

1 **Materials and Methods:**

2 **Experiment 1**

3 **Methods**

4 **Participants.** Our final sample included 24 healthy, full-term infants (15 female,
5 $M_{age}=9.95m$, range=9.43-10.53). Eight more infants were tested, excluded, and replaced (1 for
6 fussiness that prevented study completion, 1 for technical failure, 5 for inattentiveness during test
7 events, and 1 for experimental error). Sample size and exclusion criteria were fixed prior to the
8 start of data collection, and decisions concerning exclusions were made by researchers unaware
9 of the order of the test events viewed by the infant (and therefore blind to the data that the infant
10 provided). All participants were recruited from the greater Boston area and tested at the
11 Laboratory for Developmental Studies at Harvard University with parental informed consent.
12 Families received a small thank-you gift (e.g. a t-shirt or toy) for participating. All study
13 protocols were approved by the Committee on the Use of Human Subjects at Harvard University.

14 **Materials and Design.** All animated events were created in Blender (38), synchronized
15 with a custom audio track in iMovie, and presented using Keynote on a 101.6cm by 132.1cm
16 LCD projector screen. Two speakers flanking the screen played all stimuli-related sounds.
17 Infants' looking time data were coded online using Xhab64 (39) software and offline using jHab
18 (40).

19 The experiment consisted of 3 pairs of familiarization trials, 1 pre-test trial, and 2 pairs of
20 test trials. All familiarization (Movie S1) and test trials (Movie S2) began with an attention-
21 getting animation and sound (3.0s), followed by looped sequences of events. *Familiarization*
22 sequences consisted of 4 videos (8s or 8.9s; see below) and *test* sequences consisted of a single
23 video (5.2s), each played on a loop with black screens (0.5s) interspersed between the events.
24 Events featured three agents with eyes and a smiling mouth: a central red spherical agent, and
25 two target agents (a cone and cylinder). All familiarization events featured the agent paired with
26 one of the two targets at a time; all test events featured both targets together. The timing of all
27 actions was held constant within familiarization and test blocks. The location of the targets was
28 constant across participants, and the identity of the higher value target was counterbalanced
29 across participants.

30 During *familiarization* (Movie S1), the agent responded to the call of one of the targets
31 by accepting or refusing to jump over a small (1 units tall), medium (6 units), or large (10 units)
32 barrier that fell with a thud between the agent and target. In events where the agent accepted the
33 cost (8.9s), it looked up at the barrier, made a positive "Mmmm!" sound, and leapt over it to
34 reach the target. The agent always acted efficiently (Gergely et al., 1995) adapting the height of
35 its jump to the height of the barrier; all jumps were accompanied with a popping sound. In events
36 in which the agent refused the cost (8.0s), it looked up at the barrier, made a mildly negative
37 "Hmmm..." sound, and backed away, returning to the center of the screen. Each block of
38 familiarization events consisted of 4 videos, wherein the agent accepted a small cost and refused
39 a medium cost for one target, and accepted a medium cost and refused a large cost for the other
40 target. (For convenience, we describe events where the agent does not take a costly action as
41 "refusing" or "declining" to take an action, which is how these events appear to adults. But for
42 the purposes of our hypothesis about infants' action understanding, and for our computational
43 model, what matters in these cases is only that the agent is presented with a costly action,
44 considers it and does not take it.) Each familiarization trial consisted of looped blocks either in
45 the above order (small, medium, medium, large) or the opposite order. Across the 6
46 familiarization trials, both orders were presented 3 times in an ABABAB pattern. The identity of

1 the higher value target and the first familiarization block were counterbalanced across
2 participants.

3 A single non-looped *pre-test* event following familiarization featured a still image of the
4 two targets without the agent. During *test* (Movie S4), the agent reappeared, rotated left then
5 right while saying “Hmmm...”, and then approached one of the targets: either the higher- or
6 lower value agent. The same sound accompanied each approach. Across 2 pairs of test trials
7 presented in alternation, the agent approached the Triangle target twice and the Square target
8 twice. The first test trial (higher- or lower value approach) was counterbalanced across
9 participants.

10 **Procedure.** Infants were seated on their caregivers’ laps approximately 1.5m away from
11 the screen. Caregivers were instructed to keep their eyes closed and to refrain from interacting
12 with their infants throughout the experiment, and were monitored for compliance.

13 After calibrating infants to the screen using a toy, the researcher began the experiment.
14 The researcher had access to a video feed of the infants’ faces, a computer screen indicating the
15 current trial, and a third screen indicating when to conclude a trial. The researcher ran the
16 experiment and coded looking time online while unaware of the order of events (and therefore
17 unaware of the infant’s differential reactions to the displays), but could determine the start of
18 each trial as well as the timing of actions (e.g. when the central agent approached one of the two
19 targets) based on auditory cues.

20 Across both the familiarization and test phases of the experiment, the researcher began
21 coding a trial immediately following the attention getter, and concluded the trial once the infant
22 had attended to the screen for 60s cumulatively or looked away for 2s consecutively. During pre-
23 test, the researcher waited until the infant looked towards each target agent at least once, and
24 then began the test trials. These criteria were fixed prior to the start of data collection.

25 **Coding and analysis.** Videos of all test sessions were coded offline by observers who
26 were unaware of the order of events that infants viewed, using the same thresholds as online
27 coding, and reviewed for predetermined exclusion criteria (fussiness that prevented study
28 completion, online coding error, experimenter error, technical failure, and parental interference).
29 Further, if infants were determined to have missed a critical part of the test trial (i.e., looked
30 away for the entirety of the approach at test), then that test pair was marked and excluded from
31 subsequent analyses. If infants missed a critical portion of both test pairs, then they were dropped
32 from the sample and replaced. To assess the reliability of the offline-coded data, 100% of the test
33 trials were recoded independently by an additional researcher who was unaware of test pair
34 order. The two coders agreed on trial cutoffs for 95% of the test trials, and the intraclass
35 correlation (ICC) between the two raters was 0.994, 95% CI [0.991, 0.996]. Thus, the primary
36 offline coding data were used in our analyses.

37 The primary dependent measure was log-transformed looking time (32) but plots and
38 descriptive statistics feature raw values for ease of interpretation.

39 All models were fit in R (41). Linear mixed models were fit using the lme4 package (42).
40 Detection of influential observations was conducted using the influence.ME package (43). Plots
41 were produced using the ggplot2 package (44). To explicitly take into account repeated
42 measures, all mixed models included participant identity as a random intercept. Three classes of
43 models were fit: (1) null models, featuring participant identity as the only predictor, (2)
44 hypothesis-driven models, which included additional manipulated factor(s), and (3) exploratory
45 models, which included additional non-hypothesis driven factors. We leveraged likelihood ratio
46 tests (LRTs) and the Aikake Information Criterion (AIC) to evaluate model fit and parsimony..

1 All degrees of freedom from mixed effects models were calculated using the Satterthwaite
2 approximation method. Bracketed values indicate 95% confidence intervals.

3 We predicted that if infants can infer value from effort and use this information to predict
4 the agents' subsequent actions, then they will differentiate between the more probable outcome
5 (when the agent chooses the higher value goal) and the less probable one (when the agent
6 chooses the lower value goal) by showing a looking preference in either direction.

7 **Results**

8 **Hypothesis-driven Results.** A model with the single predictor of test trial (higher- or
9 lower value) revealed that infants looked longer at the lower value action ($M=28.41s$, $SD=14.85$)
10 than the higher value action ($M=21.79s$, $SD=12.29$), $B=0.327$, $SE=0.130$, $\beta=0.502$, $t(24)=2.523$,
11 $p=.019$, $[0.062, 0.591]$. This model outperformed a null model by a LRT, $X^2(1)=5.648$, $p=.017$. A
12 leverage analysis using Cook's Distance revealed 1 influential observation in this model.
13 Removal of this case produced an inferentially equivalent result, $B=0.263$, $SE=0.119$, $\beta=0.454$,
14 $t(23)=2.221$, $p=.037$, $[0.021, 0.506]$.

15 **Exploratory Results.** A first exploratory analysis tested for an effect of test pair order by
16 including an interactive effect between test trial presentation order and trial type. Infants
17 discriminated between the test events to a similar degree regardless of whether they were
18 assigned to watch the lower- or higher value event first, $B=0.212$, $SE=0.255$, $\beta=0.325$,
19 $t(24)=0.829$, $p=0.415$, $[-0.310, 0.733]$. Removing one influential case produced an inferentially
20 equivalent result, $B=0.354$, $SE=0.226$, $\beta=0.610$, $t(23)=1.569$, $p=0.130$, $[-0.107, 0.816]$. A second
21 model was fit with the additive effects of test pair order and test trial type, summed looking time
22 during familiarization, sex, the identity of the higher value character, and the order of the first
23 block of familiarization. Infants' looking preferences were not predicted by any of the
24 exploratory factors, with all CIs containing 0, $ps>0.1$. Removal of one influential case produced
25 inferentially equivalent results. The best model out of the above 4 was the hypothesis-driven
26 model ($AIC=92.500$).

27 **Experiment 2**

28 **Methods**

29 **Participants.** Our final sample included 24 healthy, full-term infants (15 female,
30 $M_{age}=9.88m$, $range=9.47-10.43$). Ten more infants were tested, excluded, and replaced (1 for
31 fussiness that prevented study completion, 1 for technical failure, 2 for online coding errors, 2 for
32 parental interference, and 4 for inattentiveness during test events).

33 **Materials, design, and procedure.** All materials, design, and procedure were identical to
34 those from Exp. 1 except that during familiarization (Movie S2), each target appeared at the top
35 of a ramp and the agent either refused or accepted to climb the ramp to reach it. To manipulate
36 action cost while controlling for path length, the angle of the ramps varied (11.51° , 39.26° , and
37 64.09°) such that the agent began 10 Blender units away from the top of the ramp. When the
38 agent accepted the cost (6.8s), the agent moved up the ramp once, slid back down, and then
39 moved all the way up to the target. When the agent refused the cost (5.5s), the agent moved up
40 the ramp once, slid back down, and then turned away from the target, back towards the center for
41 the screen.

42 **Coding and analysis.** All coding and analysis procedures were identical to those from
43 Exp. 1. To assess the reliability of the offline-coded data, 100% of the test events were recoded
44 independently by an additional researcher who was unaware of test pair order. The two coders
45 agreed on trial cutoffs for 98% of the test trials, and the intraclass correlation (ICC) between the
46 two raters was 0.978, 95% CI $[0.967, 0.985]$. Thus, the primary offline coding data were used in

our analyses. Because of our strong directional prediction, all reported p-values in hypothesis-driven results in Exp. 2 are one-tailed. All other reported p-values are two-tailed.

Results

Hypothesis-driven Results. As in Exp. 1, infants looked longer at the lower value action ($M=30.84$, $SD=13.79$) than the higher value action ($M=27.05s$, $SD=17.55$), $B=0.250$, $SE=0.109$, $\beta=0.408$, $t(24)=2.294$, $p=.015$, $[0.028, 0.472]$. This model outperformed a null model by LRT, $X^2(1)=4.760$, $p=.029$. No influential cases were detected.

Exploratory Results. An exploratory model testing explicitly for presentation order revealed that infants differentiated between the test events differently depending on whether they were assigned to watch the lower value versus the higher value approach first, $B=0.521$, $SE=0.190$, $\beta=0.851$, $t(24)=2.741$, $p=.011$, $[0.133, 0.909]$. No influential cases were detected. which revealed that whereas infants who saw the lower value test event first looked longer at the lower value ($M=28.70s$, $SD=11.12$) versus higher value ($M=20.01s$, $SD=14.83$) test trials, $B=0.511$, $SE=0.134$, $t(24)=3.797$, $p=.001$, $[0.233, 0.788]$, infants who saw the higher value choice did not differentiate between the lower- ($M=32.97s$, $SD=15.39$) and higher value test events ($M=34.08$, $SD=17.78$), $B=-0.011$, $SE=0.134$, $t(24)=0.079$, $p=.938$, $[-0.288, 0.267]$. An additional model testing for effects of summed attention during test, sex, the identity of the higher value target, and the first familiarization loop revealed that no further effects other than one of first familiarization loop, where infants assigned to watch a sequence of low to high cost first looked longer overall at test, $B=0.440$, $SE=0.185$, $\beta=0.851$, $t(24)=2.374$, $p=.023$, $[0.062, 0.819]$. Removal of 3 influential observations from this model yielded inferentially equivalent results. The best model of the four reported was the simpler exploratory model with the single interactive effect ($AIC=78.219$).

Experiment 3

Methods

Participants. Our final sample included 32 healthy, full-term infants (15 female, $M_{age}=10.03m$, range=9.57-10.50). Six more infants were tested, excluded, and replaced (1 for fussiness that prevented study completion, 3 for online coding errors, 1 for parental interference, and 1 for inattentiveness during test events). Sample size was determined from a simulation power analysis over data from Exp 1-2. The design, procedure, and analyses of this experiment were pre-registered via the Open Science Framework (<https://osf.io/k7yjt/>).

Materials, design, and procedure. All materials, design, and procedure were identical to those from Exp. 1-2 except as follows. During familiarization (Movie S3), each target appeared at the far end of a trench (15 units deep), and the agent either accepted this cost by jumping across it or refused to jump. To manipulate action cost while controlling for movement against gravity, the width of the trench varied (5, 10, and 15 units) such that the agent began 13, 18, or 23 units away from the target. When the agent accepted the cost (8.4s for the small cost, 9.3s for the medium), the agent looked down at the bottom of the trench, looked at its target, and backed up and leapt across the trench. When the agent refused the cost (7.2s), it looked down at the bottom of the trench, looked at its target, and backed up the turned away, back towards the center for the screen. The agent backed away a longer distance and accelerated more when accepting the medium versus the small cost. During test (Movie S5), infants saw the same choice events as in Exp. 1-2, but the events took place on a platform at the same height used for familiarization. Prior to familiarization, infants watched an additional video (Movie S6), in which a ball rolled across the platform and off the small, medium, and wide trench (5.8s, 5.7s, and 5.0s), identical in width and depth to those during familiarization and situated at the center of the screen. During

1 each segment of the video, the ball traveled in a parabolic trajectory after it rolled off the edge
 2 and shattered against the far wall of the trench. The location of impact depended on the width of
 3 the trench, such that the ball shattered at a lower point on the far wall for wider trenches. The
 4 side of the screen that the ball emerged from (left vs. right) was consistent across all videos
 5 within participants and was counterbalanced across participants.

6 **Coding and analysis.** All coding and analysis procedures were identical to those from
 7 Exp. 1-2. To assess the reliability of the offline-coded data, 100% of the test events were recoded
 8 independently by an additional researcher who was unaware of test pair order. The two coders
 9 agreed on trial cutoffs for 96% of the test trials, and the intraclass correlation (ICC) between the
 10 two raters was 0.995, 95% CI [0.993, 0.996]. Thus, the primary offline coding data were used in
 11 our analyses. All reported *p*-values in hypothesis-driven results in Exp. 3 are one-tailed. All other
 12 reported *p*-values are two-tailed.

13 Results

14 **Hypothesis-driven Results.** As in Exp. 1-2, infants looked longer at the lower value
 15 action ($M=23.05s$, $SD=13.58$) than the higher value action ($M=17.47$, $SD=10.69$), $B=0.260$,
 16 $SE=0.119$, $\beta=0.403$, $t(32)=2.185$, $p=.018$, [0.020, 0.501]. This model outperformed a null model
 17 by LRT, $\chi^2(1)=4.452$, $p=.035$. Removal of 1 influential case produced an inferentially equivalent
 18 result, $B=0.286$, $SE=0.120$, $\beta=0.466$, $t(31)=2.379$, $p=.011$, [0.043, 0.529].

19 **Exploratory Results.** An exploratory model testing explicitly for presentation order
 20 revealed that like in Exp. 2, infants differentiated between the test events differently depending
 21 on whether they were assigned to watch the lower value versus the higher value approach first,
 22 $B=0.770$, $SE=0.195$, $\beta=1.194$, $t(32)=3.941$, $p<.001$, [0.357, 1.165]. We detected 1 influential
 23 observation in this model and removed it from subsequent pairwise comparisons, which revealed
 24 that whereas infants who saw the lower value test event first looked longer at the lower value
 25 ($M=27.95s$, $SD=14.11$) versus higher value ($M=15.05s$, $SD=9.06$) test trials, $B=0.645$, $SE=0.117$,
 26 $t(31)=5.519$, $p<.001$, [0.407, 0.883], infants who saw the higher value choice showed a weak
 27 preference for the higher-value ($M=20.62$, $SD=11.93$) relative to the lower-value ($M=16.56s$,
 28 $SD=9.84$) test events, $B=-0.236$, $SE=0.121$, $t(31)=-1.959$, $p=0.059$, [-0.483, 0.010]. An additional
 29 model testing for effects of summed attention during test, sex, the identity of the higher value
 30 target, and the first familiarization loop revealed that no further effects other than one of total
 31 attention in seconds during familiarization, where infants who were more attentive during
 32 familiarization looked marginally longer overall at test, $B=0.002$, $SE=0.001$, $\beta=0.261$,
 33 $t(32)=1.948$, $p=.060$, [0.000, 0.005]. Removal of two influential cases from this model yielded an
 34 inferentially equivalent result, and revealed an additional finding that infants randomly assigned
 35 to conditions where the higher value target was on the right looked longer overall at test,
 36 $B=0.368$, $SE=0.149$, $\beta=0.610$, $t(30)=2.464$, $p=.020$, [0.066, 0.670] The best model of the four
 37 reported was the simpler exploratory model with the single interactive effect ($AIC=114.24$).

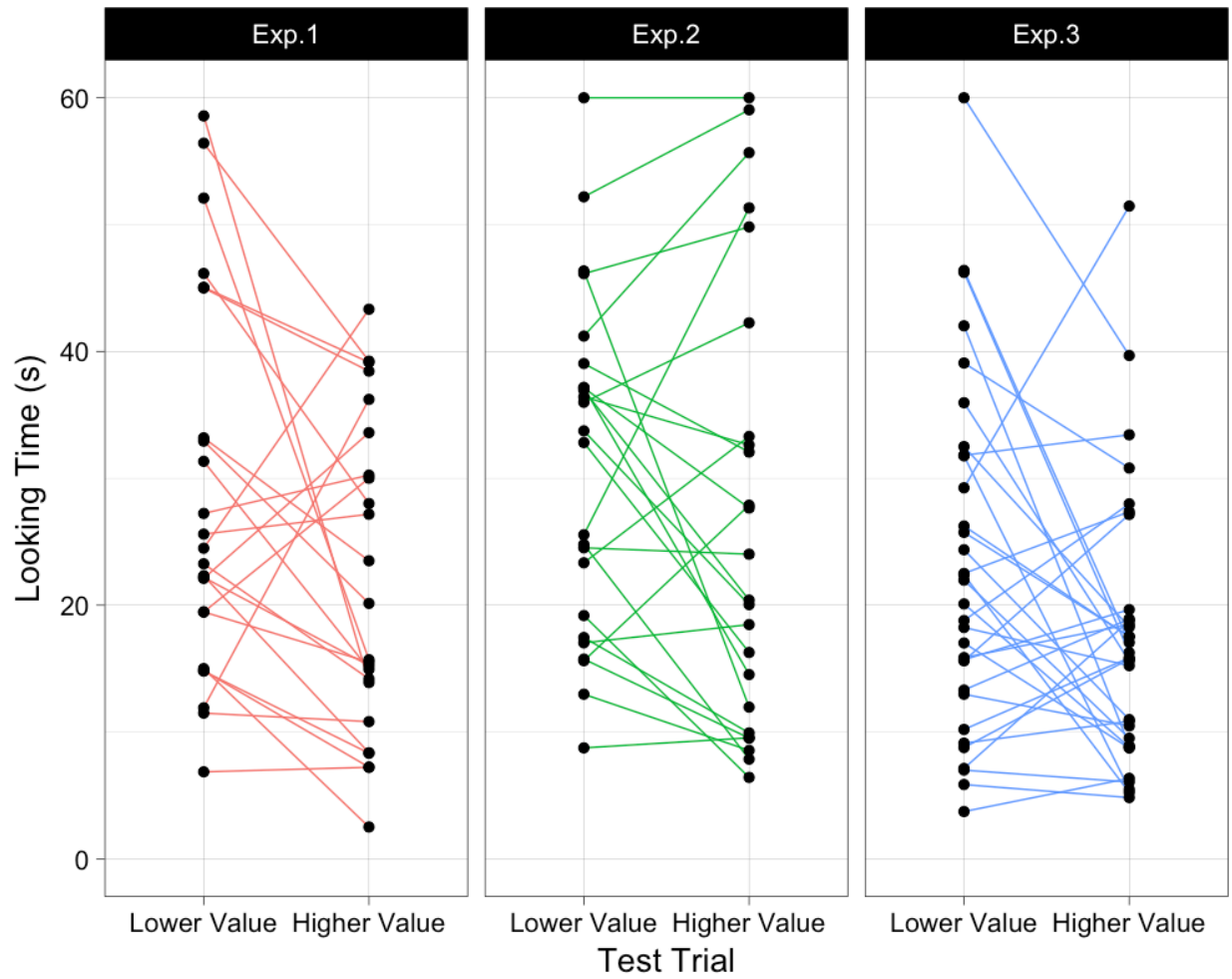
38 **Results Across Exp 1-3.** As reported in the main text, across all experiments, infants
 39 looked longer at the lower value action ($M=26.99s$, $SD=14.13$) than the higher value action
 40 ($M=21.64s$, $SD=13.94$), $B=0.277$, $SE=0.070$, $\beta=0.424$, $t(80)=3.975$, $p<.001$, one-tailed, [0.139,
 41 0.415], mixed effects model with random intercepts for participant and experiment. Removal of
 42 one influential case produced an inferentially equivalent result, $B=0.258$, $SE=0.068$, $\beta=0.406$,
 43 $t(79)=3.799$, $p<.001$, one-tailed, [0.123, 0.393]. To test explicitly for differences in responses to
 44 test events across experiments, an additional model with an interactive effect between
 45 experiment and test event was fit and revealed no differences in looking preference across

1 experiments, and removal of one influential observation produced an inferentially equivalent
2 result.

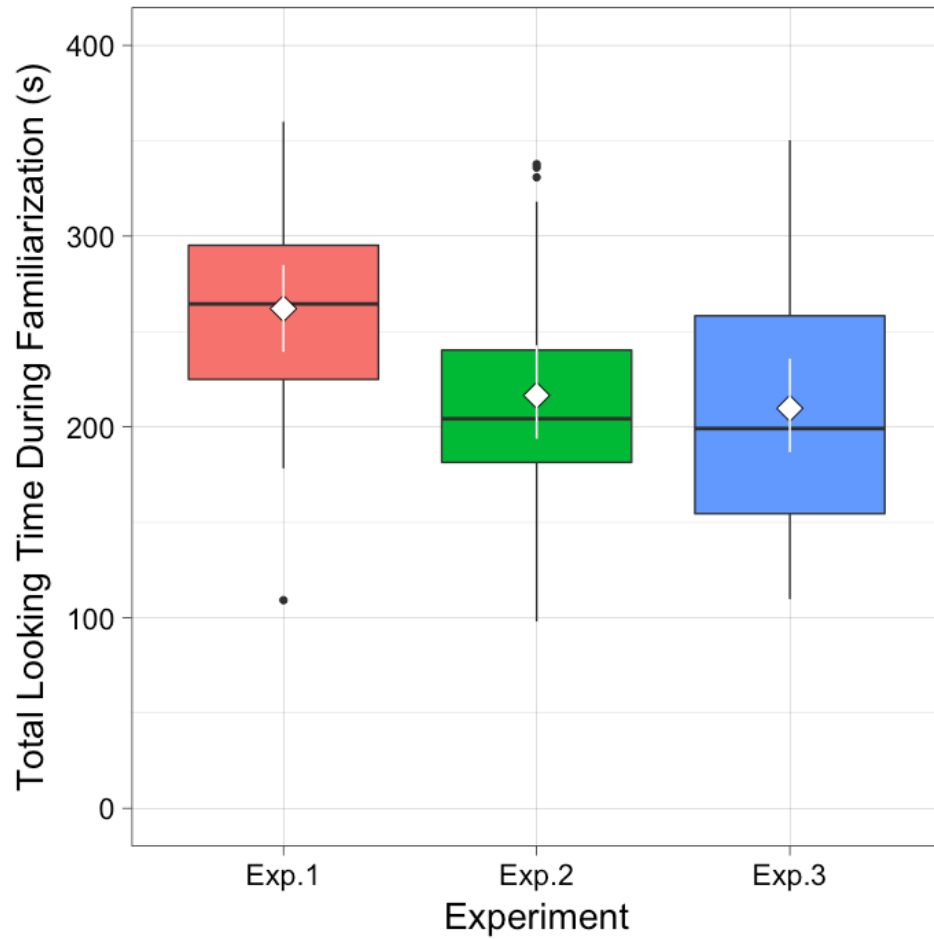
3 Two further models tested explicitly for the effect of presentation order on attention to
4 the lower- vs. higher value test events, and additional predictors of summed attention during test,
5 sex, the identity of the higher value target, and the first familiarization loop. These models
6 revealed a robust effect of presentation order, $B=0.528$, $SE=0.126$, $\beta=0.808$, $t(80)=4.179$, $p<.001$,
7 $[0.277, 0.778]$, and for an effect of attention during familiarization, $B=0.002$, $SE=0.001$, $\beta=0.209$,
8 $t(79.65)=2.221$, $p=.030$, $[0.000, 0.006]$. mixed effects model with random intercepts for
9 participant and experiment. Removal of influential cases produced inferentially equivalent
10 results. A comparison of model fit and parsimony revealed that the model including the single
11 interactive effect provided the best description of the data (AIC=280.90).

12 Fifty out of the 79 infants who demonstrated a preference in either direction (63.29%; 1
13 infant in Experiment 2 looked for 60 seconds on all test trials) looked longer on average at the
14 lower value test events, $p=.033$, 95% CI $[0.269, 0.490]$, exact binomial test. The proportion of
15 infants who looked in the predicted direction did not differ across Experiments 1 (17/24,
16 70.83%), Experiment 2 (14/23, 60.87%), and Experiment 16 (19/32, 59.38%), $X^2(2)=1.022$,
17 $p=.600$. See Figure S3. Further non-parametric analyses on raw looking times in seconds
18 supported the finding that infants looked longer at the lower value test events across all
19 experiments, $[2.165, 8.185]$, $V=2241$, $p=.001$, Wilcoxon signed rank test, and $[0.720, 8.235]$,
20 bootstrapped median difference in looking times across test events with 10,000 samples.

21 To compare attention during familiarization across the experiments, we fit a linear model
22 including summed looking time in seconds as the dependent variable and experiment (1, 2, or 3)
23 as a predictor. We found that infants looked longer during familiarization in Exp. 1 ($M=261.96s$,
24 $SD=61.29$) than in Exp. 2 ($M=216.50$, $SD=63.49$), $B=45.47$, $SE=19.22$, $\beta=0.653$, $t(77)=2.366$,
25 $p=.021$, $[7.196, 83.74]$, and in Exp. 3 ($M=209.70$, $SD=72.33$), $B=52.26$, $SE=17.98$, $\beta=0.751$,
26 $t(77)=2.907$, $p=.005$, $[16.460, 88.056]$. See Figure S2.



1
 2 **Fig S1.** Looking times averaged across two test pairs towards the lower value and higher value
 3 test event from all participants in Experiments 1-3.



1
2
3
4
5
6

Fig S2. Total looking time in seconds during familiarization across Exp 1-3. Boxes indicate middle quartiles, vertical lines indicate points within 1.5 times the interquartile range from the 25th and 75th percentiles, and horizontal lines indicate medians. Means and 95% confidence intervals are plotted in white.

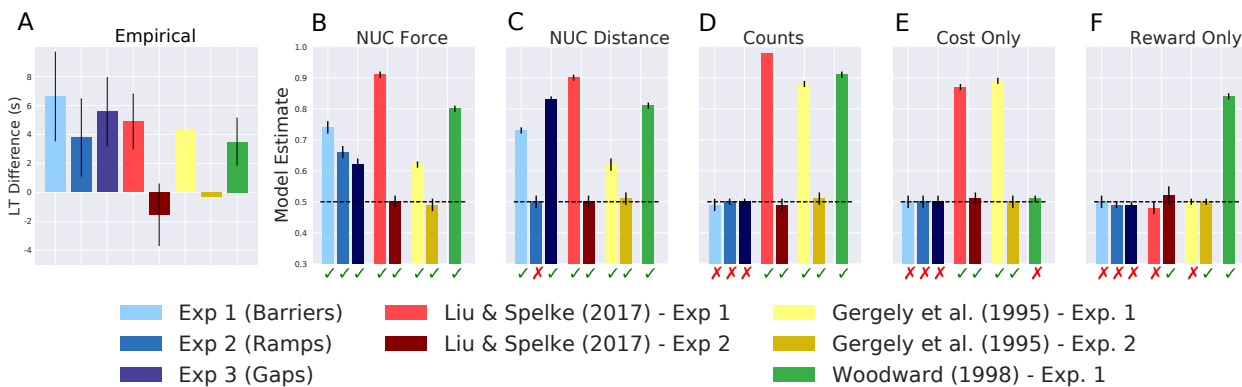
- 1 **Movie S1.** Sample familiarization event for Experiment 1.
- 2 **Movie S2.** Sample familiarization event for Experiment 2.
- 3 **Movie S3.** Sample familiarization event for Experiment 3.
- 4 **Movie S4.** Sample test event for Experiments 1-2.
- 5 **Movie S5.** Sample test event for Experiment 3.
- 6 **Movie S6.** Pre-familiarization event for Experiment 3.
- 7

1 **Computational Modeling Details**

2
3 We constructed several alternative models for predicting the actions of a central red
4 agent, given previous stimuli showing its behavior. We evaluated the models both as accounts of
5 our present three experiments testing infants' ability to integrate costs and rewards in action
6 understanding, and as accounts of five previous experiments that tested infants' abilities to make
7 inferences about rewards (6) or costs (7, 13) individually. As our main proposed model, we
8 constructed a probabilistic program for inferring rewards over possible goals through a Bayesian
9 computation, based on the Naïve Utility Calculus (NUC) and inverse planning. The cost-function
10 of this model was either based on distance (NUC-distance), or physical effort through the
11 application of force (NUC-force). We next constructed three alternative models that were lesions
12 of various aspects of this model. Only one of our two NUC models (NUC Force; Fig. S3B)
13 accounts for all of these findings, including Experiment 2 in which infants saw the agent cover
14 the same total distance at the same speed but on trajectories requiring different amounts of
15 physical effort due to the force of gravity. In contrast, two models that do not integrate costs and
16 rewards (Costs Only, Fig. S3E and Rewards Only, Fig. S3F), as well as two integrative but
17 simpler models using only perceptual cues for value (Counts, Fig. S3D) or distance-based
18 representations of cost (NUC Distance, Fig. S3C) instead of abstract effort-reward tradeoffs,
19 account for just a subset of this body of empirical findings.

20 The eight experiments shown in Fig. S3 only begin to test the explanatory scope of our
21 computational framework. This approach naturally extends to explain other findings (23),
22 including infants' rational imitation of other agents' goal-directed actions (45, 46), understanding
23 of the transitivity of agents' preferences (47), inferences of agents' preferences from their non-
24 random selections of objects (48), and inferences about the existence and location of an occluded
25 object from the cost of an agent's action (49).

26 In the rest of this section we describe the full NUC model, then the alternative proposals.
27 We then compare the predictions of the different models to the empirical results of eight
28 different experiments with young children that all relate to reasoning about cost, reward and
29 efficient behavior.
30
31
32



1
2

3 **Fig S3.** Comparison of model predictions (B-F) to the empirical results (A) of Experiments 1-3
 4 in this paper, as well as additional experiments on infants' sensitivity to costs and rewards (6, 7,
 5 13). Model estimates were generated from 5,000 MH samples with a 10-step interval and burn-in
 6 of 1,000. Error bars indicate 95% CIs (model estimates) and standard errors (where available
 7 from empirical findings). Check marks and crosses respectively indicate consistency or
 8 inconsistency with empirical findings.

9

Decision-making framework

Following the logic of “Bayesian Theory of Mind” (4, 23, 50, 51), the primary model assumes agents have a planning procedure for generating actions that maximize their utility or expected utility, and that agents then use observed actions to invert this procedure to reason about the utilities and constraints of the planning agents.

Following a standard decision-making framework for rational planning (e.g. 48) we assume that a rational planning agent:

- (i) Divides the world into possible states S , such that for each state $s \in S$ there exists a set of possible actions A_s .
- (ii) Has access to a transition function T such that:

$$T(S, A_s) \rightarrow P(S'), \quad (1)$$

where $P(S')$ is the probability distribution over the states of the worlds S' that can result from taking action A_s in state S .

- (iii) Has a utility function U such that:

$$U(A, S) = \text{Reward}(S) - \text{Cost}(A), \quad (2)$$

where *Reward* and *Cost* are functions that map from states and actions to real numbers. This separation of the utility into independent components follows recent work by Jara-Ettinger et al. (2016), who showed that both components can be separate targets of inference, and used as explanatory variables by young children across different scenarios.

- (iv) Is guided by a decision-function $D(S, A_s) \rightarrow P(A)$ that selects a given action in a state, in order to maximize the expected utility U . For our simple scenarios, we assume a soft-max decision function (53):

$$P(A_s = a | U, S = s) = \frac{e^{\beta U(a, T(s, a))}}{\sum_j e^{\beta U(a, T(s, a))}}, \quad (3)$$

where β is a noise-parameter (inverse-temperature) adjusting the agent’s determinism. As $\beta \rightarrow 0$, the agent will behave in a more random fashion. As $\beta \rightarrow \infty$, the agent will tend to use greedy action selection. Intermediate values of β (around 1) approximate probability matching behavior. Thus, we place an equal prior probability on β taking one of two values, $\beta = 3$ (representing a near-optimal rational agent) and $\beta = 0$ (representing a random, non-rational agent). The exact value of the parameter for a rational agent is not important for the qualitative pattern of results discussed in the following sections.

The Specific Decision-making Environment

States: In our studies infants saw the agent making a single choice between staying put and moving towards one of two possible goal agents. The full set of possible states in Experiments 1-3 are $s \in \{Start, Target_A, Target_B\}$, where $Target_i$ can be the Square or Triangle goal agent. Not all states are available in every situation, depending on the stimuli (for example, if the stimuli shows the decision making agent near a ramp with the Square goal agent at the top, the sub-set of states for this environment does not include reaching the Triangle goal agent).

Actions: In Experiment 1 the actions in any situation were a subset of $\{Nothing, Left, Right, Jump\ tall, Jump\ medium, Jump\ short\}$.

1 In Experiment 2 the actions were a subset of
2 $\{\text{Nothing}, \text{Left}, \text{Right}, \text{Climb steep}, \text{Climb moderate}, \text{Climb shallow}\}$.

3 In Experiment 3 the actions were a subset of
4 $\{\text{Nothing}, \text{Left}, \text{Right}, \text{Jump wide}, \text{Jump medium}, \text{Jump narrow}\}$.

5 The set of possible actions and states was similarly altered for the other 5 experiments we
6 considered, that were not part of the current study.

7 **Cost functions:** We do not know the specific cost that infants believe the agent incurs for
8 jumping over the barriers or climbing the inclines, but we assume that the cost for
9 jumping/climbing is greater than staying put, and that cost generally scales with distance. There
10 are at least two potential models for how to take this distance into account when calculating cost.
11 The first model (NUC-distance) ignores physical forces and effort, and uses only total distance
12 traveled between states as the input to the cost function:

$$13 \quad \text{Distance}_{\text{cost}}(S_1 \rightarrow S_2) \propto \int_C a \, ds, \quad (4)$$

14 where C is the trajectory from S_1 to S_2 , and a is some constant factor. We considered each
15 Blender unit of distance to be one unit of cost for the model.

16 The second model (Force) considers the physical work required to get from one state to
17 the next. Specifically, we consider the work done in the presence of a conservative force field
18 (gravity):

$$19 \quad \text{Work}_{\text{cost}}(S_1 \rightarrow S_2) = \int_C F \cdot ds = \int_{t_1}^{t_2} F \cdot dt. \quad (5)$$

20 For Experiments 1&2, we considered the force of over-coming friction as well as gravity,
21 meaning $F = f \cdot g \cdot \cos(\alpha) + g \cdot \sin(\alpha)$, where α is the angle opposite the vertical part of the
22 incline, and f is the friction coefficient. The mass (m) was set arbitrarily to 1 and so does not
23 appear. Without the friction component the non-accelerated movement of an agent on a non-
24 inclined plane is effortless. We set the friction coefficient arbitrarily to $f = 0.5$, similar to wood
25 sliding on wood or metal sliding on wood, though we note that the exact choice of f in the range
26 0 to 1 does not affect the qualitative results. For experiment 3 we considered in addition the
27 impulse forces that accelerate and decelerate the agent. This instantaneous change in velocity can
28 be related to the work done as a change in kinetic energy, meaning $W = \Delta E_k = \frac{1}{2}(V_1^2 - V_2^2)$,
29 where V_1 and V_2 are the velocities before and after the application of the impulse. Both models
30 would predict similar cost differences for Experiments 1 and 3, but they would diverge with
31 regards to Experiment 2. $\text{Distance}_{\text{cost}}$ of the steep and shallow inclines is the same, while $\text{Work}_{\text{cost}}$
32 is sensitive to the incline of the ramp.

33 **Rewards:** We do not make a-priori assumptions about the reward associated with
34 reaching the goal agents, except that they can potentially exceed the range of the costs:

$$35 \quad \text{Reward}(\text{Target A}) \sim \text{Uniform}(0, \text{Reward}_{\text{max}}),$$

$$36 \quad \text{Reward}(\text{Target B}) \sim \text{Uniform}(0, \text{Reward}_{\text{max}}).$$

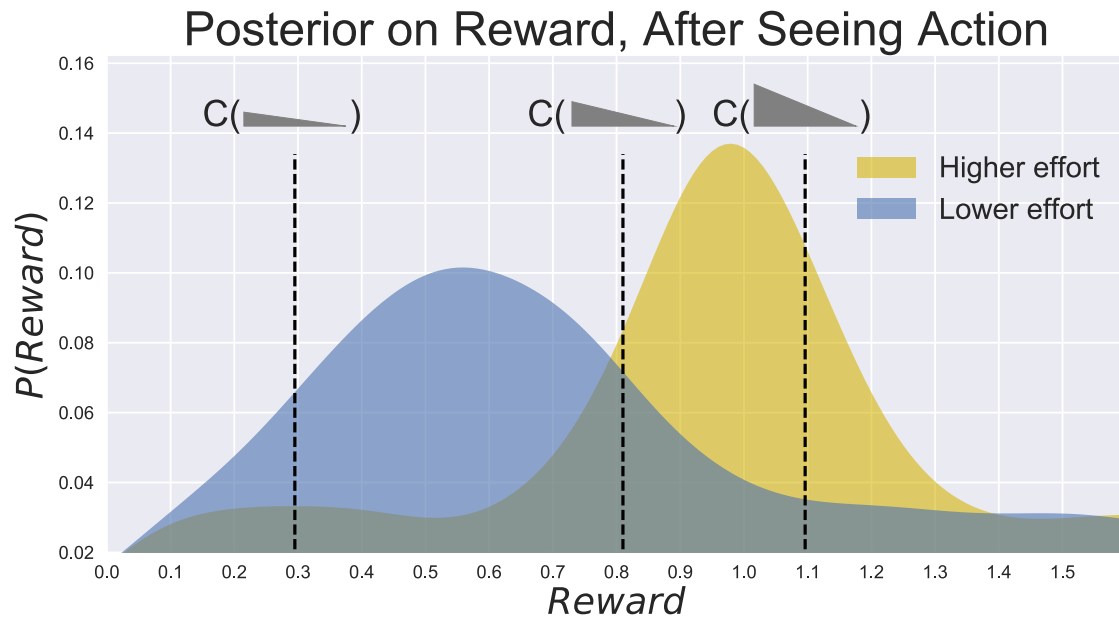
38 Inference

39
40 Using the previous assumptions, we can calculate the posterior distribution over the
41 rewards and costs conditioned on an observed action A. Applying Bayesian reasoning, this is a
42 combination of the likelihood of the observed action (given by the planning procedure), and the
43 prior distributions over costs and rewards:

$$44 \quad P(\text{Reward}, \text{Cost} \mid A, S) \propto P(A \mid S, \text{Cost}, \text{Reward}) P(\text{Reward}, \text{Cost}) \quad (6)$$

1 We constructed a probabilistic program that samples rewards and costs from this
2 posterior, conditioned on observed actions matching the actions of the agent . The program was
3 written in Church (54), a probabilistic programming language based on Scheme.

4 This program can be used to infer the value of the reward for the different experiments
5 (cases 1-8 as detailed above). Below we consider the specific case of Experiment 2 from the
6 current paper (Ramps), using the NUC-force model. The sampling procedure in Church uses the
7 Metropolis-Hastings algorithm, and Fig. S4 shows the resulting approximation to the posterior in
8 Eq. 6 for Experiment 2 and the NUC-force model, with 5,000 samples at 10-step intervals with a
9 burn-in of 1,000 samples. As can be seen in Fig. S4, the model shifts probability away from a
10 uniform distribution for both goals, such that the reward for the higher-effort target is higher in
11 expectation than the lower-effort one. For both the targets, most of the probability mass is in
12 between the accepted and rejected costs, as expected.
13



1
2
3
4
5
6
7

Fig S4. The model's inferred posterior probability distribution over possible reward values, after seeing the actions taken in Experiment 2 when using a Force-based cost function. Note that for the force-based model the cost of the different ramps is a combination of the effort of working against gravity and overcoming friction.

1 Action Prediction

2
3 The ‘forward planning’ part of the program implements Eq. 1-6 and the decision-making
4 framework of the first section. This forward planning can use the posterior distribution inferred
5 from observing the agent’s action to predict its next actions, $P(\text{Action} \mid \text{Previous stimuli})$. This
6 probability can be calculated for cases in which there is only one reward and multiple possible
7 actions (as in Gergely et al. 1995, Experiments 1&2) or multiple rewards (the other experiments
8 considered). In all model comparisons we used 5,000 samples at 10-step intervals with a burn-in
9 of 1,000 samples for the inference of the reward distributions, and for the sampling of action
10 predictions. We discounted the ‘Do Nothing’ action predictions, as the infants never saw these in
11 the test stimuli, and were faced rather with a direct comparison between two actions.

12 To give an example of a particular way in which action prediction proceeds, we consider
13 again Experiment 2 from the current paper under the NUC-force model. For the test stimuli, the
14 situation includes no barriers, and both targets are obtainable. The possible actions are: going to
15 target A, going to target B, or doing nothing. We assume that:

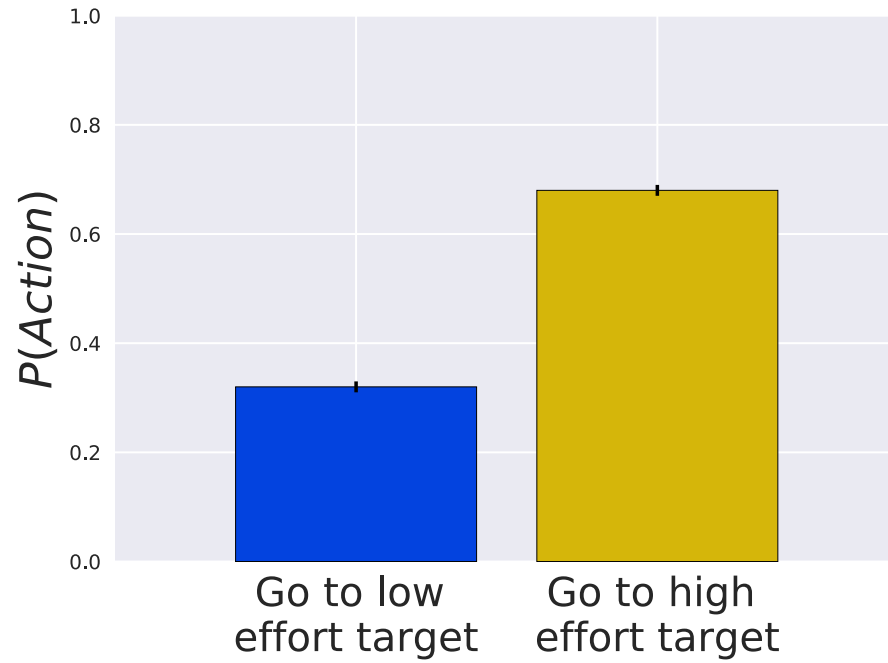
$$16 \quad \text{Reward}(\text{Target A}) \sim P(\text{Reward}(\text{Target A}) \mid \text{Previous Stimuli}), (7)$$

$$17 \quad \text{Reward}(\text{Target B}) \sim P(\text{Reward}(\text{Target B}) \mid \text{Previous Stimuli}), (8)$$

18 where $P(\text{Reward}(\text{Target}) \mid \text{Previous Stimuli})$ is calculated using the program approximating Eq.
19 6.

20 Fig. S5 shows the resulting prediction for the probability distribution over the agent’s
21 action. As can be seen in Fig. S5, the model predicts the agent will move towards the Target that
22 it expended higher effort to reach earlier. Consequently, when the agent moves towards the lower
23 effort Target, this goes against the prediction.

24 The amount of surprise the model predicts can be directly related to the inverse of the
25 probability of a given event or action (34).
26



1
2 **Fig S5.** Probability of predicted next action by the NUC-force model for Experiment 2 (Ramps),
3 for the test case in which both targets are equally accessible. Black lines show bootstrapped 95%
4 confidence intervals.

5
6

1 Alternative Models

2
3 The NUC-force and NUC-distance models outlined so far assume that infants are able to
4 use Bayesian reasoning to integrate inferences about cost, reward, and agent rationality to predict
5 the next action of an agent. The NUC-distance model assumes that cost is proportional to
6 perceptually observable distance, while the NUC-force model assumes that cost is proportional
7 to physical effort.

8 But an alternative option is that infants are performing a much more low-level analysis of
9 the scene, or reasoning about costs and rewards in a non-integrated way. Consider for example
10 the findings from Experiment 1 in Woodward (1998), showing that 9-month old infants expected
11 an agent to reach for a goal object A over B, after A and B had switched positions compared to
12 habituation. A full NUC mental model could be deployed to infer that the reward for A is higher
13 than the reward for B, and then correctly predict the action. But it is also possible that infants are
14 reasoning along the lines of “A goal previously reached for will be reached for again, regardless
15 of its position”.

16 We thus consider 3 alternatives to the NUC model, all of which involve different forms of
17 lower level or non-integrated reasoning, as follows:

- 18 1) **Count-based:** This is a perceptual-based account that relies on directly observable cues:
19 reaching an item, and the distance to the item. This model is similar to the NUC in that
20 agents are assumed to be rational planners that act to achieve goals under constraints.
21 However, rather than using full Bayesian inverse planning to reason about rewards, it
22 uses an easy-to-calculate proxy for the reward of a target: The model tallies the number
23 of times a goal object has been approached or reached, and considers that tally in
24 proportion to the total number of times any goals have been reached. Furthermore, the
25 model considers the cost as proportional to the distance. Thus, for this model:
26

$$27 \quad P(\text{Reward}_i | \text{Stimuli}) = \frac{\#(\text{Reward}_i \text{ approached})}{\sum_j \#(\text{Reward}_j \text{ approached})} \cdot (9)$$

28
29 The cost is calculated as in Eq. 4, and the planning proceeds as in Eq. 1-3.

- 30
31 2) **Cost Only:** This model assumes all goal states are equally rewarding, but is sensitive to
32 the cost in the form of distance as in Eq. 4. Such a model is able to reason about agents
33 acting efficiently to reach goals, in the sense of minimizing distance as in Gergely et al.
34 (1995), for example.
- 35 3) **Reward Only:** This model assumes rewards can vary and performs the correct Bayesian
36 updating on the probability distribution over reward as in Eq. 6, but without the cost
37 factor. Such a model is sufficient for correctly predicting behavior when reasoning about
38 the choices of agent that shows a simple preference for one target over another, as in
39 Woodward (1998), for example.

41 Comparison to Empirical Data from Eight Experiments

42

1 In this section we compare the outputted prediction of the NUC models as well as the
2 alternative models, to eight different experiments examining infants' expectations about cost,
3 reward and efficient behavior. The experiments considered were:

- 4 1) Experiment 1 in the current paper (Barriers).
- 5 2) Experiment 2 in the current paper (Ramps).
- 6 3) Experiment 3 in the current paper (Gaps).
- 7 4) Experiment 1 in Liu and Spelke (2017), in which six-month-old infants first view an
8 agent jumping over barriers of varying sizes to get to a goal. In the test phase, the
9 agent is blocked by a previously unseen small barrier and either makes a small jump
10 (expected) or a large jump (unexpected) over the barrier to reach the goal.
- 11 5) Experiment 2 in Liu and Spelke (2017), which is identical to (3), except that the
12 barriers are placed behind the goal such that they are not blocking the path of the
13 agent. Infants' looking time was at chance when shown a small or large jump,
14 presumably inferring that the agent is irrational based on previous behavior.
- 15 6) Experiment 1 in Gergely et al. (1995), in which 12-month-olds first view an agent
16 jumping over a barrier to reach a goal. In the test phase, the barrier is removed and
17 the agent either moves straight to the goal (expected) or performs the previously seen
18 jump on the way to the goal (unexpected).
- 19 7) Experiment 2 in Gergely et al. (1995), which is identical to (5) except the barrier is
20 placed behind the goal such that it is not blocking the path of the agent. Infants'
21 looking time was at chance when shown the agent moving straight or making a jump
22 on the way to the goal.
- 23 8) Experiment 1 in Woodward (1998), in which 9-month-old infants first saw an agent
24 reaching for and grasping one of two possible goals. In the test phase, the positions of
25 the goals were switched. The agent either reached to the same goal using a new
26 trajectory (expected) or follow the same trajectory to reach a new goal (unexpected).

27
28 We use the same model set-up and parameters for all experiments, varying only the actions and
29 stimuli observed, and the different cost and reward functions used depending on the model.

30 Following (34), we relate the amount of surprise in infants' looking time to $1-P(\text{outcome})$
31 as predicted by the model. Specifically, we relate this measure to the difference in looking time
32 (measured in seconds) between the unexpected event and the expected event. The results of the
33 models and the comparison to the eight experiments are shown in Fig. S3.

34 35 36 **Code**

37
38 All the code implementing the computational models and the analysis of their output is available
39 on Open Science Framework (<https://osf.io/crx4d/>).

40
41
42
43
44