

## Explaining Constrains Causal Learning in Childhood

Caren M. Walker

*University of California, San Diego*

Joseph J. Williams

*Harvard University*

Tania Lombrozo

*University of California, Berkeley*

Anna N. Rafferty

*Carleton College*

Alison Gopnik

*University of California, Berkeley*

Three experiments investigate how self-generated explanation influences children's causal learning. Five-year-olds ( $N = 114$ ) observed data consistent with two hypotheses and were prompted to *explain* or to *report* each observation. In Study 1, when making novel generalizations, explainers were more likely to favor the hypothesis that accounted for more observations. In Study 2, explainers favored a hypothesis that was consistent with prior knowledge. Study 3 pitted a hypothesis that accounted for more observations against a hypothesis consistent with prior knowledge. Explainers were more likely to base generalizations on prior knowledge. Findings suggest that attempts to explain drive children to evaluate hypotheses using features of "good" explanations, or those supporting generalizations with broad scope, as informed by children's prior knowledge and observations.

Since Piaget, researchers have commonly regarded children's explanations as a window into cognitive development, revealing their understanding of the world (Frazier, Gelman, & Wellman, 2009; Gopnik, 2000; Hickling & Wellman, 2001; Keil, 2006; Legare, Wellman, & Gellman, 2009). More recently, the very process of seeking, generating, and evaluating explanations has additionally been proposed as a powerful mechanism for learning and generalization, scaffolding knowledge acquisition, and contributing to theory change (e.g., Fonseca & Chi, 2010; Legare, 2012; Lombrozo, 2006, 2012; Wellman & Liu, 2007; Williams & Lombrozo, 2010, 2013). Here, we investigate the role of explanation in children's causal learning, focusing on whether and

how explaining influences how children evaluate new observations in light of their current theories.

Discovering the underlying causal structure in the world is one of the major inductive problems that learners face as they construct and revise early intuitive theories. Researchers have proposed that the acquisition of this causal knowledge is supported by powerful learning mechanisms that allow children (and adults) to effectively integrate novel evidence with prior beliefs (e.g., Gopnik et al., 2004; Griffiths, Sobel, Tenenbaum, & Gopnik, 2011). For example, 5-year-olds can implicitly track covariation between events to infer a novel causal relationship but require much stronger evidence to endorse a cause that conflicts with their prior beliefs (e.g., a psychological state causes a tummy ache) than one that is consistent with their prior beliefs (e.g., a particular food causes a tummy ache; Schulz, Bonawitz, & Griffiths, 2007; see also; Griffiths et al., 2011).

The integration of new evidence and prior beliefs can be naturally represented by normative Bayesian models (e.g., Griffiths, Kemp, & Tenenbaum, 2008). However, adults sometimes ignore prior probabilities in the face of compelling evidence (e.g.,

---

The research was funded by grants from the McDonnell Foundation to Alison Gopnik, the National Science Foundation (Grant DRL-1056712) to Tania Lombrozo, and the American Psychological Foundation (Elizabeth Munsterberg Koppitz) to Caren M. Walker. The authors would like to thank the parents and children who participated in this research. We are grateful to Christine Legare for her feedback on related collaborative work. We also thank the University of California, Berkeley Early Childhood Centers, the Center for Childhood Creativity at the Bay Area Discovery Museum, the Lawrence Hall of Science, and Habitat for facilitating recruitment. Finally, we thank Anna Akullien, Brynna Ledford, Sierra Eisen, and Rosie Aboody for assisting in data collection.

Correspondence concerning this article should be addressed to Caren M. Walker, Department of Psychology, University of California, San Diego, 9500 Gillman Drive, La Jolla, CA 92093-0109. Electronic mail may be sent to [carenwalker@ucsd.edu](mailto:carenwalker@ucsd.edu).

© 2016 The Authors

Child Development © 2016 Society for Research in Child Development, Inc. All rights reserved. 0009-3920/2016/xxxx-xxxx

DOI: 10.1111/cdev.12590

Kahneman & Tversky, 1973) and sometimes are overly reliant on prior beliefs (e.g., Chapman & Chapman, 1967; Wason, 1960). How might *explaining* their observations influence the relative contributions of evidence and prior knowledge in children's causal learning?

One possibility is that explaining is a process by which learners approximate Bayesian updating (i.e., compute a posterior probability), which is a function of both novel evidence and prior beliefs. Broadly consistent with this idea, the philosopher Peter Lipton suggests that "explanatory considerations may help enquirers to determine prior probabilities, to move from prior to posterior probabilities, and to determine which data are relevant to the hypothesis under investigation" (Lipton, 2001, p. 94). Going beyond Lipton's proposal to a descriptive claim, it could be that considering whether and how something can be explained contributes to the identification and integration of evidence and prior beliefs. As a result, explaining could improve learning by generating judgments that are more consistent with the results of Bayesian inference, relative to the judgments that would have been reached in the absence of explaining.

Although the relationship between children's explanations and Bayesian inference has never been tested directly, there are reasons to suspect that the process of explaining could influence children's sensitivity to both novel evidence and prior beliefs. For example, explaining could draw attention to anomalous observations (Legare, 2012; Legare, Gelman, & Wellman, 2010), thus making it more likely that prior beliefs inconsistent with that evidence will be revised (e.g., Brown & Kane, 1988; Rittle-Johnson, 2006; Siegler, 1995). On the other hand, explaining could encourage children to accommodate what they are trying to explain to what they already believe (e.g., Chi, de Leeuw, Chiu, & LaVancher, 1994; Kuhn & Katz, 2009; Lombrozo, 2006) and to "explain away" anomalous observations to preserve their current theory (Bonawitz, van Schijndel, Friel, & Schulz, 2012; Chinn & Brewer, 1993). How and why might explanation have these effects?

In the present article, we investigate a proposal about why explaining increases sensitivity to both evidence (Study 1) and prior beliefs (Study 2), with implications for how evidence and prior beliefs are negotiated when they come into conflict (Study 3). Specifically, we suggest that explaining recruits a set of evaluative criteria for what constitutes a *good* explanation. As a result, explaining could encourage learners to formulate and privilege hypotheses

that exhibit certain "explanatory virtues" (Lipton, 2001, 2004), hypotheses that they may not have otherwise considered.

The explanatory virtue on which we focus here is "scope," which we define as the amount of data that a candidate hypothesis explains (Khemlani, Sussman, & Oppenheimer, 2011; Lombrozo, 2012). Judgments of scope can be informed both by newly observed data (Williams & Lombrozo, 2010) and by past observations reflected in prior beliefs (Williams & Lombrozo, 2013). In general, hypotheses with greater scope also allow more general and wide-ranging inferences (Lombrozo & Carey, 2006), so scope is likely to be an especially powerful and important feature of explanations.

Research with adults supports the idea that explaining recruits scope as a constraint on learning. For example, Williams and Lombrozo (2010) had adult participants learn to classify novel robots into two categories by either *explaining* each robot's category membership or by engaging in a control task, such as thinking aloud. Compared to participants in other conditions, those prompted to explain were more likely to discover a pattern in category membership that accounted for 100% (as opposed to 75%) of the novel observations and thus had broader scope. Subsequent work (Williams & Lombrozo, 2013) found conditions under which prompting learners to explain made them more likely to identify and favor patterns that were consistent with prior beliefs, which have broader scope if one considers both current and past observations. These findings indicate that adult learners who are asked to explain privilege hypotheses with greater scope, including current and past observations (i.e., evidence and prior knowledge). If children are similarly sensitive to scope, we predict that a prompt to explain will similarly increase their sensitivity to each of these cues in the context of causal learning.

Recent developmental findings suggest that by age 5, children possess the basic prerequisites to test this proposal. For instance, we know that as early as 2 children have a sense of what counts as an explanation (Frazier et al., 2009), that by 4 they offer domain-appropriate explanations (Hickling & Wellman, 2001; Schulz et al., 2007), and that by 5 they prefer some kinds of explanations to others. For example, they prefer explanations that are simpler (Bonawitz & Lombrozo, 2012) and, at least in some domains, that offer a purpose or goal (Keil, 1995; Kelemen, 1999).

We also know that by age 5, children have developed abstract, coherent representations of

causal relationships in a variety of domains (e.g., Carey, 1985; Gelman & Wellman, 1991; Gopnik & Meltzoff, 1997; Inagaki & Hatano, 1993; Perner, 1991) and can reason successfully in a variety of causal inference tasks (e.g., Gopnik & Sobel, 2000). Children of this age are also capable of engaging in probabilistic reasoning (Gopnik et al., 2004; Griffiths, et al., 2004; Griffiths et al., 2011; Schulz & Gopnik, 2007; Schulz et al., 2007) and have developed the basic capacities for forming novel inferences based on covariation data (e.g., Kushnir, Xu, & Wellman, 2010)—two capacities that are needed for tracking and using information about the scope of a hypothesis.

Finally, two recent studies support the idea that prompting preschool-aged children to explain can actually change causal learning and inference relative to a control condition. First, Legare and Lombrozo (2014) presented 3- to 6-year-old children with a novel toy involving interconnected gears. They found that children who were prompted to explain how the toy worked (Study 1) or who spontaneously explained in response to a broader prompt (Study 2) outperformed their peers when it came to measures of causal mechanism learning (e.g., the shape of gears) but not on measures involving causally irrelevant perceptual details (e.g., the color of gears). Second, Walker, Lombrozo, Legare, and Gopnik (2014) found that when 3- to 5-year-olds were prompted to explain why various blocks did or did not activate a novel machine, they were more likely to focus on the internal properties and category membership of the blocks. These findings not only support the proposal that prompts to explain can systematically change learning and inference but additionally point to the idea that explaining directs learners to privilege certain types of hypotheses and aspects of their environment that are most likely to support future generalizations.

In sum, prior work is consistent with the idea that explaining constrains learning and inference, with evidence for effects of scope in adults. However, the evidence for these effects in young children is absent or indirect, and many questions remain unanswered. Scope is of special interest in young children not only because it is linked to explaining in adults (Williams & Lombrozo, 2010, 2013) but because it gets at the heart of learning as a function of both evidence and prior beliefs: consistency with present and past observations. We thus set out to address the following questions. First, when children's prior beliefs are held constant, does explaining make them more likely to favor hypotheses that account for more of the evidence they observe

(Study 1)? Second, when the evidence is held constant, does explaining make children more likely to favor hypotheses that are more consistent with prior beliefs (Study 2)? Third, does explaining influence how children balance evidence and prior beliefs when the two conflict (Study 3)?

In these experiments children learned about a novel causal system in which some objects generated an effect and others did not, where the objects were designed to support two candidate causal hypotheses. In Study 1, the two hypotheses were equally consistent with children's prior knowledge, so the only factor differentiating them was their consistency with the current evidence: One hypothesis accounted for the causal efficacy of *most* objects, whereas the other accounted for the causal efficacy of *all* objects. In Study 2, both hypotheses accounted for the causal efficacy of all objects, but one hypothesis was more consistent with children's prior knowledge. In Study 3, the hypothesis that accounted for fewer observations was also more consistent with prior knowledge, thereby pitting evidence against prior knowledge.

In all studies, half of the children were prompted to explain the presence or absence of the effect after each observation, and the other half were prompted to *report* the outcome. Reporting was selected as a control task because it shares several commonalities with explanation: It draws children's attention to the evidence, it requires them to verbalize in a social context, and it roughly matches children's time engaging with each outcome across conditions. In the test phase of all studies, children were then asked to generate predictions about the causal efficacy of novel objects.

These studies are designed to test the claim that explaining directs children to the hypothesis with the broadest scope (the one consistent with the greatest amount of current and past data). We predicted that children prompted to explain would be more likely than controls to favor the hypothesis that accounted for more evidence in Study 1 and the hypothesis more consistent with prior knowledge in Study 2. Study 3 allowed us to rule out alternative explanations for Study 1 and to investigate the extent to which explaining led children to balance evidence and prior beliefs in a manner consistent with Bayesian inference.

### Study 1

In Study 1, each observation involved a wooden block being placed on a "machine," where the

machine either did or did not play music (e.g., Gopnik & Sobel, 2000). Two features varied across the blocks: the color of a Lego affixed to the top and the color of a Lego affixed to the front. One of these features covaried perfectly with whether the machine played music across eight observations, whereas the other predicted six of the eight activations. These features were the basis for two candidate causal hypotheses about what led the machine to activate, which we refer to as the 100% color hypothesis and 75% color hypothesis, with the former being more consistent with the evidence. These proportions were selected following previous research conducted with adults (see Williams & Lombrozo, 2010) as well as developmental evidence that the difference in frequencies was large enough to be appreciated (Sobel, Tenenbaum, & Gopnik, 2004).

After children observed each of the blocks placed on the machine and either explained or reported the outcome, novel test blocks were placed inside a “hiding box” that obscured one of the two features. By asking children which of the test blocks would activate the machine, we were able to assess which hypothesis they favored in making novel causal judgments.

### Method

#### Participants

Forty-two 5-year-olds ( $M = 64.2$  months,  $SD = 3.6$ , range = 59.9–72.7; 25 girls) were included in Study 1, with 21 children randomly assigned to each of two conditions (*explain* and *control*). There was no significant difference in age between conditions, and there were approximately equal numbers of boys and girls assigned to each group. One additional child was tested, but excluded for failing to complete the study. Data were collected from January 2011 to June 2011.

Approximately half of the children tested were recruited from university preschools, and the other half from a local museum. Although we did not collect specific demographic information for each child, the following demographic information describes the population in each recruitment location. The preschools include the following racial/ethnic groups: 58% Caucasian, 15% Asian, 4% Native American or Alaskan Native, 2% Latino or Hispanic, and 21% mixed racial/ethnic background. Tuition for preschools ranged from \$15,000 to \$40,000 per year and the neighborhoods had a median household income ranging from \$80,000 to

\$131,000 per year. The museum visitors include the following racial/ethnic groups: 60% Caucasian, 28% Asian, 1% Native American or Alaskan Native, 14% Latino or Hispanic, 4% African American, and 13% mixed racial/ethnic background. The average income for museum visitors is between \$100,000 and \$150,000 per year.

#### Materials

*Machine.* The machine used in the training phase of all studies was similar to the “blicket detectors” used in past research on causal reasoning (e.g., Gopnik & Sobel, 2000). The machine consisted of a  $10 \times 6 \times 4$  in. opaque box constructed from cardboard and painted white with blue borders. The box contained a wireless doorbell that was not visible to the participant. When an object “activated” the machine, the doorbell played a melody. The machine was in fact surreptitiously activated by a remote control that was held out of view by the experimenter.

*Training blocks.* There were eight training blocks (four *causal* blocks and four *inert* blocks) made of 2 in. wooden cubes, illustrated in Figure 1 (top row). A plastic Lego plate was affixed to the top and front of each cube, and each Lego plate had one small, rectangular Lego. Thus, each training block had two Legos: one attached to the top and one to the front.

For half of the children, the feature corresponding to the 100% color hypothesis (the blue or yellow Lego) appeared on the top of the block and the feature corresponding to the 75% color hypothesis (the red or white Lego) appeared on the front of the block, and for half of the children, these positions were reversed. To create features consistent with the 100% color hypothesis, all four causal blocks (i.e., blocks that activated the machine) had a blue Lego on the top (front), and all four inert blocks (i.e., blocks that did not activate the machine) had a yellow Lego on the top (front). To create features consistent with the 75% color hypothesis, three of the four causal blocks had a red Lego on the front (top) and one had a white Lego, and three of the four inert blocks had a white Lego on the front (top) and one had a red Lego.

Small cards were constructed to serve as memory aids during the experiment. One card had an image of a black music note (placed in front of the causal objects, which activated the machine), and the other had an image of a black music note crossed out with a red “X” (placed in front of the inert objects, which did not activate the machine).

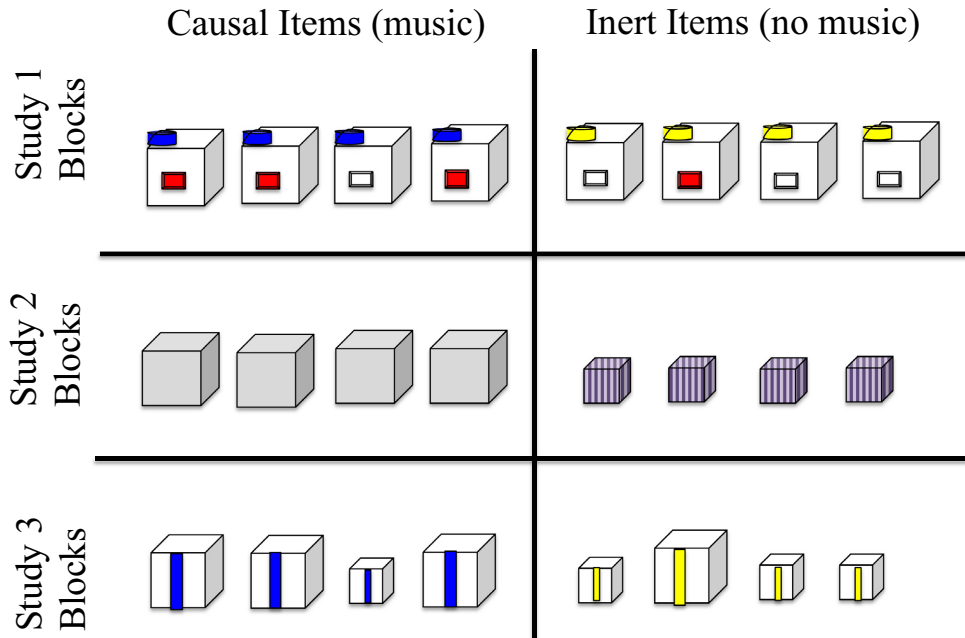


Figure 1. The complete set of eight training blocks used in Studies 1 (top row), 2 (center row), and 3 (bottom row). In Study 1, the 100% color hypothesis is represented by the blue (causal) and yellow (inert) Legos attached to the top of the blocks and the 75% color hypothesis is represented by the red (causal) and white (inert) Legos attached to the front of the blocks. The placement of the features corresponding to the 75% and 100% hypotheses were counterbalanced across participants. In Study 2, the 100% color hypothesis is represented by the solid silver (causal) and purple striped (inert) blocks and the 100% size hypothesis is represented by the relative size of the blocks (large [causal] and small [inert]). In Study 3, the 100% color hypothesis is represented by the blue (causal) and yellow (inert) bands on the blocks and the 75% hypothesis is represented by the relative size of the blocks (large [causal] and small [inert]).

*Test blocks.* The testing phase involved four additional blocks that were identical to the training blocks, but each included a single Lego attached to one side, with one block with each Lego color (blue, red, yellow, and white). These blocks could be partially obscured in a “hiding box” that was constructed from cardboard and painted black. The box included four cut-out windows that were covered with black felt flaps. In order to have children base their predictions on one cause (but not the other), this box was used to obscure parts of the test blocks. The experimenter could place two blocks inside the hiding box and lift two flaps to show a single Lego on each block, where the block position and lifted flaps determined whether the Lego appeared at the top or the front.

#### Procedure

*Training phase.* Following a warm-up period in which the child was familiarized with the experimenter, the machine was placed on the table. The experimenter said,

This is my machine. Some things make my machine play music and some things do not. We

are going to put all of the things that make my machine play music over here with this music note (experimenter places the *causal* memory card on one side of the table) and all of the things that do not make my machine play music over here with this crossed out music note (experimenter places the *inert* memory card on the other side of the table).

The experimenter then produced the first training block and placed it on the machine. After children observed the outcome, they were asked for a verbal response. In the explain condition, children were asked to explain the outcome: “Why did/didn’t this block make my machine play music?” In the control condition, children were asked to report the outcome: “What happened to my machine when I put this block on it? Did it play music?” Afterward, the experimenter asked the child to sort the object to one side of the table or the other (i.e., with the causal or inert memory card). This process was repeated for all eight training blocks.

The order in which the eight training blocks were presented was semirandom. The first four blocks included two causal and two inert blocks that were consistent with both the 100% and 75%

color hypotheses (i.e., blue and red blocks activate the machine and yellow and white blocks did not). The fifth and sixth blocks introduced the anomalous cases for both the causal and inert sets (i.e., a blue and white block activated the machine and a yellow and red block did not), which was consistent with the 100% color hypothesis, but not with the 75% color hypothesis. Finally, the seventh and eighth blocks were again consistent with both hypotheses (i.e., a blue and red block activated the machine and a yellow and white block did not).

To eliminate memory demands, all blocks remained visible and were grouped on the table with their corresponding causal or inert memory cards throughout the training and test phases of the experiment. Because we were interested in whether prompts to explain would influence children's inferences and generalization, we did not want working memory constraints to impact performance. (By providing these memory aids, it is possible that children would generate a hypothesis about the entire set of evidence after they have seen all blocks on the toy. Given that we are interested in whether explaining highlights the candidate cause with the broadest scope, it does not matter whether this inference is formed during the training trials or post hoc, in response to the test questions.)

*Test phase.* For the test phase, the machine was removed and the "hiding box" was placed on the table. The experimenter explained: "This is my hiding box. I am going to put two new blocks into my hiding box, and lift these flaps so you can only see part of each block." The experimenter demonstrated lifting each flap to familiarize the child with the apparatus. The child was then given the following instructions: "One of the blocks I put in my hiding box will make my machine play music, and one of the blocks will not. I want you to tell me which one you think *will* make my machine play music." Each test question was later repeated a second time in which the experimenter asked the child to indicate which block *would not* activate the machine.

In the first set of questions, the 100% *no conflict test items* and the 75% *no conflict test items*, the causal features corresponding to each hypothesis were pit against the inert features for that hypothesis (i.e., the 100% *no conflict test items* involved blue vs. yellow, and the 75% *no conflict test items* involved red vs. white). These questions ensured that children noted the covariation between color and the machine's activation corresponding to each hypothesis. In the next set of questions, the experimenter presented *conflict items*, in which the causal feature corresponding to the 100% color hypothesis

(blue) was pitted against the causal feature corresponding to the 75% color hypothesis (red). These questions presented a conflict between the two candidate hypotheses to examine which one children would favor in predicting a causal outcome. There was a total of six test questions in this format: four 100% and 75% *no conflict test items* (two blue vs. yellow, two red vs. white) and two *conflict items* (two blue vs. red), where one item of each type was asked in the positive format ("which *will* make my machine play music?") and one in the negative format ("which *will not* make my machine play music?").

*Coding.* Responses on test items were scored according to accuracy in tracking each property. For the 100% and 75% *no conflict test items*, accuracy reflected the correct selection of the block with a causal feature when asked for an item that would activate the machine, and the block with an inert feature when asked for an item that would not activate the machine. The two judgments of each type were averaged to create a single score corresponding to accuracy on the 100% *no conflict test items* and another for accuracy on the 75% *no conflict test items*.

For the *conflict items*, we coded children's predictions as "1" if they conformed to the 100% color hypothesis (selecting the blue Lego over the red Lego) or as "0" if they agreed with the 75% color hypothesis (selecting the reverse). This produced a *conflict items* score—the proportion of judgments (of two) consistent with the 100% color hypothesis.

Children's explanations were transcribed and coded from the videos. We analyzed the frequency with which different types of explanations were produced for each of the eight training blocks. All explanations were coded as belonging to one of four categories: (a) 75% color hypothesis (e.g., "It made the machine play because it has a red thing on it"), (b) 100% color hypothesis (e.g., "It made the machine play because it has a blue part"), (c) combined 75% and 100% color hypotheses (e.g., "It made the machine play because it has a blue and red one"), and (d) other/uninformative (e.g., "I don't know," "Because it is heavy," "Because it wants to").

Children's responses were recorded by a second researcher during the testing session, and most sessions were video recorded for independent coding by a third researcher who was naive to the hypotheses of the experiment. All available videos (95%) were independently coded to establish reliability. Interrater reliability was very high; the two coders agreed on 98% of children's responses to the

test questions and 96% of explanation codings. The few minor discrepancies were resolved by a third party.

*Results and Discussion*

*No Conflict Test Items*

To examine whether children noted the covariation between features and the machine’s activation, we conducted a repeated measures analysis of variance (ANOVA) with the two test item types (100% and 75% no conflict test items) as the repeated measure and condition (explain vs. control) as a between-subjects variable. This analysis did not reveal significant differences in children’s performance on the two question types,  $F(1, 40) = 0.195$ ,  $p = .661$ , nor across conditions,  $F(1, 40) = 1.84$ ,  $p = .182$ . The interaction between question type and condition was also not significant,  $F(1, 40) = 0.780$ ,  $p = .382$ .

Children in both the explain condition ( $M = 3.48$  of 4,  $SD = 0.92$ ) and the control condition ( $M = 3.08$  of 4,  $SD = 0.88$ ) learned about the 100% and 75% color hypotheses during training,  $t(20) = 7.29$ ,  $p < .0001$  and  $t(20) = 5.65$ ,  $p < .0001$ , respectively, and used this information when generalizing to novel blocks.

*Conflict Test Items*

To examine whether a prompt to explain influenced which causal hypothesis children favored, the conflict items score was analyzed with a one-way ANOVA with condition (explain vs. control) as the between-subjects variable (see Figure 2). There was a significant difference between conditions,  $F(1, 40) = 5.79$ ,  $p < .02$ ,  $\eta_p^2 = .13$ , with children in the explain condition more likely ( $M = 1.34$  of 2,  $SD = 0.66$ ) than children in the control condition ( $M = 0.98$  of 2,  $SD = 0.76$ ) to favor the 100% color hypothesis. These results support the claim that explaining helps young learners to discover and extend observed patterns that are consistent with the greatest number of cases (i.e., to privilege the hypothesis with broader scope).

*Qualitative Analysis of Explanations*

There was little variation in the types of explanations that children provided during the training phase. In fact, 80% of all explanations were categorized as combined 100% and 75% color hypotheses for training blocks 1–8, demonstrating that children were attending to the features relevant for both hypotheses. Given the limited variance, we do not provide additional analyses of the

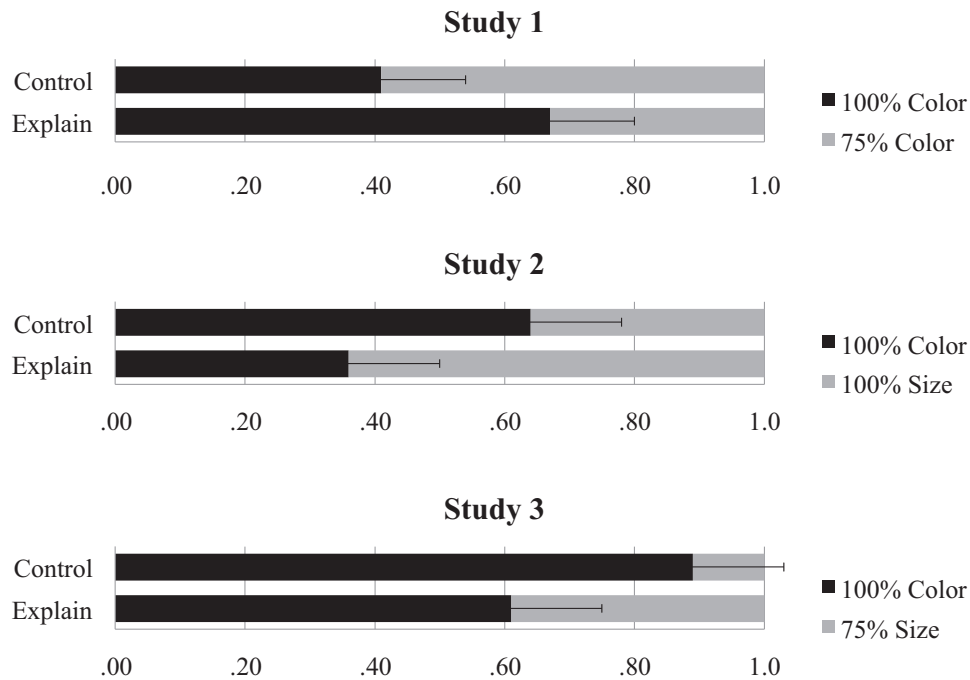


Figure 2. Mean proportion of responses on conflict items for Study 1 (100% color hypothesis vs. 75% color hypothesis), Study 2 (100% color hypothesis vs. 100% size hypothesis), and Study 3 (100% color hypothesis vs. 75% size hypothesis).

qualitative data for Study 1. Interestingly, 33% of children prompted to explain also spontaneously mentioned the weight of the objects, revealing the existence of prior beliefs about plausible causal mechanisms.

The findings from Study 1 suggest that explaining can increase children's responsiveness to evidence, and we propose that this is a consequence of privileging broader scope, in the sense that the hypothesized cause accounted for more observations. However, it is also possible that explaining drew attention to the presence of the single anomalous observation (Legare, 2012; Legare et al., 2010), penalizing the 75% rule without children ever engaging in a comparative assessment of scope. We return to this point in motivating Study 3.

## Study 2

Study 1 examined whether prompting children to explain could influence causal learning and inference, leading them to privilege a causal hypothesis that was more consistent with the evidence. The results suggest that it can: Children who explained were more likely than controls to favor a hypothesis that accounted for 100% of observations over a hypothesis that accounted for only a subset of observations (75%). This finding is consistent with the adult literature (Williams & Lombrozo, 2010). In Study 2, we consider whether explaining can also lead children to privilege hypotheses that are more consistent with prior knowledge. We therefore match the evidence for the two regularities, but make one more consistent with prior knowledge.

Study 2 again presented children with two novel hypotheses. However, both hypotheses (*color* and *size*) accounted for 100% of the data: Blocks of one color activated the machine 100% of the time, whereas those of another color always failed to do so, and larger blocks activated the machine 100% of the time while smaller blocks always failed to do so.

Size was selected as the feature for our second hypothesis because we anticipated that children would favor size as a more plausible causal factor in activating the machine than color. This expectation was informed by previous research examining children's beliefs about density (Esterly, 2000) and additionally verified during pilot testing, in which children often appealed to the weight of the objects, even when they did not vary in size. In addition, one third of the children prompted to explain in

Study 1 spontaneously mentioned the weight of the objects, even though this information did not differentiate the blocks. This suggests that children's appeal to weight was driven by prior beliefs and not by evidence from the task. Moreover, our assumption that children's prior beliefs favor weight as a candidate causal factor was confirmed in the computational model described in Study 3, which suggests that children assign the highest prior probability to the hypothesis that larger blocks activate the machine.

## Method

### Participants

Thirty-six 5-year-olds ( $M = 65.7$  months,  $SD = 5.1$ , range = 74.1–60.2; 20 girls) were included in Study 2, with 18 children randomly assigned to each of two conditions (explain and control). There was no significant difference in age between conditions, and there were approximately equal numbers of boys and girls assigned to each. Three additional children were tested but excluded due to experimenter error. Recruitment procedures and population demographics were equivalent to Study 1. Data were collected from February 2013 to October 2013.

### Materials

Study 2 used the same machine from Study 1. An illustration of the set of training blocks appears in Figure 1 (middle row). The causal blocks were made with four large (3 in.) wooden cubes and were painted metallic silver. The inert blocks were made with four small (1 in.) cubes and were covered in purple corduroy fabric. As in Study 1, small cards served as memory aids. Two additional blocks were used in the testing phase: a large block covered in purple corduroy and a small block painted silver. In place of the hiding box used in Study 1, test blocks were completely hidden in an opaque bag.

### Procedure

*Training phase.* The procedure for the training phase in Study 2 was identical to Study 1 but used the new set of eight training blocks.

*Test phase.* The procedure for the testing phase of Study 2 was similar to Study 1, with the following changes. Rather than placing the test objects in the hiding box, the experimenter looked



inside an opaque bag and described one feature for each of two new blocks, saying, for example, "I see a silver one and I see a purple one. Which one will [will not] make my machine play music, the silver one or the purple one?" As in Study 1, the first four *no conflict test items* contrasted the causal and inert features for each hypothesis (i.e., big vs. small, silver vs. purple) and the two *conflict test items* pit the causal features of one hypothesis against the other. For the conflict items, children were shown two blocks with a novel combination of features, a big purple block versus a small silver block, and asked to select the one that would (would not) activate the machine. This selection forced children to choose between size and color.

*Coding.* Response coding for Study 2 was identical to Study 1. To facilitate comparison across studies, *conflict items* were always coded in the same way: the *conflict items* score was equivalent to the proportion of judgments consistent with the 100% color hypothesis.

Explanations were coded as belonging to one of four categories: (a) color/texture (e.g., "It made the machine play because it is silver/smooth/shiny/sparkly"), (b) size/weight (e.g., "It made the machine play because it is big and heavy"), (c) insides/hidden properties (e.g., "It made the machine play because it has electricity inside of it"), and (d) other/uninformative (e.g., "I don't know," "Because it played music").

Seventy-eight percent of the videos were coded to assess reliability. Coders agreed on all but one child's response to test questions and 95% of explanations. Discrepancies were resolved by a third party.

## Results and Discussion

### No Conflict Test Items

To test whether children learned the covariation between the features (color and size) and the machine's activation, we conducted a one-way ANOVA to assess accuracy for the four *no conflict test items*, with condition (explain vs. control) as the between-subjects variable. The analysis revealed no difference between conditions,  $F(1, 34) = .92$ ,  $p = .345$ . As in Study 1, children in both the explain condition ( $M = 3.32$  of 4  $SD = 0.68$ ) and the control condition ( $M = 3.56$  of 4,  $SD = 0.72$ ) were able to track the relationship between the machine's activation and the features corresponding to the two hypotheses, and to use this information when

generalizing,  $t(17) = 8.25$ ,  $p < .0001$  and  $t(17) = 9.37$ ,  $p < .0001$ , respectively.

### Conflict Test Items

To analyze children's performance on the two *conflict test items*, a one-way ANOVA was conducted with condition (explain vs. control) as the between-subjects variable (see Figure 2). There was a significant difference between conditions,  $F(1, 34) = 4.46$ ,  $p < .05$ ,  $\eta_p^2 = .12$ , with children in the explain condition significantly less likely to choose the 100% color hypothesis ( $M = 0.72$  of 2,  $SD = 0.90$ ) than children in the control condition ( $M = 1.28$  of 2,  $SD = 0.90$ ). Instead, children who explained were more likely to choose the 100% size hypothesis. In other words, when the evidence was held constant (i.e., both hypotheses accounted for 100% of the data), children who explained were more likely than controls to privilege a hypothesis consistent with their prior knowledge (i.e., size).

These results provide further support for the proposal that explaining prompts children to favor hypotheses with broader scope. Not only does the search for explanations with broad scope direct children to favor a candidate cause that accounts for the greatest number of current observations (as shown in Study 1), it also leads them to favor a candidate cause that is more consistent with their *prior* observations, which are captured by their prior beliefs.

### Qualitative Analysis of Explanations

The majority of children's explanations appealed to size/weight (63% of all explanations), consistent with our expectations about children's prior beliefs. The 12 children who appealed to size/weight most often (i.e., as their modal response) were also more likely to privilege size in their responses to the conflict items ( $M = 1.42$  of 2,  $SD = 0.80$ ) than the remaining six children ( $M = 1.0$  of 2,  $SD = 1.05$ ), who provided other explanation types as their modal responses. However, perhaps due to the small sample sizes, this difference was not statistically significant,  $p = .36$ .

## Study 3

In Study 1, when prior knowledge was matched across hypotheses, explaining prompted children to select the hypothesis that accounted for a greater proportion of the evidence (i.e., the 100% color

hypothesis over the 75% color hypothesis). In Study 2, when the evidence was matched, explaining prompted children to select the hypothesis that was more consistent with their prior knowledge (i.e., the 100% size hypothesis over the 100% color hypothesis). Both findings are consistent with the idea that explaining prompts children to favor hypotheses with broader scope, whether scope is computed in terms of present observations (as in Study 1) or past observations as reflected in prior beliefs (as in Study 2). However, evidence and prior beliefs can come into conflict when the hypothesis that accounts for the most current evidence is not the one most consistent with prior beliefs. In Study 3, we investigate how prompting children to explain affects the balance between evidence and prior beliefs.

Study 3 had two additional aims: ruling out an alternative interpretation of Study 1 and providing data to compare against a Bayesian model. First, one criticism of the findings from Study 1 is that explaining could lead children to pay more attention to the presence of the single counterexample to the 75% rule and therefore to favor the 100% rule without evaluating scope as such. In Study 3, children were presented with two hypotheses: a 100% color hypothesis, which was consistent with a greater proportion of the evidence, and a 75% size hypothesis, which was more consistent with prior knowledge but involved a single counterexample, like the 75% color hypothesis from Study 1. If explaining simply draws attention to this counterexample, then children prompted to explain should favor the 100% color hypothesis over the 75% size hypothesis. In contrast, if explaining is related to scope—which is informed by both evidence and prior beliefs—then it is plausible that the single counterexample would not outweigh children’s antecedent commitments, and explaining will result in judgments that favor the 75% size hypothesis over the 100% color hypothesis.

Second, Study 3 also allowed us to investigate the correspondence between children’s judgments and the predictions of a Bayesian model. If explaining helps children approximate Bayesian inference, then we would expect children who explain to respond more like the model predictions than controls. In contrast, if explaining leads children to favor evidence or prior knowledge more than they “ought” to, we might expect controls to look more like the model predictions than those in the explain condition. By considering a case in which prior knowledge and evidence come into

conflict, we can more clearly differentiate these possibilities.

## Method

### Participants

Thirty-six 5-year-olds ( $M = 64.4$ ,  $SD = 3.8$ , range = 60.1–71.7; 20 girls) were included in Study 3, with 18 children randomly assigned to each condition (explain and control). There were no significant differences in age between the conditions, and there were approximately equal numbers of boys and girls in each. Four additional children were tested but excluded: two for failing to complete the study and two as a result of experimenter error. Recruitment procedures and population demographics were equivalent to Studies 1 and 2. Data were collected from June 2011 to November 2011.

### Materials

Study 3 used the same machine from Studies 1 and 2. An illustration of the complete set of training blocks appears in Figure 1 (bottom row). The causal blocks were made with three large (3 in.) wooden cubes and one small (1 in.) cube. The inert blocks were made with three small (1 in.) cubes and one large (3 in.) cube. For the 100% color hypothesis, a half-inch band of colored electrical tape was affixed to each of the eight blocks. The four causal blocks had a blue band and the four inert blocks had a yellow band. As in Study 2, test blocks were hidden in an opaque bag and several small cards served as memory aids. One additional large block with a blue band was used for the testing phase.

### Procedure

*Training phase.* The procedure for the training phase in Study 3 was identical to Studies 1 and 2 but used the new set of eight training blocks.

*Test phase.* The procedure for the testing phase of Study 3 was identical to Study 2.

*Coding.* Coding for Study 2 was identical to Studies 1 and 2. All explanations were coded as belonging to one of the same four categories identified in Study 2. Interrater reliability was very high; of the (75%) of videos coded, the two coders agreed on > 99% of the children’s responses to the test questions and 92% of the explanation coding. The few minor discrepancies were resolved by a third party.

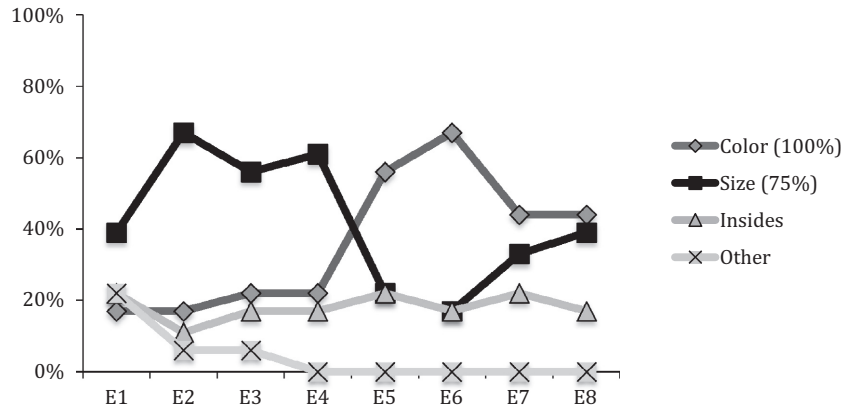


Figure 3. Proportion of explanations of each type (color, size, insides, other) provided for each of the eight training blocks by children in the explain condition in Study 3. Observed events are presented in order and indicated by E1–E8. The anomalous training blocks for both the causal and inert sets were presented in E5 and E6.

### Results and Discussion

#### No Conflict Test Items

To test whether children noted the covariation between features and the machine’s activation, we conducted a repeated measures ANOVA on accuracy for the test items, using question type (100% and 75% no conflict test items) as the repeated measure and condition (explain vs. control) as the between-subjects variables. The analysis revealed no main effect of condition,  $F(1, 34) = 1.81, p = .19$ , and no interaction between question type and condition,  $F(1, 34) = .05, p = .82$ . As in the previous studies, children in both the explain condition ( $M = 3.28$  of 4,  $SD = 0.63$ ) and the control condition ( $M = 3.67$  of 4,  $SD = 0.34$ ) were able to track the relationships between the machine’s activation and the features corresponding to each novel hypothesis and to use this information when generalizing to novel blocks with the same sets of features,  $t(17) = 4.30, p < .001$  and  $t(17) = 10.31, p < .0001$ , respectively.

#### Conflict Test Items

To analyze performance on the conflict items, a one-way ANOVA was conducted with condition (explain vs. control) as the between-subjects variable (see Figure 2). There was a significant difference between conditions,  $F(1, 34) = 6.64, p < .02$ ,  $\eta_p^2 = .16$ , with children in the explain condition significantly *less* likely to respond in line with the 100% color hypothesis ( $M = 1.22$  of 2,  $SD = 0.81$ ) than controls ( $M = 1.78$  of 2,  $SD = 0.43$ ).

These results provide evidence against the possibility that explaining simply prompts children to note counterexamples in the data they observe. They also weigh against a potential alternative interpretation of Studies 1 and 2 in which explanation is leading to greater overall engagement with the task, where greater engagement predicts responses that are more consistent with the observations they were prompted to explain.

#### Qualitative Analysis of Explanations

Figure 3 illustrates the frequency of each type of explanation provided for training blocks 1–8. The majority of explanations for the first four consistent blocks (ranging from 39% to 67% of the total number of explanations) focused on size, consistent with our expectations about children’s prior beliefs. Once children observed the anomalous events in blocks 5–6 (i.e., the small, blue causal block and the big, yellow inert block), 44% of all children changed their explanations to focus on the color of the blocks, consistent with the anomalous observations. As additional events consistent with prior knowledge were presented (events 7–8), the frequency of explanations that focused on size increased, while the frequency of explanations that focused on color decreased. These patterns are suggestive, but the sample sizes prevent a more systematic analysis. Interestingly, the minority of children ( $N = 4$ ) who began with an explanation that focused on insides or hidden properties (e.g., batteries, magnets, electricity) maintained this explanation when confronted with the anomalous training blocks. Because inside parts are invisible, they provide a

plausible explanation for all outcomes regardless of the size of the blocks, so there was no need for belief revision.

These qualitative data provide evidence that children's explanations were not simply determined by prior knowledge but rather informed by both observations and prior knowledge, even when these were in conflict. As in previous work (Walker et al., 2014), receiving a prompt to explain impacted children's inferences, even when the explanations that were generated did not appeal to the relevant causal property. Although generating the relevant explanation seems to be the most direct path to forming a corresponding inference, the process of explaining could itself yield other cognitive effects, such as promoting comparison and abstraction, that could lead children to the relevant inference by more indirect means. We return to this phenomenon in the General Discussion.

In sum, our qualitative data support the proposal that prompts to explain increase children's reliance on scope as a basis for inference and further suggest that effects of explanation are not restricted to children who happen upon the "correct" explanation for the task.

#### *Comparing Children's Performance to a Bayesian Model*

The findings reported above are consistent with more than one interpretation. One possibility is that explaining led children to integrate prior beliefs with novel evidence in a way that more closely approximated Bayesian inference. A second possibility is that explaining made children *less* likely to accurately approximate Bayesian inference, instead leading them to privilege prior beliefs more than they "should" have according to Bayes' rule.

Bayesian inference provides a formal account of how a learner should update her prior belief in some hypothesis,  $h$ , in light of new evidence,  $d$  (e.g., Gopnik et al., 2004; Griffiths et al., 2011). Specifically, the learner evaluates the posterior probability of the hypothesis,  $p(h|d)$ , by applying Bayes' rule:  $p(h|d) = p(h) \times p(d|h)/p(d)$ , where  $p(d|h)$  is the "likelihood" of the data given the hypothesis, and  $p(d)$  is the probability of the data under all hypotheses in question,  $h$ , and alternatives,  $\sim h$ . With Bayesian inference as a normative standard against which to assess children's performance in Study 3, we can ask whether children who explained responded in line with the Bayesian posterior *more* often or *less* often than those in the control condition.

We present the details of our model in the Supporting Information. Briefly, we first conducted a

behavioral study in which an additional eighteen 5-year-olds observed the same data as participants in Study 3 but generated a *prediction* about whether each block would activate the machine before seeing it placed on top. We assumed that children assigned some prior probability  $p$  to the hypothesis that large blocks activate the machine, with the remaining probability divided equally between the hypotheses that small, blue, or yellow blocks activate the machine. We additionally assumed that the system was somewhat noisy, with a probability  $\epsilon$  that an observation would depart from the actual rule (e.g., that a small block would activate the machine even though the correct hypothesis was that large blocks activate the machine).

The maximum likelihood estimates for these parameters, given the eight distinct prediction judgments from children during the task, corresponded to a prior probability that large blocks activate the machine of 51% and a "noise" parameter of 6%. When these numbers were updated by an ideal Bayesian learner in light of the actual observations from Study 3, the result was a posterior probability of 98.7% favoring the 100% color hypothesis, not the 75% size hypothesis. This number is much closer to the judgments observed for children in the control condition (83% favoring the 100% color hypothesis) than those in the explain condition (64% favoring the 100% color hypothesis), suggesting that the prompt to explain made children respond *less* normatively, overweighting their prior beliefs relative to the strength of the evidence.

### General Discussion

In three studies, prompting children to explain influenced the hypothesis they privileged when generalizing a causal relationship with novel cases. In Study 1, one hypothesis accounted for a greater number of observations. In this case, children who explained were more likely than controls to generalize according to the hypothesis that accounted for more of their observations. In Study 2, both hypotheses accounted for all observations, but one was more consistent with children's prior beliefs. In this case, children who explained favored the hypothesis that was consistent with their prior beliefs more often than those in the control condition. Finally, in Study 3, a candidate cause that accounted for all observations was pitted against an alternative that accounted for fewer observations but was consistent with prior knowledge. When presented with this conflict between current

evidence and prior commitments, children who explained were *less* likely than controls to favor the hypothesis that accounted for more of the observed events, instead making judgments consistent with prior beliefs more often than children in the control condition. The results of Study 3 provide evidence against the possibility that explaining simply increases overall attention to the task, responsiveness to the evidence alone, or sensitivity to the presence of a counterexample. The modeling results from Study 3 (see the Supporting Information) additionally suggest that explaining did not make children “more Bayesian.”

Taken together, these studies shed light on the mechanisms by which explanation informs and constrains causal learning in early childhood and help explain the conflicting results of earlier studies. In particular, explanation leads children to consider hypotheses that capture the explanatory virtue of broad scope, which results in (at least) two distinct effects on learning: Explanation can make learners more sensitive to evidence (Study 1) or more likely to rely on prior beliefs (Study 2). Depending on which effect dominates, explanation can lead to either an increase (e.g., Brown & Kane, 1988; Rittle-Johnson, 2006; Siegler, 1995) or a decrease (Bonawitz et al., 2012; Chi et al., 1994; Chinn & Brewer, 1993; Lombrozo, 2006) in belief revision, relative to children who do not explain.

Second, our findings can help us understand how the balance between evidence and prior beliefs is negotiated when the two conflict. In Study 3, we found that children who explained privileged consistency with prior knowledge over consistency with the data. Presumably, children favored the size hypothesis initially because it fit with a general principle that larger (or heavier) objects have a stronger causal impact than smaller ones, a principle that applies to many folk physical cases, particularly those involving contact between objects (in cartoons, it is an anvil that flattens roadrunner not a thumbtack). Children may have maintained the size hypothesis in the face of exceptions or reverted to it when observations were once again consistent with this idea, not only because of its broad scope but also because it offered other explanatory virtues, including a sense for the causal mechanism (Lipton, 2000). Study 3, then, may have pit one explanatory virtue (scope) against several (scope plus a sense of mechanism), with the result that explaining tipped the balance toward prior beliefs over evidence.

Although many questions about the role of explanation in early learning remain open, our

findings do provide evidence *against* several possibilities. First, given the results of Study 1, it cannot be that the effects of explanation are restricted to accommodating new observations in the context of prior knowledge, as children who explained were also more likely to generalize even when prior beliefs were held constant. Second, Studies 2 and 3 rule out the possibility that explaining always leads children to greater responsiveness to the evidence or to notice and respond to counterexamples. Third, given the results of Study 3, we can also rule out the possibility that explanation leads to a uniform boost in attention to the observations being explained (see also, Walker et al., 2014), as that would predict judgments more consistent with the evidence. Finally, one interesting question introduced in the discussion of Study 3 is whether explanation causes children to become more or less optimally Bayesian. Our model and analysis found that children who were prompted to explain were *less* likely to conform to the model predictions than children in the control condition. That is, children prompted to explain seemed to maintain prior beliefs more strongly than they should have, given their priors and the evidence observed.

Although explaining may not increase fidelity to Bayesian conditionalization, it may nonetheless be possible to provide a formal account of explanation’s effects in Bayesian terms. In the spirit of the quote from Peter Lipton in the introduction, “explanatory considerations” could influence how Bayesian inference is approximated, if not always leading to greater accuracy. In particular, recent work has explored the idea that, at an algorithmic level, both children and adults approximate ideal Bayesian inference by “sampling” procedures. In these procedures, learners generate a few hypotheses to test at a time, adjusting the probabilities of those hypotheses as they acquire more data (Bonawitz et al. 2014a, 2014b). Explaining could influence how this sampling process occurs, especially at the stage of hypothesis *generation* (e.g., Bonawitz & Griffiths, 2010; Ullman, Goodman, & Tenenbaum, 2012). For example, Bonawitz and Griffiths (2010) designed a causal learning experiment in which half of the participants were primed with the correct causal rule, and the other half were given a neutral prime. Priming participants changed the probability with which they *generated* those hypotheses when asked to describe the rule that best captured the evidence but did not influence their *evaluation* of hypotheses that were subsequently provided.

Another possibility is that children who explain *are* in fact responding more normatively but that

Bayesian inference is not the right standard against which to assess their performance. For instance, Douven and Schupbach (2015) suggest that Bayesian inference may be used when the goal is to minimize expected inaccuracy in the long run but that probabilistic versions of “inference to the best explanation”—which use explanatory considerations as a guide to inference, as we have suggested—could be appropriate when the goal is to get things “mostly right” in the short term. Another possibility is that explaining involves considerations of utility. Errors of overgeneralization may be preferable to errors of undergeneralization (Williams, Lombrozo, & Rehder, 2013), for instance. Finally, it is worth noting that a more complex model than the one we report could reveal a different picture. We hope that our initial steps toward a formal analysis of children’s explanation-based reasoning prompts further study.

#### *Additional Questions for Future Research*

The interpretation outlined thus far has focused primarily on the impact of explanation on the formation of children’s causal beliefs. However, it is also worthwhile to consider the performance of children who were *not* prompted to explain. In particular, why did the children in the control condition in Study 2 not spontaneously consult their prior knowledge? Williams and Lombrozo (2013) report similar findings in adults: Learners who explained were more likely to be influenced by labels that cued prior knowledge. Williams and Lombrozo suggest that explanation can guide learners to consult prior knowledge that would otherwise remain inert or underutilized. These results can also be interpreted consistent with Rozenblit and Keil’s (2002) “illusion of explanatory depth,” or the bias to overestimate one’s own explanatory understanding of causal mechanisms (e.g., how a bicycle works), which has also been found in young children (Mills & Keil, 2004). If children erroneously believe that they already possess an adequate explanation, they may not feel it necessary to explain presented observations and therefore fail to capitalize on the resources that are recruited by explaining.

Another possibility, of course, is that explaining may not always be beneficial. In fact, Study 3 suggests that under some conditions, explaining may result in causal inferences that are less normative. Relatedly, other research has suggested that explanation can have associated costs: Children prompted to explain why blocks activate a machine

are less likely to remember superficial properties of each block than are those in a control condition (Walker et al., 2014), and children prompted to explain how a gear toy works are less likely than controls to remember the colors of particular gears (Legare & Lombrozo, 2014). With adults, Williams et al. (2013) report cases in which prompting adults to explain can impair learning by leading to errors of overgeneralization. Finally, it is also possible that explanation may lead to verbal overshadowing (Schooler & Engstler-Schooler, 1990), in that generating an explanation could interfere with nonverbal processing, such as perceptual memory. Although we think that this phenomenon is unlikely to account for the present findings, as children’s inferences did not depend on accurate visual memory (i.e., children were provided with memory aids), verbal overshadowing could play an important role in explanation’s effects more generally.

There are also interesting open questions regarding the role of anomalous data in the current experiments. In particular, we used a fixed pattern of data, in which children were introduced to the anomalous observations midway (in Trials 5 and 6). This method of presentation was chosen to avoid order effects found in causal learning tasks (e.g., Abbott & Griffiths, 2011). However, it is possible that the placement of anomalous data may have influenced responses. Previous research has found that the order of explanations has an effect on learning (Ihme & Wittwer, 2015): Adults tend to prefer the first explanation provided. It is therefore reasonable to assume that the order of presentation—and in particular the relative placement of anomalies—could influence children’s causal learning as well.

Open questions also remain regarding the development of explanation as a mechanism for learning over time. Although the reported results demonstrate the presence of an early effect of explanation, we are unable to make claims regarding the developmental progression of these particular effects before 5 years of age. There do appear to be developments in children’s ability to engage in diagnostic inference in the face of uncertainty in the preschool period (Fernbach, Macris, & Sobel, 2012). However, previous research indicates uniform effects of explanation between 3 and 5 years (Walker et al., 2014). Our findings are also broadly consistent with previous research on the effect of explanation on adult category learning (Williams & Lombrozo, 2010, 2013)—A potentially surprising correspondence given children’s comparatively impoverished language skills, immature metacognitive abilities, and lower levels of prior knowledge.

Future research should also explicitly consider how these findings relate to previous proposals regarding the role of explanation for learning. Much of the evidence for the benefits of explanation comes from research on the “self-explanation effect,” the finding from educational psychology that prompting students to explain can improve learning (e.g., Fonseca & Chi, 2010). Researchers have proposed a variety of plausible mechanisms that could underlie the effect. For example, Siegler (2002) suggests (among other things) that one consequence of explaining is a general increase in engagement, and several researchers have suggested that explanations invoke prior beliefs (e.g., Ahn, Brewer, & Mooney, 1992; Chi et al., 1994; Lombrozo, 2006). Additional proposals include the ideas that explaining improves metacognitive monitoring, encourages learners to draw novel inferences, and helps form effective procedures (e.g., Chi, Bassok, Lewis, Reimann, & Glaser, 1989; Fonseca & Chi, 2010; Johnson-Laird, Girotto, & Legrenzi, 2004; Siegler, 2002).

The current work builds upon these accounts by demonstrating how explaining can moderate the relative contributions of prior beliefs and novel observations. Considering these mechanisms in conjunction can also help explain why we observe reliable effects of explaining even when the content of children’s explanations do not always map onto their judgments. For example, although it is clear why producing an explanation that refers to a block’s size would lead children to generalize according to size, it is much less clear why producing an explanation that appeals to a block’s color would lead to this same generalization. Previous research by Chi and colleagues (1994) suggest that incorrect explanations can make a (potentially implicit) belief explicit, thereby highlighting conflicts between the content of the explanation and the evidence. By noting the mismatch between an explanation and the data, children may be in a better position to reject an incorrect hypothesis and seek better alternatives.

Consistent with this idea, Wilkenfeld and Lombrozo (2015) propose that the process of explaining can have epistemic value, regardless of whether it results in a correct explanation—an idea that they dub “explaining for the best inference,” in contrast to the more typical “inference to the best explanation.” In addition to the potential metacognitive benefits noted above, the process of explaining could encourage typically beneficial cognitive processes, such as comparison and abstraction (e.g., Edwards, Williams, Lombrozo, & Gentner, under

review; Williams & Lombrozo, 2010). For example, in a study by Walker et al. (2014), children prompted to explain in a causal learning task exhibited more mature patterns of inference—privileging inductively rich, but hidden features over perceptually salient ones—even when the content of their explanations focused on perceptual features.

### Conclusion

In three experiments, we provide evidence for the role of explaining in guiding causal learning in early childhood. Our findings support the idea that generating explanations prompts young learners to favor hypotheses with broad scope, where assessments of scope are informed both by current evidence and by prior beliefs. Although the current findings contribute to our understanding of the role of explanation for learning, in particular, they also shed light on the nature of learning in general. When learning by explaining, the learner gains “new” knowledge by engaging with information that she already has. This phenomenon of “learning by thinking” (Lombrozo & Walker, in prep) challenges a simple data-driven view of knowledge acquisition, in which learning is simply a function of observations and testimony. Instead, these findings provide evidence for a more complex picture, one in which processes such as explaining to oneself—which does not involve new data—influence how the data and currently held theories inform judgments. Understanding how engaging in explanation influences early learning therefore contributes to a more complete understanding of how knowledge is acquired and revised.

### References

- Abbott, J. T., & Griffiths, T. L. (2011). Exploring the influence of particle filter parameters on order effects in causal learning. In L. A. Carlson, C. Hölscher, & T. F. Shipley (Ed.), *Proceedings of the 33rd annual conference of the Cognitive Science Society* (pp. 2950–2956). Austin, TX: Cognitive Science Society.
- Ahn, W. K., Brewer, W. F., & Mooney, R. J. (1992). Schema acquisition from a single example. *Journal of Experimental Psychology: Learning, Memory, & Cognition*, *18*, 391–412. doi:10.1037/0278-7393.18.2.391
- Bonawitz, E. B., Denison, S., Gopnik, A., & Griffiths, T. (2014). Win-Stay, Lose-Sample: A simple sequential algorithm for approximating Bayesian inference. *Cognitive psychology*, *74*, 35–65. doi:10.1016/j.cogpsych.2014.06.003

- Bonawitz, E., Denison, S., Griffiths, T. L., & Gopnik, A. (2014). Probabilistic models, learning algorithms, and response variability: sampling in cognitive development. *Trends in cognitive sciences*, 18(10), 497–500. doi:10.1016/j.tics.2014.06.006
- Bonawitz, E. B., & Griffiths, T. (2010). Deconfounding hypothesis generation and evaluation in Bayesian models. In N. Miyake, D. Peebles, & R. P. Cooper (Ed.), *Proceedings of the 32nd annual conference of the Cognitive Science Society* (pp. 2260–2265). Austin, TX: Cognitive Science Society.
- Bonawitz, E. B., & Lombrozo, T. (2012). Occam's rattle: Children's use of simplicity and probability to constrain inference. *Developmental Psychology*, 48, 1156–1164. doi:10.1037/a0026471
- Bonawitz, E. B., van Schijndel, T. J. P., Friel, D., & Schulz, L. (2012). Children balance theories and evidence in exploration, explanation, and learning. *Cognitive Psychology*, 64, 215–234. doi:10.1016/j.cogpsych.2011.12.002
- Brown, A. L., & Kane, M. J. (1988). Preschool children can learn to transfer: Learning to learn and learning from example. *Cognitive Psychology*, 20, 493–523. doi:10.1016/0010-0285(88)90014-X
- Carey, S. (1985). *Conceptual change in childhood*. Cambridge, MA: MIT Press/Bradford Books.
- Chapman, L. J., & Chapman, J. P. (1967). Genesis of popular but erroneous psycho-diagnostic observations. *Journal of Abnormal Psychology*, 72, 193–204. doi:10.1037/h0024670
- Chi, M. T. H., Bassok, M., Lewis, M., Reimann, P., & Glaser, R. (1989). Self-explanations: How students study and use examples in learning to solve problems. *Cognitive Science*, 13, 145–182. doi:10.1207/s15516709cog1302\_1
- Chi, M. T. H., de Leeuw, N., Chiu, M. H., & LaVancher, C. (1994). Eliciting self-explanations improves understanding. *Cognitive Science*, 18, 439–477. doi:10.1207/s15516709cog1803\_3
- Chinn, C. A., & Brewer, W. F. (1993). The role of anomalous data in knowledge acquisition: A theoretical framework and implications for science instruction. *Review of Educational Research*, 63, 1–49. doi:10.3102/00346543063001001
- Douven, I., & Schupbach, J. N. (2015). Probabilistic alternatives to Bayesianism: The case of explanationism. *Frontiers in Psychology*, 6, 459. doi:10.3389/fpsyg.2015.00459
- Edwards, B. J., Williams, J. J., Lombrozo, T., & Gentner, D. (under review). Explanation recruits comparison: Insights from a category-learning task.
- Esterly, J. B. (2000). "Its size is really wood": The development of buoyancy and material kind understanding in children between three and seven years of age. *Dissertation Abstracts International: Section B: The Sciences and Engineering*, 61, 3878.
- Fernbach, P., Macris, D. M., & Sobel, D. M. (2012). Which one made it go?: The emergence of diagnostic reasoning in preschoolers. *Cognitive Development*, 27, 39–53. doi:10.1016/j.cogdev.2011.10.002
- Fonseca, B., & Chi, M. T. H. (2010). The self-explanation effect: A constructive learning activity. In R. Mayer & P. Alexander (Eds.), *Handbook of research on learning and instruction* (pp. 296–321). New York, NY: Routledge Press. doi:10.4324/9780203839089
- Frazier, B. N., Gelman, S. A., & Wellman, H. M. (2009). Preschoolers' search for explanatory information within adult-child conversation. *Child Development*, 80, 1592–1611. doi:10.1111/j.1467-8624.2009.01356.x
- Gelman, S. A., & Wellman, H. M. (1991). Insides and essence: Early understandings of the non-obvious. *Cognition*, 38, 213–244. doi:10.1016/0010-0277(91)90007-Q
- Gopnik, A. (2000). Explanation as orgasm and the drive for causal knowledge: The function, evolution, and phenomenology of the theory formation system. In F. Keil & R. Wilson (Eds.), *Explanation and cognition* (pp. 299–323). Cambridge, MA: MIT Press.
- Gopnik, A., Glymour, C., Sobel, D., Schulz, L., Kushnir, T., & Danks, D. (2004). A theory of causal learning in children: Causal maps and Bayes nets. *Psychological Review*, 111, 1–31. doi:10.1037/0033-295X.111.1.3
- Gopnik, A., & Meltzoff, A. N. (1997). *Words, thoughts and theories*. Cambridge, MA: MIT Press. doi:10.5860/CHOICE.34-6547
- Gopnik, A., & Sobel, D. (2000). Detecting blickets: How young children use information about novel causal powers in categorization and induction. *Child Development*, 71, 1205–1222. doi:10.1111/1467-8624.00224
- Griffiths, T. L., Kemp, C., & Tenenbaum, J. B. (2008). Bayesian models of cognition. In R. Sun (Ed.), *Cambridge handbook of computational psychology* (pp. 59–100). New York, NY: Cambridge University Press. doi:10.1017/CBO9780511816772.006
- Griffiths, T. L., Sobel, D., Tenenbaum, J. B., & Gopnik, A. (2011). Bayes and blickets: Effects of knowledge on causal induction in children and adults. *Cognitive Science*, 35, 1407–1455. doi:10.1111/j.1551-6709.2011.01203.x
- Hickling, A. K., & Wellman, H. M. (2001). The emergence of children's causal explanations and theories: Evidence from everyday conversation. *Developmental Psychology*, 37, 668–683. doi:10.1037/0012-1649.37.5.668
- Ihme, N., & Wittwer, J. (2015). The role of consistency, order, and structure in evaluating and comprehending competing scientific explanations. *Instructional Science*, 43(4), 507–526. doi:10.1007/s11251-015-9349-6
- Inagaki, K., & Hatano, G. (1993). Young children's understanding of the mind-body distinction. *Child Development*, 64, 1534–1549. doi:10.2307/1131551
- Johnson-Laird, P. N., Girotto, V., & Legrenzi, P. (2004). Reasoning from inconsistency to consistency. *Psychological Review*, 111, 640–661. doi:10.1037/0033-295X.111.3.640
- Kahneman, D., & Tversky, A. (1973). On the psychology of prediction. *Psychological Review*, 80, 237–251. doi:10.1037/h0034747
- Keil, F. (1995). The growth of causal understandings of natural kinds. In D. Sperber, D. Premack, & A. J.



- Premack (Eds.), *Causal cognition: A multi-disciplinary debate* (pp. 234–262). Oxford: Clarendon Press. doi:10.1093/acprof:oso/9780198524021.003.0009
- Keil, F. (2006). Explanation and understanding. *Annual Review of Psychology*, 57, 227–254. doi:10.1146/annurev.psych.57.102904.190100
- Kelemen, D. (1999). Functions, goals, and intentions: Children's teleological reasoning about objects. *Trends in Cognitive Sciences*, 12, 461–468. doi:10.1016/S1364-6613(99)01402-3
- Khemlani, S., Sussman, A., & Oppenheimer, D. (2011). Harry Potter and the sorcerer's scope: Latent scope biases in explanatory reasoning. *Memory & Cognition*, 39, 527–535. doi:10.3758/s13421-010-0028-1
- Kuhn, D., & Katz, J. (2009). Are self-explanations always beneficial? *Journal of Experimental Child Psychology*, 103, 386–394. doi:10.1016/j.jecp.2009.03.003
- Kushnir, T., Xu, F., & Wellman, H. M. (2010). Young children use statistical sampling to infer the preferences of others. *Psychological Science*, 21, 1134–1140. doi:10.1177/0956797610376652
- Legare, C. H. (2012). Exploring explanation: Explaining inconsistent information guides hypothesis-testing behavior in young children. *Child Development*, 83, 173–185. doi:10.1111/j.1467-8624.2011.01691.x
- Legare, C. H., Gelman, S. A., & Wellman, H. M. (2010). Inconsistency with prior knowledge triggers children's causal explanatory reasoning. *Child Development*, 81, 929–944. doi:10.1111/j.1467-8624.2010.01443.x
- Legare, C. H., & Lombrozo, T. (2014). The selective benefits of explanation on learning in early childhood. *Journal of Experimental Child Psychology*, 126, 198–212. doi:10.1016/j.jecp.2014.03.001
- Legare, C. H., Wellman, H. M., & Gellman, S. A. (2009). Evidence for an explanation advantage in naïve biological reasoning. *Cognitive Psychology*, 58, 177–194. doi:10.1016/j.cogpsych.2008.06.002
- Lipton, P. (2000). Inference to the best explanation. In W. H. Newton-Smith (Ed.), *A companion to the philosophy of science* (pp. 184–193). London, England: Blackwell.
- Lipton, P. (2001). Is explanation a guide to inference? In G. Hon & S. Rackover (Eds.), *Explanation: Theoretical approaches* (pp. 93–120). Dordrecht: Kluwer. doi:10.1007/978-94-015-9731-9\_4
- Lipton, P. (2004). *Inference to the best explanation*. London: Routledge.
- Lombrozo, T. L. (2006). The structure and function of explanations. *Trends in Cognitive Sciences*, 10, 464–470. doi:10.1016/j.tics.2006.08.004
- Lombrozo, T. L. (2012). Explanation and abductive inference. In K. J. Holyoak & R. G. Morrison (Eds.), *Oxford handbook of thinking and reasoning* (pp. 260–276). Oxford, UK: Oxford University Press. doi:10.1093/oxfordhb/9780199734689.013.0014
- Lombrozo, T., & Carey, S. (2006). Functional explanation and the function of explanation. *Cognition*, 99, 167–204. doi:10.1016/j.cognition.2004.12.009
- Mills, C., & Keil, F. C. (2004). Knowing the limits of one's understanding: The development of an awareness of an illusion of explanatory depth. *Journal of Experimental Child Psychology*, 87, 1–32. doi:10.1016/j.jecp.2003.09.003
- Perner, J. (1991). *Understanding the representational mind*. Cambridge, MA: MIT Press.
- Rittle-Johnson, B. (2006). Promoting transfer: The effects of direct instruction and self-explanation. *Child Development*, 77, 1–15. doi:10.1111/j.1467-8624.2006.00852.x
- Rozenblit, L. R., & Keil, F. C. (2002). The misunderstood limits of folk science: An illusion of explanatory depth. *Cognitive Science*, 26, 521–562. doi:10.1207/s15516709cog2605\_1
- Schooler, J. W., & Engstler-Schooler, T. Y. (1990). Verbal overshadowing of visual memories: Some things are better left unsaid. *Cognitive Psychology*, 22, 36–71. doi:10.1016/0010-0285(90)90003-M
- Schulz, L. E., Bonawitz, E. B., & Griffiths, T. L. (2007). Can being scared make your tummy ache? Naive theories, ambiguous evidence and preschoolers' causal inferences. *Developmental Psychology*, 43, 1124–1139. doi:10.1037/0012-1649.43.5.1124
- Schulz, L. E., & Gopnik, A. (Eds.). (2007). *Causal learning: Psychology, philosophy, & computation*. Oxford, UK: Oxford University Press. doi:10.1093/acprof:oso/9780195176803.001.0001
- Siegler, R. S. (1995). How does change occur: A microgenetic study of number conservation. *Cognitive Psychology*, 28, 225–273. doi:10.1006/cogp.1995.1006
- Siegler, R. S. (2002). Microgenetic studies of self-explanations. In N. Granott & J. Parziale (Eds.), *Microdevelopment: Transition processes in development and learning* (pp. 31–58). New York, NY: Cambridge University. doi:10.1017/CBO9780511489709.002
- Sobel, D. M., Tenenbaum, J. B., & Gopnik, A. (2004). Children's causal inferences from indirect evidence: Backwards blocking and Bayesian reasoning in preschoolers. *Cognitive Science*, 28, 303–333. doi:10.1207/s15516709cog2803\_1
- Ullman, T. D., Goodman, N. D., & Tenenbaum, J. B. (2012). Theory acquisition as stochastic search in the language of thought. *Cognitive Development*, 27, 455–480. doi:10.1016/j.cogdev.2012.07.005
- Walker, C. M., Lombrozo, T., Legare, C., & Gopnik, A. (2014). Explaining prompts children to favor inductively rich properties. *Cognition*, 133, 343–357. doi:10.1016/j.cognition.2014.07.008
- Wason, P. C. (1960). On the failure to eliminate hypotheses in a conceptual task. *Quarterly Journal of Experimental Psychology*, 12, 129–140. doi:10.1080/17470216008416717
- Wellman, H. M., & Liu, D. (2007). Causal reasoning as informed by the early development of explanations. In L. Schulz & A. Gopnik (Eds.), *Causal learning: Psychology, philosophy, & computation* (pp. 261–279). New York, NY: Oxford University Press. doi:10.1093/acprof:oso/9780195176803.001.0001

- Wilkenfeld, D., & Lombrozo, T. L. (2015). Inference to the best explanation (IBE) versus explanation to the best inference (EBI). *Science and Education, 10*, 1–19. doi:10.1007/s11191-015-9784-4
- Williams, J. J., & Lombrozo, T. (2010). The role of explanation in discovery and generalization: Evidence from category learning. *Cognitive Science, 34*, 776–806. doi:10.1111/j.1551-6709.2010.01113.x
- Williams, J. J., & Lombrozo, T. (2013). Explanation and prior knowledge interact to guide learning. *Cognitive Psychology, 66*, 55–84. doi:10.1016/j.cogpsych.2012.09.002
- Williams, J. J., Lombrozo, T., & Rehder, B. (2013). The hazards of explanation: Overgeneralization in the face

of exceptions. *Journal of Experimental Psychology: General, 142*, 1006–1014. doi:10.1037/a0030996

### Supporting Information

Additional supporting information may be found in the online version of this article at the publisher's website:

**Figure S1.** Consistency With Observed Data for Varying Parameter Values

**Appendix S1.** Details of the Bayesian Model