# Supporting Information

## Battaglia et al. 10.1073/pnas.1306572110

### SI Materials and Methods

**Subjects and Stimulus Apparatus.** All subjects volunteered in response to an advertisement posted on the Massachusetts Institute of Technology (MIT) Brain and Cognitive Sciences human subjects notification system (we did not collect background personal data, but our subject population is composed of roughly half MIT students or employees and half local community members). All gave informed consent, were treated according to protocol approved by MIT's Institutional Review Board, and were compensated $10/h for participation. All experimental sessions were 1 h long and took place in 1 d, and subjects ran in exactly one session in one experiment. All had normal or corrected-to-normal vision. Stimuli were presented on a liquid-crystal display computer monitor, which subjects free-viewed from a distance of 0.5–0.75 m. They indicated their responses by depressing a key on the keyboard or by adjusting the computer mouse and clicking to lock in their choice. The numbers of subjects per experiment are listed in Table S1.

**Exp. 1.** Subjects were presented with virtual, 3D tower scenes [rendered using Panda3D (1)], like those in Fig. 2B (main text) and Fig. S1, which contained 10 rectangular blocks (each with 3D aspect ratio 3:1:1). The blocks were stacked within a square column by a sequential random process such that when placed on the tower, no single block would fall off of its support (although, when later blocks were added, they might cause previously placed blocks to fall). This design increased the complexity of the scenes and judgments, by ensuring that pairwise comparisons of adjacent blocks would not provide information about whether the towers were stable or not. Subjects were asked, "Will this tower fall?", and responded on a 1–7 scale, in which 1, 4, and 7 corresponded to "definitely will not fall", "not sure", and "definitely will fall", respectively, with intermediate numbers indicating intermediate degrees of confidence. These stimuli were deliberately chosen to be challenging and to evoke judgments across this graded scale, to more clearly test the role of simulation and probability in physical scene understanding. For more natural everyday scenes, people's judgments and the intuitive physics engine (IPE) model's predictions would likely be more deterministic (more clustered at the endpoints of this 1–7 scale) and also more objectively accurate (more similar to ground truth physics).

On each trial, subjects first viewed the scene for a 3-s "prephysics" interval, in which gravity and other forces were absent (i.e., the scene remained static regardless of whether the blocks would fall when gravity was applied). During this interval the camera panned 180° around the tower so subjects could view the tower's full 3D geometry. (Pilot testing suggested that performance would be similar, if more variable, for subjects presented with a single static view.) At the end of this interval a cylindrical occluder dropped vertically over the tower and the subject was able to make a judgment under no time constraint. In Exp. 1 (feedback), feedback was presented immediately after the subject responded in the form of a 2-s movie in which gravity was turned on and either the tower remained upright or some (or all) of its blocks fell; the floor also changed color to distinguish between "fell" (red) and "remained standing" (green). In Exp. 1 (no feedback), the next trial began immediately after the subject responded.

Before the test session in Exps. 1–4, participants performed a 20-trial familiarization session with a set of tower stimuli different from those in the test session and with feedback, to acclimate them to the task timing and judgment, the response keys, and the parameters of the scene (e.g., the towers' general appearances and physical characteristics such as friction, etc.). Exp. 1's

test session contained 360 trials: 60 towers (half of which would fall under gravity), repeated six times each across different trial blocks with randomized stimulus orderings and randomized trial-by-trial tower block colors and initial camera angles.

**Exp. 2.** Exp. 2 was similar to Exp. 1, except that all towers always fell under gravity (at least two blocks dropped). Subjects were asked "Which direction will the tower fall in?" and reported their judgments using the computer mouse by adjusting the orientation of a line on the floor, which extended from the base of the tower to the floor's perimeter, to indicate their expectation of its dominant fall direction (Fig. 3A, main text). Exp. 2's test session contained 360 trials, consisting of 60 stimuli repeated six times each in separate trial blocks, with the same randomizations of stimulus ordering, block coloring, and initial camera angles as in Exp. 1.

**Exps. 3 and 4.** Exps. 3 and 4 used a new set of tower scenes (Fig. 3 D and G) that was similar to that in Exps. 1 and 2, except that half of the blocks were 10 times heavier than the others. Blocks of different masses were visually distinguished by a dark, stone-like texture (heavy) and a pale green striped texture (light). As in Exps. 1 and 2, subjects were asked Will the tower fall? (Exp. 3) and Which direction will the tower fall in? (Exp. 4). The towers were organized into "state pairs": two towers whose geometries (block positions and poses) were identical, but whose heavy/light assignments were different. For some state pairs, the different mass assignments caused very different outcomes with respect to the model's predictions. As in Exp. 1, half of the towers in Exp. 3 would fall under gravity; as in Exp. 2, Exp. 4's towers always fell under gravity. In both experiments there were 48 towers with unique arrangements of blocks, each with two state-pair assignments of heavy and light blocks, for a total of 96 stimuli. The test sessions each contained 384 trials, consisting of the 96 stimuli repeated 4 times each in separate trial blocks, with the same randomizations as in Exps. 1 and 2.

**Exp. 5.** Exp. 5's scenes depicted a table on which a collection of blocks were arranged (Fig. 4A), half of which were red and the other half of which were yellow. Subjects were asked, "If the table were bumped, which color would be more likely to fall off?", and responded on a 1–7 scale in which 1, 4, and 7 corresponded to "definitely yellow", "not sure", and "definitely red". Subjects were instructed to assume that the bump's force was great enough to knock off at least one block. There were 12 possible block configurations that varied in the numbers and shapes of the blocks and the numbers, heights, and positions of the stacks in which they were arranged. There were also five different tables, one with a flat surface and four others with two short walls rigidly attached to different edges, which were designed to increase the complexity of the ensuing dynamics once the table was bumped. The 60 scene stimuli included all possible combinations of the 12 block configurations and five tables. We ran two conditions that differed by what information was provided about the ensuing bump: In the "uncued" condition, no additional information was provided about the bump; in the "cued" condition, a blue arrowhead in the scene pointed in the direction from which the bump would strike (Fig. 4B). Each stimulus was shown six times (with randomized red/yellow group assignments and viewing angles), two times in the uncued condition and four times in the cued condition (two times each with two different bump directions), for a total of 360 trials. Each

block of 60 trials showed the stimulus from only one of the two conditions, ordered "no cue", "cue", "cue", "no cue", "cue", "cue". On each trial, the stimulus interval was 3 s, during which the camera panned over a 90° arc; the scene remained unoccluded during the response interval. After every fifth trial, feedback was presented in which a force was applied to the table for 200 ms with random magnitude but strong enough to cause at least one block to topple and fall off the edge. In the bump cue condition, the feedback's actual bump matched the cued direction; in the uncued direction the bump's direction was randomized. Subjects first performed a 20-trial familiarization session with feedback on each trial to become acquainted with the task and the effects of the bump. The familiarization stimuli contained different scenes with fewer blocks in simpler arrangements, so that the judgments were not challenging.

**Exp. S1.** Exp. S1 was almost identical to Exp. 1, except that the towers were selected such that all were of the same height. There were 432 experimental trials: 108 towers, repeated four times each, and trials were divided into four trial blocks, with the same randomizations over trial orderings, block colorings, and initial camera angles as in Exp. 1. Correlations between subjects' judgments and IPE model predictions are shown in Fig. S4.

**Exp. S2.** Exp. S2 was identical to Exp. 2, except that subjects were asked to judge "How far will the furthest block come to rest?". They reported their judgments using the computer mouse by adjusting the radius of a circle on the floor, centered at the tower's base, to match the distance they expected the tower's blocks to come to rest. The model's output predicate was defined as the farthest point from the tower's base on any block.

**Data Analysis.** In Exps. 1, 3, and 5, where people responded on a 1–7 scale, these judgments were rescaled linearly to lie between 0 and 1 and averaged across repetitions and subjects. We compared the model's responses with subjects', using a Pearson correlation. In Exps. 2 and 4 we computed the circular mean of subjects' directional responses across repetitions and subjects. For Exps. 2 and 4, because judgments were angular, we used circular correlations (2) to quantify their agreement; note that circular correlation coefficient values are analogous but not directly comparable to Pearson correlations.

All correlation coefficients and 95% confidence intervals (CIs) were estimated using a bootstrap resampling procedure with 10,000 resamplings. The coefficients were taken as the median of the bootstrapped distribution; the lower and upper confidence intervals were always the 2.5th and 97.5th percentiles. All $P$ values were estimated using a direct bootstrap hypothesis test over the 10,000 resamples (3).

In Exps. 2 and 4, for some towers, the model's distribution of predicted directions was dispersed broadly around the circle and thus the model's circular mean estimate was very unstable (single samples could cause the response to shift by up to 180°) (Fig. S2). In these cases the model was indecisive and its point estimates of the direction of fall were not meaningful. We thus excluded those towers for which less than 80% of the model's distribution of predicted directions fell in any one-half of the circle, which left 47 of 60 towers included in the analysis for Exp. 2 and 72 of 96 for Exp. 4 (Fig. S2). In Exp. S2 we computed the mean of subjects' radial responses across repetitions and subjects.

The reported $F$-scores comparing IPE model confidences in Exp. S2 and Exp. 1 refer to a one-way ANOVA $F$-test statistic, measuring how separable the model's distributions of predictions were across stimuli within each experiment. They were computed as the ratio of the variance of per-stimulus prediction means over the means of the per-stimulus prediction variances.

## SI Appendix: IPE Model

**Computational Theory.** The model forms a judgment, $J_q$, by computing the expected value of a physical property query, $\mathcal{Q}_q(S_{0:T})$, which is a function of the initial state $(S_0)$ and the sequence of future states $(S_{1:T})$. The model's knowledge about these states is summarized as the Bayesian posterior probability distribution given observed information about the states and latent forces $(I_{S_0:S_T}$ and $I_f)$.

**Definitions.**

- $S_t$ : Scene state at time, $t$.
- $S_{t_0:t_1} = (S_{t_0}, S_{t_0+1}, \ldots, S_{t_1-1}, S_{t_1})$ : Sequence of scene states from time $t_0$ to $t_1$.
- $f_t$ : Extrinsic force applied beginning at time $t$.
- $f_{t_0:t_1}$ : Sequence of extrinsic forces applied from time $t_0$ to $t_1$.
- $I_{S_t}$ : Observed information about $S_t$.
- $I_f$ : Observed information about $f_{t_0:t_1}$.
- $\psi(\cdot)$ : Deterministic physical dynamics from $t_0$ to $t_1$, which maps $S_{t_0}$ to a new state at time $S_{t_1}$ : $S_{t_1} = \psi(S_{t_0}, f_{t_0}, t_1 - t_0)$. The force, $f_{t_0}$, is applied for a duration $t_1 - t_0$. The dynamics can be applied recursively,

$$S_{t_2} = \psi\big(\psi\big(S_{t_0}, f_{t_0}, t_1 - t_0\big), f_{t_1}, t_2 - t_1\big).$$

- We denote the repeated application of $\psi(\cdot)$ from $(t_0 : t_n)$ as $\Psi(\cdot)$,

$$S_{t_n} = \psi\big(\ldots\psi\big(S_{t_0}, f_{t_0}, t_1 - t_0\big), \ldots, f_{t_{n-1}}, t_n - t_{n-1}\big)$$
$$= \Psi\big(S_{t_0}, f_{t_0:t_{n-1}}, t_0 : t_n\big).$$

- $\mathcal{Q}_q(\cdot)$: Output predicate corresponding to a query, $q$, which maps an initial $(S_0)$ and a future sequence of scene states $(S_{1:T})$ to a judgment, $J_q = \mathcal{Q}_q(S_{0:T})$. In our experiments, the queries were sensitive only to the initial and final scene states (i.e., for Will the tower fall?, the query reflected how many blocks dropped from $t = 0$ to $t = T$), and so $J_q = \mathcal{Q}_q(S_0, S_T)$.

**Inputs.** The model represents knowledge of $S_t$, using a probability distribution, $\Pr(S_t)$, which, if $I_{S_t}$ is available, will be a posterior distribution defined by Bayes' rule, $\Pr(S_t|I_{S_t}) = \frac{\Pr(I_{S_t}|S_t)\Pr(S_t)}{\Pr(I_{S_t})}$. Similarly, the model represents knowledge of $f_{t_0:t_1}$, using a probability distribution, $\Pr(f_{t_0:t_1})$, or, if $I_f$ is available, $\Pr(f_{t_0:t_1}|I_f)$. For brevity, the remaining formulas assume observed information is available.

**Physical inference.** The probability of a future state, $S_{t+1}$, given a previous state, $S_t$, $\Pr(S_{t+1}|S_t, f_t) = 1$, for $S_{t+1} = \psi(S_t, f_t, 1)$, and 0 for any other value of $S_{t+1}$, because $\psi(\cdot)$ is deterministic. The model represents the distribution over initial and future states, $S_0$ and $S_{1:T}$, as determined by physics:

$$\Pr\big(S_{0:T}|I_{S_0}, I_f\big) = \int_{f_{0:T-1}} \Pr(S_T|S_{T-1}, f_{T-1})\cdots\Pr(S_1|S_0, f_0)$$

$$\Pr(S_0|I_{S_0})\Pr(f_{0:T-1}|I_f)\mathrm{d}f_{0:T-1}$$

$$\begin{aligned} = \int_{f_{0:T-1}} &\Pr(\psi(S_{T-1}, f_{T-1}, 1)|S_{T-1}, f_{T-1})\cdots \\ &\Pr(\psi(S_0, f_0, 1)|S_0, f_0)\Pr(S_0|I_{S_0})\cdots \\ &\Pr(f_{0:T-1}|I_f)\mathrm{d}f_{0:T-1}. \end{aligned} \quad \textbf{[S1]}$$

In this work we focus on the influence that the latent forces, $f_{0:T-1}$, have on future scene states—not on the latent forces themselves—which is why they can be integrated out.

When only $S_0$ and $S_T$ are required, as is the case in the present work, the intermediate times can be integrated out,

$$\Pr\left(S_T, S_0 | I_{S_0}, I_f\right) = \int_{f_{0:T-1}} \int_{S_{1:T-1}} \Pr(S_T | S_{T-1}, f_{T-1}) \ldots \Pr(S_1 | S_0, f_0)$$

$$\Pr(S_0 | I_{S_0}) \Pr(f_{0:T-1} | I_f) \mathrm{d}S_{1:T-1} \mathrm{d}f_{0:T-1}$$

$$= \int_{f_{0:T-1}} \Pr(S_T | S_0, f_{0:T-1}) \Pr(S_0 | I_{S_0}) \Pr(f_{0:T-1} | I_f) \mathrm{d}f_{0:T-1}$$

$$= \int_{f_{0:T-1}} \Pr(\Psi(S_{t_0}, f_{0:T-1}, 0:T) | S_0, f_{0:T-1})$$

$$\Pr(S_0 | I_{S_0}) \Pr(f_{0:T-1} | I_f) \mathrm{d}f_{0:T-1}.$$

[S2]

Note that because $\psi(\cdot)$ can be applied recursively,

$$\Pr(S_T | S_0, f_{0:T-1}) = \Pr(\Psi(S_0, f_{0:T-1}, 0:T) | S_0, f_{0:T-1}).$$

**Outputs.** The model's output judgment for a query, $q$, is

$$J_q = \mathrm{E}\left[\mathcal{Q}_q(S_{0:T}) | I_{S_0}, I_f\right]$$

$$= \int_{S_{0:T}} \mathcal{Q}_q(S_{0:T}) \Pr(S_{0:T} | I_{S_0}, I_f) \mathrm{d}S_{0:T}, \qquad \text{[S3]}$$

where $\Pr(S_{0:T} | I_{S_0}, I_f)$ is defined in Eq. **S1**.

Again, as our experimental queries were sensitive only to $S_0$ and $S_T$,

$$J_q = \mathrm{E}\left[\mathcal{Q}_q(S_0, S_T) | I_{S_0}, I_f\right]$$

$$= \int_{S_T} \int_{S_0} \mathcal{Q}_q(S_0, S_T) \Pr(S_T, S_0 | I_{S_0}, I_f) \mathrm{d}S_0 \, \mathrm{d}S_T, \qquad \text{[S4]}$$

where $\Pr(S_T, S_0 | I_{S_0}, I_f)$ is defined in Eq. **S2**.

**Experimental IPE Implementation.** *Inputs.* As described in the main text, the IPE model tested experimentally had three numerical parameters, $(\sigma, \mu, \phi)$, which controlled the state uncertainty, physical attributes, and latent force inputs to the simulation, respectively. They were defined and implemented as follows.

The model's state representation could be separated into the geometric state, $G$, and the physical state, $p$: $S = (G, P)$. In our experiments, $G$ included the numbers, positions, poses, and shapes of the objects' states; and $P$ included the mass density of each object, represented as the vector $m$, in addition to other physical parameters such as coefficients of friction and elasticity that were fixed across all studies to values typical for everyday wooden blocks. Similarly, $I_S = (I_G, I_P)$, where $I_G$ represents observed information about $G$, e.g., the objects' visually indicated geometry in the image; and $I_P$ represents observed information about $P$, e.g., the objects' visually indicated mass assignments in Exps. 3 and 4.

Because subjects viewed scenes from many viewpoints over a continuous 180° span, and these viewpoints were randomly chosen for each trial, we represented the model's visual inference of the initial scene geometry, $G_0$, as a simple, one-parameter viewpoint-invariant approximation to the Bayesian posterior,

$$\Pr(G_0 | I_{G_0}) \approx \pi(G_0; \overline{G}_0, \sigma),$$

which represents the distribution over $G_0$ given the true geometry, $\overline{G}_0$, and the parameter, $\sigma$, which represents the magnitude of the posterior state uncertainty. The $\pi(\cdot)$ was defined by taking the true geometry, $\overline{G}_0$, and adding horizontal, zero-mean Gaussian noise (SD $\sigma$) to the ground truth object positions independently. Because the noise could cause interobject penetrations, the

objects' coordinates were then transformed by a deterministic constraint-satisfaction procedure that selected the nearest configuration for which no objects violated each others' volumes. This procedure ran very small time steps of the physics engine, resetting the objects' velocities to zero after each step, which caused all that were detected as being in collisions to shift apart until no collisions were detected. We evaluated the plausibility of this viewpoint-invariant approximation ($\pi(\cdot)$) by comparing its samples with those of a prototype Bayesian vision system that we developed, which used Markov chain Monte Carlo (MCMC) to sample directly from the Bayesian posterior of scene geometries conditioned on specific subsets of images, under a likelihood function defined by a graphics rendering package (*SI Appendix: Bayesian Vision System*). Results were highly similar, in terms of both the samples of scene geometries (Fig. S8 *A–F*) and the predictions on our experimental tasks (Fig. S8 *G* and *H*). In our simulations, $\sigma$ could take 1 of 11 possible values,

$$\sigma \in \{0.0, 0.05, 0.10, 0.15, 0.20, 0.25, 0.30, 0.35, 0.40, 0.45, 0.5\}$$

(where the short width of one of the tower's blocks in Exps. 1–4 was 1, and the cube-shaped block's width in Exp. 5 was 2).

The model's representation of objects' mass densities $m$ was controlled by the parameter $\mu$. In Exps. 3 and 4, which contained heavy and light blocks, $\mu$ represented the scalar ratio between their respective mass densities; in Exps. 1, 2, and 5, all blocks had the same density, so $\mu = 1$. The information provided by $I_P$ was approximated as deterministically indicating which blocks were heavier (in Exps. 3 and 4) or that all had the same density (in Exps. 1, 2, and 5), and the model then used its assumption about $\mu$ to set each block's individual mass density. In Exps. 1, 2, and 5, $\mu = 1$; in Exps. 3 and 4 it could take 1 of 12 values,

$$\mu \in \{0.25, 1.0, 2.0, 2.5, 3.2, 4.0, 5.0, 6.3, 8.0, 10, 13, 16\}.$$

For tower scenes (Exps. 1–4), all objects in the scene were adjusted so that the total mass was 2 kg. For Exp. 5 each block was 1 kg.

The model's latent force dynamics, $f_{0:T-1}$, represented possible vibrations and bumps that could be applied to the scene (main text). We approximated them as a horizontal force [angle $\theta$ and magnitude $\phi$ (main text)], applied from $t = 0$ ms to $t = 200$ ms (and no force after 200 ms) to the surface on which the objects were situated; for simplicity we drop the $0:T-1$ subscript and refer to $f_{0:T-1}$ as, $f = (\theta, \phi)$. The observation about the forces could be separated into two terms, $I_f = (I_\phi, I_\theta)$, where $I_\phi$ reflected language-based instructions like "If the table were bumped..." (Exp. 5). When available (Exp. 5, bump cue condition), $I_\theta$ represented the cue's indicated latent force direction; the model assumed $\theta$ was uniform over the range $[I_\theta - 45, I_\theta + 45]$. When $I_\theta$ was unavailable, the model assumed $\theta$ was uniform over the range [0, 360].

In Exps. 1–4, $\phi$ could take 12 possible values,

$$\phi \in \{0.0, 0.2, 0.4, 0.6, 0.8, 1.0, 1.2, 1.5, 2.0, 2.5, 3.0\},$$

and in our analyses we fitted the best value from Exp. 1. Fig. S3 shows the correlations between model and human judgments for a range of $\sigma$ and $\phi$ values in Exp. 1. In Exp. 5 (the "Bump?" task), $\phi$ could take 16 values,

$$\phi \in \{22.0, 25.6, 29.2, 32.8, 36.4, 40.0, 43.6, 47.2, 50.8, 54.4, 58.0,$$
$$61.6, 65.2, 68.8, 72.4, 76.0\},$$

where, in Exp. 5's scenes, $\phi = 22.0$ usually caused only one or a few blocks to fall off the table, and $\phi = 76.0$ caused many or almost all to fall off. The model's judgment took the expectation over all $\phi$ values.

**Physical simulation.** For our experiments, the model's judgment in Eq. S4 was computed via a Monte Carlo approximation

$$J_q^{MC} \approx \frac{1}{k} \sum_{i=1}^{k} \mathcal{Q}_q\Big((G_0,\mu)^{(i)}, \Psi\big((G_0,\mu)^{(i)}, f^{(i)}, 0:T\big)\Big) \qquad \textbf{[S5]}$$

that used the Open Dynamics Engine (4, 5) to approximate $\Psi(\cdot)$, over a $T$ that was always 2,000 ms with a step size of 1 ms, and where $k$ is the number of independent simulation samples, initialized by independently drawing values of $S_0$ and $f$ from the approximate posterior distributions as described above. In Exps. 1–4, $k = 48$. In Exp. 5, $k = 12$, and only the values of $S_0$ were drawn independently; rather than sampling $f$, we computed predictions for 16 equally spaced values of $\theta$, crossed with the 16 values of $\phi$ (Fig. 4 $C$–$E$).

**Query definitions.** All queries regarded an object as having "fallen" if, between $S_0$ and $S_T$, its $z$ coordinate was displaced downward greater than 0.025 unit. For Exp. 5, objects were defined as having "fallen off" the table if their $z$ coordinate was less than 0 at time $T$ (the table's surface was at $z = 0$).

The queries corresponding to the four experimental tasks were as follows:

- Exps. 1 and 3 (Will it fall?): $\mathcal{Q}_{fall}(S_0, S_T) = $ the fraction of objects that fell.
- Exps. 2 and 4 (In which direction?): $\mathcal{Q}_{dir}(S_0, S_T) = $ the angle of the center of mass of the fallen blocks, in the $x, y$ plane. For $\mathcal{Q}_{dir}$'s Monte Carlo sum, we used the circular mean across the angular query outputs.
- Exp. 5 (Bump?): $\mathcal{Q}_{bump}(S_0, S_T) = $ the ratio of red to total blocks that fell off the table.
- Exp. S2 (How far?): $\mathcal{Q}_{far}(S_0, S_T) = $ the radius of the farthest fallen block from the base of the tower, in the $x, y$ plane.

Due to ambiguities and vagueness inherent in language, more than one output predicate could be consistent with any given task query. For example, although we represented answers to the Will it fall? query as the average proportion of blocks that fell, other possibilities (e.g., the proportion of simulations for which one or more blocks fell or the proportion for which more than half of the blocks fell) would give similar results.

## SI Appendix: Model-Free Accounts

An alternative explanation for people's physical scene understanding is that they do not possess some model of the world, but instead use model-free methods that depend heavily on their experienced interactions with the world. People have been exposed to many stable and unstable arrangements of objects over their lifetime, and perhaps when presented with a new scene such as our experimental stimulus, they consult some stored representation of similar scenarios they have previously experienced and produce a response that reflects the outcomes of those scenarios. Such stored representations might be exemplar based, composed of individual examples of scene and outcome pairs, or feature based, encoding their experiences as compressed features—the gross shape of the scene, the number of vertical elements, etc.—and a function that takes as inputs these features and returns as output statistically dependent physical outcomes. Either the features or the parameters of this input–output function (or both) could be constructed on the basis of experience. To what extent could one of these model-free methods explain subjects' performance in our experiments?

One reason to be skeptical of a model-free account is the lack of significant practice effects over the course of our experiments (*SI Appendix: Analysis of Learning*), which might be expected if people's judgments were driven primarily by data-driven learning mechanisms.

However, the central argument for a model-based probabilistic physical simulation mechanism over a primarily model-free mechanism based on learning from experience is the former's generality. Although the same IPE model could be applied to each new task by using a predicate that was appropriate for the query and to each new scene type by adjusting its parameters to reflect the objects' physical attributes and the scene's patterns of applied forces, the extensive amount of variation across natural settings means that no compact set of trained features or exemplars is even applicable to all tasks or scenes.

Still, to formally evaluate an account based on memory or learned features, we built separate feature models for the Will it fall?, How far?, In which direction?, and Bump? tasks and fitted their parameters to best predict subjects' data in each experiment, using ridge regression (a linear model that prevents overfitting by subjecting the coefficients to an L2 penalty). The features used are listed in Table S2. For each experiment, the penalty parameter was selected through a cross-validated fitting procedure on subjects' data. This method is relatively generous to a feature-based account because it allows the feature weights to vary arbitrarily across experiments, taking on whatever values best fits people's performance in each particular experiment.

Because standard multivariate regression analysis methods are not available for circular data, we report the best individual circular correlation between people's responses and any single feature in the In which direction? tasks. In Exps. 3 and 4, some geometric features implicitly took the blocks' masses into account [e.g., $F_F(2)$, height of the center of mass]. To grant these features the ability to make judgments that were sensitive to the heavy/light assignments, we generated features for a range of assumed $\mu$ values and used those that best fitted people's judgments.

Across experiments, the IPE model fits were generally significantly better than those of the best feature-based models, often dramatically so (Fig. S4). This was true even allowing for features that were selected specifically for each task and multiple free parameters that were tuned to maximize their fits to each experiment separately.

The results of fitting the best feature-based account to each individual experiment were as follows. In Exp. 1, across the feedback and no-feedback conditions, one heuristic predictor, the tower's height [Table S2, $H_F(1)$], was best correlated with subjects' responses ($\rho = 0.75$, 95% CIs [0.68, 0.81]), so we conducted a controlled variant (Exp. S1) identical to Exp. 1 (feedback) except with different subjects ($n = 10$) and 108 new towers (each repeated over four blocks) that were all of the same height, to assess performance when the most dominant geometric heuristic was neutralized. The IPE, with parameters identical to those of Exp. 1 ($\sigma = 0.2$ and $\phi = 0.2$), again had a significantly higher correlation with people ($\rho = 0.81[0.74, 0.87]$) than the geometric heuristics ($\rho = 0.70[0.59, 0.78]$, $P < 0.001$). In Exp. 2, the best geometric feature had a significantly lower circular correlation with people's circular-mean responses than the IPE ($\rho_{circ} = 0.39[0.21, 0.56]$, $P < 0.001$). In Exp. 3, the regression-fit mass-sensitive feature predictions (with best-fit $\mu = 6$) had a correlation of $\rho = 0.71[0.61, 0.79]$ and the mass-insensitive features (with $\mu = 1$) had a correlation of $\rho = 0.61[0.46, 0.72]$, which were also both significantly lower than that of the IPE model ($P < 0.02$ and $P < 0.001$, respectively). The features' state-pair differences (with $\mu = 6$) correlation, $\rho = 0.54[0.34, 0.73]$, was also significantly lower ($P < 0.03$) (the mass-insensitive features did not make different state-pair predictions). In Exp. 4, the best mass-sensitive feature [$F_D(6)$, best-fit $\mu = 16$] and mass-insensitive feature [$F_D(6)$, $\mu = 16$] had correlations of $\rho_{circ} = 0.43[0.28, 0.58]$ and $\rho_{circ} = 0.31[0.17, 0.45]$, respectively, which were both significantly lower than that of the IPE model ($P < 0.001$ and $P < 0.1$). However, the state pair circular correlation between the IPE model and the subjects was not significantly higher ($P < 0.07$) than the feature correlation

($\rho_{circ} = 0.50[0.18, 0.75]$ vs. $\rho_{circ} = 0.25[0.066, 0.45]$, respectively). In Exp. 5, the regression-fit features had significantly lower correlations than the IPE model did with people's responses, in both the uncued and the cued conditions, respectively ($\rho = 0.73[0.58, 0.85]$, $P < 0.01$; $\rho = 0.68[0.58, 0.77]$, $P < 0.001$).

We do not claim that geometric features learned through experience play no role in physical scene understanding: In one experiment, Exp. S2 (How far?), a simple feature was significantly better than the IPE model at predicting people's judgments. Below (*SI Appendix: Approximations*) we further discuss why and when features may be used instead of mental simulations. However, we argue against the sufficiency of a purely memory- or feature-based account and in favor of a general capacity for simulation, based on the IPE model's distinctive ability to explain quantitatively the whole set of experimental results presented here, as well as people's inferences across a wide range of real-world scenes and tasks such as those in Fig. 1 and the examples discussed in main text under *Architecture of the IPE*. It is difficult to imagine a set of features flexible enough to capture all of these inferences, yet compact enough to be learnable from people's finite experience.

## SI Appendix: Approximations

Our IPE model as tested so far is surely incomplete in key ways. It is likely that people's capacity to represent the full physical state of a scene, to simulate the dynamics of many objects in motion over time, and to maintain faithful representations of probabilities is more limited–and also more adaptive—than our use here of a simple physics engine running several dozen stochastic samples. The human IPE likely trades off precision for efficiency much more aggressively than engineers typically do. People may summarize the individual objects in a many-object scene coarsely, substituting in aggregate representations that approximately capture how such "stuff" tends to behave. They may use simplified shape representations, such as spheroids or convex polyhedra, that allow multiobject interactions to be computed more efficiently. The temporal resolution of the IPE's simulations may be low, potentially leading to errors and biases that only more finely spaced time steps would be able to avoid. The laws of physics embedded within the human IPE may not uphold basic physical principles, such as the conservation of energy and momentum, and the simulation may blend mechanics principles that are objectively independent into combined rules; e.g., because friction is ever present, objects in simulated motion may gradually but constantly lose energy and momentum over time. The IPE may represent uncertain sets of possible states with only a few samples or instead represent a state's probability, using a continuous-valued weight rather than through its frequency of occurrence.

These approximations may have little effect in many everyday contexts, where even short simulations based on coarse representations of objects' shapes and positions are sufficient to make useful predictions. However, they could result in perceptual errors or illusions with more complex scenes or more demanding tasks, pointing to ways in which our model could be improved to better capture the approximations and shortcuts people exploit. The following subsections illustrate in more depth several types of approximations that the human IPE may make and the role they might play in people's physical intuitions.

**Approximating Physics.** One kind of approximation is motivated by the intrinsic difficulty of making certain kinds of judgments via simulation in the presence of complex dynamics. The Open Dynamics Engine (ODE) and other standard physics engines can simulate highly nonlinear systems, such as many-body collisions, for which accurate predictions over even short time intervals are computationally intensive and probably beyond what the human IPE can perform. Consider a bowling ball at the moment it leaves the bowler's hand: We can mentally extrapolate its path for several seconds as it rolls down the lane, but the instant it strikes the pins, the ensuing motions become unimaginably complex. The same dynamics apply when a stack of blocks falls or is knocked over, as in Exps. 1–5, but did not pose serious challenges there because the judgments queried were mostly insensitive to collisions occurring beyond the early stages of the simulation. To test whether people might rely on alternatives to simulation in cases where these factors matter more, we conducted a supplementary study (Exp. S2), using the same tower stimuli of Exps. 1–2, but instead asked subjects to predict "How far will the blocks come to rest?". This judgment depends sensitively on tracking each block precisely until it comes to rest, often after multiple collisions. Here the IPE model makes much less confident predictions (*F*-ratio of 66.34 vs. 172.17 in Exp. 1, where higher values indicate greater separation between the model's simulated distributions across stimuli) (*SI Appendix: Data Analysis*), and people's judgments on this task were also much less accurate, as assessed both by correlation with ground truth ($\rho = 0.38$ vs. $\rho = 0.64[0.46, 0.79]$ in Exp. 1) and by correlation with the IPE model ($\rho = 0.71[0.56, 0.81]$ vs. $\rho = 0.92[0.88, 0.94]$ in Exp. 1 using the same parameter values) in Exps. 1 and 2. Intriguingly, in this case a simple geometric feature—the height of the tower—was correlated with the IPE model's inferences and better predicts people's judgments ($\rho = 0.93[0.87, 0.96]$) than the model does (Fig. S4). People may fall back on such learned features as simple heuristics for predictions when complex dynamics make mental simulation impractical, just as a bowler learns that targeting the 1 and 3 pins from a slight angle predicts a strike even though she cannot imagine how the pins will move to produce that outcome.

**Approximating the Scene.** Another kind of challenge for mental simulation can arise under even simple dynamics but with complex object shapes. In our IPE model, as in most standard physics simulation engines, predicting an object's response to gravity or other forces depends on representing its center of mass and moments of inertia. These variables are easy to estimate for the blocks in Exps. 1–5, with their cuboidal shape and uniform density, but for objects whose mass is distributed in an unknown way over a complex 3D volume, they likely exceed people's ability to estimate with much precision. People may use strongly simplifying priors, such as taking an object's mass to be uniformly distributed over a coarse approximation to its shape (e.g., a cuboid, an ellipsoid, or a convex hull). This would be consistent with recent reports that people can accurately judge the stability of asymmetric objects when they are convex (6), while also explaining why people incorrectly predict the stability and dynamics of objects with more complex, nonconvex shapes (Fig. S5): They expect a wheel rim will roll downhill at the same rate as a filled disk (Fig. S5*A*) (7), when in fact it rolls more slowly, and they are surprised to see the dragonfly in Fig. S5*B* apparently defying gravity, when in fact it is balanced stably around its center of mass (Fig. S5*C*).

**Approximating Probabilities.** There are several important questions of approximation in our proposed IPE's sample-based representation of probabilities. How many simulation samples does the human IPE use? How faithfully do these samples represent probabilistic quantities of interest? We explored these questions by examining the variance of subjects' responses, which will decrease with the number of simulation samples that contribute to each judgment, $k$ (Eq. **S5**), in the same way that the (squared) SEM in experimental statistics will decrease as the sample size is increased. For each stimulus in each experiment, we computed the across-subject response variances (or circular variances, in Exps. 2 and 4) and model variances as follows. We first normalized the subjects' and model's responses in Exps. 1, 3, and 5

to be in [0, 1] so they would have common coordinates (note that Exps. 2 and 4's model and subject responses did not need to be normalized because they were both already in radians). We then shifted their values so that they would have a mean of zero across repetitions for each stimulus to remove the component of variance due to mean differences. Fig. S6 A–E shows the relationships between these model and subject response variances; the colored lines correspond to the relationships between the variances for different values of $k$ (Fig. S6 legend for details). The relationship between the model and subject variances favors an effective value of $k = 1$ for Exps. 1–4 and $k = 5$ or 6 for Exp. 5. However, these estimates are based on the assumption that the sampling variance was the only source of trial-by-trial variance and thus represents a lower bound: There are surely other sources of variance that are typically present in psychophysical data, such as general decision-making noise. This may be why in Exps. 1–4 the subjects' variances are elevated above the $k = 1$ line. To separate the sampling variance from these other stimulus-independent noise sources and better estimate the true value of $k$, we performed the following linear regression analysis. If we assume that the sampling variation is statistically independent of the other sources of variance, then the variance in people's responses is the sum of these two variance components and the slope of the fit line will correspond to $\frac{1}{k}$, whereas the intercept will correspond to $\omega^2$, where $\omega$ is the SD of the other noise sources. The estimated $k$ and $\omega$ values for Exps. 1–5, respectively, were

$$k_1 = 4.8[2.2, 66], \omega_1 = 0.27[0.25, 0.29]$$
$$k_2 = 2.8[1.9, 5.1], \omega_2 = 0.56[0.51, 0.60]$$
$$k_3 = 4.5[2.9, 8.9], \omega_3 = 0.21[0.19, 0.22]$$
$$k_4 = 3.1[2.1, 5.4], \omega_4 = 0.47[0.42, 0.51]$$
$$k_5 = 7.4[4.8, 13], \omega_5 = 0.09[0.07, 0.10].$$

This suggests that subjects typically formed their judgments based on 3–7 simulation samples per trial. A possible reason for the small differences between Exps. 1–4's $k$ values of 3–5 and Exp. 5's of 7 is that Exp. 5's model predictions included uncertainty in the bump's direction and magnitude, in addition to the state uncertainty ($\sigma$) also present in Exps. 1–4: It may be that in such cases of greater uncertainty the IPE increases its numbers of samples to more accurately approximate the physical outcome. The differences between the $\omega$ terms in Exps. 1 and 3 and in Exp. 5 may reflect different effects that the tasks' demands caused in subjects' decision processes (Exps. 2 and 4's responses were in radians and so not comparable to the other experiments' $\omega$ terms).

That these estimated values of $k$ were around 3–7 is consistent with the "one and done" account of Vul et al. (8), which claims that many judgments and behaviors might be the product of taking a small number of probabilistic samples from an internal posterior distribution. To examine how the number of samples affects the IPE model's ability to make accurate probabilistic inferences, for each experiment we created several IPE model variants that ran only a small number of samples and compared their judgments to those of the original IPE model whose judgments were based on many samples. Fig. S6 F–J shows that across experiments, the IPE model needs only a small number of samples to make predictions that correlate well with the predictions of the original IPE model; the correlations associated with model variants that used the (rounded) numbers of samples estimated above were 0.93, 0.83, 0.89, 0.84, and 0.93 for Exps. 1–5, respectively. This demonstrates that for our tasks, small numbers of samples are generally sufficient to make similar predictions to those of a model that has a more complete representation of probabilities.

Besides the simple Monte Carlo methods that we used here, there are other more sophisticated ways to use sampling in simulation to represent physical uncertainty over time. These are especially useful in cases when new observations are being continually collected (unlike in most of our experimental tasks, where all relevant observations are made at the beginning of each trial). Sequential importance sampling and particle filtering can use time-dependent samples that have weights associated with them, where the weight is proportional to an estimate of the posterior probability of that sample's state given the observations up to that point in time. The unscented Kalman filter uses a procedure to select a representative nonrandom subsample of states to which to apply the deterministic physical dynamics, to capture uncertainty over time. Future work should pursue this question of how exactly samples are used to represent and update dynamically changing probabilities within a simulation-based framework.

## SI Appendix: Analysis of Learning

Across all experiments, our analyses treated the subjects' data as stationary and constant by collapsing across multiple stimulus repetitions, yet subjects performed hundreds of trials with repeated stimulus presentations and often received feedback. One possibility is that subjects arrived at or substantially improved their behavior through learning over the course of the trials, rather than drawing primarily on a fixed internalized model of physics. We examined each subject's responses for evidence of practice or learning effects by computing their judgments' differences from the IPE model and its variants as a function of trials completed (Fig. S7). Fig. S7 A and B shows subjects' time-averaged differences as a function of trial number for Exp. 1's feedback and no-feedback conditions, respectively; Fig. S7 C and D shows the judgment differences averaged across subjects. There appear to be at most negligible changes over the course of the trials for most individual subjects, as well as for the average over subjects. An almost-perfect correlation ($\rho = 0.95$, 95% CIs [0.95, 0.95]) between subjects' mean responses in Exp. 1's no-feedback and feedback conditions (Fig. S7E) also suggests that feedback did little to alter people's judgments.

We also fitted regression lines to every subject's judgment differences as a function of trials completed and found that the largest shift in any single subject's judgments away from the model was +20%, whereas the largest shift toward any model variant was −17% (mean −1.6%, SD 6.6%). Fig. S7 F–L shows histograms of these slopes per experiment; most shifted by less than 10%, which indicates that subjects' responses were largely constant across trials. Fig. S7M shows these histograms (for shifts with respect to the probabilistic IPE model) pooled across all experiments and subjects, as well as a bootstrap resampled set of all experiments' data with the trial orderings randomized to express the hypothesis that there was no true shift in subjects' responses across trials. A bootstrapped hypothesis test of the difference between the empirically observed distribution of shifts and this no-shift reference distribution was not significant ($P = 0.38$). Taken together, these analyses suggest that feedback played a minimal role and that there was little if any effect of practice or learning across trials.

## SI Appendix: Bayesian Vision System

Our IPE model uses a simplified input representation of the scene: a sample of its 3D geometric state that approximates a Bayesian posterior distribution on scenes' given images. To explain physical scene understanding more fully, however, we ultimately need to capture all of the factors that govern how people infer underlying scenes from observed images, including effects of viewpoint, occlusion, and so on. Although a full treatment of the visual inference problem is beyond our scope here and would stretch the bounds of most conventional machine vision systems, here we offer an initial attempt both to validate our model's assumptions about the input geometry distribution and as a proof of concept to motivate future research integrating Bayesian vision with probabilistic physical reasoning.

We implemented an image-based vision system for approximating the Bayesian posterior distribution over $G$ given $I_G$ by MCMC in an inverse-graphics model and compared its judgments (Fig. S8 $D$ and $F$) to those of our IPE's viewpoint-invariant, $\pi$-based samples (Fig. S8 $C$ and $E$) for the tower stimuli in Exps. 1 and 2 (Fig. S8 $A$ and $B$). The system's geometric state estimate $G$ was defined as the positions and poses of all 10 blocks in the tower, and its observed evidence $I_G$ was three $256 \times 256$ 8-bit (256 grayscale levels) images of the tower rendered under perspective projection (using OpenGL) from a series of three viewpoints rotated by 45° and then filtered using a normalized 2D Gaussian blur kernel with $x,y$ scale parameter values of 11 pixels. The blur kernel's scale was selected so that the variance of the Bayesian system's samples would roughly match the viewpoint-invariant model's variances. The prior over $G$ was assumed to be uniform over all possible positions and poses of the blocks. The likelihood (probability of the image data given the geometric scene state) was defined as the product of per-pixel normal distributions (whose SD was 128 grayscale levels) between the observed image and an OpenGL-rendered image of the latent scene.

We used a Metropolis–Hastings (MH) sampling algorithm (9) to draw 5,000 geometric state samples (Fig. S8 $D$ and $F$), with samples grouped into 10-sample scans. Within a scan, each block's state was updated once, and the order of block updates was chosen uniformly with replacement, independently per scan. The MH algorithm's proposals consisted of updating each block's horizontal position, drawing a value from a horizontal bivariate normal distribution (with a diagonal covariance matrix whose SD terms were 0.2), and adding that vector to the block's current position, keeping its vertical position fixed. We initialized the sampler's latent state at the true state to minimize "burn-in" overhead and additionally discarded the first 1,000 samples of each run as burn-in as well. These simplifications make the system not generally applicable for performing vision in arbitrary scenes. However, they allowed us to conveniently draw approximate posterior samples representing how people could parse our experimental stimuli, under the assumption that we can accurately perceive the numbers of objects and their general positions.

Fig. S8 $G$ and $H$ shows that when the Bayesian system's samples are input to the IPE model, the resulting judgments are highly correlated (Exp. 1, $\rho = 0.90[0.84, 0.94]$; Exp. 2 [circular], $\rho = 0.94[0.91, 0.97]$) with judgments of the viewpoint-invariant ($\pi$-based) model used in the main text.

1. CMU Entertainment Technology Center (2010) Panda3D. Available at www.panda3d.org. Accessed October 7, 2013.
2. Fisher N, Lee A (1983) A correlation coefficient for circular data. *Biometrika* 70(2):327–332.
3. Efron B, Tibshirani R (1993) *An Introduction to the Bootstrap* (Chapman & Hall/CRC, New York and London), Vol 57.
4. Smith R (2010) Open Dynamics Engine. Available at www.ode.org. Accessed October 7, 2013.
5. Baraff D (2001) Physically based modeling: Rigid body simulation. SIGGRAPH Course Notes, *Association for Computing Machinery SIGGRAPH* 2:2–1.
6. Cholewiak SA, Fleming RW, Singh M (2013) Visual perception of the physical stability of asymmetric three-dimensional objects. *J Vis* 13(4):12.
7. Proffitt DR, Kaiser MK, Whelan SM (1990) Understanding wheel dynamics. *Cognit Psychol* 22(3):342–373.
8. Vul E, Goodman N, Griffiths T, Tenenbaum J (2009) One and done? Optimal decisions from very few samples. *Proceedings of the 31st Conference of the Cognitive Science Society*, eds Taatgen N, van Rijn H (Cognitive Science Society, Austin, TX), pp 66-72.
9. MacKay D (2003) *Information Theory, Inference, and Learning Algorithms* (Cambridge Univ Press, Cambridge, UK).
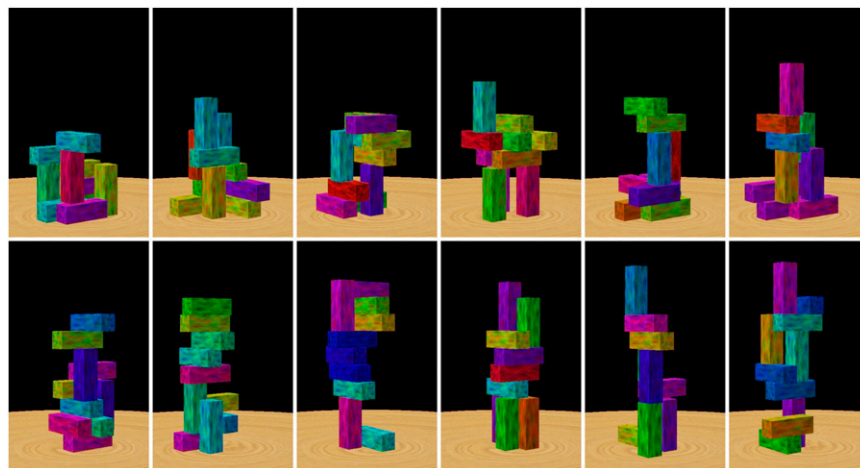
**Fig. S1.** Example of tower stimuli from Exp. 1.

**Fig. S2.** Exp. 2 IPE model's decisiveness. Each ring plot shows a sample of one tower's model-predicted directions. The towers are sorted by how concentrated the predictions are (numbers indicate rank). The 47th tower and lower (blue) were included in the analysis but the 48th and higher (red) were not.



**Fig. S3.** Correlations between Exp. 1's model and human judgments as a function of IPE parameters, $\sigma$ (x axis) and $\phi$ (y axis). The green and orange boxes/lines indicate the probabilistic and ground truth best fit parameters, corresponding to Fig. 2 *C* and *D*, respectively.



**Fig. S4.** Summary of model fits from Exps. 1–5 and Exps. S1 and S2. Bar graphs show the correlations (including 95% CIs) with people's judgments for both the IPE model and the best feature-based alternative account, across each experiment. Exps. 5uc and 5c refer to Exp. 5 uncued and cued conditions, respectively. Different features were required to fit behavior in different experiments. The graph (*Left*) shows which sets of features were applicable (colored boxes) or inapplicable (white boxes) to each experiment.
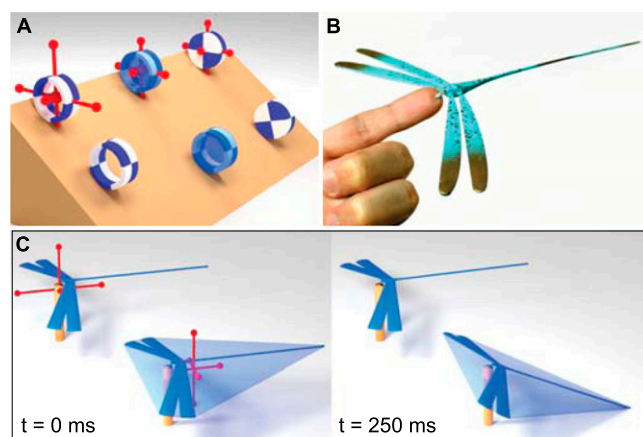
**Fig. S5.** Physics illusions and errors may arise from how the IPE approximates objects' dynamical properties. Moments of inertia along each object's principal axes are indicated by the length (inverse) of the red bars, centered at the estimated center of mass. (*A*) A wheel rim (leftmost object) has a different inertial tensor than a disk (rightmost object), such that it rolls downhill more slowly, but these inertial tensors are nontrivial to calculate. Treating each object's mass distribution as uniform over its convex hull gives a cognitively plausible approximation in general and here predicts (as people naively expect) that the wheel rim (middle object) has the same inertial properties and rolls at the same rate as the disk. (*B*) A toy dragonfly perches dramatically on its nose and is surprisingly balanced. (*C*) A model dragonfly with a physically correct center of mass and inertial tensor (upper) remains balanced after 250 ms in a simulation, but modeling it with the same convex-hull approximation as in *A* (lower) locates the dragonfly's center of mass slightly behind its nose (as people do) and leads to the intuitive expectation that it will tip from its perch.
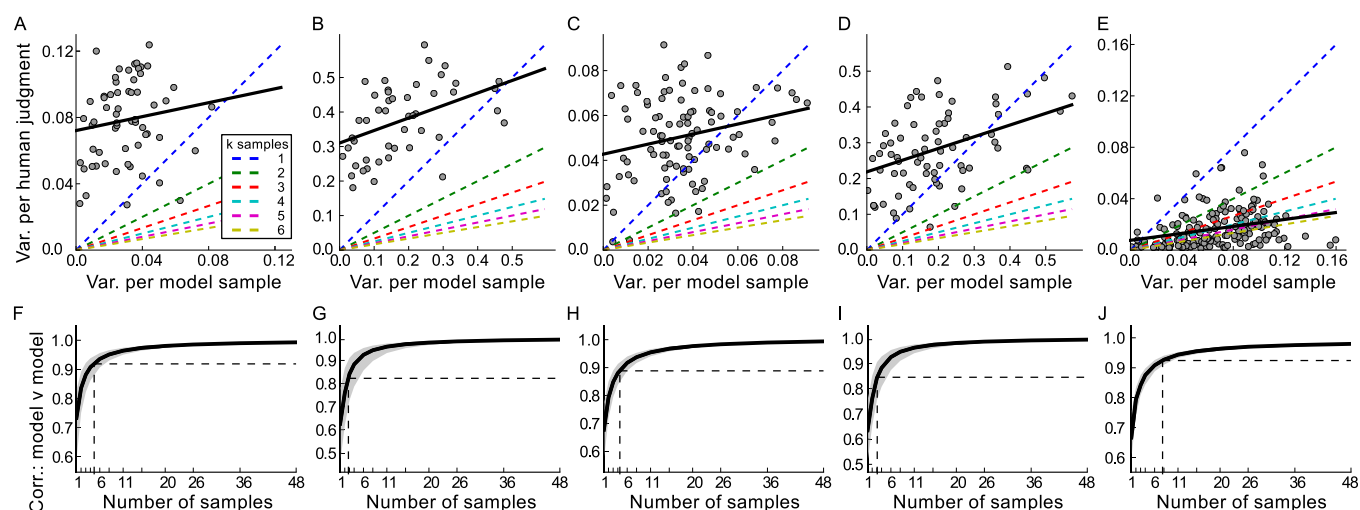


**Fig. S6.** Analysis of the number of IPE simulation samples. The columns of subplots correspond to Exps. 1–5, respectively. (*A–E*) Each point indicates the variance across human responses for a single stimulus (*y* axis) vs. the variance across model samples for that same stimulus (*x* axis). The blue dashed line (with the greatest slope) indicates a one-to-one correspondence between the model's sample variances and people's judgment variances, which would be consistent with each person's judgment being based on a single simulation sample. If instead each person formed judgments by taking a mean across $k > 1$ samples (drawn from the same probabilistic model), then we would expect the variances of people's judgments to be smaller than the model's sample variances by a factor of $k$ (analogous to the squared SEM). The other dashed lines, with decreasing slopes ($\frac{1}{k}$, for $k = 2 \ldots 6$), depict expected correspondences between the human response variances and the model's sample variances if human judgments were based on means of 2–6 simulation samples, respectively. The black solid lines show best-fit regression lines, whose intercepts reflect stimulus-independent trial-by-trial variance due to sources other than sampling variability. (*F–J*) The thick black line depicts the correlation (*y* axis) between the IPE model's predictions based on $k$ samples (*x* axis) and the model's predictions with the full set of samples from the original IPE; the gray ranges are 95% CIs (estimated by bootstrapped resampling of the $k$ samples); at $k = 48$ samples the correlations converge to 1. The vertical dashed line indicates the best-fit value of $k$ estimated from the linear regression analysis (*SI Appendix: Approximating Probabilities*), and the horizontal dashed line indicates the correlation level at this value of $k$.
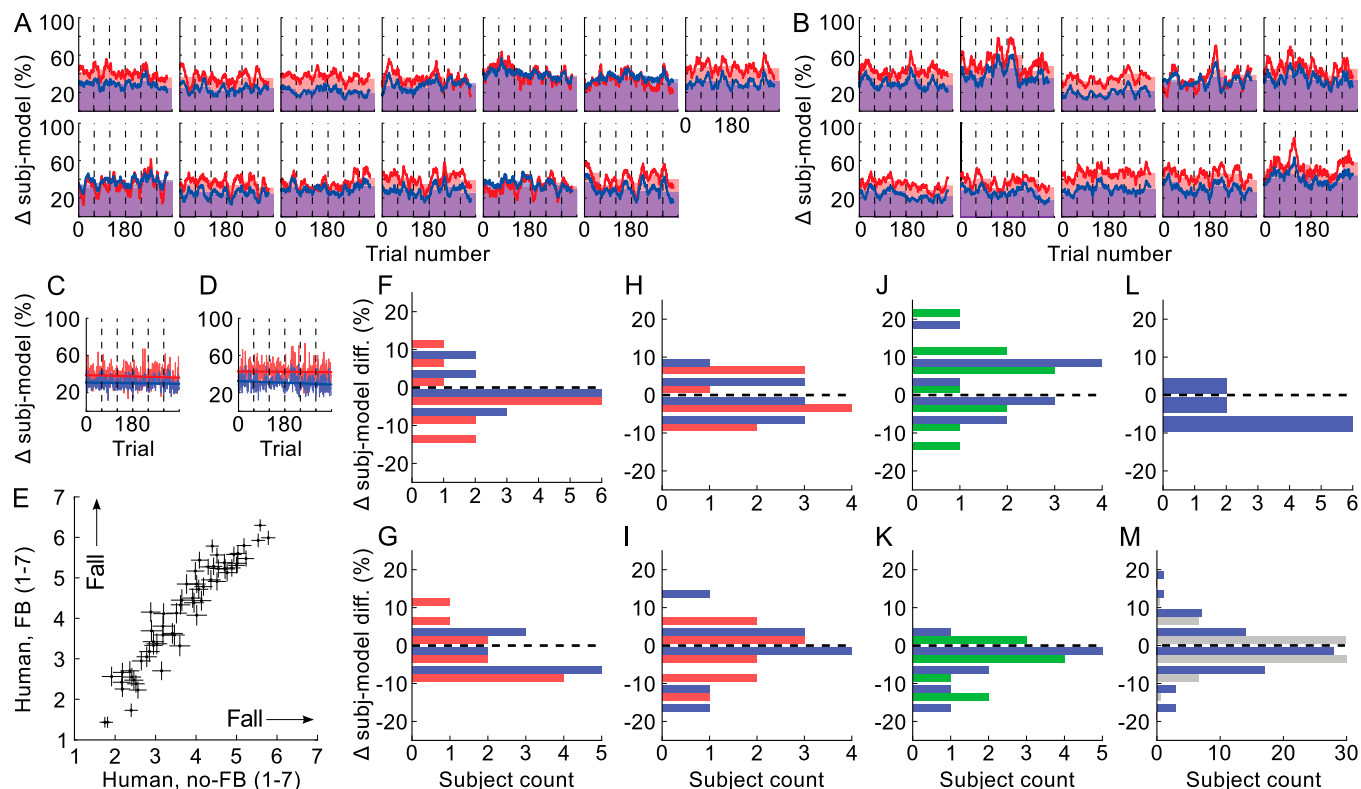
**Fig. S7.** Analysis of learning. (*A*) Individual subjects' average (over 20-trial sliding window) response differences (percentage, *y* axis) from ground truth (red: $\sigma = 0$, $\phi = 0$) and probabilistic IPE model (blue: $\sigma = 0.2$, $\phi = 0.2$) as a function of trial number (*x* axis) for Exp. 1's feedback condition. The vertical dashed lines indicate different trial blocks. (*B*) The same as *A* for Exp. 1's no-feedback condition. (*C*) Average response differences across all subjects (not a sliding window) for Exp. 1's feedback condition. The overlaid colored lines are best-fit regression fits. (*D*) The same as *C*, but for Exp. 1's no-feedback condition. (*E*) Correlation between subjects' average judgments (raw 1–7 responses, with SEM) in Exp. 1's no-feedback condition (*x* axis) vs. Exp. 1's feedback condition (*y* axis). (*F*) Histogram of slopes (*y* axis) of best-fit regression lines fitted to individual subjects' response differences from ground truth (red) and the probabilistic IPE model (blue) for Exp. 1's feedback condition. The slopes' units are change in response differences (percentage) across the experimental session; 0.0 indicates no change and −100% indicates a change from 100% error to 0% error. (*G*) The same as *F* for Exp. 1's no-feedback condition. (*H*) The same as *F* for Exp. 1 (same height). (*I*) The same as *F* for Exp. 2. (*J*) The same as *F* for Exp. 3. The blue and green bars are the slopes from fits to the mass-sensitive and mass-insensitive IPE models, respectively. (*K*) The same as *J* for Exp. 4. (*L*) The same as *F* for Exp. 5. The bars are slopes from fits to the full IPE model. *A–D* and *F–L* show minor changes in subjects' response differences from either the ground truth or the probabilistic IPE model across trials, indicating minimal effects of practice or learning. (*M*) Histograms of all slopes across experiments with respect to the probabilistic IPE model (blue) and the (rescaled) bootstrapped distribution of slopes for subjects' data with randomly permuted trial orderings (gray).

**Fig. S8.** Bayesian vision system. (*A* and *B*) Two original tower images. (*C* and *E*) Viewpoint-invariant ($\pi$-based) scene samples. (*D* and *F*) Bayesian vision system's scene samples. (*G*) Correlation between Exp. 1's model judgments for IPE model variants that input the Bayesian vision system's scene samples (*x* axis) and the viewpoint-invariant ($\pi$) samples used in the main text (*y* axis). (*H*) The same as *G*, but for Exp. 2.

**Table S1. Numbers of subjects per experiment**

| Experiment no. | Judgment | No. subjects |
| --- | --- | --- |
| 1, feedback | Fall? | 13 |
| 1, no feedback | Fall? | 10 |
| S1, same height | Fall? | 10 |
| 2 | Direction? | 10 |
| 3 | Fall?, mass | 11 |
| 4 | Direction?, mass | 10 |
| 5 | Bump? | 10 |
| S2 | Far? | 10 |
| Total | | 74 |

**Table S2. Geometric features**

| $F_F(\cdot)$ | Exps. 1 and 3 and Exps. S1 and S2, Will it fall?/How far? |
|---|---|
| $F_F(1)$ | Tower's height |
| $F_F(2)$ | Height of the tower's center of mass |
| $F_F(3)$ | Minimum critical angle* of the tower |
| $F_F(4)$ | Minimum critical angle across subtowers† |
| $F_D(\cdot)$ | Exps. 2 and 4, In which direction? |
| $F_D(1)$ | $x,y$ angle of minimum critical angle |
| $F_D(2)$ | $x,y$ angle of minimum critical angle across subtowers |
| $F_B(\cdot)$ | Exp. 5, Bump? |
| $F_B(1)$ | Average distance from the table's center |
| $F_B(2)$ | Average $x,y$ distance from the table's center |
| $F_B(3)$ | Average height |
| $F_B(4)$ | Average distance from the nearest edge |
| $F_B(5)$ | Average $x,y$ distance from the nearest edge |
| $F_B(6)$ | Minimum $x,y$ distance from the nearest edge + height |
| $F_B(7)$ | Minimum $x,y$ distance from the nearest edge |
| $F_B(8)$ | Like 1, except with maximum, instead of average |
| $F_B(9)$ | Like 2, except with maximum, instead of average |
| $F_B(10)$ | Like 3, except with maximum, instead of average |
| $F_B(11)$ | Like 4, except with minimum, instead of average |
| $F_B(12)$ | Like 5, except with minimum, instead of average |
| $F_B(13)$ | Like 6, except with minimum, instead of average |
| $F_B(14)$ | Like 7, except with minimum, instead of average |

*Critical angle is defined as the angle of center-of-mass of the tower about the nearest edge (in the horizontal plane) of the convex hull around the tower's base. Negative critical angles mean the center-of-mass is outside the convex hull (more unstable), and positive values mean it is inside. For $F_B$, each property was computed over the red blocks and over all of the blocks, and their ratio was used as the corresponding feature's value.

†Subtower is defined as a disjoint (noncontacting) subset of the blocks in a tower. Multiple subtowers existed in some stimuli as a natural consequence of the random procedure by which towers were generated; these subtowers were supported by the ground, but were not in contact with each other.