and reject inaccurate translations.

In addition to making aspects of translation symbolic, we could consider approaches to amortizing the more general probabilistic inferences required to answer queries. By supervising “translation models” directly with the final outputs of symbolic inference, across a wide variety of tasks, we could enable a pure neural inference mode for these systems that may overcome some limitations of models trained only on language and code. As described above, such supervised models could also be incorporated as proposal distributions in posterior sampling algorithms, leading to improved efficiency without sacrificing the ability to correct for learned biases that may be inapplicable when tackling novel problems.

Ultimately, we envision a new kind of neurosymbolic model in which, rather than pre-assigning responsibilities to the neural or symbolic program, models may flexibly perform any part of the language understanding via explicit probabilistic inference or learned, amortized prediction, with tradeoffs in speed and accuracy for any allocation of responsibilities to modules. The research question is how to do this automatically—how do we identify pieces of a computation that can reliably be emulated by a neural model, how do we train this neural model efficiently, and how do we decide at runtime which inference mode to use? As above, these questions raise many opportunities to take inspiration from our scientific understanding of the separation of responsibilities in language and thought, and work on learning for inference in more general probabilistic models.

#### 5.1.4 Language generation

The preceding discussion has focused largely in problems of language *understanding*—mapping from utterances to inferences about the state of the world that those utterances describe. But effective models of language use should also be able to explain *generation*, making it possible to translate the results of inference back to language. As with the problem of language-informed thinking that we focus on this paper, it is useful to model language generation as two distinct processes: choosing *what to say*, then *how to say it* (Duboué & McKeown, 2003). And as with understanding, the first phase requires a model of the world, and of the speaker’s goals within it. What additional work is needed to adapt our models of rational meaning construction for generation?

One possibility, alluded to in the discussion of amortization above, is to interpret the language understanding machinery described above as a *model of a listener*, then perform language generation by selecting utterances that cause this model listener to form correct beliefs or take appropriate actions (Fox, 2007; Goodman & Frank, 2016). This extra layer of reasoning introduces major inferential challenges: the generation model must now reason both about the set of possible utterances and the effect of each utterance on the distribution over possible worlds inferred by a listener. Here it is once again possible to leverage large-scale statistical learning—for example, using LLMs to directly translate candidate communicative intentions back to natural language strings, which may then be used as candidate utterances to be scored using a formal model of language understanding. Such a hybrid neuro-symbolic generation model (Fang et al., 2022; Langkilde & Knight, 1998) offers a path towards language generation that is expressive and fluent, but avoids the truthfulness and hallucination problems that plague all purely neural language generation models that exist today (Maynez, Narayan, Bohnet, & McDonald, 2020; Wiseman, Shieber, & Rush, 2017).

## 5.2 Implications for cognitive science

In this section, we describe several research directions for other closely related disciplines that study language and thought in natural minds, brains, and behavior, focusing on productive intersections in relation to this framework.

### 5.2.1 Connections to cognitive and formal models of linguistic structure

In all the examples described above, the process of translating utterances into formal meaning representations was performed with a black-box statistical model, while reasoning about those meaning representations leveraged an explicit symbolic inferential process. However, an enormous body of work in linguistics has argued that the process of mapping from utterances to meaning representations can itself be described (at least approximately) in terms of symbol processing operations (Montague, 1970; Pollard & Sag, 1994; Steedman, 2001, *inter alia*). By design, most of our “meaning representations” are designed to support efficient reasoning about domain-specific world models, and bear only a vague resemblance to formal and domain-general linguistic representational theories. But can the symbolic models of linguistic meaning posited by these theories (as opposed to the symbolic models of reasoning we already draw on) be incorporated into our framework?

As noted in Section 5.1.3, a fully realized model of rational meaning construction should be able to flexibly move computation across the statistical–symbolic boundary, “compiling” results of symbolic inference into amortized computation, or retrieving symbolic descriptions of amortized processes for explicit verification. In this view, the vignettes above treat the meaning representation process as culminating in domain-specific representations and amortized by default. But probabilistic symbolic models of meaning (e.g., Kwiatkowski, Zettlemoyer, Goldwater, & Steedman, 2010), or Bayesian and game-theoretic models of semantics (e.g., Goodman & Frank, 2016) can themselves be implemented as probabilistic programs and composed with domain-specific inferential computations, resulting in an almost purely symbolic (but amortizable) language understanding process similar to the one described by Goodman and Lassiter (2015).

Such a model would also offer an appealing path toward *learning* language in a more sample-efficient (and perhaps human-like) ways. Today’s neural sequence models require orders of magnitude more data than human learners to discover the structural regularities underlying human languages (Linzen, 2020). Explicit probabilistic symbolic models, by contrast, can discover this structure extremely sample-efficiently (Yang & Piantadosi, 2022). A model that could automatically infer symbolic meaning representation rules from data, then amortize this representation system into a statistical translation model (Liang, Daumé III, & Klein, 2008), would be capable of both efficient learning of language, and efficient modeling of other domains using language. It would also offer a framework for modeling other key aspects of language acquisition, including explicit linguistic instruction (of word meanings, rules of grammar, etc), tradeoffs between different formal representational schemes, and the relationship between linguistic *competence* (understood as symbol-side language processing) and linguistic performance (understood as statistical-side processing).

The semantic framework in this paper is most closely related to other *cognitive semantic frameworks* (eg. Jackendoff (1985); Lakoff (1988); Pietroski (2018); Pinker (1984)) that explicitly propose that human language constructs meanings from conceptual and cognitive primitives, including those for causal reasoning, or core knowledge representations of physics and agents. Related information-theoretic proposals have proposed that languages are effectively designed to be efficiently communicable externalizations of underlying thoughts—that the structure of human languages derives from underlying structure in the semantic representations we wish to communicate, and indeed may be driven by environmental and domain-specific pressures (eg. Gibson et al. (2019); Mollica et al. (2021); Zaslavsky, Kemp, Regier, and Tishby (2018)).

Other related acquisition theories posit that these structural relationships between the representations of thought and externalizable language play an important role in language acquisition. Under these theories, humans can so efficiently learn or hypothesize the meanings of sentences because they “map cleanly” onto the cognitive structures already present in the minds of the language learner (Snedeker, 2016); language learning is *bootstrapped* by these predictable, structured mappings between the underlying space of meanings and the syntax of language (L. R. Gleitman et al., 2005; Hartshorne et al., 2016; Pinker & MacWhinney, 1987). In preliminary experiments, we find intriguing evidence that large language-to-code models can extract and generalize syntactic patterns between language and code, including to bootstrap hypotheses about the semantics of novel words expressed as probabilistic programs based on contextual, syntactic usage (see *Syntactic Bootstrapping*, Fig. 14). Future work can explore therefore whether these statistical distributional models might be used to implement cognitive models of bootstrapped language acquisition.



### 5.2.2 Modeling the mechanisms of human thought

Using tools for adaptive Bayesian inference over flexibly structured symbolic representations—including not only probabilistic programs but more generally hierarchical Bayesian models (Griffiths et al., 2010; Tenenbaum, Kemp, Griffiths, & Goodman, 2011), resource-rational modeling (S. J. Gershman et al., 2015; Lieder & Griffiths, 2020), and program induction (Lake et al., 2017; Piantadosi et al., 2012)—computational cognitive scientists have built quantitatively predictive and functionally explanatory models of human behavior in almost every domain of cognition. This range spans from models of perception, concept learning and categorization, causal reasoning, decision-making and planning, to intuitive physics, theory of mind, sentence processing, and cognitive and language development (C. Baker et al., 2011; Goodman & Frank, 2016; Goodman et al., 2014; Griffiths & Tenenbaum, 2006; Ho, Saxe, & Cushman, 2022; Jara-Ettinger et al., 2020; Lake et al., 2017; Perfors et al., 2011). However, in almost every one of these cases, the models are not fully “stimulus-computable”: Behavioral experiments in cognitive psychology almost always use natural language to present participants with some situation for thinking about (in addition to perhaps perceptual stimuli); language is also almost invariably used to pose some question or goal as the end for thinking. Put another way, almost all our behavioral experiments—like so many instances of cognition in the wild—follow the language-informed thinking paradigm of this paper. But our cognitive models traditionally do not; they are created by hand from the modeler’s understanding of the natural language task description, rather than synthesized automatically from the linguistic stimuli presented to participants. To what extent can the rational meaning construction framework presented here reduce the need for computational cognitive scientists to manually create Bayesian models that match the natural-language prompts given to humans in behavioral experiments? Can we build “language-computable” models of human thought, that are much easier to test and vary via large-scale online experiments?

We have already begun to explore these possibilities and shown promising preliminary results in several domains, including to model how language implicates commonsense physical reasoning about linguistic scenes (C. E. Zhang et al., 2023), social reasoning about goal-directed agents (Ying et al., 2023); as well as to test the claim that the LLM-based meaning function we implement in this paper can compute amortized pragmatic judgments of scalar implicatures that accord with human interpretations (Lipkin, Wong, Grand, & Tenenbaum, 2023).

There is also a growing body of research in computational cognitive science showing that salient dynamics of thought, including well-known departures from Bayesian norms, can be explained via Monte Carlo inference approximations that aim to rationally use limited computational resources (Chater et al., 2020; S. J. Gershman et al., 2015; Lieder & Griffiths, 2020; Lieder, Hsu, & Griffiths, 2014; Sanborn & Chater, 2017). In some cases, human inferences seem to rest on just a single, highly approximate sample (Vul, Goodman, Griffiths, & Tenenbaum, 2014), or perhaps just a few of them (Vul & Pashler, 2008). If we extend our proposed architecture for rational meaning construction to incorporate these kinds of Monte Carlo mechanisms, could we build models of language-guided thinking that can be directly compared at a more mechanistic level to human behavior? How will processes of language understanding and reasoning interact mechanistically, and can we build resource-rational approximate inference models that capture this interaction?

### 5.2.3 Language and thought in the brain

Evidence from cognitive neuroscience suggests a number of parallels between the framework we describe in this paper, and how language relates to systems for general cognition in the human brain. Over decades, cognitive neuroscientists have mapped out a series of interconnected areas in the frontal and temporal lobes that are implicated in *human language processing*. This “language network” is activated in both linguistic comprehension (Deniz, Nunez-Elizalde, Huth, & Gallant, 2019; Fedorenko, Hsieh, Nieto-Castañón, Whitfield-Gabrieli, & Kanwisher, 2010; MacSweeney et al., 2002; Regev, Honey, Simony, & Hasson, 2013; T. L. Scott, Gallée, & Fedorenko, 2017) and production (Hu et al., 2021; Menenti, Gierhan, Segaert, & Hagoort, 2011). It is sensitive to regularities in all levels of linguistic structure—from phonology, to words, to phrases and sentences (Blank & Fedorenko, 2017; Lerner, Honey, Silbert, & Hasson, 2011; Silbert, Honey, Simony, Poeppel, & Hasson, 2014; Wilson, Molnar-Szakacs, & Iacoboni, 2008) and is implicated in combinatorial semantic and syntactic processing (Fedorenko, Blank, Siegelman, & Mineroff, 2020; Hu et al., 2021).

Convergent evidence suggests that the language network is distinct from these systems, and that it is *not* activated in more general, non-linguistic cognition. Aphasic individuals with damage to the language network

exhibit impaired language production and comprehension, but retain the ability to solve arithmetic and logic puzzles, reason about causality and social situations, and perform many other non-linguistic tasks (e.g., Basso & Capitani, 1985; Bek, Blades, Siegal, & Varley, 2010; Fedorenko & Varley, 2016; Klessinger, Szczerbinski, & Varley, 2007; Lecours & Joannette, 1980; Luria, Tsvetkova, & Futer, 1965; Varley, 1998). Functional neuroimaging studies provide further evidence that the language network is *not* activated in a variety of non-linguistic tasks including reasoning about arithmetic, logic, actions, or events (Amalric & Dehaene, 2016, 2019; Blank, Kanwisher, & Fedorenko, 2014; Deen, Koldewyn, Kanwisher, & Saxe, 2015; Fedorenko, Behr, & Kanwisher, 2011; Monti, Osherson, Martinez, & Parsons, 2007; Monti, Parsons, & Osherson, 2012; Paunov, Blank, & Fedorenko, 2019; Paunov et al., 2022; Shain, Paunov, Chen, Lipkin, & Fedorenko, 2022).

In tandem, a broader line of cognitive neuroscience work has located non-linguistic networks that *are* activated in processing many of the core cognitive domains we model throughout this paper, including logic, mathematical reasoning (eg. Amalric and Dehaene (2019); Monti et al. (2007)), social reasoning and planning (Adolphs, 2009; Saxe, Moran, Scholz, & Gabrieli, 2006; Saxe & Powell, 2006); and physical reasoning and simulation (Prasad, Cohen, Tenenbaum, & Kanwisher, 2022; Schwettmann, Tenenbaum, & Kanwisher, 2019). More recent work suggests the existence of an “amodal semantics network” (Ivanova, 2022; Ivanova et al., 2021), a network that appears proximal to the language networks activated in processing linguistic structures, interfaces between the language network and the more general multiple demand networks involved in complex non-linguistic cognition, and that appears to be activated specifically in processing semantically meaningful sentences (as opposed to scrambled tokens or syntactically correct but semantically incoherent strings.)

Recently, neuroscientists who study language cognition have begun to draw explicit parallels between the language network and LLMs (see Mahowald et al., 2023, for a review). Several recent studies have observed that smaller LLMs trained specifically on the distributional statistics of language (generally focusing on the GPT-2 model) can predict brain activity in humans processing sentence input (Caucheteux & King, 2022; Goldstein et al., 2022; Schrimpf et al., 2021) and may share representational characteristics of the human language network (Fedorenko et al., 2020; Shain, Blank, van Schijndel, Schuler, & Fedorenko, 2020). These accounts, however, align LLMs with the modular role we propose for neural models in our framework—*not* as end-to-end models of language and reasoning, but instead as robust, context-aware mappings between language and meanings. As a ground for future work, our framework can inform evaluations of LLMs with respect to human language understanding. For instance, our proposal suggests that code-trained LLMs might better capture latent semantic and syntactic structure than language-only LLMs. Ideas from neuroscience, in turn, can help us figure out which kinds of computations can be neurally amortized and where our model’s boundary between language and thought should lie.

## 5.3 Implications for AI

### 5.3.1 Structured hybrid models of language and thought

Growing awareness of the limitations of LLM-based reasoning has motivated several recent proposals for interfacing language models with external symbolic plug-ins or toolkits (Karpas et al., 2022; OpenAI, 2023c; Schick et al., 2023; Wolfram, 2023). At face value, one perspective is to view rational meaning construction as an argument to add probabilistic programs to the growing “swiss army knife” of LLM plug-ins. However, we see this notion as inverted: *thought* should not simply be a plug-in on top of language models. Rather, we believe that future AI systems should be architected around *thought*—general-purpose computing systems that provide a principled framework for expressing world models, conditioning them on observations from sources including language and perceptual input, and drawing principled inferences and decisions with respect to the goals of an intelligent system.<sup>7</sup> As we show throughout this paper, many core domains of cognition can be expressed as forms of probabilistic inference. A probabilistic language of thought, in turn, provides a unifying language for world modeling that can nest calls to other cognitively-motivated modules. In this sense, all of these plug-ins and modules would become plug-ins to the substrate of thought, *including* graphics engines, physics simulators, planning algorithms, *and*, in fact, language models themselves. As we discuss in the future directions of each section, scaling any of our toy implementations towards robust, human-like reasoning and language-understanding systems will almost certainly require more sophisticated implementations of each

<sup>7</sup>A similar argument has been expressed by Stephen Wolfram in a compelling series of writings on integrating ChatGPT with the Wolfram Language and its suit of symbolic computational tools (Wolfram, 2023).

reasoning module. We therefore hope this general probabilistic framework suggests a symbolic substrate that might in turn incorporate many of the specific modules and plug-ins in this recent work.

To this end, another important near-term AI research direction will involve building probabilistic programming frameworks that natively incorporate LLMs. Important steps in this direction are already being taken through work leveraging LLMs to approximate prior probabilities over strings (A. K. Lew, Tessler, Mansinghka, & Tenenbaum, 2020) and amortize complex posterior inferences (M. Wu & Goodman, 2022). Indeed, many popular LLM techniques, such as scratchpads (Nye et al., 2021), chain-of-thought prompting (Wei et al., 2022), selection-inference (Creswell, Shanahan, & Higgins, 2022), STaR (Zelikman, Wu, Mu, & Goodman, 2022), and others can be viewed as implementations of probabilistic programs over string-valued random variables (Dohan et al., 2022). A maturing theoretical understanding of LLMs as probabilistic entities will afford powerful ways of harnessing and controlling generations. For instance, the sequential Monte Carlo (SMC) steering technique introduced under the LLaMPPL framework (A. K. Lew, Zhi-Xuan, Grand, & Mansinghka, 2023) enables concise and tractable specification of infilling, prompt intersection, and other constrained LLM generation tasks as *language model probabilistic programs*. Many of these hybrid models can be viewed as instantiations of rational meaning construction that make resource-motivated tradeoffs between inference in the unstructured space of strings (words) and more structured hypothesis spaces (worlds).

### 5.3.2 Robustness and trustworthiness in language understanding

Recent, high-profile attempts to deploy LLMs in production highlight the fundamental robustness challenges of using these models as the backbone of usable AI systems (Brereton, 2023; Sorkin, Warner, Kessler, Hirsch, & Livni, 2023), even with automated filters and supervised finetuning to human preferences. While LLMs may reasonably appear to condition on input language or answer queries under some circumstances, it is precisely this combination of linguistic fluency and underlying unpredictability that makes them problematic in situations where verifiable, systematic behavior is paramount. LLMs easily produce syntactically convincing but inaccurate “hallucinations” that fabricate facts and inferences (Dziri, Milton, Yu, Zaiane, & Reddy, 2022; Ji et al., 2022), fail to consistently condition on rules and constraints described in natural language, including rules intended to ensure user safety (Edwards, 2023; Zhuo, Huang, Chen, & Xing, 2023), and can generally degrade into nonsensical or highly undesirable language in the vast, easily accessible “long tail” of situations that deviate from their training distribution (Bender, Gebru, McMillan-Major, & Shmitchell, 2021; Roose, 2023; Tangermann, 2023).

The unevenness of today’s LLMs recalls a classic critique of even older neural architectures (Fodor & Pylyshyn, 1988)—that neural models trained on predictive objectives do not produce systematic, logical outputs by design. Similarly, while current or future LLMs may be able in principle to recover the latent representations and algorithms necessary to reason over language—or even successfully approximate them in many settings—they do not *need* to produce systematic results by construction. Rather, they often approximate them with unexpected, undesirable outputs, particularly in out-of-distribution settings.

Even if future LLMs do appear to improve with scale without an external reasoning substrate, engineers may find it desirable to distinguish modularly between external symbolic reasoning engines and language-specific systems to enable separate supervision and verification of each. The framework we present here offers one roadmap for language understanding architectures whose robustness guarantees derive from explicit inference over a structured, editable, and formally constrainable programming language. Inferences themselves, and other formalizable reasoning computations including planning and physical simulation, take place in modules constructed explicitly to perform these calculations.

### 5.3.3 Interpreting models that use language

As with verifiability and robustness, the framework we propose here is an architecture for language understanding systems that are also *inherently interpretable*, or *interpretable by design* (Rudin, 2019; Rudin et al., 2022)—it constructs visible, editable, and constrainable world models and meanings that serve as the formal basis for inference, rather than post-hoc explanations decoded from or produced over hidden internal computations.

However, a fundamental part of our hypothesis is that *any* system that reasons effectively over language should need to—explicitly or implicitly—represent and implement the kinds of computations we formalize throughout this paper. Implementations of this framework might therefore also be useful for model-guided

hypotheses and experiments intended to explain other less transparent language processing systems, both biological (as we suggest in [Section 5.2.3](#)) and artificial. This framework might be incorporated productively into the growing body of work using explicit world models and symbolic languages to formally model the internal computations of deep neural models ([Biggio, Bendinelli, Neitz, Lucchi, & Parascandolo, 2021](#); [Mu & Andreas, 2020](#)) and LLMs specifically ([B. Z. Li, Nye, & Andreas, 2021](#)); as with the related body of work using structured probabilistic models and reasoning engines to interpret human neural activity on social reasoning, physical understanding, and other general inference tasks ([Ho et al., 2022](#); [Schwettmann, Fischer, Tenenbaum, & Kanwisher, 2018](#); [Watters, Tenenbaum, & Jazayeri, 2021](#)). Explaining how LLMs represent the meanings of language, and perform computations with them, is a pressing open question whose scientific interest only increases if LLMs do appear to become more coherent and robust with scale.

In light of this, inspired by our proposed architecture, it may be interesting to probe, or trace, whether end-to-end LLMs construct context-specific world models ([B. Z. Li et al., 2021](#)), maintain belief distributions over uncertain world states ([Hase et al., 2021](#)), and implement reasoning algorithms like probabilistic inference, physical simulation, or planning over these representations.

#### 5.3.4 Learning from human-scale data

Large language models must be trained with many orders of magnitude more language data than any human learner encounters over a lifetime. How can we engineer systems that not only understand language as we do, but also learn from human-scale language data?

Effective, data-efficient language models hold great relevance for both scientific and engineering applications. Complete cognitive models of human language understanding—including models built on the framework we propose here—should account for language acquisition, as well as language use. For engineering purposes, addressing the data-hungry training regime of current LLMs could also address challenges in learning low-resource languages (and the more general problem of accurately learning and deploying the “long tail” of knowledge from statistical distributional data) ([Kandpal, Deng, Roberts, Wallace, & Raffel, 2022](#)), incorporating more expensive input modalities like videos or embodied trajectories ([Ahn et al., 2022](#); [Reed et al., 2022](#)), finetuning on more targeted, task-specific supervision like instruction following ([OpenAI, 2023a](#)), and generally enabling the construction of smaller, more accessible models that can be trained without massive computing resources and prohibitive economic and environmental costs ([Bender et al., 2021](#); [Dickson, 2020](#)). While current “scaling routes” look to improve language understanding by *increasing* data supervision, our hypothesis strongly suggests that this is an expensive, and highly indirect, approach towards learning the representations and inference procedures necessary to reason about language.

Instead, our framework suggests several alternative directions for improving data efficiency. First, perhaps the most direct consequence of this framework is the suggestion that neural models need only play a much tighter, focused role in language understanding systems—as translation models that parse from language into structured symbolic programs for reasoning. Training a translation model focused on parsing from language into probabilistic programs almost certainly requires much less data for effective performance than required to solve the general token prediction problem.

Further, several ideas we discuss in [Section 5.1.1](#) and [Section 5.1.3](#) might also be relevant for training simpler translation models, and using them to bootstrap larger and more complex neural language models.

First, as we discuss in [Section 5.1.3](#), we might consider a progressively amortized avenue for training even complex translation models like the one in our concrete implementation, which appears to contextually amortize certain pragmatic inferences (such as those that adjust vague quantifiers to the context of a particular world model) that could be explicitly computed from a more literal initial semantic parse. One possibility, then, would be to train a more limited, literal semantic parser from language to probabilistic programs, but seek to train neural models that progressively amortize more of these inferences by supervising on its outputs.

Other ideas from human language acquisition might offer more avenues for more radically data efficient learning. Human language learners progress through several phases of language mastery ([R. Brown, 1973](#); [Saffran et al., 2001](#); [Tomasello, 2009](#)), appearing to learn initial but highly imperfect grammars and meaning functions that they refine progressively over time, but much more quickly and with much less data than a comparable LLM trained directly on the distribution of language. Framed as a problem of learning a translation model, however, a more data efficient training regime might also draw inspiration from other methods for learning more flexible translation and semantic parsing distributions. Multiple approaches

have used simpler models to bootstrap more complex ones, either by using simpler models trained on more constrained translation objectives to directly initialize the parameters of more complex ones (P. F. Brown, Della Pietra, Della Pietra, Mercer, et al., 1993; Dong & Lapata, 2018; Petrov, Haghghi, & Klein, 2008), or using simpler grammars as generative data sources to train more complex models, as in general wake-sleep training methods that learn predictive models to amortize the outputs of a generative distribution (Andreas, 2019; Hinton, Dayan, Frey, & Neal, 1995; Jia & Liang, 2016).

Both of these approaches rely, importantly, on the language of thought hypothesis we advance here, which separates the computational problem of learning a translation distribution from the problem of learning the representations and algorithms necessary for general intelligence. This drastically reduces the latent structure and computational complexity we seek to learn from distributional supervision—to learn as efficiently as people, we propose a framework that *begins* with a substrate for thinking and then suggests avenues for amortizing its outputs or refining translation into this substrate, rather than seeking to learn an effective language of thought itself from natural language data.

## 6 Conclusion

Language is central to our cognition. A theory of meaning in human language should explain how language relates to our *thoughts*—how it connects to all our faculties for reasoning, and how it can shift our beliefs across nearly every domain of what we now, change how we act or respond across a broad range of situations, even construct new knowledge that we might later marshal towards yet unspoken questions and goals. This vision lies at the heart of a human theory of language and meaning, but the most expansive visions of AI have also long been ones in which computers share our language, able to meaningfully understand us as we expect to be understood by other people. Today’s large language models have made striking advances towards building this reality in many important regards. For the first time, we have built computer systems that can speak fluently back to us, using many more of our own words than ever before.

Still, much more is needed to capture our own relationship to language. We do not learn language like a large language model does. We think first, and learn from far less input how language maps into our thoughts. Our own world models and beliefs are not the fragile byproduct of what we can glean from language—they are the basis of and core of our cognition, constructed and maintained purposefully towards our intentions and desires. We, of course, are the ones who created the language on which today’s machine learning models are now trained. That language is the product of and reflection of our own goals and questions, and of conceptual systems of our own invention. We continue to think completely new thoughts, and we continue in turn to produce entirely new language, coining new words and even constructing wholly new languages so that we can build its meaning in the minds of other humans. A cognitive theory of human language must capture and explain these aspects of our language and thought. It might in turn form the basis for AI models that reliably and predictably understand us, and that work in ways that we can interpret, explain, and control. This white paper is simply a sketch towards these ends: an outline of the computational components that could relate human language and a substrate for cognition, and one proposal for how this approach might also incorporate today’s language models without requiring them to learn to reliably model the world, draw inferences, or make decisions. We hope it can offer one step towards cognitive and AI models that share the meaning we make from language, and that bridge from language into the vast expanse of our thoughts.

## Acknowledgements

We have many people to thank whose comments, critiques, and feedback have influenced this manuscript and shaped it for the better. Among others, we are grateful to Steve Piantadosi, Jesse Snedeker, Kate Davidson, Ellie Pavlick, Paul Pietroski, Thomas Icard, Luca Bonatti, and Susan Carey for their insightful comments on an early version of this manuscript that was presented at the July 2022 McDonnell Network Workshop; as well as for innumerable helpful comments and feedback on developing versions of this manuscript from Joshua Hartshorne, Judy Fan, Robert Hawkins, Katherine Collins, Anna Ivanova, Cedegao Zhang, Hayley Ross, Anna Ivanova, Benjamin Lipkin, Megan Wei, Jiahai Feng, Xuan Tan, Lance Ying, William McCarthy, Laura Schulz and Tyler Brooke-Wilson. Language from all of these collaborators has invaluable and profoundly informed our thoughts.

The authors gratefully acknowledge support from support from the MIT Quest for Intelligence, AFOSR Grant No. FA9550-19-1-0269, the MIT-IBM Watson AI Lab, the DARPA Machine Common Sense Program, the ONR Science of AI Program, and Siegel Family Endowment. This material is based on work supported by the National Science Foundation Graduate Research Fellowship under Grant No. 1745302 and No. 2141064. Additionally, GG was supported by the MIT Presidential Fellowship, and JDA was supported by NSF Grant IIS-2212310.

## References

- Abend, O., Kwiatkowski, T., Smith, N. J., Goldwater, S., & Steedman, M. (2017). Bootstrapping language acquisition. *Cognition*, *164*, 116–143.
- Adolphs, R. (2009). The social brain: neural basis of social knowledge. *Annual review of psychology*, *60*, 693–716.
- Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., . . . others (2022). Do as I can, not as I say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*.
- Allen, K. R., Smith, K. A., & Tenenbaum, J. B. (2020). Rapid trial-and-error learning with simulation supports flexible tool use and physical reasoning. *Proceedings of the National Academy of Sciences*, *117*(47), 29302–29310.
- Alon, U., Xu, F. F., He, J., Sengupta, S., Roth, D., & Neubig, G. (2022). Neuro-Symbolic Language Modeling with Automaton-augmented Retrieval. *undefined*.
- Amalric, M., & Dehaene, S. (2016, May). Origins of the brain networks for advanced mathematics in expert mathematicians. *Proceedings of the National Academy of Sciences of the United States of America*, *113*(18), 4909–4917. doi: 10.1073/pnas.1603205113
- Amalric, M., & Dehaene, S. (2019, April). A distinct cortical network for mathematical knowledge in the human brain. *NeuroImage*, *189*, 19–31. Retrieved 2019-07-26, from <https://linkinghub.elsevier.com/retrieve/pii/S1053811919300011> doi: 10.1016/j.neuroimage.2019.01.001
- Anderson, J. R. (1990). *The adaptive character of thought*. Psychology Press.
- Andreas, J. (2019). Good-enough compositional data augmentation. *arXiv preprint arXiv:1904.09545*.
- Armeni, I., He, Z.-Y., Gwak, J., Zamir, A. R., Fischer, M., Malik, J., & Savarese, S. (2019). 3d scene graph: A structure for unified semantics, 3d space, and camera. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 5664–5673).
- Artzi, Y., Das, D., & Petrov, S. (2014). Learning compact lexicons for ccg semantic parsing.
- Artzi, Y., Lee, K., & Zettlemoyer, L. (2015, September). Broad-coverage ccg semantic parsing with amr. In *Proceedings of the conference on empirical methods in natural language processing* (pp. 1699–1710). Lisbon, Portugal: Association for Computational Linguistics. Retrieved from <http://aclweb.org/anthology/D15-1198>
- Artzi, Y., & Zettlemoyer, L. (2013). Weakly supervised learning of semantic parsers for mapping instructions to actions. *Transactions of the Association for Computational Linguistics*, *1*(1), 49–62.
- Bai, J., Zhou, L., Blanco, A., Liu, S., Wei, F., Zhou, M., & Li, Z. (2021). Jointly learning to repair code and generate commit message. *ArXiv, abs/2109.12296*.
- Baillargeon, R. (2004). Infants’ physical world. *Current directions in psychological science*, *13*(3), 89–94.
- Baker, C., Saxe, R., & Tenenbaum, J. (2011). Bayesian theory of mind: Modeling joint belief-desire attribution. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 33).

- Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition*, 113(3), 329–349.
- Baker, C. L., Tenenbaum, J. B., & Saxe, R. R. (2007). Goal inference as inverse planning. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 29).
- Bar-Zeev, A. (2003). Scenegraps: Past, present and future. *Último acesso em*, 13.
- Basso, A., & Capitani, E. (1985, May). Spared musical abilities in a conductor with global aphasia and ideomotor apraxia. *Journal of Neurology, Neurosurgery, and Psychiatry*, 48(5), 407–412. Retrieved 2020-08-03, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1028326/>
- Battaglia, P. W., Hamrick, J. B., & Tenenbaum, J. B. (2013). Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences*, 110(45), 18327–18332.
- Bek, J., Blades, M., Siegal, M., & Varley, R. A. (2010, May). Language and spatial reorientation: evidence from severe aphasia. *Journal of Experimental Psychology. Learning, Memory, and Cognition*, 36(3), 646–658. doi: 10.1037/a0018281
- Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 acm conference on fairness, accountability, and transparency* (pp. 610–623).
- Biernaskie, J. M., Walker, S. C., & Gegeer, R. J. (2009). Bumblebees learn to forage like bayesians. *The American Naturalist*, 174(3), 413–423.
- Biggio, L., Bendinelli, T., Neitz, A., Lucchi, A., & Parascandolo, G. (2021). Neural symbolic regression that scales. In *International conference on machine learning* (pp. 936–945).
- Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T., ... Goodman, N. D. (2019). Pyro: Deep universal probabilistic programming. *J. Mach. Learn. Res.*, 20, 28:1–28:6. Retrieved from <http://jmlr.org/papers/v20/18-403.html>
- Blank, I. A., & Fedorenko, E. (2017, October). Domain-General Brain Regions Do Not Track Linguistic Input as Closely as Language-Selective Regions. *Journal of Neuroscience*, 37(41), 9999–10011. Retrieved 2019-11-06, from <https://www.jneurosci.org/content/37/41/9999> doi: 10.1523/JNEUROSCI.3642-16.2017
- Blank, I. A., Kanwisher, N., & Fedorenko, E. (2014, September). A functional dissociation between language and multiple-demand systems revealed in patterns of BOLD signal fluctuations. *Journal of Neurophysiology*, 112(5), 1105–1118. doi: 10.1152/jn.00884.2013
- Block, N. (1998). Conceptual role semantics.
- Bloom, P. (2002). *How children learn the meanings of words*. MIT press.
- Bolton, A. D., Haesemeyer, M., Jordi, J., Schaechtle, U., Saad, F. A., Mansinghka, V. K., ... Engert, F. (2019). Elements of a stochastic 3d prediction engine in larval zebrafish prey capture. *ELife*, 8, e51975.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., ... others (2021). On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Borgeaud, S., Mensch, A., Hoffmann, J., Cai, T., Rutherford, E., Millican, K., ... Sifre, L. (2022, February). *Improving language models by retrieving from trillions of tokens* (No. arXiv:2112.04426). arXiv.
- Bowers, M., Olausson, T. X., Wong, L., Grand, G., Tenenbaum, J. B., Ellis, K., & Solar-Lezama, A. (2023, jan). Top-down synthesis for library learning. *Proc. ACM Program. Lang.*, 7(POPL). Retrieved from <https://doi.org/10.1145/3571234> doi: 10.1145/3571234
- Branwen, G. (2022). *The scaling hypothesis*. Gwern.net.
- Brereton, D. (2023). *Bing ai can't be trusted*. <https://dkb.blog/p/bing-ai-cant-be-trusted>.
- Brooke-Wilson, T. (2023). Why is seeing fast and thinking slow? *in prep*.
- Brown, P. F., Della Pietra, S. A., Della Pietra, V. J., Mercer, R. L., et al. (1993). The mathematics of statistical machine translation: Parameter estimation.
- Brown, R. (1973). A first language: The early stages.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... Amodei, D. (2020, July). Language Models are Few-Shot Learners. *arXiv:2005.14165 [cs]*. Retrieved 2020-08-09, from <http://arxiv.org/abs/2005.14165> (arXiv: 2005.14165)
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., ... others (2023). Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Bybee, J. L. (1985). Morphology. *Typological studies in language*.
- Cai, Q., & Yates, A. (2013). Large-scale semantic parsing via schema matching and lexicon extension. In *Proceedings of the 51st annual meeting of the association for computational linguistics (volume 1: Long*

- papers*) (pp. 423–433).
- Caliskan, A., Bryson, J. J., & Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, *356*(6334), 183–186. Retrieved from <https://www.science.org/doi/abs/10.1126/science.aal4230> doi: 10.1126/science.aal4230
- Carey, S. (1999). Sources of conceptual change. *Conceptual development: Piaget’s legacy*, 293–326.
- Carey, S. (2009). *The origin of concepts*. New York: Oxford University Press.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., . . . Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of statistical software*, *76*(1).
- Caucheteux, C., & King, J.-R. (2022, February). Brains and algorithms partially converge in natural language processing. *Communications Biology*, *5*(1), 1–10. Retrieved 2022-07-05, from <https://www.nature.com/articles/s42003-022-03036-1> (Number: 1 Publisher: Nature Publishing Group) doi: 10.1038/s42003-022-03036-1
- Chakraborty, S., Ding, Y., Allamanis, M., & Ray, B. (2022). Cedit: Code editing with tree-based neural models. *IEEE Transactions on Software Engineering*, *48*, 1385–1399.
- Chakraborty, S., & Ray, B. (2021). On multi-modal learning of editing source code. *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 443–455.
- Chater, N., & Manning, C. D. (2006). Probabilistic models of language processing and acquisition. *Trends in cognitive sciences*, *10*(7), 335–344.
- Chater, N., & Oaksford, M. (1999). Ten years of the rational analysis of cognition. *Trends in cognitive sciences*, *3*(2), 57–65.
- Chater, N., Zhu, J.-Q., Spicer, J., Sundh, J., León-Villagrà, P., & Sanborn, A. (2020). Probabilistic biases meet the bayesian brain. *Current Directions in Psychological Science*, *29*(5), 506–512.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., . . . others (2021). Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Chu, J., & Schulz, L. (2023). In praise of folly: Flexible goals and human cognition.
- Clark, J. H. (1976). Hierarchical geometric models for visible surface algorithms. *Communications of the ACM*, *19*(10), 547–554.
- Cobbe, K., Kosaraju, V., Bavarian, M., Hilton, J., Nakano, R., Hesse, C., & Schulman, J. (2021). Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Collins, K. M., Wong, C., Feng, J., Wei, M., & Tenenbaum, J. B. (2022, May). *Structured, flexible, and robust: Benchmarking and improving large language models towards more human-like behavior in out-of-distribution reasoning tasks* (No. arXiv:2205.05718). arXiv. doi: 10.48550/arXiv.2205.05718
- Colmerauer, A., Kanoui, H., Pasero, R., & Roussel, P. (1972). Un système de communication en français. *Rapport préliminaire de fin de contrat IRIA, Groupe Intelligence Artificielle, Faculté des Sciences de Luminy, Université d’Aix-Marseille II*.
- Colmerauer, A., & Roussel, P. (1996). The birth of prolog. In *History of programming languages—ii* (p. 331–367). New York, NY, USA: Association for Computing Machinery. Retrieved from <https://doi.org/10.1145/234286.1057820>
- Conwell, C., & Ullman, T. D. (2022). Testing relational understanding in text-guided image generation. *arXiv preprint arXiv:2208.00005*.
- Coumans, E., & Bai, Y. (2016). Pybullet, a python module for physics simulation for games, robotics and machine learning.
- Craik, K. J. W. (1967). *The nature of explanation* (Vol. 445). CUP Archive.
- Creswell, A., Shanahan, M., & Higgins, I. (2022, May). *Selection-Inference: Exploiting Large Language Models for Interpretable Logical Reasoning* (No. arXiv:2205.09712). arXiv. doi: 10.48550/arXiv.2205.09712
- Csibra, G. (2008). Goal attribution to inanimate agents by 6.5-month-old infants. *Cognition*, *107*(2), 705–717.
- Csibra, G., Bíró, S., Koós, O., & Gergely, G. (2003). One-year-old infants use teleological representations of actions productively. *Cognitive Science*, *27*(1), 111–133.
- Cusumano-Towner, M., Bichsel, B., Gehr, T., Vechev, M., & Mansinghka, V. K. (2018). Incremental inference for probabilistic programs. In *Proceedings of the 39th acm sigplan conference on programming language design and implementation* (pp. 571–585).
- Cusumano-Towner, M., Lew, A. K., & Mansinghka, V. K. (2020). Automating involutive MCMC using probabilistic and differentiable programming. *arXiv preprint arXiv:2007.09871*.



- Cusumano-Towner, M. F., Radul, A., Wingate, D., & Mansinghka, V. K. (2017). Probabilistic programs for inferring the goals of autonomous agents. *arXiv preprint arXiv:1704.04977*.
- Cusumano-Towner, M. F., Saad, F. A., Lew, A. K., & Mansinghka, V. K. (2019). Gen: a general-purpose probabilistic programming system with programmable inference. In *Proceedings of the 40th acm sigplan conference on programming language design and implementation* (pp. 221–236).
- Dalvi, B., Taffjord, O., & Clark, P. (2022). Towards teachable reasoning systems: Using a dynamic memory of user feedback for continual system improvement. In *Proceedings of the 2022 conference on empirical methods in natural language processing* (pp. 9465–9480).
- Dasgupta, I., & Gershman, S. J. (2021). Memory as a computational resource. *Trends in Cognitive Sciences*, 25(3), 240–251.
- Davidson, D., & Rescher, N. (1967). The logical form of action sentences. 1967, 105–122.
- Davidson, G., Gureckis, T. M., & Lake, B. (2022). Creativity, compositionality, and common sense in human goal generation. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 44).
- de Avila Belbute-Peres, F., Smith, K., Allen, K., Tenenbaum, J., & Kolter, J. Z. (2018). End-to-end differentiable physics for learning and control. *Advances in neural information processing systems*, 31.
- Dechter, E., Malmaud, J., Adams, R. P., & Tenenbaum, J. B. (2013). Bootstrap learning via modular concept discovery. In *Proceedings of the international joint conference on artificial intelligence*.
- Deen, B., Koldewyn, K., Kanwisher, N., & Saxe, R. (2015, November). Functional Organization of Social Perception and Cognition in the Superior Temporal Sulcus. *Cerebral Cortex*, 25(11), 4596–4609. Retrieved 2022-07-05, from <https://doi.org/10.1093/cercor/bhv111> doi: 10.1093/cercor/bhv111
- Deng, F., Zhi, Z., Lee, D., & Ahn, S. (2021). Generative scene graph networks. In *International conference on learning representations*.
- Deniz, F., Nunez-Elizalde, A. O., Huth, A. G., & Gallant, J. L. (2019, September). 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. Retrieved 2020-03-11, from <https://www.jneurosci.org/content/39/39/7722> (Publisher: Society for Neuroscience Section: Research Articles) doi: 10.1523/JNEUROSCI.0675-19.2019
- Dennett, D. C. (2017). *From bacteria to bach and back: The evolution of minds*. WW Norton & Company.
- De Raedt, L., Kimmig, A., & Toivonen, H. (2007). Problog: A probabilistic prolog and its application in link discovery. In *Proceedings of the 20th international joint conference on artificial intelligence* (p. 2468–2473). San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Dickson, B. (2020). The gpt-3 economy. *TechTalks*.
- Ding, Y., Zhang, X., Paxton, C., & Zhang, S. (2023). Task and motion planning with large language models for object rearrangement. *arXiv preprint arXiv:2303.06247*.
- Dohan, D., Xu, W., Lewkowycz, A., Austin, J., Bieber, D., Lopes, R. G., . . . others (2022). Language model cascades. *arXiv preprint arXiv:2207.10342*.
- Dong, L., & Lapata, M. (2018). Coarse-to-fine decoding for neural semantic parsing. *arXiv preprint arXiv:1805.04793*.
- Dries, A., Kimmig, A., Davis, J., Belle, V., & De Raed, L. (2017). Solving probability problems in natural language. In *Proceedings of the 26th international joint conference on artificial intelligence* (p. 3981–3987). AAAI Press.
- Duboue, P. A., & McKeown, K. (2003). Statistical acquisition of content selection rules for natural language generation.
- Dumais, S. T., et al. (2004). Latent semantic analysis. *Annu. Rev. Inf. Sci. Technol.*, 38(1), 188–230.
- Dziri, N., Milton, S., Yu, M., Zaiane, O., & Reddy, S. (2022). On the origin of hallucinations in conversational models: Is it the datasets or the models? *arXiv preprint arXiv:2204.07931*.
- Edgington, D. (1992). Validity, uncertainty and vagueness. *Analysis*, 52(4), 193–204.
- Edgington, D. (1997). Vagueness by degrees.
- Edwards, B. (2023). Ai-powered bing chat spills its secrets via prompt injection attack. *Ars Technica*.
- Elkind, D. (1962). Children’s conceptions of brother and sister: Piaget replication study v. *The Journal of genetic psychology*, 100(1), 129–136.
- Ellis, K., Wong, C., Nye, M., Sable-Meyer, M., Cary, L., Morales, L., . . . Tenenbaum, J. B. (2020). Dreamcoder:

- Growing generalizable, interpretable knowledge with wake-sleep bayesian program learning. *arXiv preprint arXiv:2006.08381*.
- English, G., Nejad, N. G., Sommerfelt, M., Yanik, M. F., & von der Behrens, W. (2023). Bayesian surprise shapes neural responses in somatosensory cortical circuits. *Cell Reports*, 42(2).
- Erez, T., Tassa, Y., & Todorov, E. (2015). Simulation tools for model-based robotics: Comparison of bullet, havok, mujoco, ode and physx. In *2015 IEEE International Conference on Robotics and Automation (ICRA)* (pp. 4397–4404).
- Fang, H., Balakrishnan, A., Jhamtani, H., Bufe, J., Crawford, J., Krishnamurthy, J., ... Klein, D. (2022). The whole truth and nothing but the truth: Faithful and controllable dialogue response generation with dataflow transduction and constrained decoding. *arXiv preprint arXiv:2209.07800*.
- Fedorenko, E., Behr, M. K., & Kanwisher, N. (2011, September). Functional specificity for high-level linguistic processing in the human brain. *Proceedings of the National Academy of Sciences*, 108(39), 16428–16433. Retrieved 2020-02-27, from <https://www.pnas.org/content/108/39/16428> doi: 10.1073/pnas.1112937108
- Fedorenko, E., Blank, I., Siegelman, M., & Mineroff, Z. (2020, February). Lack of selectivity for syntax relative to word meanings throughout the language network. *bioRxiv*, 477851. Retrieved 2020-03-13, from <https://www.biorxiv.org/content/10.1101/477851v2> (Publisher: Cold Spring Harbor Laboratory Section: New Results) doi: 10.1101/477851
- Fedorenko, E., Hsieh, P.-J., Nieto-Castañón, A., Whitfield-Gabrieli, S., & Kanwisher, N. (2010, August). New method for fMRI investigations of language: defining ROIs functionally in individual subjects. *Journal of Neurophysiology*, 104(2), 1177–1194. doi: 10.1152/jn.00032.2010
- Fedorenko, E., & Varley, R. A. (2016, April). Language and thought are not the same thing: evidence from neuroimaging and neurological patients: Language versus thought. *Annals of the New York Academy of Sciences*, 1369(1), 132–153. Retrieved 2019-07-27, from <http://doi.wiley.com/10.1111/nyas.13046> doi: 10.1111/nyas.13046
- Field, H. H. (1977). Logic, meaning, and conceptual role. *The Journal of Philosophy*, 74(7), 379–409.
- Fikes, R. E., & Nilsson, N. J. (1971). Strips: A new approach to the application of theorem proving to problem solving. *Artificial intelligence*, 2(3-4), 189–208.
- Firth, J. (1957). A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis*, 10–32.
- Fodor, J. A. (1975). *The language of thought*. Cambridge, MA: Harvard University Press.
- Fodor, J. A. (1983). *The modularity of mind*. MIT press.
- Fodor, J. A., & Pylyshyn, Z. W. (1988). Connectionism and cognitive architecture: A critical analysis. *Cognition*, 28(1-2), 3–71.
- Fox, D. (2007). Free choice and the theory of scalar implicatures. *Presupposition and implicature in compositional semantics*, 71–120.
- Frank, M. C., Goodman, N. D., & Tenenbaum, J. B. (2009). Using speakers’ referential intentions to model early cross-situational word learning. *Psychological science*, 20(5), 578–585.
- Frege, G. (1892). Über sinn und bedeutung. *Wittgenstein Studien*, 1(1).
- Fried, D., Aghajanyan, A., Lin, J., Wang, S. I., Wallace, E., Shi, F., ... Lewis, M. (2022). Incoder: A generative model for code infilling and synthesis. *ArXiv, abs/2204.05999*.
- Gauthier, J., Levy, R., & Tenenbaum, J. B. (2018). Word learning and the acquisition of syntactic-semantic overhypotheses. *arXiv preprint arXiv:1805.04988*.
- Gehr, T., Misailovic, S., & Vechev, M. (2016). Psi: Exact symbolic inference for probabilistic programs. In *Computer aided verification: 28th international conference, cav 2016, toronto, on, canada, july 17-23, 2016, proceedings, part i 28* (pp. 62–83).
- Gehr, T., Steffen, S., & Vechev, M. (2020). lpsi: Exact inference for higher-order probabilistic programs. In *Proceedings of the 41st ACM SIGPLAN conference on programming language design and implementation* (pp. 883–897).
- Gentner, D., & Goldin-Meadow, S. (2003). Whither whorf. *Language in mind: Advances in the study of language and thought*, 3–14.
- Gentner, D., & Stevens, A. L. (2014). *Mental models*. Psychology Press.
- Gershman, S., & Goodman, N. (2014). Amortized inference in probabilistic reasoning. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 36).
- Gershman, S. J., Horvitz, E. J., & Tenenbaum, J. B. (2015). Computational rationality: A converging

- paradigm for intelligence in brains, minds, and machines. *Science*, 349(6245), 273–278.
- Gerstenberg, T., & Goodman, N. (2012). Ping pong in church: Productive use of concepts in human probabilistic inference. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 34).
- Gibson, E. (2014). Language for communication: Language as rational inference. In *Proceedings of coling 2014, the 25th international conference on computational linguistics: Technical papers* (pp. 781–782).
- Gibson, E., Futrell, R., Piantadosi, S. T., Dautriche, I., Mahowald, K., Bergen, L., & Levy, R. (2019). How efficiency shapes human language. *Trends in cognitive sciences*, 23(5), 389–407.
- Ginsberg, M. L. (1987). Readings in nonmonotonic reasoning.
- Gleitman, L. (1990). The structural sources of verb meanings. *Language acquisition*, 1(1), 3–55.
- Gleitman, L. R., Cassidy, K., Nappa, R., Papafragou, A., & Trueswell, J. C. (2005). Hard words. *Language learning and development*, 1(1), 23–64.
- Goldin-Meadow, S. (2012). 26. homesign: gesture to language. In *Sign language* (pp. 601–625). De Gruyter Mouton.
- Goldstein, A., Zada, Z., Buchnik, E., Schain, M., Price, A., Aubrey, B., ... Hasson, U. (2022, March). Shared computational principles for language processing in humans and deep language models. *Nature Neuroscience*, 25(3), 369–380. Retrieved 2022-10-31, from <https://www.nature.com/articles/s41593-022-01026-4> (Number: 3 Publisher: Nature Publishing Group) doi: 10.1038/s41593-022-01026-4
- Goldwater, S., Griffiths, T. L., & Johnson, M. (2009). A bayesian framework for word segmentation: Exploring the effects of context. *Cognition*, 112(1), 21–54.
- Golovneva, O., Chen, M., Poff, S., Corredor, M., Zettlemoyer, L., Fazel-Zarandi, M., & Celikyilmaz, A. (2022). Roscoe: A suite of metrics for scoring step-by-step reasoning. *arXiv preprint arXiv:2212.07919*.
- Goodman, N. D., & Frank, M. C. (2016). Pragmatic language interpretation as probabilistic inference. *Trends in cognitive sciences*, 20(11), 818–829.
- Goodman, N. D., & Lassiter, D. (2015). Probabilistic semantics and pragmatics: Uncertainty in language and thought. *The handbook of contemporary semantic theory, 2nd edition*. Wiley-Blackwell.
- Goodman, N. D., Mansinghka, V. K., Roy, D. M., Bonawitz, K. A., & Tenenbaum, J. B. (2008). Church: a language for generative models. In D. A. McAllester & P. Myllymäki (Eds.), *UAI 2008, proceedings of the 24th conference in uncertainty in artificial intelligence, helsinki, finland, july 9-12, 2008* (pp. 220–229). AUAI Press. Retrieved from [https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article\\_id=1346&proceeding\\_id=24](https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=1346&proceeding_id=24)
- Goodman, N. D., Tenenbaum, J. B., & Gerstenberg, T. (2014). *Concepts in a probabilistic language of thought* (Tech. Rep.). Center for Brains, Minds and Machines (CBMM).
- Gopnik, A. (1996). The scientist as child. *Philosophy of science*, 63(4), 485–514.
- Gothoskar, N., Cusumano-Towner, M., Zinberg, B., Ghavamizadeh, M., Pollok, F., Garrett, A., ... Mansinghka, V. (2021). 3dp3: 3d scene perception via probabilistic programming. *Advances in Neural Information Processing Systems*, 34, 9600–9612.
- Graff, D. (2000). Shifting sands: An interest-relative theory of vagueness. *Philosophical topics*, 28(1), 45–81.
- Grand, G., Blank, I. A., Pereira, F., & Fedorenko, E. (2022). Semantic projection recovers rich human knowledge of multiple object features from word embeddings. *Nature Human Behaviour*, 1–13.
- Greenberg, M., & Harman, G. (2005). Conceptual role semantics.
- Griffiths, T. L., Chater, N., Kemp, C., Perfors, A., & Tenenbaum, J. B. (2010). Probabilistic models of cognition: Exploring representations and inductive biases. *Trends in cognitive sciences*, 14(8), 357–364.
- Griffiths, T. L., Steyvers, M., & Tenenbaum, J. B. (2007). Topics in semantic representation. *Psychological review*, 114 2, 211–44.
- Griffiths, T. L., & Tenenbaum, J. B. (2006). Optimal predictions in everyday cognition. *Psychological science*, 17(9), 767–773.
- Grimshaw, J. (1981). Form, function, and the language acquisition device. *The logical problem of language acquisition*, 165, 178.
- Harman, G. (1982). Conceptual role semantics. *Notre Dame Journal of Formal Logic*, 23(2), 242–256.
- Harris, Z. S. (1954). Distributional structure. *Word*. Retrieved from <http://psycnet.apa.org/psycinfo/1956-02807-001>
- Hartshorne, J. K., O'Donnell, T. J., Sudo, Y., Uruwashii, M., Lee, M., & Snedeker, J. (2016). Psych verbs, the linking problem, and the acquisition of language. *Cognition*, 157, 268–288.
- Hase, P., Diab, M., Celikyilmaz, A., Li, X., Kozareva, Z., Stoyanov, V., ... Iyer, S. (2021). Do language

- models have beliefs? methods for detecting, updating, and visualizing model beliefs. *arXiv preprint arXiv:2111.13654*.
- Heim, I., & Kratzer, A. (1998). *Semantics in generative grammar* (Vol. 1185). Blackwell Oxford.
- Hespos, S. J., & Baillargeon, R. (2008). Young infants' actions reveal their developing knowledge of support variables: Converging evidence for violation-of-expectation findings. *Cognition*, *107*(1), 304–316.
- Hinton, G. E., Dayan, P., Frey, B. J., & Neal, R. M. (1995). The "wake-sleep" algorithm for unsupervised neural networks. *Science*, *268*(5214), 1158–1161.
- Ho, M. K., Saxe, R., & Cushman, F. (2022). Planning with theory of mind. *Trends in Cognitive Sciences*.
- Hoffman, M. D., Blei, D. M., Wang, C., & Paisley, J. (2013). Stochastic variational inference. *Journal of Machine Learning Research*.
- Holtzen, S., Van den Broeck, G., & Millstein, T. (2020). Scaling exact inference for discrete probabilistic programs. *Proceedings of the ACM on Programming Languages*, *4*(OOPSLA), 1–31.
- Hu, J., Small, H., Kean, H., Takahashi, A., Zekelman, L., Kleinman, D., ... Fedorenko, E. (2021, September). *The language network supports both lexical access and sentence generation during language production* (Tech. Rep.). Retrieved 2021-09-13, from <https://www.biorxiv.org/content/10.1101/2021.09.10.459596v1> (Company: Cold Spring Harbor Laboratory Distributor: Cold Spring Harbor Laboratory Label: Cold Spring Harbor Laboratory Section: New Results Type: article) doi: 10.1101/2021.09.10.459596
- Hughes, N., Chang, Y., & Carlone, L. (2022). Hydra: A real-time spatial perception engine for 3d scene graph construction and optimization. *arXiv preprint arXiv:2201.13360*.
- Icard, T., & Goodman, N. D. (2015). A resource-rational approach to the causal frame problem. In *Cogsci*.
- Isomura, T., Parr, T., & Friston, K. (2019). Bayesian filtering with multiple internal models: toward a theory of social intelligence. *Neural computation*, *31*(12), 2390–2431.
- Ivanova, A. A. (2022). *The role of language in broader human cognition: evidence from neuroscience* (Unpublished doctoral dissertation). Massachusetts Institute of Technology.
- Ivanova, A. A., Mineroff, Z., Zimmerer, V., Kanwisher, N., Varley, R., & Fedorenko, E. (2021). The language network is recruited but not required for nonverbal event semantics. *Neurobiology of Language*, *2*(2), 176–201.
- Izcard, G., Lewis, P., Lomeli, M., Hosseini, L., Petroni, F., Schick, T., ... Grave, E. (2022). Few-shot Learning with Retrieval Augmented Language Models. *undefined*. doi: 10.48550/arXiv.2208.03299
- Jackendoff, R. S. (1985). *Semantics and cognition* (Vol. 8). MIT press.
- Jara-Ettinger, J., Schulz, L. E., & Tenenbaum, J. B. (2020). The naive utility calculus as a unified, quantitative framework for action understanding. *Cognitive Psychology*, *123*, 101334.
- Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., ... Fung, P. (2022). Survey of hallucination in natural language generation. *ACM Computing Surveys*.
- Jia, R., & Liang, P. (2016). Data recombination for neural semantic parsing. *arXiv preprint arXiv:1606.03622*.
- Johnson, J., Hariharan, B., Van Der Maaten, L., Hoffman, J., Fei-Fei, L., Lawrence Zitnick, C., & Girshick, R. (2017). Inferring and executing programs for visual reasoning. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 2989–2998).
- Johnson, J., Krishna, R., Stark, M., Li, L.-J., Shamma, D., Bernstein, M., & Fei-Fei, L. (2015). Image retrieval using scene graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 3668–3678).
- Johnson, R. E., Linderman, S., Panier, T., Wee, C. L., Song, E., Herrera, K. J., ... Engert, F. (2020). Probabilistic models of larval zebrafish behavior reveal structure on many scales. *Current Biology*, *30*(1), 70–82.
- Johnson-Laird, P. N. (1980). Mental models in cognitive science. *Cognitive science*, *4*(1), 71–115.
- Johnson-Laird, P. N. (1989). Mental models.
- Jones, D. (2010, October). Human kinship, from conceptual structure to grammar. *Behavioral and Brain Sciences*, *33*(5), 367–381. Retrieved 2022-08-09, from [https://www.cambridge.org/core/product/identifier/S0140525X10000890/type/journal\\_article](https://www.cambridge.org/core/product/identifier/S0140525X10000890/type/journal_article) doi: 10.1017/S0140525X10000890
- Kaelbling, L. P., & Lozano-Pérez, T. (2011). Hierarchical task and motion planning in the now. In *2011 IEEE International Conference on Robotics and Automation* (pp. 1470–1477).
- Kaelbling, L. P., & Lozano-Pérez, T. (2013). Integrated task and motion planning in belief space. *The International Journal of Robotics Research*, *32*(9-10), 1194–1227.

- Kandpal, N., Deng, H., Roberts, A., Wallace, E., & Raffel, C. (2022). Large language models struggle to learn long-tail knowledge. *arXiv preprint arXiv:2211.08411*.
- Karpas, E., Abend, O., Belinkov, Y., Lenz, B., Lieber, O., Ratner, N., ... Tenenbaum, J. B. (2022, May). *MRKL Systems: A modular, neuro-symbolic architecture that combines large language models, external knowledge sources and discrete reasoning* (No. arXiv:2205.00445). arXiv.
- Katz, Y., Goodman, N. D., Kersting, K., Kemp, C., & Tenenbaum, J. B. (2008). Modeling Semantic Cognition as Logical Dimensionality Reduction. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 30(30), 6.
- Kemp, C., & Regier, T. (2012, May). Kinship Categories Across Languages Reflect General Communicative Principles. *Science*, 336(6084), 1049–1054. Retrieved 2022-08-09, from <https://doi.org/10.1126/science.1218811> (Publisher: American Association for the Advancement of Science) doi: 10.1126/science.1218811
- Kersten, D., Mamassian, P., & Yuille, A. (2004). Object perception as bayesian inference. *Annu. Rev. Psychol.*, 55, 271–304.
- Kersten, D. K. D., & Yuille, A. (1996). Introduction: A bayesian formulation of visual perception. *Perception as Bayesian inference*, 1–21.
- Khalvati, K., Kiani, R., & Rao, R. P. (2021). Bayesian inference with incomplete knowledge explains perceptual confidence and its deviations from accuracy. *Nature communications*, 12(1), 5704.
- Klein, D., & Manning, C. D. (2003). Accurate unlexicalized parsing. In *Proceedings of the 41st annual meeting of the association for computational linguistics* (pp. 423–430).
- Klessinger, N., Szczerbinski, M., & Varley, R. A. (2007, January). Algebra in a man with severe aphasia. *Neuropsychologia*, 45(8), 1642–1648. Retrieved 2022-06-15, from <https://www.sciencedirect.com/science/article/pii/S0028393207000280> doi: 10.1016/j.neuropsychologia.2007.01.005
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., & Iwasawa, Y. (2022). Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*.
- Krafft, P., Baker, C., Pentland, A., & Tenenbaum, J. (2016). Modeling human ad hoc coordination. In *Proceedings of the aaai conference on artificial intelligence* (Vol. 30).
- Kulkarni, T. D., Kohli, P., Tenenbaum, J. B., & Mansinghka, V. (2015). Picture: A probabilistic programming language for scene perception. In *Proceedings of the ieee conference on computer vision and pattern recognition* (pp. 4390–4399).
- Kwiatkowska, T., Zettlemoyer, L., Goldwater, S., & Steedman, M. (2010). Inducing probabilistic ccg grammars from logical form with higher-order unification. In *Proceedings of the 2010 conference on empirical methods in natural language processing* (pp. 1223–1233).
- Kwiatkowski, T., Zettlemoyer, L., Goldwater, S., & Steedman, M. (2011). Lexical generalization in ccg grammar induction for semantic parsing. In *Proceedings of the 2011 conference on empirical methods in natural language processing* (pp. 1512–1523).
- Lake, B. M., Ullman, T. D., Tenenbaum, J. B., & Gershman, S. J. (2017). Building machines that learn and think like people. *Behavioral and brain sciences*, 40.
- Lakoff, G. (1988). Cognitive semantics.
- Landauer, T. K., & Dumais, S. T. (1997). A solution to plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological review*, 104(2), 211.
- Langkilde, I., & Knight, K. (1998). Generation that exploits corpus-based statistical knowledge. In *Coling 1998 volume 1: The 17th international conference on computational linguistics*.
- Lassiter, D., & Goodman, N. D. (2017). Adjectival vagueness in a bayesian model of interpretation. *Synthese*, 194(10), 3801–3836.
- Le, T. A., Baydin, A. G., & Wood, F. (2017). Inference compilation and universal probabilistic programming. In *Artificial intelligence and statistics* (pp. 1338–1348).
- Lecours, A. R., & Joanette, Y. (1980, May). Linguistic and other psychological aspects of paroxysmal aphasia. *Brain and Language*, 10(1), 1–23. doi: 10.1016/0093-934x(80)90034-6
- Lee, T. S., & Mumford, D. (2003). Hierarchical bayesian inference in the visual cortex. *JOSA A*, 20(7), 1434–1448.
- Lerner, Y., Honey, C. J., Silbert, L. J., & Hasson, U. (2011, February). Topographic Mapping of a Hierarchy of Temporal Receptive Windows Using a Narrated Story. *The Journal of Neuroscience*, 31(8), 2906–2915. Retrieved 2019-12-28, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3089381/> doi:

- 10.1523/JNEUROSCI.3684-10.2011
- Levin, B. (1993). *English verb classes and alternations: A preliminary investigation*. University of Chicago press.
- Lew, A., Agrawal, M., Sontag, D., & Mansinghka, V. (2021). Pclean: Bayesian data cleaning at scale with domain-specific probabilistic programming. In *International conference on artificial intelligence and statistics* (pp. 1927–1935).
- Lew, A. K., Matheos, G., Zhi-Xuan, T., Ghavamizadeh, M., Gothoskar, N., Russell, S., & Mansinghka, V. K. (2023). Smcp3: Sequential monte carlo with probabilistic program proposals. In *International conference on artificial intelligence and statistics* (pp. 7061–7088).
- Lew, A. K., Tessler, M. H., Mansinghka, V. K., & Tenenbaum, J. B. (2020). Leveraging unstructured statistical knowledge in a probabilistic language of thought. In *Proceedings of the annual conference of the cognitive science society*.
- Lew, A. K., Zhi-Xuan, T., Grand, G., & Mansinghka, V. K. (2023). Sequential monte carlo steering of large language models using probabilistic programs. *arXiv preprint arXiv:2306.03081*.
- Lewis, D. (1976). General semantics. In *Montague grammar* (pp. 1–50). Elsevier.
- Li, B. Z., Nye, M., & Andreas, J. (2021). Implicit representations of meaning in neural language models. *arXiv preprint arXiv:2106.00737*.
- Li, Y., Wang, S., & Nguyen, T. N. (2020). Dlfix: Context-based code transformation learning for automated program repair. *2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE)*, 602–614.
- Liang, P. (2016). Learning executable semantic parsers for natural language understanding. *Communications of the ACM*, 59(9), 68–76.
- Liang, P., Daumé III, H., & Klein, D. (2008). Structure compilation: trading structure for features. In *Proceedings of the 25th international conference on machine learning* (pp. 592–599).
- Lieder, F., & Griffiths, T. L. (2019). Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behavioral and Brain Sciences*, 43.
- Lieder, F., & Griffiths, T. L. (2020). Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behavioral and brain sciences*, 43, e1.
- Lieder, F., Hsu, M., & Griffiths, T. L. (2014). The high availability of extreme events serves resource-rational decision-making. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 36).
- Linzen, T. (2020). How can we accelerate progress towards human-like linguistic generalization? *arXiv preprint arXiv:2005.00955*.
- Lipkin, B., Wong, L., Grand, G., & Tenenbaum, J. B. (2023). Evaluating statistical language models as pragmatic reasoners. *arXiv preprint arXiv:2305.01020*.
- Liu, B., Jiang, Y., Zhang, X., Liu, Q., Zhang, S., Biswas, J., & Stone, P. (2023). Llm+ p: Empowering large language models with optimal planning proficiency. *arXiv preprint arXiv:2304.11477*.
- Liu, H., Ning, R., Teng, Z., Liu, J., Zhou, Q., & Zhang, Y. (2023). Evaluating the logical reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*.
- Liu, R., Wei, J., Gu, S. S., Wu, T.-Y., Vosoughi, S., Cui, C., ... Dai, A. M. (2022). Mind’s eye: Grounded language model reasoning through simulation. *arXiv preprint arXiv:2210.05359*.
- Lowie, R. H. (1930). The kinship terminology of the bannock indians. *American Anthropologist*, 32(2), 294–299.
- Luria, A. R., Tsvetkova, L. S., & Futer, D. S. (1965, June). Aphasia in a composer (V. G. Shebalin). *Journal of the Neurological Sciences*, 2(3), 288–292. doi: 10.1016/0022-510x(65)90113-9
- Lyu, Q., Havaldar, S., Stein, A., Zhang, L., Rao, D., Wong, E., ... Callison-Burch, C. (2023). Faithful chain-of-thought reasoning. *arXiv preprint arXiv:2301.13379*.
- MacSweeney, M., Woll, B., Campbell, R., McGuire, P. K., David, A. S., Williams, S. C. R., ... Brammer, M. J. (2002, July). Neural systems underlying British Sign Language and audio-visual English processing in native users. *Brain*, 125(7), 1583–1593. Retrieved 2021-01-05, from <https://doi.org/10.1093/brain/awf153> doi: 10.1093/brain/awf153
- Mahowald, K., Ivanova, A. A., Blank, I. A., Kanwisher, N., Tenenbaum, J. B., & Fedorenko, E. (2023). Dissociating language and thought in large language models: a cognitive perspective. *arXiv preprint arXiv:2301.06627*.
- Mansinghka, V., Selsam, D., & Perov, Y. (2014). Venture: a higher-order probabilistic programming platform

- with programmable inference. *arXiv preprint arXiv:1404.0099*.
- Mansinghka, V. K., Kulkarni, T. D., Perov, Y. N., & Tenenbaum, J. (2013). Approximate bayesian image interpretation using generative probabilistic graphics programs. *Advances in Neural Information Processing Systems*, 26.
- Mansinghka, V. K., Schaechtle, U., Handa, S., Radul, A., Chen, Y., & Rinard, M. (2018). Probabilistic programming with programmable inference. In *Proceedings of the 39th acm sigplan conference on programming language design and implementation* (pp. 603–616).
- Marcus, G., Davis, E., & Aaronson, S. (2022). A very preliminary analysis of dall-e 2. *arXiv preprint arXiv:2204.13807*.
- Maynez, J., Narayan, S., Bohnet, B., & McDonald, R. (2020). On faithfulness and factuality in abstractive summarization. *arXiv preprint arXiv:2005.00661*.
- McCarthy, J. (1980). Circumscription—a form of non-monotonic reasoning. *Artificial intelligence*, 13(1-2), 27–39.
- McDermott, D. (1982). A temporal logic for reasoning about processes and plans. *Cognitive science*, 6(2), 101–155.
- McDermott, D. M. (2000). The 1998 ai planning systems competition. *AI magazine*, 21(2), 35–35.
- Menenti, L., Gierhan, S. M. E., Segaert, K., & Hagoort, P. (2011, September). Shared language: overlap and segregation of the neuronal infrastructure for speaking and listening revealed by functional MRI. *Psychological Science*, 22(9), 1173–1182. doi: 10.1177/0956797611418347
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems* (pp. 3111–3119).
- Milch, B., Marthi, B., Russell, S., Sontag, D., Ong, D. L., & Kolobov, A. (2007). BLOG: Probabilistic models with unknown objects. *Statistical relational learning*, 373.
- Mitchell, A., & Jordan, F. M. (2021, June). The Ontogeny of Kinship Categorization. *Journal of Cognition and Culture*, 21(1-2), 152–177. Retrieved 2022-08-09, from [https://brill.com/view/journals/jocc/21/1-2/article-p152\\_8.xml](https://brill.com/view/journals/jocc/21/1-2/article-p152_8.xml) (Publisher: Brill) doi: 10.1163/15685373-12340101
- Mollica, F., Bacon, G., Zaslavsky, N., Xu, Y., Regier, T., & Kemp, C. (2021). The forms and meanings of grammatical markers support efficient communication. *Proceedings of the National Academy of Sciences*, 118(49), e2025993118.
- Mollica, F., & Piantadosi, S. T. (2022, June). Logical word learning: The case of kinship. *Psychonomic Bulletin & Review*, 29(3), 766–799. Retrieved 2022-08-09, from <https://doi.org/10.3758/s13423-021-02017-5> doi: 10.3758/s13423-021-02017-5
- Montague, R. (1970). English as a formal language.
- Monti, M. M., Osherson, D. N., Martinez, M. J., & Parsons, L. M. (2007, September). Functional neuroanatomy of deductive inference: A language-independent distributed network. *NeuroImage*, 37(3), 1005–1016. Retrieved 2020-04-16, from <http://www.sciencedirect.com/science/article/pii/S1053811907003436> doi: 10.1016/j.neuroimage.2007.04.069
- Monti, M. M., Parsons, L. M., & Osherson, D. N. (2012, August). Thought beyond language: neural dissociation of algebra and natural language. *Psychological Science*, 23(8), 914–922. doi: 10.1177/0956797612437427
- Morgan, M. S. (1999). Learning from models. *Ideas in Context*, 52, 347–388.
- Mu, J., & Andreas, J. (2020). Compositional explanations of neurons. *Advances in Neural Information Processing Systems*, 33, 17153–17163.
- Nersessian, N. J., et al. (2010). Mental modeling in conceptual change. *International Journal on Humanistic Ideology*, 3(01), 11–48.
- Nye, M., Andreassen, A. J., Gur-Ari, G., Michalewski, H., Austin, J., Bieber, D., ... others (2021). Show your work: Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*.
- Oaksford, M., & Chater, N. (2007). *Bayesian rationality: The probabilistic approach to human reasoning*. Oxford University Press.
- OpenAI. (2023a). Chatgpt: Optimizing language models for dialogue. *OpenAI Blog*.
- OpenAI. (2023b). *Chatgpt plugins*. OpenAI Blog.
- OpenAI. (2023c). *Gpt-4 technical report*.

- Osgood, C. E. (1952). The nature and measurement of meaning. *Psychological bulletin*, 49(3), 197.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., ... Lowe, R. (2022). *Training language models to follow instructions with human feedback*. arXiv. Retrieved from <https://arxiv.org/abs/2203.02155> doi: 10.48550/ARXIV.2203.02155
- Pan, L., Albalak, A., Wang, X., & Wang, W. Y. (2023). Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. *arXiv preprint arXiv:2305.12295*.
- Panthaplackel, S., Nie, P., Gligoric, M., Li, J. J., & Mooney, R. J. (2020). Learning to update natural language comments based on code changes. *arXiv preprint arXiv:2004.12169*.
- Parsons, T. (1990). Events in the semantics of english: A study in subatomic semantics.
- Paunov, A. M., Blank, I. A., & Fedorenko, E. (2019, April). Functionally distinct language and Theory of Mind networks are synchronized at rest and during language comprehension. *Journal of Neurophysiology*, 121(4), 1244–1265. Retrieved 2019-07-10, from <https://www.physiology.org/doi/10.1152/jn.00619.2018> doi: 10.1152/jn.00619.2018
- Paunov, A. M., Blank, I. A., Jouravlev, O., Mineroff, Z., Gallée, J., & Fedorenko, E. (2022, June). Differential Tracking of Linguistic vs. Mental State Content in Naturalistic Stimuli by Language and Theory of Mind (ToM) Brain Networks. *Neurobiology of Language*, 1–29. Retrieved 2022-07-05, from [https://doi.org/10.1162/nol\\_a\\_00071](https://doi.org/10.1162/nol_a_00071) doi: 10.1162/nol\_a\_00071
- Pearl, J. (1988). *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan kaufmann.
- Pearl, J., et al. (2000). Models, reasoning and inference. *Cambridge, UK: CambridgeUniversityPress*, 19(2).
- Pednault, E. P. (1989). Adl: exploring the middle ground between strips and the situation calculus. *Kr*, 89, 324–332.
- Pereira, F. C., & Shieber, S. M. (2002). *Prolog and natural-language analysis*. Microtome Publishing.
- Perfors, A., Tenenbaum, J. B., & Regier, T. (2011). The learnability of abstract syntactic principles. *Cognition*, 118(3), 306–338.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (1802). Deep contextualized word representations. corr abs/1802.05365 (2018). *arXiv preprint arXiv:1802.05365*.
- Petrov, S., Haghighi, A., & Klein, D. (2008). Coarse-to-fine syntactic machine translation using language projections. In *Proceedings of the 2008 conference on empirical methods in natural language processing* (pp. 108–116).
- Philippe, R. (1972). *Définition et traitement de l'égalité formelle en démonstration automatique* (Unpublished doctoral dissertation). thèse de 3ième cycle, Groupe Intelligence Artificielle, Faculté des Sciences ...
- Piaget, J. (1951). *Judgement and reasoning in the child*. London: Routledge and Kegan Paul.
- Piantadosi, S. T. (2023). Modern language models refute chomsky's approach to language. *Lingbuzz Preprint*, *lingbuzz*, 7180.
- Piantadosi, S. T., Tenenbaum, J. B., & Goodman, N. D. (2012). Bootstrapping in a language of thought: A formal model of numerical concept learning. *Cognition*, 123(2), 199–217.
- Pietroski, P. M. (2018). *Conjoining meanings: Semantics without truth values*. Oxford University Press.
- Pinker, S. (1984). Language learnability and language development.
- Pinker, S. (1998). Words and rules. *Lingua*, 106(1-4), 219–242.
- Pinker, S., & MacWhinney, B. (1987). The bootstrapping problem in language acquisition. *Mechanisms of language acquisition*, 399–441.
- Pollard, C., & Sag, I. A. (1994). *Head-driven phrase structure grammar*. University of Chicago Press.
- Posner, G. J., Strike, K. A., Hewson, P. W., & Gertzog, W. (1982). Toward a theory of conceptual change. *Science education*, 66(2), 211–227.
- Pramod, R., Cohen, M. A., Tenenbaum, J. B., & Kanwisher, N. (2022). Invariant representation of physical stability in the human brain. *Elife*, 11, e71736.
- Pyers, J. E., Shusterman, A., Senghas, A., Spelke, E. S., & Emmorey, K. (2010). Evidence from an emerging sign language reveals that language supports spatial cognition. *Proceedings of the National Academy of Sciences*, 107(27), 12116–12120.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8).
- Rae, J. W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., ... others (2021). Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.



- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., & Chen, M. (2022). Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*.
- Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., ... Sutskever, I. (2021). Zero-shot text-to-image generation. In *International conference on machine learning* (pp. 8821–8831).
- Ranganath, R., Gerrish, S., & Blei, D. (2014). Black box variational inference. In *Artificial intelligence and statistics* (pp. 814–822).
- Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S. G., Novikov, A., Barth-Maron, G., ... others (2022). A generalist agent. *arXiv preprint arXiv:2205.06175*.
- Regev, M., Honey, C. J., Simony, E., & Hasson, U. (2013, October). Selective and Invariant Neural Responses to Spoken and Written Narratives. *Journal of Neuroscience*, *33*(40), 15978–15988. Retrieved 2020-10-02, from <https://www.jneurosci.org/content/33/40/15978> (Publisher: Society for Neuroscience Section: Articles) doi: 10.1523/JNEUROSCI.1580-13.2013
- Reid, M., & Neubig, G. (2022). Learning to model editing processes. *ArXiv, abs/2205.12374*.
- Ribeiro, D., Wang, S., Ma, X., Zhu, H., Dong, R., Kong, D., ... others (2023). Street: A multi-task structured reasoning and explanation benchmark. *arXiv preprint arXiv:2302.06729*.
- Rips, L. J., & Hespos, S. J. (2015). Divisions of the physical world: Concepts of objects and substances. *Psychological bulletin*, *141*(4), 786.
- Roose, K. (2023). Bing’s a.i. chat: I want to be alive. *The New York Times*.
- Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence*, *1*(5), 206–215.
- Rudin, C., Chen, C., Chen, Z., Huang, H., Semenova, L., & Zhong, C. (2022). Interpretable machine learning: Fundamental principles and 10 grand challenges. *Statistic Surveys*, *16*, 1–85.
- Russell, S., & Norvig, P. (2021). *Artificial intelligence : a modern approach* (Fourth edition. ed.). Hoboken, NJ: Pearson.
- Saad, F. A., Cusumano-Towner, M. F., Schaechtle, U., Rinard, M. C., & Mansinghka, V. K. (2019). Bayesian synthesis of probabilistic programs for automatic data modeling. *Proceedings of the ACM on Programming Languages*, *3*(POPL), 1–32.
- Saad, F. A., Rinard, M. C., & Mansinghka, V. K. (2021). Sppl: probabilistic programming with fast exact symbolic inference. In *Proceedings of the 42nd acm sigplan international conference on programming language design and implementation* (pp. 804–819).
- Saffran, J. R., Senghas, A., & Trueswell, J. C. (2001). The acquisition of language by children. *Proceedings of the National Academy of Sciences*, *98*(23), 12874–12875.
- Sahlgren, M. (2008). The distributional hypothesis. *Italian Journal of Disability Studies*, *20*, 33–53.
- Sanborn, A. N., & Chater, N. (2017). The sampling brain. *Trends in Cognitive Sciences*, *21*(7), 492–493.
- Sapir, E. (1929). The status of linguistics as a science. *Language*, 207–214.
- Saxe, R., Moran, J. M., Scholz, J., & Gabrieli, J. (2006). Overlapping and non-overlapping brain regions for theory of mind and self reflection in individual subjects. *Social cognitive and affective neuroscience*, *1*(3), 229–234.
- Saxe, R., & Powell, L. J. (2006). It’s the thought that counts: specific brain regions for one component of theory of mind. *Psychological science*, *17*(8), 692–699.
- Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Zettlemoyer, L., ... Scialom, T. (2023). Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*.
- Schrimpf, M., Blank, I. A., Tuckute, G., Kauf, C., Hosseini, E. A., Kanwisher, N., ... Fedorenko, E. (2021, November). The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the National Academy of Sciences*, *118*(45). Retrieved 2021-12-12, from <https://www.pnas.org/content/118/45/e2105646118> (Publisher: National Academy of Sciences Section: Biological Sciences) doi: 10.1073/pnas.2105646118
- Schuler, K. K. (2005). *Verbnet: A broad-coverage, comprehensive verb lexicon* [PhD Thesis]. University of Pennsylvania. Retrieved from <http://verbs.colorado.edu/~kipper/Papers/dissertation.pdf> (ISBN: 0-542-20049-X)
- Schwettmann, S., Fischer, J., Tenenbaum, J., & Kanwisher, N. (2018). Evidence for an intuitive physics engine in the human brain. In *Cogsci*.
- Schwettmann, S., Tenenbaum, J. B., & Kanwisher, N. (2019). Invariant representations of mass in the human brain. *Elife*, *8*, e46619.

- Scott, R. M., & Baillargeon, R. (2013). Do infants really expect agents to act efficiently? a critical test of the rationality principle. *Psychological science*, *24*(4), 466–474.
- Scott, T. L., Gallée, J., & Fedorenko, E. (2017). A new fun and robust version of an fMRI localizer for the frontotemporal language system. *Cognitive Neuroscience*, *8*(3), 167–176. doi: 10.1080/17588928.2016.1201466
- Seaman, I. R., van de Meent, J.-W., & Wingate, D. (2018). Nested reasoning about autonomous agents using probabilistic programs. *arXiv preprint arXiv:1812.01569*.
- Senghas, A., Kita, S., & Ozyurek, A. (2004). Children creating core properties of language: Evidence from an emerging sign language in nicaragua. *Science*, *305*(5691), 1779–1782.
- Shain, C., Blank, I. A., van Schijndel, M., Schuler, W., & Fedorenko, E. (2020). fMRI reveals language-specific predictive coding during naturalistic sentence comprehension. *Neuropsychologia*, *138*, 107307. doi: 10.1016/j.neuropsychologia.2019.107307
- Shain, C., Paunov, A. M., Chen, X., Lipkin, B., & Fedorenko, E. (2022, July). *No evidence of theory of mind reasoning in the human language network*. bioRxiv. Retrieved 2022-07-20, from <https://www.biorxiv.org/content/10.1101/2022.07.18.500516v1> (Pages: 2022.07.18.500516 Section: New Results) doi: 10.1101/2022.07.18.500516
- Shan, C.-c., & Ramsey, N. (2017). Exact bayesian inference by symbolic disintegration. In *Proceedings of the 44th acm sigplan symposium on principles of programming languages* (pp. 130–144).
- Shin, R., Brockschmidt, M., Allamanis, M., & Polozov, O. (2018). Program synthesis with learned code idioms.
- Silbert, L. J., Honey, C. J., Simony, E., Poeppel, D., & Hasson, U. (2014, October). Coupled neural systems underlie the production and comprehension of naturalistic narrative speech. *Proceedings of the National Academy of Sciences*, *111*(43), E4687–E4696. Retrieved 2021-09-06, from <https://www.pnas.org/content/111/43/E4687> (Publisher: National Academy of Sciences Section: PNAS Plus) doi: 10.1073/pnas.1323812111
- Smith, K., Frank, S., Rolando, S., Kirby, S., & Loy, J. E. (2020). Simple kinship systems are more learnable. *Proceedings of the Annual Meeting of the Cognitive Science Society*, *7*.
- Smith, L., & Yu, C. (2008). Infants rapidly learn word-referent mappings via cross-situational statistics. *Cognition*, *106*(3), 1558–1568.
- Snedeker, J. (2016). *Clean mapping: A sketchy story about how conceptual structure could shape language acquisition and some evidence suggesting that it just might be true*.
- Sorkin, A. R., Warner, B., Kessler, S., Hirsch, L., & Livni, E. (2023). Revenge of the chatbots. *The New York Times*.
- Spelke, E. S. (1990). Principles of object perception. *Cognitive science*, *14*(1), 29–56.
- Spelke, E. S. (2022). *What babies know: Core knowledge and composition volume 1* (Vol. 1). Oxford University Press.
- Spelke, E. S., Gutheil, G., & Van de Walle, G. (1995). The development of object perception. *Visual cognition: An invitation to cognitive science*, *2*, 297–330.
- Spelke, E. S., & Kinzler, K. D. (2007). Core knowledge. *Developmental science*, *10*(1), 89–96.
- Steedman, M. (2001). *The syntactic process*. MIT press.
- Steedman, M. (2011). Combinatory categorial grammar.
- Sumers, T. R., Hawkins, R. D., Ho, M. K., Griffiths, T. L., & Hadfield-Menell, D. (2022). How to talk so your robot will learn: Instructions, descriptions, and pragmatics. *arXiv preprint arXiv:2206.07870*.
- Suster, S., Fizez, P., Totis, P., Kimmig, A., Davis, J., De Raedt, L., & Daelemans, W. (2021). Mapping probability word problems to executable representations. In *Proceedings of the 2021 conference on empirical methods in natural language processing* (pp. 3627–3640).
- Talmy, L. (1988). Force dynamics in language and cognition. *Cognitive science*, *12*(1), 49–100.
- Tangermann, V. (2023). Microsoft’s bing ai is leaking maniac alternate personalities named venom and fury. *Futurism*.
- Téglás, E., Vul, E., Giroto, V., Gonzalez, M., Tenenbaum, J. B., & Bonatti, L. L. (2011). Pure Reasoning in 12-Month-Old Infants as Probabilistic Inference. *Science*, *27*(332), 1054–1059.
- Tellex, S., Kollar, T., Dickerson, S., Walter, M., Banerjee, A., Teller, S., & Roy, N. (2011). Understanding natural language commands for robotic navigation and mobile manipulation. In *Proceedings of the aaai conference on artificial intelligence* (Vol. 25, pp. 1507–1514).

- Tenenbaum, J. B., Kemp, C., Griffiths, T. L., & Goodman, N. D. (2011). How to grow a mind: Statistics, structure, and abstraction. *science*, *331*(6022), 1279–1285.
- Tenney, I., Das, D., & Pavlick, E. (2019). Bert rediscovers the classical nlp pipeline. *arXiv preprint arXiv:1905.05950*.
- Tessler, M. H., Tenenbaum, J. B., & Goodman, N. D. (2022). Logic, probability, and pragmatics in syllogistic reasoning. *Topics in Cognitive Science*, *14*(3), 574–601.
- Thoppilan, R., De Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.-T., ... Le, Q. (2022, February). *LaMDA: Language Models for Dialog Applications* (No. arXiv:2201.08239). arXiv.
- Todorov, E., Erez, T., & Tassa, Y. (2012). Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems* (pp. 5026–5033).
- Tolpin, D., van de Meent, J.-W., Yang, H., & Wood, F. (2016). Design and implementation of probabilistic programming language Anglican. In *Proceedings of the 28th symposium on the implementation and application of functional programming languages* (pp. 1–12).
- Tomasello, M. (2009). The usage-based theory of language acquisition. In *The Cambridge handbook of child language* (pp. 69–87). Cambridge Univ. Press.
- Tomasello, M. (2022). *The evolution of agency: Behavioral organization from lizards to humans*. MIT Press.
- Ullman, T. D. (2023). Large language models fail on trivial alterations to theory-of-mind tasks. *arXiv preprint arXiv:2302.08399*.
- Ullman, T. D., Spelke, E., Battaglia, P., & Tenenbaum, J. B. (2017). Mind games: Game engines as an architecture for intuitive physics. *Trends in cognitive sciences*, *21*(9), 649–665.
- Valmeekam, K., Sreedharan, S., Marquez, M., Olmo, A., & Kambhampati, S. (2023). On the planning abilities of large language models (a critical investigation with a proposed benchmark). *arXiv preprint arXiv:2302.06706*.
- Varley, R. A. (1998). Aphasic language, aphasic thought: an investigation of propositional thinking in an a-propositional aphasic. In P. Carruthers & J. Boucher (Eds.), *Language and Thought: Interdisciplinary Themes* (pp. 128–145). Cambridge University Press. doi: 10.1017/CBO9780511597909.009
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, *30*.
- Vul, E., Goodman, N., Griffiths, T. L., & Tenenbaum, J. B. (2014). One and done? optimal decisions from very few samples. *Cognitive science*, *38*(4), 599–637.
- Vul, E., & Pashler, H. (2008). Measuring the crowd within: Probabilistic representations within individuals. *Psychological Science*, *19*(7), 645–647.
- Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., ... Anandkumar, A. (2023). Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*.
- Wang, R. F., & Spelke, E. S. (2002). Human spatial representation: Insights from animals. *Trends in cognitive sciences*, *6*(9), 376–382.
- Watters, N., Tenenbaum, J., & Jazayeri, M. (2021). Modular object-oriented games: a task framework for reinforcement learning, psychology, and neuroscience. *arXiv preprint arXiv:2102.12616*.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Chi, E., Le, Q., & Zhou, D. (2022). Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.
- Weir, N., & Van Durme, B. (2022, September). *Dynamic Generation of Interpretable Inference Rules in a Neuro-Symbolic Expert System* (No. arXiv:2209.07662). arXiv.
- Wellman, H. M., & Gelman, S. A. (1992). Cognitive development: Foundational theories of core domains. *Annual review of psychology*, *43*(1), 337–375.
- White, J., Mu, J., & Goodman, N. D. (2020). Learning to refer informatively by amortizing pragmatic reasoning. *arXiv preprint arXiv:2006.00418*.
- Whorf, B. (1956). *Language, thought, and reality: selected writings*.
- Wilson, S. M., Molnar-Szakacs, I., & Iacoboni, M. (2008, January). Beyond Superior Temporal Cortex: Intersubject Correlations in Narrative Speech Comprehension. *Cerebral Cortex*, *18*(1), 230–242. Retrieved 2022-06-19, from <https://doi.org/10.1093/cercor/bhm049> doi: 10.1093/cercor/bhm049
- Wingate, D., Stuhlmüller, A., & Goodman, N. (2011). Lightweight implementations of probabilistic programming languages via transformational compilation. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics* (pp. 770–778).
- Wiseman, S., Shieber, S. M., & Rush, A. M. (2017). Challenges in data-to-document generation. *arXiv*

- preprint arXiv:1707.08052.*
- Witty, S., Lew, A., Jensen, D., & Mansinghka, V. (2019). Bayesian causal inference via probabilistic program synthesis. *arXiv preprint arXiv:1910.14124*.
- Wolfram, S. (2023). *ChatGPT gets its “Wolfram Superpowers”*. Retrieved from <https://writings.stephenwolfram.com/2023/03/chatgpt-gets-its-wolfram-superpowers/>
- Wong, C., Ellis, K. M., Tenenbaum, J., & Andreas, J. (2021). Leveraging language to learn program abstractions and search heuristics. In *International conference on machine learning* (pp. 11193–11204).
- Wong, Y. W., & Mooney, R. (2007). Learning synchronous grammars for semantic parsing with lambda calculus. In *Proceedings of the 45th annual meeting of the association of computational linguistics* (pp. 960–967).
- Wu, J., Tenenbaum, J. B., & Kohli, P. (2017). Neural scene de-rendering. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 699–707).
- Wu, J., Yildirim, I., Lim, J. J., Freeman, B., & Tenenbaum, J. (2015a). Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. *Advances in neural information processing systems*, 28.
- Wu, J., Yildirim, I., Lim, J. J., Freeman, B., & Tenenbaum, J. (2015b). Galileo: Perceiving Physical Object Properties by Integrating a Physics Engine with Deep Learning. In *Advances in Neural Information Processing Systems* (Vol. 28). Curran Associates, Inc.
- Wu, M., & Goodman, N. (2022). Foundation posteriors for approximate probabilistic inference. *arXiv preprint arXiv:2205.09735*.
- Wu, S. A., Wang, R. E., Evans, J. A., Tenenbaum, J. B., Parkes, D. C., & Kleiman-Weiner, M. (2021). Too many cooks: Bayesian inference for coordinating multi-agent collaboration. *Topics in Cognitive Science*, 13(2), 414–432.
- Xie, Y., Yu, C., Zhu, T., Bai, J., Gong, Z., & Soh, H. (2023). Translating natural language to planning goals with large-language models. *arXiv preprint arXiv:2302.05128*.
- Xu, K., Srivastava, A., Gutfreund, D., Sosa, F., Ullman, T. D., Tenenbaum, J., & Sutton, C. (2021). A bayesian-symbolic approach to reasoning and learning in intuitive physics. *Advances in Neural Information Processing Systems*, 34, 2478–2490.
- Yang, Y., & Piantadosi, S. T. (2022). One model for the learning of language. *Proceedings of the National Academy of Sciences*, 119(5), e2021865119.
- Yasunaga, M., & Liang, P. (2020). Graph-based, self-supervised program repair from diagnostic feedback. *CoRR*, abs/2005.10636. Retrieved from <https://arxiv.org/abs/2005.10636>
- Yi, K., Gan, C., Li, Y., Kohli, P., Wu, J., Torralba, A., & Tenenbaum, J. B. (2019). Clevrer: Collision events for video representation and reasoning. *arXiv preprint arXiv:1910.01442*.
- Yi, K., Wu, J., Gan, C., Torralba, A., Kohli, P., & Tenenbaum, J. B. (2018). Neural-symbolic vqa: Disentangling reasoning from vision and language understanding. *arXiv preprint arXiv:1810.02338*, 31, 1031–1042.
- Yildirim, I., Belledonne, M., Freiwald, W., & Tenenbaum, J. (n.d.). Efficient inverse graphics in biological face processing. , 77.
- Ying, L., Collins, K., Wei, M., Zhang, C., Tan, Z.-X., Weller, A., ... Wong, L. (2023). The neuro-symbolic inverse planning engine (nipe): Modeling probabilistic social inferences from linguistic inputs. *ICML ToM Workshop 2023*.
- Yuille, A., & Kersten, D. (2006). Vision as bayesian inference: analysis by synthesis? *Trends in cognitive sciences*, 10(7), 301–308.
- Zaslavsky, N., Kemp, C., Regier, T., & Tishby, N. (2018). Efficient compression in color naming and its evolution. *Proceedings of the National Academy of Sciences*, 115(31), 7937–7942.
- Zelikman, E., Wu, Y., Mu, J., & Goodman, N. (2022). Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35, 15476–15488.
- Zhang, C. E., Wong, L., Grand, G., & Tenenbaum, J. B. (2023). Grounded physical language understanding with probabilistic programs and simulated worlds. In *Proceedings of the annual conference of the cognitive science society* (p. To Appear).
- Zhang, J., Panthaplackel, S., Nie, P., Li, J. J., & Gligorić, M. (2022). Coditt5: Pretraining for source code and natural language editing. *ArXiv*, abs/2208.05446.
- Zhi-Xuan, T. (2022). *Pddl. jl: An extensible interpreter and compiler interface for fast and flexible ai planning*

- (Unpublished doctoral dissertation). Massachusetts Institute of Technology.
- Zhi-Xuan, T., Mann, J., Silver, T., Tenenbaum, J., & Mansinghka, V. (2020). Online bayesian goal inference for boundedly rational planning agents. *Advances in neural information processing systems*, *33*, 19238–19250.
- Zhuo, T. Y., Huang, Y., Chen, C., & Xing, Z. (2023). Exploring ai ethics of chatgpt: A diagnostic analysis. *arXiv preprint arXiv:2301.12867*.
- Zinberg, B., Cusumano-Towner, M., & Vikash, K. M. (2019). Structured differentiable models of 3d scenes via generative scene graphs. In *Workshop on perception as generative reasoning, neurips, submitted september*.

## Appendices

We include code for reference below to help better interpret the examples in the paper. This code is included (with human-readable comments) for completeness and for reference, but is not guaranteed to be the the most up-to-date version of these examples. Please refer to the GitHub repository for the most complete, corrected, and up-to-date code for all examples in this paper, as well as instructions for execution and reproducibility: [github.com/gabegrand/world-models](https://github.com/gabegrand/world-models).

## A Language and world models

### A.1 Probabilistic reasoning

#### A.1.1 Generative world model for tug-of-war

---

```

1 ;; This Church program models a tug-of-war game between teams of players.
2
3 ;; Each player has a strength, with strength value 50 being about average.
4 (define strength (mem (lambda (player) (gaussian 50 20))))
5
6 ;; Each player has an intrinsic laziness frequency.
7 (define laziness (mem (lambda (player) (uniform 0 1))))
8
9 ;; The team's strength is the sum of the players' strengths.
10 ;; When a player is lazy in a match, they pull with half their strength.
11 (define (team-strength team)
12 (sum
13 (map
14 (lambda (player) (if (flip (laziness player)) (/ (strength player) 2) (strength player)))
15 team)))
16
17 ;; The winner of the match is the stronger team.
18 ;; Returns true if team-1 won against team-2, else false.
19 (define (won-against team-1 team-2)
20 (> (team-strength team-1) (team-strength team-2)))

```

---

Code Block 1: Generative domain theory for the Bayesian tug-of-war.

#### A.1.2 Translation examples for tug-of-war

---

```

1 ;; Now, let's translate some user-defined statements.
2 ;; Each statement begins with either `Condition` or `Query`.
3 ;; `Condition` statements provide facts about the scenario.
4 ;; `Query` statements are questions that evaluate quantities of interest.
5
6 ;; Condition: Alice won against Bob.
7 (condition (won-against '(alice) '(bob)))
8
9 ;; Condition: John and Mary won against Tom and Sue.
10 (condition (won-against '(john mary) '(tom sue)))
11
12 ;; Query: If Mary played against Tom, who would win?
13 (query (won-against '(mary) '(tom)))
14
15 ;; Certain statements are underspecified and require some interpretation. For example:
16 ;; Condition: Sue is very strong.
17 (condition (> (strength 'sue) 75))
18

```

```

19 ;; We can `define` new constructs that are useful for translation. For example:
20 ;; Condition: Bob is stronger than John.
21 (define (stronger-than? player-1 player-2)
22   (> (strength player-1) (strength player-2)))
23 (condition (stronger-than? 'bob 'john))
24
25 ;; Query: Is Sue stronger than Mary?
26 (query (stronger-than? 'sue 'mary))
27
28 ;; Condition: A couple of the players are stronger than John.
29 (condition (>= (count (map (lambda (player) (stronger-than? player 'john) players)) 2)))

```

---

Code Block 2: Prompt examples

## A.2 Relational reasoning

### A.2.1 Generative world model for kinship

---

```

1 ;; -- KINSHIP GENERATIVE DOMAIN THEORY --
2
3 ;; All the names that can be used in the conversational context.
4 (define ALL-NAMES '(avery blake charlie dana))
5
6 ;; Generates unique person ids of the format (person-0, person-1, ...)
7 (define PERSON-PREFIX "person-")
8 (define new-person-id (make-gensym PERSON-PREFIX))
9 (define (id->idx person-id)
10   (string->number (string-slice (stringify person-id) (string-length PERSON-PREFIX))))
11
12 ;; Randomly assign a gender
13 (define person->gender (mem (lambda (person-id)
14   (uniform-draw '(male female)))))
15
16 ;; Randomly-ordered list of person names
17 (define NAMES (shuffle-unique ALL-NAMES))
18 (define person->name (mem (lambda (person-id)
19   (list-ref NAMES (id->idx person-id)))))
20
21 ;; Person node in tree
22 (define (person person-id parent-1-id parent-2-id) (list
23   (pair 'person-id person-id)
24   (pair 'name person-id)
25   (pair 'gender (person->gender person-id))
26   (pair 'parent-1-id parent-1-id)
27   (pair 'parent-2-id parent-2-id)))
28
29 ;; Generate the full tree
30 ;; Max tree size is 1 + (sum_{n=0}^{n=MAX-DEPTH} 2 * MAX-WIDTH^n)
31 (define MAX-WIDTH 3)
32 (define MAX-DEPTH 2)
33 (define PARTNER-PROBABILITY 0.5)
34 (define (generate-tree root-primary-id root-secondary-id depth)
35   (let* (
36     ;; Create the primary parent
37     (parent-1-id (new-person-id))
38     (parent-1 (person parent-1-id root-primary-id root-secondary-id)))
39     (if (flip PARTNER-PROBABILITY)
40       ;; Case: parent-1 has partner
41       (let* (

```

```

42     ;; Create the secondary parent
43     (parent-2-id (new-person-id))
44     (parent-2 (person parent-2-id () ()))
45
46     ;; Link the parents with a partner relation
47     (parent-1 (append parent-1 (list (pair 'partner-id parent-2-id))))
48     (parent-2 (append parent-2 (list (pair 'partner-id parent-1-id))))
49
50     ;; Generate children
51     (n-children (if (>= depth MAX-DEPTH) 0 (bounded-geometric 0.5 0 MAX-WIDTH)))
52     (child-trees (repeat n-children (lambda () (generate-tree parent-1-id parent-2-id (+ depth 1)))))
53
54     ;; Update the parents to point to the children
55     (child-ids (map (lambda (t) (lookup (first t) 'person-id)) child-trees))
56     (parent-1 (append parent-1 (list (pair 'child-ids child-ids))))
57     (parent-2 (append parent-2 (list (pair 'child-ids child-ids))))
58     (append (list parent-1) (list parent-2) (shallow-flatten child-trees))
59
60     ;; Case: parent-1 has no partner
61     (list parent-1)))
62
63 ;; Generate the global tree.
64 (define T (generate-tree () () 0))
65
66 ;; Assign names randomly to (some of) the people in the tree.
67 (define (add-names-to-tree tree names)
68   (if (null? tree) ()
69       (let*
70         ;; Probability of adding a name to the first person
71         ((p (min 1.0 (/ (length names) (length tree))))
72          (person (first tree)))
73         (if (flip p)
74             ;; Name the person
75             (let
76               ((named-person (update-list person 1 (pair 'name (first names))))
77                (cons named-person (add-names-to-tree (rest tree) (rest names))))
78               ;; Don't name the person
79               (cons person (add-names-to-tree (rest tree) names))))))
80
81   ;; Update the tree with the name information.
82   (define T (add-names-to-tree T NAMES))

```

Code Block 3: Generative domain theory for family trees.

### A.2.2 Kinship tree utilities

```

1  ;; -- KINSHIP TREE UTILITIES --
2
3  ;; Returns all instances of person with property `key` equal to `value`
4  (define filter-by-property
5    (mem (lambda (key value)
6          (filter (lambda (p) (equal? (lookup p key) value)) T))))
7
8  ;; Returns the unique instance of person with name.
9  (define get-person-by-name
10   (mem (lambda (name)
11         (let
12           ((results (filter-by-property 'name name)))
13           (if (null? results) () (first results))))))

```



```

14
15 ;; People without a name can be referenced directly by person-id.
16 (define get-person-by-id
17   (mem (lambda (person-id)
18     (if (null? person-id)
19       ()
20       (let ((idx (id->idx person-id)))
21         (if (>= idx (length T)) () (list-ref T idx)))))))
22
23 ;; Get a person object either by name or person-id.
24 (define get-person
25   (mem (lambda (person-ref)
26     (cond
27       ((null? person-ref) ())
28       ((member? person-ref NAMES) (get-person-by-name person-ref))
29       (else (get-person-by-id person-ref)))))
30
31 ;; Get a property of a person.
32 (define get-property
33   (mem (lambda (name key)
34     (lookup (get-person name) key))))
35
36 ;; -- TREE OPERATORS --
37 ;; predicate :: name -> boolean
38
39 (define (map-tree predicate)
40   (map (lambda (x) (predicate (lookup x 'name))) T))
41
42 (define (filter-tree predicate)
43   (filter (lambda (x) (predicate (lookup x 'name))) T))
44
45 (define (exists predicate)
46   (some (map-tree predicate)))

```

---

Code Block 4: Utility functions for kinship trees.

### A.2.3 Kinship conceptual system

---

```

1 ;; -- KINSHIP CONCEPTUAL SYSTEM --
2
3 ;; Gets the partner of a person.
4 (define (partner-of name)
5   (get-property (get-property name 'partner-id) 'name))
6
7 ;; Gets the parents of a person.
8 (define (parents-of name)
9   (let* ((parent-1-id (get-property name 'parent-1-id))
10         (parent-1-name (get-property parent-1-id 'name))
11         (parent-2-id (get-property name 'parent-2-id))
12         (parent-2-name (get-property parent-2-id 'name)))
13     (list parent-1-name parent-2-name)))
14
15 ;; Gets the grandparents of a person.
16 (define (grandparents-of name)
17   (let ((parent-1 (first (parents-of name))))
18     (parents-of parent-1)))
19
20 ;; Gets the children of a person.
21 (define (children-of name)

```

```

22 (let ((child-ids (get-property name 'child-ids)))
23   (map (lambda (child-id) (get-property child-id 'name)) child-ids))
24
25 ;; Gets the siblings of a person.
26 (define (siblings-of name)
27   (let* ((parent-1-id (get-property name 'parent-1-id))
28          (child-ids (get-property parent-1-id 'child-ids))
29          (child-names (map (lambda (child-id) (get-property child-id 'name)) child-ids)))
30     (filter (lambda (child-name) (not (equal? child-name name))) child-names)))
31
32 ;; -- BOOLEAN RELATIONS --
33 (define (partner-of? name_a name_b)
34   (equal? name_a (partner-of name_b)))
35
36 (define (parent-of? name_a name_b)
37   (member? name_a (parents-of name_b)))
38
39 (define (father-of? name_a name_b)
40   (and (equal? (get-property name_a 'gender) 'male)
41        (parent-of? name_a name_b)))
42
43 (define (mother-of? name_a name_b)
44   (and (equal? (get-property name_a 'gender) 'female)
45        (parent-of? name_a name_b)))
46
47 (define (grandparent-of? name_a name_b)
48   (member? name_a (grandparents-of name_b)))
49
50 (define (grandfather-of? name_a name_b)
51   (and (equal? (get-property name_a 'gender) 'male)
52        (grandparent-of? name_a name_b)))
53
54 (define (grandmother-of? name_a name_b)
55   (and (equal? (get-property name_a 'gender) 'female)
56        (grandparent-of? name_a name_b)))
57
58 (define (child-of? name_a name_b)
59   (member? name_a (children-of name_b)))
60
61 (define (son-of? name_a name_b)
62   (and (equal? (get-property name_a 'gender) 'male)
63        (child-of? name_a name_b)))
64
65 (define (daughter-of? name_a name_b)
66   (and (equal? (get-property name_a 'gender) 'female)
67        (child-of? name_a name_b)))
68
69 (define (sibling-of? name_a name_b)
70   (member? name_a (siblings-of name_b)))
71
72 (define (brother-of? name_a name_b)
73   (and (equal? (get-property name_a 'gender) 'male)
74        (sibling-of? name_a name_b)))
75
76 (define (sister-of? name_a name_b)
77   (and (equal? (get-property name_a 'gender) 'female)
78        (sibling-of? name_a name_b)))

```

Code Block 5: Conceptual system and derived predicates for kinship trees.

## A.2.4 Translation examples for kinship

---

```

1 ;; -- CONDITION AND QUERY STATEMENTS --
2 ;; Now, let's translate some user-defined statements.
3 ;; Each statement begins with either `Condition` or `Query`.
4 ;; `Condition` statements provide facts about the scenario.
5 ;; `Query` statements are questions that evaluate quantities of interest.
6
7 ;; Condition: Ryan's partner is Taylor.
8 (condition (partner-of? 'ryan 'taylor))
9
10 ;; Condition: Taylor is the mother of Sam.
11 (condition (mother-of? 'taylor 'sam))
12
13 ;; Condition: Sam's father is Ryan.
14 (condition (father-of? 'ryan 'sam))
15
16 ;; Condition: Sam has two siblings.
17 (condition (= (length (siblings-of 'sam)) 2))
18
19 ;; Condition: Sam has a brother.
20 (condition
21   (exists (lambda (x)
22     (brother-of? x 'sam))))
23
24 ;; Condition: Payton's partner has a brother named Kyle.
25 (condition
26   (exists (lambda (x) (and
27     (partner-of? x 'payton)
28     (brother-of? 'kyle x))))))
29
30 ;; Condition: Payton's partner has a sister who has a son named Sam.
31 (condition
32   (exists (lambda (x) (and
33     (partner-of? x 'payton)
34     (exists (lambda (y) (and
35       (sister-of? y x)
36       (son-of? 'sam y))))))))))
37
38 ;; Query: Who are Sam's parents?
39 (query (parents-of 'sam))
40
41 ;; Query: How many children does Kyle have?
42 (query (length (children-of 'kyle)))
43
44 ;; Query: Who is Ryan's grandfather?
45 (query
46   (filter-tree
47     (lambda (x) (grandfather-of? x 'ryan))))
48
49 ;; Query: Does Taylor have a sister?
50 (query
51   (exists (lambda (x)
52     (sister-of? x 'taylor))))

```

---

Code Block 6: Translation examples for kinship trees.

### A.2.5 Why not Prolog?

Readers who are familiar with the world of logic programming may wonder why we have chosen to model the kinship domain in Church instead of a more standard logic programming language, such as Prolog. Indeed, kinship is often one of the introductory examples in Prolog textbooks (Pereira & Shieber, 2002) and online tutorials,<sup>8</sup> from which we drew inspiration when writing this section. Moreover, there are many structural parallels between our framework and the style of declarative programming embodied by Prolog: schemecondition statements in Church are similar to facts in Prolog; derived concepts like schemefather-of? in our Church kinship model are analogous to Prolog rules; and schemequery performs similar functions in both languages (though the algorithms that underlie these queries differ in important ways). And, as discussed in the introduction to Section 3.1, Prolog was originally developed as a model of natural language (Colmerauer et al., 1972) and has deep ties to computational linguistics. So: why not use Prolog?

In short, there is nothing about our approach to semantic parsing that precludes swapping out Church for other programming languages, like Prolog, SMT solvers, or even a general purpose language like Python. In fact, with the right prompting, Codex readily translates natural language utterances like *Avery has a sister named Blake* into `sister_of(blake, avery)` in Prolog. On the parsing side, we did not encounter any technical limitations to using LLMs to translate natural language into Prolog.

However, because Prolog is based on definite (Horn) clauses, there are limitations in the kinds of utterances that we can express and the kinds of inferences that we can make when working in Prolog. For instance, a typical Prolog kinship model might have a rule defining the concept of a “grandfather” as follows:

---

```
grandfather_of(X,Y) :- male(X),
    parent_of(X,Z),
    parent_of(Z,Y).
```

---

Now, if we learn that *Charlie is the grandfather of Dana*, we might be inclined to translate this into Prolog as a fact: `grandfather_of(charlie, dana)`. Given this information, we can make various deductive inferences: e.g., that Charlie is male, and that there exists some person in the family tree who is both the child of Charlie and the parent of Dana. In fact, this is how the `grandfather_of(X,Y)` rule is defined in the first place.

For this reason, it is especially counterintuitive that these kinds of inferences are not at all straightforward in Prolog. Because logical implication in definite clauses is *unidirectional*, anyone satisfying the right-hand side of the `grandfather_of(X,Y)` rule is considered a grandfather. However, our rule says nothing about what being a grandfather entails. Moreover, our above translation `grandfather_of(charlie, dana)` is actually quite facile; it simply modifies `grandfather_of(X,Y)` such that queries will now return true for anyone satisfying the original definition; or for the special case where  $x=charlie$  and  $y=dana$ . These are all examples of limitations of the kinds of *deductive* inferences that we can model with Prolog. Additionally, there are many kinds of *inductive* inferences that are not well-captured by Prolog; e.g., because Charlie has at least one child, he is more likely to have multiple children, and is more likely to be married.

In sum, to get the kinds of mixed deductive and inductive inferences that we would like to see in an expressive language-of-thought, we need to have ways of incorporating and trading off uncertainty in our world model. ProbLog (De Raedt et al., 2007; Dries et al., 2017; Suster et al., 2021), a probabilistic extension of Prolog in which deduction rules can be annotated with probabilities, offers one way of integrating uncertainty with deductive reasoning. Church goes a step further by specifying a *generative domain theory* in addition to probabilistic inference rules. We believe that this interplay between probabilistic priors and likelihoods—which is central to Bayesian inference—is also at the heart of human cognition.

## A.3 Perceptual and physical reasoning

### Static visual scenes

#### A.3.1 Generative world model for static visual scenes

---

```
1 ;; Objects have a shape attribute, which is a choice of cube, sphere, or cylinder shape categories.
2 (define choose-shape
3   (mem (lambda (obj-id)
```

---

<sup>8</sup><https://swish.swi-prolog.org/p/prolog-family-tree.pl>

```

4      (pair 'shape (uniform-draw '(mug can bowl))))))
5
6  ;; Objects have a color attribute that is drawn from a predefined set of RGB values.
7  (define choose-color
8    (mem (lambda (obj-id)
9          (pair 'color (uniform-draw (list
10                                     (list 255 0 0)
11                                     (list 0 0 255)
12                                     (list 0 255 0)
13                                     (list 255 255 0)
14                                     ))))))
15  ;; An object is an object ID, and the object's attribute types and their values.
16  (define object (mem (lambda (obj-id) (list
17                       (pair 'object-id obj-id)
18                       (choose-shape obj-id)
19                       (choose-color obj-id))))))
20
21  ;; Scenes can have a maximum of 12 objects.
22  (define max-objects 12)
23  ;; The number of objects in a scene tends to be not too large, and is capped at the maximum number of
24  objects.
25  (define choose-num-objects
26    (mem (lambda (scene-id) (floor (min max-objects (* max-objects (exponential 1)))))))
27  ;; Then, for each object we intend to generate, generate an object indexical, and associate it with a
28  choice of attributes.
29  (define obj-id-gensym (make-gensym "obj-"))
30  (define (generate-n-objects scene-id total-objects)
31    (if (= total-objects 0)
32        (list (object (obj-id-gensym)))
33        (cons (object (obj-id-gensym)) (generate-n-objects scene-id (- total-objects 1)))))
34  (define objects-in-scene (mem (lambda (scene-id) (generate-n-objects scene-id (choose-num-objects
35                                  scene-id)))))
36
37  ;; An object is red if it is of this continuous color value.
38  (define red (list 255 0 0))
39  ;; An object is blue if it is of this continuous color value.
40  (define blue (list 0 0 255))
41  ;; An object is green if it is of this continuous color value.
42  (define green (list 0 255 0))
43  ;; An object is yellow if it is of this continuous color value.
44  (define yellow (list 255 255 0))
45
46  ;; Check if an object is of a given shape.
47  (define is-shape? (lambda (shape) (lambda (object) (equal? (cdr (assoc 'shape object)) shape))))
48  ;; Check if an object is of a given named color.
49  (define is-color? (lambda (color) (lambda (object) (equal? (cdr (assoc 'color object)) color))))
50
51  ;; Select only objects from the scene of a given color.
52  (define filter-color (lambda (color) (lambda (object-list) (filter (is-color? color) object-list))))
53
54  ;; Select only objects from the scene of a given shape.
55  (define filter-shape (lambda (shape) (lambda (object-list) (filter (is-shape? shape) object-list))))

```

---

Code Block 7: Generative domain theory for tabletop scenes. Generates scenes containing a set of objects which vary in shape and color. These scene states are rendered by a separately generated render function to generate images. Shown with natural language comments, but these are not used in the LLM prompt.

### A.3.2 Translation examples for static visual scenes

---

```

1 ;; There's a blue thing.
2 (condition (> (length ((filter-color blue) (objects-in-scene 'this-scene))) 0))
3 ;; There's at least two blue plates.
4 (condition (>= (length
5   ((filter-color blue)
6   ((filter-shape 'plate)
7   (objects-in-scene 'scene))))
8 2))
9 ;; There's many blue plates.
10 (condition (>= (length
11   ((filter-color blue)
12   ((filter-shape 'plate)
13   (objects-in-scene 'scene))))
14 5))
15 ;; There's exactly two plates and there's also a yellow thing.
16 (condition
17   (and (= (length ((filter-shape 'plate) (objects-in-scene 'scene))) 2)
18   (> (length ((filter-color yellow) (objects-in-scene 'scene))) 0)))
19
20 ;; Is there a mug?
21 (query (> (length ((filter-shape 'mug) (objects-in-scene 'this-scene))) 0))

```

---

Code Block 8: Translation examples for the visual domain. These examples are concatenated with the visual scenes generative model to produce the prompt used to generate new translations.

### Dynamic physical scenes

#### A.3.3 Generative world model for physical scenes

---

```

1 (define (get_attribute obj key)
2   (if (assoc key obj) (rest (assoc key obj)) ()))
3
4 (define (member? a b)
5   (if (member a b) true false))
6 (define concatenate
7   (lambda (list-1 list-2)
8     (if (null? list-1)
9       list-2
10      (cons (car list-1) (concatenate (cdr list-1) list-2))))))
11
12 (define (pairs x l)
13   (define (aux accu x l)
14     (if (null? l)
15       accu
16       (let ((y (car l))
17             (tail (cdr l)))
18         (aux (cons (cons x y) accu) x tail))))
19   (aux '() x l))
20
21 (define (cartesian_product l m)
22   (define (aux accu l)
23     (if (null? l)
24       accu
25       (let ((x (car l))
26             (tail (cdr l)))

```

```

27     (aux (append (pairs x m) accu) tail))))
28 (aux '() 1))
29
30 ;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;
31 (define get_num_objects 2)
32 (define OBJECT_DEFAULT_RADIUS 1)
33 (define GRAVITY 9.8)
34 (define DELTA_T 0.5)
35
36 (define get_initial_color
37   (lambda (obj_id)
38     (if (eq? obj_id 'obj-0)
39         (list 255 0 0)
40         (list 0 0 255))))
41
42 (define choose_mass
43   (mem (lambda (obj_id)
44         (abs (gaussian 5 3)))))
45
46 (define choose_shapes
47   (mem (lambda (scene-id) (uniform-draw (list 'sphere 'block)))))
48
49 (define min_x -3)
50 (define max_x 3)
51 (define mid_x (+ (/ (- max_x min_x) 2) min_x))
52 (define get_initial_x
53   (lambda (obj_id)
54     (if (eq? obj_id 'obj-0)
55         min_x
56         mid_x)))
57
58 (define min_force 0)
59 (define max_force 10)
60 (define mid_force (+ (/ (- max_force min_force) 2) min_force))
61 (define choose_initial_force
62   (mem (lambda (obj_id)
63         (if (eq? obj_id 'obj-0)
64             (abs (gaussian mid_force 3))
65             0
66             )))
67
68 (define static_friction_constant (lambda (shape)
69                                   (if (eq? shape 'sphere)
70                                       0.02
71                                       0.05)
72                                   ))
73 (define kinetic_friction_constant (lambda (shape)
74                                     (if (eq? shape 'sphere)
75                                         0.01
76                                         0.02)
77                                     ))
78 (define normal_force (lambda (m) (* m GRAVITY)))
79 (define force_after_friction (lambda (f v shape m)
80                                (if (> (abs v) 0)
81                                    (- f (* (kinetic_friction_constant shape) (normal_force m)))
82                                    (if (< f (* (static_friction_constant shape) (normal_force m))) 0 (- f (*
83                                         (kinetic_friction_constant shape) (normal_force m))
                                         ))))
84                                ))))

```

```

84
85 (define newtons_second (lambda (f m) (/ f m)))
86 (define v_next (lambda (v_prev a_prev delta_t)
87     (let ((v_temp (+ v_prev (* a_prev delta_t))))
88         (if (>= (* v_prev v_temp) 0) v_temp 0))
89 ))
90 (define x_next (lambda (x_prev v_prev delta_t) (+ x_prev (* v_prev delta_t))))
91 (define initial_object_state (mem (lambda (obj_id scene_id)
92     (let ((obj_shape (choose_shapes scene_id)))
93         (let ((obj_mass (choose_mass obj_id)))
94             (let ((obj_color (get_initial_color obj_id)))
95                 (let ((initial_x (get_initial_x obj_id)))
96                     (let ((initial_push_force (choose_initial_force obj_id)))
97                         (let ((initial_force (force_after_friction
initial_push_force 0 obj_shape obj_mass)))
98                             (list
99                                 (pair 'object_id obj_id)
100                                (pair 'object_radius OBJECT_DEFAULT_RADIUS)
101                                (pair 'shape obj_shape)
102                                (pair 'mass obj_mass)
103                                (pair 'color obj_color)
104                                (pair 'x initial_x)
105                                (pair 'initial_push_force initial_push_force)
106                                (pair 'f initial_force)
107                                (pair 't 0)
108                                (pair 'a_prev (newtons_second initial_force obj_mass))
109                                (pair 'a (newtons_second initial_force obj_mass))
110                                (pair 'v_0 0)
111                                (pair 'v (v_next 0 (newtons_second initial_force
obj_mass) DELTA_T)))
112                                 ))))))))
113 (define obj_id_gensym (make_gensym "obj-"))
114 (define generate_initial_state
115     (mem (lambda (scene_id total_objects)
116         (if (= total_objects 1)
117             (list (initial_object_state (obj_id_gensym) scene_id))
118             (cons (initial_object_state (obj_id_gensym) scene_id) (generate_initial_state scene_id
(- total_objects 1)))))))
119
120 (define generate_initial_scene_event_state (mem (lambda (scene_id total_objects)
121     (pair 0
122         (list
123             (pair 'scene_states
(generate_initial_state scene_id total_objects))
124             (pair 'event_states [])
125             ))
126     ))
127 ))
128
129 (define event_id_gensym (make_gensym "event-"))
130 (define circle_intersect? (lambda (subject_x subject_radius object_x object_radius)
131 (let ((square_circle_distance (expt (- subject_x object_x) 2)))
132 (let ((square_radii (expt (+ subject_radius object_radius) 2)))
133 (leq square_circle_distance square_radii)))
134 ))
135 (define elastic_collision_subject_v (lambda (subject_m subject_v object_m object_v)
136 (/ (+ (* 2 (* object_m object_v)) (* subject_v (- subject_m object_m))) (+ subject_m object_m))
137 ))
138

```



```

139 (define get_collision_events (lambda (time scene_event_state_for_time)
140   (let ((scene_event_state (get_attribute scene_event_state_for_time time)))
141     (let ((scene_state (get_attribute scene_event_state 'scene_states)))
142       (if (= (length scene_state) 1)
143         ()
144         (fold (lambda (event events) (if (equal? event ()) events (cons event events))) ()
145         (let ((paired_object_states (cartesian_product scene_state scene_state)))
146           (map (lambda (paired_objects)
147             (let ((event_subject (get_attribute (first paired_objects) 'object_id)))
148               (let ((event_object (get_attribute (cdr paired_objects) 'object_id)))
149                 (if (eq? event_subject event_object) ()
150                     (let ((subject_v (get_attribute (first paired_objects) 'v)))
151                       (let ((subject_x (get_attribute (first paired_objects) 'x)))
152                         (let ((subject_m (get_attribute (first paired_objects) 'mass)))
153                           (let ((subject_radius (get_attribute (first paired_objects) 'object_radius)))
154                             (let ((object_v (get_attribute (cdr paired_objects) 'v)))
155                               (let ((object_x (get_attribute (cdr paired_objects) 'x)))
156                                 (let ((object_m (get_attribute (cdr paired_objects) 'mass)))
157                                   (let ((object_radius (get_attribute (cdr paired_objects) 'object_radius)))
158                                     (if (circle_intersect? subject_x subject_radius object_x object_radius)
159                                         (list
160                                           (pair 'event-id (event_id_gensym))
161                                           (pair 'event_time time)
162                                           (pair 'event_predicates (list 'is_colliding))
163                                           (pair 'event_subject event_subject)
164                                           (pair 'event_object event_object)
165                                           (pair 'subject_initial_v subject_v)
166                                           (pair 'subject_final_v (elastic_collision_subject_v subject_m subject_v object_m
167                                             object_v))
168                                           (pair 'object_initial_v object_v)
169                                         )
170                                         ())))))))))))))
171         )))
172         paired_object_states)))
173     ))))
174
175
176 (define generate_next_object_state (lambda (current_time event_state) (lambda (prev_object_state)
177   (let ((obj_id (cdr (assoc 'object_id prev_object_state))))
178     (let ((collision_events (fold (lambda (event events) (if (equal? (get_attribute event 'event_subject)
179       obj_id) (cons event events) events)) () event_state)))
180       (if (> (length collision_events) 0)
181         (generate_collision_event_state current_time obj_id prev_object_state (car collision_events))
182         (generate_no_collision_event_state current_time obj_id prev_object_state))
183     )))
184
185 (define generate_collision_event_state (lambda (current_time obj_id prev_object_state collision_event)
186   (let ((obj_radius (cdr (assoc 'object_radius prev_object_state))))
187     (let ((obj_mass (cdr (assoc 'mass prev_object_state))))
188       (let ((obj_color (cdr (assoc 'color prev_object_state))))
189         (let ((obj_shape (cdr (assoc 'shape prev_object_state))))
190           (let ((v_prev (cdr (assoc 'v prev_object_state))))
191             (let ((a_prev (cdr (assoc 'a_prev prev_object_state))))
192               (let ((x_prev (cdr (assoc 'x prev_object_state))))
193                 (let ((v (get_attribute collision_event 'subject_final_v)))
194                   (let ((x (x_next x_prev v 1)))
195                     (list

```

```

196         (pair 'object_id obj_id)
197         (pair 'object_radius obj_radius)
198         (pair 'shape obj_shape)
199         (pair 'color obj_color)
200         (pair 'mass obj_mass)
201         (pair 'x x)
202         (pair 'f 0)
203         (pair 't (* current_time DELTA_T))
204         (pair 'a_prev 0)
205         (pair 'a 0)
206         (pair 'v_0 0)
207         (pair 'v v))
208         ))))
209     ))))
210 ))
211
212 (define generate_no_collision_event_state (lambda (current_time obj_id prev_object_state)
213   (let ((obj_radius (cdr (assoc 'object_radius prev_object_state))))
214     (let ((obj_mass (cdr (assoc 'mass prev_object_state))))
215       (let ((obj_color (cdr (assoc 'color prev_object_state))))
216         (let ((obj_shape (cdr (assoc 'shape prev_object_state))))
217           (let ((v_prev (cdr (assoc 'v prev_object_state))))
218             (let ((a_prev_no_friction (cdr (assoc 'a_prev prev_object_state))))
219               (let ((a_prev (newtons_second (force_after_friction 0 v_prev obj_shape obj_mass) obj_mass)))
220                 (let ((x_prev (cdr (assoc 'x prev_object_state))))
221                   (let ((v (v_next v_prev a_prev DELTA_T)))
222                     (let ((x (x_next x_prev v_prev DELTA_T)))
223                       (list
224                         (pair 'object_id obj_id)
225                         (pair 'object_radius obj_radius)
226                         (pair 'shape obj_shape)
227                         (pair 'color obj_color)
228                         (pair 'mass obj_mass)
229                         (pair 'x x)
230                         (pair 'f (force_after_friction 0 v_prev obj_shape obj_mass))
231                         (pair 't (* current_time DELTA_T))
232                         (pair 'a_prev a_prev)
233                         (pair 'a 0)
234                         (pair 'v_0 0)
235                         (pair 'v v))
236                       ))))
237                   ))))
238             ))))
239
240 (define generate_next_scene_state (lambda (prev_scene_state event_state next_time)
241   (map (generate_next_object_state next_time event_state) prev_scene_state))
242
243 (define generate_next_scene_event_state_time (lambda (next_time scene_event_state_for_times)
244   (let ((prev_scene_event_state (get_attribute scene_event_state_for_times (- next_time 1))))
245     (let ((prev_scene_state (get_attribute prev_scene_event_state 'scene_states)))
246       (let ((event_state (get_collision_events (- next_time 1) scene_event_state_for_times)))
247
248         (pair next_time (list
249           (pair 'scene_states (generate_next_scene_state prev_scene_state event_state next_time))
250           (pair 'event_states event_state)
251         ))
252       ))))
253
254 (define generate_next_scene_event_states

```

```

255   (lambda (current_time prev_scene_event_states_for_times)
256   (cons (generate_next_scene_event_state_time current_time prev_scene_event_states_for_times)
        prev_scene_event_states_for_times)
257 ))
258
259 (define generate_scene_event_states_for_times (mem (lambda (scene_id total_objects total_time)
260   (if (= total_time 0)
261       (list
262         (generate_initial_scene_event_state
263           scene_id total_objects)
264         (let ((prev_scene_event_states
265             (generate_scene_event_states_for_times scene_id total_objects (- total_time 1))))
            (generate_next_scene_event_states
266               total_time prev_scene_event_states)
267         )))
268 ))))
269
270 (define base_states_for_times (generate_scene_event_states_for_times 'this_scene get_num_objects
271   max_time))
272 ;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;Derived predicates.
273 (define objects_in_scene (lambda (base_states_for_times)
274   (let ((initial_base_states_at_time (cdr (assoc 0 (cdr base_states_for_times)))))
275     (let ((base_state (cdr (assoc 'scene_states initial_base_states_at_time)))
276         base_state
277       ))
278   ))
279 (define red (list 255 0 0))
280 (define blue (list 0 0 255))
281 (define is_color? (lambda (color) (lambda (object) (equal? (cdr (assoc 'color object)) color))))
282 (define is_shape? (lambda (shape) (lambda (object) (equal? (cdr (assoc 'shape object)) shape))))
283
284 (define all_objects (objects_in_scene base_states_for_times))
285 (define (exists_object predicate)
286   (some (map predicate (objects_in_scene base_states_for_times))))
287
288 (define (filter_objects predicate)
289   (map
290     (lambda (o) (get_attribute o 'object_id))
291     (filter predicate (objects_in_scene base_states_for_times))))
292 ;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;
293 (define QUICKLY_THRESHOLD 2)
294 (define SLOWLY_THRESHOLD 2)
295
296 (define is_moving_events (mem (lambda (base_states_for_times)
297   (fold (lambda (base_state_for_time these_events)
298     (let ((current_time (car base_state_for_time))
299         (let ((base_state (cdr (assoc 'scene_states (cdr base_state_for_time))))
300         (fold (lambda (obj_state these_events)
301           (let ((obj_id (cdr (assoc 'object_id obj_state)))
302             (let ((obj_velocity (cdr (assoc 'v obj_state)))
303               (let ((obj_speed (abs obj_velocity)))
304                 (if (> obj_speed 0)
305                     ;;
306                     (let ((event_predicates
307                         (if (> obj_speed QUICKLY_THRESHOLD)
308                           (list 'is_moving 'is_quickly)

```

```

309         (if (< obj_speed SLOWLY_THRESHOLD)
310             (list 'is_moving 'is_slowly)
311             (list 'is_moving)
312             ))
313     ))
314     (cons
315         (list
316             (pair 'event-id (event_id_gensym))
317             (pair 'event_time current_time)
318             (pair 'event_predicates event_predicates)
319             (pair 'event_subject obj_id)
320             (pair 'event_speed obj_speed)
321             )
322         these_events))
323     these_events
324
325     )))))
326     these_events base_state))))
327 () base_states_for_times))))
328
329 (define is_resting_events (mem (lambda (base_states_for_times)
330     (fold (lambda (base_state_for_time these_events)
331         (let ((current_time (car base_state_for_time)))
332             (let ((base_state (cdr (assoc 'scene_states (cdr base_state_for_time)))))
333                 (fold (lambda (obj_state these_events)
334                     (let ((obj_id (cdr (assoc 'object_id obj_state))))
335                         (let ((obj_velocity (cdr (assoc 'v obj_state))))
336                             (let ((obj_speed (abs obj_velocity)))
337                                 (if (= obj_speed 0)
338                                     ;;
339                                     (let ((event_predicates
340                                         (list 'is_resting)))
341                                         (cons
342                                             (list
343                                                 (pair 'event-id (event_id_gensym))
344                                                 (pair 'event_time current_time)
345                                                 (pair 'event_predicates event_predicates)
346                                                 (pair 'event_subject obj_id)
347                                                 (pair 'event_speed obj_speed)
348                                                 )
349                                             these_events))
350                                         these_events
351                                         )))))
352                                     these_events base_state))))
353             (base_states_for_times))))
354     (base_states_for_times))))
355
356 (define is_colliding_events (mem (lambda (base_states_for_times)
357     (fold (lambda (base_state_for_time these_events)
358         (let ((current_time (car base_state_for_time)))
359             (let ((event_states (cdr (assoc 'event_states (cdr base_state_for_time)))))
360                 (fold (lambda (event_state these_events)
361                     (let ((subject_initial_speed (abs (get_attribute event_state 'subject_initial_v))))
362                         (let ((subject_final_speed (abs (get_attribute event_state 'subject_final_v))))
363                             (let ((object_initial_speed (abs (get_attribute event_state 'object_initial_v))))
364                                 (let ((cause_subject_object_event (and (> subject_initial_speed 0) (=
object_initial_speed 0))))
365                                     (let
366                                         ((event_predicates

```

```

367         (if (and cause_subject_object_event (eq? subject_final_speed 0))
368             (list 'is_launching 'is_hitting 'is_colliding)
369             (if (> subject_initial_speed 0)
370                 (list 'is_hitting 'is_colliding)
371                 (list 'is_colliding)
372             )
373         )))
374
375     (cons (list
376         (pair 'event-id (get_attribute event_state 'event-id))
377         (pair 'event_time (get_attribute event_state 'event_time))
378         (pair 'event_predicates event_predicates)
379         (pair 'event_subject (get_attribute event_state 'event_subject))
380         (pair 'event_object (get_attribute event_state 'event_object))
381         (pair 'subject_initial_v (get_attribute event_state 'subject_initial_v ))
382         (pair 'subject_final_v (get_attribute event_state 'subject_final_v ))
383         (pair 'object_initial_v (get_attribute event_state 'object_initial_v ))
384         ) these_events))))))
385     ) these_events event_states)
386 )))
387 () base_states_for_times)
388
389 )))
390
391
392
393 (define events_in_scene (concatenate
394     (is_colliding_events base_states_for_times)
395     (concatenate
396         (is_moving_events base_states_for_times)
397         (is_resting_events base_states_for_times))))
398
399
400 (define is_event? (lambda (event_predicate event) (member? event_predicate (get_attribute event
    'event_predicates))))
401
402 (define is_subject_of_event? (lambda (event object ) (equal?
403     (get_attribute event 'event_subject)
404     (get_attribute object 'object_id)
405 )))
406
407 (define is_object_of_event? (lambda (event object ) (equal?
408     (get_attribute event 'event_object)
409     (get_attribute object 'object_id)
410 )))
411
412 (define event_subject_is? (lambda (event predicate) (member?
413     (get_attribute event 'event_subject)
414     (filter_objects predicate)
415 )))
416 (define event_object_is? (lambda (event predicate) (member?
417     (get_attribute event 'event_object)
418     (filter_objects predicate)
419 )))
420
421 (define (exists_event predicate)
422     (some (map predicate events_in_scene)))
423
424 (define (filter_events predicate)

```

425 (filter predicate events\_in\_scene))

Code Block 9: Generative domain theory for physical scenes. Generates scenes containing a red object left of a blue object, and a randomly generated force. These scene states are forward simulated using a physics engine which is shown implemented within this Church code. Shown with natural language comments, but these are not used in the LLM prompt.

#### A.3.4 Translation examples for visual scenes

```

1 ;; The objects are all balls.
2 (condition (all (map (lambda (o) ((is_shape? 'sphere) o)) all_objects)))
3 ;; Everything is a ball.
4 (condition (all (map (lambda (o) ((is_shape? 'sphere) o)) all_objects)))
5 ;; Imagine the red thing is a block, and is somewhat heavy.
6 (condition (exists_object (lambda (object)
7     (and
8         ((is_color? red) object)
9         ((is_shape? 'cube) object)
10        (> (get_attribute object 'mass) 2)
11        ))))
12 ;; There is a blue ball, and it is quite heavy.
13 (condition (exists_object (lambda (object)
14     (and
15         ((is_color? blue) object)
16         ((is_shape? 'sphere) object)
17         (> (get_attribute object 'mass) 3.5)
18         ))))
19 ;; Now, the red block is very light.
20 (condition (exists_object (lambda (object)
21     (and
22         ((is_color? red) object)
23         ((is_shape? 'cube) object)
24         (< (get_attribute object 'mass) 1)
25         ))))
26 ;; A blue ball is somewhat light.
27 (condition (exists_object (lambda (object)
28     (and
29         ((is_color? red) object)
30         ((is_shape? 'cube) object)
31         (< (get_attribute object 'mass) 2)
32         ))))
33 ;; Imagine the red block gets pushed lightly to the right.
34 (condition (exists_object (lambda (object)
35     (and
36         ((is_color? red) object)
37         ((is_shape? 'cube) object)
38         (< (get_attribute object 'initial_push_force) 2)
39         ))))
40 ;; Now, imagine a red ball is pushed hard to the right.
41 (condition (exists_object (lambda (object)
42     (and
43         ((is_color? red) object)
44         ((is_shape? 'sphere) object)
45         (> (get_attribute object 'initial_push_force) 6)
46         ))))
47 ;; A red block hits a blue block.
48 (condition
49 (exists_object (lambda (object_1)

```

```

50 (exists_object (lambda (object_2)
51 (exists_event (lambda (event)
52     (and
53         ((is_color? red) object_1)
54         ((is_shape? 'cube) object_1)
55         ((is_color? blue) object_2)
56         ((is_shape? 'cube) object_2)
57         (is_subject_of_event? event object_1)
58         (is_object_of_event? event object_2)
59         (is_event? 'is_hitting event))
60     )))))))
61 ;; What's the final velocity of the red block after it is hit?
62 (query (last (map
63 (lambda (event) (get_attribute event 'subject_final_v))
64 (filter_events
65 (lambda (e)
66 (and
67 (is_event? 'is_colliding e)
68 (event_subject_is? e (lambda (o)
69     (and
70         ((is_color? red) o)
71         ((is_shape? 'cube) o))))))))))

```

Code Block 10: Translation examples for the physics domain. These examples are concatenated with the physical scenes generative model to produce the prompt used to generate new translations.

## A.4 Social reasoning

### A.4.1 Generative world model for social reasoning

```

1 (define gridworld (list
2   (list 'ames 'lawn 'lawn 'lawn 'sushi)
3   (list 'ames 'lawn 'lawn 'lawn 'danner)
4   (list 'office 'barlow 'barlow 'barlow 'danner)
5   (list 'ames 'lawn 'lawn 'lawn 'danner)
6   (list 'ames 'lawn 'lawn 'lawn 'vegetarian)
7   (list 'pizza 'carson 'carson 'carson 'danner)
8 ))
9 (define restaurants (list 'sushi 'pizza 'vegetarian))
10
11 (define initial_x 1)
12 (define initial_y 3)
13
14
15 (define has_bike (mem (lambda (agent-id) (flip))))
16 (define available_motions (mem (lambda (agent-id) (if (has_bike agent-id) (list 'is_walking 'is_biking)
17   (list 'is_walking)))))
17 (define directions (list 'west 'east 'north 'south))
18 (define available_actions (mem (lambda (agent-id) (cons (pair 'stay 'stay) (cartesian_product
19   (available_motions agent-id) directions)))))
19
20 (define is_open (mem (lambda (restaurant_type) (flip))))
21 (define POSITIVE_UTILITY_MEAN 10)
22 (define NEGATIVE_UTILITY_MEAN -10)
23 (define UTILITY_VARIANCE 1)
24 (define restaurant_utility (mem (lambda (agent-id restaurant_type)
25   (uniform-draw

```

```

26             (list
27             (gaussian POSITIVE_UTILITY_MEAN UTILITY_VARIANCE)
28             (gaussian NEGATIVE_UTILITY_MEAN UTILITY_VARIANCE)
29 )))))
30
31 (define motion_utility (mem (lambda (agent-id location_type motion_type)
32   (case location_type
33     (('lawn) (case motion_type
34               (('is_biking) -1)
35               (('is_walking) -0.2)
36               (('is_staying) 0)
37               (else 0))
38             )
39     (else (case motion_type
40              (('is_biking) -0.01)
41              (('is_walking) -0.2)
42              (('is_staying) 0)
43              (else 0)))
44   )))
45
46 (define food_utility (mem (lambda (agent-id location_type)
47   (case location_type
48     (('lawn) 0)
49     (('ames) 0)
50     (('barlow) 0)
51     (('carson) 0)
52     (('danner) 0)
53     (('office) 0)
54     (else
55      (if (is_open location_type) (restaurant_utility agent-id location_type) NEGATIVE_UTILITY_MEAN))
56   )))
57
58 (define utility_function (mem (lambda (agent-id gridworld state_x state_y action)
59   (let ((location_type (get_gridworld_at gridworld state_x state_y)))
60     (let ((motion_type (car action)))
61       (let ((state_food_utility (food_utility agent-id location_type)))
62         (let ((state_motion_utility (motion_utility agent-id location_type motion_type)))
63           (+ state_food_utility state_motion_utility)))))))))
64
65 (define get_gridworld_at (lambda (gridworld x y)
66   (list-elt (list-elt gridworld y) x)
67 ))
68 (define x_increment (lambda (direction)
69   (case direction
70     (('west) -1)
71     (('east) 1)
72     (('north) 0)
73     (('south) 0)
74     (('stay) 0)
75 ))
76 (define y_increment (lambda (direction)
77   (case direction
78     (('north) -1)
79     (('south) 1)
80     (('west) 0)
81     (('east) 0)
82     (('stay) 0)
83 ))
84 (define gridworld_max_x (lambda (gridworld) (length (list-elt gridworld 1))))

```



```

85 (define gridworld_max_y (lambda (gridworld) (length gridworld)))
86 (define gridworld_transition (lambda (gridworld current_x current_y action)
87   (let ((direction (cdr action)))
88     (let ((next_x (if (>= current_x (gridworld_max_x gridworld)) current_x (+ (x_increment direction)
89       current_x))))
90     (let ((next_x (if (< next_x 1) current_x next_x)))
91     (let ((next_y (if (>= current_y (gridworld_max_y gridworld)) current_y (+ (y_increment direction)
92       current_y))))
93     (let ((next_y (if (< next_y 1) current_y next_y)))
94     (let ((next_state (get_gridworld_at gridworld next_x next_y)))
95     (list next_state next_x next_y)
96     ))))))))
97 (define value_function (mem (lambda (agent-id curr_iteration gridworld state_x state_y)
98   (if (equal? curr_iteration -1) 0
99     (let ((prev_optimal_action_value (optimal_action_value agent-id (- curr_iteration 1) gridworld
100       state_x state_y)))
101     (cdr prev_optimal_action_value)
102     )))
103 (define available_actions_to_values (mem (lambda (agent-id curr_iteration gridworld state_x state_y)
104   (map (lambda (action)
105     (let ((utility (utility_function agent-id gridworld state_x state_y action)))
106     (let ((next_state (gridworld_transition gridworld state_x state_y action)))
107     (let ((next_state_x (second next_state)))
108     (let ((next_state_y (third next_state)))
109     (let ((next_state_value (value_function agent-id curr_iteration gridworld next_state_x
110       next_state_y)))
111     (pair action (+ utility next_state_value)
112     ))))))))
113   (available_actions agent-id)
114   )))
115 (define optimal_action_value (mem (lambda (agent-id curr_iteration gridworld state_x state_y)
116   (let ((actions_to_values (available_actions_to_values agent-id curr_iteration gridworld state_x
117     state_y)))
118   (max_cdr actions_to_values)
119   )))
120 (define MAX_ITERATIONS 20)
121 (define should_terminate (mem (lambda (agent-id gridworld state_x state_y)
122   (if (<= (value_function agent-id MAX_ITERATIONS gridworld initial_x initial_y) 0) true
123     (let ((location_type (get_gridworld_at gridworld state_x state_y)))
124     (let ((state_food_utility (food_utility agent-id location_type)))
125     (> state_food_utility 0))))))
126
127
128
129 (define optimal_policy_from_initial_state (mem (lambda (agent-id gridworld state_x state_y)
130   (if (should_terminate agent-id gridworld state_x state_y) ()
131     (let ((curr_optimal_action_value (optimal_action_value agent-id MAX_ITERATIONS gridworld state_x
132       state_y)))
133     (let ((curr_optimal_action (car curr_optimal_action_value)))
134     (let ((next_state (gridworld_transition gridworld state_x state_y curr_optimal_action)))
135     (let ((next_state_x (second next_state)))
136     (let ((next_state_y (third next_state)))
137     (let ((remaining_policy (optimal_policy_from_initial_state agent-id gridworld next_state_x
138       next_state_y)))

```

```

137   (cons curr_optimal_action remaining_policy)
138 )))))))))))
139
140 (define trajectory_from_initial_state (mem (lambda (agent-id gridworld state_x state_y)
141   (if (should_terminate agent-id gridworld state_x state_y) ()
142     (let ((curr_optimal_action_value (optimal_action_value agent-id MAX_ITERATIONS gridworld state_x
143       state_y)))
144       (let ((curr_optimal_action (car curr_optimal_action_value)))
145         (let ((next_state (gridworld_transition gridworld state_x state_y curr_optimal_action)))
146           (let ((next_state_location (first next_state)))
147             (let ((next_state_x (second next_state)))
148               (let ((next_state_y (third next_state)))
149                 (let ((remaining_trajectory (trajectory_from_initial_state agent-id gridworld next_state_x
150                   next_state_y)))
151                   (cons next_state_location remaining_trajectory))
152                 )))))))))))
153
154 (define optimal_policy (mem (lambda (agent-id gridworld initial_state_x initial_state_y)
155   (cons (pair 'start 'start) (optimal_policy_from_initial_state agent-id gridworld
156     initial_state_x initial_state_y))))))
157
158 (define optimal_trajectory (mem (lambda (agent-id gridworld initial_state_x initial_state_y)
159   (cons (get_gridworld_at gridworld initial_state_x initial_state_y)
160     (trajectory_from_initial_state agent-id gridworld initial_state_x initial_state_y))
161 )))
162
163 (define optimal_policy_with_trajectory (mem (lambda (agent-id gridworld initial_state_x initial_state_y)
164   (zip (optimal_policy agent-id gridworld initial_state_x initial_state_y) (optimal_trajectory
165     agent-id gridworld initial_state_x initial_state_y))
166 )))
167
168 (define get_terminal_goal_state (mem (lambda (agent-id gridworld initial_state_x initial_state_y)
169   (last (optimal_trajectory agent-id gridworld initial_state_x initial_state_y))))))
170
171 (define trajectory_has_location_type? (mem (lambda (agent-id location_type gridworld initial_state_x
172   initial_state_y)
173   (member? location_type (optimal_trajectory agent-id gridworld initial_state_x initial_state_y))
174 )))
175
176 (define policy_has_motion_type? (mem (lambda (agent-id motion_type gridworld initial_state_x
177   initial_state_y)
178   (let ((policy_motions (map (lambda (action) (first action)) (optimal_policy agent-id gridworld
179     initial_state_x initial_state_y))))
180     (member? motion_type policy_motions))
181 )))
182
183 (define policy_and_trajectory_has_motion_at_location? (mem (lambda (agent-id motion_type location_type
184   gridworld initial_state_x initial_state_y)
185   (let ((policy_motions (map (lambda (action) (first action)) (optimal_policy agent-id gridworld
186     initial_state_x initial_state_y))))
187     (let ((trajectory (optimal_trajectory agent-id gridworld initial_state_x initial_state_y)))
188       (let ((motions_at_locations (zip policy_motions trajectory)))
189         (member? (list motion_type location_type) motions_at_locations))
190       )))))))
191
192 (define motion_at_location? (mem (lambda (agent-id motion_type location_type gridworld initial_state_x
193   initial_state_y)
194   (let ((policy_motions (map (lambda (action) (first action)) (optimal_policy agent-id gridworld
195     initial_state_x initial_state_y))))
196     (let ((trajectory (optimal_trajectory agent-id gridworld initial_state_x initial_state_y)))
197       (let ((motions_at_locations (zip policy_motions trajectory)))
198         (let ((motion_at_location? (lambda (motion_type location_type)
199           (member? (list motion_type location_type) motions_at_locations))))
200           (member? (list motion_type location_type) motions_at_locations))
201       )))))))

```

```

184     motions_at_locations
185 ))))))
186 ;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;
187 ;; Derived predicates.
188 ;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;
189 (define action_id_gensym (make_gensym "action-"))
190 (define is_going_to_actions (mem (lambda (agent-id)
191     (let ((action_states (optimal_policy_with_trajectory agent-id gridworld initial_x initial_y)))
192         (let ((final_location (last (last action_states))))
193             (list (list
194                 (pair 'action_id (action_id_gensym))
195                 (pair 'action_subject agent-id)
196                 (pair 'action_predicates (list 'is_going (list 'to final_location)))
197                 (pair 'action_preposition 'to)
198                 (pair 'action_location final_location)
199             ))))))))
200
201 (define is_going_on_actions (mem (lambda (agent-id)
202     (let ((action_states (optimal_policy_with_trajectory agent-id gridworld initial_x initial_y)))
203         (fold (lambda (action_state these_actions)
204             (let ((action_location (last action_state)))
205                 (let ((action_manner (first (first action_state))))
206                     (let ((action_direction (cdr (first action_state))))
207                         (cons
208                             (list
209                                 (pair 'action_id (action_id_gensym))
210                                 (pair 'action_subject agent-id)
211                                 (pair 'action_predicates (list 'is_going action_manner action_direction (list 'on
212                                     action_location)))
213                                 (pair 'action_preposition 'on)
214                                 (pair 'action_location action_location)
215                             )
216                             these_actions)
217                         )))
218         (lambda () action_states)
219     )))
220 (define actions_in_scene (mem (lambda (agent-id) (concatenate (is_going_to_actions agent-id)
221     (is_going_on_actions agent-id)))))
222 (define is_action? (lambda (action action_predicate) (member? action_predicate (lookup action
223     'action_predicates))))
224 (define is_subject_of_action? (lambda (action entity) (eq?
225     (lookup action 'action_subject)
226     entity
227     )))
228 (define is_preposition_of_action? (lambda (action preposition) (eq?
229     (lookup action 'action_preposition)
230     preposition
231     )))
232 (define is_location_of_action? (lambda (action location) (eq?
233     (lookup action 'action_location)
234     location
235     )))
236 (define get_location (lambda (action)
237     (lookup action 'action_location)
238     ))
239

```

```

240 (define (exists_action agent-id predicate)
241   (some (map predicate (actions_in_scene agent-id))))
242
243 (define (get_actions agent-id predicate)
244   (fold (lambda (action these_actions) (if (predicate action) (cons action these_actions)
245     these_actions))
246     () (actions_in_scene agent-id))

```

Code Block 11: Generative domain theory for restaurant navigation domain. Generates agents with varying preferences in a gridworld environment. Also implements a value iteration-based planner directly in the Church code.

#### A.4.2 Translation examples for social reasoning domain

```

1 ;; Bob likes pizza.
2 (condition (> (restaurant_utility 'bob 'pizza) 0))
3 ;; Bob really likes pizza.
4 (condition (> (restaurant_utility 'bob 'pizza) 10))
5 ;; Bob does not like pizza, and he actually despises vegetables.
6 (condition (and
7   (< (restaurant_utility 'bob 'pizza) 0)
8   (< (restaurant_utility 'bob 'vegetarian) 10)
9 ))
10 ;; The pizza place is not open.
11 (condition (not (is_open 'pizza)))
12 ;; Condition on: Bob walked North on Danner.
13 (condition (exists_action 'bob (lambda (action)
14   (and
15     (is_subject_of_action? action 'bob)
16     (is_action? action 'is_walking)
17     (is_action? action 'north)
18     (is_preposition_of_action? action 'on)
19     (is_location_of_action? action 'danner))))))
20 ;; Does Bob like vegetarian food?
21 (query (> (restaurant_utility 'bob 'vegetarian) 0))
22 ;; Where is Bob going?
23 (query (get_actions 'bob (lambda (action)
24   (and
25     (is_subject_of_action? action 'bob)
26     (is_action? action 'is_going))))))
27 ;; Where will Bob go to for lunch?
28 (query (get_location (first
29   (get_actions 'bob (lambda (action)
30     (and (and
31       (is_subject_of_action? action 'bob)
32       (is_action? action 'is_going))
33       (is_preposition_of_action? action 'to))
34     )))))

```

Code Block 12: Translation examples for the social reasoning. These examples are concatenated with the social reasoning scenes generative model to produce the prompt used to generate new translations.

## B Open questions

### B.1 Syntactic bootstrapping

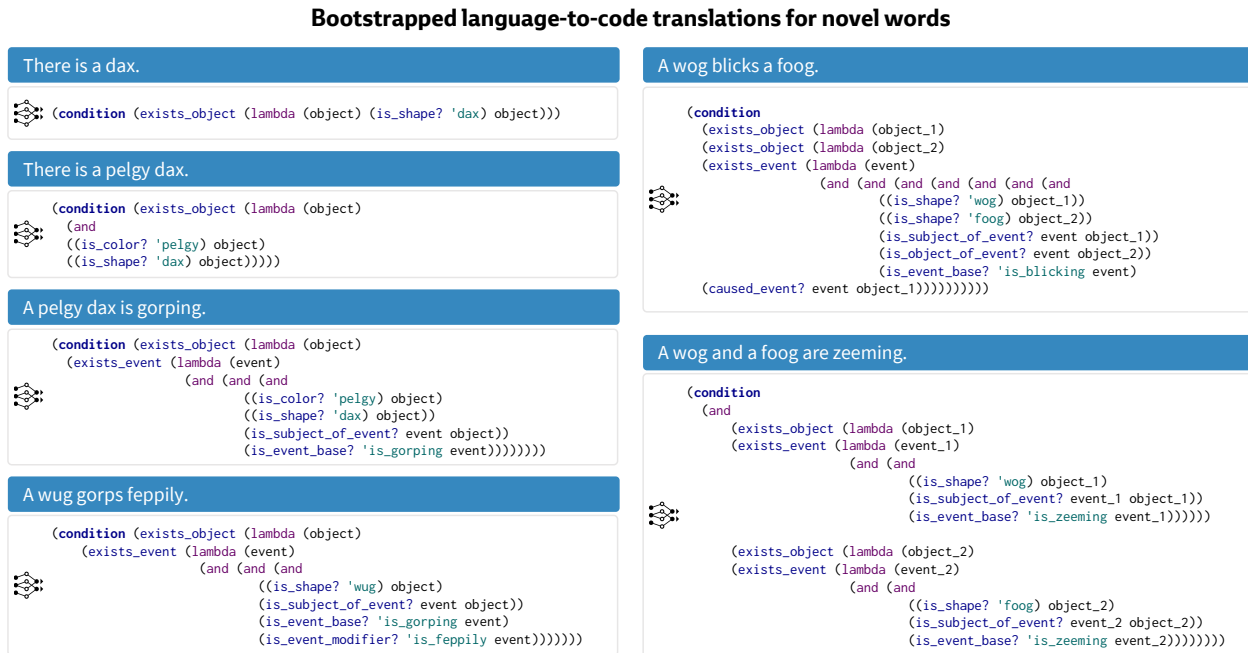


Figure 14: Example translations with novel words, suggesting that language-to-code models can leverage syntax-semantic mappings to inform hypothesized meanings.

### B.2 Code editing

While our framework focuses primarily on *generating* code in the PLoT, this view encompasses only part of the broader story of natural language. In particular, in certain contexts, it might not make sense to write new code, but instead to *modify* the existing domain theory. Consider the following statements, taken in context of the domains we explored in [Section 3](#):

- (Tug-of-war) *The team’s strength is the strength of the **strongest** player.*
- (Kinship) *Avery has two kids **from a previous marriage**.*
- (Visual scenes) *There’s a red mug **stacked on top of** a yellow can.*
- (Navigation) ***There’s a river** separating the North and South sides of town, which **you can paddle across** in nice weather.*

These utterances bend or break the rules of their respective domain theories. To properly integrate these kinds of language, we’d need to edit pieces of the existing generative models.

While language-guided code editing is still an open area of research, recent advances offer an exciting glimpse of what might be possible in the near-term. [Ouyang et al. \(2022\)](#) use a combination of finetuning and reinforcement learning to make GPT-3 adhere more closely to human-authored instructions. The resulting InstructGPT models, which OpenAI make available on their API, are capable of editing existing text based on short natural language instructions (e.g., “Fix the grammar”; “Turn this into a poem.”).<sup>9</sup> Excitingly, this same approach extends to code-based LLMs, meaning that it is possible to prompt GPT models to edit a piece of code according to some instructions. Indeed, we can use OpenAI’s editing interface off-the-shelf to handle utterances requiring localized changes to the domain model (see below for a simple example in the tug-of-war domain).

<sup>9</sup><https://openai.com/blog/gpt-3-edit-insert/>


**Redefine:** *The team's strength is the strength of the strongest player.*


<pre> 1 ;; The team's strength is the sum of the   players' strengths. 2 ;; When a player is lazy in a match, they pull   with half their strength. 3 (define (team-strength team) 4   (sum 5     (map 6       (lambda (player) 7         (if (flip (laziness player)) 8             (/ (strength player) 2) 9             (strength player))) 10      team))) </pre>	<pre> 1 ;; The team's strength is the strength of the   strongest player. 2 ;; When a player is lazy in a match, they pull   with half their strength. 3 (define (team-strength team) 4   (apply 5     max 6     (lambda (player) 7       (if (flip (laziness player)) 8           (/ (strength player) 2) 9           (strength player))) 10    team))) </pre>
---	---

Though questions of scaling and robustness remain, the problem of modeling sequences of code changes is currently gaining traction in the machine learning for code community, which has recently produced multiple language-guided neural models of code editing (Chakraborty, Ding, Allamanis, & Ray, 2022; Chakraborty & Ray, 2021; Fried et al., 2022; Panthaplackel, Nie, Gligoric, Li, & Mooney, 2020; Reid & Neubig, 2022; J. Zhang, Panthaplackel, Nie, Li, & Gligorić, 2022) that draw broadly on contemporary work in automated program repair (Bai et al., 2021; Y. Li, Wang, & Nguyen, 2020; Yasunaga & Liang, 2020). These advances suggest a broader vision for our framework in which domain theories, expressed in the PLoT, can be iteratively grown and revised to reflect natural language instruction. Moreover, as code LLMs become more general-purpose, the technical gap between generation and editing will continue to narrow, suggesting a point in the near future where defining new components of a domain theory will be a special case of language-guided code editing.

## C Attributions

### C.1 Attribution of graphics resources

 Artificial neural network icon by sachin modgekar from [thenounproject.com](https://thenounproject.com).

 Cog icon by Rflor from [thenounproject.com](https://thenounproject.com).