# Modeling habituation in infants and adults using rational curiosity over perceptual embeddings

**Gal Raz**[*]
MIT
galraz@mit.edu

**Anjie Cao**[*]
Stanford University
anjiecao@stanford.edu

**Rebecca Saxe**
MIT
saxe@mit.edu

**Michael C. Frank**
Stanford University
mcfrank@stanford.edu

## Abstract

From birth, human infants engage in intrinsically motivated, open-ended learning, mainly by deciding what to attend to and for how long. Yet, existing formal models of the drivers of looking are very limited in scope. To address this, we present a new version of the Rational Action, Noisy Choice for Habituation (RANCH) model. This version of RANCH is a stimulus-computable, rational learning model that decides how long to look at sequences of stimuli based on expected information gain (EIG). The model captures key patterns of looking time documented in the literature, habituation and dishabituation. We evaluate RANCH quantitatively using large datasets from adult and infant looking time experiments. We argue that looking time in our experiments is well described by RANCH, and that RANCH is a general, interpretable and modifiable framework for the rational analyses of intrinsically motivated learning by looking.

## 1 Introduction

Human infants have limited motor capacity, so they engage in intrinsically motivated open-ended learning mainly by deciding what to attend to and for how long. Developmental psychologists have long capitalized on this fact, probing infants' mental representations through their looking behavior [1, 2, 5]. In looking time experiments, infants are repeatedly shown one stimulus until their looking time decreases significantly (i.e. habituation), and then shown a novel test stimulus. Infants look longer at the novel stimulus (i.e. dishabituation). Why do infants look longer at the novel stimulus? One intuition is that infants look longer when they recognize learning opportunities. In this paper, we offer a formal model of this connection between looking and learning: Rational Action, Noisy Choice for Habituation Model (RANCH). We validate the model using behavioral datasets from both adults and infants.

Existing models of looking behaviors in infants leverage event probabilities to connect information theoretic measures with looking behaviors [8, 9]. For example, Poli et al. [13] utilized a paradigm in which infants were shown sequences of events until they looked away. A rational learning model of event probabilities (a Dirichlet-Multinomial model) computed various information theoretic metrics such as surprise and KL-divergence associated with each event. The results suggested that infants were looking longest at optimally informative stimuli. However these models only retrospectively fit infants' behaviors, without modelling the online learning and decision making underlying infant looking. Also, these models only apply to learning the probabilities of events, assuming that the conceptual content and boundaries of events are already given.

---

[*]Equal contribution.

The Rational Action, Noisy Choice for Habituation (RANCH) model was proposed to address these limitations [3]. The initial version of RANCH was a Bayesian concept learner that made moment-by-moment sampling decisions based on its expected information gain, a metric used for the rational analysis of information sampling [11, 12] and artificial agent behavior [15]. This model successfully predicted habituation and dishabituation patterns of adult participants. Nevertheless, as in prior work, the first version lacked a principled way of representing learning from the actual stimuli in the experiment. In the previous version of RANCH, stimuli were represented by ad hoc binary feature vectors. Without principled stimulus representations, any linkage to specific experimental results is necessarily mediated by ad-hoc, experiment-specific stipulations about how stimuli are encoded.

In this paper, we present a new version of RANCH that explicitly models learning of visual concepts, represented as convex regions in a continuous perceptual feature space [4, 6]. Importantly, RANCH is now fully stimulus-computable, so that it can generate predictions from raw pixels. It therefore instantiates a formal hypothesis about how humans go from perceiving stimuli to an abstract representation over which the ideal learner forms its representations. The learner's goal is to form a simple perceptual concept, but our framework lends itself to representing more open-ended types of learning. We evaluated the predictions of RANCH using behavioral datasets collected from adults and infants performing a simple perceptual learning task. Our results show that RANCH captures the attentional patterns of humans across development, and a developmental comparison between best-fitting parameters provides deep insights about the different priors that adults and infants bring to bear on perceptual learning.

## 2 Model

RANCH is a Bayesian perception/action model in which a learner makes optimal perceptual sampling decisions [3]. The learner learns the location and variance of a Gaussian category in a perceptual space by observing a series of noisy perceptual samples from a sequence of stimuli; the model makes decisions about how many samples to receive of each stimulus before disengaging. We describe RANCH's perceptual representation, learning model, and decision model in turn.

**Perceptual representation** The current version of RANCH extends previous versions by using stimulus-computable perceptual embeddings obtained from a model presented recently by Lee & DiCarlo [10]. This deep neural network uses ResNet50 for encoding and then projects the final layer onto a lower-dimensional embedding. This final projection is "perceptually-aligned", in that it was trained to match perceptual dissimilarity matrices derived from human adult reaction times in a 2-AFC match-to-sample task. We use these projections into a perceptually-aligned embedding space as a principled low-dimensional representation of stimuli, over which our learning model can form perceptual concepts. A visualization of experimental stimuli in the embedding space can be seen in Figure 1A.
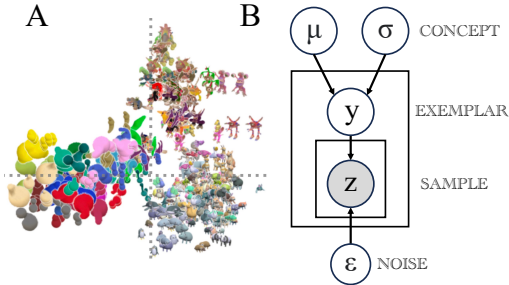


Figure 1: Panel A: Stimulus embeddings in PC-space. Panel B: Plate diagram of the learning model, which learns distributions over a mean and standard deviation from noisy perceptual samples.

**Learning model** RANCH's goal is to learn a concept in the perceptual embedding space described above, through noisy perceptual samples from a stimulus. The concept is parameterized by $\mu, \sigma$, which represents beliefs about the location and variance of the presented concept in the embedding space. This concept $\mu, \sigma$ generates exemplars $y$: exemplars of the concept. RANCH observes repeated noisy samples $\bar{z}$ from each exemplar. For any sample $z$ from an exemplar $y$, the model expects the observation to get corrupted by zero-mean noise, represented by $\epsilon$. A plate diagram is shown in Figure 1B. We used a normal-inverse-gamma prior on the concept, the conjugate prior for a normal with unknown mean and variance, on the concept parameterized as $\mu_p, \nu_p \, \alpha_p, \beta_p$. Still, applying perceptual noise to $y$ breaks the conjugate relation, so we computed approximate posteriors using grid approximation over $\mu, \sigma$ and $\epsilon$.

**Decision model** To decide whether to request an additional sample from the same stimulus, RANCH computes expected information gain (EIG) of the next sample. EIG is computed as the product of the posterior predictive probability of the next sample and the information gained conditioned on that next sample, via a grid approximation of possible subsequent samples. RANCH then makes a softmax choice (with temperature = 1) between next-sample EIG and a constant "environmental EIG" assumed to be the amount of information to be gained via looking away from the stimulus.

## 3 Behavioral data

We evaluated the predictions of RANCH using adapted versions of two previously-published behavioral datasets: adults (N = 380) from [3] and infants (N = 92) from [14]. The adult behavioral dataset was collected using an online self-paced looking time paradigm, where participants were instructed to watch blocks of six animations, consisting mostly of one animation (the background), and a second animation (the deviant) being shown on the 2nd, 4th or 6th trial, or not at all. Adults indicated when they wanted to continue to the next stimulus with a keypress. The infant dataset was collected using a novel online looking time paradigm. Infants watched blocks consisting of familiarization to one animation for different exposure durations (between 5 and 45 seconds), followed by a test trial which either showed the same stimulus again or a new stimulus. We measured looking time as the total time infants looked at a test trial until the first 2-second lookaway. The experiments used distinct stimulus sets.

## 4 Model fitting methods

We tailored the model predictions to each dataset by creating an "adult experiment" and an "infant experiment" for the model. In the adult experiment, the model decided after each sample whether to keep looking at the current stimulus, or move on to the next trial, for six successive trials. A deviant stimulus was presented on the 2nd, 4th, or 6th trial, or was absent. In the infant experiment, we created a "familiarization phase" where the model was presented with a fixed number of samples of the background stimulus; and then a "test phase" of either the background of the deviant stimulus where the model decided after each sample whether to keep looking or move on.

We also compared the RANCH model with lesioned models to test model assumptions' relevance. In the "No Learning" lesion, the model makes sampling decision randomly rather than based on learning. In the "No Noise" lesion, the model assumes that each observation is noiseless. These two lesioned models were used as comparisons in evaluating the earlier version of RANCH [3].

For each dataset, we conducted an iterative grid search across free parameters for each dataset to select the best-fitting parameters. The free parameters we searched over were the priors over $\mu$, $\sigma$ and $\epsilon$, as well as the actual noise $\epsilon$. For the adult dataset, we used 10% of the data as the training set to select parameters. For each set of parameters, we calculated a Pearson's r between the model outputs and the training dataset. We selected the parameter set with the highest Pearson's r between condition means and participant means as the best-fitting model. We compared results under these parameters to the remaining 90% behavioral dataset for adults. Given the sparsity of infant data, we used a leave-one-out cross-validation procedure, where we iteratively fit parameters to all but one infant, and then generated trial-wise predictions for the left-out infant. For each fold of the cross-validation, we fit a linear model using model predictions and block number as predictors (an experimental variable not accounted for by RANCH, but associated with fatigue and decreasing looking times in infants). We then used the resulting coefficients to predict looking times of the held-out infant. We also computed a noise-ceiling using a linear model that had access to the experimental conditions themselves (test type and prior exposure duration), following the same cross-validation procedure.

## 5 Results

**Adult fit to data** RANCH provided reasonable qualitative fits with the behavioral data from adults (Figure 2A). The best fitting parameters were ($\mu_p = 0$, $\nu_p = 1$, $\alpha_p = 1$, $\beta_p = 1$, $\epsilon = 0.0001$). Quantitatively, the model achieved a good fit with the behavioral data ($r = 0.95$ [0.90, 0.98], $RMSE$ = 3573.076 [2789.96, 4970.69]). The model fit was significantly better than both the No Learning model ($r = -0.16$ [-0.53, 0.26], $RMSE$ = 3583.762 [2798.3, 4985.56]) and the No Noise model ($r = 0.69$ [0.39, 0.85], $RMSE$ = 3584.49 [2798.87, 4986.57]). Interestingly, RANCH did not capture the complexity effect reported in the original publication [3]. Adults looked longer at the more complex visual stimuli ($\beta = -0.05$, $SE = 0.02$, $t = -2.3$, $p = 0.02$) but the effect of complexity was
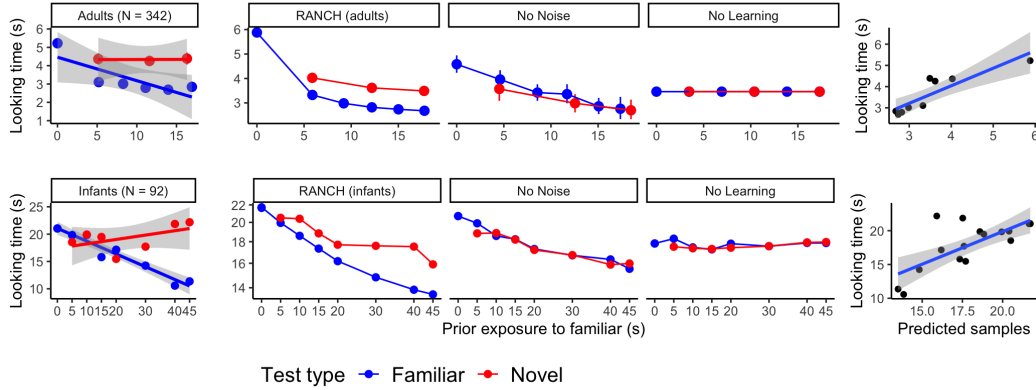
Figure 2: Behavior and model fits of RANCH and baseline models for adults and infants. The first row and the second row represents the behavioral dataset and the model fits for the adult dataset and the infant dataset respectively. For the four columns on the left, X-axis shows the accumulated prior exposure to the familiar stimuli in seconds. Y-axis represents the looking time to the test stimuli. Blue dots represent the test stimuli being familiar and the red dots represent the test stimuli being novel. The rightmost show the fit between the RANCH predicted results on the X-axis and the looking time results on the y-axis.

not significant in the model outputs ($\beta$ = -0.001, $SE$ = 0.002, $t$ = -0.42, $p$ = 0.67). This discrepancy suggests that stimulus complexity is not currently represented in our low dimensional embedding.

**Infant data fit** For the infant data, the best parameters were different from adults ($\mu_p = 0$, $\nu_p = 2$, $\alpha_p = 5$, $\beta_p = 15$, $\epsilon = 0.00001$). Using these parameters, RANCH provided a good fit to infants' looking times (Figure 2B; $r$ = 0.72 [0.36, 0.90], $RMSE$ = 2357.992 [1756.16, 3588.70]), better than the No Learning ($r$ = 0.48 [-0.02, 0.79], $RMSE$ = 3422.515 [2548.99, 5208.83]) and No Noise model ($r$ = 0.11 [-0.40, 0.58], $RMSE$ = 3015.594 [2245.92, 4589.52]). The main behavioral patterns reported in [14], habituation and dishabituation, were reproduced by RANCH. In our leave-one-out cross-validation analysis, we found that our trialwise predictions showed weaker differentiation between RANCH ($r$ = 0.34 [0.25, 0.42], $RMSE$ = 13.74 [12.01, 16.06]) and the two baseline models (No Learning: $r$ = 0.32 [0.22, 0.39], $RMSE$ = 13.90 [12.14, 16.24]; No Noise: $r$ = 0.32 [0.23, 0.40], $RMSE$ = 13.86 [12.11, 16.20]). However, the noise ceiling was also quite low ($r$ = 0.35 [0.27, 0.43], $RMSE$ = 13.67 [11.94, 15.97]), suggesting that trial-wise data is generally noisier in infants compared to adults. In other words, RANCH's performance is close to maximal, and mostly limited by the data.

## 6  Discussion

We present a modular, computational framework for the rational analysis of intrinsically motivated learning through looking. In the current implementation, we made specific decisions about perceptual representation, learning model, and decision model. This modular setup of RANCH easily lends itself to modification of any of these components and investigating the effects on sampling behavior.

One unique advantage of the RANCH model is the parameter interpretability. The priors on $\mu$ and $\sigma$ were parameterized by a normal inverse-gamma prior with $\mu$, $\nu$, $\alpha$ and $\beta$, the conjugate prior to a normal distribution with unknown mean and variance. While the mean was fixed to be 0 ($\mu$), variation in the the other parameters express distinct hypotheses about the precision of the location of the concept ($\nu$), as well as the variance of the concept ($\alpha$ and $\beta$). The best-fitting priors for adult and infants therefore lend themselves to comparison. The most striking difference was in the prior on the concept variance: While the adult version of RANCH achieved the highest fit with parameters $\alpha = 1$ and $\beta = 1$, infants' achieved the highest fit with $\alpha = 5$ and $\beta = 15$, indicating far wider prior variance. This result is consistent with previous proposals that infants bring less refined prior knowledge to bear on learning, and that the role of development is to refine those priors [7, 16].

Despite the generally good fit between model predictions and data, there is one qualitative mismatch to both the adult and infant datasets. RANCH predicts that looking to deviant stimuli will gradually decrease (i.e. a negative slope) as a function of prior exposure to the background stimulus. In both

infants and adults, looking to the deviant stimuli stayed the same (i.e. flat slope). Such qualitative deviations from the data point to differences in the computation underlying attentional decision-making between humans and the current version of RANCH. In this case, one possibility is that humans use different decision models to link learning and looking. In previous work we explored the effect of using different decision models including both the optimal, forward-looking EIG used here and several simpler-to-compute proxies, surprisal and KL-divergence [3]. Surprisal, unlike EIG and KL-divergence, typically results in an increasing looking toward deviant stimuli as a function of previous exposure to the background stimulus, and thus might better approximate human behavior in our paradigm. Alternatively, it may be that the description of humans as learning single concepts may be oversimplified. Formulating the learning problem as hierarchical, where a learner attempts to understand how many concepts are present, and which concept to attribute the current observation, may result in a closer fit to the data.

Furthermore, the imperfect differentiation between model fits in the infant cross-validation analysis points to two issues: First, noise in infant data makes it hard to achieve good fit with trial-wise predictions. This is corroborated by the low noise ceiling, which suggests that even with perfect information about the experimental condition, predictive power is generally low in infant looking time data. Second, it is likely that an overall model fit to the entire dataset is not the most sensitive measure of RANCH's performance. RANCH's predictions are only importantly different from the baseline models in a subset of the data: after longer exposures, when looking to novel items is higher than to familiar items. Other aspects of the prediction, like the effect of block number (a strong predictor of a decrease in looking time due to fatigue), and looking time for shorter prior exposures, are shared between RANCH and baseline models. Future analyses should therefore consider more sensitive human-model comparison methods which focus their assessment on the crucial parts of RANCH's behavior.

Overall, RANCH instantiates a hypothesis about how looking achieves rational intrinsically motivated learning in a simple perceptual learning task. Moving forward, its modular nature would lend itself to capturing how humans "look" at more open-ended learning problems. By changing the perceptual representation and learning model to reflect more open-ended learning problems, RANCH provides a general framework in which one could instantiate hypotheses about how humans use looking to learn, in general. We believe that using this framework can move the study of looking towards a predictive science, formally linking it to its underlying learned representations.

## References

[1] Richard N Aslin. What's in a look? *Developmental Science*, 10(1):48–53, 2007.

[2] Renee Baillargeon, Elizabeth S Spelke, and Stanley Wasserman. Object permanence in five-month-old infants. *Cognition*, 20(3):191–208, 1985.

[3] Anjie Cao, Gal Raz, Rebecca Saxe, and Michael C Frank. Habituation reflects optimal exploration over noisy perceptual samples. *Topics in Cognitive Science*, 15(2):290–302, 2023.

[4] Shimon Edelman. Representation is representation of similarities. *Behavioral and Brain Sciences*, 21(4):449–467, 1998.

[5] Robert L Fantz. Pattern vision in newborn infants. *Science*, 140(3564):296–297, 1963.

[6] Martin N Hebart, Charles Y Zheng, Francisco Pereira, and Chris I Baker. Revealing the multidimensional mental representations of natural objects underlying human similarity judgements. *Nature Human Behaviour*, 4(11):1173–1185, 2020.

[7] Frank C Keil. *Concepts, Kinds, and Cognitive Development*. MIT Press, 1992.

[8] Celeste Kidd, Steven T Piantadosi, and Richard N Aslin. The goldilocks effect: Human infants allocate attention to visual sequences that are neither too simple nor too complex. *PloS ONE*, 7(5):e36399, 2012.

[9] Celeste Kidd, Steven T Piantadosi, and Richard N Aslin. The goldilocks effect in infant auditory attention. *Child Development*, 85(5):1795–1804, 2014.

[10] Michael J Lee and James J DiCarlo. An empirical assay of view-invariant object learning in humans and comparison with baseline image-computable models. *bioRxiv*, pages 2022–12, 2023.

[11] Doug Markant and Todd Gureckis. Does the utility of information influence sampling behavior? In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 34, 2012.

[12] Mike Oaksford and Nick Chater. A rational analysis of the selection task as optimal data selection. *Psychological Review*, 101(4):608, 1994.

[13] F Poli, G Serino, RB Mars, and S Hunnius. Infants tailor their attention to maximize learning. *Science Advances*, 6(39):eabb5053, 2020.

[14] Gal Raz, Anjie Cao, Minh Khong Bui, Michael C Frank, and Rebecca Saxe. No evidence for familiarity preferences after limited exposure to visual concepts in preschoolers and infants. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 45, 2023.

[15] Jürgen Schmidhuber. A possibility for implementing curiosity and boredom in model-building neural controllers. In *Proc. of the international conference on simulation of adaptive behavior: From animals to animats*, pages 222–227, 1991.

[16] Elizabeth Spelke. Initial knowledge: Six suggestions. *Cognition*, 50(1-3):431–445, 1994.