

What happened? Reconstructing the past through vision and sound

Tobias Gerstenberg, Max H. Siegel & Joshua B. Tenenbaum

{tger, maxs, jbt}@mit.edu

Brain and Cognitive Sciences, Massachusetts Institute of Technology

Abstract

We introduce a novel experimental paradigm for studying multi-modal integration in causal inference. Our experiments feature a physically realistic Plinko machine in which a ball is dropped through one of three holes and comes to rest at the bottom after colliding with a number of obstacles. We develop a hypothetical simulation model which postulates that people figure out what happened by integrating visual and auditory evidence through mental simulation. We test the model in a series of three experiments. In Experiment 1, participants only receive visual information and either predict where the ball will land, or infer in what hole it was dropped based on where it landed. In Experiment 2, participants receive both visual and auditory information – they hear what sounds the dropped ball makes. We find that participants are capable of integrating both sources of information, and that the sounds help them figure out what happened. In Experiment 3, we show strong cue integration: even when vision and sound are individually completely non-diagnostic, participants succeed by combining both sources of evidence.

Keywords: mental simulation; causal inference; intuitive physics; vision; audition; multi-sensory integration.

Introduction

Forensic scientists are experts at figuring out what happened. From little evidence they reconstruct intricate stories for how the crime must have played out. Often the only source of evidence is visual, but sometimes there is auditory evidence, too. When JFK was murdered, some witnesses reported having heard one shot, whereas others reported to have heard two shots from different directions, leading to different theories about what happened. Forensic scientists are not the only ones trying to figure out what happened. We do it all the time. And in doing so, we often draw on both visual and auditory information (Gaver, 1993). In this paper, we bring together research on people’s intuitive understanding of physics with research on multi-sensory integration to better understand human causal inference.

Recent work on intuitive physics has made popular the idea that people’s predictions about the future (Bates, Yildirim, Tenenbaum, & Battaglia, 2015; Battaglia, Hamrick, & Tenenbaum, 2013; Smith & Vul, 2013), inferences about the past (Smith & Vul, 2014, but see Carroll & Kemp, 2015), and causal judgments (Gerstenberg, Goodman, Lagnado, & Tenenbaum, 2015; Gerstenberg, Peterson, Goodman, Lagnado, & Tenenbaum, 2017) are guided by an approximate mental simulation of the underlying physical processes. However, work on intuitive physics has almost exclusively relied on presenting participants with visual information (although see Ecker & Heller, 2005; Sekuler, Sekuler, & Lau, 1997; Siegel, Magid, Tenenbaum, & Schulz, 2014).

Work on multi-sensory integration (Ernst & Banks, 2002; Körding et al., 2007; Sekuler et al., 1997) has shown how people combine evidence from multiple senses to arrive at a better understanding of the world. For example, Körding et al. (2007) have shown that people optimally combine visual

and auditory cues in a spatial localization task.

Here, we introduce a novel experimental paradigm for studying the integration of visual and auditory evidence in causal inference. We argue that people combine evidence from multiple senses by running hypothetical simulations on a mental model of the situation to predict what will happen in the future, and infer what happened in the past. We study the process of multi-sensory integration in an engaging physical setup in which participants’ task is to predict where a ball will end up that is dropped into a Plinko Machine (see Figure 1), or infer in which hole the ball was dropped.

The rest of the paper is organized as follows. We first describe the novel experimental paradigm we have developed for studying the integration of vision and sound for prediction and inference. We then describe the *hypothetical simulation model* which we use to explain participants’ judgments. In Experiment 1, participants are asked to make predictions and inferences based on visual information only. In Experiment 2, we look at how participants integrate both visual and auditory information. Finally, in Experiment 3, we demonstrate strong cue integration: even when visual and auditory information are individually completely non-diagnostic, together they reveal what happened.

Experimental paradigm

All of the experiments reported in this paper used the Plinko Machine setup shown in Figure 1. A ball is dropped into one of the three holes of the machine and lands in the sand at the bottom. The configuration of the three obstacles differs between trials. In the prediction task, participants saw the initial position of the ball and had to guess where it will land (Figure 1a). In the inference task, participants saw the final position of the ball and had to guess in which hole it was dropped (Figure 1b). In Experiments 2 and 3, the box was initially fully occluded by a cover and participants merely heard the sounds that the ball makes when being dropped. In the inference task of Experiment 2b, the cover was then removed so that participants could see where the ball landed. In Experiment 3, only part of the cover was removed so that the obstacles were visible but not the final position of the ball (Figure 1c).

Model

We modeled the Plinko Machine with the 2D physics engine Pymunk (<http://www.pymunk.org>) and rendered the stimuli in Unity (<https://unity3d.com>). The Pymunk physics engine allows us to generate physically realistic simulations of how the dropping ball interacts with the obstacles and walls.

Hypothetical simulation model

In line with the work on people’s intuitive understanding of physics discussed above, we assume that people mentally

simulate where the ball would end up if it was dropped in either of the holes, and that upon seeing a ball in the sand, they can use these simulations to infer in which hole it was dropped. Moreover, we assume that people not only mentally simulate the ball’s trajectory but also the sounds that it makes. We further assume that people have uncertainty both about the ball’s trajectory, as well as what sounds it makes. We describe how the model captures visual and auditory uncertainty in turn.

Vision Our hypothetical simulation model incorporates two sources of visual uncertainty. First, we assume that participants may be uncertain about exactly how the ball is dropped into a hole. We capture this uncertainty by initiating the ball at the center of a hole and uniformly varying the angle at which the ball is dropped within a small range. Second, we assume that participants have dynamic uncertainty about collisions (cf. Smith & Vul, 2013). We model this uncertainty by adding Gaussian noise to the magnitude of the ball’s velocity after each collision with an obstacle or a wall.

Figure 2 row 3 shows the predictions of the hypothetical simulation model. We refer to each of the selected trials as a world. The colored circles at the bottom of each world show where the ball would land if it was dropped in the different holes according to the noiseless physical simulation. The densities indicate where the hypothetical simulation predicts the ball would end up. For each hole, we ran 110 forward simulations. In world 3, the model is uncertain about where the ball would end up if it was dropped in hole 1. It could either end up to the left of the pentagon, or to the right of it. This bimodal distribution is due to the dynamic uncertainty that affects the ball’s magnitude of velocity after the initial collision with the triangle. When the ball is dropped in hole 2, the model believes that the ball will most likely end up to the right of the pentagon. However, it also predicts that it may end up to the left of the pentagon. This bimodality is due to the uncertainty about how exactly the ball is dropped.

We can use the hypothetical simulation model to infer in which hole a ball was dropped given its final position. To do so, we first consider the likelihood that the ball would end up in this position for each of the three holes, and then use Bayesian inference to infer a posterior distribution over holes starting from a uniform prior. To calculate the likelihood, we used a Gaussian linking function on the distance between

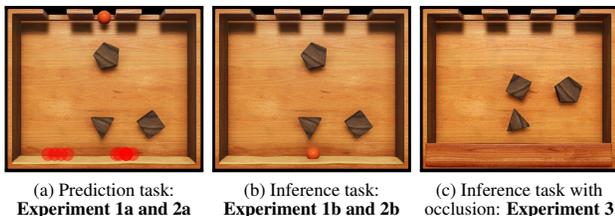


Figure 1: Illustration of (a) the prediction task, (b) the inference task, and (c) the inference task with occlusion. In the prediction task, participants had to guess where the ball will land. The red points at the bottom indicate a single participant’s guesses about where the ball will land. In the inference task, participants had to guess in which hole the ball was dropped.

the final location of the ball in a simulated trajectory, and the actual ball’s position: the smaller the distance, the greater the likelihood of a sampled trajectory.

Sound In some of the experiments, participants heard the sounds that the ball makes when being dropped. We used the physics engine to modulate pre-recorded impact sounds according to the difference in the magnitude of the ball’s velocity before and after a collision. Participants heard impact sounds whenever the ball collided with an obstacle or a wall. While the loudness of the impact sounds was modulated, the quality of each sound was the same. For example, there was no difference in tone between the ball colliding with an obstacle or the wall. Participants also heard a beep sound that indicated the time at which the ball was dropped, and another sound when the ball landed in the sand.

Forward simulations of the model not only generate trajectories of the ball’s motion but also a spectrogram of sounds. To do inference based on the sounds that the ball made, we only consider the time points in the simulation at which the impact sounds occurred. So, for each simulation we record the final position of where the ball lands, as well as a vector with the time points at which impact sounds happened.

To do inference on the auditory information, we first reject all samples that contain the wrong number of impact sounds. For the remaining samples, we score the likelihood of the vector of time points at which the sounds happened by again using a Gaussian linking function. A sample’s likelihood increases the closer its vector of time points at which collisions happened matches the vector of collision time points in the actual observation.

Integration of vision and sound The hypothetical simulation model combines visual and auditory information by assuming a multiplicative relationship between the visual and auditory likelihoods, and updating its posterior accordingly. Figure 2 row 6 shows the model’s predictions of where the ball will drop given both auditory and visual information. For many of the cases, the predictions are similar to a model that only considers visual information (row 4). However, in world 4, when the ball is dropped in the middle, a model that considers sound, rules out that the ball may have ended up to the right of the rectangle. When the ball is dropped in the middle, it first collides with the triangle, then with the rectangle, then with the wall, and then lands in the sand. Because the trajectory where the ball ends up to the right of the rectangle implies fewer impact sounds than actually observed, this possibility is ruled out.

Heuristic model

We compare the hypothetical simulation model to a heuristic model which only considers visual information and infers the hole that the ball was dropped in by proximity between the ball’s final location to the center of each hole. For example, in world 2 the heuristic model considers it most likely that the ball was dropped in hole 3, less likely that it was dropped in hole 2, and even less likely that it was dropped in hole 3 (due to its greater distance to the ball’s position).

Stimulus selection

We used the physics engine to randomly generate a large set of worlds, and then used the hypothetical simulation model to select 40 worlds of varying difficulty for the inferential task ranging from easy worlds in which the model was certain in which hole the ball was dropped to difficult worlds in which the model was uncertain about where the ball had been dropped. Worlds differed only in where the obstacles were positioned and in which hole the ball had been dropped. We generated the worlds by first defining a 3×3 grid of possible locations for three obstacles (triangle, rectangle, and pentagon). We then randomly assigned the three objects to the nine possible initial positions, jittered each obstacle in the x- and y-direction, and randomly rotated each object. Figure 2 row 1 shows a selection of worlds.¹

Experiment 1a: Prediction (vision)

In this experiment, participants see the initial location of the ball (cf. Figure 1a) and their task is to predict where it will land. This experiment tests whether participants are able to use their intuitive understanding of physics to simulate the ball's future trajectory. The ability to forward simulate hypothetical future trajectories is critical for inference which we will look at in Experiment 1b.

Methods

Participants Participants for all experiments were recruited via Amazon Mechanical Turk using Psiturk (Gureckis et al., 2016). 45 participants ($M_{\text{age}} = 37$, $SD_{\text{age}} = 11$, 21 female) took 26 minutes on average to complete this experiment.

Design and procedure Participants first learned about the Plinko Machine and that their task would be to indicate where they think the ball will land. A set of comprehension checks made sure that participants understood the instructions. Participants then watched four videos of the ball being dropped into Plinko Machines with the three obstacles placed in different positions. Each video was shown twice. Afterwards, participants saw 122 trials. First, two practice trials, and then 120 test trials in randomized order (3 trials, one for each hole, for each of the 40 worlds that we had selected for the inference task).

In each trial, the ball was positioned at the center of one of the three holes, and participants were asked to indicate where they think the ball will land by clicking ten times. For each click, a red semi-transparent circle appeared at the location of the click (see Figure 1a). Participants were told that they can indicate their confidence by clicking on the same location several times.

Results and discussion

Figure 2 row 2 shows participants' predictions for a selection of trials. For example, in world 1, participants' estimates of where the ball would land were quite accurate, although they underestimated how close the ball would end up to the right wall when being dropped in the middle hole. For

world 4, participants predicted that when the ball is dropped in the middle hole, it would land to the right of the rectangle, whereas in fact it lands on the left. Row 3 shows the predictions of the hypothetical simulation model.

To evaluate how well the model captures participants' predictions, we calculated participants' average prediction of where the ball lands for each world and hole, and compared that to the averaged prediction based on the model's simulations. Figure 3a shows that the model correlates well with participants' predictions.

Experiment 1b: Inference (vision)

The results of Experiment 1a show that people use their intuitive understanding of physics to mentally simulate where the ball would end up when being dropped into the Plinko Machine. In this experiment, we test how well participants can infer what happened.

Methods

Participants 46 participants ($M_{\text{age}} = 39$, $SD_{\text{age}} = 12$, 21 female) took 14 minutes on average to complete this experiment.

Design and procedure The instructions as well as the rest of the experimental procedure were similar to those of Experiment 1. Participants again first watched practice videos and then saw 42 trials including two practice trials at the beginning. The forty test trials were presented in randomized order. On each trial, participants saw the final location of the ball (see Figure 1b) and they were asked to indicate where they think the ball was dropped by entering percentage values into text boxes displayed on top of each of three holes. Participants were only able to proceed to the next trial if the percentage values of the three text boxes added up to 100% ($\pm 1\%$).

Results and discussion

Figure 2 row 5 shows participants' beliefs about where the ball was dropped for a selection of cases. For example, in world 1, participants are uncertain about whether the ball was dropped in hole 2 or hole 3, but ruled out hole 1. In world 4, they incorrectly believed that the ball was dropped in hole 1 whereas in fact it was dropped in hole 2. Note that in this case, the ball almost lands at the exact same spot whether it's dropped in hole 1 or hole 2 according to the ground truth.

Figure 4 shows how well the different models account for participants' inferences in this experiment. The hypothetical simulation model which uses the physics engine (Figure 4a) predicts participants' inferences based on the noisy forward samples that were simulated for the ball being dropped into each of the three holes (see Figure 2 row 3). Figure 4b shows the predictions of the hypothetical simulation model that uses participants' predictions from Experiment 1a (see Figure 2 row 4) to infer what hole the ball was dropped in. The heuristic model (Figure 4c) predicts inferences based on the ball's closeness to the three holes. The heuristic model cannot account for participants' judgments. For example, in world 3, the ball lands at a position that is close to the center

¹All materials including videos and sound files are available here: https://github.com/tobiasgerstenberg/what_happened

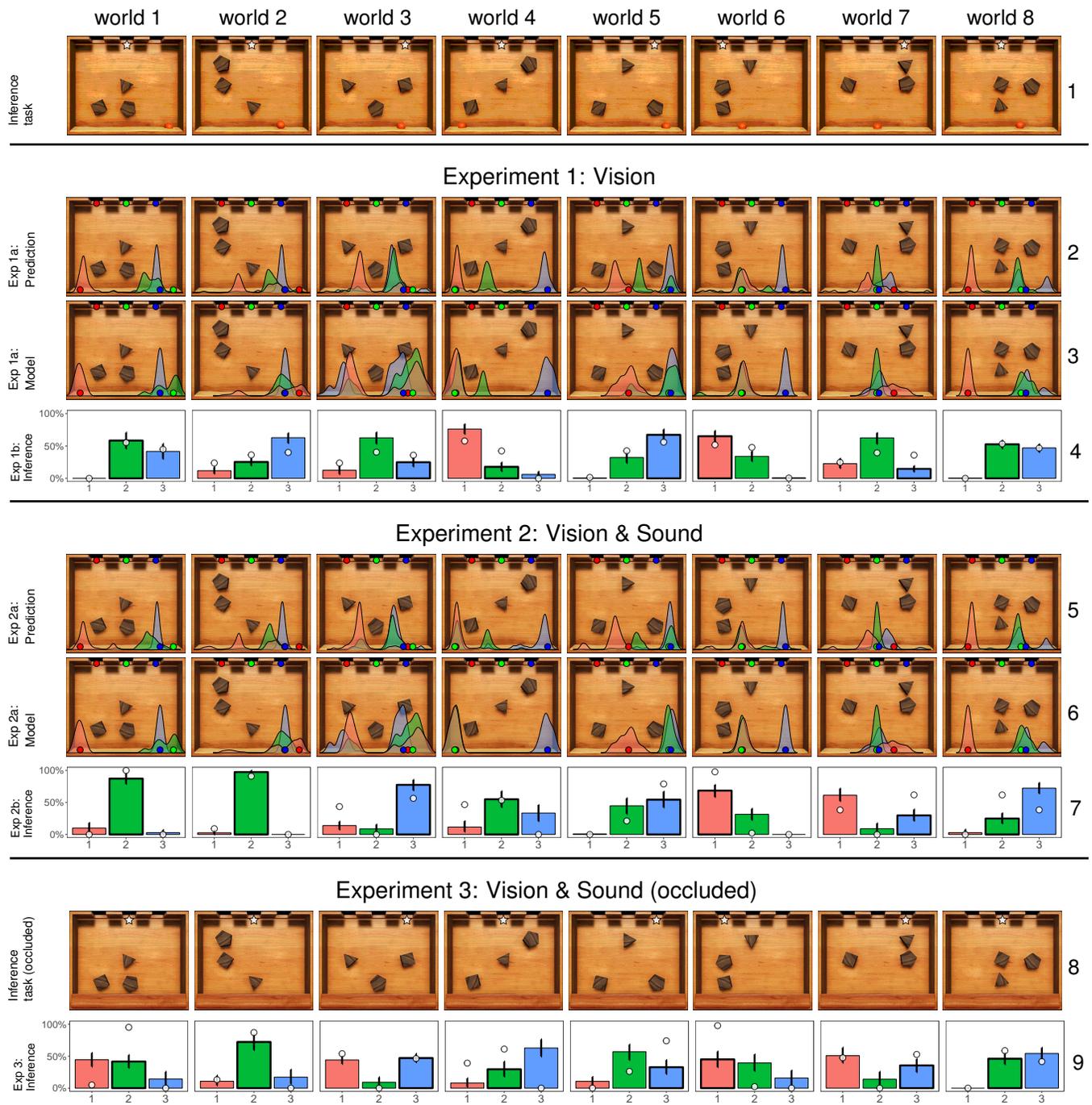


Figure 2: Summary of results for an illustrative selection of cases. These cases were chosen to illustrate how auditory information may improve performance (world 1–4), make no difference (world 5–7), or decrease performance (world 8). **Row 1:** Stimuli in the inference task. The star indicates through which hole the ball was dropped. **Row 2:** Densities showing participants’ predictions in **Experiment 1a** of where the ball would land if it was dropped in hole 1 (red), hole 2 (green), or hole 3 (blue). The colored circles at the bottom show for each hole where the ball would land according to the ground truth physics model. **Row 3:** Predicted densities generated by the hypothetical simulation model. **Row 4:** Inferences in **Experiment 1b** about which hole the ball was dropped in. Thick outlines indicate the correct response. Error bars indicate bootstrapped 95% confidence intervals. The white circles show the predictions of the hypothetical simulation model. **Rows 5–7:** Analog to rows 2–4 for **Experiment 2** in which participants also heard the ball dropping. The hypothetical simulation model shown in row 6 (densities) and row 7 (white circles) integrates both vision and sound. **Row 8:** Stimuli in the inference task with occlusion. **Row 9:** Inferences in **Experiment 3** in which the final ball position was occluded. *Note:* Error bars in all figures indicate bootstrapped 95% confidence intervals.

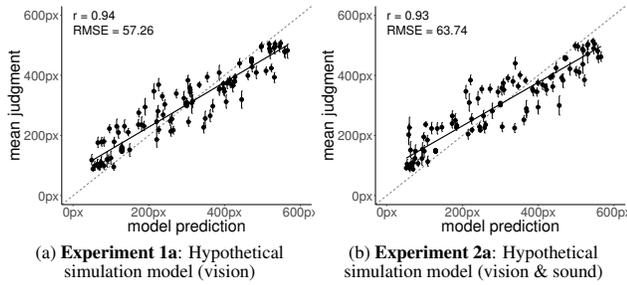


Figure 3: Relationship between the hypothetical simulation model and participants' judgments of where the ball will land.

of hole 3 but participants are more likely to believe the ball was dropped in hole 2.

Overall, the hypothetical simulation model based on people's predictions explains participants' inferences best, suggesting they approach the inference task by simulating forward where the ball would land if it was dropped in each of the three holes, and then placing most of the belief on holes that best explain why the ball ended up where it did. The results of Experiments 1a and 1b show that people can mentally simulate what will happen, as well as figure out what happened based on visual information about the ball's final position. In Experiment 2, we will look at how people integrate both visual and auditory information in prediction and inference.

Experiment 2a: Prediction (vision and sound)

In this experiment, participants predicted where the ball would land based on auditory and visual information.

Methods

Participants 36 participants ($M_{age} = 34$, $SD_{age} = 10$, 12 female) took 49 minutes on average to complete this experiment.

Design and procedure Design and procedure were identical to Experiment 1a with one important change. Participants now received auditory information about what happened. The practice videos now featured sounds. Participants heard a beep sound when the ball was dropped, impact sound when the ball collided with an obstacle or a side wall, and another sound when it landed in the sand. In the test trials, participants first saw a fully occluded box that only revealed in which hole the ball was dropped. They then listened twice to the sounds that the ball makes when being dropped. Afterwards, the occluder was removed to reveal the obstacles. Par-

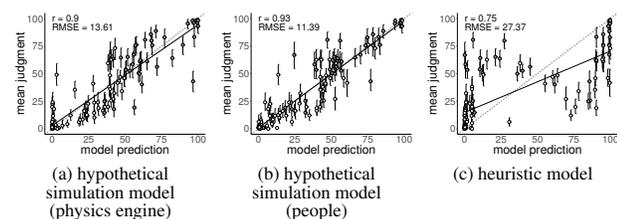


Figure 4: **Experiment 1b:** Scatter plots showing how well the different models can account for participants' inferences about where the ball was dropped. Gray circles indicate the correct holes.

ticipants then clicked ten times to indicate where they think the ball would end up like in Experiment 1a.

Results and discussion

Figure 2 row 4 shows participants' predictions based on having heard the sounds that the ball makes when being dropped and, subsequently, seeing where the obstacles are positioned. Overall, participants' predictions are similar to those from Experiment 1a where participants didn't have any auditory information. However, there are situations in which the auditory information affected participants' predictions. For example, in world 4 when the ball is dropped in hole 2, participants who had heard what sounds the ball made correctly predicted that it would land to the left of the rectangle at the bottom. In contrast, participants who hadn't heard the sounds believed that it would land to the right of the rectangle. The audio reveals that the ball collided three times, and thus suggests that it must have collided with the left wall before landing in the sand. The alternative path according to which the ball lands on the right of the rectangle would have only produced two impact sounds. Figure 3b shows that the hypothetical simulation model which takes into account both visual and auditory information provides a close fit to participants' predictions.

Experiment 2b: Inference (vision and sound)

Experiment 2b looks at how well participants integrate visual and auditory information to figure out what happened.

Methods

Participants 47 participants ($M_{age} = 36$, $SD_{age} = 11$, 21 female) took 21 minutes on average to complete this experiment.

Design and procedure Design and procedure were analog to Experiment 1b with the addition of the auditory evidence. For each trial, the box was first fully covered and participants heard the sounds that the ball makes when being dropped. Then the cover was revealed so that participants could see the obstacles as well as where the ball had landed in the sand.

Results and discussion

Figure 2 row 7 shows participants' inferences. In worlds 1–4, participants benefitted from having heard the sound. They were more likely to infer the correct hole compared to the condition in which participants didn't have any auditory information (row 4). For example, in worlds 1 and 2, participants ruled out hole 3 because of having heard an impact sound. In world 3, participants heard two impact sounds and thus ruled out that it was dropped in hole 2.

Figure 5 shows how well the different models account for participants' inferences in this experiment. The hypothetical simulation model again provides a good account of participants' inferences whether its predictions are based on the physics engine, or on participants' judgments in Experiment 2a. Figure 6 compares the accuracy in both the prediction and inference task between Experiments 1 and 2. Having auditory information improved performances only somewhat

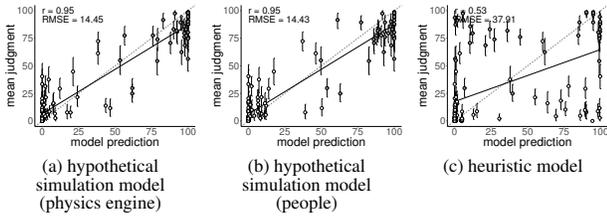


Figure 5: **Experiment 2b**: Scatter plots showing how well the different models account for participants' inferences about where the ball was dropped. Gray circles indicate the correct holes.

for the prediction task, whereas it helped considerably in the inference task. Taken together, the results of Experiment 1 and 2 show that people integrate visual and auditory information to predict what will happen as well as infer what happened. In Experiment 3 we use the Plinko Machine paradigm to demonstrate strong multi-sensory integration by considering a task in which vision and sound individually provide no information at all about what happened.

Experiment 3: Inference with occlusion

In this experiment, we made a minimal change to the Plinko Machine setup which renders the visual information alone non-diagnostic about what happened. We simply covered up the final position of the ball by adding an occluder in front of the bottom of the Plinko Machine (see Figure 1c). Participants still saw the position of the obstacles but not where the ball was. Thus, the visual cue itself is completely non-diagnostic about which hole the ball was dropped in. The auditory information itself is also completely non-diagnostic for where the ball was dropped. However, when both auditory and visual information are combined, it is possible to infer what happened by mentally simulating which of the possible paths that the ball could have taken are consistent with the sounds that one has heard.

Methods

Participants 43 participants ($M_{\text{age}} = 36$, $SD_{\text{age}} = 10$, 23 female) completed the experiment in 22 minutes on average.

Design and procedure Design and procedure were identical to Experiment 2b, with the only difference being that the final position of the ball is occluded (see Figure 1c). Figure 2 row 8 shows what the different worlds look like when the final position of the ball is occluded.

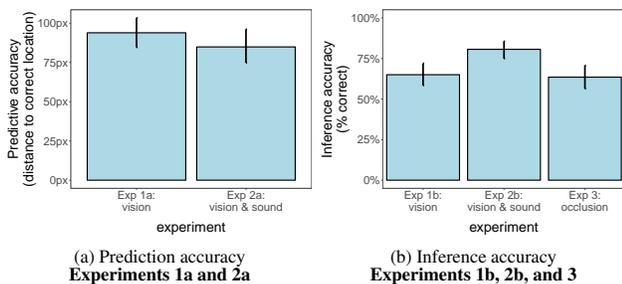


Figure 6: Accuracy in (a) predictive task as measured by the root mean squared error between judged and actual position of the ball (lower values indicate better accuracy), and (b) inference task as measured by the averaged percentage assigned to the correct hole.

Results and discussion

Figure 2 row 9 shows participants' inferences as well as model predictions. Overall, the hypothetical simulation model again provided a good account of participants' inferences ($r = 0.85$, $RMSE = 20.15$). Even without seeing the final position of the ball, participants were still able to infer what must have happened. Indeed, as Figure 6 shows, participants' accuracy in this experiment was similar to the experiment in which participants saw the final position of the ball but had no auditory evidence.

These results show a striking form of multi-sensory integration: through the process of mental simulation that operates over an intuitive physical understanding of a situation, people can integrate two sources of evidence that are individually completely non-diagnostic to figure out what happened.

General discussion

In this paper, we introduced a new experimental paradigm for studying multi-sensory integration of visual and auditory information through mental simulation. In the Plinko Machine task, participants either predict where a ball will land, or infer where it was dropped. The results of Experiment 1 show that participants can use their intuitive understanding of physics to predict what will happen, and infer what happened based on visual information only. In Experiment 2, we show that participants' inferences about what happened improve if, in addition to seeing where the ball is, they also heard how it got there. Finally, in Experiment 3, we illustrate a striking form of cue integration through mental simulation. In a setup where both visual and auditory evidence are individually completely non-diagnostic, participants were able to combine the information from both senses to figure out what happened.

Acknowledgments We thank Jeremy Schwartz for help with generating the stimuli. This work was supported by the Center for Brains, Minds & Machines (CBMM), funded by NSF STC award CCF-1231216.

References

Bates, C. J., Yildirim, I., Tenenbaum, J. B., & Battaglia, P. W. (2015). Humans predict liquid dynamics using probabilistic simulation. In *CogSci Proceedings*.

Battaglia, P. W., Hamrick, J. B., & Tenenbaum, J. B. (2013). Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences*, 110(45), 18327–18332.

Carroll, C. D., & Kemp, C. (2015). Evaluating the inverse reasoning account of object discovery. *Cognition*, 139, 130–153.

Ecker, A. J., & Heller, L. M. (2005). Auditory – visual interactions in the perception of a ball's path. *Perception*, 34(1), 59–75.

Ernst, M. O., & Banks, M. S. (2002). Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415(6870), 429–433.

Gaver, W. W. (1993). What in the world do we hear?: An ecological approach to auditory event perception. *Ecological Psychology*, 5(1), 1–29.

Gerstenberg, T., Goodman, N. D., Lagnado, D. A., & Tenenbaum, J. B. (2015). How, whether, why: Causal judgments as counterfactual contrasts. In *CogSci Proceedings*.

Gerstenberg, T., Peterson, M. F., Goodman, N. D., Lagnado, D. A., & Tenenbaum, J. B. (2017). Eye-tracking causality. *Psychological Science*, 28(12), 1731–1744. doi: 10.1177/0956797617713053

Gureckis, T. M., Martin, J., McDonnell, J., Rich, A. S., Markant, D., Coenen, A., ... Chan, P. (2016). psturk: An open-source framework for conducting replicable behavioral experiments online. *Behavior research methods*, 48(3), 829–842.

Körding, K., Beierholm, U., Ma, W., Quartz, S., Tenenbaum, J., & Shams, L. (2007). Causal inference in multisensory perception. *PLoS one*, 2(9), e943.

Sekuler, R., Sekuler, A. B., & Lau, R. (1997). Sound alters visual motion perception. *Nature*, 385(6614), 308–308.

Siegel, M. H., Magid, R., Tenenbaum, J. B., & Schulz, L. E. (2014). Black boxes: Hypothesis testing via indirect perceptual evidence. In *CogSci Proceedings*.

Smith, K. A., & Vul, E. (2013). Sources of uncertainty in intuitive physics. *Topics in Cognitive Science*, 5(1), 185–199.

Smith, K. A., & Vul, E. (2014). Looking forwards and backwards: Similarities and differences in prediction and retrodiction. In *CogSci Proceedings*.