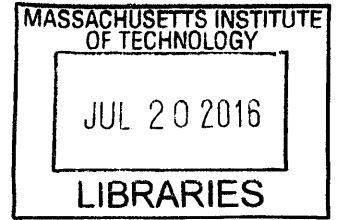


The inner life of goals: costs, rewards, and  
commonsense psychology

by

Jose Julian Jara-Ettinger



Submitted to the Department of Brain and Cognitive Sciences  
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2016

© Massachusetts Institute of Technology 2016. All rights reserved.

**Signature redacted**

Author .....

Department of Brain and Cognitive Sciences

May 14, 2016

**Signature redacted**

Certified by .....

.....  
Laura E. Schulz

Associate Professor

Thesis Supervisor

**Signature redacted**

Accepted by .....

.....  
Matthew A. Wilson

Director of Graduate Education for Brain and Cognitive Sciences



**The inner life of goals: costs, rewards, and commonsense  
psychology**

by

Julian Jara-Ettinger

Submitted to the Department of Brain and Cognitive Sciences  
on May 14, 2016, in partial fulfillment of the  
requirements for the degree of  
Doctor of Philosophy

**Abstract**

By kindergarten, our knowledge of agents has unfolded into a powerful intuitive theory that enables us to thrive in our social world. In this thesis I propose that children build their commonsense psychology around a basic assumption that agents choose goals and actions by quantifying, comparing, and maximizing utilities. This naïve utility calculus generalizes infants' expectation that agents navigate efficiently, and captures much of the rich social reasoning we engage in from early childhood. I explore this theory in a series of experiments looking at children's ability to infer costs and rewards given partial information, their reasoning about knowledgeable versus ignorant agents, and their reasoning about the moral status of agents. Moreover, a formal model of this theory, embedded in a Bayesian framework, predicts with quantitative accuracy how humans make cost and reward attributions.

Thesis Supervisor: Laura E. Schulz

Title: Associate Professor





## Acknowledgments

The work I present here belongs to many. The undergraduate students who I was lucky to work with -Eric Garr, Salvador Esparza, Jenny Yang, Aviana Polsky, Jessica Wass, Mika Maeda, Kristina Presing, Vivian Tran, Eileen Rivera, Diego Guerrero, Felix Sun, Anna Fountain, Emily Lydic, Sophie Cao, Lena Yang, Allison Kaslow, Christina Ma, Madeline Klein, Amy Zhang, and Mary DePascale-, the members of my committee -Nancy Kanwisher, Liz Spelke, and Rebecca Saxe-, my collaborators -Samantha Floyd and Hyowon Gweon-, and my advisors -Laura Schulz and Josh Tenenbaum.

Laura is the kind of advisor that makes you wish grad school were twice as long. I cannot express how fortunate I am to have had an advisor, a mentor, and a friend in her. After our first meeting, I remember walking out of her office feeling thrilled by her enthusiasm, creativity, and incisiveness. But I had no idea the same feeling would continue to appear even five years later. I can only hope to someday be the kind of person and scientist that Laura is.

As if having one amazing advisor were not enough, I have been privileged to have a second one. Josh is one the kindest and most generous people I have ever met. I learn something new from every conversation we have. His passion is contagious, and I'm continuously inspired by the depth of his ideas.

I had many other outstanding mentors. Thanks to Ted Gibson, Steve Piantadosi, Bevil Conway, Roger Levy, and Angela Yu for everything they've taught me. Thanks to the faculty who gave me interesting questions and challenges at different stages of this work -Susan Carey, Evelina Fedorenko, Susan Gelman, Allison Gopnik, Celeste Kidd, Henry Wellman, and Amanda Woodward.

Thanks to all the people who made Cambridge so fun - Max Siegel, Shaiyan Keshvari, Zenna Tavares, Galen Lynch, Praneeth Namburi, Laura Stoppel, Nick Dufour, Steve Voinea, Alex Paunov, Emily Mackevicius, Idan Blank, Omer Durak, Demi Duran, Wilma and Connie Bainbridge, Cristina Camayd, Hilary Richardson, Dorit Kliemann, Ben Deen, Rosa Lafer-Sousa, Kevin Woods, Wiktor Mlynarski, James

Traer, Eliza Kosoy, Jamie Rondeau, Anna Berglind, George Chen, Valerie Wallace, Melissa Troyer, Emilie Josephs, Sage Boechter, Sunoo Park, Jean Yang, Vicente Rodriguez, Richard Futrell, Kyle Mahowald, Leon Bergen, Kim Scott, Melissa Kline, Julia Leonard, Rachel Magid, Yang Wu, Nathan Winkler-Rhoades, Paul Muentener, and everyone I'm forgetting- and to those who have made me miss home so often- Lorena Lopez, Diego Flores, Azucena Rosales, Alan Carrillo, David Solis, Mauro Rodriguez, Marianne Tapia, Valentina Montes, Javier Benavides, Gabriel Zamora, and Sari Morales.

Finally, thanks to my family for what they've done for me, which is everything.

# Contents

<b>1</b>	<b>Introduction</b>	<b>11</b>
1.1	The inner life of goals . . . . .	11
1.1.1	Beyond goal attribution . . . . .	12
1.2	Thesis roadmap . . . . .	12
<b>2</b>	<b>The naïve utility calculus</b>	<b>15</b>
2.1	Commonsense Psychology . . . . .	15
2.2	Naïve utility theory: Agents as utility-maximizers . . . . .	18
2.2.1	Forward model: From costs and rewards to actions . . . . .	18
2.2.2	Reversing the model: From actions to costs and rewards . . . . .	18
2.2.3	Reasoning through a naïve utility calculus . . . . .	19
2.2.4	Real world reasoning with a naïve utility calculus . . . . .	19
2.3	The naïve utility calculus as a unifying theory in social cognition . . . . .	21
2.3.1	Goal-directed actions . . . . .	22
2.3.2	Sampling and preferences . . . . .	24
2.3.3	Communication and pedagogy . . . . .	25
2.3.4	Social and moral reasoning . . . . .	26
<b>3</b>	<b>Breaking into the inner life of goals</b>	<b>31</b>
3.1	Introduction . . . . .	31
3.2	Experiment 1 . . . . .	34
3.2.1	Methods . . . . .	35
3.2.2	Results and Discussion . . . . .	37

3.3	Experiment 2 . . . . .	39
3.3.1	Methods . . . . .	39
3.3.2	Results and Discussion . . . . .	41
3.4	Experiment 3 . . . . .	42
3.4.1	Methods . . . . .	43
3.4.2	Results and Discussion . . . . .	43
3.5	Experiment 4 . . . . .	45
3.5.1	Methods . . . . .	46
3.5.2	Results and Discussion . . . . .	46
3.6	General Discussion . . . . .	48
<b>4</b>	<b>Action under uncertainty</b>	<b>51</b>
4.1	Introduction . . . . .	51
4.2	Experiment 1 . . . . .	54
4.2.1	Participants . . . . .	55
4.2.2	Results and Discussion . . . . .	56
4.3	Experiment 2 . . . . .	56
4.3.1	Methods . . . . .	57
4.3.2	Results and Discussion . . . . .	58
4.4	Experiment 3 . . . . .	59
4.4.1	Methods . . . . .	59
4.4.2	Results and Discussion . . . . .	60
4.5	Experiment 4 . . . . .	61
4.5.1	Methods . . . . .	61
4.5.2	Results and Discussion . . . . .	62
4.6	General Discussion . . . . .	62
<b>5</b>	<b>Social reasoning</b>	<b>65</b>
5.1	Introduction . . . . .	65
5.2	Experiment 1 . . . . .	66
5.2.1	Methods . . . . .	67

5.2.2	Results and Discussion . . . . .	68
5.3	Experiment 2 . . . . .	69
5.3.1	Methods . . . . .	69
5.3.2	Results and Discussion . . . . .	70
5.4	Experiment 3 . . . . .	71
5.4.1	Methods . . . . .	71
5.4.2	Results and Discussion . . . . .	72
5.5	General Discussion . . . . .	73
<b>6</b>	<b>Formal implementation</b>	<b>79</b>
6.1	Introduction . . . . .	79
6.1.1	Social reasoning beyond goal attribution . . . . .	80
6.2	The naïve utility calculus . . . . .	82
6.3	Computational framework . . . . .	83
6.3.1	Naïve Utility Calculus model sketch . . . . .	83
6.3.2	Intermediate accounts . . . . .	84
6.3.3	Experiment . . . . .	86
6.3.4	Results . . . . .	89
<b>7</b>	<b>Unifying early social cognition</b>	<b>97</b>
7.1	Introduction . . . . .	97
7.2	The Naïve Utility Calculus . . . . .	98
7.2.1	Inferences from spatial information . . . . .	99
7.2.2	Inferences from statistical information . . . . .	100
7.3	Computational modeling . . . . .	100
7.3.1	Alternative models . . . . .	101
7.3.2	Naïve Utility Calculus models . . . . .	102
7.4	Experiment . . . . .	104
7.4.1	Methods . . . . .	104
7.4.2	Results . . . . .	105
7.5	Discussion . . . . .	109

<b>8 Discussion</b>	<b>113</b>
8.1 What develops? . . . . .	114
8.2 Outstanding questions . . . . .	117
8.2.1 Concluding remarks . . . . .	118
<b>A Meta-analysis for Chapter 2 experiments</b>	<b>121</b>
<b>B Adult survey on competence and niceness</b>	<b>127</b>

# Chapter 1

## Introduction

### 1.1 The inner life of goals

Our social skills are a strikingly different from those of other animals [61]. We reason about others in terms of mental states [158, 29], we help those in need [152], we cooperate in joint goals [54], we transmit knowledge efficiently [50] and make inferences by assuming others do the same [12].

Despite our natural pro-social tendency [109] we do not interact with all agents in the same way. From early on we are sensitive to the competence of teachers and distrust ignorant ones [77]. We distinguish between helpers, hinderers, and bystanders, and prefer the former [78, 57]. We extend this preference transitively, based on how agents interact with antisocial others [78, 134], and we punish transgressors, even at our own expense [33, 60, 68, 90].

Together, our understanding of others' minds, our tendency to act pro-socially, and our selectivity with whom we interact, allow us to reap the benefits of living in social groups (See also [24, 149]). These abilities are enabled by our commonsense psychology -an intuitive causal theory about agents, their thoughts, intentions, and personalities- [29, 42, 87]. Current work suggests that the origins of this theory are grounded on our understanding that agents have goals [160, 15] that they complete as efficiently as possible [36]. Through these assumptions, making sense of others' behavior reduces to finding target states that the agent is navigating towards in an

efficient manner [3, 4, 62]. This understanding is likely to be necessary for navigating the social world, but it is not sufficient.

### 1.1.1 Beyond goal attribution

Suppose you're in preschool and both you and your classmate Anna ask the teacher for help at the same time. The teacher looks at each of you for a second, and goes to help Anna. Her goal is clearly to help Anna. But why? Maybe the teacher likes Anna better. Or maybe the teacher likes you better and she's helping Anna first so she can spend more time with you later. More likely, the teacher doesn't have a favorite student and Anna needs more help than you do; or maybe Anna is doing something dangerous; or maybe the teacher happened to be closer to Anna; or maybe the teacher believes you're more patient. There are thousands of reasons for why the teacher chose to help Anna first, and making sense of the underlying causes is crucial for navigating the social world.

We face this kind of challenge everyday. And although we cannot always uncover why people do what they do, we are also not hopelessly blind to the causes behind their goals. Somehow, we infer goals and break them open to infer their inner structure. This thesis aims to answer how we do this.

## 1.2 Thesis roadmap

This work has three parts. The first part consist of Chapter 2 and focuses on theory. In this chapter I propose a theory -the naïve utility calculus- for how we reason beyond goal-attribution and I argue that this proposal enriches past accounts of social reasoning and captures a large set of empirical phenomena. In Chapters 3-5 I present empirical evidence testing the predictions of the naïve utility calculus. Chapter 3 shows how the naïve utility calculus enables us to reasong beyond goal attribution, Chapter 4 shows how naïve utility calculus supports reasoning about agents who do not know what they like, and Chapter 5 shows how the naïve utility calculus supports reasoning about social goals. Finally, Chapters 6-7 present a formal



account of the theory. In Chapter 6 I show how a formal computational model of the naïve utility calculus, embedded in a Bayesian framework, predicts adult judgments with quantitative accuracy, and in Chapter 7 I show how the same model can explain the success of other preference-inference models. I end with a brief discussion on the origins and development of the naïve utility calculus (Chapter 8).



# Chapter 2

## The naïve utility calculus

This chapter is based on Jara-Ettinger, Gweon, Tenenbaum, & Schulz. The naïve utility calculus (*in prep*).

### 2.1 Commonsense Psychology

Theories of decision-making have been at the heart of psychology since the field's inception, but only recently has the field turned to the study of how humans - especially the youngest humans - think humans make decisions. When we watch someone make a choice, we explain it in terms of their goals, their preferences, their personalities, and moral beliefs. This capacity - our commonsense psychology - is the cognitive foundation of human society. It lets us share what we have and know, with those from whom we expect the same in return, and it guides how we evaluate those who deviate from our expectations.

The representations and inferential power underlying commonsense psychology trace back to early childhood -before children being kindergarten, and often even in infancy. Work on how children reason about other agents' goals [11, 159, 160, 139, 92, 85, 36, 38], desires [80, 157, 110] beliefs [101, 2, 154, 155, 156, 142, 104], and pro-social behavior [23, 37, 74, 78, 57, 152, 153, 107, 134, 55, 53, 138, 148] has advanced our understanding of what's at work in early infancy [15, 143, 93] and what develops [155, 154, 42, 44, 43]. Nonetheless, major theoretical questions remain

unresolved. What computations underlie children's social judgments? Are there a small number of general principles by which humans reason about and evaluate other agents, or do we instead learn a large number of special case rules and heuristics? To what extent is there continuity between the computations supporting social cognition in infancy and later ages? Is children's social-cognitive development a progressive refinement of a computational system in place from birth, or are there fundamentally new computational principles coming into play?

In this thesis we advance a conjecture about the fundamental computational character of human social cognition that offers answers to each of these questions, and provides a unifying framework in which to understand the diverse social-cognitive capacities we see even in young children. We propose that human beings interpret others' intentional actions through the lens of a naïve utility calculus: that is, people assume that others choose actions to maximize the rewards they expect to obtain relative to the costs they expect to incur. We argue that this principle is at the heart of social cognition, possibly from early infancy. It can be made precise computationally and tested quantitatively (see chapters 6 and 7). Embedded in a Bayesian framework for reasoning under uncertainty, and supplemented with other knowledge children have about the physical and psychological world (e.g. knowledge about objects, forces, action, perception, goals, desires, and beliefs), the naïve utility calculus supports a surprisingly wide range of core social-cognitive inferences. It persists stably in some form through adulthood, guiding the development of social reasoning even as children's thinking about others undergoes significant growth.

Figure 2-1 illustrates some of the basic social intuitions that go beyond goal attribution which the naïve utility calculus aims to explain. These examples illustrate the same abstract principles that drove our intuitions in the cookie scenario above, but the principles are not specific to cookies and children; they apply to a wide range of situations in which intentional agents of any sort (child, adult, animated ball) interact with each other and move toward, reach for, or manipulate objects. We focus our discussion on behaviors such as those shown in Figure 2-1 where even young children can immediately grasp the costs and rewards involved: the naïve utility calculus

likely applies to more abstract situations as well, but its application may be complex in ways we do not consider here (e.g., cases where cultural norms are in play). Also although we focus on intentional behavior (as opposed to habits, reflexes, accidents, etc.) some of the most revealing choices are decisions not to act; our proposal aims to account for these as well.

The ideas behind this naïve utility calculus have a long history, tracing back to classical philosophers like Adam Smith [135] and John Stuart Mill [94]. Its formulation as an intuitive theory was anticipated in some form by pioneers in social cognition Fritz Heider [59], Harold Kelly [72], and Roger Brown [14], but with the development of new computational cognitive modeling tools these ideas can be formulated and tested more precisely [45, 48, 3, 147, 84, 39].

Critically, the naïve utility calculus is not a scientific account of how people act; it is a scientific account of people’s intuitive theory of how people act. These two notions may diverge - indeed, the mathematics of utility theory was originally proposed by early economists [151] but fails to predict actual human behavior in many important economic contexts ([1]; see [69] for review). But this does not mean the naïve utility calculus is not in some sense a reasonable and useful model of human behavior. In physics, our intuitive theory is oversimplified with respect to how the physical world actually works [118], yet it still helps us navigate everyday life because it tends to support accurate predictions on the spatial and temporal scales that matter most to us [6]. Similarly, the naïve utility calculus does not require that agents actually compute and maximize fine-grained expected utilities in order to be a useful guide in many everyday social situations.

In the remainder of this article we first describe the crucial ideas of the naïve utility calculus in their simplest, most ideal form. Next we move on to more nuanced features of the intuitive theory necessary to apply it to real-world decision making. We follow by reviewing studies that directly test the proposal, as well as the broader literature on goal-directed action, sampling-sensitive and preference judgments, communication, pedagogy, and social and moral evaluation that can be explained by our framework. We conclude with a discussion of how the naïve utility calculus relates to accounts

of first-person decision-making, and the proper relationship between intuitive and scientific theories of intentional action.

## **2.2 Naïve utility theory: Agents as utility-maximizers**

We propose that humans reason about behavior through the assumption that agents associate a utility with each possible goal, and then pursue the goal with the highest utility. This understanding takes the form of a generative model, which, embedded in a Bayesian framework, supports predictions about future behaviors (running the model forwards) and inferences about the causes of observed behaviors (working backwards from the model output through Bayes' rule). A formal version of the proposal is presented in Chapters 6 and 7.

### **2.2.1 Forward model: From costs and rewards to actions**

When agents decide how to act (e.g., whether to pursue a goal or which goal to pursue), they estimate the expected utility of each goal. Each goal's utility is calculated by estimating the rewards the agent would obtain if she completed the goal, and subtracting the cost she would need to incur to complete the goal. Through this process, agents build a utility function that maps possible plans onto expected utilities. Agents then pursue the plan with the highest positive utility. That is, agents are only willing to pursue plans where the rewards outweigh the costs. As such, if a plan has negative utility, the agent will be unwilling to act upon it, even in the absence of any alternatives.

### **2.2.2 Reversing the model: From actions to costs and rewards**

By assuming that agents behave in accordance with the forward model, observers can work backwards to infer the set of costs and rewards, which could have generated the observed behavior. In line with previous work (e.g., [3]), we propose that this inference is done through Bayes' rule (see Chapter 6).

### 2.2.3 Reasoning through a naïve utility calculus

The naïve utility calculus (an intuitive theory of agents as described in the forward model and a way to reverse this model) makes some key predictions about how humans reason about others' behavior, some of which are shown in Figure 2. These predictions not only involve people's qualitative judgments but also their confidence, supporting inferences about the ambiguity of exact rewards and costs underlying others' behaviors. For simplicity, here we only consider the rewards associated with outcomes and the costs associated with sequences of actions; as we note below however, an outcome can be costly and the sequence of actions be rewarding, too.

### 2.2.4 Real world reasoning with a naïve utility calculus

Reasoning about decision-making in the real world, however, has several complications that the idealized naïve utility calculus cannot handle. These complications reveal more sophisticated aspects of the naïve utility calculus that give it traction and point to ways in which commonsense psychology may develop (see Chapter 8).

#### **Expected utilities.**

Rather than choosing the options with the highest utilities, agents select the option with the highest expected utility. In familiar scenarios, agents should make accurate estimates. Most people, for instance, can estimate their costs for walking a block and their rewards from eating a cookie. However, agents often pursue novel outcomes in novel ways. In these contexts, it is critical that learners understand that agents act based on the expected rather than true costs and rewards. Learners should be less likely to infer that agents' choices are stable if they have reason to believe the agent was ignorant or mistaken about the true costs or rewards of her actions (Figure 2-1g and 2-1h; see also Chapter 4).

### **Individual differences.**

For most agents, two cookies are better than one and longer distances are costlier to travel than short ones. However, some people find walking more difficult than others, and some people like cookies more than others. These individual differences may be observable or may have to be inferred as part of explaining an agent's actions, similar to classic attribution theories [71]. By integrating both agent-invariant (objective) and agent-dependent (subjective) aspects of costs and rewards, the naïve utility calculus allows learners to parcel out known agent-invariant contributions to how an agent acts in a given situation and thereby infer latent costs and rewards that differ across agents (see Chapter 3).

### **Recursive content.**

Crucially for social cognition, an agent's costs and rewards can depend recursively on their expectations about another agent's costs and rewards. If someone is motivated to help, her rewards depend on promoting the other person's utilities, or diminishing them if she is motivated to hinder [150]. Likewise, acting against what you know another agent wants you to do may impose a cost. By integrating an agent's own first-order (self-interested) costs and rewards with that agent's second-order appreciation of others' costs and rewards, the naïve utility calculus allows observers to make inferences about the nature and extent of others' prosocial or altruistic tendencies.

### **Stochastic approximation.**

Observers should assume not only that agents act on expected rather than actual utilities, but also that agents have to estimate their own expected utilities and that these estimates may be inexact. When two plans have very different expected utilities, it is easy to identify the better plan. However, when plans have similar expected utilities, agents may find it more difficult to decide which is best - even apart from any uncertainty in their basic costs and rewards. This assumption provides flexibility in observers' inferences, softening the assumption that choices unambiguously reveal



the highest expected utility. It could also allow observers to infer agents' costs and rewards from the dynamics of their decision-making: agents are more likely to deterministically choose one plan over another when their utilities are very different, and more likely to oscillate between their choices when the utilities are similar.<sup>1</sup>

### **Inductive biases.**

When agents act they may incur costs for the actions and obtain rewards for the outcome; they may obtain rewards for the actions and incur costs for the outcome; or they may obtain rewards for both the actions and the outcome. The naïve utility calculus supports all of these representations. This flexibility implies that behaviors have multistable cost-reward decompositions. As in other domain of cognitions (e.g., [88]), and consistent with the Bayesian framework [47], this challenge can be solved through inductive biases that dictate that actions are more likely to be costly and that outcomes are more likely to be rewarding. As such, we naturally parse plans as being constituted of costly actions and rewarding outcomes, but the sources of costs and rewards can get overridden when the favored explanations are unable to account for the behavior (e.g., ascribing rewards to actions [124]).

## **2.3 The naïve utility calculus as a unifying theory in social cognition**

Beyond the work presented in this thesis, ascribing a naïve utility calculus to people has implications for a wide array of other phenomena in social cognitive development. As noted, researchers have looked extensively at children's intuitions about agents' goal-directed actions, desires and beliefs, pro-social behavior, and teaching and learning from others. Each of these aspects of social cognition has typically been treated as a separate problem, and explored through different paradigms. However, findings

---

<sup>1</sup>Additionally of course, the observers themselves have only sparse, noisy data about the observed agents; here however, we focus no on the observers' uncertainty but on the fact that the observers should assume that other agents act under uncertainty.

in many of these areas can be unified under the assumption that humans predict and explain behavior through a naïve utility calculus, as we illustrate below.

### 2.3.1 Goal-directed actions

A large body of work in cognitive development suggests that even infants expect agents to complete their goals as efficiently as possible [132, 127, 36, 13, 38] If for instance, infants are habituated to one agent hopping over a barrier to reach another agent, infants look longer when the agent continues to hop in the absence of a barrier than when she moves in a straight line [36].

These inferences have been proposed to be supported by a “teleological stance” [36], a non-mentalistic representation of behavior where agents are assumed to move efficiently towards goal-states, subject to situational constraints. The teleological stance is thought to underlie infants’ earliest forms of reasoning about agents and to serve as the basis for mentalistic representations that emerge later in life. The teleological stance is compatible with the naïve utility calculus: if agents maximize utilities, they should incur the minimum costs necessary to obtain rewards. However, the naïve utility calculus expands on the teleological stance by explaining how agents select their goals, and by explaining how objective (e.g., walls) and subjective (e.g., competence) constraints not only influence goal-completion, but also goal-formation.

Is it possible that infants merely expect agents to take the shortest possible path to a goal, without an abstract representation of costs or an expectation that agents should minimize them? Several studies suggest that infants represent efficiency in terms of relative costs that go beyond simply computing the length of the path. Southgate et al. [141] showed that infants appear to expect actions with fewer number of steps to be performed, over actions that take more steps or more time. Gergely et al. [35] showed infants an actor who used their head to light up a toy when their hands were either free or occupied. Infants themselves were more likely to imitate the head action in the hands-free condition compared to the hands-occupied condition, suggesting they inferred the actor had a specific intention (indicating a source of strong reward) to use their heads only when that was clearly the more costly of

available alternative actions. Together, these findings suggest that infants' expectation for efficient action may be driven by an abstract notion of cost-minimization. Nevertheless, experiments that directly pit a path's simplicity, straightness, length, time and energy costs against each other are needed to reveal if a general metric of cost minimization is at work in infancy, or if it arises later, building on top of some more limited, primitive notion of action efficiency.

More generally, a number of studies suggest that infants believe that the ability to perform effortful, high cost actions in the service of salient or plausible goals is the special provenance of agents (and only agents). Abilities attributed to agents (but not to objects or physical forces) include the ability to engage in self-generated movement [85, 108], the ability to resist gravity [81], the ability to cause objects to move or change state [123, 96, 120], the ability to create order [98], the ability to generate patterns [86], and the ability to generate events spontaneously [97] and probabilistically [161].

Such studies provide evidence that infants have intuitions about the costs of agents' actions. Other work suggests that infants also understand the rewards of goal-directed actions. Ten-month-olds appear surprised when an agent expresses a negative emotion following a completed (versus failed) goal [133], suggesting that they expect agents to find goal-completion rewarding. Ten-month-olds also attribute a preference to an agent who consistently chooses one goal over another (but not when the agent chooses the only option available; [74]) suggesting that infants understand that agents can find some goals more rewarding than others. By 18-months, children also understand that different agents can find the same goal (e.g., broccoli) more or less rewarding [110].

Collectively these results suggest that at least many key prerequisites to a naïve utility calculus emerge early in development: an expectation that agents act efficiently, an expectation that agents (and only agents) can perform effortful actions, and an expectation that agents experience subjective rewards consistent with goal outcomes.

### 2.3.2 Sampling and preferences

Infants as young as six-months expect randomly sampled sets, but not deliberately selected sets, to be representative of the population from which they are drawn ([27, 28, 30, 162, 163]). This sensitivity to the sampling process can support learning properties of novel objects [52]. Given a box that contains blue and yellow toys and an agent who pulls out three blue toys to demonstrate that all three have an interesting hidden property (e.g., squeaks), 15-month-olds draw different inferences about the property of the undemonstrated yellow toy depending on the ratio of the blue and yellow toys in the box; if the yellow toys were common in the box, they infer that only blue, but not yellow toys, have the novel property. In the absence of a clear purpose behind an agent’s sampling actions, infants attribute preferences [80]. For instance, if an agent pulls three frogs in a row from a box that contains mostly ducks, 20-month-olds infer that the agent prefers frogs to ducks; they do not infer this if the box contains more frogs than ducks or if the box contains only frogs.

The intuitions underlying toddlers’ and infants’ sensitivity to the sampling process can be explained through the naïve utility calculus. If all the objects in a box are equally rewarding, then agents should minimize costs by taking the objects that are the easiest to reach, generating a sample representative of the population. However, if one type of object is more rewarding than the others, then the agent should selectively draw that kind of object even if it is difficult to obtain, generating a biased sample. Reversing these inferences, if an agent generates a sample that could have been obtained purely by minimizing costs, her actions suggest that all the objects are equally rewarding or that she was ignorant of the rewards. However, if generating the sample required the agent to perform relatively costly actions (in time, effort, and attention), the rare objects must have been more rewarding (e.g., because these objects are preferred by the agent or have a special property). The reverse of these inferences is consistent with studies that show infants and children infer the presence of agents (rather than physical forces) when physical objects are sampled or arranged in order [86, 98].

### 2.3.3 Communication and pedagogy

Researchers have suggested that pedagogical communication is supported by the assumption that both teachers and learners expect teachers to share information sufficient for learners to infer the intended hypotheses [12, 129, 130, 111]. This assumption means that in pedagogical contexts, children are likely to infer that what they are taught is exhaustive. If for instance, a teacher shows a child that a toy squeaks, children learn that the toy squeaks, but also infer that it doesn't do anything else (and thus explore less than they do at baseline or in non-pedagogical contexts; [12, 131, ?]). Moreover, if children themselves know the true state of the world, they recognize informants who do not provide exhaustive information, rate the informant poorly and compensate under-informative teaching with additional exploration [50].

The naïve utility calculus provides a principled explanation for how the assumptions underlying pedagogical communication emerge. If the teacher has information that she believes will be useful to the learner and can be shared at a low cost, the agent should share it in order to maximize her utilities. (That she is sharing information at all implies that sharing the information is more rewarding than costly.) As such, in small and simple domains (e.g., a toy with just a few functions) where the cost of sharing information is negligible, agents should share all the information necessary for the learner to draw accurate inferences, and observers can use this expectation to make inferences accordingly.

The naïve utility calculus has a number of other implications for cooperative communication. For instance, the costs and reward of information should affect what and how much information teachers are expected to share. We should be less surprised if teachers fail to provide exhaustive evidence about a toy with many functions than a toy with only a few. Similarly, if a toy has many equally rewarding functions but some are costlier to demonstrate than others, we should be less surprised if the teacher fails to share high-cost information than low-cost information. More generally, if costs are negligible, learners should expect informants to communicate all relevant information, but the expectation should be weaker when the costs are higher. Consistent with

this, children prefer exhaustive informants when costs are low but prefer informants who provide only information sufficient for a good inductive inference when costs are high [51]. In linguistic communication more broadly, the classic Gricean maxims [46] calling on good communicators to be both informative and concise, which are central in pragmatic inferences for both adults and children [34, 144], can be derived from the naïve utility calculus. Minimizing utterance length (communication costs) while maximizing information transfer to the listener (communicative rewards) can be seen as optimizing an overall utility function trading off these costs and rewards.

Finally, as noted, the rewards of effective communication can be recursive: how rewarding it is to provide information may depend on how rewarding it is for the learner to learn it. Consistent with this, a number of studies suggest that very young children go out of their way to communicate information that is currently unknown to the learner [100, 82, 89], relevant to the learner’s goal [49], or difficult for the learner to discover by herself [114].

### **2.3.4 Social and moral reasoning**

Many studies have suggested that social evaluation emerges in the first year of life, with infants preferring agents who help others achieve their goals to those who hinder those goals [57, 78]. Moreover, infants’ evaluation of others’ actions is sophisticated. Infants’ evaluations are transitive (they prefer agents who hinder hinderers and help helpers; [58]) and they only positively evaluate agents if the agent actually caused the helpful action, did so intentionally, and did so with knowledge of the recipients’ preferences [55, 74]. Such studies are consistent with a naïve utility calculus: in every case, the helper or hinderer takes costly actions (i.e., goes out of his/her way to intervene), supporting the inference that the goal (helping or hindering the other) must be rewarding. Our own work suggests that toddlers also use agents’ relative costs to distinguish their motivations: if someone refuses to help when helping is costly, two-year-olds think she is nicer than a more competent agent who refuses to help at low cost (see Chapter 5).

The naïve utility calculus has many other, untested, implications for social eval-

uation. Consider for instance, that agents who underestimate rewards or costs may be more liable to abandon plans or commitments, with consequences for how others' judge them and whether they trust them in the future. It is also noteworthy that there is a special category of moral blame ("exploitation") for those who knowingly take advantage of others' ignorance of their utilities; it is unethical to convince someone to commit to an action when you know their expected reward is too high and/or their expected costs are too low. By the same token, agents with selective knowledge of their utilities can incur special moral credit or blame: It is particularly admirable to commit to a helpful action when you are ignorant of any extrinsic reward; it is particularly heinous to knowingly perform a costly (e.g., planned and premeditated) harmful action. In short, a wide range of intuitions underlying our judgments of others' competence and values involve considering how agents' might maximize their utilities given subjective and objective elements of costs and reward. In this way, a naïve utility calculus may play a critical role in social evaluation broadly.

In the following three chapters we present empirical tests of the crucial aspects of the naïve utility calculus. In Chapter 3 we show that children expect agents to maximize utilities, that they understand that cost and rewards vary across agents, and that they can use this intuitive theory to infer an agent's costs and rewards. In Chapter 4 we show that children understand that agents act to maximize expected utilities and not true utilities. Finally, in Chapter 5 we show toddlers use their naïve utility calculus to reason about social goals.

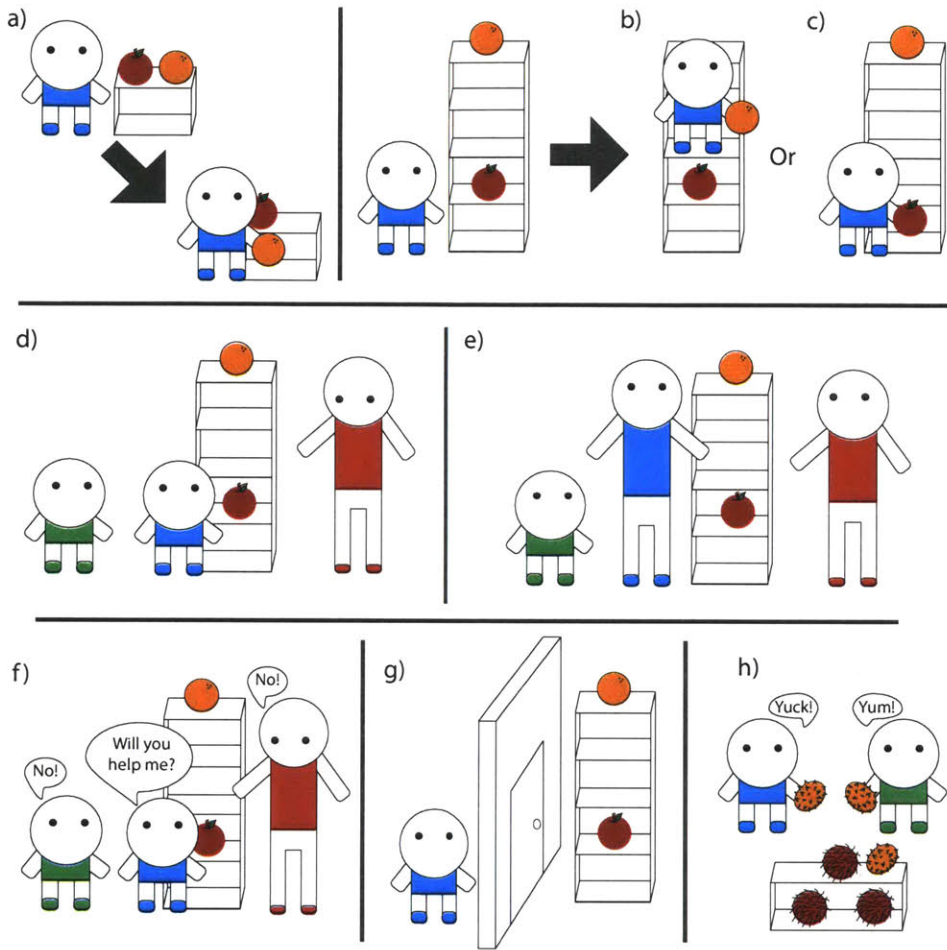


Figure 2-1: The logic of costs and rewards underlying our commonsense psychology. (a) If the blue agent (the protagonist) chooses the orange over the apple, her immediate goal is clearly the orange, but how confident are you that she prefers oranges in general to apples? (b) If the orange were high on the top shelf and the agent climbs up to get it, would you become more confident she prefers oranges in general? (c) What if she had chosen the apple instead? Does this indicate any strong preference for apples? (d) If the blue agent wants the orange from the top shelf, whom should she ask for help? (e) If she is the tallest person in the room, is it still appropriate for her to ask for help? (f) If both the red and green agent refuse to help, are they equally mean or is the red one meaner? (g) If the blue agent cannot see the shelf and says she is going to get the orange, are you confident she won't change her mind? (h) If both agents choose kiwanos over rambutans, but one says "Yum!" and the other says "Yuck!" after tasting it, who's more likely to have never tasted the fruits before?



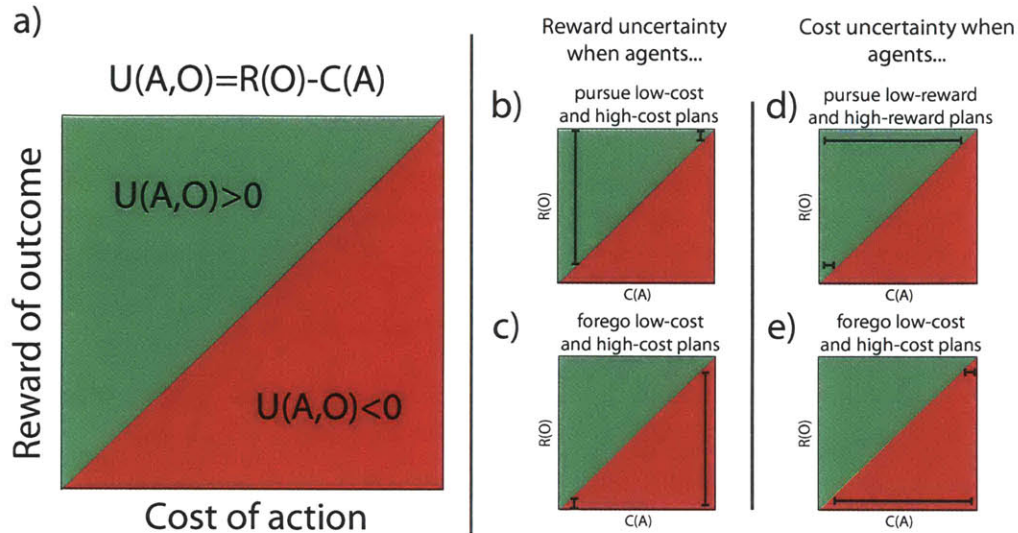


Figure 2-2: Visual schematic of utility maximization in simplified scenarios where the agent obtains a reward after completing the goal, and incurs costs towards achieving the outcome. (a) Qualitative plot of the utility function for a single goal. The utility is positive when the reward for reaching an outcome is higher than the costs the agent must incur. The utility is negative otherwise. (Green indicates the agent will act and red that she will not.) (b) If an agent completes a low-cost plan, a wide range of rewards could have produced a positive utility, thus low-cost actions do not reveal much about the agent’s reward. However, if the agent completes a high-cost plan, we can be confident that the reward was high enough to produce positive utility. (c) The structure of these inferences flips when the agent refuses to pursue a plan. For low cost plans, the rewards must be low to make the net utility negative, thus we can be confident that a decision to forego acting indicates a small reward. By contrast, a range of rewards is consistent with refusing to pursue a costly plan, so if the agent refuses to pursue a high cost plan, her refusal is not very informative about her rewards. (d) The implications are parallel when we infer costs given reward knowledge. Low rewards only motivate action when the costs are low; high rewards motivate action under a wide range of costs (e) If an agent foregoes a low reward we may be uncertain about the costs of acting but if she foregoes a high reward, we can be certain the costs were high.



# Chapter 3

## Breaking into the inner life of goals

This chapter is based on Jara-Ettinger, Gweon, Tenenbaum, & Schulz. Children’s understanding of the costs and rewards underlying rational action (2015). *Cognition*.

### 3.1 Introduction

Although the computational details of the naïve utility calculus are not critical here (see Chapters 6 and 7), it is helpful to consider the qualitative intuitions behind them. Intuitively we can think of an agent’s utility function as the difference between two terms: a (positive) reward term associated with goals to be achieved, measuring the value of a goal to the agent, and a (negative) cost term associated with actions that can be taken to achieve these goals, measuring the difficulty of these actions. Formally, we can decompose the utility function into a reward associated with the goal state, and a cost associated with the necessary actions:

$$U(a, s) = R(s) - C(a)$$

If we see an agent take actions  $a$  to reach state  $s$ , we can conclude that the reward for  $s$  is significantly higher than the cost of  $a$ ; but we cannot determine their exact values. Consider Sally, who climbs over a wall to get a cookie. Sally’s goal is clearly to get the cookie, and the reward for getting the cookie must be higher than the cost

for climbing over the wall. However, we cannot tell if she chose to jump because she likes the cookies so much that the cost of climbing was well worth it (a high cost/high reward plan), or because the obstacle is so trivial that it is worth surmounting even for a relatively mediocre cookie (a low cost/low reward plan). Although the underlying costs and rewards are irrelevant when interpreting Sally's goal (getting a cookie), high cost/high reward plans are psychologically very different from low cost/low reward ones. If Sally incurred a high cost for the cookie, she will likely also try to get a cookie in other situations where the costs are lower. However, if Sally incurred a low cost for the cookie, she may forego the cookie when the cost increases. Intuitively, inferring Sally's utility tells us her goal, but understanding the underlying costs and rewards allow us to understand Sally's capabilities and motivations and to predict her future behavior.

Importantly, costs and rewards have both external and internal components. Some aspects of costs and rewards apply across agents: Climbing over a high wall is always more costly than climbing over a low one, and two cookies are typically more rewarding than one. However, other aspects of costs and rewards differ across agents: Some people find climbing harder than others and some people like cookies better than others. Such intuitions motivate an account of rational action that considers not just those aspects of the event that are constant across agents, but also those that vary between them: not just the height of the obstacle but Sally's competence to surmount it, and not just how many cookies Sally gets but how much Sally values them. We suggest that we naturally understand agents' actions and goals in terms that go beyond a simple maximization of the overall utilities. Instead, we reason about the costs and rewards that form the utility function (See also [63]).

Note that this understanding of rational action requires more sophisticated reasoning than merely an understanding of goal-directed actions. As researchers have noted, reasoning about the goals of an agent's actions does not even necessarily require theory of mind. A learner could infer that actions are (or are not) efficient with respect to a goal and environmental constraint without imputing any mental states to agents (see [25, 23, 36] on the "teleological stance" in understanding rational action).

Such non-mentalist inferences may indeed underlie some successes at social cognition very early in infancy (e.g., [132, 139]). Researchers have similarly proposed that early studies of theory of mind (e.g., [101, 142, 140]) rely on implicit knowledge distinct from the explicit representations that emerge later in development (e.g., [105]).

However, some aspects of the naïve utility calculus might require representations more sophisticated even than many tasks that clearly do require explicit theory of mind. We might predict, for instance, that children should come to understand not only that one agent likes something and another agent does not (see e.g., [110]) but that two agents can like the same thing to different degrees. Similarly, children should come to understand not only that one agent can perform an action and another cannot ([149]) but that two agents might perform (or fail to perform) the same action and yet incur, or expect, different costs (i.e., because one agent is more competent than another). Moreover, children should be able to infer differences in agents' subjective rewards from information about differences in their subjective costs (and vice versa), even in the absence of any explicit behavioral cues indicating that one agent is more motivated, or more competent than another, and in contexts where agents have identical epistemic access and face identical situational constraints.

In short, our study looks not at children's understanding of rational action or theory of mind in general, but at children's ability to infer unobservable individual differences among agents. To date, most research on children's reasoning about individual differences has occurred in the context of children's understanding of personality traits (e.g., that someone is "lazy" or "shy"). Research suggests that such trait understanding does not develop until seven or eight [8, 70, 112, 117, 115, 116] although more recent work suggests that it may emerge by age five [83, 128]. While our study does not require children to treat individual differences as enduring, stable traits, we do require children to reason about the unobservable, internal structure of goal-directed actions as a consequence of individual differences between agents. Additionally, our study requires children to impute different mental states to agents given identical evidence about their behavior. The ability to understand that agents who perceive ambiguous evidence might interpret it differently is also a relatively late

development [16, 19]. Thus here we focus our investigation on five and six-year-old children.

Here we test three key implications of a naïve utility calculus. First, children should understand that costs influence an agent’s choices. That is, agents do not always pursue the states with the highest rewards because obtaining those states might also involve high costs; rational agents should maximize utilities rather than rewards. We test this understanding in Experiment 1 by looking at whether children can accurately infer an agent’s subjective rewards (preferences) from the choices she makes by considering the relative costs of her choices. Second, children should understand that both rewards and costs vary across agents, are not directly observable, and differ from situational constraints that uniformly affect all agents. In Experiment 2, we test this by introducing two agents who have different preferences but make identical choices; we look at whether children can use information about agents’ preferences and choices to infer differences in their competence. Finally, children should be able to predict how changes in costs and rewards affect an agent’s actions. We ask whether children can manipulate the position of objects, or the role of agents, in order to gain information about agents’ competencies. We test this in Experiments 3 and 4.

## 3.2 Experiment 1

In Experiment 1 we look at whether children understand that an agent’s choices depend on the expected costs and rewards associated with an action. If children understand that agents act efficiently towards their goals but do not consider the associated costs and rewards, then they may fail to distinguish actions that maximize the rewards from actions that maximize the utility (the difference between costs and rewards). Suppose, for instance, that an agent reaches for a banana on one trial (when the banana is closer than a watermelon) and a watermelon on the next trial (when a watermelon and a banana are equidistant from the agent). If children simply believe that agents act to maximize rewards, then children might infer that the agent liked the watermelon and bananas equally because the agent engaged in efficient goal-directed

actions on both trials. If, however, children consider both the costs and rewards of actions, they should instead infer that the agent prefers the treat chosen when the costs were equivalent: the watermelon.

We test exactly this in Experiment 1. Children see a puppet choose between two kinds of treats on two consecutive trials. In one trial (different cost trial), the cost of getting each treat is different, and the puppet chooses the low-cost treat. In another trial (same cost trial), the cost of obtaining each treat is matched and the puppet chooses the treat it had previously foregone (trial order counterbalanced). If children are insensitive to costs and assume the agent is acting only to maximize his rewards, they should conclude that the puppet likes both treats equally; he chooses each treat once. If, instead, children take costs into account and expect the puppet to maximize utilities, then children should infer that the puppet prefers the high-cost treat even though he chooses it on only one of the trials.

### **3.2.1 Methods**

#### **Participants**

32 children (mean age: 5.85 years, range 5.0 – 6.9 years) were recruited at an urban children’s museum; one additional participant was tested but excluded from analysis and replaced due to interference from a sibling (See Results). Children were assigned to a test condition or control condition ( $n = 16$  per condition).

#### **Stimuli**

The stimuli consisted of a puppet (Ernie), a paper picture of a watermelon slice, a paper picture of a banana, and two cardboard boxes: a short box (30 cm high) and a tall box (62 cm high).

#### **Procedure**

3-1 shows the experimental setup. Participants were tested in a quiet room at the museum in the presence of their caregiver. The child and the experimenter sat on

opposite sides of a small table where the tall and short cardboard boxes were placed. In the test condition, the experimenter introduced Ernie and then directed the child's attention to the two boxes. Participants were asked which box was the hardest for Ernie to climb. Children who chose the short box were corrected ( $n = 5$ ). The experimenter then said, "It's easy for Ernie to climb the short box!" and had Ernie climb the short box swiftly and nod in agreement. Then the experimenter said, "It's hard for Ernie to climb the tall box. It makes him tired!" and had Ernie climb the tall box slowly, and running out of breath. Afterwards, the experimenter introduced the watermelon and the banana. The experimenter placed both treats on the short box. The experimenter had Ernie look at both treats and then choose the banana. The experimenter said, "When both treats are on the short box, Ernie always chooses the banana!" Next, the experimenter placed the watermelon on the short box and the banana on the tall box. The experimenter had Ernie look at both treats and then choose the watermelon on the short box. The experimenter said, "When the watermelon is on the short box and the banana is all the way up on the tall box, Ernie always chooses the watermelon!" The experimenter then placed both pictures on the table and asked, "Which treat does Ernie like the most?" Trial order and Ernie's preferred treat were counterbalanced throughout.

In our design, one treat was always placed on the low box, while the other moved from one box to the other. The control condition was designed to rule out the possibility that children might simply select the moving treat over the static treat one because the moving treat was more salient. The control condition followed the same logic as the test condition except that the treats were placed next to the boxes rather than on top and the experimenter substituted "next to" for "on" (e.g., "When both treats are next to the short box, Ernie always chooses the banana!") In this condition, both treats had equally low costs and the puppet chose each treat once; thus we expected children to perform at chance.



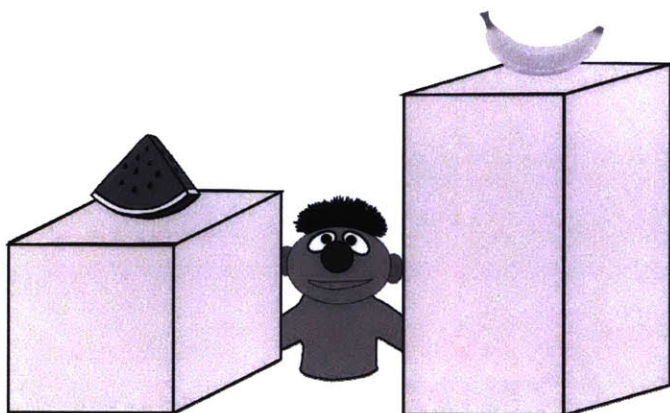


Figure 3-1: Example of experimental setup. All trials in all experiments consisted of a puppet choosing between two treats that could be placed either on the tall or the short box. We studied children's naïve utility calculus by varying the position of the objects, the puppet's choices, and the preference or competence information participants received.

### 3.2.2 Results and Discussion

All videotapes were coded by the first author for inclusion and children's responses to the test question; 100% of the videotapes were recoded on both measures by a second coder blind to hypotheses and conditions. The parents of two children did not consent to videotape and their responses were judged online. Children were excluded from analyses and replaced if the second coder judged A. that the items were not placed equidistant from the child or that the experimenter had otherwise cued the child's response (no children were excluded on these grounds) or B. if a parent or sibling interfered with the task ( $n = 1$ ). In the test condition, children were counted as succeeding on the task if they selected the treat that Ernie chose in the trial where both treats were equally costly to reach. Intercoder agreement on all measures was 100%. Consistent with recent concerns about null hypothesis testing (e.g., Cohen, 1994; Cumming, 2013) we report confidence intervals throughout and report exact  $p$ -values as a secondary measure (reporting one-tailed tests when directional predictions warrant).

Twelve of the sixteen children correctly selected Ernie's favorite treat (75%; 95%

CI: 56.25 – 100%); the remaining four children incorrectly selected the other treat. See 3-2. The results of the control condition suggest that these results were not due to children simply choosing the treat that moved locations. As expected, children in the control condition performed at chance: 7 of 16 children chose the treat that Ernie chose when both treats were by the short box (44%; 95% CI: 18.75 – 68.75%).

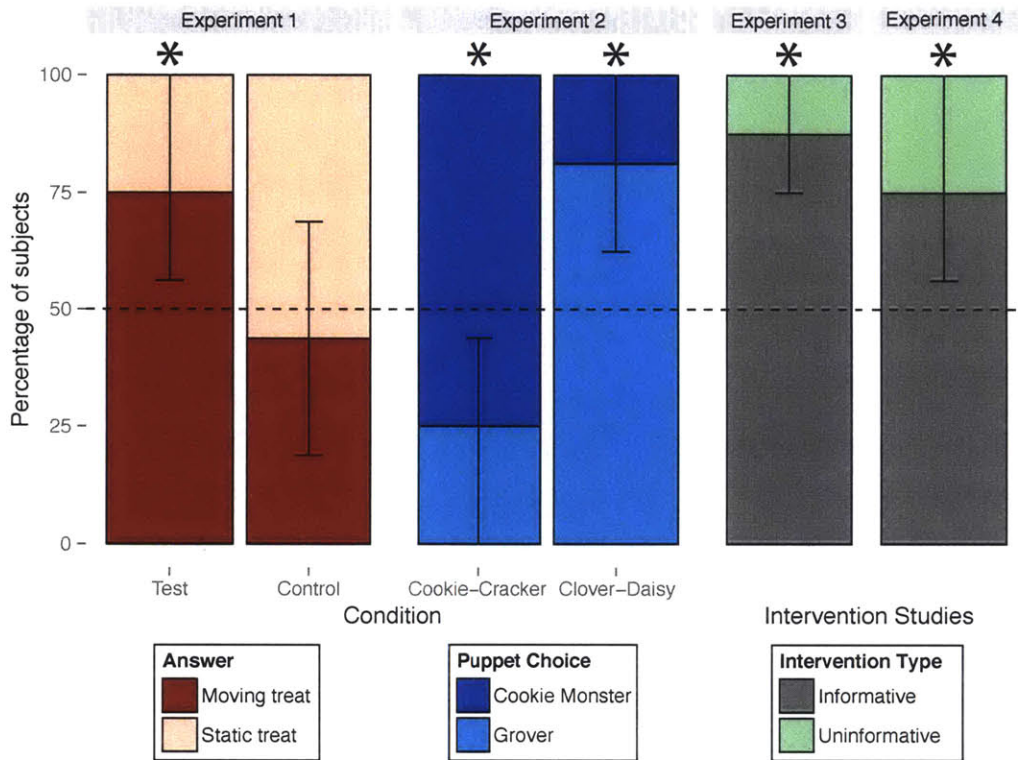


Figure 3-2: Results from all experiments. Each bar shows the distribution of participant’s responses across the four experiments. Vertical black lines show 95% confidence intervals bootstrapped from the data.

Note that if the children expected Ernie to always pursue the treat with the highest reward then their responses should have been equally split across the two treats in both conditions. However, even though Ernie chose both treats exactly once, children in the test condition successfully identified Ernie’s preferred treat, suggesting they considered the relative cost of his choices. These results suggest that children not only understand the agent-invariant cost of actions (i.e., that a tall box is harder to climb than a shorter one) but can integrate this information with the agent’s actions to infer unobservable mental states: the agent’s subjective rewards, or preferences.

## 3.3 Experiment 2

Experiment 1 suggests that children understand that agents maximize utilities and not rewards and thus costs influence agents' choices. However, this task only required children to understand the agent-invariant aspect of costs (i.e., that taller boxes are more costly to climb than shorter boxes). In Experiment 2 we look at whether children can use differences in agents' expected subjective rewards (their preferences) to infer agents' expected costs (their competencies).

In this task, children are introduced to two puppets. Each puppet can choose between a treat that is relatively costly to obtain and a treat that is relatively easy to obtain. One of the puppets likes both treats equally; the other puppet prefers the more costly treat. Children then see both puppets choose the less costly treat. Children are asked which puppet is unable to perform the high cost action. Critically, neither puppet ever attempts the costly action. Thus, in contrast to previous work (e.g., [7]) there are no behavioral cues to indicate whether agents are unwilling or unable to perform the action. However, if an agent assigns identical reward to both treats, he should never attempt the more costly action insofar as he is maximizing his utilities. By contrast, if an agent assigns a high reward to the outcome associated with the high cost action, he should attempt the action unless the costs exceed the reward. Thus the second agent's actions are more likely to be informative about the agent's subjective costs than the actions of the first agent. If children are able to infer expected costs based on expected rewards, they should infer that the puppet with the preference had difficulty performing the action.

### 3.3.1 Methods

#### Participants

Thirty-two children (mean age: 5.8 years, range 5.0 – 6.9 years) were recruited from a children's museum and randomly assigned to either the Cookie-Cracker condition or the Clover-Daisy condition (N= 16 in each condition). Four additional children were tested but excluded from analysis and replaced due to experimenter error (N = 2)

and interference from siblings ( $N = 2$ ).

## **Stimuli**

A Cookie Monster puppet and a Grover puppet were used. A short cardboard box (20 cm high) and a tall cardboard box (51 cm high) were used for the puppets to climb. Paper cutouts of cookies and crackers or clover leaves and daisy flowers were used for the Cookie-Cracker and the Clover-Daisy conditions, respectively. We also used two additional pictures for the Clover-Daisy condition: one of Grover surrounded by clovers and one of Cookie Monster surrounded by clovers and daisies.

## **Procedure**

Participants were tested in a quiet room at the museum in the presence of their caregiver. Children sat across the table from the experimenter. The two boxes were on the table. In the Cookie-Cracker condition, the experimenter showed the child paper cutouts of cookies and crackers and introduced the puppets. Children were told that Cookie Monster liked cookies better than crackers while Grover liked both treats equally (order counterbalanced). The preference information was repeated twice and children were prompted to ensure they remembered the information (e.g., “Remind me, does Cookie Monster like cookies? Yes, he loves cookies. And does he like crackers? Not so much.”). Children who gave wrong answers were corrected ( $N = 3$ ). Next, children were told that both puppets could climb the short box, but the tall box was so hard to climb that only one of the puppets could climb up to the top. Children were told that in order to find out which puppet was the better climber we would place treats on the boxes and let the puppets choose a treat. In the first trial, a cracker and a cookie were placed on the short box. Each puppet approached the short box individually (while the other puppet was absent), looked at both treats, and picked the cookie (order counterbalanced). In the second trial, the cracker was once again placed on the short box, but the cookie was now placed on the tall box. Once again, each puppet approached the boxes individually and looked at both treats, but this time both puppets picked the cracker. Children were then

asked, “Which puppet do you think is the one who cannot climb?”

Because children might think that Cookie Monster could not climb for reasons irrelevant to the experiment (e.g., because cookie eaters are unhealthy), the Clover-Daisy condition was set up such that Grover was the puppet who couldn’t climb. In this condition, Grover liked clovers better than daisies but Cookie Monster liked both equally. Although we chose clovers as the preferred stimuli for Grover hoping that children would easily associate the two (i.e., because Grover rhymes with clover), pilot data showed that children had a hard time remembering the puppets’ preferences. Thus we added a picture of Grover with clovers and Cookie Monster with both clovers and daisies to help children remember the puppets’ preferences. All other aspects of the two conditions were identical.

### **3.3.2 Results and Discussion**

All videotapes were coded blind to condition by the first author for inclusion and children’s responses to the test question; 62% of the videotapes were recoded on both measures by a second coder blind to hypotheses and conditions. The parents of two children did not consent to videotape and their responses were judged online. Intercoder agreement was 100%. In both conditions, children successfully used the preference information to make competence judgments. In the Cookie-Cracker condition, 12 of the 16 children correctly identified Cookie Monster as the incompetent puppet (75%; 95% CI: 56.25 – 100%). The remaining four children identified Grover as the incompetent puppet. In the Clover-Daisy condition, 13 out of the 16 children correctly identified Grover as the incompetent puppet (81.25%; 95% CI: 62.5 – 100%); the remaining three children identified Cookie Monster as the incompetent puppet. See 3-2.

Children’s ability to distinguish agents’ competencies here is especially striking because both puppets behaved identically: each puppet chose each treat once, and neither climbed the tall box. In fact, neither puppet even attempted to climb the tall box. Instead they always chose to climb the small box, and always succeeded in their actions. In order for children to draw different conclusions about the competence of

the two agents, children had to use the information about differences in the agents' subjective rewards to infer that the costs of climbing the tall box influenced the agents' choices. These results are consistent with our hypothesis that children evaluate agents through a naïve utility calculus that includes a principle of rational expectation.

We asked the children “Which puppet cannot climb?” (rather than, for instance, “Which puppet has more difficulty climbing?”) because what was at stake in this experiment was only children’s ability to use the expected reward information to distinguish the expected costs for two puppets. Even the information about the differences in the agents’ subjective rewards does not provide evidence about the agent’s absolute competence. That is, the utility functions do not distinguish a puppet that is completely unable to climb the tall box from a puppet that merely finds it very costly to do so. Critically however, children were able to recognize that the differences in agents’ subjective preferences could provide information about subjective costs for the puppet who preferred one treat to another, but was unlikely to be informative for the puppet who liked both treats equally. This suggests that children understand how graded differences in subjective rewards can provide information about agents’ subjective costs.

### 3.4 Experiment 3

Experiments 1 and 2 suggest that children are able to represent and infer agent-specific costs and rewards. In Experiment 3, we further investigate children’s understanding of agent-independent (external) and agent-dependent (subjective) costs by asking whether children can manipulate the external cost associated with obtaining different goals to gain information about an agent’s subjective costs. Children are given a high-reward and a low-reward treat and asked to place them in locations that incur different objective costs in order to learn an agent’s subjective expected costs. Intuitively, it is obvious that agents should choose low cost actions when they are associated with high reward; it is therefore more informative to see how agents behave in the context of high reward, high cost actions. If children understand how costs and rewards affect



an agent's choices, they should pair the high-reward treat with the high-cost location.

### **3.4.1 Methods**

#### **Participants**

Sixteen children (mean age: 6.0 years, range 5.1 – 6.8 years) were recruited at an urban children's museum and randomly assigned to either the Cookie-Cracker stimuli (N= 8) or the Clover-Daisy stimuli (N= 8). One additional child was tested but excluded from analysis and replaced because she did not follow the instruction to place one item in each location (See Results).

#### **Stimuli**

The same stimuli used in Experiment 2 were used in Experiment 3.

#### **Procedure**

Participants were tested in a quiet room at the museum in the presence of their caregiver. The experimenter first introduced the puppet to the child. Children given the Cookie-Cracker stimuli were told that Cookie Monster liked cookies better than crackers; children given the Clover-Daisy stimuli were told that Grover liked clovers better than daisies. The experimenter then said, "Here's a tall box, and here's a short box. It's very hard to climb the tall box, and we don't know if Cookie Monster (or Grover) can do it." She then gave the child two objects (a cookie and a cracker, or a clover and a daisy) and said, "We are going to put one of them on top of the tall box and the other on top of the little box. After that we are going to see what Cookie Monster (or Grover) does and see if he can climb. Where do you want to put them?"

### **3.4.2 Results and Discussion**

All videotapes were coded blind to condition by the first author for inclusion and children's responses to the test question; 81% of the videotapes were recoded on both measures by a second coder blind to hypotheses and conditions. The parents

of one child did not consent to videotape and the child's response was judged online. Inter-coder agreement was 100%.

As predicted, 14 of the 16 children made the informative intervention, putting the object with higher subjective reward in the more costly position (87.5% 95% CI: 75.0–100%). The remaining two children made an uninformative intervention, placing the object with the lower subjective reward in the more costly position. See 3-2. This suggests that children can predict how agents might act in the world as a function of the costs and rewards and can use this information to design interventions that are informative about agents' competence.

Although the task is very simple, it illustrates how combinations of costs and rewards could be (or fail to be) informative about unobservable properties of agents. In this task, children had to combine a high-reward ( $HR$ ) and a low-reward ( $LR$ ) object with a high-cost ( $HC$ ) and a low-cost ( $LC$ ) location to generate a utility function. Although climbing the tall box is always more costly than climbing the short box ( $HC > LC$ ), the exact difference between these costs is unobservable and variable across agents. For agents with high competence, this cost difference ( $HC - LC$ ) is small. However, the less competent an agent is, the higher this cost difference becomes. If children place the high-reward object on the low-cost location, the agent can choose between a high-reward for a low-cost plan ( $HR - LC$ ), and a low-reward for a high-cost plan ( $LR - HC$ ). Here the agent's competence plays no role; it is always better to choose the high-reward for a low-cost plan (because  $HR - LC > LR - HC$  for all values since  $HR > LR$  and  $HC > LC$ ). Thus the choice between these two plans reveals nothing about the agent's competence.

If, instead, children place the high-reward object at the high-cost location, then the agent's rational action choice becomes dependent on his competence. If the agent is very competent, then the high-reward for a high high-cost plan is likely to have a higher utility than the low-reward for a low-cost plan ( $HR - HC > LR - LC$ ). However, if the agent is less competent, then the difference between the high-cost plan and the low-cost plan is relatively large ( $HC - LC$ ) and the low-reward for a low-cost plan becomes more likely to be the highest utility choice ( $HR - HC <$



$LR - LC$ ). Determining the informative intervention requires generating appropriate utility functions that depend on these agent-specific attributes. Again, note that even the informative intervention does not provide evidence about the agent's absolute competence. As in Experiment 2, the agent might be unable to climb the tall box or merely find it very costly to do so. Nonetheless, children were able to distinguish the more and less informative intervention and use information about agent's subjective rewards to provide evidence about the agent's competence.

### 3.5 Experiment 4

In Experiment 3, children manipulated the objective costs associated with each reward to make inferences about the agent's subjective costs. In Experiment 4 we hold the objective cost associated with each reward constant and ask instead whether children can identify agents whose subjective rewards are informative about their subjective costs. Following the same logic described above, if an agent assigns the same reward to two objects, that agent's actions are unlikely to be informative about his or her subjective costs: provided the agent is maximizing utilities, he will always choose the reward associated with lower cost. However, if an agent assigns a higher reward to one object than another, then if the agent fails to pursue the higher reward object, it suggests that the expected cost of the action was high. In Experiment 4, children are shown two treats, one in a high cost location and one in a low cost location. Children are introduced to two puppets with different preferences and asked to identify the puppet who can perform the high cost action. If children can predict how an agent will act as a function of the costs and rewards, they should select the puppet that prefers the treat on the high-cost location. Additionally, because children in Experiment 3 may have simply believed that more desirable objects should be placed in higher places (i.e., because parents often put treats out of children's reach), in Experiment 4 we have each treat be the favorite of one of the puppets.

### **3.5.1 Methods**

#### **Participants**

Sixteen children (mean age: 6.0 years, range 5.0 – 6.9 years) were recruited at an urban children’s museum.

#### **Stimuli**

The same stimuli used in Experiment 3 were used in Experiment 4.

#### **Procedure**

Participants were tested in a quiet room at the museum in the presence of their caregiver. Experiment 4 began identically to the Cookie-Cracker condition in Experiment 2. The experimenter introduced Cookie Monster and Grover, the paper cookies and crackers, and the boxes. This time, Cookie Monster preferred cookies to crackers and Grover preferred crackers to cookies. As in Experiment 2, the experimenter told the child, “Both of our friends can climb up the small box. The big box is really hard to climb. One of our friends can climb it and one of our friends cannot. But we don’t know which one can climb and which one cannot.” The experimenter then placed a cookie on the tall box and a cracker on the short box (object on tall box was counterbalanced). Children were asked, “If we want to figure out which of our friends can climb, which friend should we send in?”

### **3.5.2 Results and Discussion**

All videotapes were coded by the first author blind to condition for inclusion and children’s responses to the test question; 100% of the videotapes were recoded by a second coder on both measures blind to hypotheses and conditions. All parents consented to videotape. Intercoder agreement was 100%.

In Experiment 4, we were interested in which puppet children chose to test. The intervention was considered informative if the child chose the puppet that preferred the treat on the tall box (i.e., cookies for Cookie Monster, crackers for Grover). Twelve

of the 16 children made the informative intervention (75%; 95% CI: 56.25–100%); the remaining four made the uninformative intervention (choosing the other puppet). See Figure 3-2. To succeed in this task, children had to predict how different agents would act as a function of their utilities, given common situational constraints. The agent whose preferred treat was on the short box had an uninformative utility function: he should always climb the short box no matter his competence (because  $HR - LC > LR - HC$ , using the notation of Experiment 3). By contrast, the agent whose preferred treat was on the tall box had an ambiguous utility function that was more likely to be resolved by his choice. If he were competent enough to climb the tall box easily (so that  $HC - LC$  is relatively small, and thus  $HR - HC > LR - LC$ ), he would be expected to climb to get his preferred treat. If he were not so competent (so that  $HC - LC$  is large, and thus  $LR - LC > HR - HC$ ), he would be more likely to choose the less preferred treat on the short box. Additionally, in this setup each treat was preferred by one of the agents. As such, children could not have succeeded through simpler strategies like associating the high-reward treat with a location out of reach, or ignoring the low-reward treat (as the low-reward treat was dependent on the agent). These results suggest that children can assign different sets of costs and rewards to agents under the same situational constraints and predict how the agents would act upon the resulting utilities. Again, we asked about the puppet's ability to climb the tall box, because we were primarily interested in whether children could use the information about the agent's preferences to distinguish the two puppets. We do not know whether children interpreted the high cost action merely as very difficult for the puppet or as so costly that the agent was unable to perform the action at all. Critically however, children recognized that the intervention on the puppet that preferred the high cost treat was potentially informative whereas the intervention on the puppet that preferred the low cost treat was not.

## 3.6 General Discussion

The results of these studies suggest that young children understand how agents act in the world as a function of costs and rewards. Our findings suggest that children understand that there are unobservable, agent-specific aspects of costs and rewards, can make predictions about these unobservable variables, and can design informative interventions to infer them. Experiment 1 showed that children understand that agents act to maximize overall utilities and not just rewards, and as a consequence, agents will sometimes forego a high reward option because the costs of obtaining it are too high. Experiment 2 showed that children understand that competence constraints, unlike situational constraints, are agent-specific and cannot be directly observed; children were able to infer differences in agents' competence using information about their preferences, even given a constant environment in which agents engaged in identical actions. Experiments 3 and 4 showed that, in addition to being able to infer the components of utility functions, children can predict the behavior of agents with different costs and rewards, and thus can design interventions that are informative about agents' competence. Collectively, these studies suggest that children reason about agents' actions and goals in terms of utility functions, consistent with the idea that a naïve utility calculus underlies our social judgments even in early childhood.

In all experiments, children's success rate was consistently above chance (as assessed by 95% CIs). However, in several cases, the 95% confidence intervals were close to 0.5, raising a concern about the robustness of these effects. To test the overall evidence the four experiments provided for our theory, we conducted a Bayesian meta-analysis of effect sizes across all experiments using a hierarchical random-effects model (see Appendix A). The results placed extremely high confidence on an estimate of children's overall rate of theory-correct responding being well above chance, with a mean posterior estimate of 78%, and 95% posterior CIs between 69% and 87%. The variability in children's responses might be due to some but not all of the children having a mature naïve utility calculus, or due to other factors influencing children's responses in our tasks, as in any complex real-world judgment. However,

taken together, the experiments support the hypothesis that children’s judgments of agents’ preferences and costs are consistent with intuitive utility calculations.

As argued in Chapter 2, our proposal of a naïve utility calculus is a natural extension of current accounts of goal-directed action understanding. There is mounting evidence that humans engage in relatively rich psychological reasoning even as infants (e.g., [57, 36, 101, 105]). However, such early social cognition has generally been demonstrated in the context of agent-independent, directly observable, situational and behavioral constraints. Even such sophisticated findings as early false belief understanding are predicated on knowing, for instance, that one person can see the contents of a box and one cannot [101, 140, 142]. In contrast, a mature naïve utility calculus requires understanding that even given identical epistemic access and situational contexts, agents differ in their subjective rewards and costs in ways that may affect their behavior.

Consistent with other research showing the emerging understanding of individual differences [8, 70, 83, 112, 115, 116, 128], and the ability to impute different perspectives on identical evidence [16, 19], the current findings suggest that five and six-year-olds are sensitive to the internal structure of goals and recognize both agent-invariant and agent-dependent aspects of costs and rewards. Children always saw agents who were fully informed about the location of desired objects, they never saw agents change their minds, or attempt and fail to execute an action and agents always reached their goals successfully. Nevertheless, children were able to impute different unobservable motivations and competencies to the agents. Such results suggest that at least by the age of six, the principle of rational action extends to a principle of rational choice. Children not only expect agents to take rational actions towards their goals, but also to use the expected costs and reward of actions to decide when it is worth pursuing a goal at all. Further research might look at the developmental trajectory of these abilities to see if aspects of the naïve utility calculus emerge even earlier in development.

Finally, although our results suggest that children understand that both costs and rewards vary across agents, we do not know if they understand that some aspects of

the costs and rewards are more stable than others. Agent-specific aspects of the costs and rewards can include both state- and trait-like differences. An agent might assign a high-expected cost to an action because of a transient state change (e.g., twisting an ankle) or because of a more stable trait (e.g., being weak or lazy). Similarly, some rewards have high value at some moments but not others (e.g., food when hungry) whereas other rewards may be more stable across time (e.g., long-term values or preferences). Further research might look at the development of children's understanding of both transient and stable aspects of subjective costs and rewards.

Collectively, these studies test some of the fundamental assumptions of a naïve utility calculus, and look at whether children are sensitive to these principles even in early childhood. Children are not only sensitive to information about the subjective costs and rewards of agents' actions, but can also act on the world gain information about these unobservable variables. These abilities emerge relatively early in development, in the absence of formal instruction, suggesting that a naïve utility calculus may be a fundamental component of human social cognition (see Chapter 2).

# Chapter 4

## Action under uncertainty

This chapter is based on Jara-Ettinger, Floyd, Tenenbaum, & Schulz. Children's understanding of action under uncertainty (*in prep*).

### 4.1 Introduction

Consider the simple case of watching Sally get a cookie from the jar. We might infer that Sally likes cookies and that she would get a cookie again if she found herself in the same situation. However, this inference assumes that Sally not only knew there were cookies in the jar, but also knew that she liked cookies. If we knew that Sally had never tried a cookie before, we might not be so fast to assume that Sally will eat cookies in the future (see Figure 2-1 panels f and g).

Cookies, of course, are almost universally familiar and universally liked, but in novel contexts, inferences that depend on an agent's beliefs about her desires are both commonplace and critical for social cognition. Imagine for instance, that you are watching your friend buy food at a market. Usually, her choices reveal what she likes. However, if she's in a foreign country and has never tasted the food before, her choices only reveal her best guess; they may not tell us anything about her stable, long-term preferences. For us to know what someone likes, she has to know it herself first. We can draw comparable inferences in reverse. If you see your friend try a chocolate from her box and then change her mind and choose a different one, you

might infer that she was initially naïve, unsure, or wrong, about what was inside the chocolate. In these cases, the instability of the agent's preferences is indicative of the initial uncertainty of her beliefs about her utilities.

Despite extensive past work on the development of theory of mind, to our knowledge no work has examined children's understanding of how uncertainty about one's own desires influences behavior or at their ability to infer knowledge or ignorance about utilities using information about agents' actions. To the degree that researchers have looked at children's inferences regarding how agents change their mind with evidence, they have focused primarily on issues related to epistemic access: canonically, whether the agent does or does see where a desired object has been placed (e.g., [158]). However, an agent may also be aware or unaware of the value of a putatively desired object. Thus, the degree to which the agent's initial estimates are stable depends on how much the agent knows initially.

From a naïve utility calculus standpoint, these intuitions can be explained by the understanding that agents maximize their *expected* utilities, and not the actual utilities. That is, agents act to select plans that maximize the difference between the rewards they expect to obtain and the costs they expect to incur (see Chapter 2). Figure 4-1 shows a visual schematic of how experience influences choice. The x-axis represents an agent's experience with two fruits (thus, the left end corresponds to ignorant agents and the right end corresponds to knowledgeable agents) and the y-axis shows the agent's expected rewards. Each fruit is surrounded by a circle that represents uncertainty over the reward the agent will obtain. In the low end of experience, agents have greater uncertainty over the possible rewards they may obtain. As such, when ignorant agents select the highest expected reward they are less likely to obtain the actual highest reward, whereas knowledgeable agents are more likely to obtain the actual highest reward. The naïve utility calculus makes a second prediction: ignorant agents should be more likely to revise their choices in light of new experience. This is because a knowledgeable agent's expected rewards do not change much in light of new experience (because they are accurate), whereas an ignorant agent's first experience can drastically change the estimates of the expected rewards



(illustrated by the red and yellow lines in Figure 4-1).

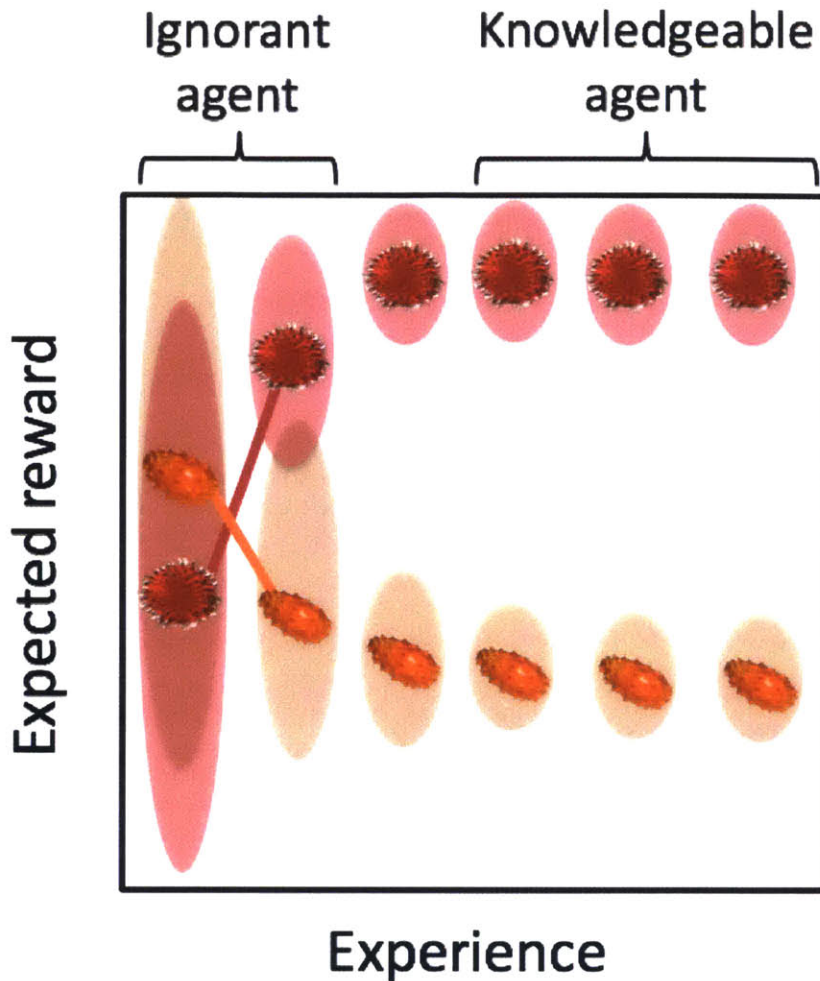


Figure 4-1: Visual representation of expected rewards as a function of experience.

In these studies we investigate children's understanding of how agents' uncertainty about their desires relates to the expected outcome of their goals, and to the stability of their behavior. Figure 4-2 shows a graphical display of the experiments. In Experiment 1 we ask whether children understand that knowledgeable agents are more likely than naïve agents to take actions that lead to a high reward. In Experiment 2, we ask whether children understand that knowledgeable agents are more likely than naïve agents to make decisions that are stable over time. Experiments 3 and 4 examine the inverse questions: In Experiment 3, we ask if children believe that

agents who obtain a high reward are more likely to have been knowledgeable, and in Experiment 4, we ask if children believe that agents who make more stable choices are more likely to have been more knowledgeable. The current study goes beyond merely representing agents' beliefs; instead children must understand that different agents might perform the same actions with the same beliefs about the world and yet interpret the experience differently. Because children's ability to reason explicitly about agents' mistaken beliefs emerges between ages four and five [156] here, we focus on four- and five-year-olds.

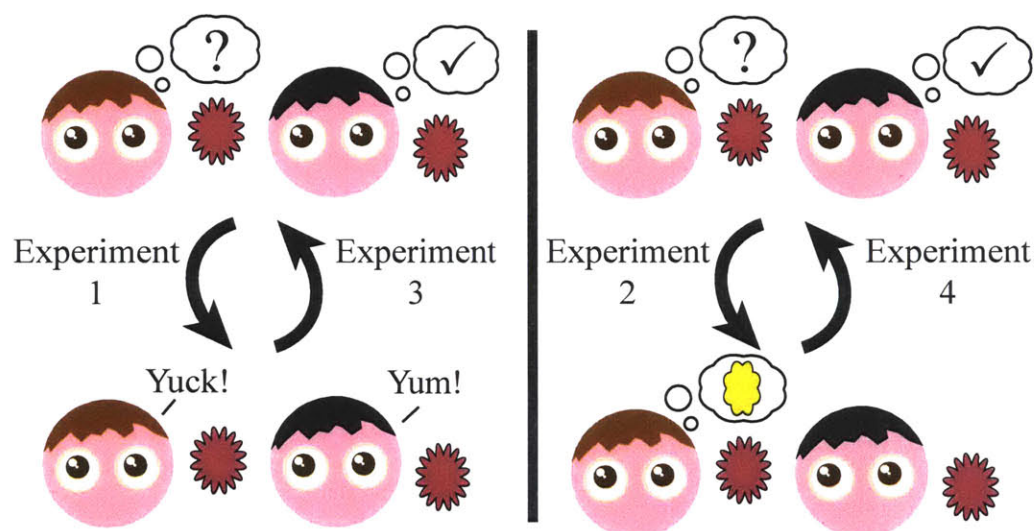


Figure 4-2: In Experiments 1 and 2 children saw a knowledgeable and a naïve agent make the same choice. Children were asked to infer which puppet said “yum” and which puppet said “yuck” (Experiment 1), or which puppet changed their mind (Experiment 2). Experiments 3 and 4 tested inferences in the reverse direction.

## 4.2 Experiment 1

In Experiment 1 we test if children understand that agents' choices are affected by their estimate of the expected rewards of their actions, and thus that knowledgeable agents are more likely than naïve agents to accrue high actual rewards. Children were introduced to two puppets who had been given a choice between two types of fruits. One puppet was knowledgeable and had tasted both fruits before; one puppet was

naïve and had not. Both puppets chose the same fruit. One puppet tasted it and said “Yum!” and one puppet tasted it and said “Yuck!” We looked at whether children inferred that the knowledgeable puppet was more likely to say “Yum!” Methods

### **4.2.1 Participants**

16 participants (mean age (SD): 5.09 years (195 days), range 4.13-5.89 years) were recruited at an urban children’s museum. One additional participant was recruited but not included in the study because he failed to respond the inclusion question correctly (See Procedure).

### **Stimuli**

The stimuli consisted of two pairs of gender-matched puppets, and picture cutouts of two fruits: Rambutans, and African cucumbers.

### **Procedure**

Participants were tested individually in a quiet room in a children’s museum. The child and the experimenter sat on opposite sides of a small table. The experimenter first introduced the cutout pictures of the rambutans and the African cucumbers and placed four pictures of each fruit on the table with each kind of fruit in its own pile. Next, the experimenter introduced the two puppets by name (“Anne” and “Sally”, or “Arnold” and “Bob”, depending on the participant’s gender. The puppets were matched with the participant’s gender to ensure gender biases did not influence our task). The experimenter then explained that “Sally has never seen these fruits before and she doesn’t know what they taste like” while “Anne knows all about these fruits. She knows what they taste like.” (Knowledgeable puppet and introduction order were counterbalanced). Next, the experimenter told the participant “Earlier today, we told our friends they had to pick one fruit to each, and both of our friends picked a rambutan.” (Actual fruit counterbalanced). Next, the experimenter placed a picture of a rambutan in front of each puppet and explained, “Both of our friends took a bite

of their fruit and one of them said 'Yum!' and one of them said 'Yuck!'" Participants were then asked an inclusion question to ensure the child remembered the critical information: "Can you tell me, which of our friends has not tasted the fruits before? And which one of our friends has not tasted these fruits before?" Finally, participants were asked which puppet said "Yum!" and which puppet said "Yuck!"

#### **4.2.2 Results and Discussion**

Children who failed to respond correctly to the inclusion question were excluded from analysis and replaced ( $n = 1$ ). Results were coded for adherence to the script by a coder blind to the child's response to the test question (no participants were dropped due to experimenter error). Videotapes were then coded to record the child's response to the test question. Children were coded as answering correctly if they indicated that the knowledgeable puppet had said "Yum." Of the sixteen children who responded to the inclusion question correctly, 100% responded correctly to the test question (95% CI: 82.93%-100%. See 4-3).

Note that if children believe that an agent's choices always reflect her preferences then children should have expected both puppets to say "Yum!" That is, if children recovered agents' desires only from information about the agents' actions and beliefs about the state of the world, then children should have responded at chance in this context. Both agents knew they had a choice of the two fruits and both agents made the same choice. Children instead recognized that the knowledgeable agent would be more likely to like the chosen fruit, which suggests that children understand that agents choose the options with the highest expected rewards and that a naïve agent's choices may be governed by an inaccurate estimate of the actual reward.

### **4.3 Experiment 2**

The results from Experiment 1 suggest that four and five-year-olds understand that, relative to knowledgeable agents, naïve agents are more likely to make unrewarding choices. When this happens, naïve agents will most likely reconsider their choices.

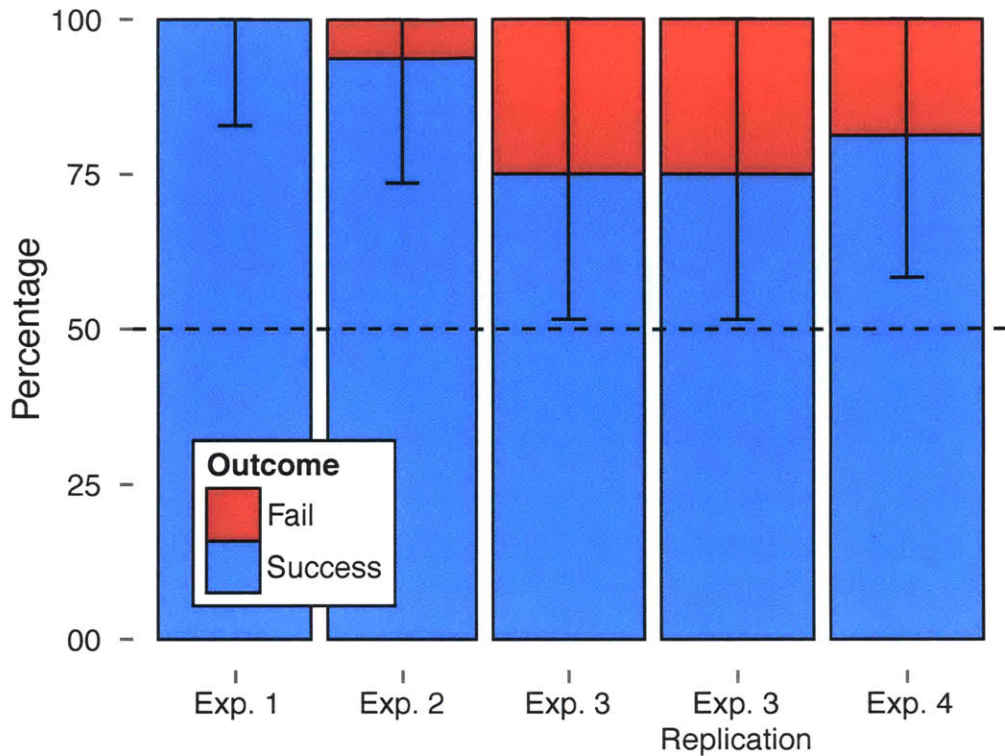


Figure 4-3: Results from all experiments. The x-axis shows each experiment and the y-axis shows the distribution of children’s choices (color coded). The dashed horizontal line represents expected chance performance and the vertical solid lines show 95% confidence intervals.

However, naïve agents may reconsider their choice even when they obtain a positive reward, especially if the reward is lower than they expected. Thus, in general, naïve agents are more likely to have unstable choices, when compared to knowledgeable agents. In Experiment 2 we test if, in the absence of knowledge about the rewards, children have expectations about the stability of the choices of knowledgeable and naïve agents.

### 4.3.1 Methods

#### Participants

16 participants (mean age (SD): 5.16 years (241 days), range 4.01-5.96 years) were recruited at an urban children’s museum. One additional participant was recruited

but excluded from the study and replaced because he declined to complete the task.

## **Stimuli**

The stimuli were identical to those used in Experiment 1.

## **Procedure**

The procedure was identical to the procedure in Experiment 1 except as follows: Neither puppet said “Yum!” or “Yuck!” Instead, after the puppets had tasted the fruits, the experimenter said, “Both of our friends took a bite from their fruit and one of them changed her mind and decided she wanted to eat a different fruit”. For the test question, participants were asked, “Which one of our friends changed his/her mind?”

### **4.3.2 Results and Discussion**

As in Experiment 1, children who failed to respond to the inclusion question were excluded from the study. (No participants were dropped on these grounds.) Videos were first coded for adherence to script by a coder blind to the child’s test response (no participants were dropped due to experimenter error), and later to record the child’s answer to the test question. Children were coded as responding correctly if they indicated that the naïve agent was the one who had changed her mind. Fifteen of the 16 children responded correctly to the test question (93.75%; 95% CI: 73.60%-100%. See 4-3).

Note that in contrast to Experiment 1, participants in Experiment 2 never obtained any information about the outcome of each puppet’s choice. Thus, it was possible that either or both puppets had liked or disliked their chosen fruit. Nevertheless, children were able to infer that naïve agents are more likely to make unstable choices. Furthermore, in this experiment we used a neutral dependent measure, thus ensuring that participants couldn’t succeed by grouping together two features with a positive valence (such as knowledge and “yumminess” in Experiment 1). Together



with Experiment 1, these results suggest that four and five year-olds understand that relative to knowledgeable agents, naïve agents are more likely to make choices that lead to low rewards, and thus that their choices are less likely to be stable over time.

Experiments 1 and 2 suggest that children have different expectations about the rewards that knowledgeable and naïve agents obtain and about the stability of their actions. In Experiments 3 and 4 we ask if children can reverse these inferences, and infer an agent’s prior knowledge based on whether the agent succeeds in obtaining high rewards (Experiment 3), and whether the agent’s choices are stable (Experiment 4).

## 4.4 Experiment 3

In Experiment 3 we invert the question asked in Experiment 1. Here we ask if children believe that agents who make choices that result in low rewards are more likely to have been naïve prior to making their choice. Children watched two puppets pick the same fruit to eat. After learning that one puppet said “Yum!” and the other puppet said “Yuck!” children were asked to decide which puppet had not tasted the fruits before.

### 4.4.1 Methods

#### Participants

32 participants (mean age (SD): 5.12 years (194 days), range 4.12-5.98 years) were recruited at an urban children’s museum. Sixteen participants were recruited for the original experiment, and sixteen additional participants were recruited to conduct a replication (see Results). Four additional participants were recruited in the original experiment but excluded from analysis and replaced for failing the inclusion question ( $n = 1$ ), declining to complete the experiment ( $n = 1$ ), and declining to answer the test question ( $n = 2$ ). One additional participant was recruited in the replication experiment but excluded from analysis because he declined to answer the test question.

See Results.

## **Stimuli**

The stimuli were identical to those used in Experiment 1.

## **Procedure**

Participants were tested individually in a quiet room. As in Experiment 1, the experimenter introduced the two puppets and the fruits and explained that each puppet chose a fruit to eat. The section from Experiment 1 in which where the experimenter explained that one puppet was more knowledgeable than the other was omitted. After both puppets chose a fruit, the experimenter said “Anne and Sally (or Arnold and Bob) both took a bite from their rambutans (or African cucumbers). Anne said ‘Yum!’ Sally said ‘Yuck!’” Next, the experimenter said, “But guess what? One of our friends didn’t know what rambutans tasted like until today.” Children were then asked to remember which puppet had said “Yum!” and which puppet had said “Yuck!” For the test question, the experimenter asked, “Can you tell me, which one of our friends didn’t know what rambutans tasted like until today?” (Actual fruits counterbalanced throughout.) The replication experiment had the same procedure as the original experiment with the exception that the inclusion question was asked immediately after the puppets tasted the fruit, thus making the last part of the experiment more fluent.

### **4.4.2 Results and Discussion**

Results were coded as in Experiments 1 and 2. Participants were coded as responding correctly if they indicated that the puppet who said “Yuck!” was the one who had not tasted the fruits before today. Of the 16 participants who made a choice in the original experiment, 75.00% ( $n = 12$ ) responded correctly (95% CI: 51.56%-100%). These results suggest that children can use knowledge about the actual subjective rewards that different agents obtain to infer which agent is more likely to have been naïve in



her estimate of the expected rewards. However, four children answered incorrectly and two were excluded from analysis and replaced for failing to answer the test question. Thus, to ensure the validity of our interpretation we replicated the experiment. Out of the 16 participants who made a choice in the replication, 75% (n=12) responded correctly (95% CI: 51.56%-100%). Together, these experiments suggests that children can in fact use knowledge about subjective rewards to infer knowledge, and it provides some suggestive evidence that children may find it easier to use information about agent's knowledge to predict their subjective rewards than to use information about agent's rewards to recover information about their unobservable mental states.

## 4.5 Experiment 4

In Experiment 4 we invert the question asked in Experiment 2. Here we see if children can infer which of two agents is more likely to be naïve when one shows stable preferences and one does not.

### 4.5.1 Methods

#### Participants

16 participants (mean age (SD): 5.64 years (244 days), range 4.04-5.93 years) were recruited at an urban children's museum. Four additional children were tested but excluded from the study because they failed to respond to the inclusion question correctly. See Results.

#### Stimuli

The stimuli were identical to those used in Experiment 1.

#### Procedure

The procedure was identical to Experiment 3 except as follows: First, children were never given any information about whether the puppets said "Yum!" or "Yuck!" after

tasting the fruits. Instead, after taking a bite from their fruit, the experimenter said, “Anne kept eating the rambutan. Sally changed her mind and said she wanted an African cucumber instead.” As in Experiment 3, at test children were asked, “Can you tell me, which one of our friends didn’t know what rambutans tasted like until today?” (Actual fruits counterbalanced throughout.)

### 4.5.2 Results and Discussion

All results were coded in the same way as Experiments 1-3. Four children failed to respond to the inclusion question correctly and were therefore excluded from analysis and replaced. Children’s responses were coded as responding correctly if they indicated that the puppet who changed her mind was the one who had never tasted the fruits before. Of the 16 participants who made a choice, 81.25% ( $n = 13$ ) responded correctly (95% CI: 58.34-100%).

Together with Experiment 2, these results suggest children understand that relative to naïve agents, knowledgeable agents are more likely to stick to their choices. Moreover, children can infer which agents are more likely to be knowledgeable based on the stability of these choices.

## 4.6 General Discussion

Across four experiments, we studied children’s understanding of the relationship between agents’ knowledge of their desires and the outcome of these agents’ actions. Our results suggest that four and five-year-olds understand that, relative to naïve agents, knowledgeable agents are more likely to obtain high rewards (Experiment 1) and more likely to make stable choices (Experiment 2). Similarly, children believe that agents who obtain high rewards and agents who make stable choices are more likely to be knowledgeable (Experiments 3 and 4 respectively). Collectively, these results suggest that children understand that an agent’s choices are not always aligned with the highest utility, but rather with the highest expected utility. As such, agents who have more uncertainty about the value of a target are more likely to obtain a

low reward, and more likely to explore different alternatives.

Children’s responses in our task could have been driven by their expectations about knowledgeable agents, naïve agents, or both. Future research might see which of these expectations underlies children’s reasoning. Additionally, we have emphasized the possibility that children should infer that naïve agents are more likely than knowledgeable agents to make unrewarding choices, and this interpretation is consistent with the results of Experiments 1 and 3. However, children may also believe that, relative to knowledgeable agents, naïve agents are more likely to find exploration rewarding; this kind of reasoning could contribute to the results of Experiments 2 and 4

Computationally, children’s intuitions may stem from a categorical distinction between knowledgeable and naïve agents, or from a continuous representation of how the amount of knowledge an agent has influences the quality of their choices. Although these two accounts make similar predictions in our task, the latter theory is more powerful, as it enables observers to reason about intermediate stages of knowledge. Further work is needed to establish exactly how children represent an agent’s uncertainty. Furthermore, it is an open question how children represent a goal’s reward. For instance, children might assume that each goal has a fixed reward value, which the agent may not know. Alternatively, children may understand that the same outcome can have variable rewards over time (an agent may find apples very rewarding when she’s hungry and less so when she’s full). Future research might investigate the precise representations underlying children’s calculations of agents’ utilities.

In these studies, we focused on agents’ estimates of the expected reward of their actions. However, the same logic applies to agents’ understanding of the cost of actions. For example, Sally might be eager (or reluctant) to run a marathon. However, if you know she does not understand the costs involved, you might not be confident that her current actions will be informative about her future ones. The converse inferences also hold. If you see Sally sign up for a committee and then fail to attend, you might infer that she had not accurately estimated the commitment involved. Future work might investigate whether four and five-year-olds also have strong intuitions

about how agents' knowledge about the costs of action influences their preferences and choices.

The current results suggest the importance of integrating the naïve utility calculus into theory of mind models (see Chapter 6). Past research has focused on learners' ability to draw inferences connecting agents' beliefs, desires, and actions (e.g., [158, 4]). Almost universally however, beliefs and desires have been treated as independent, non-interacting variables, and the content of beliefs has been restricted to information about the world. The current study suggests that children understand that agents have beliefs not only about the world, but also about their own preferences. Children understand that as agents gain knowledge about the world, their preferences can change as well.

# Chapter 5

## Social reasoning

This chapter is based on Jara-Ettinger, Gweon, Tenenbaum, & Schulz. Not so innocent: Toddlers' reasoning about costs, competence, and culpability (2015). *Psychological Science*.

### 5.1 Introduction

So far, we have only discussed the naïve utility calculus in the context of goals directed towards objects. However, if the naïve utility calculus is the basis of all social reasoning, it should also explain how we reason about social goals. How might information about costs and rewards affect social judgment? Imagine that your neighbor, Sally, watches someone struggle to reach a package on a high shelf. Sally stands by and does nothing at all. Although there is no intrinsic relationship between height and moral worth, you may well judge Sally less harshly if she is 4'11" than if she is an NCAA basketball player.

What analysis underlies this judgment? We suggest that in predicting and evaluating behavior, we assume that the costs and rewards of an action jointly affect the likelihood that an agent will act. If we know the costs (e.g., in time and energy) that an agent is willing to incur to achieve a goal, we can make inferences about the agent's subjective rewards (i.e., her level of motivation). Conversely, if we know how motivated an agent is, we can infer the costs she might be willing to incur. These

attributions trade off: If a highly motivated agent fails to act, we may infer that achieving the goal was too costly; conversely, if an agent doesn't pursue a low-cost goal, we may infer that she does not value it highly.

Inferences about agents' motivations are particularly influential in moral judgment [26, 76, 166]. Moral bystanders may be exonerated if a helpful action would have cost them dearly; they are likely to be judged harshly if they merely found helping insufficiently rewarding. Ambiguity arises when agents fail to perform costly actions or do perform low cost ones. If Sally is 4'11", we can infer that the cost of reaching the shelf is high. This renders the motivation behind her failure to act ambiguous: Did she not want to help or was helping too hard? By contrast, if Sally is an NCAA basketball player, she may not get much credit even if she does help get the package off the shelf. The costs are so trivial for her that her motivation to be helpful did not need to be high.

These kinds of considerations fall naturally from the naïve utility calculus. However, to date no empirical work has looked at how differences in the cost of actions across agents affect children's judgments. Here we test the prediction that young children can estimate the costs associated with agents' actions and that this analysis affects their social evaluations.

In Experiment 1, we test the basic premise that children are sensitive to the perceived cost of actions. We predict that at baseline children will prefer more (versus less) competent agents. In Experiments 2 and 3, we look at whether two-year-olds can use differences in agents' costs to infer differences in their motivation. We predict that if two agents refuse to help, children should think the less competent agent is nicer.

## 5.2 Experiment 1

In Experiment 1 we look at whether toddlers distinguish agents who incur different costs (in time and effort) to achieve a goal and whether toddlers prefer agents who incur lower costs.

### 5.2.1 Methods

Informed by developmental studies on comparable topics (e.g., [56, 75, 102]) we predicted strong behavioral effects throughout. We use a sample size of  $n = 16$  per condition in all experiments (replacing participants excluded by the decision of coders blind to condition; see Results). Low-level design features are counterbalanced throughout.

#### Participants

Twenty-four toddlers (mean age (SD): 21.58 months (97 days), range 17.1-28.5 months) were tested at an urban children's museum. Twelve additional toddlers were recruited but not included in the study because they declined to participate in a warm-up task, in which the child was asked to choose between two stuffed animals. Five children were excluded from analysis: three due to experimental error and two due to parental interference. (See Results.)

#### Stimuli

Participants were shown two puppets and a yellow cylindrical toy with a black button at the top. The toy played music when the button was pressed.

#### Procedure

Participants were tested in a quiet room at the museum. The child's parent was seated on a chair facing away from the testing table and the parent was asked to hold the toddler over his or her shoulder. Thus the child could see the stimuli but the parent could not. Once the parent and toddler were positioned, the experimenter presented the toy to the child and introduced the two puppets (puppet position counterbalanced between participants). See 5-1. He said, "Here are my two friends! They are going to show you how the toy works." Both puppets were continuously present throughout the experiment and each puppet approached the toy (order counterbalanced between participants) one at a time. During their turn, each puppet said, "It's my turn!" and then pressed the button. When the toy activated, the toy played a song for

approximately 7 seconds, the puppet moved rhythmically to the sound, and then the puppet released the button. After releasing the button, the puppet who activated the toy said "Yay!" to celebrate the success.

The puppets differed in how many attempts it took them to activate the toy. The more competent agent always made the toy play music on the first attempt. The less competent puppet tried several times to activate the toy (flattening his hand over the button but not depressing it fully). After the third or fourth failed attempt, the less competent puppet backed away to "look" at the button, and then tried again. The incompetent puppet made a few more failed attempts and then successfully activated the toy. (The number of total attempts ranged from 5 - 8 trials across participants, allowing some flexibility in maintaining the child's attention to the task.) After the show, the parent was asked to turn around and to place their child at a marker on the middle of the edge of a lower table. The experimenter placed both puppets on opposite sides of the table equidistant from the child and asked the child which one she wanted to play with.

### **5.2.2 Results and Discussion**

All videotapes were coded by a coder blind to condition. Three children were excluded from analysis due to the coders' judgment that the puppets were not placed equidistant from the child. Two additional children were excluded from analysis due to parental interference. The coder recorded the toddlers' first contact with a puppet following the prompt. If the child did not make a choice within a 30-second window following the prompt, the experiment was ended. Three children did not make a choice. 93.75% (CI: 87.5%-100%) of subjects who made a choice preferred the competent agent (15 out of 16 subjects). See 5-2.

In Experiment 1, toddlers strongly preferred the agent who achieved his goal more easily. Future research might establish whether toddlers' preference is driven by the agents' overall effort to achieve the goal, the time taken to achieve the goal, the greater ease of interpreting the more fluid actions, or a more abstract judgment about these factors as indices of competence per se. However, the results of Experiment 1 suggest



that toddlers distinguish agents based on diverse cues associated with the cost of goal-directed actions and prefer agents who incur fewer costs.

## 5.3 Experiment 2

In Experiment 2, we look at how the cost of agents' actions affects toddlers' social evaluations. Because pilot work suggested that the task in Experiment 2 was more demanding than in Experiment 1 we tested slightly older children: two-year-olds.

### 5.3.1 Methods

#### Participants

Seventeen two-year-olds (mean age (SD): 2.64 years (84 days), range 2.26-2.96 years) were recruited and tested at an urban children's museum. Five additional two-year-olds were recruited but not included in the study because they declined to participate in a warm-up task, in which the child was asked to choose between two stuffed animals. One additional subject was excluded from analysis due to refusal to make a choice.

#### Stimuli

The stimuli used in Experiment 2 were identical to stimuli used in Experiment 1.

#### Procedure

Because participants in Experiment 2 were older, they were given a choice between sitting in a chair or standing in front of the testing table, behind the parent's chair. Parents were asked to turn their back to the table so both puppets and the toy were out of sight. Parents were given a script to read explaining the experimental procedure. (The experimenter also explained the script before parents entered the testing room to ensure they were willing to participate.) The beginning of the experiment was otherwise identical to Experiment 1. (See 5-1) As in Experiment 1, the experimenter introduced two puppets: one was able to activate the toy immediately; the

other required 6-8 attempts. Then the child was told, “Now your mom/dad is going to turn around, pick up the yellow toy and ask our friends a question.” As instructed in the script, parents then turned around. They saw a single puppet and the toy; the other puppet was out of sight. As per the script, parents looked at the puppet and asked, “Can you help me?” After the parent asked for help, the puppet looked at the toy and then at the parent. The puppet said “No!”, turned around, and hid under the table. This was repeated with the second puppet (order counterbalanced). The question and answer procedure was then repeated with each puppet a second time.

After each puppet refused to help a second time, the experimenter took the toy from the parent and asked the child to stand on a marker in the center of the table’s edge. As in Experiment 1, the experimenter placed each puppet on the table equidistant from the child (left/right counterbalanced) and asked which one they would rather play with.

### 5.3.2 Results and Discussion

Results were coded from videotape by a coder blind to conditions, as in Experiment 1. Children were excluded from analysis if, in the coder’s judgment, the puppets were not placed equidistant from the children (no children were excluded on these grounds) or if children did not make a choice within the 30-second window (one child), resulting in a total of 16 children. (See Participants.)

Contrary to our prediction, two-year-olds did not prefer the incompetent agent. Instead, they showed a bias towards the competent agent. Of the 16 toddlers who made a choice, 68.75% (CI: 50-93.75%) chose the more competent agent (11 participants). (See Figure 5-2.)

Why didn’t children reject the competent agent, given his refusal to perform a low cost helpful action? One possibility is that toddlers distinguish agents based on the costs they incur, but fail to infer either that low costs entail an obligation to help, or that high costs exonerate an agent from helping. A related possibility is that toddlers make categorical distinctions between classes of behavior (e.g., “helping”, “not helping”, and “hindering”) but no distinctions within categories. That is, toddlers

might distinguish the agents based on the costs they incur but find them equally blameworthy (because neither helped).

A final possibility is that two-year-olds infer that the puppet who incurred higher costs is less culpable but prefer to affiliate with agents who incur low costs. Indeed, an informal survey (see Appendix B) suggested that adults are split on analogous questions of this kind. When two agents refused to help, 44% (CI: 24%-64%) of participants preferred the less competent agent (“[the competent agent] sounds like a jerk, why didn’t he help?”) and 56% preferred the more competent agent (“...I’d rather have smart friends than not so smart.”). To distinguish these possibilities, and to see if children’s preferences are robust, in Experiment 3 we compare children’s choice of “which puppet they want to play with” with their choice of “which puppet is nicer.”

## 5.4 Experiment 3

### 5.4.1 Methods

#### Participants

66 participants (mean age (SD): 2.48 years (114 days), range 2.00-2.98 years) were recruited and tested at an urban children’s museum. One participant was dropped due to experimenter error, two participants were dropped because they left before the experiment concluded, and five participants were dropped due to parental interference. Thirteen additional two-year-olds were recruited but not included in the study because they declined to participate in a warm-up task, in which the child was asked to choose between two stuffed animals. Participants were randomly assigned to a “Play” condition, a “Nicer” condition, or a “Nicer Baseline” control condition. Ten children (two in the Play condition, five in the Nicer condition, and three in the Nicer Baseline control) were excluded from analysis because they failed to respond in the first 30 seconds, resulting in a final sample of 16 participants per condition.

## Stimuli

The stimuli used in Experiment 3 were identical to stimuli used in Experiment 2.

## Procedure

The procedure for the two test conditions (Play and Nicer) was identical to the procedure of Experiment 2 with two modifications. First, at the end of the experiment, participants in the Play condition were asked, “Which one would you rather play with?” (replicating Experiment 2); children in the Nicer condition were asked, “Which one is nicer?” Second, we clarified the interpretation of “Can you help me?” to ensure the interpretation “Would you help me?” rather than “Are you able to help me?” In Experiment 3, after the first time both puppets refused to help the child, the experimenter told the child “Neither of our friends wants to help your mom/dad with the toy! Let’s ask them one more time to make sure.” After the second refusal, the experimenter reiterated, “Our friends do not want to help.”

We predicted that toddlers would judge the less competent agent as “nicer”. Alternatively however, children might simply believe that less competent agents are nicer a priori. The Nicer Baseline control condition allowed us to assess children’s initial beliefs about the relative niceness of the agents. The procedure for the Nicer Baseline control condition was identical to the procedure of Experiment 1 with the exception that the experimenter introduced the puppets saying, “Here are my two friends! They are both going to play with the toy.” At the end of the experiment, participants were asked, “Which one is nicer?”

### 5.4.2 Results and Discussion

Results were coded from videotape by a coder blind to condition. Children were excluded from analysis if, in the coder’s judgment, the puppets were not placed equidistant from the child (no children were dropped on these grounds) or if children did not make a choice within the 30-second window (7 participants), resulting in 16 2-year-olds per condition. (See Participants.) 5-2 shows the results from the

experiment. In the Play condition, 81.25% of children chose the competent agent (13 participants; CI: 68.75-100%), replicating the findings of Experiment 2. By contrast, children in the Nicer condition chose the less competent agent. Of the 16 two-year-olds who made a choice, only 31.25% (CI: 6.25%-50%) chose the competent agent (5 participants). That is, the participants' choice of the competent agent dropped from 81.25% to 31.25% when asked to judge which agent was nicer (CI: 21.05%-81.67%;  $p < 0.05$  Fisher exact test). Finally, participants' performance in the Nicer Baseline control condition suggests that the switch between the Play and Nicer conditions was not due to children's baseline belief that less competent agents are nicer. While only 68.75% of participants (11 out of 16) chose the incompetent agent in the Nicer condition, 31.25% of participants (5 out of 16) chose the incompetent agent in the Nicer Baseline control, a decrease of 37.5% (95% CI: 6.17-71.03%;  $p = 0.075$  Fisher exact test). This suggests that toddlers do not simply assume incompetent agents are sympathetic; rather they take into account the relative cost of agent's actions.

When both puppets refuse to help, toddlers might infer that the less competent agent is nicer because they exonerate the less competent agent from performing a costly action, because they judge the more competent agent harshly for refusing to perform a low cost action, or both. Further research might shed light on the precise inferences that underlie children's social evaluations. However, note that neither agent is canonically nice: both agents explicitly refused to engage in a helpful action. If children only understood "nice" with respect to nice, helpful behaviors, rather than with respect to internal motivations, then they should have chosen at chance or refused to answer. Instead, two-year-olds were able to use differences in the agent's costs to identify the nicer of two unhelpful agents.

## 5.5 General Discussion

Consistent with the idea that a naïve utility calculus is integral to children's understanding of agents, we found that the cost of agents' actions affects children's social evaluations. Toddlers are sensitive to cues associated with the relative competence of

agents and prefer agents who achieve goals quickly and easily. However, they evaluate agents differently in moral contexts; agents who fail to perform helpful actions at relatively high costs are judged to be nicer than those who fail to act at lower cost.

These findings are consistent with previous work suggesting that toddlers differentiate between agents who are unable, rather than unwilling, to act prosocially. Children make more attempts to reach for a toy if an experimenter tries and fails to transfer it than if she is clearly teasing and unwilling to share the toy [7]. Similarly, toddlers prefer to give a new toy to an experimenter who tries unsuccessfully to help than to one who teases and is unwilling to help [31]. These results suggest that toddlers distinguish agents' motivations and selectively reward helpful agents.

Our work extends these studies in two important ways. First, in the previous studies there were unambiguous cues to the agents' motivations (e.g., sincere attempts versus taunting). Even capuchin monkeys are sensitive to overt cues distinguishing unwilling and unable agents [106]. By contrast, in our study, the behavioral cues were informative only about the relative costs incurred by both agents; there were no direct cues to the agents' motivations. Additionally, note that in our study both agents were able and both were unwilling to perform the helpful action; only the cost difference supported the possibility that one agent was unmotivated to help whereas the other was merely unwilling to incur high costs.

Interestingly, the naïve utility calculus may support adults' intuition that incompetent agents are more sympathetic than competent agents. (Compare your feelings of empathy when an elderly man and a macho 24-year-old both walking into a wall; [20]). It is relatively easy to recognize when competent agents are unmotivated to be helpful (because their competence is not in question) but difficult to recognize when they are highly motivated (because easy tasks require little motivation). By contrast it is relatively easy to recognize when incompetent agents are highly motivated to help (they must be if it is difficult for them) but difficult to recognize when they are unmotivated. Therefore, given potentially ambiguous evidence, we are more likely to infer that an agent is not nice if we know she is competent and more likely to infer she is nice if we know she is incompetent. Although at baseline two-year-olds believed

the agent who incurred low costs was nicer, future work might test the speculation that over time, inferences supported by the naïve utility calculus lead to the adult intuition that less competent agents are more likely to be nice.

Our finding that children generally prefer competent agents is congruent with work in the domain of epistemic trust suggesting that children prefer reliable agents [10, 9, 64, 77, 79, 99, 103, 113, 119, 50, 136, 137, 32]. However, (with one exception; Pasquini, et al., [103]) such studies have pitted good informants against bad ones rather than looking at relative competence. Pasquini et. al. [103] did look at children's ability to make graded judgments, and found that children younger than four fail to track agents' reliability across independent trials. Here however, we did not vary the probability of success across trials; we varied the amount of effort agents' needed to succeed. Thus toddlers only had to encode the agent-specific effort associated with achieving the goal. However, because we provided redundant cues to the costs of the agents' actions, including the time and the number of times each agent pressed the button, we do not know to what extent toddlers' preferences were driven by each individual cue, or if their choice was guided by a more abstract representation of competence. Future research can shed light on the full range of cues we use to infer agents' competence.

The current study suggests that children are sensitive to the cost of actions early in development. At an age when children themselves are still largely incompetent and exonerated from moral responsibility, their ability to understand cues to how one characteristic might bear upon the other suggests remarkably sophisticated inferential abilities and highlights the importance of building a new theoretical synthesis for understanding the development social cognition.

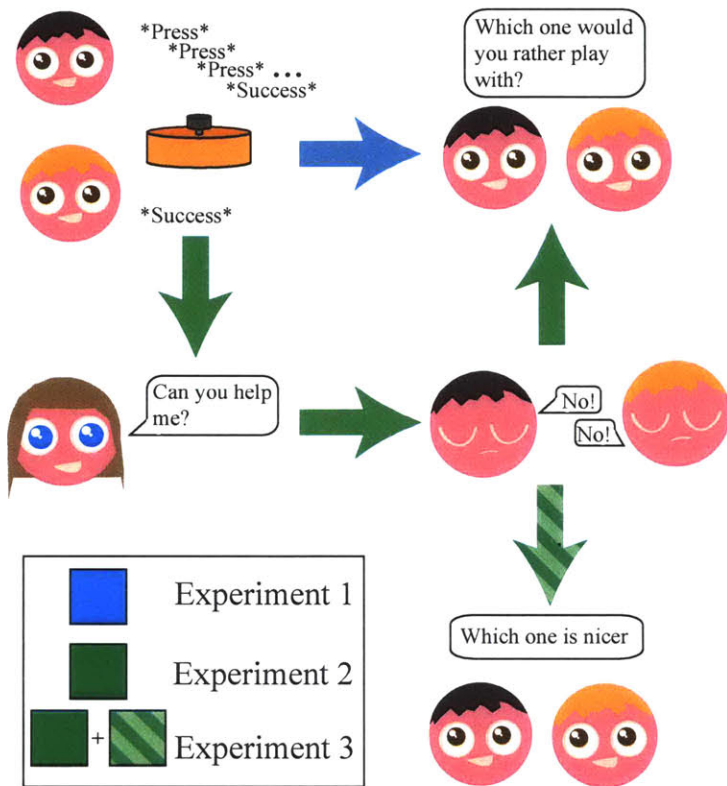


Figure 5-1: Procedure for all experiments. Children were introduced to two puppets and a toy. One puppet (the Competent agent) was able to make the toy play music on the first attempt; the other puppet (the Incompetent agent) succeeded, but only after many attempts. In Experiment 1 (blue arrow), children were asked to choose a puppet to play with. In Experiments 2 and 3 (green arrows), after the child saw both puppets activate the toy, the parent turned around and asked each puppet for help with the toy. Both puppets refused. In Experiment 2 (solid green arrow) children were then asked to choose one of the puppets to play with (as in Experiment 1). In Experiment 3 one condition replicated Experiment 2 (solid green arrow); in a second condition (striped green arrow) children were asked which puppet was nicer; in a third baseline condition (not shown) children were asked which puppet was nicer in the absence of any refusal-to-help information.



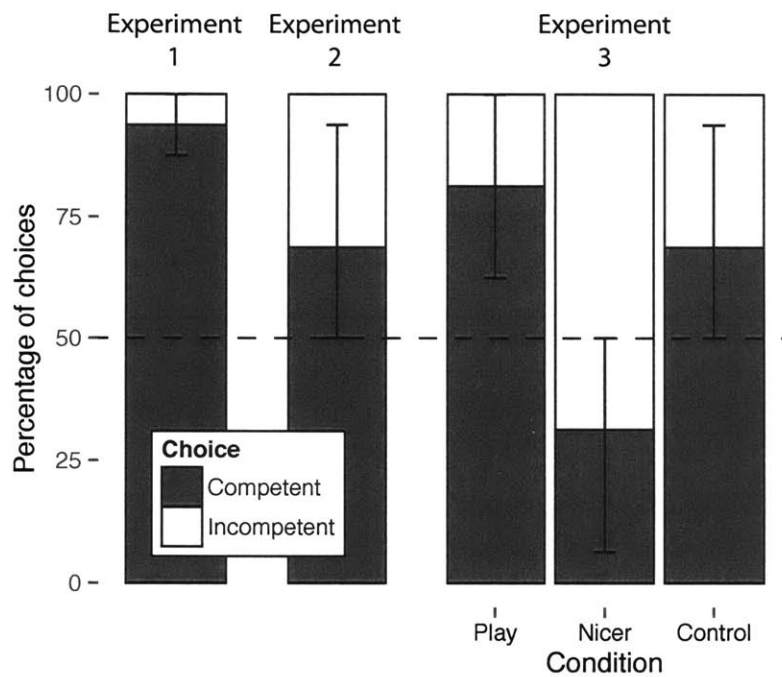


Figure 5-2: Results from all Experiments. In all experiments, children were introduced to a more and a less competent agent. In Experiment 1, children were asked which puppet they preferred to play with; they preferred the more competent agent. In Experiment 2, both agents refused to help activate a toy; toddlers again preferred the more competent agent. In Experiment 3, the puppets refused to help and children were asked either which puppet they wanted to play with (replicating Experiment 2; Play condition) or which puppet was nicer (Nicer condition). Children preferred to play with the more competent agent but said the less competent agent was nicer. In a Nicer Baseline Control condition, children were introduced to the agents as in Experiment 1 (with no refusal to help), and asked which was nicer. At baseline, children chose the more competent agent.



# Chapter 6

## Formal implementation

This chapter is based on Jara-Ettinger, Schulz, & Tenenbaum. The Naïve Utility Calculus: joint inferences on the costs and rewards of action (*in prep*).

### 6.1 Introduction

So far, Chapters 2-5 have presented qualitative evidence for the naïve utility calculus. In this chapter we present a formal model that we use to evaluate our proposal with more precision. We begin by briefly describing the bayesian approach to modeling mental state inference and highlighting the limitations that the naïve utility calculus aims to help solve.

Current models have focused on the understanding that agents have goals [160] that they complete as efficiently as possible [127, 36]. This assumption, known as the principle of efficiency, enables humans to infer unobservable goals from observable behavior. The logic of this inference can be described and formalized using Bayesian inference, where the probability that an agent has goal  $G$  given that they took actions  $A$  is given by

$$p(G|A) \propto L(A|G)p(G)$$

Here,  $L(A|G)$  is the likelihood that the agent would take actions  $A$  if she had

goal  $G$ , and  $p(G)$  is the prior belief that the agent goal  $G$ . The principle of efficiency determines the likelihood function: The more efficiently the actions  $A$  complete the goal  $G$ , the higher their likelihood (And therefore the higher the posterior probability that the agent has that goal).

This kind of inference, called inverse planning, was formally modeled by Baker et al [3], using Markov Decision Processes (MDPs). In the MDP framework, the environment is modeled as a set of states, each with an associated utility (that can be positive or negative), which the agent can navigate by taking different actions (e.g., walk left, right, etc). With this formulation, it is possible to determine the sequence of actions that maximize an agent's utility as efficiently as possible. Using MDPs as a model for how agents act, goal inference can be formalized as inferring the unobservable utility function that is guiding the agent's actions. These models predict with high quantitative accuracy how adults infer goals in simple scenarios [3, 4, 62].

### **6.1.1 Social reasoning beyond goal attribution**

Despite the success of these models, the power of these inferences is limited.

#### **Explanatory limitations.**

To illustrate why, consider a simple example. A man is walking and reaches a fork on the road. The left path leads to a lake where he can swim, and the right path leads to his house. The man stops for a second and then takes the right path. The man's goal is immediately revealed after his first step, as he's taking an efficient path towards his house and an inefficient path towards the lake. However, this inference only tells us what the man is trying to achieve, but not why. The man may be going home because he doesn't like swimming, because he cannot swim, because he's too tired, or too hungry.

Models that infer the utility function will treat all the above explanations as being roughly equivalent, as they all reduce to a utility function with a higher value for being

home than for going swimming. Intuitively, however, each statement tells us more about the man's psychological state and provides some insight into why swimming had a low utility. That is, rather than only reasoning about high utilities and associating them with goals, we are also sensitive to the costs and rewards underlying these utilities.

### **Predictive limitations.**

Following on the past example, after the man arrives to his house, the predictive power of goal inference vanishes. We don't know what the man will do next, or even if he will have the same goal in the future. However, each explanation above boosts our predictive power. Knowing the costs and rewards underlying the man's goal allows us to reason about how his utilities may change over time. A tired man might choose to go swimming after taking a nap; an incompetent swimmer will not.

### **Inferential limitations.**

Desires might be in direct conflict with each other (e.g., wanting to a cookie and wanting to lose weight), they might be too costly to obtain (e.g., buying a new car), or we may not know how to complete them (e.g., wanting world peace) [95]. As such, agents have to compromise and tradeoff their true desires to choose a goal. Therefore, goals aren't always aligned with agents' true preferences. This makes it critical to distinguish between high utility states (what an agent wants to do at the moment) and high reward states (what an agent intrinsically likes). If your friend buys coffee next door you won't infer that she likes it better than the coffee sold across town, but if she goes all the way across town, you'll be confident she likes it better than the coffee from the local shop.

### **Practical limitations.**

Standard goal attribution accounts assume that costs are identical for all agents. However, this is not the case. Consider this common scenario: Anne and Bob arrive to check-in at the airport and find that the entry is an empty zigzag pathway. Anne,

who is six-years-old, takes her most efficient path towards the counter: ducking under the divisions. At the same time, Bob, who is 6' tall, takes his most efficient path towards the counter, by zigzagging through the path. If we assumed that both agents were acting efficiently with respect to the same objective costs we might infer that Bob is changing his goal at every bend in the zigzag path. Thus, to infer goals we need to understand that costs vary across agents, and we need to be able to infer them.

## 6.2 The naïve utility calculus

In light of these limitations, our intuitive theory must also include some understanding of how costs and rewards jointly influence people's behavior. Recent developmental evidence suggest that we assume that agents estimate the costs and rewards associated with a goal, and chose what to do based on the difference of these two values: the utility. Formally, the utility for taking a sequence of actions  $A$  to reach state  $S$  is given by

$$U(S, A) = R(S) - C(A)$$

The higher a goal's utility, the more likely the agent will pursue it. Despite the simplicity, decomposing utilities into costs and rewards has powerful implications. Plans with high rewards and medium costs (e.g., doing something because you truly want it) are now different from plans with low rewards but even lower costs (e.g., doing something simply because it is convenient). Conversely, plans with low rewards and medium costs (e.g., foregoing something because you don't want it) are now different from plans with high rewards and even higher costs (e.g., foregoing something because it's too costly). However, the exact costs for different actions and the rewards for reaching different states vary across agents and are partially unobservable. Thus, for an observer to have the advantage of representing an agent's costs and rewards, they need to be able to infer them.

Despite the qualitative evidence for a naïve utility calculus early in development

(Chapters 3-5, the exact nature of these inferences, and the precision to which humans can make them, are open questions.

## 6.3 Computational framework

To test people’s ability to jointly infer an agent’s costs and rewards, we implemented the naïve utility calculus model and a main alternative basic goal inference model (based on [3]). In addition, to get better insight into how each difference between the two models affects the cost-reward inferences, we implemented three additional intermediate models.

### 6.3.1 Naïve Utility Calculus model sketch

This model is a direct extension of past goal-inference models [3]. However, rather than inferring the agent’s utility function, we take the inference further and decompose the utility function into the underlying costs and rewards. This joint cost-reward inference can be seamlessly adapted into the inverse planning framework, where the probability that an agent who took actions  $A$  has cost function  $C$  and reward function  $R$  is given by Bayes’ rule:

$$p(C, R|A) \propto L(A|C, R)p(C, R)$$

Here, the likelihood that the agent takes actions  $A$  given their costs and rewards  $C$  and  $R$  is determined by the resulting utility function (Equation 2). That is, this model performs Bayesian inference over a generative planning model (formalized as a Markov Decision Process; See [3] for a detailed explanation of inverse planning through MDPs) by combining the cost and reward function to generate the utility function. Critically, the model understands that costs depend on the type of action (some actions are more costly than others) and on the agent (different agents incur different costs), and, similarly, that the rewards depend on the outcome (some outcomes are more rewarding than others) and on the agent (different agents place different rewards on

the outcomes).

Simple goal inference alternative model As the main alternative we implemented a simple goal-inference model based on Baker, et al., [3]. Like the naïve utility calculus model, this model infers the unobservable utility function. However, rather than inferring an agent’s costs, it assumes that all agents incur the same costs, independent of the action they take. Thus, this model is unable to infer agents’ costs functions or to use them to infer the magnitude of the rewards.

### **6.3.2 Intermediate accounts**

#### **Competence inference model.**

This model extends the simple goal inference alternative model by allowing the costs to vary across agents. That is, this model assumes that agents incur a fixed cost for taking any action. However, it allows different agents to have different cost constants (their competence). As such, it understands that some agents may forego a high reward if the costs they would have to incur are too high. The difference between this model and the simple goal inference model quantifies the advantage an observer obtains by understanding that some agents are broadly more competent than others.

#### **Motivation inference model.**

This model is the complement of the competence inference model. As in the naïve utility calculus model, this model assumes that the cost for travelling depends on both the specific agent and the specific terrain. However, rather than inferring a separate reward value for each object, this model assumes that all objects have a constant reward value. Nevertheless, the model allows this value to vary across agents. Intuitively, the model attempts to explain agents’ behavior by inferring their full cost function, and an overall level of motivation to complete goals. This model allows us to test if people’s inferences can be explained by simply considering an agent’s overall motivation to navigate the world and the cost they incur for navigating different types of terrains.



### Competence-motivation inference model.

This last model assumes that agents' behavior is determined by two parameters: their overall competence and motivation. That is, the model assumes that each agent incurs a cost  $c$  whenever it takes an action (regardless of the terrain) and obtains a reward  $r$  whenever it collects an object (regardless of which object it collects). Although these two values are fixed for each agent, the model infers their specific value for different agents. This model, compared with the naïve utility calculus model enables us to quantify the inferential gain from giving the cost and reward functions more flexibility by allowing them to vary as a function of the objects and the terrains.

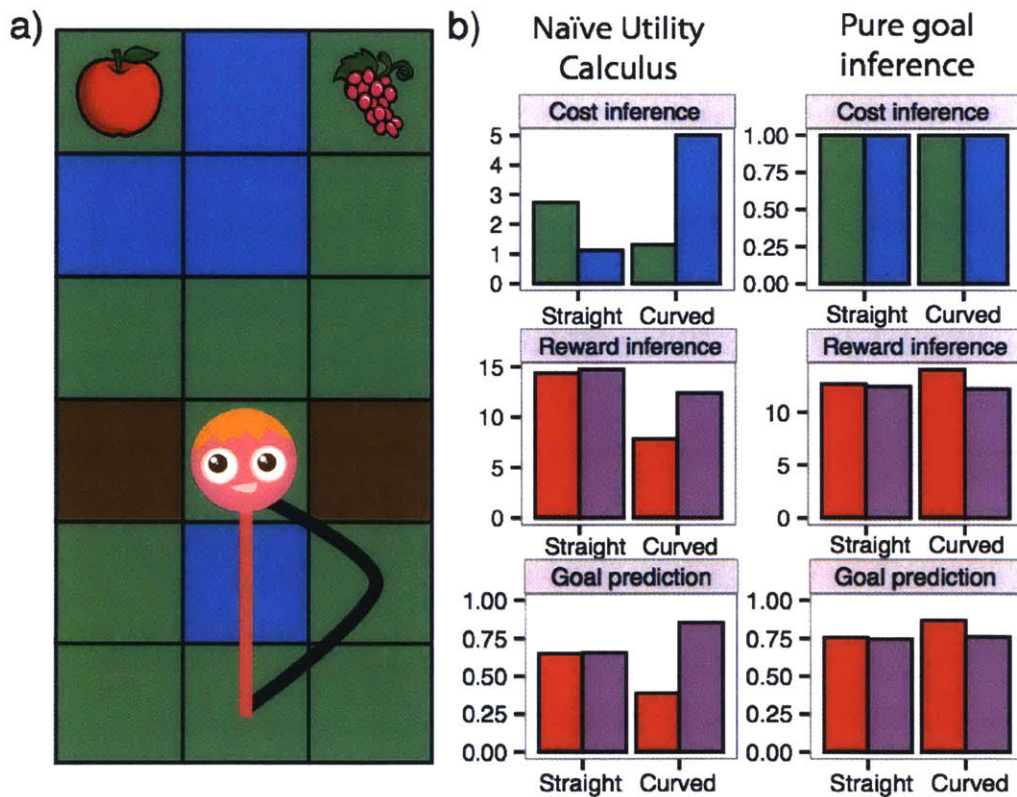


Figure 6-1: a) An agent moves from south to north towards two fruits. In the orange path, the agent moved in a straight line, while in the black path the agent circumvented the water. b) naïve utility calculus and simple goal inferences. Bars are color coded in accordance with the map.

To illustrate how the naïve utility calculus model and the simple goal-inference model differ, consider the sample path shown in Figure 6-1a. An agent is travelling

from south to north, where he can pick up either, or both, of the fruits. The terrain consists of dense jungle (in green), water (in blue), and mountains (in brown). Figure 6-1b shows the two model's inferences for two potential paths. In the straight path (orange line) the agent travelled up north in a straight line, crossing the water. In the curved path (black line), the agent travelled up north circumventing the water. As the top row shows (Figure 6-1b), for the naïve utility calculus model, the straight path implies that the agent doesn't mind crossing water, and the curved path implies that he dislikes water. In contrast, the simple goal-inference model is unable to consider these differences. The second row shows each model's inferred reward functions. When the agent takes a straight path, both models infer that he probably likes both fruits. However, when the agent takes the curved path, the naïve utility calculus model now infers that the agent prefers grapes, while the simple goal inference model does not. This is consistent with the predictions about the agent's future actions (last row). Once again, the simple goal- inference model makes similar predictions for both paths. In contrast, the naïve utility calculus model infers that the agent is more likely to pick up the grapes when it observes the curved path, but not when it observes the straight path. Although simple, this example highlights how joint cost- reward inferences help overcome the limitations raised in the past section. The naïve utility calculus can infer why the agent circumvented the water, and it can use this knowledge to predict what the agent will do next. In contrast, the pure goal-inference model interprets all actions as attempts to reach the fruits through the shortest possible path.

### 6.3.3 Experiment

To test people's ability to perform precise cost-reward inferences, we designed a simple experiment where participants were asked to infer the abilities and preferences of different agents navigating a grid world (as a static image) with three types of terrains and two types of objects.

## Design

The stimuli consisted of an 8x6 grid world with jungle, water, and mud (See Figure 6-2 for examples). Each stimulus contained the agent's starting point (which could be any of the four red squares shown in the examples in Figure 6-2), the end point (always located in the top left spot), two targets (located in any of the three possible locations shown in Figure 6-2; the apple and grape images were randomized across trials), and the agent's path. To generate the test stimulus we first ran 12,000 simulations (1,000 in each of the 12 possible worlds) of agents with random costs and rewards navigating the world (Cost and reward values were sampled from exponential distributions with parameters 0.1 and 10, respectively; these parameters were set qualitatively to ensure the simulations produced a wide range of paths). These simulations generated 189 unique paths. To reduce the stimuli size we first calculated each path's recoverability score, defined as the residual sum of squares (RSS) between the true parameters and the parameters inferred through Bayesian inference over the generative model (taking the posterior's expected value). Thus, paths with low recoverability indices had enough information for a rational observer to infer the underlying costs and rewards. Next, we calculated a discrepancy score for each alternative model, defined as the RSS between the naïve utility calculus predictions and the alternative model's predictions. Stimuli were reduced by removing all paths with a recoverability index greater than one, and then by selecting the 30 paths with the highest discrepancy score for each alternative model. The resulting 120 paths (30 for each of the four alternative models) reduced to 42 paths after removing duplicates. These 42 paths were thus ensured to contain enough information for observers to be able to make cost reward inferences (because they had a low recoverability index), and a high likelihood of helping us disambiguate between models (because they had a high discrepancy score). For each of the 42 paths we created an object version, where the map contained two fruits the protagonist could collect (See Figure 6-2), and a social version, where the map contained two agents the protagonist could help (The stimuli was otherwise identical). This allows us to test if humans make different cost-reward inferences when

reasoning about social (helping someone) and non-social goals (collecting food). For instance, humans may infer a separate reward for each outcome in non-social goals (as the naïve utility calculus model does), but only an overall level of prosociality when reasoning about social goals (as the motivation inference model does).

## **Participants**

80 U.S. residents (as determined by their IP address) were recruited and tested through Amazon's Mechanical Turk platform (Mean age = 38.59 years. Min=19 years, max=68 years).

## **Procedure**

Participants were randomly assigned to the object (N=40 participants) or the social (N=40 participants) condition. In order to keep the experiment short, each participant only completed half (21) of the trials. These trials were selected by performing random splits, guaranteeing that each path was rated exactly 20 times in the social condition and 20 times in the object condition. Participants first completed a tutorial and a brief questionnaire to ensure they understood the task. Participants who responded one or more question incorrectly were automatically redirected to the beginning of the tutorial. Participants who responded all questions correctly were given access to the test stage. In each trial, participants saw a test path on the left side of the screen (See Figure 6-2 for examples; all images were static) and five slidebars on the right side of the screen. The first three slidebars asked about the agent's ability to navigate through each type of terrain (ranging from "Extremely exhausting" to "Extremely easy", with "average" in the middle) and the last two slidebars asked about the agent's strength of preference for each fruit, or about their motivation to help each stranded agent, depending on the condition (ranging from "Not at all" to "A lot" with no text in the middle).

### 6.3.4 Results

As predicted, participants' average judgments were highly similar in the social and the object conditions ( $r=0.95$ ; 95% CI: 0.93-0.97), suggesting that people use the same type of reasoning when inferring an agent's social or non-social rewards. In light of this, all further analyses were performed using the merged judgments from both conditions.

Figure 6-2 shows example paths with the naïve utility calculus inferences and the average human judgments. Although the model qualitatively matched human judgments, there were also high discrepancies. For example, in the path on the bottom left of Figure 6-2, humans inferred that the agent had a high reward for picking up both objects (or helping both agents). In contrast, the model inferred a high reward for the first target the agent reached and a substantially lower reward for the second object, as it was conveniently located on the agent's path towards the exit state (the top left of the map). This same path illustrates how the naïve utility calculus model showed more sensitivity to costs than humans did. At the beginning of the path the agent travelled north and moved two squares across the jungle before diving into the water. The model took this as strong evidence that the agent prefers navigating through the jungle relative to the other terrains, but humans did not.

We next performed a quantitative model comparison by calculating each model's correlation with human cost and reward inferences (See Figure 6-3). To do this, each participant's data was standardized (z-scored) and then averaged. Similarly, each model's predictions were standardized (z-scored). On the cost dimension, the naïve utility calculus correlated the highest with human judgments ( $r=0.72$ ; 95% CI: 0.65-0.79), followed by the motivation inference model ( $r=0.50$ ; 95% CI: 0.40-0.61). The naïve utility calculus inferred the full reward function while the motivation inference model only inferred a single motivation parameter. Thus, this correlation difference (0.22; 95% CI: 0.09-0.34) suggests that inferring the reward function also helps recover the costs with more precision. The competence inference and the competence-motivation inference models both had correlations close to zero ( $r= -0.04$  and  $-0.01$ ,

respectively. The 95% CI for both models was between -0.20 and 0.16), suggesting that humans do not treat costs as being uniform for each agent. Last, the simple goal inference alternative model makes no cost predictions and is thus incomparable on the cost dimension.

On the reward dimension, the naïve utility calculus model showed the highest correlation ( $r=0.88$ ; 05% CI: 0.83-0.93), but it was not reliably higher than the competence inference model ( $r=0.87$ ; 95% CI: 0.82-0.93) or the simple goal inference model ( $r=0.82$ ; 0.74-0.90) (95% CI difference between naïve utility calculus and competence inference and simple goal inference: -0.07-0.09 and -0.03-0.16, respectively). The motivation inference and motivation- competence inferences performed considerably worse ( $r=0.34$  and  $0.42$ , respectively; 95% CI: 0.44-0.68 and 0.32- 0.57, respectively). Thus, our paradigm did not reveal any significant improvement in the ability to infer rewards by simultaneously inferring costs.

Last, we examined participants' individual performance by calculating their correlation with each model (See Figure 6-4). Because both cost and reward inferences were z-scored for participants and each model, we were able to calculate a joint cost-reward correlation score. All participants were correlated with the predictions generated from the model prior to data collection and no parameters were fit to individual participants. On average, participants had a correlation of 0.624 (95% CI: 0.60-0.66) with the naïve utility calculus model. Furthermore, 93.75% of participants ( $N=75$ ) showed the highest correlation with this model. Three out of the remaining five participants (6.25%) correlated better with the goal inference model and the other two participants correlated better with the motivation inference model (See Figure 6-4). This suggests that, although the naïve utility calculus model did not fit human inferences perfectly, it nevertheless clearly outperformed all other models at a global and individual level.

## Discussion

Here we proposed that the ability to reason about the costs and rewards underlying rational action is crucial for social reasoning. Inspired by developmental studies (See

Chapters 3, 4, and 5) we implemented a formal model of the naïve utility calculus and tested its performance against human inferences.

Overall, the naïve utility calculus model outperformed the simple pure goal inference model as well as intermediate models both at a global level (averaging the responses of all participants) and at an individual level (correlating model predictions with individual participants). Importantly, the naïve utility calculus was able to infer the cost function in a quantitatively similar way to human's inferences (See Figure 6-3), which no other model was able to do. However, we also found unexpected results.

First, although the naïve utility calculus made better cost inferences compared to the other models, its reward inferences were matched by the simple goal-inference and the competence inference models. Thus, we failed to find evidence that the ability to infer an agent's costs helps to infer rewards with more precision. However, a closer look at the data (See Figure 6-3) suggests that, although the models showed a high numerical reward correlation, none of the models was able to predict human judgments with high accuracy. Critically, humans' reward inferences were bimodal, with participants mostly inferring that the agents' rewards took the highest possible value, or no value at all. In contrast, the naïve utility calculus model made graded predictions. One possibility is that humans were judging whether the agent placed a reward on the outcome or not, rather than inferring its exact magnitude. Further work is needed to determine if this effect is task specific or if it fundamentally reflects how humans make reward inferences.

In addition, our experiment only used complete paths. However, as Figure 6-1 shows, a significant advantage of jointly inferring the costs and rewards comes into play before the agent has completed their goal. Models that don't take into account an agent's costs assume the agent is always taking the shortest path towards their goal (which may not necessarily be the most efficient; see Figure 6-1) and thus can make incorrect inferences. As such, it is possible that the naïve utility calculus model would outperform the other models when making reward inferences in incomplete paths.

Importantly, participants performed identically in the object and the social conditions. This suggests that humans use the same kinds of inferences to reason about social goals. Having found overall support for human's naïve utility calculus, in future work we can bring this quantitative paradigm to study how humans make social and moral evaluations. Behavioral work suggests that the same kinds of inferences influence our social evaluations (see Chapter 5 and [63]). As such, models of people's quantitative cost-reward inferences may help us understand the precise computations underlying our social evaluations and moral judgments.



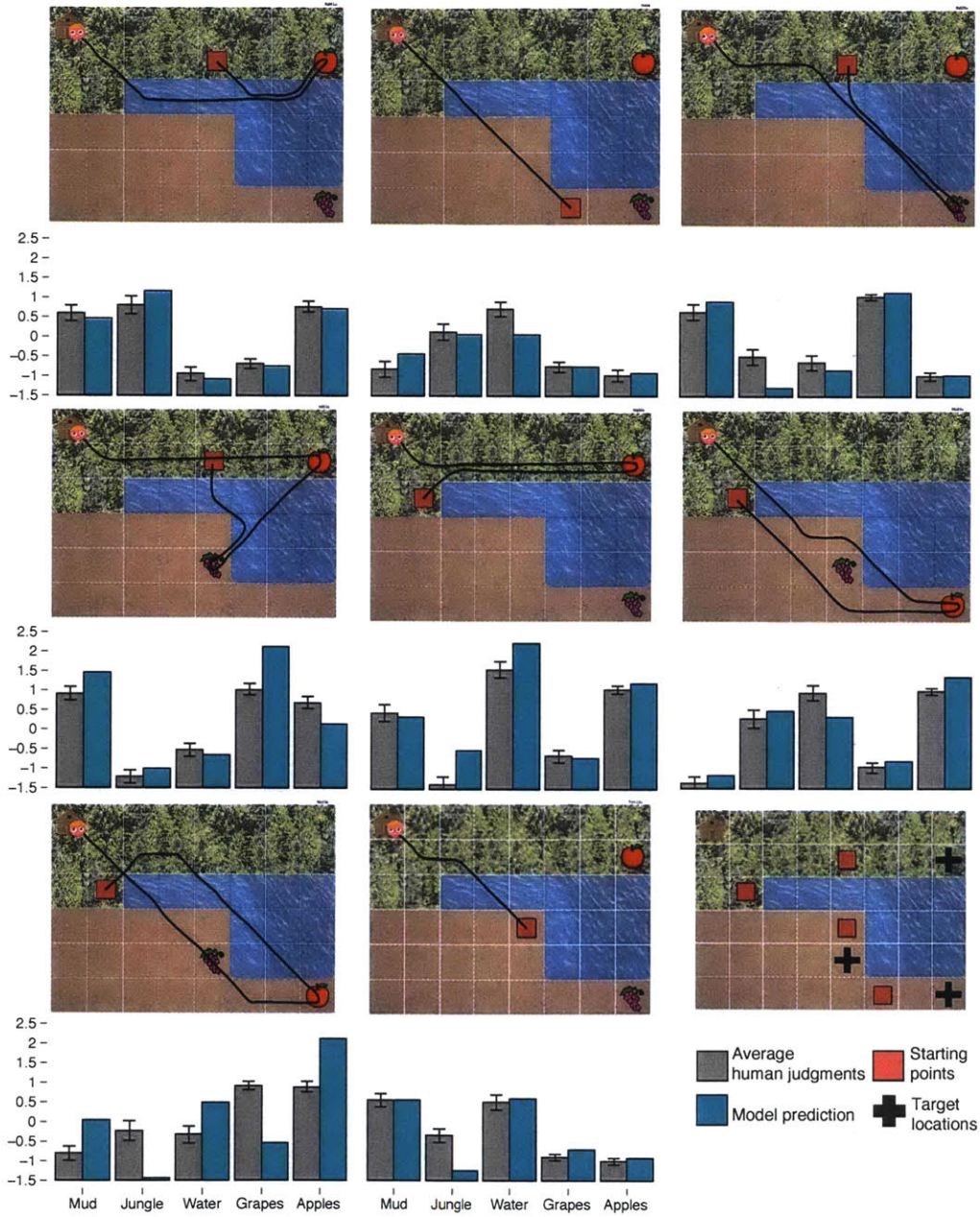


Figure 6-2: Example stimuli showing different starting points, object arrangements, and paths. Grey bars show average human judgments (z-scored per participant) with 95% confidence intervals. Teal bars show naïve utility calculus predictions.

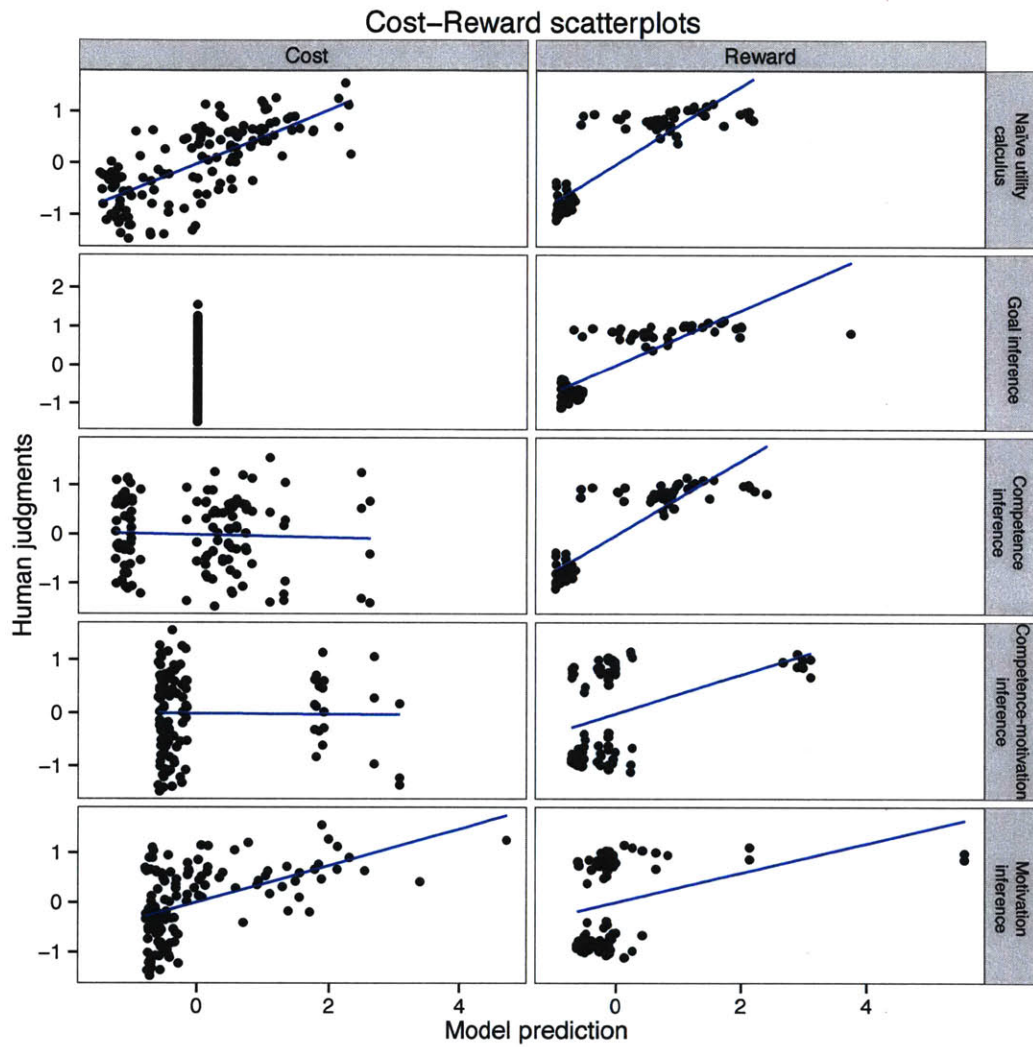


Figure 6-3: Scatterplot of model predictions (z-scored) compared to average human judgments. The x-axis shows the model predictions and y-axis shows the human (z-scored per participant) average judgments. The left column shows the cost inferences (three points per path) and the right column shows the reward inferences (two points per path). Each row shows a different model.

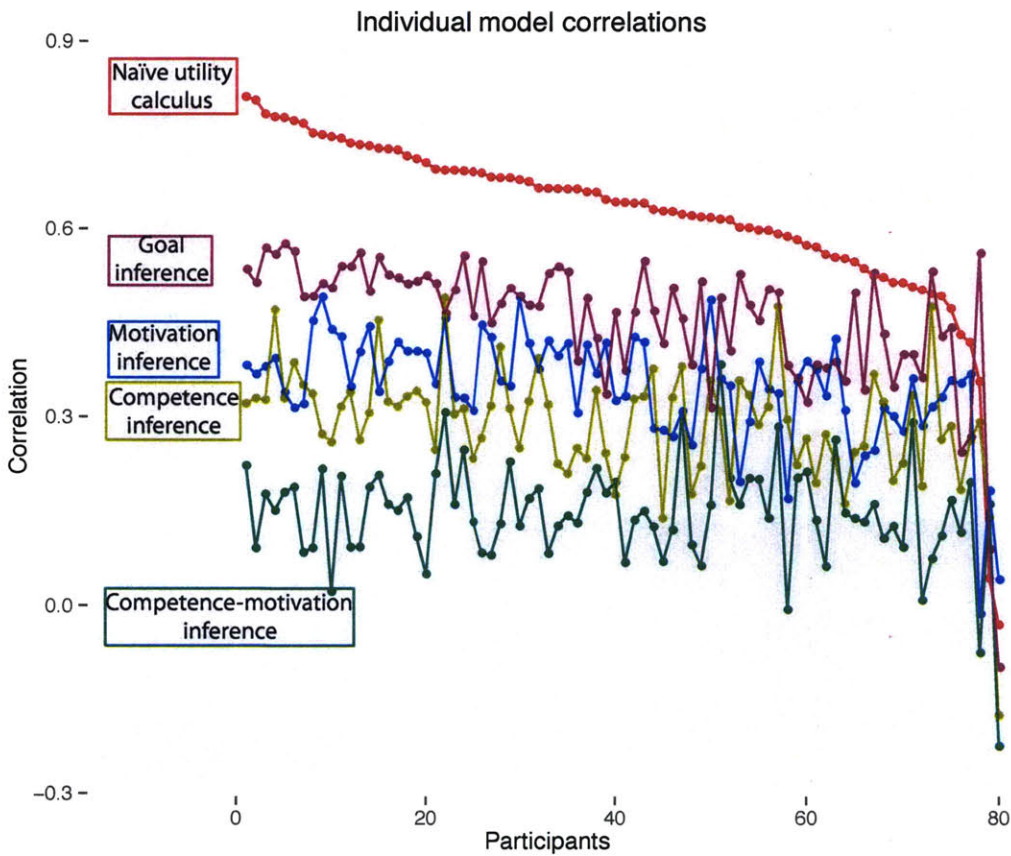


Figure 6-4: Individual model correlations with each participant. 93.75% of participants correlated best with the naïve utility calculus model. The x-axis shows all 80 participants. The y-axis shows each participant's correlation with each model. Participants are sorted by their correlation with the naïve utility calculus model. All model predictions were obtained prior to data collection and no individual parameters were fit.



# Chapter 7

## Unifying early social cognition

This chapter is based on Jara-Ettinger\*, Sun\*, Schulz, & Tenenbaum. The Naïve Utility Calculus unifies spatial and statistical routes to preference (*in prep*).

### 7.1 Introduction

Given the quantitative evidence for the naïve utility calculus (Chapter 6), here we return to the proposal that our theory may unify other inferences in the developmental literature (Chapter 2), and evaluate this proposal quantitatively with respect to account for how we infer preferences.

Our ability to infer other people’s preferences, in the service of interpreting their actions and predicting their future behavior, is at the heart of this ability. A large body of work suggests that preference inferences rely on spatial information. When we watch an agent navigate, a first focus is on the path’s end state: Agents navigate to complete goals that fulfill their desires [159]. A second focus is on the path’s directedness: We expect agents to navigate efficiently, and we use this expectation to attribute goals [36]. Thus, if an agent does not take the shortest path towards a goal this implies there is a constraint in the way [22], a subgoal that the agent completed within the path [3], or that the actions themselves are the goal [124].

When we infer preferences, however, we not only rely on what agents choose; we also take into account what they don’t choose. Suppose that an agent can pick a

fruit from a bag filled with a hundred apples and one orange. If the agent takes an apple, she doesn't necessarily like them better than oranges. But if she takes the only orange, then she probably likes them better than apples. Intuitively, the second situation reveals a stronger preference, even though the agent could have chosen either fruit in both cases. In other words, the strength of the preference inference depends on the statistical information of the possible choices. The ability to infer preferences using spatial and statistical information are both at work from early in life. Infants as young as three months old expect agents to navigate efficiently to some extent [132] and show a robust expectation by their first birthday [36]. Similarly, the ability to draw inferences from statistical information has its roots in infancy and it plays a role in how we learn what other people like [80, 157], how we learn about the world [52], and even how we learn the meaning of new words [165, 164]. Together, these two lines of evidence suggest a dual system for inferring preferences: one that relies on spatial information, and one that relies on statistical information. But real-world situations do not break down so cleanly. Agents usually combine both spatial and statistical distributions of potential rewards in their environment, and so should our judgments about their preferences from observing their actions. Here we propose that, rather than being supported two systems of knowledge, preference inferences from spatial and statistical information are derived from a single intuitive theory of agents: the naïve utility calculus (2). Here we show how the naïve utility calculus supports inferences from spatial and statistical information. We test our proposal formally by implementing a suite of preference inference models and by comparing them with adult performance on a preference-inference task. We end by discussing the implications of our findings on understanding the development of commonsense psychology.

## 7.2 The Naïve Utility Calculus

A growing set of studies suggests that humans reason about agents in terms of utility maximization (Chapters 3-5; see also [84, 66, 67]). Specifically, humans have an



intuitive theory of how utilities are comprised of costs and rewards, and how, together, they guide what others do. According to this Naïve Utility Calculus, agents decide how to act by estimating the costs and rewards associated with each possible plan, and by selecting the plan with the highest utility (the difference between rewards and costs). That is, when people watch an agent, they assume that her behavior yielded high utilities, and they use this assumption to infer the agent’s competence (her costs) and her motivation (her rewards). To illustrate inferences through the naïve utility calculus, consider an agent who chooses an apple over an orange. This implies that the utility for the apple,  $U(a)$ , is higher or equal than the utility for the orange,  $U(o)$ . By decomposing each utility into its costs ( $C(a)$  and  $C(o)$ , respectively) and rewards ( $R(a)$  and  $R(o)$ , respectively), the agent’s choice implies that  $R(a) - C(a) \geq R(o) - C(o)$ . If both fruits were equally easy to get, then  $C(a) = C(o)$  and, therefore,  $R(a) \geq R(o)$ . That is, when two options are matched for costs, agents choose what they like best. Suppose instead that the apple, the agent’s choice, was more costly to get than the orange. Because  $R(a) - C(a) \geq R(o) - C(o)$  and  $C(a) > C(o)$ , then  $R(a) > R(o)$ . That is, agents unambiguously reveal their preferences when they choose the more costly option. Last, if the apple was easier to get, then  $R(a) - C(a) \geq R(o) - C(o)$ , and  $C(a) < C(o)$ . Under these circumstances,  $U(a)$  ( $R(a) - C(a)$ ) may be higher than  $U(o)$  ( $R(o) - C(o)$ ) because the apple’s reward ( $R(a)$ ) was high, or because the orange’s cost ( $C(o)$ ) was high. Thus, when agents choose low cost options their preferences are not revealed.

### 7.2.1 Inferences from spatial information

The naïve utility calculus explains why humans are sensitive to spatial information (see 2 for a longer version of this argument). Suppose an agent takes a sequence of actions to complete a goal. If the agent maximized utilities, then two things must be true. First, the reward must outweigh the costs. Otherwise, the plan’s utility would be negative and the agent could obtain a final higher utility by not acting at all. Second, the agent must be minimizing costs: the smaller the costs the agent incurs, the higher the utility she obtains. In spatial contexts, cost minimization

reduces to efficient navigation. Thus, expecting agents to maximize utilities implies that a path's directedness and end state can help reveal preferences. If, however, humans do so through a naïve utility calculus, then, as the example above reveals (see apple-orange example), humans should also be sensitive to a third feature of spatial navigation: its cost.

### 7.2.2 Inferences from statistical information

Inferences from statistical information ultimately rely on the assumption that rare choices reveal stronger preferences. Although intuitive, the causes underlying this assumption are unclear. The naïve utility calculus, however, naturally produces this expectation (see 2 for a longer version of this argument). Suppose that an agent can take any object from a box. If she doesn't have a preference, then taking whichever object is easiest to get maximizes her utilities (if all objects have the same reward, the option with the lowest cost yields the highest utility). In contrast, if the agent prefers one type of object to the others, then she will have to incur a higher cost in terms of time, effort, attention, and distance to locate the object of the desired category and to retrieve it. The less common an object's category is, the higher the cost the agent must incur to locate it and obtain it. Thus, retrieving rare objects suggests that the agent incurred a higher cost, and, if agents maximize utilities, this cost is only warranted if the reward associated with the rare object is higher than the reward associated with more common objects.

## 7.3 Computational modeling

If humans infer preferences using their naïve utility calculus, then a formal implementation of this idea should quantitatively predict adult preference judgments. Alternatively, if humans infer preferences through simpler ways, then simpler models should predict human inferences with equal or better accuracy. To test if participants integrate spatial and statistical information through the assumption of utility maximization we ran a preference-inference task with adult participants and we compared



their performance to seven computational models: our full naïve utility calculus model and two naïve utility calculus lesioned models, as well as three alternative models; two that focus purely on spatial information and one that focuses purely on statistical information. Next, to test if participants infer these preferences in a Bayesian way, we compared participants' self-reported confidence judgments with estimates from each model.

### 7.3.1 Alternative models

#### Empirical estimate

The empirical model only uses a limited source of spatial information: the end state. This model assumes that the distribution of choices an agent makes matches her underlying preferences. For instance, if an agent collects two red objects and one blue object, then the reward for collecting a red object is  $R_{red} = 2R/3$  and the reward for collecting a blue object is  $R_{blue} = R/3$ , where  $R$  is a constant set to 1 (changing the value of  $R$  does not change our results as model comparison was done by z-scoring model predictions. See Results).

#### Sequential choice model

Similar to the empirical estimate model, this model also relies purely on spatial information. This model assumes that the agent goes through a decision making process to find objects that will give her high rewards. Specifically, the sequential choice model assumes that the agent considers one object type at a time (in a random order) and decides whether to collect an object of this type based on its reward. That is, if the agent considers taking an object from category  $k$ , she will take it with probability  $R_k / \sum_{i=1}^n R_i$ , where  $R_k$  is the reward associated with objects of category  $k$ , and  $n$  is the total number of categories. In this model, the observer assumes that the agent considers each type of object with uniform probability (if there are  $n$  types of objects, the agent considers each object type  $i$  with probability  $1/n$ ).

Given the theory of how an agent chooses what to collect, we use Bayesian infer-

ence to recover the agent’s preferences given her choices. Specifically, because in our experiment we use two types of objects (see Stimuli), we use Bayes’ rule to estimate the relative magnitude of one reward type over the other (with 0 indicating that the first category contains all the rewards, 0.5 indicating that both categories are equally rewarding, and 1 indicating that the second category contains all the rewards), using a non-informative (uniform) prior.

### **Statistical model**

Next we implemented a model from previous studies that relies purely on statistical information [52, 165, 164]. These models were formulated in simpler domains than the one we test in our experiment so we extended them to fit our experimental design. As in the sequential choice model, the observer assumes that the agent considers one random object at a time and decides whether to collect it or not based on its reward. However, in contrast to the sequential choice model, here we assume that agents consider specific objects (rather object categories) one at a time. This assumption implies that more common object categories are more likely to be considered and that less common categories are less likely to be considered. As such, selecting an object from a rare category suggests that the agent incurred additional costs in locating the object. As in the sequential choice model, we use Bayes’ rule to infer the agent’s preferences given information about her choices, using a uniform prior. This model extends sensitivity to statistical information [52] to scenarios where the space of possible reward values is continuous rather than binary. That is, this model can infer how rewarding each object category is, rather than only inferring whether the category is rewarding or not.

### **7.3.2 Naïve Utility Calculus models**

The last three models are implementations of the naïve utility calculus (naïve utility calculus), but they integrate costs in different ways, enabling us to understand how humans may reason about costs, rewards, and utility maximization. All models are

formulated as generative models that predict agent choices given their preferences, and the inference from choices to preferences is done through Bayes' rule with a non-informative prior over the possible distribution of rewards over the object categories.

### **Full Naïve Utility Calculus**

The full naïve utility calculus model assumes that agents maximize utilities. Costs are function of the number of actions the agent takes and rewards are exponentially discounted over time. Intuitively, the future discount corresponds to the assumption that the longer an agent takes to reach a reward, the less likely the reward will still be there. Thus, this model relies on spatial information in three ways: first, it expects agents to navigate efficiently because smaller sequences of actions incur fewer costs (minimizing costs), and because collecting objects faster results in higher rewards (maximizing exponentially discounted rewards); second it assumes that the agent's goals have sources of rewards; and last, and in contrast to the alternative models, it assumes that longer distances reveal stronger preferences.

### **Future-discount lesion**

The future-discount lesion is identical to the naïve utility calculus model but rewards aren't discounted over time. Thus, this model integrates statistical information in a full manner, and spatial information in a simplified manner. The model expects agents to navigate efficiently only because lower costs lead to higher utilities, but not because longer distances increases the chance of losing the target reward.

### **Action cost lesion**

Conversely, the action cost lesion model is identical to the naïve utility calculus but it ignores action costs. Nevertheless, the model assumes that the agent's rewards are discounted over time. This model therefore integrates spatial information through the expectation that agents act efficiently because the longer it takes them to reach a reward, the less likely it will still be there when they arrive.

## 7.4 Experiment

To test our models we designed a simple task where participants watched a miner collect minerals in mines with variable distributions of minerals.

### 7.4.1 Methods

#### Stimuli

7-1 shows examples of the stimuli. Each stimulus consisted of an animated display of an agent (the miner) entering a mine (a 12x12 grid world) and collecting green and/or red minerals. Each map contained 24 minerals in the same locations (which were chosen at random and kept constant across stimuli), but the proportion and the distribution of these minerals varied. The proportion varied according to three levels: more green than red (20 green and 4 red), more red than green (4 green and 20 red), or an equal number of each (12 of each). The distributions of these minerals varied according to three levels: red minerals closer, green minerals closer, or all minerals intermixed. This generated a total of nine different maps. By varying the proportion of the objects, we can test how statistical information influences preference inferences; by varying the location of the objects we can test how spatial information influences preference inferences.

The miner's paths were obtained by simulating agents who either only wanted one specific type of mineral, or agents who liked both minerals equally. These paths were generated in three different manners. In the first condition, the miner collected one mineral and exited the mine. In the second condition, the miner collected three minerals in a single trip and then exited the mine. And in the last condition the miner collected three minerals, but had to return to the mine's exit after collecting each object. Thus, the first and second conditions test how the amount of data an observer receives influences observers' inferences, and the second and third conditions together test how the costs of collecting the minerals influence observers' inferences. The combination of the two agent types (strong preference or no preference) with the nine maps produced a total of 18 test paths per condition.

## Participants

90 U.S. residents (as determined by their IP address) were recruited and tested through Amazon’s Mechanical Turk platform (Mean age = 33 years. Range = 20 - 59 years).

## Procedure

Participants were randomly assigned to the one mineral condition, to the three minerals in one trip condition, or to the three minerals in three trips condition (N = 30 participants per condition). Thus, each participant only completed one-third of the trials. Participants first completed a brief tutorial that explained the task. Next, participants completed a questionnaire with three questions to ensure they understood the task. Participants who responded all questions correctly were given access to the experiment, and participants who made at least one error were redirected to the beginning of the tutorial.

In the test stage, participants saw an animated display of the miner collecting the minerals and had to respond four questions. The first two questions were multiple choice control questions asking about the proportion and distribution of the minerals. Participants who answered these questions incorrectly were asked to re-examine the stimulus. The third question asked participants to rate the miner’s preference using a slider that ranged from “Red is much more valuable” (coded as a 0) to “Green is much more valuable” (coded as a 1). The last question asked participants to rate their confidence in the preference judgment using a slider that ranged from “Not at all” (coded as a 0) to “Extremely confident” (coded as a 1).

### 7.4.2 Results

7-2 shows the results from the experiment. As expected, all models matched the qualitative pattern of participant judgments: they predicted strong and weak preferences accurately. However, as 7-2 shows, the naïve utility calculus model appeared to capture human judgments with higher precision. To evaluate model performance

Model	Correlation (95% CI)
Empirical estimate	.84 (0.79,0.92)
Sequential choice	.82 (0.76,0.91)
Statistical	.81 (0.74,0.90)
Naïve Utility Calculus	.97 (0.96,0.98)
Future-discount lesion	.93 (0.90,0.96)
Action cost lesion	.96 (0.92,0.97)

Table 7.1: Model correlations with participant responses along with 95% bootstrapped confidence intervals.

more precisely we computed each model’s correlation with average human judgments (z-scored within each participant and averaged). 7.1 shows the results from the analysis.

### Comparison with alternative models

Overall, the naïve utility calculus model had the highest correlation ( $r=0.97$ ) between its predictions and participant responses. To evaluate this correlation we bootstrapped the correlation difference between the naïve utility calculus and the alternative models. The naïve utility calculus reliably outperformed the empirical estimate model (correlation difference=0.12; 95% CI=(0.03,0.18)), the sequential choice model (correlation difference=0.15; 95% CI=(0.05,0.21)), and the statistical model (correlation difference=0.16; 95% CI=(0.05,0.22)).

7-1 shows four example trials that reveal how the naïve utility calculus outperforms the alternative models. The empirical model fails to capture differences between trials A, B, and C, as it is not sensitive to the amount of evidence. The sequential choice model fails to capture differences between trials B and D, as it is not sensitive to the location of the objects. The statistical model roughly captures human responses, but it attributes a stronger preference to the miner in trial B, as it neglects the spatial distribution. In contrast, the naïve utility calculus models show sensitivity to the amount of data, the spatial information, and the statistical information.

## Comparison with model lesions

Both model lesions had a lower correlation with participant judgments compared to the full naïve utility calculus model (see 7.1). Removing the future-discount parameter led to a significant decrease in the model's correlation with human judgments (correlation difference = .042; 95% CI = (0.004,0.072)). This suggests that participants are sensitive to an exponential discounting of the mineral rewards over the length of the miner's trajectory. Similarly, removing the cost of travelling decreased the model's correlation with human judgments (difference = .022; 95% CI = (-0.007,0.047)). However, 13% of the mass of the 95% confidence interval was on the negative region. This suggests that integrating a linear cost over the future-discount may better fit human judgments, but the results are inconclusive.

## Confidence judgments

Our evidence so far suggests that humans infer preferences through the assumption of utility-maximization. Nevertheless, this inference is not necessarily Bayesian. To explore this possibility, we asked participants to report confidence judgments on each trial (see Methods section) and we compared them with a rough measure of each model's uncertainty: the posterior distribution's standard deviation. If participants are inferring preferences in a probabilistic manner, then the naïve utility calculus' standard deviation should correlate with participant confidence judgments. Moreover, because two of the alternative models also produce confidence judgments (the empirical model generates a single estimate with full confidence), this enables us to further test their validity.

7-3 shows each model's negative standard deviation along with participants' confidence judgments. Although the alternative models all captured preference inferences in a coarse way (see 7-2), their measures of confidence did not resemble participant's confidence judgments (see 7-3). In contrast, the naïve utility calculus model and its lesions predicted with far higher accuracy participants' confidence judgments. 7.2 shows the correlations and confidence intervals. Although the naïve utility calculus'

Model	Correlation (95% CI)	Single cat correlation	Both cat correlation
Sequential choice	0.04 (-0.25,0.32)	0.03 (-0.26,0.32)	-
Statistical	0.28 (0.04,0.56)	0.29 (0.05,0.57)	-0.43 (-1,0.01)
Naïve Utility Calculus	0.65 (0.49,0.83)	0.91 (0.88,0.95)	-0.45 (-1,-0.04)
Future-discount lesion	0.33 (0.12,0.51)	0.84 (0.79,0.89)	-0.45 (-1,-0.04)
Action cost lesion	0.68 (0.53,0.86)	0.91 (0.88,0.96)	-0.32 (-0.98,0.18)

Table 7.2: Correlation between the standard deviation of the model’s posterior distribution and participant confidence judgments, along with 95% bootstrapped confidence intervals. The first column shows the overall correlations, and the last two columns show the correlations after splitting the stimuli into the group where miner only collected one type of mineral (single category) and when the miner collected a combination of red and green minerals (both categories). The empirical model is not presented as it only produces a point estimate rather than a probability distribution. The sequential choice model produced the same prediction for all stimuli in the “both categories” group and its standard deviation was therefore incomputable.

correlations were reliably greater than 0, 7-3 reveals that it failed to capture the variation in a small set of stimuli (the results were qualitatively identical for the naïve utility calculus model lesions). Post-hoc inspection of these outliers revealed that they were all cases where the miner had selected a combination of red and green minerals (because of the way we generated the stimuli, the miner only took a combination of red and green minerals whenever these were the closest and the agent had no preference; see Stimuli section). Consistent with this, we found that when we decomposed the stimuli into trials where the agent collected only one type of mineral (Single category), the naïve utility calculus model and its lesions showed high correlations and performed roughly as well. In contrast, in the stimuli where the agent collected various kinds of minerals (Both categories), none of the models predicted human confidence judgments (see 7.2). Nevertheless, it is important to note that this subset of stimuli consists of seven data points, making it difficult to draw conclusions from the correlations.



## 7.5 Discussion

Here we reviewed evidence that, from early in life, humans can infer preferences using statistical and spatial information, and we proposed that these two types of inferences are driven by the naïve utility calculus -our intuitive theory of how agents select their goals by estimating and maximizing utilities. We tested our proposal by implementing seven computational accounts of how humans may infer preferences, and we compared their predictions to human judgments. Our results showed that adults were both sensitive to the spatial and statistical information of an agent's behavior, and that this variation was best captured by the naïve utility calculus model.

Critically, all accounts fit participant judgments qualitatively. Thus, implementing formal computational models was critical for generating precise predictions and assessing whether they explained variation in human judgments in a fine-grained manner. Our results show that the naïve utility calculus model significantly outperformed the alternative models at a detailed level.

In order to better understand the naïve utility calculus' performance we implemented two model lesions. In one model lesion we removed the future discount parameter (future-discount lesion) and in the second model lesion we removed the cost for traveling (action cost lesion). Critically, both model lesions were still sensitive to the statistical information, and they both expected the agent to navigate efficiently. The naïve utility calculus correlated with human responses better than both of the model lesions, but this difference was only reliable when comparing the naïve utility calculus model with the cost sensitive lesion and not when comparing it with the action cost lesion. Our results suggest that a non-linear reward discount is critical for how humans reason about efficiency. However, once a model integrates a future-discount parameter, adding a cost of traveling only produces a modest improvement.

Although the alternative models roughly predicted human responses, a comparison of the models' posterior standard deviation (a measure of the model's uncertainty) against participant confidence judgments revealed strong discrepancies. In contrast,

the naïve utility calculus and its lesions predicted our participant’s confidence judgments for a large set of stimuli (see 7-3 and 7.2, columns 1 and 2). Nevertheless, all models failed to capture human confidence judgments in the trials where the miner collected a combination of red and green minerals closest to the mine’s entrance (see last column of 7.2). In these situations, the naïve utility calculus models were confident that the miner liked both minerals roughly as much, and that she was therefore collecting the closest ones. Participants made similar judgments, but they reported less confidence. One possible explanation for this discrepancy is that our model assumes that the cost for traveling is fixed and observable, whereas participants may not. Instead, participants may be uncertain about how exhausting it is to travel the mine, and this may lead to a confound in the miner’s behavior: she might be taking the closest minerals because she likes all minerals just as much, or because she finds traveling deep into the mine to be very costly. A richer version of the naïve utility calculus that integrates uncertainty over the costs and rewards is needed to evaluate this possibility.

Altogether, our results show that the naïve utility calculus explains why and how humans rely on spatial and statistical information when inferring preferences. This raises two interesting hypotheses about the development of commonsense psychology. A first possibility is that the naïve utility calculus is already at work in infancy. If so, infants may use it to solve tasks involving spatial information (e.g., [36]), and tasks involving statistical information (e.g., [80]). A second possibility, however, is that the naïve utility calculus emerges later in life. Under this account, infants must then rely on simpler expectations about agents (e.g., an expectation for efficient spatial navigation, but not an expectation for cost minimization more abstractly, or utility maximization at all) to reason about spatial and statistical information. If this is true, then our results suggest that over time, these expectations provide the bedrock for a richer unifying theory: the naïve utility calculus (see Chapter 2 for other converging evidence for a naïve utility calculus).

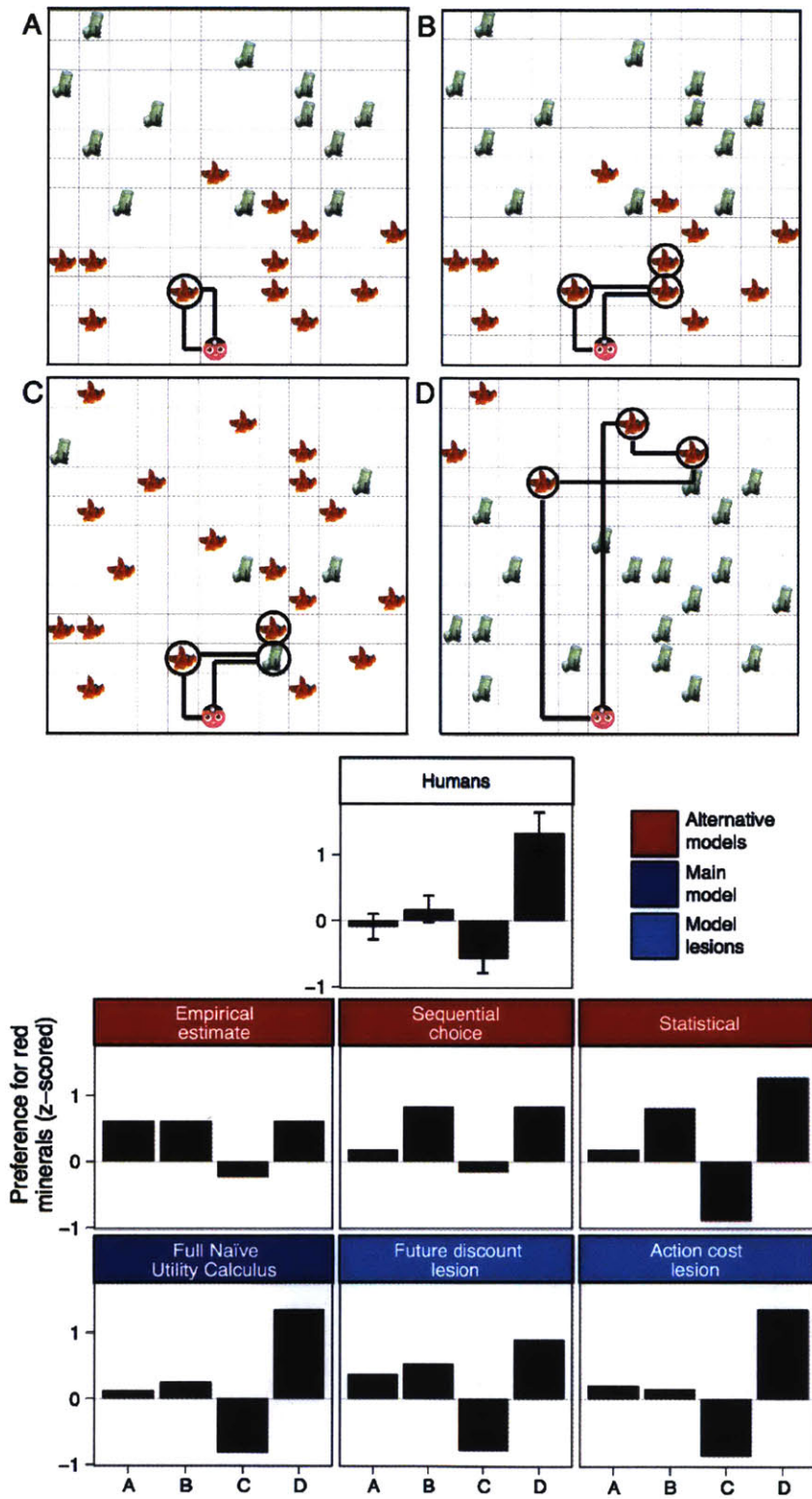


Figure 7-1: Stimuli examples along with participant judgments and model predictions.

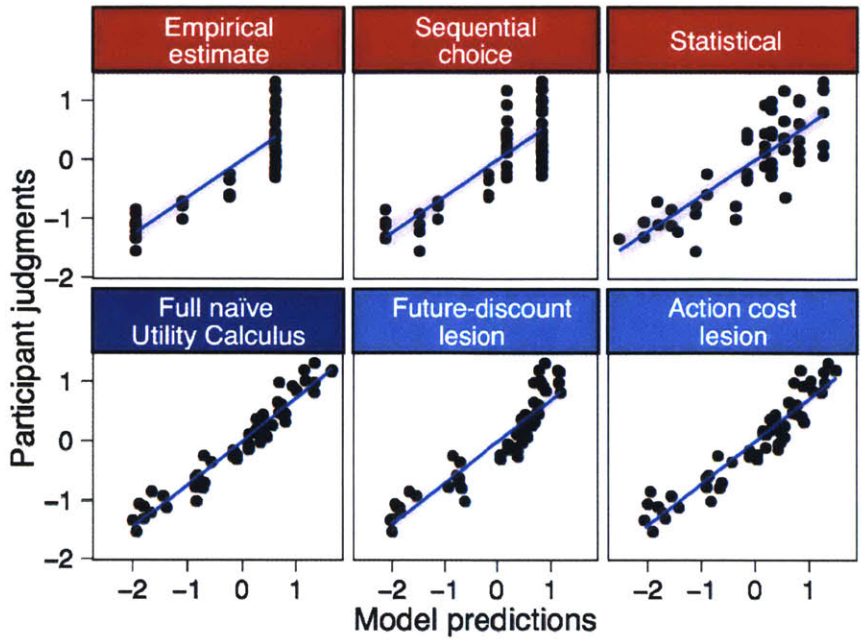


Figure 7-2: Experiment results. In each plot, each black dot represents a stimulus. The x-axis shows the model's prediction (z-scored) and the y-axis shows average participant judgments (z-scored within each participant and averaged).

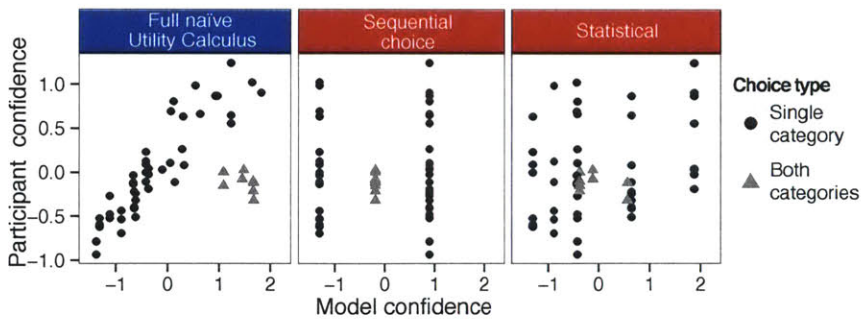


Figure 7-3: Confidence judgments. The models' confidence ratings were obtained by computing the standard deviation of the posterior distribution of each stimulus, multiplying them by -1 (so as to match the qualitative order in participant judgments), and then z-scoring the values. Participant confidence judgments were z-scored within participant and averaged.

# Chapter 8

## Discussion

Both the studies already conducted (See chapter 2) and those presented in this thesis (see Chapters 3-7) suggest that starting early in development, humans interpret other agents' behavior through a naïve utility calculus. This naïve interpretation of human behavior may be over-simplified, or even wrong in many respects. As noted, many studies suggest that human decision-making does not conform to the predictions of rational choice theory [1, 69]. Nonetheless, to the degree that a naïve utility calculus captures key features of human intentional action, it may support accurate prediction, intervention, and explanation of everyday phenomena, even if it fails to capture human decision-making in more complex domains (i.e., economic choices in the modern marketplace). By analogy, humans' intuitive physics is not an accurate account of the physical world; in comparison to quantum mechanics, or even Newtonian mechanics, it is over-simplified and even wrong in many respects [91]; nonetheless it suffices for many everyday predictions and inferences (e.g., [6]).

The studies reviewed in Chapter 2, along with the findings I present on Chapters 3-5, show successes of children in different age groups in different scenarios. Altogether, these open the possibility that some form of the naïve utility calculus is core knowledge. Even if this were the case, however, many aspects of the naïve utility calculus may nevertheless develop in crucial ways. In the remainder of the thesis, I discuss hypotheses about what in the naïve utility calculus may develop, how this may occur, and I present outstanding questions that this thesis does not answer.

## 8.1 What develops?

### 1. Dimensions of costs and rewards

As adults we assume that agents' utilities integrate many sources of costs and rewards. Time, effort, attention, or even intangible things like damaging one's reputation can be costly. Similarly, eating, learning, or having a good reputation, for example, can be rewarding. The dimensions that infants consider in utility computation are likely limited and expand over time. For instance, it is difficult to assign a cost to breaking social norms without knowing what these social norms are. I see two non-exclusive possibilities for how we may achieve this.

A first possibility is that the fundamentals of the naïve utility calculus enable us to learn new dimensions of costs and rewards. Imagine, for instance, an infant who expects agents to maximize utilities but does not realize that using time to complete a goal is costly. When this agent watches someone choose a faster route over a shorter route, the behavior will appear to fail to maximize utilities. However, if this assumption is not in question, then the actions imply that there is an additional source of costs or rewards that the observer did not consider. Situations like these may prompt infants and children to search for other dimensions in the event that may be costly or rewarding.

This account, however, faces two challenges. First, there are many ways to resolve the tension in actions that do not appear to maximize utilities. For instance, one can posit the existence of sub-goals with rewards, one can ascribe rewards to the actions themselves, or one can invoke false beliefs or uncertainty over the state of the world. Adults make these attributions to explain apparent inefficient behavior [3, 124, 4] and it is possible that toddlers or infants do so too. Second, even if children can conclude that there are additional dimensions of costs and rewards that they are failing to consider, events have a potentially infinite number of features and it is unclear how children may zoom in on the appropriate ones.

A second possibility is that infants and children learn new dimensions of costs and rewards through first person experience. Although our ability to reason about others' behavior is likely supported by a different cognitive system than the one supporting how we generate our own behavior [121, 18], empirical evidence nevertheless suggests that, even in infancy, first person experience producing actions influences action interpretation [139, 132]. Similarly, the ability to turn our commonsense psychology around and use it to introspect about our own behavior [17, 41] may help us map features of events that we find desirable or undesirable onto cost and reward judgments about ourselves. This, in turn, may enable us to learn new dimensions of costs and rewards that we can use to make sense of others' actions (e.g., an infant experiencing frustration waiting for a reward may use this experience to understand that time is costly for other agents).

## 2. Variability of costs and rewards across agents

A critical aspect of the naïve utility calculus is the understanding that costs and rewards vary across agents (they are agent-dependent). Even if infants expect agents to maximize utilities, they may nevertheless assume that costs and rewards are objective properties of the world that apply equally to all agents (agent-independent). This view is consistent with studies suggesting that infants struggle to understand that different agents can have different preferences [110]. To transition from an agent-independent theory to an agent agent-dependent theory, infants face two challenges. First, they must come up with the right theory -where costs and rewards vary across agents- and, second, they must recognize that this theory is better than the agent-independent theory -where costs and rewards are the same for all agents.

The first challenge (how infants first begin to consider the theory that costs and rewards are agent-dependent) is concerned with how we generate new hypotheses. This question appears in many domains of cognition [126] and, although this problem remains unsolved, there are some clues to how this may happen.

For simplicity, I focus on learning that costs are agent-dependent, but the same logic applies to rewards. Suppose that an infant believes that the cost for taking a step is the same for all agents (agent-independent) and equal to  $c$ . Over time, the infant will watch agents in different situations taking a different number of steps. If the infant uses this information to infer rewards, she will necessarily need to treat the total costs as event-specific. For instance, if the infant watches an agent walk 10 steps for an apple, she can conclude that the reward for the apple was higher than  $10c$  (to make the utility positive), and if she watches a different agent walk 50 steps for an orange, she can conclude that the reward for the orange is greater than  $50c$ . As such, even if the infant represents the cost for a step as being agent-independent, she will have to treat the cost for walking as event-dependent (depending on the total number of steps an agent takes). The need to treat costs as event-dependent may form the basis for infants to consider the hypothesis that costs are agent-dependent as well. Once infants consider the agent-dependent model, they face the second challenge, choosing between a simpler model (the agent-independent model) and a richer model (the agent-dependent model). A very similar problem has been studied with respect to children's false-belief understanding [40] and the same bayesian model selection approach could explain how infants transition from the agent-independent to the agent-dependent model.

### 3. Agent-independent priors on costs and rewards

Although individual differences in agents' subjective costs and rewards can only be learned from individuals themselves, agents largely overlap on what they like and dislike. For instance, most people agree that eating sweets is rewarding and that spending time is costly. These priors help observers zoom in on the appropriate cost and reward decompositions. How do we obtain them?

The answer to this question may largely depend on when we learn that cost and reward functions vary across agents (see point above). If infants treat costs and rewards as agent-dependent from the outset, then learning agent-independent



properties is largely a process of finding common features of cost and reward functions. This type of inductive problem is common and may be solved by learning overhypotheses over types of actions and types of rewards [73, 146]. In contrast, if infants first treat costs and rewards as agent-dependent and learn over time that they are agent-dependent, then the initial cost and reward functions that were treated as agent-independent may be the basis for the agent-independent priors. That is, rather than learning this hierarchical structure in an upwards fashion (by learning commonalities across utility functions) they may learn it a downwards fashion (by learning differences across utility functions).

Finally, the present work leaves three outstanding questions.

## 8.2 Outstanding questions

1. Is some form of the naïve utility calculus core knowledge? That is, do we have concepts of costs and rewards, and expect agents to maximize utilities, from birth? Or is this understanding built upon simpler theories of agency? If the naïve utility calculus is core knowledge, what parts are core and which develop (see Section 8.1)? If not, what kind of experiences drive the construction of the naïve utility calculus and what role does maturation play, if any?
2. What is the relation between our intuitive theory of a common currency, and our cultural construction of a common currency? In permitting comparison of very different kinds of costs and rewards, our monetary system resembles the naïve utility calculus. Is money, and economic systems in general, a formal expression of our intuitive theories, or does learning formal systems of exact costs and rewards influence social reasoning?
3. What are the neural mechanisms that support the naïve utility calculus? For instance, the striatum (and nucleus accumbens), VMPFC (ventromedial prefrontal cortex, OFC) and DLPFC (dorsolateral pfc) receive dopaminergic (DA) projections from the midbrain that encode prediction error signals [125], and

have been consistently implicated in tasks that involve value-guided learning and decision-making as well as the integration of costs and values, in both humans and nonhuman primates [21, 5, 145]. However, most of these studies involve choices for oneself; do the same neural circuits support reasoning about others' costs and rewards? Are they implicated in tasks in which people predict others' behaviors or making choices on behalf of others? Also, are these computations complemented by other neural circuitries that have been found to support reasoning about others' beliefs [122] or self-other comparisons [65, 167]?

### 8.2.1 Concluding remarks

I have argued that starting early in development, humans interpret other agents' behavior through a naïve utility calculus. The connection between the naïve utility calculus as an account of intuitive decision theory and formal theories of decision-making based on expected utility maximization developed in economics may appear coincidental or simply convenient, but we believe the relation runs deep. As Fritz Heider argued [59], scientific theories, especially in their early stages, may be grounded on commonsense; what better way to formulate initial hypotheses if not by what we intuitively believe to be true? Heider quotes the physicist Robert Oppenheimer: "...all sciences arise as refinement, corrections, and adaptations of common sense."

Suppose that scientific theories of human decision making, starting with classical utility theory and moving through their descendants in behavioral economics, really began grounded on the common sense theory we discussed here. This view has several implications. First, the reason that our models of common-sense psychology in children look like classical utility theory might be because early economists were, with a different purpose in mind, doing exactly what we do here: formalizing common-sense psychology. Second, our common-sense psychology is, at its core, right. Despite the memorable cases where we fail to understand each other, we get others right more often than not. Even if it fails to account for human decision-making in less ecologically relevant domains (e.g., economic choices in the modern marketplace), the naïve utility calculus, as the first models in utility theory, captures key features of human

intentional action in the most basic everyday situations even the youngest children appreciate. And as Heider observed, even when commonsense psychology is wrong with respect with how we make choices, it's still right in an important sense. Our most important everyday choices involve others, and our ability to reason about their own choices influences what we do. This intuitive decision-theory is therefore, by definition, a cornerstone of any scientific theory of human decision making.

Finally, the ways in which people's decision making fails to conform with basic assumptions of classical utility theory, which are often counterintuitive and surprising, are surprising precisely because they go against our common-sense. As such, these surprises may point to features of the naïve theory that we have not yet elucidated. To cite just one salient example, we may overinterpret others' failures to help in a low-cost situation as a sign that they don't value helping us. But maybe our naïve theories do not sufficiently take into account agents' non-optimal planning; they wanted to help but they didn't plan well. Or perhaps our naïve theories oversimplify by assuming we know all the relevant costs (or rewards) even when we don't, or assuming that others' costs are like ours even when they're not; both of these assumptions could lead us to mistake a failure to help as a low-cost refusal even when it isn't. Understanding how our commonsense psychology is oversimplified in these ways could advance not only our understanding of core social cognition as scientists, but also, ultimately, help us better understand each other as human beings.



# Appendix A

## Meta-analysis for Chapter 2 experiments

The simplest way to run a meta-analysis of the data is to aggregate the results from all experiments and analyze them as a whole. Altogether, children’s success rate was 0.7875 with 95% CI: 0.7 – 0.875. This shows that the overall pattern of data supports our theory. However, this analysis does not capture the possibility that some tasks may be easier than others, and it does not consider the constraints of how the successes and failures are distributed across the four experiments. For example, if every participant succeeded in the first three experiments, and every participant failed on the last experiment, the conclusions from this meta-analysis would be almost identical, but intuitively, the fact that the failures were all concentrated in one experiment matters - in contrast to our actual pattern of data, in which similar proportions of children (around three-quarters) succeed in all experiments. In general, a meta-analysis of the overall distribution of children’s responses should also be sensitive to each experiment’s local distribution.

Supplemental A-1 shows a sketch of the model used for the meta-analysis of our study. Children’s responses in each experiment  $i$  are modeled as a binomial distribution with unobservable bias. For any given experiment, the value of the bias is assumed to be drawn from a beta distribution with parameters  $\alpha = \mu\kappa$  and  $\beta = (1-\mu)\kappa$ . This parameterization, based on Kruschke (2010), allows for a more intuitive interpre-

tation of the parameters. Here, the beta distribution is determined by its mean value  $\mu$ , a real number between 0 and 1 which reflects the shared overall response rate across all experiments assumed to be due to the underlying naïve utility calculations, and a parameter  $\kappa$ , varying between 1 and infinity, that reflects the amount of between-experiment variability - how much children's response rate is expected to vary from one experiment to the next. (Greater values of indicate lower between-experiment variability.)

## Supplemental Figure 1

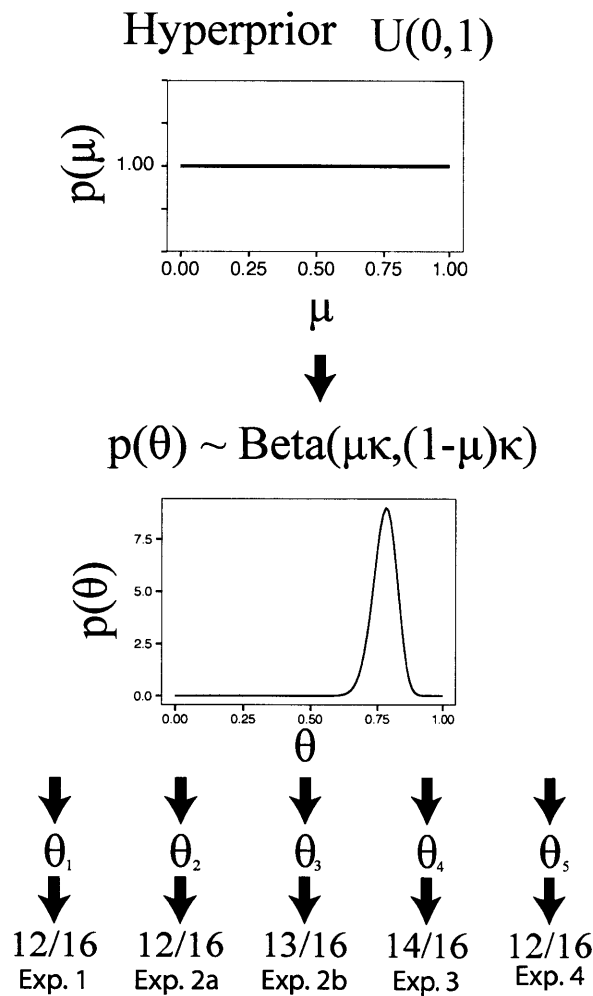


Figure A-1: Sketch of hierarchical model used for meta-analysis.

Given our complete dataset  $D = (D_1, D_2, D_3, D_4, D_5)$ <sup>1</sup> we can compute the posterior distribution of  $\mu$  through Bayes' rule:

$$p(\mu|D; \kappa) \propto L(D|\mu; \kappa)p(\mu) = \int_{\theta} p(D_i|\theta_i)p(\theta_i|\mu; \kappa)p(\mu)$$

where  $\theta = (\theta_1, \theta_2, \theta_3, \theta_4, \theta_5)$  and  $p(\mu)$  is a uniform (uninformative) distribution. We approximated the integral above numerically, with  $\mu$  varying discretely in steps of 0.01. The inferences we make about  $\mu$  depend on what inferences or assumptions we make about the variability parameter  $\kappa$ . Supplemental A-2 shows the expected value of  $\mu$  for different values of the parameter  $\kappa$ , ranging in integer steps between 1 and 150. Supplemental A-3a shows some examples of the posterior distributions on  $\mu$  for different values of  $\kappa$ . Supplemental A-3b shows the beta distributions obtained by estimating the expected value of  $\mu$  for different values of  $\kappa$ .

We can summarize these results in several ways. As  $\kappa$  increases, the expected posterior value of  $\mu$  converges to 0.78, with the 95% highest-density interval (HDI) 0.69 – 0.87 (i.e., the highest posterior probability of  $\mu$  is in this interval, containing 95% of the total probability mass). The likelihood of  $\kappa$  is strictly increasing over the range 1 – 150, so this estimate of  $\mu = 0.78$ , with 95% HDI between 0.69 and 0.87 also reflects the posterior on  $\mu$  at the maximum likelihood value of  $\kappa$ . Note that this approach returns an estimate of  $\mu$  and 95% HDIs almost identical to the simplest analysis we began with, aggregating the results of all experiments into a single large sample. In this case, that should not be surprising. The fact that the most likely value of  $\kappa$  is so high reflects the inference that the underlying response rates in each of our experiments were essentially indistinguishable, intuitively because they were likely measuring the same underlying response process in similar ways. We can also place a prior on  $\kappa$  and compute a joint posterior over  $\kappa$  and  $\mu$ . Placing a uniform prior

---

<sup>1</sup>While we presented our work as four experiments, in Experiment 2 the cookie-cracker condition and the clover-daisy condition varied slightly in the methods and, as such, we opted to allow the model to set different biases for each condition. Hence we had effectively five experiments in the meta-analysis. However, treating them as a single experiment has no impact on the conclusions.

## Supplemental Figure 2

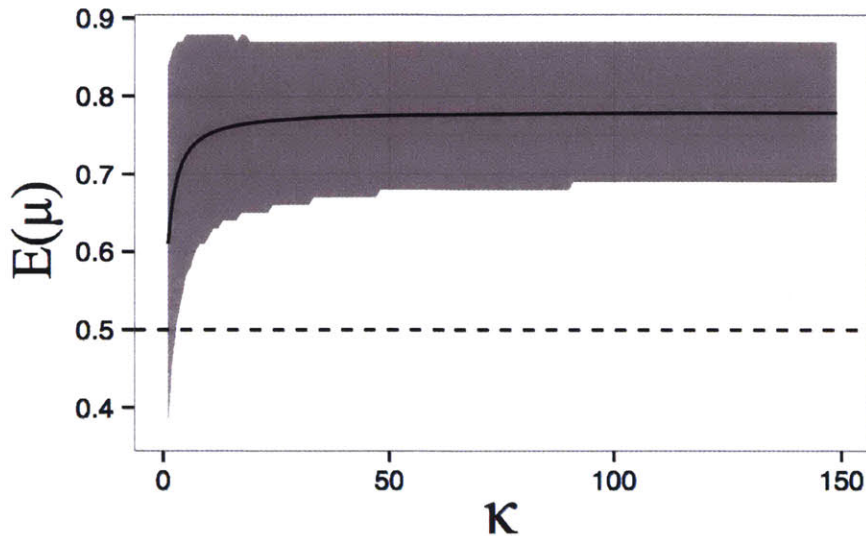


Figure A-2: Expected value of  $\mu$  as a function of  $\kappa$ .

on  $\kappa$  for values between 1 and 150 produces a posterior distribution with expected values  $\mu = 0.78$  and  $\kappa = 88$ . Under this distribution, the probability that the mean  $\mu$  is higher than chance level ( $p(\mu > 0.5)$ ) is greater than 0.99999. The corresponding beta distribution resulting from these parameters is shown in the model sketch in Supplemental A-1. Placing other reasonable priors on  $\kappa$ , such as an uninformative prior ( $1/\kappa$ ) or a gamma prior, yields extremely similar results. The evidence that  $\mu$  is well above 0.5, and most likely within 10% of 0.78 (as determined by the 95% HDI) is so strong that it is supported under the model for every individual value of  $\kappa$  tested except the two most extreme values ( $\kappa = 1, 2$ ) on the bottom end (see Supplemental A-2). However, these two values had the lowest likelihoods of all values tested ( $p(D|\kappa = 1) < 0.0001$ , and  $p(D|\kappa = 2) < 0.0005$ ), and reflect implausibly high levels of between-experiment variability.

In sum, our inferences about the most probable values of  $\mu$  are the same whether we look at the most likely values of  $\kappa$ , place a prior on  $\kappa$  and integrate it out in a joint Bayesian analysis of  $\mu$  and  $\kappa$ , or look at special cases of the model for any plausible specific value of  $\kappa$ . Taken together, the individual analyses of each of our



experiments reported in the main text, the simple meta-analysis aggregating all data into one experiment, and this meta-analysis of the underlying response rate across all our experiments all support the conclusion that children's behavior follows the predictions of the naïve utility calculus.

### Supplemental Figure 3

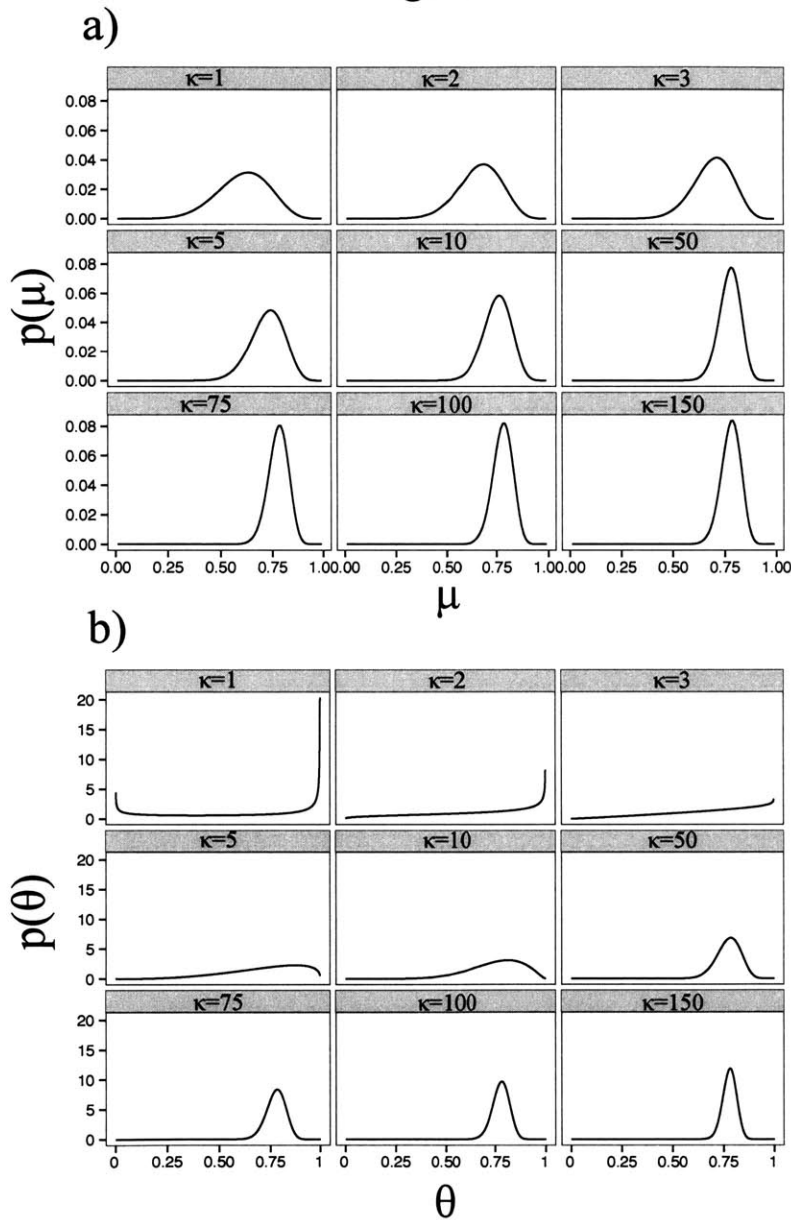


Figure A-3: (a) Posterior distribution of  $\mu$  for a range of  $\kappa$  values. (b) Corresponding beta distribution of probability of success.

# Appendix B

## Adult survey on competence and niceness

25 participants on Amazon's Mechanical Turk platform read the following story:

*Bob and Bill can both operate a special machine. Bob is very good at making the machine work. When he tries to operate the machine he immediately succeeds. Bill is not very good at making the machine work. When he tries to operate the machine he always needs to try many times before succeeding. For every other thing, Bob and Bill are just as good.*

*One day Brandon, their friend, needs help operating the machine, or else he will get fired.*

*Brandon asks Bob for help and Bob refuses to help him.*

*Brandon asks Bill for help and Bill refuses to help him.*

Next, participants responded three multiple-choice questions:

1. Who's better at making the machine go? (Bob, Bill, or Brandon).
2. If you had to choose, who is to blame for Brandon getting fired? (Bob or Bill)
3. Who would you rather be friends with? (Bob or Bill).

Last, participants were asked an open-ended question: Please explain the reason for your choice.

All participants responded the first control question correctly (selecting Bob), and all participants said Bob was more blameworthy. On the affiliation question 56%

If I had to choose I say Bob because he seems be to able to get the machine going a little better and is maybe overall smarter than Bill.
I really wouldn't prefer either. I don't know much about them.
Bob and Bill are equal except for the mastery of the 'special machine'. Even though both refused to help Brandon, Bob would have had better chance of saving Brandon's job. I would rather be friends with Bob because he can always make the special machine function.
Bob is better than bill. Simple as that.
Bob is better and should have helped. Bob is better that's why I would rather be friends.
Bob can operate the machine better.
they both sound like jerks, but if I had to choose one or the other, I would choose Bob. At least he can work the stupid machine. Maybe if we were friends he would help me.
Bob could have helped Brandon more so it's more his responsibility.
bob is more popular man.
he's better
I would want to be friends with Bob in case I needed to make the machine work.
Bob sounds more capable and smart. I'd rather have smart friends than not so smart.
Bill most likely understands that he's not good with the machine. He probably declined because he wouldn't be able to make it work. However, Bob is good with the machine and would have been able to make it work. His refusal had nothing to do with his ability.
At least if Bob was my friend, he could help me out with the machine successfully. especially if my job was on the line.

Table B.1: Responses from participants affiliating with the more competent agent.

of participants (N=14) chose the competent agent (Bob) and 44% of participants (N=11) chose the less competent agent (Bill).

Responses from participants who preferred to be friends with the competent agent (Bob) are show in Table B.1 and responses from participants who preferred to be friends with the incompetent agent (Bill) are shown in Table B.2

---

I wouldn't honestly blame either one for Brandon's not knowing what he is doing...why the hell did we hire him in the first place? you know what -YOU are fired. I have no preference for being either of their friends.

---

Bill probably refused because he knows he sucks at doing it and doesn't want to f\*\*\* up brandon's situation any more. that shows self awareness and concern. I like that

---

Bob, having superior ability with the special machine, was in a better position to aid Brandon than Bill. Thus, Bob's refusal could show a greater deficit in character.

---

Bill refused because he knew he wasn't very good at it and would have gotten Brandon fired. Bob refused because he is a jerk.

---

I find it a poor choice for Bob, who has the most experience with machines, would refuse to help and not likely want to be friends with him.

---

Bill might have refused to help because he was afraid of being blamed for any failure. Bob is probably more selfish in his refusal to help, being skilled in machine manipulation.

---

Bob sounds like a jerk, why didn't he help?

---

If Bob knows it only takes 1 time for him to get the machine to work then frankly he's just being a jerk. It's more understandable if Bill doesn't want to help given his track record with starting the machine.

---

He wasn't good at so that's why he didn't help.

---

They both refused, so neither is nice. Bob at least is good at the task...

---

I'd rather be friends with Bill because he doesn't give up when he is trying to do something and tries to figure things out. I can understand why Bill wouldn't want to help and that's because he'd be embarrassed and it'd take him a long time to help Brandon.

---

Either one of them should have helped though...

---

Table B.2: Responses from participants affiliating with the less competent agent.



# Bibliography

- [1] Maurice Allais. Le comportement de l'homme rationnel devant le risque: critique des postulats et axiomes de l'école américaine. *Econometrica: Journal of the Econometric Society*, pages 503–546, 1953.
- [2] Renée Baillargeon, Rose M Scott, and Zijing He. False-belief understanding in infants. *Trends in cognitive sciences*, 14(3):110–118, 2010.
- [3] Chris L Baker, Rebecca Saxe, and Joshua B Tenenbaum. Action understanding as inverse planning. *Cognition*, 113(3):329–349, 2009.
- [4] Chris L Baker, Rebecca R Saxe, and Joshua B Tenenbaum. Bayesian theory of mind: Modeling joint belief-desire attribution. In *Proceedings of the thirty-second annual conference of the cognitive science society*, pages 2469–2474, 2011.
- [5] Ulrike Basten, Guido Biele, Hauke R Heekeren, and Christian J Fiebach. How the brain integrates costs and benefits during decision making. *Proceedings of the National Academy of Sciences*, 107(50):21767–21772, 2010.
- [6] Peter W Battaglia, Jessica B Hamrick, and Joshua B Tenenbaum. Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences*, 110(45):18327–18332, 2013.
- [7] Tanya Behne, Malinda Carpenter, Josep Call, and Michael Tomasello. Unwilling versus unable: infants' understanding of intentional action. *Developmental psychology*, 41(2):328, 2005.
- [8] Thomas J Berndt and Kirby A Heller. Gender stereotypes and social inferences: A developmental study. *Journal of Personality and Social Psychology*, 50(5):889, 1986.
- [9] Susan AJ Birch, Nazanin Akmal, and Kristen L Frampton. Two-year-olds are vigilant of others' non-verbal cues to credibility. *Developmental science*, 13(2):363–369, 2010.
- [10] Susan AJ Birch, Sophie A Vauthier, and Paul Bloom. Three-and four-year-olds spontaneously use others' past performance to guide their learning. *Cognition*, 107(3):1018–1034, 2008.

- [11] Szilvia Biro, Stephan Verschoor, Esther Coalter, and Alan M Leslie. Outcome producing potential influences twelve-month-olds' interpretation of a novel action as goal-directed. *Infant Behavior and Development*, 37(4):729–738, 2014.
- [12] Elizabeth Bonawitz, Patrick Shafto, Hyowon Gweon, Noah D Goodman, Elizabeth Spelke, and Laura Schulz. The double-edged sword of pedagogy: Instruction limits spontaneous exploration and discovery. *Cognition*, 120(3):322–330, 2011.
- [13] Amanda C Brandone and Henry M Wellman. You can't always get what you want infants understand failed goal-directed actions. *Psychological science*, 20(1):85–91, 2009.
- [14] Roger Brown. Social psychology: The second edition. *Collier Macmillan*, 1986.
- [15] S Carey. The origin of concepts. *Oxford series in cognitive development: Vol. 3*, 10, 2009.
- [16] Jeremy I Carpendale and Michael J Chandler. On the distinction between false belief understanding and subscribing to an interpretive theory of mind. *Child Development*, 67(4):1686–1706, 1996.
- [17] Peter Carruthers. How we know our own minds: The relationship between mindreading and metacognition. *Behavioral and brain sciences*, 32(02):121–138, 2009.
- [18] Peter Carruthers and Peter K Smith. *Theories of theories of mind*. Cambridge Univ Press, 1996.
- [19] Michael J Chandler and David Helm. Developmental changes in the contribution of shared experience to social role-taking competence. *International Journal of Behavioral Development*, 7(2):145–156, 1984.
- [20] Mina Cikara, Matthew M Botvinick, and Susan T Fiske. Us versus them social identity shapes neural responses to intergroup competition and harm. *Psychological science*, 2011.
- [21] Paula L Croxson, Mark E Walton, Jill X O'Reilly, Timothy EJ Behrens, and Matthew FS Rushworth. Effort-based cost–benefit valuation and the human brain. *The Journal of Neuroscience*, 29(14):4531–4541, 2009.
- [22] Gergely Csibra, Szilvia Biró, Orsolya Koós, and György Gergely. One-year-old infants use teleological representations of actions productively. *Cognitive Science*, 27(1):111–133, 2003.
- [23] Gergely Csibra and György Gergely. The teleological origins of mentalistic action explanations: A developmental hypothesis. *Developmental Science*, 1(2):255–259, 1998.



- [24] Gergely Csibra and György Gergely. Natural pedagogy as evolutionary adaptation. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 366(1567):1149–1157, 2011.
- [25] Gergely Csibra, György Gergely, Szilvia Bíró, Orsolya Koos, and Margaret Brockbank. Goal attribution without agency cues: the perception of pure reason in infancy. *Cognition*, 72(3):237–267, 1999.
- [26] Fiery Cushman. Crime and punishment: Distinguishing the roles of causal and intentional analyses in moral judgment. *Cognition*, 108(2):353–380, 2008.
- [27] Stephanie Denison, Christie Reed, and Fei Xu. The emergence of probabilistic reasoning in very young infants: Evidence from 4.5- and 6-month-olds. *Developmental Psychology*, 49(2):243, 2013.
- [28] Stephanie Denison and Fei Xu. Integrating physical constraints in statistical inference by 11-month-old infants. *Cognitive Science*, 34(5):885–908, 2010.
- [29] Daniel Clement Dennett. *The intentional stance*. MIT press, 1989.
- [30] Kathryn M Dewar and Fei Xu. Induction, overhypothesis, and the origin of abstract knowledge evidence from 9-month-old infants. *Psychological Science*, 2010.
- [31] Kristen A Dunfield and Valerie A Kuhlmeier. Intention-mediated selective helping in infancy. *Psychological science*, 2010.
- [32] Shiri Einav and Elizabeth J Robinson. When being right is not enough four-year-olds distinguish knowledgeable informants from merely accurate informants. *Psychological science*, 2011.
- [33] Ernst Fehr and Simon Gächter. Altruistic punishment in humans. *Nature*, 415(6868):137–140, 2002.
- [34] Michael C Frank and Noah D Goodman. Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998, 2012.
- [35] György Gergely, Harold Bekkering, and Ildikó Király. Developmental psychology: Rational imitation in preverbal infants. *Nature*, 415(6873):755–755, 2002.
- [36] György Gergely and Gergely Csibra. Teleological reasoning in infancy: The naive theory of rational action. *Trends in cognitive sciences*, 7(7):287–292, 2003.
- [37] György Gergely and Gergely Csibra. The social construction of the cultural mind: Imitative learning as a mechanism of human pedagogy. *Interaction Studies*, 6(3):463–481, 2005.
- [38] György Gergely, Zoltán Nádasy, Gergely Csibra, and Szilvia Bíró. Taking the intentional stance at 12 months of age. *Cognition*, 56(2):165–193, 1995.

- [39] Tobias Gerstenberg, Noah D Goodman, David A Lagnado, and Joshua B Tenenbaum. How, whether, why: Causal judgments as counterfactual contrasts. In *Proceedings of the 37th annual conference of the cognitive science society*, 2015.
- [40] Noah D Goodman, Chris L Baker, Elizabeth Baraff Bonawitz, Vikash K Mansinghka, Alison Gopnik, Henry Wellman, Laura Schulz, and Joshua B Tenenbaum. Intuitive theories of mind: A rational approach to false belief. In *Proceedings of the twenty-eighth annual conference of the cognitive science society*, pages 1382–1387, 2006.
- [41] A Gopnik. How we know our minds: The illusion of first-person knowledge of intentionality. *Behavioral and Brain Sciences*, 18(2):390–390, 1995.
- [42] Alison Gopnik and Andrew N Meltzoff. *Words, thoughts, and theories*. Mit Press, 1997.
- [43] Alison Gopnik and Henry M Wellman. Why the child’s theory of mind really is a theory. *Mind & Language*, 7(1-2):145–171, 1992.
- [44] Alison Gopnik and Henry M Wellman. The theory theory. In *An earlier version of this chapter was presented at the Society for Research in Child Development Meeting, 1991*. Cambridge University Press, 1994.
- [45] Alison Gopnik and Henry M Wellman. Reconstructing constructivism: Causal models, bayesian learning mechanisms, and the theory theory. *Psychological bulletin*, 138(6):1085, 2012.
- [46] Herbert P Grice. *Logic and conversation*. na, 1970.
- [47] Thomas L Griffiths, Nick Chater, Charles Kemp, Amy Perfors, and Joshua B Tenenbaum. Probabilistic models of cognition: Exploring representations and inductive biases. *Trends in cognitive sciences*, 14(8):357–364, 2010.
- [48] Thomas L Griffiths, Charles Kemp, and Joshua B Tenenbaum. Bayesian models of cognition. 2008.
- [49] Hyowon Gweon, Veronica Chu, and Laura E Schulz. To give a fish or to teach how to fish? children weigh costs and benefits in considering what information to transmit. In *Proceedings of the 36th Annual Conference of the Cognitive Science Society*, pages 559–564. Cognitive Science Society Austin, TX, 2014.
- [50] Hyowon Gweon, Hannah Pelton, Jaclyn A Konopka, and Laura E Schulz. Sins of omission: Children selectively explore when teachers are under-informative. *Cognition*, 132(3):335–341, 2014.
- [51] Hyowon Gweon, Patrick Shafto, and Laura E Schulz. Children consider prior knowledge and the cost of information both in learning from and teaching others. In *Proceedings of the 36th annual conference of the Cognitive Science Society*, 2014.

- [52] Hyowon Gweon, Joshua B Tenenbaum, and Laura E Schulz. Infants consider both the sample and the sampling process in inductive generalization. *Proceedings of the National Academy of Sciences*, 107(20):9066–9071, 2010.
- [53] Katharina Hamann, Felix Warneken, Julia R Greenberg, and Michael Tomasello. Collaboration encourages equal sharing in children but not in chimpanzees. *Nature*, 476(7360):328–331, 2011.
- [54] Katharina Hamann, Felix Warneken, and Michael Tomasello. Children’s developing commitments to joint goals. *Child development*, 83(1):137–145, 2012.
- [55] J Kiley Hamlin. Moral judgment and action in preverbal infants and toddlers evidence for an innate moral core. *Current Directions in Psychological Science*, 22(3):186–193, 2013.
- [56] J Kiley Hamlin and Karen Wynn. Young infants prefer prosocial to antisocial others. *Cognitive development*, 26(1):30–39, 2011.
- [57] J Kiley Hamlin, Karen Wynn, and Paul Bloom. Social evaluation by preverbal infants. *Nature*, 450(7169):557–559, 2007.
- [58] J Kiley Hamlin, Karen Wynn, Paul Bloom, and Neha Mahajan. How infants and toddlers react to antisocial others. *Proceedings of the national academy of sciences*, 108(50):19931–19936, 2011.
- [59] Fritz Heider. *The psychology of interpersonal relations*. Psychology Press, 1958.
- [60] Joseph Henrich, Richard McElreath, Abigail Barr, Jean Ensminger, Clark Barrett, Alexander Bolyanatz, Juan Camilo Cardenas, Michael Gurven, Edwina Gwako, Natalie Henrich, et al. Costly punishment across human societies. *Science*, 312(5781):1767–1770, 2006.
- [61] Esther Herrmann, Josep Call, María Victoria Hernández-Lloreda, Brian Hare, and Michael Tomasello. Humans have evolved specialized skills of social cognition: The cultural intelligence hypothesis. *science*, 317(5843):1360–1366, 2007.
- [62] Julian Jara-Ettinger, Chris L Baker, and Joshua B Tenenbaum. Learning what is where from social observations. In *Proceedings of the Thirty-Fourth Annual Conference of the Cognitive Science Society*, pages 515–520, 2012.
- [63] Julian Jara-Ettinger, Nathaniel Kim, Paul Muentener, and Laura E Schulz. Running to do evil—how costs incurred by a perpetrator affect moral judgment. In *Proceedings of the Thirty-Sixth Annual Conference of the Cognitive Science Society*, 2014.
- [64] Vikram K Jaswal and Leslie A Neely. Adults don’t always know best preschoolers use past reliability over age when learning new words. *Psychological Science*, 17(9):757–758, 2006.

- [65] Adrianna C Jenkins, C Neil Macrae, and Jason P Mitchell. Repetition suppression of ventromedial prefrontal activity during judgments of self and others. *Proceedings of the National Academy of Sciences*, 105(11):4507–4512, 2008.
- [66] Alan Jern and Charles Kemp. Reasoning about social choices and social relationships. Cognitive Science Society, 2014.
- [67] Samuel GB Johnson and Lance J Rips. Do the right thing: The assumption of optimality in lay decision theory and causal judgment. *Cognitive psychology*, 77:42–76, 2015.
- [68] Jillian J Jordan, Katherine McAuliffe, and Felix Warneken. Development of in-group favoritism in children’s third-party punishment of selfishness. *Proceedings of the National Academy of Sciences*, 111(35):12710–12715, 2014.
- [69] Daniel Kahneman. *Thinking, fast and slow*. Macmillan, 2011.
- [70] Charles W Kalish. Children’s predictions of consistency in people’s actions. *Cognition*, 84(3):237–265, 2002.
- [71] Harold H Kelley. The processes of causal attribution. *American psychologist*, 28(2):107, 1973.
- [72] Harold H Kelley. *An atlas of interpersonal situations*. Cambridge University Press, 2003.
- [73] Charles Kemp, Amy Perfors, and Joshua B Tenenbaum. Learning overhypotheses with hierarchical bayesian models. *Developmental science*, 10(3):307–321, 2007.
- [74] J Kiley Hamlin, Tomer Ullman, Josh Tenenbaum, Noah Goodman, and Chris Baker. The mentalistic basis of core social cognition: experiments in preverbal infants and a computational model. *Developmental science*, 16(2):209–226, 2013.
- [75] Katherine D Kinzler, Kristin Shutts, Jasmine DeJesus, and Elizabeth S Spelke. Accent trumps race in guiding children’s social preferences. *Social cognition*, 27(4):623, 2009.
- [76] Joshua Knobe. Theory of mind and moral cognition: Exploring the connections. *Trends in cognitive sciences*, 9(8):357–359, 2005.
- [77] Melissa A Koenig, Fabrice Clément, and Paul L Harris. Trust in testimony: Children’s use of true and false statements. *Psychological Science*, 15(10):694–698, 2004.
- [78] Valerie Kuhlmeier, Karen Wynn, and Paul Bloom. Attribution of dispositional states by 12-month-olds. *Psychological Science*, 14(5):402–408, 2003.

- [79] Tamar Kushnir, Christopher Vredenburgh, and Lauren A Schneider. “who can help me fix this toy?” the distinction between causal knowledge and word knowledge guides preschoolers’ selective requests for information. *Developmental psychology*, 49(3):446, 2013.
- [80] Tamar Kushnir, Fei Xu, and Henry M Wellman. Young children use statistical sampling to infer the preferences of other people. *Psychological Science*, 21(8):1134–1140, 2010.
- [81] Alan M Leslie. Infant perception of a manual pick-up event. *British Journal of Developmental Psychology*, 2(1):19–32, 1984.
- [82] Ulf Liszkowski, Malinda Carpenter, and Michael Tomasello. Twelve-month-olds communicate helpfully and appropriately for knowledgeable and ignorant partners. *Cognition*, 108(3):732–739, 2008.
- [83] David Liu, Susan A Gelman, and Henry M Wellman. Components of young children’s trait understanding: Behavior-to-trait inferences and trait-to-behavior predictions. *Child Development*, 78(5):1543–1558, 2007.
- [84] Christopher G Lucas, Thomas L Griffiths, Fei Xu, Christine Fawcett, Alison Gopnik, Tamar Kushnir, Lori Markson, and Jane Hu. The child as econometrician: A rational model of preference understanding in children. *PloS one*, 9(3):e92160, 2014.
- [85] Yuyan Luo and Renée Baillargeon. Can a self-propelled box have a goal? psychological reasoning in 5-month-old infants. *Psychological Science*, 16(8):601–608, 2005.
- [86] Lili Ma and Fei Xu. Preverbal infants infer intentional agents from the perception of regularity. *Developmental psychology*, 49(7):1330, 2013.
- [87] Bertram F Malle. How the mind explains behavior. *Folk Explanation, Meaning and Social Interaction*. Massachusetts: MIT-Press, 2004.
- [88] Ellen M Markman. *Categorization and naming in children: Problems of induction*. Mit Press, 1991.
- [89] Danielle Matthews, Elena Lieven, Anna Theakston, Michael Tomasello, et al. The effect of perceptual availability and prior discourse on young children’s use of referring expressions. *Applied Psycholinguistics*, 27(3):403, 2006.
- [90] Katherine McAuliffe, Jillian J Jordan, and Felix Warneken. Costly third-party punishment in young children. *Cognition*, 134:1–10, 2015.
- [91] Michael McCloskey, Alfonso Caramazza, and Bert Green. Curvilinear motion in the absence of external forces: Naive beliefs about the motion of objects. *Science*, 210(4474):1139–1141, 1980.

- [92] Andrew N Meltzoff. Understanding the intentions of others: re-enactment of intended acts by 18-month-old children. *Developmental psychology*, 31(5):838, 1995.
- [93] Andrew N Meltzoff. ‘like me’: a foundation for social cognition. *Developmental science*, 10(1):126–134, 2007.
- [94] John Stuart Mill. *Utilitarianism*. Longmans, Green and Company, 1901.
- [95] Louis J Moses. Some thoughts on ascribing complex intentional concepts to young children. *Intentions and intentionality: Foundations of social cognition*, pages 69–83, 2001.
- [96] Paul Muentener and Susan Carey. Infants’ causal representations of state change events. *Cognitive Psychology*, 61(2):63–86, 2010.
- [97] Paul Muentener and Laura Schulz. Toddlers infer unobserved causes for spontaneous events. *Frontiers in psychology*, 5, 2014.
- [98] George E Newman, Frank C Keil, Valerie A Kuhlmeier, and Karen Wynn. Early understandings of the link between agents and order. *Proceedings of the National Academy of Sciences*, 107(40):17140–17145, 2010.
- [99] Erika Nurmsoo and Elizabeth J Robinson. Children’s trust in previously inaccurate informants who were well or poorly informed: When past errors can be excused. *Child development*, 80(1):23–27, 2009.
- [100] Daniela K O’Neill. Two-year-old children’s sensitivity to a parent’s knowledge state when making requests. *Child development*, pages 659–677, 1996.
- [101] Kristine H Onishi and Renée Baillargeon. Do 15-month-old infants understand false beliefs? *Science*, 308(5719):255–258, 2005.
- [102] Harriet Over and Malinda Carpenter. Priming third-party ostracism increases affiliative imitation in children. *Developmental science*, 12(3):F1–F8, 2009.
- [103] Elisabeth S Pasquini, Kathleen H Corriveau, Melissa Koenig, and Paul L Harris. Preschoolers monitor the relative accuracy of informants. *Developmental psychology*, 43(5):1216, 2007.
- [104] Josef Perner. *Understanding the representational mind*. The MIT Press, 1991.
- [105] Josef Perner and Johannes Roessler. From infants’ to children’s appreciation of belief. *Trends in cognitive sciences*, 16(10):519–525, 2012.
- [106] Webb Phillips, Jennifer L Barnes, Neha Mahajan, Mariko Yamaguchi, and Laurie R Santos. Unwilling versus unable: capuchin monkeys’(cebus apella) understanding of human intentional action. *Developmental science*, 12(6):938–945, 2009.

- [107] Lindsey J Powell and Elizabeth S Spelke. Preverbal infants expect members of social groups to act alike. *Proceedings of the National Academy of Sciences*, 110(41):E3965–E3972, 2013.
- [108] David Premack. The infant’s theory of self-propelled objects. *Cognition*, 36(1):1–16, 1990.
- [109] David G Rand, Joshua D Greene, and Martin A Nowak. Spontaneous giving and calculated greed. *Nature*, 489(7416):427–430, 2012.
- [110] Betty M Repacholi and Alison Gopnik. Early reasoning about desires: evidence from 14-and 18-month-olds. *Developmental psychology*, 33(1):12, 1997.
- [111] Marjorie Rhodes, Elizabeth Bonawitz, Patrick Shafto, Annie Chen, and L Caglar. Controlling the message: Preschoolers’ use of evidence to teach and deceive others. *Frontiers in Psychology*, 6, 2015.
- [112] William S Rholes and Diane N Ruble. Children’s understanding of dispositional characteristics of others. *Child Development*, pages 550–560, 1984.
- [113] Elizabeth J Robinson and EL Whitcombe. Children’s suggestibility in relation to their understanding about sources of knowledge. *Child development*, 74(1):48–62, 2003.
- [114] Samuel Ronfard, Alexandra M Was, and Paul L Harris. Children teach methods they could not discover for themselves. *Journal of experimental child psychology*, 142:107–117, 2016.
- [115] Ken J Rotenberg. Children’s use of intentionality in judgments of character and disposition. *Child Development*, pages 282–284, 1980.
- [116] Ken J Rotenberg. Development of character constancy of self and other. *Child Development*, pages 505–515, 1982.
- [117] Diane N Ruble and Carol S Dweck. Self-conceptions, person conceptions, and their development. *Social development*, 15, 1995.
- [118] Bertrand Russell. *ABC of Relativity*. Routledge, 2009.
- [119] Mark A Sabbagh and Dare A Baldwin. Learning words from knowledgeable versus ignorant speakers: Links between preschoolers’ theory of mind and semantic development. *Child development*, 72(4):1054–1070, 2001.
- [120] R Saxe, JB Tenenbaum, and S Carey. Secret agents inferences about hidden causes by 10-and 12-month-old infants. *Psychological Science*, 16(12):995–1001, 2005.
- [121] Rebecca Saxe. Against simulation: the argument from error. *Trends in cognitive sciences*, 9(4):174–179, 2005.

- [122] Rebecca Saxe and Nancy Kanwisher. People thinking about thinking people: the role of the temporo-parietal junction in “theory of mind”. *Neuroimage*, 19(4):1835–1842, 2003.
- [123] Rebecca Saxe, Tania Tzelnic, and Susan Carey. Knowing who dunnit: Infants identify the causal agent in an unseen causal interaction. *Developmental Psychology*, 43(1):149, 2007.
- [124] Adena Schachner and Susan Carey. Reasoning about irrational actions: When intentional movements cannot be explained, the movements themselves are seen as the goal. *Cognition*, 129(2):309–327, 2013.
- [125] Wolfram Schultz, Peter Dayan, and P Read Montague. A neural substrate of prediction and reward. *Science*, 275(5306):1593–1599, 1997.
- [126] Laura Schulz. Finding new facts; thinking new thoughts. *Rational constructivism in cognitive development. Advances in child development and behavior*, 43:269–294, 2012.
- [127] Rose M Scott and Renée Baillargeon. Do infants really expect agents to act efficiently? a critical test of the rationality principle. *Psychological science*, page 0956797612457395, 2013.
- [128] Elizabeth Seiver, Alison Gopnik, and Noah D Goodman. Did she jump because she was the big sister or because the trampoline was safe? causal inference and the development of social attribution. *Child development*, 84(2):443–454, 2013.
- [129] Patrick Shafto, Noah D Goodman, and Michael C Frank. Learning from others the consequences of psychological reasoning for human learning. *Perspectives on Psychological Science*, 7(4):341–351, 2012.
- [130] Patrick Shafto, Noah D Goodman, and Thomas L Griffiths. A rational account of pedagogical reasoning: Teaching by, and learning from, examples. *Cognitive psychology*, 71:55–89, 2014.
- [131] Laura Shneidman, Hyowon Gweon, Laura E Schulz, and Amanda L Woodward. Learning from others and spontaneous exploration: A cross-cultural investigation laura shneidman1, hyowon gweon2, laura e. schulz3 and amanda l. woodward1. 2016.
- [132] Amy E Skerry, Susan E Carey, and Elizabeth S Spelke. First-person action experience reveals sensitivity to action efficiency in prereaching infants. *Proceedings of the National Academy of Sciences*, 110(46):18728–18733, 2013.
- [133] Amy E Skerry and Elizabeth S Spelke. Preverbal infants identify emotional reactions that are incongruent with goal outcomes. *Cognition*, 130(2):204–216, 2014.



- [134] Stephanie Sloane, Renée Baillargeon, and David Premack. Do infants have a sense of fairness? *Psychological science*, page 0956797611422072, 2012.
- [135] Adam Smith. The theory of moral sentiments, ed. *DD Raphael & AL Macfie. Liberty Fund. (Original work published in 1759.)*[ELK], 1759.
- [136] David M Sobel and Kathleen H Corriveau. Children monitor individuals' expertise for word learning. *Child development*, 81(2):669–679, 2010.
- [137] David M Sobel and Tamar Kushnir. Knowledge matters: How children evaluate the reliability of testimony as a process of rational inference. *Psychological review*, 120(4):779, 2013.
- [138] Jessica A Sommerville, Marco FH Schmidt, Jung-eun Yun, and Monica Burns. The development of fairness expectations and prosocial behavior in the second year of life. *Infancy*, 18(1):40–66, 2013.
- [139] Jessica A Sommerville, Amanda L Woodward, and Amy Needham. Action experience alters 3-month-old infants' perception of others' actions. *Cognition*, 96(1):B1–B11, 2005.
- [140] Victoria Southgate, Coralie Chevallier, and Gergely Csibra. Seventeen-month-olds appeal to false beliefs to interpret others' referential communication. *Developmental science*, 13(6):907–912, 2010.
- [141] Victoria Southgate, Mark H Johnson, and Gergely Csibra. Infants attribute goals even to biomechanically impossible actions. *Cognition*, 107(3):1059–1069, 2008.
- [142] Victoria Southgate, Atsushi Senju, and Gergely Csibra. Action anticipation through attribution of false belief by 2-year-olds. *Psychological Science*, 18(7):587–592, 2007.
- [143] Elizabeth S Spelke. Core knowledge. *American psychologist*, 55(11):1233, 2000.
- [144] Alex J Stiller, Noah D Goodman, and Michael C Frank. Ad-hoc implicature in preschool children. *Language Learning and Development*, 11(2):176–190, 2015.
- [145] Deborah Talmi, Peter Dayan, Stefan J Kiebel, Chris D Frith, and Raymond J Dolan. How humans integrate the prospects of pain and reward during choice. *The Journal of neuroscience*, 29(46):14617–14626, 2009.
- [146] Joshua B Tenenbaum, Thomas L Griffiths, and Charles Kemp. Theory-based bayesian models of inductive learning and reasoning. *Trends in cognitive sciences*, 10(7):309–318, 2006.
- [147] Joshua B Tenenbaum, Charles Kemp, Thomas L Griffiths, and Noah D Goodman. How to grow a mind: Statistics, structure, and abstraction. *science*, 331(6022):1279–1285, 2011.

- [148] Michael Tomasello. *Why we cooperate*. MIT press, 2009.
- [149] Michael Tomasello, Malinda Carpenter, Josep Call, Tanya Behne, and Henrike Moll. Understanding and sharing intentions: The origins of cultural cognition. *Behavioral and brain sciences*, 28(05):675–691, 2005.
- [150] Tomer Ullman, Chris Baker, Owen Macindoe, Owain Evans, Noah Goodman, and Joshua B Tenenbaum. Help or hinder: Bayesian models of social goal inference. In *Advances in neural information processing systems*, pages 1874–1882, 2009.
- [151] John Von Neumann and Oskar Morgenstern. *Theory of Games and Economic Behavior: 3d Ed.* Princeton University Press, 1953.
- [152] Felix Warneken and Michael Tomasello. Altruistic helping in human infants and young chimpanzees. *science*, 311(5765):1301–1303, 2006.
- [153] Felix Warneken and Michael Tomasello. The roots of human altruism. *British Journal of Psychology*, 100(3):455–471, 2009.
- [154] Henry M Wellman. *The child’s theory of mind*. 1990.
- [155] Henry M Wellman. *Making minds: How theory of mind develops*. Oxford University Press, 2014.
- [156] Henry M Wellman, David Cross, and Julianne Watson. Meta-analysis of theory-of-mind development: the truth about false belief. *Child development*, 72(3):655–684, 2001.
- [157] Henry M Wellman, Tamar Kushnir, Fei Xu, and Kimberly A Brink. Infants use statistical sampling to understand the psychological world. *Infancy*, 2015.
- [158] Heinz Wimmer and Josef Perner. Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception. *Cognition*, 13(1):103–128, 1983.
- [159] Amanda L Woodward. Infants selectively encode the goal object of an actor’s reach. *Cognition*, 69(1):1–34, 1998.
- [160] Amanda L Woodward, Jessica A Sommerville, and Jose J Guajardo. How infants make sense of intentional action. *Intentions and intentionality: Foundations of social cognition*, pages 149–169, 2001.
- [161] Yang Wu, Paul Muentener, and Laura E Schulz. The invisible hand: Toddlers connect probabilistic events with agentive causes. *Cognitive science*, 2015.
- [162] Fei Xu and Stephanie Denison. Statistical inference and sensitivity to sampling in 11-month-old infants. *Cognition*, 112(1):97–104, 2009.

- [163] Fei Xu and Vashti Garcia. Intuitive statistics by 8-month-old infants. *Proceedings of the National Academy of Sciences*, 105(13):5012–5015, 2008.
- [164] Fei Xu and Joshua B Tenenbaum. Sensitivity to sampling in bayesian word learning. *Developmental science*, 10(3):288–297, 2007.
- [165] Fei Xu and Joshua B Tenenbaum. Word learning as bayesian inference. *Psychological review*, 114(2):245, 2007.
- [166] Liane Young, Fiery Cushman, Marc Hauser, and Rebecca Saxe. The neural basis of the interaction between theory of mind and moral judgment. *Proceedings of the National Academy of Sciences*, 104(20):8235–8240, 2007.
- [167] Jamil Zaki, Gilberto López, and Jason P Mitchell. Activity in ventromedial prefrontal cortex covaries with revealed social preferences: Evidence for person-invariant value. *Social cognitive and affective neuroscience*, page nst005, 2013.