The potential for effective reasoning guides children's preference for small group discussion over crowdsourcing

## Emory Richardson & Frank C. Keil

### Yale University

Address for : Emory Richardson

correspondence Department of Psychology

Yale University 2 Hillhouse Ave

New Haven, CT, 06520-8205

Email : emory.richardson@yale.edu

Word Count : 3795 (Main Text), 1046 (Methods), 1628 (Results): 6469 Total

Version : Revised Submission

### Acknowledgments

We thank the members of the Yale Cognition & Development lab for helpful feedback. This research was supported by NSF grant DRL 1561143 awarded to Frank C. Keil.

### **Author Contributions**

ER developed the study concept; ER designed the experiment with input from FCK; ER collected and analyzed the data. ER drafted the manuscript, and FCK provided critical revisions. All authors approved the final version of the manuscript for submission.

#### **Abstract**

Communication between social learners can make a group collectively "wiser" than any individual, but conformist tendencies can also distort collective judgment. We asked whether intuitions about when communication is likely to improve or distort collective judgment could allow social learners take advantage of the benefits of communication while minimizing the risks. In three experiments (n=360), 7- to 10-year old children and adults decided whether to refer a question to a small group for discussion or "crowdsource" independent judgments from individual advisors. For problems affording the kind of 'demonstrative' reasoning that allows a group member to reliably correct even errors made by a majority, all ages preferred to consult the discussion group, even compared to a crowd ten times as large — consistent with past research suggesting that discussion groups regularly outperform even their best members for reasoning problems. In contrast, we observed a consistent developmental shift towards crowdsourcing independent judgments when reasoning by itself was insufficient to conclusively answer a question. Results suggest sophisticated intuitions about the nature of social influence and collective intelligence may guide our social learning strategies from early in development.

When is advice from multiple people more likely to clarify than confound a learner's understanding? Consider two ways one could learn from multiple people at once: by eliciting a consensus judgment from a small group discussion, or by "crowdsourcing" many independent answers. Discussion may enable groups to correct mistakes and combine insights, producing an accurate consensus answer that no individual could have found alone. However, without an objective method of evaluating solutions, discussion may drag on endlessly, be misled by charismatic leaders or groupthink, and ultimately only create an illusion of consensus around a wrong answer. In contrast, crowdsourcing may still include mistakes that discussion could have corrected, but particularly in *large* crowds of *independent* responders, the majority, plurality, and average response can all be surprisingly accurate [1, 2]. Nevertheless, shared culture and cognitive biases can create illusions of consensus even without direct communication between individuals [2-4]. The empirical advantages of discussion and various crowdsourcing strategies are well-documented. Less attention has been given to laypeople's own intuitions about the tradeoffs between them (note that our use of "intuition" does not refer to the intuitive-deliberative distinction in dual systems theory; rather, it follows the frequent use of "intuitive theories" in developmental psychology to describe the untaught assumptions about the world that help learners structure their experience; for a recent review, see Gerstenberg & Tenenbaum, 2017 [5]). Here, we investigate early developing intuitions about when group discussion or crowdsourcing is a more effective use of collective intelligence.

While debate over whether crowds can be trusted is at least as old as philosophy itself [6-7], mathematical models suggest that under certain conditions, crowds can be "wise". Given a set of options to "vote" for, majority and plurality accuracy increase to near certainty as crowd size increases [8-9], and evolutionary simulations suggest that conformist learning strategies are often more adaptive than alternatives [10,1]. Similarly, averaging crowd members' individual judgments can produce a collective estimate that is more accurate than the crowd's most accurate member [11-14, 2]. Interestingly, many species faced with the problem of learning from multiple sources at once rely on similar heuristics to evaluate collective opinion. Even in early childhood, people trust majority over minority judgment, and give more weight to stronger majorities [15]. By adulthood, people also trust pluralities, and give more weight to the judgments of larger crowds [16-18].

However, crowdsourcing heuristics share a common weakness: for crowds to be "wise", individual judgments must be independent. Social influence can compound individual error, particularly when a large proportion of the population are conformist learners [19-20]. Yet, while popular concerns about echo chambers and media bias suggest that laypeople intuitively recognize some of the risks of social influences, it remains unclear how well people compensate for them in practice. For example, while adults and even children as young as six prefer firsthand knowledge over hearsay in an eyewitness memory context [3, 21], adults are just as trusting of an economic forecast repeated by five news articles citing a single primary source as they are of the same forecast citing five different primary sources [3,22]. Similarly, while children as young as four expect randomly sampled evidence to cause others to revise their beliefs more than evidence from a biased sampling process [23], when the source of the sampling bias is the selection of informants itself, children are sometimes insensitive to bias even late in development—particularly when the degree of consensus is high [24-29]. Indeed, even as adults, people frequently mistake the frequency of a belief in their local networks for its frequency in the population as a whole [30-32]. In short, people's trust and distrust of consensus seems to selectively disregard one of the proposed preconditions of consensus' accuracy — independent sources.

One reason for people's occasional indifference to their sources' independence may be that social influence frequently makes their judgments *more* accurate [33-38]. For instance, while open discussion may risk groupthink by sacrificing individuals' independence, it also allows individuals to pool their knowledge and generate new insights; discussion can also ease the cognitive load on individuals, increase a groups' capacity for processing information, and allow the group to correct individual mistakes [39-41]. This division of cognitive labor means that discussion groups may be able to quickly generate solutions that most individuals would never produce alone, and may make discussion an attractive learning strategy for a wide variety of problems, particularly as the evidence load increases [42-45]. Most notably, to the extent that discussion enables even a single group member to correct a *majority* that has made a mistake, discussion may also have an advantage over crowdsourcing heuristics like majority rule [46]. Studies of group problem-solving have suggested that this 'truth wins' effect occurs when a shared conceptual system enables individuals to conclusively *demonstrate* that a given answer is correct or incorrect — and it is the

strength of their argument, rather than the individual's confidence or simply the presentation of the correct answer, which predicts whether the majority will be persuaded. Importantly, these studies suggest that "demonstrability" is a matter of degree, ranging from mathematics as the "preeminent domain of demonstrability" to purely judgmental tasks such as attitudes and preferences, with a variety of evidence-based reasoning and insight problems also being high in "demonstrability" [47-49]. Note the implication of the "truth wins" effect for social learners: if a minority is able to demonstrate that their judgment is accurate, the majority is not simply *influenced* by the judgment of the minority, they will *defer* to it. Thus, when demonstrations are possible, discussion groups may offer substantially more accurate collective judgments than a "crowdsourced" majority, with little risk of distorting an accurate majority judgment. Indeed, recent accounts suggest that reasoning itself may be most naturally deployed in service of argumentation and function most effectively in interpersonal contexts [50-51].

Note that we are not claiming that group discussion is *only* beneficial for questions that afford demonstrative reasoning, or that demonstration and reasoning are synonymous. Rather, we focus on discussion and crowdsourcing as flexible, commonsense approaches to a fundamental problem for any social learner: integrating information from multiple sources without inheriting their errors. Our intent is to examine laypeople's intuitions about their tradeoffs. Though crowdsourcing heuristics like majority rule can be remarkably accurate, they also presuppose independent judges — an unrealistic assumption about human societies. Meanwhile, work on group problem-solving has repeatedly found that discussion not only allows groups to outperform heuristics like majority rule, but that their ability to do so depends on the "demonstrability" of the problem [57-49, 52]. Of course, demonstration is possible without reasoning (e.g., by physically demonstrating how an artifact works or showing the location of an object), and reasoning cannot always conclusively demonstrate a that a solution is optimal. However, reasoning may be a reliable means of correcting errors even when physical demonstrations are not feasible, and when a correct answer cannot be simply deduced. For example, knowing the distance from New York to Chicago won't allow a group to deduce the distance from New York to Cleveland, but it may enable them correct some over- and underestimates without needing to actually measure the distance. Indeed, in a recent comparison of group discussion with the wisdom of crowds on a numerical estimation task, the average collective estimate of

four small-group discussions was more accurate than the average of 1,400 individual estimates, and participants reported arriving at their estimates by "sharing arguments and reasoning together" [53]. In short, to the extent that people expect to be able to rely on demonstrative reasoning to minimize the risks groupthink, it may be intuitive to disregard the importance of independent judgment, even if they favor crowdsourcing heuristics in other cases.

In the present work, we asked whether people would favor different social learning strategies for problems that afford demonstrative reasoning than those that do not. Crowdsourcing independent judgments may be more valuable when the potential for reasoning is less salient, particularly when the crowd is large. Discussion may be more valuable when demonstrative reasoning provides a reliable means of analyzing problems and identifying errors, even if the discussion group is small. Because past work suggests that sophisticated social learning strategies emerge in early childhood but also that children appear to underestimate some risks of social influence even in late childhood [21, 26], we focused on adults and children ages 7-10. Understanding how the ability to balance the risks and benefits of social influence develops could shed light on the incongruence of our remarkable capacity for collective problem-solving and our apparent susceptibility to groupthink. It may also provide clues as to where interventions to thwart misinformation may be most effective.

In each experiment (Fig 1), participants were shown eight questions (4 *Reasoning* and 4 *Non-Reasoning*), and for each question, they rated whether crowdsourcing or discussion would be more helpful in answering, on a 4-point scale. In Experiment 1, this meant that participants rated whether it would be more helpful to ask five people to each answer independently, or to ask the same five people to give a single group answer after discussing. In Experiments 2 and 3, we contrasted the five-person group discussion with a crowd of 50 people answering alone. For the *Reasoning* questions, we chose a set of constraint-satisfaction problems that would challenge adults' capacities, but still be understandable to children (e.g., Sudoku). Because the solutions to these questions must satisfy a mutually understood set of explicit constraints, discussion can help groups generate potential solutions and and reduce processing demands on individuals while relying on demonstrative reasoning to correct errors. In Experiments 1 and 2, we contrasted the *Reasoning* questions with *Population Preference* questions (e.g., most popular fruit in the world). Though individuals' intuitions may sway as the

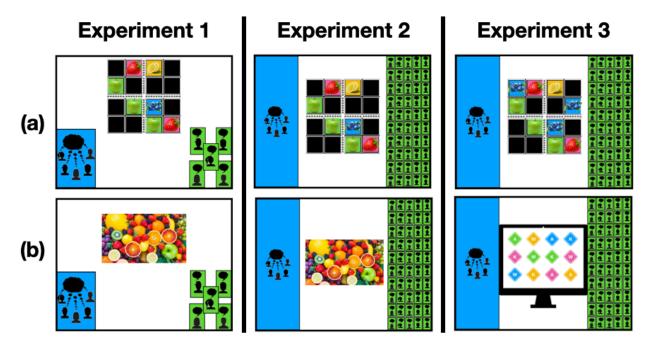
discussion generates potential answers, discussion provides no objective means of adjudicating disagreement; thus, it may distort intuitions rather than sharpen them. In Experiment 3, we contrasted easy versions of the same *Reasoning* questions with a set of challenging Perceptual Discrimination questions (e.g., fastest rotating item in an array), which a separate sample had rated as more difficult than the *Reasoning* questions. This allowed us to test the role of perceived difficulty against the potential for effective reasoning. If participants simply favor discussion for questions that feel more difficult — regardless of whether discussion can reliably adjudicate disagreement — then the preference for group discussion will be stronger for the *Perceptual Discrimination* questions than the *Reasoning* questions. Our general prediction in all three experiments was that sensitivity to the contrast between reasoning and intuitive judgment would lead all ages to prefer group discussion for reasoning questions. However, because past work has suggested that children may underestimate the risks of social influence until between the ages of 6 and 9 [21, 23-24, 26], we predicted that a robust preference for crowdsourcing non-reasoning questions would emerge only among older children (ages 9-10) and adults, while younger children (ages 7-8) would favor group discussion for both kinds of questions in Experiments 1 and 3. All experiments were preregistered, and the <u>data</u>, <u>materials</u>, <u>and power analyses</u> (<u>https://osf.io/6pw5n/?</u> <u>view\_only=1b5c7b4316e74c028c67eae0c9350b86</u>) are available on the OSF repository. All experiments were approved by the Yale University Institutional Review Board and conducted according to their guidelines. Written informed consent was obtained from all adult participants. Because children participated online, parents were recorded reading the informed consent form aloud.

# **Experiment 1**

#### Method

**Participants**. We recruited 40 adults through MTurk, as well as 80 children (40 Younger, M=8.01, SD=.56; 40 Older, M=9.92, SD=.56; 39 girls). Children participated through an online platform for developmental research that allows researchers to video chat with families using pictures and videos on slides [54]. Sample size was chosen based on the estimated effect size from pilot results.

**Materials.** We asked eight test questions (Fig 1), four from each of two question types: *Reasoning* and *Popularity*. Questions were presented from the perspective of a



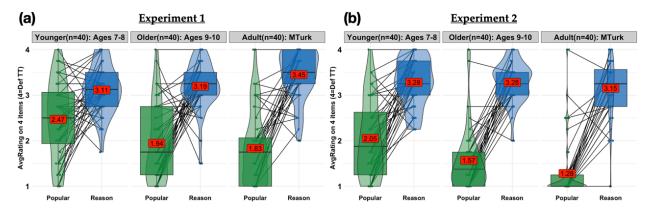
**Figure 1. Example stimuli: Groups and Crowds.** Each participant saw 8 questions. Experiments 1 and 2 used the same questions (4 *Reasoning*, 4 *Population Preference*). Experiment 3 contrasted *Easy* versions of the *Reasoning* questions with *Hard Perceptual Discrimination* questions. **(a)** Top Row: Example *Reasoning* question. Two-by-two Fruit Sudoku from Experiments 1 & 2 & partially completed "Easy" version in Experiment 3. **(b)** Bottom Row: Example *Population Preference* question (most popular fruit in the world) and *Perceptual Discrimination* question (which shape is spinning the fastest).

protagonist (Jack). The *Reasoning* questions were chosen to be simple enough to explain to children, but challenging enough that the answer would not be immediately obvious to adults. (1) A 4x4 Sudoku puzzle adapted for children. (2) A vehicle routing problem that required a MarioKart character to find the shortest road through all the treasures on a map without taking "two in a row that are the same color, or two in a row that are the same shape". (3) A single-heap game of Nim ("Each side takes turns picking up pencils. Each turn, you can pick up either 1, 2, or 3 pencils. The winner is the person who picks up the last pencil. There are 5 pencils left in this game; how many pencils should Jack pick up?"). (4) An "impossible object" puzzle that requires the solver to remove a dowel held in place by a nut and bolt from inside a bottle without breaking the dowel or the bottle. The *Popularity* questions concerned the most common subjective preferences in a population. (1) Whether pizza or hot dogs were more preferred by students in Jack's school. (2) What most people in the world say their favorite fruit is. (3) What most people in the world say their favorite fruit is.

their favorite color is. Questions were written to have approximately equal word counts  $(M_{Reas} = 73.75, M_{Pop} = 67.75)$ . Three counterbalances were created to vary the order of the questions — Forward, Reverse, and a Shuffle. Color coding of answer choice and left/right presentation were also counterbalanced between participants.

**Procedure.** Children were introduced to the protagonist, Jack (a silhouette). They were told that Jack was unsure of the answers to the questions, and could ask five people for help. The five people could either help by *Talking Together* (giving Jack a single answer as a group), or by Answering Alone (each giving Jack their own answer after thinking about the question without consulting others). For each item, children and adults rated whether "talking together" or "answering alone" was "probably more helpful, or definitely more helpful", producing a 4-point scale of relative preference, where 1 corresponds to "definitely answering alone", and 4 corresponds to "definitely talking together." Adults used the scale directly; children's responses were staggered: they first chose the more helpful strategy, and then were asked for a "probably/definitely" judgment. After answering the eight test items, participants were asked the two comprehension check questions (these were not counterbalanced: Comp\_TT was always presented first). Two features of the procedure are important to keep in mind. First, participants could not evaluate the content of any answer to any question, because none was given: they were asked to choose a *means* of advice, not evaluate the quality of the advice itself. Secondly, they could not make judgments based on degree or quality of consensus — they only knew that the group would have to give one answer, while the crowd would have to give 5 independent answers which could differ or not.

**Results.** For the primary test, the four responses within each question domain (Fig. 2a) were averaged to create a single score for each domain. A repeated measures ANOVA revealed a significant effect of question Type (F(1,117)=132.87, p<.001,  $\eta_p^2$  = .532) and an  $AgeGroup^*Type$  interaction (F(2,117)=7.83, p<.001,  $\eta_p^2$  = .118), and a marginal but non-significant effect of AgeGroup (F(2,117)=2.82, p=.064,  $\eta_p^2$  = .046). Multiple comparisons suggested that intuitions about how to manage collective wisdom appear by at least age 7: consistent with the empirical literature suggesting that group reasoning outperforms individual reasoning, all age groups believed that Talking Together would be more helpful than Answering Alone for Reasoning questions, both as compared to Popularity questions (Bonferroni corrected, Younger: t(117) = 3.66 p=0.0057, Older: t(117) = 7.105, p<.0001, Adult: t(117) = 9.201, p<.0001), and compared to chance



**Figure 2. Experiments 1-2: Results.** Preference for group discussion or independent crowd poll for (a) Experiment 1 (b) Experiment 2, averaged across four *Reasoning* questions (blue boxplots) and four *Population Preference* questions (green boxplots). Group discussion was 5 people in both experiments; crowd was 5 people in Experiment 1, and 50 people in Experiment 2. Higher ratings indicate stronger preference for group discussion. Boxplots showing median and interquartile range overlay violin plots; red labels show means; black lines show within-subject differences for the average rating by question *Type*.

(Younger: M=3.11, SD=.52, t(39) = 7.42, p<.0001, Older: M=3.19, SD=.51, t(39) = 8.53, p<.0001, Adult: M=3.45, SD=.54, t(39) = 11.237, p<.0001). Moreover, both Older children and Adults favored *Answering Alone* over *Talking Together* for *Popularity* questions, though Younger children's answers for *Popularity* did not differ significantly from chance (Younger: M=2.46, SD=.85 t(39) = -0.232, p=n.s., Older: M=1.94, SD=.86, t(39) = -4.076, p<.001, Adult: M=1.83, SD=.81, t(39) = -5.24, p<.0001). The preference for group reasoning did not differ by age (all ps>.4), though Older children and Adults showed a stronger preference for crowdsourcing *Popularity* questions than Younger children (Bonferroni corrected: Adult vs. Older: t(78)=0.718, p=ns; Adult vs. Younger: t(78)=4.067, p<0.001; Older vs. Younger: t(78)=3.350, p<0.0143). This developmental shift towards *Answering Alone* when discussion provides no objective criteria for evaluating accuracy is slightly earlier than we had predicted, but consistent with past work on children's evaluation of non-independent testimony [26].

Finally, all ages agreed that a teacher who wanted a group of five students to answer test questions accurately should have the students *Talk Together*, while a teacher who wanted to know which students had done their homework should have the students *Answer Alone* (Comp\_AA:  $M_{Young}=65\%$ , p=.04,  $M_{Old}=87.5\%$ , p<.0001,  $M_{Adult}=92.5\%$ , p<.0001, Comp\_TT:  $M_{Young}=70\%$ , p=.008,  $M_{Old}=85\%$ , p<.0001,  $M_{Adult}=87.5\%$ , p<.0001). This suggests that by age 7, children recognize that discussion could undermine

inferences about individuals' "independent" beliefs, but expect group discussion to either generate or disseminate accurate answers.

Taken together, these two tasks suggest that sophisticated intuitions about the risks and benefits of social influence may guide decisions about how to learn from collective judgment. Notably, these intuitions are consistent with empirical findings documenting the a group advantage over individuals for reasoning questions, and the value of independent responding when discussion is likely to bias collective judgment.

### **Experiment 2**

Could Experiment 1 have underestimated the value of crowdsourcing? Crowdsourcing may be most valuable with large crowds: larger crowds are more likely to include at least one accurate individual, and better represent the relative frequency of beliefs in the population. Moreover, in large enough crowds, even a minimal plurality will easily outnumber the unanimous consensus of a small group. Thus, if a belief's frequency is a cue to its accuracy, a large crowd will always be more informative than a small group. In Experiment 2, we contrasted the 5-person group with a larger 50-person crowd. We predicted that since the *Popularity* questions simply ask the group or crowd to estimate what *most* people in a population prefer, all age groups would find it intuitive to ask *more* people — i.e., the crowd. The benefit of large crowds is less clear for *Reasoning* questions. If few individuals can solve a problem alone, identifying the correct answer in the crowd may be akin to finding a needle in a haystack; indeed, if individual accuracy is known to be rare, the most common answer may be a widelyshared misconception [53]. Yet, if many individuals can solve the problem alone, large crowds are redundant and a learner can outsource evaluating accuracy to a group discussion. We therefore predicted that adults and older children would continue to favor group deliberation over crowdsourcing for *Reasoning* questions. However, we saw two plausible alternatives for younger children. First, younger children could show the mature pattern. Alternatively, younger children's preference for reasoning in groups could be attenuated by a "more is better" bias. Additionally, since the only difference between Experiments 1 and 2 was the increased crowd size, our design also allows us to explore the effects of crowd size itself by comparing the two experiments directly.

### Method

**Participants**. We recruited 40 adults through MTurk, as well as 80 children (40 Younger, M=8.01, SD=.56; 40 Older, M=9.92, SD=.56; 39 girls). As in Experiment 1, children participated through an online platform for developmental research that allows researchers to video chat with families using pictures and videos on slides [53]. One additional child was excluded and replaced because the family lost internet connection partway through the experiment and could not rejoin.

**Materials & Procedure.** The materials and procedure were identical to Experiment 1, but participants were first shown a large crowd of people, and told that Jack could either ask 5 of them to Talk Together, or 50 of them to answer alone. The answer choices from Experiment 1 were altered to display fifty cartoon icons for *Answering Alone* instead of five.

**Results.** As before, the four responses for each question *Type* (Fig. 2b) were averaged to create a single score for each *Type*. A repeated measures ANOVA revealed a significant effect of *Type* (F(1,117)=376.88, p<.001,  $\eta_p^2 = .763$ ) and *AgeGroup*  $(F(2,117)=9.63, p<.001 \eta_{p^2}=.141)$ , and an AgeGroup\*Type interaction (F(2,117)=5.39, p<.01, $\eta_p^2 = .084$ ). Despite the crowd having ten times as many sources as the group, participants were not swayed by a "more is better" bias; all age groups continued to prefer the group discussion for *Reasoning* questions, both as compared to *Popularity* questions (Bonferroni corrected, Younger: t(117) = 8.60 p < .0001, Older: t(117) = 11.97, p < .00010001, Adult: t(117) = 13.06, p<.0001), and compared to chance responding (Younger: M=3.28, SD=.53, t(39) = 9.29, p<.0001, Older: M=3.28, SD=.42, t(39) = 11.75, p<.0001, Adult: M=3.15, SD=.65, t(39) = 6.37, p<.0001). Moreover, even younger children in Experiment 2 favored *Answering Alone* for *Popularity* questions, suggesting that they recognized that a large crowd would provide a better estimate of population preferences than a small group (Younger: M=2.05, SD=.84, t(39) = -3.38, p=.0017, Older: M=1.57, SD=.69, t(39) = -8.52, p<.0001, Adult: M=1.28, SD=.66, t(39) = -11.75, p<.0001). As in Experiment 1, the preference for group reasoning did not differ by age (all ps > .9), though Older children and Adults again showed a stronger preference for crowdsourcing *Popularity* questions than Younger children (Bonferroni corrected, Adult vs. Older: *t*(78)=1.995, *p*=*ns*; Adult vs. Younger: *t*(78)=5.334, *p*<0.001; Older vs. Younger: t(78)= 3.339, p< 0.0147). We also conducted two exploratory analyses of the effect of crowd size. Our preregistered prediction in Experiment 2 was that participants would favor the crowd for population preference questions, but continue to favor the group for reasoning questions. However, because the only difference between Experiments 1 and 2 was the increase in crowd size from 5 to 50 people, our data also enables us to test the crowd-size effect directly. We ran separate ANOVAs for each *QuestionType* using *AgeGroup & Experiment* as predictors. The tenfold increase in crowd size had no impact on participants' preference for discussing *Reasoning* questions in small groups (F(1, 234)=0.045, p=.8320); an *AgeGroup\*ExpNum* interaction was significant (F(2,234)=4.434, p=.0129), but post-hoc comparisons revealed only a marginal difference between younger children's and adults' preferences for reasoning in groups in Exp 1, with no other differences. However, participants were significantly more likely to crowdsource *Popularity* questions in Experiment 2 than Experiment 1 (F(1, 234)=19.303, p<0.0001), with no differences between age groups.

As in Experiment 1, responses to the comprehension questions at the end of the task suggested even the youngest children recognized that talking together would make it impossible for the teacher to know which students had done their homework (Comp\_AA:  $M_{Young}=67.5\%$ , p=.019,  $M_{Old}=92.5\%$ , p<.0001,  $M_{Adult}=90\%$ , p<.0001). However, while older children and adults agreed that the students would do better on the test if they could discuss their answers, younger children were at chance (Comp\_TT:  $M_{Young}=52.5\%$ , p=.4373,  $M_{Old}=90\%$ , p<.0001,  $M_{Adult}=90\%$ , p<.0001). Younger children may be less confident in the value of discussion than their responses to the main task questions in Experiments 1 and 2 would suggest; however, informal questioning of participants after the experiment suggested that younger children in Exp 2 may have simply rejected talking together on a test as cheating, even though the question specified that the teacher could choose to allow students to talk together.

In short, Experiment 2 suggests that not only are young children's intuitions about the value of group discussion consistent with empirical demonstrations of a group advantage for reasoning questions and the value of large crowds for intuitive estimations. Moreover, directly comparing Experiments 1 and 2 suggests that while children's preference for reasoning in small groups is stable even in the face of a much larger crowd, they also recognize that for some questions, larger crowds are more helpful than smaller crowds.

#### **Experiment 3**

Using *Population Preferences* as the Non-Reasoning questions in Experiments 1 and 2 leaves two points unclear. First, since a culture's preferences are intuitive for most

people, the *Popularity* questions may have simply seemed easier to answer than the Reasoning questions. Second, because individual preferences are literally constitutive of the population preference, children's responses could reflect an understanding of the nature of preference polling as much as an understanding of the potential for groupthink. To test these two alternatives, we contrasted easy versions of the reasoning questions with challenging perceptual discrimination questions. Disagreement about a challenging perceptual discrimination task would leave a group of laypeople little to discuss beyond confidence, which may be sufficient to filter out obviously wrong answers [56-57], but is generally an unreliable proxy for accuracy. In contrast, polling a large crowd has been shown to increase the accuracy of a collective decision for perceptual tasks [58]. The relative difficulty of the Easy Reasoning and Hard Percept items was confirmed in a pre-test (Supplemental Materials). If participants' preference for group discussion in Experiments 1 and 2 was driven by perceived question difficulty, then they will prefer group discussion for *Hard Percept* questions more than for *Easy* Reasoning questions. If group discussion was preferred because of its perceived benefits for reasoning, participants will prefer group discussion more for Easy Reasoning than *Hard Percept*. If participants recognize the risks of social influence when discussants cannot rely on demonstrative reasoning, they will prefer crowdsourcing for Hard Percept questions. We predicted that adults and older children would recognize the tradeoffs, but because children under 8 frequently fail to recognize the potential for motivational biases even in simpler cases [59], we predicted that younger children would prefer group discussion for both question types. As an exploratory analysis, we also compare the results in Experiment 3 directly to Experiment 2, but given that both the perceived difficulty and the subtype of Non-Reasoning question differ between Experiments, direct comparisons should be interpreted with caution.

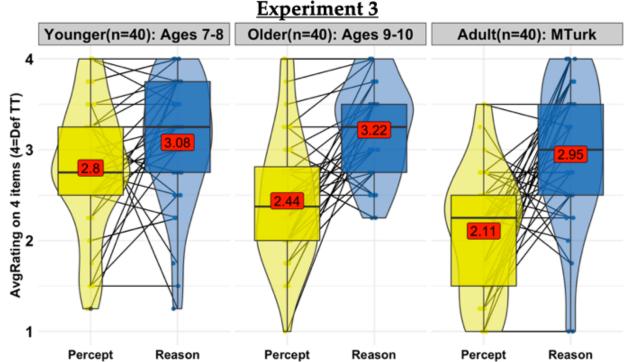
#### Method

**Participants**. We recruited 40 adults through MTurk, as well as 80 children (40 Younger, M=8.01, SD=.62; 40 Older, M=10.00, SD=.53; 37 girls). As in Experiments 1 & 2, children participated through our online platform [54]. Two children were excluded and replaced when database records identified them afterwards as having already participated in Experiment 2. Two adults were excluded and replaced as well; though our preregistered plan was to accept all MTurkers who passed the basic attention screening, two worker identification codes appeared multiple times in the data, passing

the attention screen after failing and being screened out two and three times, respectively, in violation of MTurk policies.

Materials & Procedure. Methods were identical to Experiment 2, with the exception of the following changes made to the questions themselves. First, we presented four new Non-Reasoning questions, replacing the four *Popularity* questions with four *Percept* questions: (1) decide which of two pictures of a face "at the tipping point of animacy" is a photo and which is a photorealistic drawing [60], (2) decide whether an opaque box contains 30 or 40 marbles by listening to a recording of it being shaken [61], (3) identify which of twelve colored squares in a visual array is rotating the fastest, and (4) rank the 25 brightest stars in a photo of the night sky in order of brightness. Second, we simplified the four *Reasoning* questions (see Supplemental Materials) by (1) completing most of the Sudoku, (2) reducing the number of treasures Mario was required to pick up in the vehicle routing problem, and (3) replacing the "impossible object bottle" with an analog of the "floating peanut" task, which requires the learner to extract an object from a jar of water without touching the jar or object [62]. The fourth Reasoning question, Nim, remained the same, as adults rated the 5-item Nim heap as easy to solve.

**Results.** For the primary test, the four responses within each question domain (Fig. 3) were again averaged to create a single score for each *Type*. A repeated measures ANOVA again revealed a significant effect of question Type (F(1,117)=56.12, p<.0001,  $\eta_p^2$ = .324) and AgeGroup (F(1,117)=7.01, p=.0012,  $\eta_p^2$  = .108) an AgeGroup\*Type interaction  $(F(2,117)=4.40, p=.0143, \eta_{p^2}=.070)$ . The perceived difficulty of the questions had no discernible effect on participant judgments: participants of all ages again rejected the large crowd in favor of the small group discussion for the Easy Reasoning questions (Younger: M=3.08, SD=.69 t(39) = 5.35, p < .0001, Older: M=3.22, SD=.51, t(39) = 8.97, p < .00010001, Adult: M=2.95, SD=.78, t(39) = 3.64, p=.0008). We also observed the predicted developmental shift towards Answering Alone when reasoning was insufficient to answer the question; however, in Experiment 3 the shift occurred later than expected instead of earlier. While Adults favored Answering Alone and Younger children favored Talking Together for the Hard Percept questions as predicted, older children did not show the adult pattern, instead not differing from chance for Hard Percept (Younger: M=2.80, SD=.71, t(117) = 2.66 p=0.011; Older: M=2.44, SD=.71, t(117) = -0.55, p=0.7535; Adult: M=2.11, SD=.69, t(117) = -3.62, p=.0008). Moreover, while Older children and Adults distinguished between the two question Types, Younger children did not (Bonferroni



**Figure 3. Experiment 3: Results.** Preference for 5-person group discussion or independent crowd poll of 50 people, averaged across four *Easy Reasoning* questions (blue boxplots) and four *Hard Perceptual Discrimination* questions (yellow boxplots). Higher ratings indicate stronger preference for group discussion. Boxplots showing median and interquartile range overlay violin plots; red labels show means; black lines show within-subject differences for the average rating by question *Type*.

corrected, Younger: t(117) = -1.91 p=.8701, Older: t(117) = -5.318, p<.0001, Adult: t(117) = 5.74, p=.0001). As in Experiments 1 and 2, the preference for group reasoning did not differ by age (all ps > .9), though Younger children showed a weaker preference than Adults for crowdsourcing *Percept* questions (Bonferroni corrected, Adult vs. Older: t(78)=2.156, p=ns; Adult vs. Younger: t(78)=4.516, p<0.0002; Older vs. Younger: t(78)=2.360, p=ns). Indeed, while participants of all ages were just as confident in the small group discussion for *Easy Reasoning* questions in Experiment 3 as they were for *Reasoning* questions in Experiment 2, all ages were less confident in polling a crowd of 50 for *Hard Percept* questions in Experiment 3 than for *Population Preferences* in Experiment 2 (Supplemental Materials). However, since Experiment 3 was designed contrast *Easy Reasoning* questions with *Hard Percept* questions, rather than *Hard Percept* with *Population Preferences*, these direct comparisons with Experiment 2 should be interpreted with caution: for example, the weaker preference for crowdsourcing *Hard Percept* questions than *Population Preference* questions may be due the difference in Non-

Reasoning subtype or an effect of difficulty that is specific to Non-Reasoning questions. We explore these possibilities further in the General Discussion.

#### **General Discussion**

We asked children and adults to choose between two social learning strategies: soliciting a consensus response from a small discussion group, and "crowdsourcing" many independent opinions. Though discussion can sometimes lead to groupthink, by affording individuals opportunity to correct each others' mistakes and combine insights while also reducing individual processing load, discussion can also allow small groups to outperform even their best member. In contrast, the value of crowdsourcing is fundamentally limited by the distribution of individual competence in the crowd relative to its size. The less competent individuals are on average, the larger the crowd needs to be to produce a reliably accurate estimate. Thus, when individual competence is low, crowdsourcing may be costly; when individual competence is high, the value added by crowdsourcing may have little advantage over discussion — for problems where discussion is more likely to improve accuracy than diminish it. Our results suggest that the decision to crowdsource or discuss may in part turn on learners' beliefs about the efficacy of demonstrative reasoning for a given question.

Analogously to young children's failures on false belief tasks, our results suggest that the default expectation for group judgments may be that "truth wins": though individuals may initially disagree, discussion allows groups to ultimately see the truth. As an understanding of how conscious and unconscious biases can influence people's judgments develops, learners can preempt potential biases by crowdsourcing independent judgments. Though even the youngest children in our experiments expected discussion to improve accuracy on reasoning questions, the preference for crowdsourcing non-reasoning questions underwent a developmental shift in all three experiments. Indeed, in Experiment 3, the youngest children favored discussion for both kinds of questions, suggesting that they may have failed to recognize when discussion can promote groupthink. The timing of the developmental shift is consistent with past work suggesting that between the ages of 6 and 9, children begin to use informational dependencies [21-26] and the potential for motivational bias in individual reports [29, 59] to adjudicate cases of conflicting testimony. Though recent work suggests that even preschoolers identify cases of individual bias stemming from ingroup favoritism [63], unconscious biases due to herding or groupthink may be less

obvious, particularly if people assume that informants are motivated to be accurate. For example, even though children as young as six *predict* that judges are more likely to independently give the same verdict when objective standards are available than when they are not (e.g., a footrace vs. a poetry contest), at age ten children are still no more likely to diagnose in-group favoritism as an influence on judgments in subjective contexts than objective contexts [28, 63]. In our experiments, both the reasoning and non-reasoning questions had objective answers, but only the reasoning questions afforded an objective method of finding those answers. Learning to recognize this relatively subtle distinction may allow children to take advantage of the benefits of group discussion while avoiding the risks. This is not to suggest that people expect group reasoning to be infallible — merely that they expect groups to improve individual accuracy. This is consistent with recent work asking adults to predict group and individual accuracy on a classic reasoning task: while participants radically underestimated the true group advantage, they did expect groups to be more accurate than individuals [64]. Interestingly, they also expected *dyads* to be *less* accurate than individuals. A more granular approach to intuitive beliefs about the dynamics of social influence may reveal more sophisticated intuitions: for instance, beliefs about others' conformist tendencies and the distribution of individual competence may increase confidence that "truth wins" in small groups more than in dyads.

Past work has suggested that while people dramatically underestimate crowds and overestimate their own accuracy [65,17], they defer to others more when uncertainty is high and crowds are larger [67]. While increasing the crowd size from five to fifty had no impact on *Reasoning* questions in our experiments, the larger crowd did appear to increase crowdsourcing for *Population Preference* questions. However, while we only tested *Hard Percept* questions with a crowd of fifty, confidence in crowdsourcing was lower for *Hard Percept* questions than for *Population Preferences* (Supplemental Materials). While our design licenses no firm conclusions on this point, one reason seems evident: by definition, population preferences are whatever most individuals in a population prefer, while perceptual facts like the brightness of stars are wholly independent of individual judgments. Moreover, under the right conditions, discussing perceptual judgments with a single partner *can* improve accuracy [68-69]. Thus, participants' reduced confidence in crowdsourcing *Hard Percept* questions may have been justified. The extent to which intuitive beliefs about the benefits of discussion and

crowdsourcing for different question types correspond to the empirical benefits is an open question.

Our design is limited in one important respect: the discussion group was only allowed to give a single answer, while the crowd could give multiple answers. This procedure strictly ensured that group members could not answer independently, but also entailed a unanimous consensus endorsed by a minimum of five people. Unanimous consensus can be a powerful cue: even a single dissenter can sharply reduce conformity [67, 24]. However, the meaning of dissent may vary across contexts and questions. In a crowd, a single "dissenter" may simply have made a mistake; but dissent-despite-discussion signals that the group has failed to convince them. When questions afford conclusive demonstrations of accuracy, failure to convince all discussants may reflect poorly on group accuracy. Conversely, in more ambiguous contexts, unanimity may suggest groupthink. For instance, in ancient Judea, crimes more likely to elicit widespread condemnation were tried by larger juries for the express purpose of reducing the odds of consensus, and unanimous convictions were thrown out on the grounds that a lack of dissent indicated a faulty process — an intuitive inference confirmed by modern statistical techniques [70]. A similar logic may underlie inferences about testimony that contradicts social alliances. For example, if Jenny says Jill is bad at soccer, even preschoolers give Jenny's judgment more credence if Jenny and Jill are friends than if they are enemies [63]. Our results suggest that even in early childhood, the absolute number of sources endorsing a belief may be less important than how those sources arrived at their beliefs. Indeed, the limited number of possible answers to the questions in Experiments 2 and 3 guaranteed that even a plurality of the 50-person crowd would considerably outnumber the 5-person group. Yet, participants' preference for discussion and crowdsourcing bore no relationship to the number of possible endorsers. Future work will compare explicit degrees of consensus in groups and crowds.

The last decade has produced an extensive literature describing how individual social learning heuristics and patterns of communication in social networks can improve or diminish collective learning [71-73, 33]. By focusing on population-level outcomes, much of this work has tacitly treated individuals as passive prisoners of social influence. However, the heuristics guiding social learning develop in early childhood, and recent work has shown that like other intelligent systems capable of self-

organization, people are capable of "rewiring" their social networks to improve both individual and collective learning, by "following" or "unfollowing" connections depending on their accuracy [73]. Our experiments focused on two features of communication patterns that individuals can and do control in the real world, beyond who they choose to trust: how many people to talk to, and whether to talk with those people as a group or a crowd. Our results suggest that even in early childhood, people's judgments about how to best make use of group discussion and crowdsourcing heuristics may be consistent with the empirical advantages of each strategy. An understanding of how intuitions about social influence develop may contribute to a clearer empirical picture of how people balance the benefits of learning from collective opinion with the risks of being misled by it.

### **References**

- 1. Hastie, R., & Kameda, T. (2005). The Robust Beauty of Majority Rules in Group Decisions. *Psychological Review*, 112(2), 494–508. <a href="https://doi.org/10.1037/0033-295X.112.2.494">https://doi.org/10.1037/0033-295X.112.2.494</a>
- 2. Laan, A., Madirolas, G., & de Polavieja, G. G. (2017). Rescuing Collective Wisdom when the Average Group Opinion Is Wrong. *Frontiers in Robotics and AI*, 4. <a href="https://doi.org/10.3389/frobt.2017.00056">https://doi.org/10.3389/frobt.2017.00056</a>
- 3. Yousif, S. R., Aboody, R., & Keil, F. C. (2019). The Illusion of Consensus: A Failure to Distinguish Between True and False Consensus. *Psychological Science*, *30*(8), 1195–1204. <a href="https://doi.org/10.1177/0956797619856844">https://doi.org/10.1177/0956797619856844</a>
- 4. Mercier, H., & Miton, H. (2019). Utilizing simple cues to informational dependency. *Evolution and Human Behavior*, 40(3), 301–314. <a href="https://doi.org/10.1016/j.evolhumbehav.2019.01.001">https://doi.org/10.1016/j.evolhumbehav.2019.01.001</a>
- 5. Gerstenberg, T., & Tenenbaum, J. B. (2017). Intuitive Theories. In M. R. Waldmann (Ed.), *The Oxford Handbook of Causal Reasoning* (Vol. 1). Oxford University Press. <a href="https://doi.org/10.1093/oxfordhb/9780199399550.013.28">https://doi.org/10.1093/oxfordhb/9780199399550.013.28</a>
- 6. Aristotle, ., Jowett, B., & Davis, H. W. C. (1920). Aristotle's Politics. Oxford: At the Clarendon Press.
- 7. Plato, ., & Skemp, J. B. (1952). Statesman. London: Routledge & K. Paul.
- 8. List, C., & Goodin, R. E. (2001). Epistemic Democracy: Generalizing the Condorcet Jury Theorem. *Journal of Political Philosophy*, 9(3), 277–306. <a href="https://doi.org/10.1111/1467-9760.00128">https://doi.org/10.1111/1467-9760.00128</a>
- 9. Dietrich, F., & Spiekermann, K. (2013). Epistemic Democracy with Defensible Premises. *Economics and Philosophy*, 29, 34. <a href="https://doi.org/doi:10.1017/50266267113000096">https://doi.org/doi:10.1017/50266267113000096</a>
- 10. Boyd, R., & Richerson, P. J. (1988). An evolutionary model of social learning: the effects of spatial and temporal variation. Social learning: psychological and biological perspectives, 29-48.
- 11. Galton, F. (1907). Vox populi. *Nature*, *75*, 450–451. <a href="https://doi.org/doi.org/10.1038/075450a0">https://doi.org/doi.org/doi.org/10.1038/075450a0</a>
- 12. Hong, L., & Page, S. E. (2004). Groups of diverse problem solvers can outperform groups of high-ability problem solvers. *Proceedings of the National Academy of Sciences*, 101(46), 16385–16389. <a href="https://doi.org/10.1073/pnas.0403723101">https://doi.org/10.1073/pnas.0403723101</a>
- 13. Steyvers, M., Miller, B., Hemmer, P., & Lee, M. D. (2009). The Wisdom of Crowds in the Recollection of Order Information. *Advances in Neural Information Processing Systems*, 9.
- 14. de Oliveira, S., & Nisbett, R. E. (2018). Demographically diverse crowds are typically not much wiser than homogeneous crowds. *Proceedings of the National Academy of Sciences*, 115(9), 2066–2071. https://doi.org/10.1073/pnas.1717632115

- 15. Morgan, T.J.H., Laland, K. N., & Harris, P. L. (2015). The development of adaptive conformity in young children: Effects of uncertainty and consensus. *Developmental Science*, *18*(4), 511–524. <a href="https://doi.org/10.1111/desc.12231">https://doi.org/10.1111/desc.12231</a>
- 16. Muthukrishna, M., Morgan, T. J. H., & Henrich, J. (2016). The when and who of social learning and conformist transmission. *Evolution and Human Behavior*, 37(1), 10–20. <a href="https://doi.org/10.1016/j.evolhumbehav.2015.05.004">https://doi.org/10.1016/j.evolhumbehav.2015.05.004</a>
- 17. Mannes, A. E. (2009). Are We Wise About the Wisdom of Crowds? The Use of Group Judgments in Belief Revision. *Management Science*, 55(8), 1267–1279. <a href="https://doi.org/10.1287/mnsc.1090.1031">https://doi.org/10.1287/mnsc.1090.1031</a>
- 18. Morgan, T. J. H., Rendell, L. E., Ehn, M., Hoppitt, W., & Laland, K. N. (2012). The evolutionary basis of human social learning. *Proceedings of the Royal Society B: Biological Sciences*, 279(1729), 653–662. <a href="https://doi.org/10.1098/rspb.2011.1172">https://doi.org/10.1098/rspb.2011.1172</a>
- 19. Raafat, R. M., Chater, N., & Frith, C. (2009). Herding in humans. *Trends in Cognitive Sciences*, *13*(10), 420–428. <a href="https://doi.org/10.1016/j.tics.2009.08.002">https://doi.org/10.1016/j.tics.2009.08.002</a>
- 20. Lorenz, J., Rauhut, H., Schweitzer, F., & Helbing, D. (2011). How social influence can undermine the wisdom of crowd effect. *Proceedings of the National Academy of Sciences*, *108*(22), 9020–9025. <a href="https://doi.org/10.1073/pnas.1008636108">https://doi.org/10.1073/pnas.1008636108</a>
- 21. Aboody, R., Yousif, S. R., Sheskin, M., & Keil, F. C. (2019). Says who? Children consider informants' sources when deciding whom to believe. *Poster presented at Cognitive Development Society*, Louisville, KY.
- 22. Sulik, J., Bahrami, B., & Deroy, O. (2020). Social influence and informational independence. *Proceedings of the Cognitive Science Society*, 7. <a href="https://cognitivesciencesociety.org/cogsci20/papers/0704/0704">https://cognitivesciencesociety.org/cogsci20/papers/0704/0704</a>
- 23. Magid, R. W., Yan, P., Siegel, M. H., Tenenbaum, J. B., & Schulz, L. E. (2018). Changing minds: Children's inferences about third party belief revision. *Developmental Science*, 21(2), e12553. <a href="https://doi.org/10.1111/desc.12553">https://doi.org/10.1111/desc.12553</a>
- 24. Whalen, A., Griffiths, T. L., & Buchsbaum, D. (2018). Sensitivity to Shared Information in Social Learning. *Cognitive Science*, 42(1), 168–187. <a href="https://doi.org/10.1111/cogs.12485">https://doi.org/10.1111/cogs.12485</a>
- 25. Anderson, L. R., & Holt, C. A. (1997). Information Cascades in the Laboratory. *The American Economic Review*, 87(5), 17. <a href="https://www.jstor.org/stable/2951328">https://www.jstor.org/stable/2951328</a>
- 26. Einav, S. (2018). Thinking for themselves? The effect of informant independence on children's endorsement of testimony from a consensus. *Social Development*, 27(1), 73–86. <a href="https://doi.org/10.1111/sode.12264">https://doi.org/10.1111/sode.12264</a>
- 27. Hu, J., Whalen, A., Buchsbaum, D., Griffiths, T., & Xu, F. (2015). Can Children Balance the Size of a Majority with the Quality of their Information? *Proceedings of the Annual Meeting of the Cognitive Science Society*, 35, 6.
- 28. Mills, C. M., & Keil, F. C. (2005). The Development of Cynicism. *Psychological Science*, *16*(5), 385–390. <a href="https://doi.org/10.1111/j.0956-7976.2005.01545.x">https://doi.org/10.1111/j.0956-7976.2005.01545.x</a>

- 29. Mills, C. M., & Grant, M. G. (2009). Biased decision-making: Developing an understanding of how positive and negative relationships may skew judgments. *Developmental Science*, 12(5), 784–797. <a href="https://doi.org/10.1111/j.1467-7687.2009.00836.x">https://doi.org/10.1111/j.1467-7687.2009.00836.x</a>
- 30. Marks, G., & Miller, N. (1987). Ten Years of Research on the False-Consensus Effect: An Empirical and Theoretical Review. *Psychological Bulletin*, 102(1), 19.
- 31. Lerman, K., Yan, X., & Wu, X.-Z. (2016). The "Majority Illusion" in Social Networks. *PLOS ONE*, *11*(2), e0147617. <a href="https://doi.org/10.1371/journal.pone.0147617">https://doi.org/10.1371/journal.pone.0147617</a>
- 32. Stewart, A. J., Mosleh, M., Diakonova, M., Arechar, A. A., Rand, D. G., & Plotkin, J. B. (2019). Information gerrymandering and undemocratic decisions. *Nature*, 573(7772), 117–121. <a href="https://doi.org/10.1038/s41586-019-1507-6">https://doi.org/10.1038/s41586-019-1507-6</a>
- 33. Becker, J., Brackbill, D., & Centola, D. (2017). Network dynamics of social influence in the wisdom of crowds. *Proceedings of the National Academy of Sciences*, 201615978. https://doi.org/10.1073/pnas.1615978114
- 34. Becker, J., Porter, E., & Centola, D. (2019). The wisdom of partisan crowds. *Proceedings of the National Academy of Sciences*, 116(22), 10717–10722. <a href="https://doi.org/10.1073/pnas.1817195116">https://doi.org/10.1073/pnas.1817195116</a>
- 35. Abel, M., & Bäuml, K.-H. T. (2020). Social interactions can simultaneously enhance and distort memories: Evidence from a collaborative recognition task. *Cognition*, 200, 104254. <a href="https://doi.org/10.1016/j.cognition.2020.104254">https://doi.org/10.1016/j.cognition.2020.104254</a>
- 36. Mason, W., & Watts, D. J. (2012). Collaborative learning in networks. *Proceedings of the National Academy of Sciences*, 109(3), 764–769. <a href="https://doi.org/10.1073/pnas.1110069108">https://doi.org/10.1073/pnas.1110069108</a>
- 37. Derex, M., & Boyd, R. (2015). The foundations of the human cultural niche. *Nature Communications*, *6*(1). <a href="https://doi.org/10.1038/ncomms9398">https://doi.org/10.1038/ncomms9398</a>
- 38. Barkoczi, D., & Galesic, M. (2016). Social learning strategies modify the effect of network structure on group performance. *Nature Communications*, 7(1). <a href="https://doi.org/10.1038/ncomms13109">https://doi.org/10.1038/ncomms13109</a>
- 39. Kirschner, F., Paas, F., & Kirschner, P. A. (2009a). Individual and group-based learning from complex cognitive tasks: Effects on retention and transfer efficiency. *Computers in Human Behavior*, 25(2), 306–314. <a href="https://doi.org/10.1016/j.chb.2008.12.008">https://doi.org/10.1016/j.chb.2008.12.008</a>
- 40. Kirschner, F., Paas, F., & Kirschner, P. A. (2009b). A Cognitive Load Approach to Collaborative Learning: United Brains for Complex Tasks. *Educational Psychology Review*, 21(1), 31–42. <a href="https://doi.org/10.1007/s10648-008-9095-2">https://doi.org/10.1007/s10648-008-9095-2</a>
- 41. Laughlin, P. R. (2011). Group Problem Solving. Princeton University Press.
- 42. Smith, M. K., Wood, W. B., Adams, W. K., Wieman, C., Knight, J. K., Guild, N., & Su, T. T. (2009). Why Peer Discussion Improves Student Performance on In-Class Concept Questions. *Science*, 323, 4. <a href="https://doi.org/10.1126/science.1165919">https://doi.org/10.1126/science.1165919</a>

- 43. Laughlin, P. R., Bonner, B. L., & Altermatt, T. W. (1998). Collective versus individual induction with single versus multiple hypotheses. *Journal of Personality and Social Psychology*, 75(6), 1481–1489. <a href="https://doi.org/10.1037/0022-3514.75.6.1481">https://doi.org/10.1037/0022-3514.75.6.1481</a>
- 44. Laughlin, P. R., Bonner, B. L., & Miner, A. G. (2002). Groups perform better than the best individuals on Letters-to-Numbers problems. *Organizational Behavior and Human Decision Processes*, 16. <a href="https://doi.org/10.1016/S0749-5978(02)00003-1">https://doi.org/10.1016/S0749-5978(02)00003-1</a>
- 45. Almaatouq, A., Alsobay, M., Yin, M., & Watts, D. J. (2021). Task complexity moderates group synergy. *Proceedings of the National Academy of Sciences*, 118(36), e2101062118. <a href="https://doi.org/10.1073/pnas.2101062118">https://doi.org/10.1073/pnas.2101062118</a>
- 46. Moshman, D., & Geil, M. (1998). Collaborative Reasoning: Evidence for Collective Rationality. *Thinking & Reasoning*, 4(3), 231–248. <a href="https://doi.org/">https://doi.org/</a> 10.1080/135467898394148
- 47. Trouche, E., Sander, E., & Mercier, H. (2014). Arguments, more than confidence, explain the good performance of reasoning groups. *Journal of Experimental Psychology: General*, 143(5), 1958–1971. <a href="https://doi.org/10.1037/a0037099">https://doi.org/10.1037/a0037099</a>
- 48. Laughlin, P. R., & Ellis, A. L. (1986). Demonstrability and Social Combination Processes on Mathematical Intellective Tasks. *Journal of Experimental Social Psychology*, 22, 177–189. <a href="https://doi.org/doi.org/10.1016/0022-1031(86)90022-3">https://doi.org/doi.org/10.1016/0022-1031(86)90022-3</a>
- 49. Larson, J. R. (2010). In search of synergy in small group performance. Psychology Press.
- 50. Mercier, H. (2016). The Argumentative Theory: Predictions and Empirical Evidence. *Trends in Cognitive Sciences*, 20(9), 689–700. <a href="https://doi.org/10.1016/j.tics.2016.07.001">https://doi.org/10.1016/j.tics.2016.07.001</a>
- 51. Mercier, H., & Sperber, D. (2011). Why do humans reason? Arguments for an argumentative theory. *Behavioral and Brain Sciences*, 34(2), 57–74. <a href="https://doi.org/10.1017/S0140525X10000968">https://doi.org/10.1017/S0140525X10000968</a>
- 52. Bonner, B. L., Shannahan, D., Bain, K., Coll, K., & Meikle, N. L. (2021). The Theory and Measurement of Expertise-Based Problem Solving in Organizational Teams: Revisiting Demonstrability. *Organization Science*, orsc.2021.1481. <a href="https://doi.org/10.1287/orsc.2021.1481">https://doi.org/10.1287/orsc.2021.1481</a>
- 53. Navajas, J., Niella, T., Garbulsky, G., Bahrami, B., & Sigman, M. (2018). Aggregated knowledge from a small number of debates outperforms the wisdom of large crowds. *Nature Human Behaviour*, 2(2), 126–132. <a href="https://doi.org/10.1038/s41562-017-0273-4">https://doi.org/10.1038/s41562-017-0273-4</a>
- 54. Sheskin, M., & Keil, F. (2018). The Child Lab.com A Video Chat Platform for Developmental Research. *PsyArxiv*. <a href="https://doi.org/10.31234/osf.io/rn7w5">https://doi.org/10.31234/osf.io/rn7w5</a>
- 55. Prelec, D., Seung, H. S., & McCoy, J. (2017). A solution to the single-question crowd wisdom problem. *Nature*, 541(7638), 532–535. <a href="https://doi.org/10.1038/nature21054">https://doi.org/10.1038/nature21054</a>

- 56. Massoni, S., & Roux, N. (2017). Optimal group decision: A matter of confidence calibration. *Journal of Mathematical Psychology*, 79, 121–130. <a href="https://doi.org/10.1016/j.jmp.2017.04.001">https://doi.org/10.1016/j.jmp.2017.04.001</a>
- 57. Bahrami, B., Olsen, K., Latham, P. E., Roepstorff, A., Rees, G., & Frith, C. D. (2010). Optimally Interacting Minds. *Science*, 329(5995), 1081–1085. <a href="https://doi.org/10.1126/science.1185718">https://doi.org/10.1126/science.1185718</a>
- 58. Juni, M. Z., & Eckstein, M. P. (2017). The wisdom of crowds for visual search. *Proceedings of the National Academy of Sciences*, 114(21), E4306–E4315. <a href="https://doi.org/10.1073/pnas.1610732114">https://doi.org/10.1073/pnas.1610732114</a>
- 59. Mills, C. M., & Keil, F. C. (2008). Children's developing notions of (im)partiality. *Cognition*, *107*(2), 528–551. <a href="https://doi.org/10.1016/j.cognition.2007.11.003">https://doi.org/10.1016/j.cognition.2007.11.003</a>
- 60. Looser, C. E., & Wheatley, T. (2010). The Tipping Point of Animacy: How, When, and Where We Perceive Life in a Face. *Psychological Science*, 21(12), 1854–1862. https://doi.org/10.1177/0956797610388044
- 61. Siegel, M. H., Magid, R., Tenenbaum, J. B., & Schulz, L. E. (2014). Black boxes: Hypothesis testing via indirect perceptual evidence. *Proceedings of the Cognitive Science Society*, 7.
- 62. Hanus, D., Mendes, N., Tennie, C., & Call, J. (2011). Comparing the Performances of Apes (Gorilla gorilla, Pan troglodytes, Pongo pygmaeus) and Human Children (Homo sapiens) in the Floating Peanut Task. *PLoS ONE*, *6*(6), e19555. <a href="https://doi.org/10.1371/journal.pone.0019555">https://doi.org/10.1371/journal.pone.0019555</a>
- 63. Liberman, Z., & Shaw, A. (2020). Even his friend said he's bad: Children think personal alliances bias gossip. *Cognition*, 204, 104376. <a href="https://doi.org/10.1016/j.cognition.2020.104376">https://doi.org/10.1016/j.cognition.2020.104376</a>
- 64. Mercier, H., Trouche, E., Yama, H., Heintz, C., & Girotto, V. (2015). Experts and laymen grossly underestimate the benefits of argumentation for reasoning. *Thinking & Reasoning*, 21(3), 341–355. <a href="https://doi.org/">https://doi.org/</a> 10.1080/13546783.2014.981582
- 65. Mercier, H., Dockendorff, M., Majima, Y., Hacquin, A.-S., & Schwartzberg, M. (2020). Intuitions about the epistemic virtues of majority voting. *Thinking & Reasoning*, 1–19. <a href="https://doi.org/10.1080/13546783.2020.1857306">https://doi.org/10.1080/13546783.2020.1857306</a>
- 66. Toyokawa, W., Whalen, A., & Laland, K. N. (2019). Social learning strategies regulate the wisdom and madness of interactive crowds. *Nature Human Behaviour*, 3(2), 183–193. <a href="https://doi.org/10.1038/s41562-018-0518-x">https://doi.org/10.1038/s41562-018-0518-x</a>
- 67. Asch, S. E. (1955). Opinions and social pressure. Scientific American, 193(5), 31-35.
- 68. Sorkin, R. D., Hays, C. J., & West, R. (2001). Signal-Detection Analysis of Group Decision Making. *Psychological Review*, 108(1), 21. <a href="https://doi.org/10.1037/0033-295X.108.1.183">https://doi.org/10.1037/0033-295X.108.1.183</a>
- 69. Bahrami, B., Olsen, K., Bang, D., Roepstorff, A., Rees, G., & Frith, C. (2012). What failure in collective decision-making tells us about metacognition. *Philosophical*

- *Transactions of the Royal Society B: Biological Sciences, 367*(1594), 1350–1365. https://doi.org/10.1098/rstb.2011.0420
- 70. Gunn, L. J., Chapeau-Blondeau, F., McDonnell, M. D., Davis, B. R., Allison, A., & Abbott, D. (2016). Too good to be true: When overwhelming evidence fails to convince. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 472(2187), 20150748. <a href="https://doi.org/10.1098/rspa.2015.0748">https://doi.org/10.1098/rspa.2015.0748</a>
- 71. Derex, M., & Boyd, R. (2016). Partial connectivity increases cultural accumulation within groups. *Proceedings of the National Academy of Sciences*, *113*(11), 2982–2987. <a href="https://doi.org/10.1073/pnas.1518798113">https://doi.org/10.1073/pnas.1518798113</a>
- 72. Derex, M., Perreault, C., & Boyd, R. (2018). Divide and conquer: Intermediate levels of population fragmentation maximize cultural accumulation. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 373(1743), 20170062. <a href="https://doi.org/10.1098/rstb.2017.0062">https://doi.org/10.1098/rstb.2017.0062</a>
- 73. Almaatouq, A., Noriega-Campero, A., Alotaibi, A., Krafft, P. M., Moussaid, M., & Pentland, A. (2020). Adaptive social networks promote the wisdom of crowds. *Proceedings of the National Academy of Sciences*, 201917687. <a href="https://doi.org/10.1073/pnas.1917687117">https://doi.org/10.1073/pnas.1917687117</a>

### **Supplemental Materials**

## I. Experiments 1 & 2: Comprehension Questions

After the test questions in Experiments 1 and 2, we asked two comprehension questions ("Comp\_TT" and "Comp\_AA") to test more explicitly whether participants were considering the effects of information sharing in a setting familiar to children. In these questions, Jack's teacher was giving a test to Jack's 5 informants, and participants were asked whether the 5 people should answer by Talking Together or by Answering Alone. In Comp\_TT, the teacher wanted "the 5 people to get as many answers right as possible"; in Comp\_AA, the teacher wanted to "find out which of the 5 people did their homework and which ones didn't". If children understand how discussion changes the informativeness of individual responses, they should recognize that Answering Alone is more informative to the teacher in Comp\_AA. If they understand the benefits of discussion (or at least, information sharing), they should prefer Talking Together for Comp\_TT.

In Experiment 1, children's responses to the comprehension questions suggest that even the youngest were able to choose a method of responding consistent with what the teacher wanted to learn about the students (Comp\_AA:  $M_{Young} = 65\%$ , p = .04,  $M_{Old} = 87.5\%$ , p < .0001,  $M_{Adult} = 92.5\%$ , p < .0001,  $M_{Young} = 70\%$ , p = .008,  $M_{Old} = 85\%$ , p < .0001,  $M_{Adult} = 87.5\%$ , p < .0001).

As in Experiment 1, responses to the comprehension questions at the end of Experiment 2 suggested even the youngest children recognized that talking together would make it impossible for the teacher to know who had done their homework (Comp\_AA: M<sub>Young</sub>= 67.5%, *p*=.019, M<sub>Old</sub>= 92.5%, *p*<.0001, M<sub>Adult</sub>= 90%, *p*<.0001). However, while older children and adults recognized that the students would do better on the test if they could discuss their answers, younger children were at chance (Comp\_TT: M<sub>Young</sub>= 52.5%, *p*=.4373, M<sub>Old</sub>= 90%, *p*<.0001, M<sub>Adult</sub>= 90%, *p*<.0001). Children in Experiment 2 may have been less confident in the value of discussion than their responses to the the main task questions in Experiments 1 and 2 would suggest; however, informal questioning of participants after the experiment suggested that younger children in Exp 2 may have simply rejected talking together on a test as cheating, even though the question specified that the teacher themselves could choose to allow students to talk together.

### **II. Supplementary Methods for Experiment 3**

**Norming Experiment.** In order to confirm the difficulty level of the Hard Percept and Easy Reasoning questions in Experiment 3, we first ran a norming experiment on MTurk with a separate group of 42 adult participants. Three participants were screened out for failing to answer basic comprehension questions about their job in the HIT.

We created 8 questions (4 Percept and 4 Reasoning) that we expected participants to rate as "easy" to answer and another 8 questions (4 Percept and 4 Reasoning) that we expected participants to rate as "hard" to answer. Each participant saw 8 questions: either the 4 Easy Reasoning and 4 Easy Percept questions, or the 4 Hard Reasoning and 4 Hard Percept questions. We expected the Hard Percept questions to be rated as more difficult to answer correctly than the Easy Reasoning questions. Each participant was asked "How difficult would it be to answer the question?", and rated the difficulty on a 7 point scale, from Extremely easy to Extremely difficult.

### The *Percept* questions:

**Photorealism:** decide which of two pictures of a face is a photo and which is a photorealistic drawing made by a talented artist. These materials adapted from Looser & Wheatley, 2010, which morphed faces using photographs and dolls as the anchors. We used Morph 3. The Easy version used Morph3\_052Human and Morph3\_067Human. The Hard version used Morph3\_063Human and Morph3\_065Human.

Intuitive Psychophysics (Superballs): decide how many marbles an opaque box contains by listening to it being shaken. This task was adapted from Siegel, Magid, Tenenbaum & Schulz, 2014. Two recordings were created. The Easy version asked whether the box contained 2 or 10 marbles (the recorded version contained 2). The Hard version asked whether the box contained 30 or 40 marbles (the recorded version contained 40).

<u>Brightness (Stars)</u>: decide which of the stars in a starry night sky looked the brightest. A picture of a starry night sky over a desert was used to represent the night sky, and the protagonist was said to have taken the picture so that he could "circle the brightest ones". In the Easy version, he wanted to circle the 3 brightest stars. The Hard version he wanted to circle the 25 brightest stars.

Rotation Speed: identify which of twelve colored diamonds is rotating the fastest. Each diamond had an A, a K, or a W in it to make the rotation clearer, but in the Hard version, the diamonds all had approximately the same RPM, while in the Easy version, the RPM was overall slower, and one was a clear outlier. The matrixes below show the number of rotations of each item in the Hard and Easy 4x3 arrays during the 10s display. In the Hard array, the fastest made 27 rotations in 10s, but 3 others made 26 and 2 made 25 rotations. In the Easy array, the fastest made 19 rotations in 10s, and the next closest made 12.

<u>Hard (# Rotations/10s)</u>	Easy (# Rotations/10s)
24 22 21 <u>27</u>	8 6 10 7
26 26 26 25	6 3 7 <u>19</u>
25 22 25 22	8 9 12 11

<u>The Reasoning questions:</u> The reasoning questions were adapted from Experiments 1 and 2.

<u>Sudoku</u>: Experiments 1 and 2 used a 4x4 sudoku problem rated as "easy" in a compilation, replacing the numbers with fruit to make it kid-friendly. The Easy version in Experiment 3 completed two additional moves. The Hard version used a 9x9 rated as "hard" in a compilation.

<u>Vehicle Routing Problem</u>: Experiments 1 and 2 used a custom made pathfinding puzzle which required a MarioKart find the shortest road through all the treasures on a map without taking "two in a row that are the same color, or two in a row that are the same shape". The Hard version used in these experiments had 11 treasures of different shapes and colors scattered randomly around the map. The Easy version created for Experiment 3 reduced the number of treasures to 4, of only 3 shapes and colors.

**Bottle-Jar Extraction Task**: Experiments 1 and 2 presented an "impossible object" puzzle, requiring the solver to remove a stick from a bottle without breaking the bottle or the stick. The stick was held fast inside the bottle by a nut-and-bolt. This was used as the Hard version. The Easy version substituted an analog of the "floating peanut" task (e.g., (Hanus, Mendes, Tennie, & Call, 2011), requiring the solver to remove a rubber ducky from large open-neck jar half-full of water, without touching the ducky or the jar, by pouring in the water from another jar.

Nim: In the game of Nim, each side takes turns picking up pencils. Each turn, you have to pick up either one, two, or three pencils. The winner is the person who picks up the last pencil. In Experiments 1 and 2, the we showed a game with only 5 pencils left. As adults and some older children found this 5-item version easy to solve, we created a Hard version by leave 22 pencils, and emphasizing that a wrong move would let a "super-smart computer" opponent win.

Norming Experiment: RESULTS. We fit a mixed effects model to perceived difficulty ratings, with random slopes and intercepts for each participant and question to account for repeated measures. The model confirmed that participants expected the Hard questions to be more difficult to answer than the Easy questions, ( $\beta$  = 1.95, SE = .5523, p = .0055). With the exception of the Easy version of the Percept\_Stars question, which was rated as significantly more difficult than other Easy questions ( $\beta$  = 2.45, SE = .0.4988, p < .0001), the questions within each difficulty level did not differ amongst themselves in perceived difficulty. Experiment 3 contrasted the Easy versions of the Reasoning questions with the Hard versions of the Percept questions; if participants preference for group reasoning in Experiments 2 and 3 was driven by the perceived difficulty of the question, then participants in Experiment 3 will favor group reasoning more for the *Hard Percept* questions than the *Easy Reasoning* questions.

# III. Cross-Experiment Exploratory Analyses

We conducted several exploratory analyses comparing results between experiments to examine the effects of crowd size and and question type more broadly. Experiment 1 and Experiment 2 used identical questions, but Experiment 2 increased the size of the crowd from 5 to 50 people. Our preregistered prediction was that participants would favor the crowd for population preference questions, but continue to favor the group for reasoning questions. However, we can also test the direct effect of crowd size by comparing people's judgments for reasoning and for popularity questions in Experiment 1 to their judgments in Experiment 2. Experiment 3 again used a crowd of 50 people, but contrasted easy versions of the reasoning questions from Experiments 1 and 2 with challenging perceptual discrimination tasks. This allowed us to test whether the preference for group discussion was caused by the perceived difficulty of the question. However, it also allows us to test whether the preference for

crowdsourcing observed in Experiments 1 and 2 extended to questions with a more ambiguous relationship to crowd size than population preferences.

To explore the effect of crowd size, we ran separate ANOVAs for each QuestionType using AgeGroup & Experiment as predictors (Exps 1 and 2). The tenfold increase in crowd size had no impact on participants' preference for discussing reasoning questions in small groups (F(1, 234)=0.045, p=.8320); an AgeGroup\*ExpNum interaction was significant (F(2,234)=4.434, p=.0129), but post-hoc comparisons revealed only a marginal difference between younger children's and adults' preference for reasoning in groups in Exp 1, but no other differences. However, participants were significantly more likely to crowdsource popularity questions in Experiment 2 than Experiment 1 (F(1, 234)=19.303, p<0.0001), with no differences between age groups.

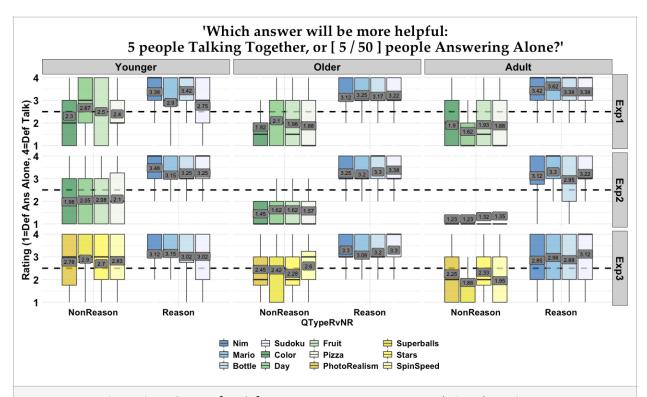
To explore whether the crowdsourcing preference was as strong for perceptual discrimination problems as population preference questions, we ran an ANOVA comparing the two types of non-reasoning questions, using AgeGroup & Experiment as predictors (Exps 2 and 3). Participants were significantly less confident that crowdsourcing would be preferable to a small group discussion for percept questions than popularity questions (F(1, 234)=76.897, p < 0.0001); the interaction was not significant (F(2,234)=0.139, p=.87). Notably however, there was no difference between participants' preference for asking a small group to discuss *Easy Reasoning* questions in Experiment 3 and *Reasoning* questions in Experiment 2, though it did approach significance (F(1, 234)=3.858, p < 0.0507).

# IV. MTurk Quality Screen

We present instructions as voice-over videos in order to prevent language bots from skimming the written text, and immediately after the videos, we simply present participants with 3 multiple choice questions about their task (*A: is their job to answer the questions themselves or decide which answer will help Jack more, B: do the people who answer alone talk together before each telling Jack their answer or not talk together, C: do the people who talk together each tell Jack their own answer after talking, or do they have to agree on a single answer to tell Jack after talking), with the correct answer being a nearly verbatim transcript from the video. Participants who get 1 or more of the attention check questions wrong have one more opportunity to answer after watching the video again;* 

if they get any questions wrong in the second round, they're blocked from taking the survey.

# V. Supplemental Plot for Exps 1-3



By-question boxplots for each of the 4 Reasoning questions (Blues) and 4 Non-Reasoning questions (Popularity - Greens; Percept - Yellows) in each experiment. Grey labels are means. (For preregistered analyses, average scores were computed for each QuestionType).

### VI. Mixed Effects Models

Our preregistered analysis plan was to compute an average score from the four questions of each QuestionType and conduct a repeated measure ANOVA on these two average scores. However, we also report mixed effects models; by including the unaveraged ratings for each question (i.e., the ratings on the 4-point scale for each of the four questions of each question type), these account for variance in the questions themselves. For each experiment, we tested the model (Ct\_Rating ~ 0+AgeGroup\*QuestionType + (1|subID), which models the responses for each of the 8

questions while treating AgeGroup and QuestionType as fixed effects, and allowing random intercepts for each subject. Centering individual ratings on 2.5 and deleting the intercept compares simple effect estimates to "chance" (i.e., 2.5 on a scale of 1 to 4) for each age group and estimates of interactions to the prior level's interaction, testing our predictions versus chance for the reference level of QuestionType and versus the magnitude of the previous age group's interaction for each interaction term; we report models with both Reasoning and Non-Reasoning questions coded as the reference level. These MEMs of raw ratings for each question produced qualitatively identical results to the repeated measures ANOVA on the averaged question ratings, with one exception: in Experiment 3, the mixed effect model suggested that while the youngest children favored group discussion for Non-Reasoning questions as well as Reasoning questions (consistent with the ANOVA), they also distinguished between the two (contrary to the ANOVA, where the difference was not significant), favoring discussion for Reasoning question more than for Non-Reasoning questions

(A) Exp 1: All age groups favored group discussion for Reasoning questions (βγουησετ = .6125, SE = .086, p = 9.33e-12; βοιdeτ = .69375, SE = .086, p = 2.15e-14; βAdult = .950, SE = .086, p < 2e-16), as well as making increasingly stronger distinctions between Reasoning and Non-Reasoning questions with age (βγουησετ = -.64375, SE = .10774, p = 3.40e-09; βοιdeτ = -.60625, SE = .15236, p = 7.52e-05; βAdult = .950, SE = .15236, p < 2.60e-10). Rerunning the regression with Non-Reasoning as the reference level showed that while younger children did not favor crowdsourcing for Non-Reasoning questions, older children and adults did (βγουησετ = -.03125, SE = .086, p = 0.717; βοιdeτ = -.55625, SE = .086, p = 4.62e-10; βAdult = -.66875, SE = .086, p < 1.47e-13)

(B) Exp 2: As in Exp 1, all age groups favored group discussion for Reasoning questions (βγοunger = .78125, SE = .086, p < 2e-16; βolder = .78125, SE = .086, p < 2e-16; βolder = .78125, SE = .086, p = 9.77e-13), as well as making increasingly stronger distinctions between Reasoning and Non-Reasoning questions with age (βγοunger = -1.23125, SE = .092, p < 2e-16; βolder = -.48125, SE = .130, p = 0.000232; βAdult = -.63750, SE = .130, p =1.16e-06). Rerunning the regression with Non-Reasoning as the reference level showed that all age groups favored crowdsourcing for Non-Reasoning questions (βγοunger = -.450, SE = .086, p < 2e-16); βAdult = -.1.21875, SE = .086, p < 2e-16).

(C) Exp 3: All age groups favored group discussion for Reasoning questions ( $\beta_{Younger} = .58125$ , SE = .096, p = 5.64e-09;  $\beta_{Older} = .71875$ , SE = .096, p = 1.45e-12;  $\beta_{Adult} = .45625$ , SE = .096, p = 3.56e-06), as well as making increasingly stronger distinctions between Reasoning and Non-Reasoning questions with age ( $\beta_{Younger} = .28125$ , SE = .106, p = 0.008342;  $\beta_{Older} = -.500$ , SE = .150, p = 0.000926;  $\beta_{Adult} = -.575$ , SE = .150, p = 0.000142). Rerunning the regression with Non-Reasoning as the reference level showed that while younger children preferred to discuss Non-Reasoning questions as well, older children had no preference, and adults preferred crowdsourcing Non-Reasoning questions ( $\beta_{Younger} = .300$ , SE = .096, p = 0.002019;  $\beta_{Older} = -.06250$ , SE = .096, p = 0.516040;  $\beta_{Adult} = -.400$ , SE = .096, p = 4.4e-05).

More complex random effects specification were overfit or failed to converge, but suggested little variance between questions themselves after accounting for the effect of QuestionType. For instance, Model 1 (below) allows for random intercepts of questions within the fixed effect of QuestionType, but fit was singular. Inspecting random effects suggested that the (1 | QuestionType:Question) explained no variance.

### Model 1:

Ct\_Rating ~ 0+QuestionType\*AgeGroup+(1 | subID)+(1 | QuestionType:Question)