

1 We thank all the reviewers for their insightful and constructive comments, and will revise the paper accordingly.

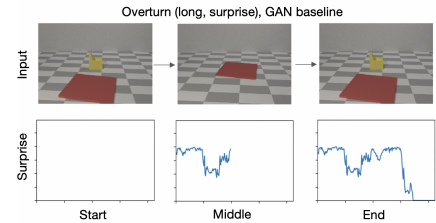
2 **(R1) Goal.** Our work has two aims. Our primary goal is to formalize human (especially infant) physical cognition, with
3 models that can be tested by developmental psychologists. But we also aim to show how modeling infant cognition can
4 inspire more robust AI vision systems that extract physical object representations from video and can detect violations
5 of physical expectations to use as learning signals.*

6 **(R1) Comparative run-times.** On a standard CPU server, baselines run at 2–4s/frame, while ADEPT runs at 6–
7 7s/frame. We currently have not optimized ADEPT for run-time performance and believe there is ample room for future
8 improvement, such as GPU acceleration via differentiable particle filtering.†

9 **(R1) Accuracy per-person.** Calculating relative accuracy requires a comparison between surprise/control pairs, but
10 each participant only observed one video from a pair (to avoid biases from seeing earlier near-identical scenes). We
11 therefore must aggregate across participants to provide this distribution.

12 **(R1) L2 loss for the GAN model.** Using an L2-based surprise score, the GAN model has a relative accuracy of 0.41
13 on our dataset (vs. 0.63 with a discriminator-based surprise score, as shown in Table 1).

14 **(R1) Baseline model images like Fig. 4.** While baselines do not have
15 interpretable internal belief states similar to the middle panels of Fig. 4,
16 we can plot model surprise over time. For example, to the right we see the
17 GAN model is not additionally surprised when the occluder fully rotates
18 over the object (vs. ADEPT in which surprise spikes; Fig. 4b).



19 **(R1, R2) Cognitive process underlying human surprise.** ADEPT re-
20 flects a plausible hypothesis that people have a probabilistic, object-based
21 model of intuitive physics (e.g., Battaglia et al., 2013‡), and that surprise is
22 driven by low probability events under that model (e.g., Teglas et al.§). We plan to further explore the human/model
23 match in future work, e.g., examining moment-by-moment surprise or neural correlates of object disappearance.

24 **(R1, R2) ADEPT hyperparameters.** We consider the difference between the sampling probabilities for surprising
25 events and their losses a form of Importance Sampling that allows rare events to be captured with small, cognitively-
26 plausible hypothesis sets by sampling those events more often than they should occur (see Lieder et al.¶ on this
27 phenomenon in humans). Empirically, we also found that performance was insensitive to the exact values or tracking
28 methods, so long as physics violations are more surprising than moderate amounts of perceptual uncertainty.

29 **(R2) Problem setup and generalization.** We designed our model to match objects based on general principles (e.g.,
30 color, shape, and size constancy). We stress that ADEPT’s training was not specific to the test dataset: there were no
31 unphysical scenes or violations in the training set. The derenderer is trained on sequences of three frames to infer the
32 instant velocity of objects, not long-term motion patterns. The training set had motion patterns similar to the test videos
33 to allow fair comparisons with baseline models, which do require longer sequences of motion to form predictions. We
34 will clarify these points in revision.

35 **(R2) Generalization to other datasets: IntPhys.** We have run ADEPT on the IntPhys dataset, only retraining the
36 derenderer to handle different visual object properties. Because the IntPhys test server is offline, we evaluate models on
37 the validation set of scenes designed to test “object permanence,” as described in Riochet et al. This set contains 90
38 matched sets of videos (2 plausible, 2 implausible within each set). ADEPT achieves an overall relative accuracy of
39 0.73, outperforming all baselines (Enc-Dec: 0.61, GAN: 0.53, LSTM: 0.65). As R2 noted, the videos of IntPhys have
40 different visual and motion patterns (e.g., complex textures and gravitational motion), so the high performance on this
41 second dataset suggests our model generalizes to situations where we have no control over the training and test data.

42 **(R2) Failure cases.** We agree and will move discussion of failure cases from supplemental to main text in the revision.

43 **(R2) Dataset.** We will release the dataset along with all code, human data, and model evaluations upon publication.

44 **(R2) Occluder modeling.** Our physics engine assumes motion changes require force, and so cannot capture the
45 up-and-down motion of occluders. We chose to model them separately to avoid producing a constant surprise signal.

46 **(R3) Additional information on metrics and fair comparisons.** We follow IntPhys by using the ‘relative accuracy’
47 metric to compare plausible and implausible videos within each matched set, and also follow their choice of plausibility
48 metrics for each baseline model. We think relative accuracy is the metric most well matched to the developmental
49 literature, which compares reactions to surprising scenes directly to those in matched controls.

*Stahl, Feigenson. Observing the unexpected enhances infants’ learning and exploration. Science 2015

†Jonschkowski et al. Differentiable particle filters: End-to-end learning with algorithmic priors. RSS 2018

‡Battaglia et al. Simulation as an engine of physical scene understanding. PNAS 2013

§Teglas et al. Pure reasoning in 12-month-old infants as probabilistic inference. Science 2011

¶Lieder et al. Over-representation of extreme events in decision making reflects rational use of cognitive resources. Psych. Rev. 2018