# Distilling Symbolic Priors for Concept Learning into Neural Networks

**Ioana Marinescu,[1] R. Thomas McCoy,[2] Thomas L. Griffiths[1,3]**

`ioanam@princeton.edu`, `tom.mccoy@yale.edu`, `tomg@princeton.edu`
[1]Department of Computer Science, Princeton University
[2]Department of Linguistics, Yale University
[3]Department of Psychology, Princeton University

## Abstract

Humans can learn new concepts from a small number of examples by drawing on their inductive biases. These inductive biases have previously been captured by using Bayesian models defined over symbolic hypothesis spaces. Is it possible to create a neural network that displays the same inductive biases? We show that inductive biases that enable rapid concept learning can be instantiated in artificial neural networks by distilling a prior distribution from a symbolic Bayesian model via meta-learning, an approach for extracting the common structure from a set of tasks. We use this approach to create a neural network with an inductive bias towards concepts expressed as short logical formulas. Analyzing results from previous behavioral experiments in which people learned logical concepts from a few examples, we find that our meta-trained models are highly aligned with human performance.

**Keywords:** inductive bias; concept learning; meta-learning

## Introduction

People can make rich inferences from remarkably little data. For instance, consider the domain that we focus on in this work: concept learning. People can learn a new concept from just a few examples (e.g., Bloom, 2002; Xu & Tenenbaum, 2007; Lake & Piantadosi, 2020). Figure 1 illustrates the concept *green or triangle*; the objects are labeled *yes* if they are an instance of this concept and *no* otherwise. Given such data, people can rapidly infer what concept underlies the labeling, for a wide range of concepts (Bruner, Goodnow, & George, 1956; Feldman, 2000; Goodman, Tenenbaum, Feldman, & Griffiths, 2008; Piantadosi, Tenenbaum, & Goodman, 2016).

Bayesian models have successfully been used to capture the human ability to learn from few examples (Goodman et al., 2008; Piantadosi et al., 2016). Although such models are effective at explaining human learning behavior, they are often built on explicitly symbolic hypotheses, which have been



Figure 1: Concept learning from examples. The concept underlying the labels is *green or triangle*: objects are labeled *yes* if they are *green* or a *triangle* and *no* otherwise. In our experiments, learners are given a set of labeled examples such as these and are then required to predict the labels for an additional set of examples.

argued against as components of mechanistic cognitive theories by advocates of artificial neural networks (McClelland et al., 2010). Recent neural network models powered by deep learning (LeCun, Bengio, & Hinton, 2015) have been shown to be extremely effective at solving a variety of problems, including learning to classify stimuli into categories. For instance, for the popular ImageNet classification challenge (Deng et al., 2009), the best-performing systems are neural networks (e.g., Yu et al., 2022; Chen et al., 2023). However, the neural networks that are used for these tasks typically require far more training examples than humans need. For instance, ImageNet contains approximately 1200 training examples per class. When standard neural networks are instead trained on smaller amounts of data, they often generalize poorly (Lake, Salakhutdinov, & Tenenbaum, 2015; Hoiem, Gupta, Li, & Shlapentokh-Rothman, 2021).

Can we use the symbolic representations from Bayesian models to enable neural networks to acquire concepts from smaller amounts of data? Answering this question is important in order to evaluate the potential of artificial neural networks as models of human learning, providing a path towards identifying more plausible cognitive mechanisms that are informed by Bayesian modeling. This problem is at its core about *inductive biases*—the factors that guide how a learner generalizes. Bayesian models are effective at capturing human inductive biases (e.g., Goodman et al., 2008; Piantadosi et al., 2016). Therefore, if neural networks can be given the same inductive biases as Bayesian models, then they could match human learning behavior as effectively as Bayesian models do. How, then, can a Bayesian model's inductive biases be given to a neural network?

In this work, we use meta-learning to distill the inductive biases from a Bayesian model into a neural network. The Bayesian model serves as our characterization of human inductive biases; in the Bayesian framework, an inductive bias can be characterized as a prior probability distribution over hypothesized concepts. Meta-learning is then used to turn this prior into a training regime for a neural network. In meta-learning, a system is trained on many related tasks in order to learn the underlying abstractions that are common across these tasks, thereby giving it inductive biases that enable it to readily learn new tasks (Thrun & Pratt, 1998; Schmidhuber, 1987). In our application of meta-learning, each task is a concept sampled from our Bayesian prior, such that meta-
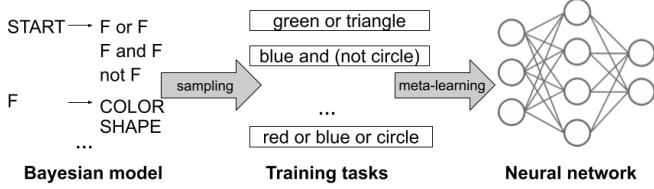
Figure 2: Learning a prior from a Bayesian model and distilling it into a neural network

$$S \rightarrow \forall x \quad l(x) \Leftrightarrow D_{\text{top}} \qquad C \rightarrow P \wedge C$$
$$D_{\text{top}} \rightarrow C_{\text{top}} \vee D \qquad\qquad C \rightarrow \text{True}$$
$$C_{\text{top}} \rightarrow P \wedge C \qquad\qquad P \rightarrow F_i$$
$$D \rightarrow C_{\text{top}} \vee D \qquad\qquad F_i \rightarrow f_i(x) = 1$$
$$D \rightarrow \text{False} \qquad\qquad F_i \rightarrow f_i(x) = 0$$

Figure 3: DNF grammar from Goodman et al. (2008). $i \in \{1, 2, ... \text{number of features}\}$

learning should result in the neural network internalizing this Bayesian prior. In this way, meta-learning serves as a formal framework for bridging probabilistic models of human concept learning and deep neural networks (see Figure 2).

In our experiments, we used this approach to distill the prior from the Rational Rules model (Goodman et al., 2008) into a neural network. Across a range of previous behavioral experiments, the resulting neural network displayed generalization behavior that was highly aligned with both humans and the Rational Rules model. In contrast, a standard neural network fared poorly across these conditions. These results demonstrate that meta-learning can successfully be used to integrate Bayesian and connectionist models of concept learning, enabling the creation of models that combine the complementary strengths of these two approaches.

## Background

### Bayesian concept learning

To define our target inductive bias, we use the prior from the Bayesian model created by Goodman et al. (2008). This prior is defined using the context-free grammar (CFG) shown in Figure 3.[1] CFGs are best known for being used to characterize the syntax of natural languages and to generate sentences in a language; in the case of Goodman et al.'s model, the CFG instead generates concept definitions.

The derivation of a concept definition starts with the symbol $S$. This symbol is then recursively expanded using the rules of the grammar until it has yielded a complete

---

[1] Goodman et al. (2008) denote their CFG in a different way that does not include the rules subscripted with *top*. However, we reran their experiments using the package Fleet (Yang & Piantadosi, 2022) and were able to replicate their results only when including the *top* rules. We therefore infer that Goodman et al.'s grammar included something akin to the *top* rules that was omitted from their diagrams for conciseness.

expression. For instance, the $S$ would first be expanded into $\forall x \; l(x) \Leftrightarrow D_{\text{top}}$. We would then expand $D_{\text{top}}$ to yield $\forall x \; l(x) \Leftrightarrow C_{\text{top}} \vee D$, after which $C_{\text{top}}$ and $D$ would each need to be expanded. Eventually this would produce a definition such as $\forall x \; l(x) \Leftrightarrow f_3(x) = 1 \wedge f_2(x) = 0 \wedge \text{True} \vee \text{False}$. This expression defines the concept characterized by feature 3 having a value of 1 and feature 2 having a value of 0 (other parts of the definition—$\forall x \; l(x)$, True, False—are included for formal reasons but do not affect the definition).

In order to make these derivations probabilistic, each rule in the grammar has a probability attached to it, and expansions are then sampled according to those probabilities.

### Neural network architectures for few-shot learning

Although standard neural network systems are not effective at learning from few examples, some alternative model architectures have been proposed that are tailored for this purpose. Examples include Matching Networks (Vinyals, Blundell, Lillicrap, Wierstra, et al., 2016) and Prototypical Networks (Snell, Swersky, & Zemel, 2017). Our approach differs from these in modifying the training regime rather than the neural network's architecture or processing mechanisms.

### Meta-learning

Besides specially-designed architectures, the other main approach for enabling neural networks to learn from few examples is meta-learning. In this work, we adopt a type of meta-learning called Model-Agnostic Meta-Learning (MAML) (Finn, Abbeel, & Levine, 2017). MAML is a machine learning technique that focuses on training a model in such a way that it can quickly adapt to new, unseen tasks with minimal data. MAML aims to learn an effective initialization for a neural network so that fine-tuning on a new task requires fewer iterations or samples. The model is trained on a diverse set of tasks (from the same task distribution) during a meta-training phase, learning the set of initial parameters. Upon encountering an unseen task, the model's learned initialization can be rapidly fine-tuned with a small amount of data from a new task. Meta-learning extracts common structure from a set of tasks, resulting in a system that has inductive biases aligned with those tasks.

MAML finds parameters $\theta$ that minimize the summed loss over tasks:

$$\underset{\theta}{\arg\min} \sum_{i=1}^{N} \mathcal{L}_{\text{val}}\left(\theta - \alpha \nabla_\theta \mathcal{L}_{\text{train}}(\theta)\right), \qquad (1)$$

where $\theta$ represents the model parameters, $N$ is the number of tasks in the meta-training set, $\mathcal{L}_{\text{val}}$ is the validation loss used to update the parameters, and $\alpha$ is the step size or learning rate for the inner loop optimization.

Prior work has established that meta-learning is a promising tool for building models of human cognition. In some cases, it is used to model ways in which humans themselves meta-learn (Griffiths et al., 2019; Wang, 2021; Kumar, Dasgupta, Cohen, Daw, & Griffiths, 2021). In other cases—as in

our work—meta-learning is used as a tool to create neural networks that have human-like inductive biases, without making the claim that meta-learning is *how* humans arrived at those biases; this framing has been used to improve how well neural networks can learn symbolic rules (Lake, 2019; McCoy, Grant, Smolensky, Griffiths, & Linzen, 2020; Lake & Baroni, 2023). Meta-learning is also powerful for connecting neural networks to other research traditions in cognitive science: Binz et al. (2023) discuss how meta-learning can connect neural networks to rational analysis, and our goal here is to connect neural networks and Bayesian models—an application of meta-learning previously used by McCoy and Griffiths (2023) in the domain of language, building upon a theoretical connection between MAML and hierarchical Bayesian models (Grant, Finn, Levine, Darrell, & Griffiths, 2018). Note that Prior-data Fitted Networks (Müller, Hollmann, Arango, Grabocka, & Hutter, 2022) also inject Bayesian priors into neural networks, but they do so in a different manner that uses standard learning rather than meta-learning.

## Distilling Priors into a Neural Network

To distill a Bayesian prior into a neural network, we use the three-step approach illustrated in Figure 2. We first use a probabilistic model to define the inductive bias that we wish to give to a neural network. Specifically, this probabilistic model is the component of a Bayesian model that defines its prior, and it gives a distribution over possible concepts. Internalizing this distribution would help a neural network learn more efficiently by shaping the hypothesis space that it searches over. In order to make this distribution accessible to a neural network, we sample a large number of concepts from the probabilistic model for use as training data. We then have a neural network meta-learn from these sampled tasks to give it an inductive bias that matches the prior distribution we started with. This approach was first used in previous work that modeled language learning (McCoy & Griffiths, 2023); the current paper applies the approach in a new domain, namely concept learning.

### Sampling data

We sample many concepts from the grammar in Figure 3 in order to serve as meta-learning data for a neural network. For each concept, we first sample the production probabilities from a Dirichlet distribution (note that, following Goodman et al. (2008), the CFG production rules are fixed, but their probabilities are not). We then sample a concept definition from the grammar with those probabilities. We must next generate both a training set and test set for this concept. The training set contains up to 20 examples, randomly sampled from the 16 possible objects. We use multi-step loss (Antoniou, Edwards, & Storkey, 2019), which has the same effect as providing the learner with training sets of various sizes, and we also allow examples to repeat within one episode. The test set is made of all 16 possible objects. Each example is labeled True if it obeys the rule defining the concept or False otherwise, except that with probability

$e^{-b}/(1+e^{-b})$ we flip the label of the example (this is how Goodman et al.'s model incorporates the possibility of outliers). We denote this outlier parameter $b \in \{1, 2, .., 8\}$.

The above procedure produces the data for a single concept. We repeat this procedure to produce data for 10,000 concepts, which will serve as the data from which the neural network will meta-learn.

### Inductive bias distillation via meta-learning

The neural network internalizes the Bayesian prior by meta-learning through experience with many similar tasks, each corresponding to learning one concept. We framed the concept learning problem as a binary classification task and assumed the desired concept was learned if specific examples were classified correctly. The architecture we used in our experiments is a multi-layer perceptron (MLP). We generated 10,000 episodes for meta-training, 100 for meta-validation, and 100 for meta-testing. We start with a baseline MLP which uses hidden size 128 and 5 layers, dropout 0.1, 1 epoch, outer learning rate 0.0005, and inner learning rate 0.1. We also use a modified version which has hidden size 256 and skip connections. We use the term **prior-trained** network to refer to a network that has undergone this distillation process; a **standard** network is one that has not undergone this process.

## Evaluating the Model

To evaluate the model we examined its performance in accounting for a set of behavioral results from experiments that had previously been used to assess the Rational Rules model (Medin & Schaffer, 1978; Medin & Schwanenflugel, 1981; Shepard, Hovland, & Jenkins, 1961; McKinley & Nosofsky, 1993). These experiments involve learning logical concepts corresponding to objects that have three or four Boolean features. For example, a concept could be represented by "$f_1(x) = 0$ and $f_3(x) = 1$" meaning that the first feature has a value of 0 and the third feature, 1. The human experiments required the participants to classify the objects into binary categories, according to the rule they inferred. We test our model against the Rational Rules model and human data.

### Category structure of Medin and Schaffer (1978)

In Table 1, we consider the category structure of Medin and Schaffer (1978), using human data from Nosofsky, Palmeri, and Mckinley (1994), and compare the behavior of humans, the Rational Rules model with $b = 1$, the prior-trained neural network, and the standard neural network.

Each object has four binary features. An object can therefore be represented by a sequence of four zeroes or ones; e.g., "0111" means that the object has a value of 0 for feature 1, a value of 1 for feature 2, a value of 1 for feature 3, and a value of 1 for feature 4. Learners (whether humans or computational models) were trained on 9 of the 16 possible objects: they were shown an object's feature values along with a label indicating whether the object belonged to category A or category B. The training examples are the ones labeled "A" or "B" in the Object column of Table 1. Participants were then

Table 1: The category structure of Medin and Schaffer (1978), with the human data of Nosofsky et al. (1994), the predictions of the Rational Rules model with $b = 1$, the predictions of the prior-trained MLP with $b = 2$, and the predictions of a standard (non-prior-trained) MLP.

| Object | Feature values | Human | RR$_{DNF}$ | Prior trained | Standard |
|--------|---------|-------|------------|---------------|----------|
| A1 | 0001 | 0.77 | 0.82 | 0.71 | 0.52 |
| A2 | 0101 | 0.78 | 0.81 | 0.76 | 0.52 |
| A3 | 0100 | 0.83 | 0.92 | 0.84 | 0.52 |
| A4 | 0010 | 0.64 | 0.61 | 0.69 | 0.52 |
| A5 | 1000 | 0.61 | 0.61 | 0.70 | 0.52 |
| B1 | 0011 | 0.39 | 0.47 | 0.40 | 0.52 |
| B2 | 1001 | 0.41 | 0.47 | 0.45 | 0.52 |
| B3 | 1110 | 0.21 | 0.21 | 0.22 | 0.52 |
| B4 | 1111 | 0.15 | 0.07 | 0.14 | 0.52 |
| T1 | 0110 | 0.56 | 0.57 | 0.56 | 0.52 |
| T2 | 0111 | 0.41 | 0.44 | 0.34 | 0.52 |
| T3 | 0000 | 0.82 | 0.95 | 0.84 | 0.52 |
| T4 | 1101 | 0.40 | 0.44 | 0.41 | 0.52 |
| T5 | 1010 | 0.32 | 0.28 | 0.39 | 0.52 |
| T6 | 1100 | 0.53 | 0.57 | 0.60 | 0.52 |
| T7 | 1011 | 0.20 | 0.13 | 0.19 | 0.52 |

shown all 16 concepts (the 9 they had been trained on, plus the 7 they had not been—labeled T1 through T7).

The last 4 columns of the table show the probability that the human or model assigned category A to a given object. For instance, humans had 0.77 probability of labeling A1 as belonging to category A. There are multiple possible rules that a learner could consider when trained on this dataset. For instance, the learner could learn a complicated rule that captures all of the training data perfectly, or it could learn a rule that misclassifies some of the training data but is simpler, perhaps under the assumption that the misclassified examples are outliers. One example of such a simpler rule would be to assume that objects with a value of 0 for feature 1 are in category A, while those with a value of 1 for feature 1 are in category B. This rule correctly classifies 7 of the 9 training examples. This interplay between multiple possible rules can be seen in the varying probability levels that are assigned to different training examples; e.g., even though learners have seen A1 through A5 as training examples labeled A, not all of them have the same probability of being labeled A by the learners; presumably, this is because some of them are consistent with more of the candidate rules than others.

The baseline architecture with a random initialization obtains $R^2 = 0$. Using the same model with the meta-learned initialization, we find $R^2 = 0.95$ for the prior-trained neural network and humans and $R^2 = 0.92$ for the prior-trained neural network and the Rational Rules model, demonstrating that meta-trained models can closely match the human results; $R^2 = 0.98$ between the Rational Rules model and humans. The specific testing examples as well as the error probability for each model are shown in Table 1. Figure 4 shows that the prior-trained neural network's predictions are highly correlated with the human predictions, while those of the standard neural network are not.
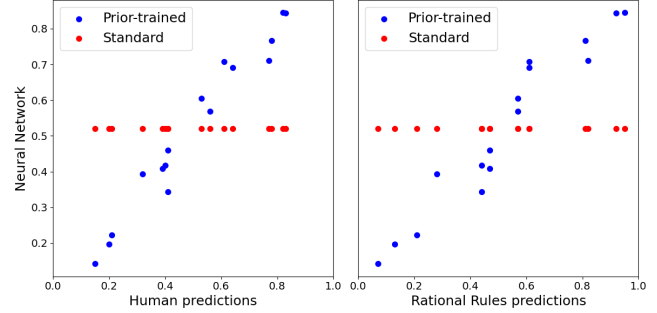


Figure 4: Predictions of prior-trained and standard neural networks vs humans and vs Rational Rules (data from Table 1).

## Linear Separability

We next consider the two concepts from Medin and Schwanenflugel (1981) which have four binary features and are shown in Table 2. Concept LS is linearly separable, meaning that it admits a linear discriminant boundary, while Concept NLS is not. Our model predicts that that Concept NLS is easier to learn, in agreement with the human and Rational Rules predictions. In Figure 5 we show the error probability: $1 - P(\text{true label} = \text{predicted label})$ in the Rational Rules model and two variations of our model; the output of our model is $P(\text{Category A})$. The two setups we study are: the modified model varying the outlier parameter $b \in \{1, 2...8\}$ with one epoch per episode at test time and the baseline model meta-trained with $b = 1$ varying $N \in \{1, 2...8\}$ epochs per episode at test time. Both neural network variants show trends similar to those found with the Rational Rules model, which in turn behaved similarly to humans (Goodman et al., 2008).

Table 2: Two concepts in Medin & Schwanenflugel (1981). Concept LS is linearly separable, Concept NLS is not.

| Concept LS | | Concept NLS | |
|------------|------------|-------------|------------|
| Category A | Category B | Category A | Category B |
| 1000 | 0111 | 0011 | 1111 |
| 0001 | 1000 | 1100 | 1010 |
| 0110 | 1001 | 0000 | 0101 |

Table 3: The six three-feature concepts with four positive and four negative examples, studied by Shepard et al. (1961).

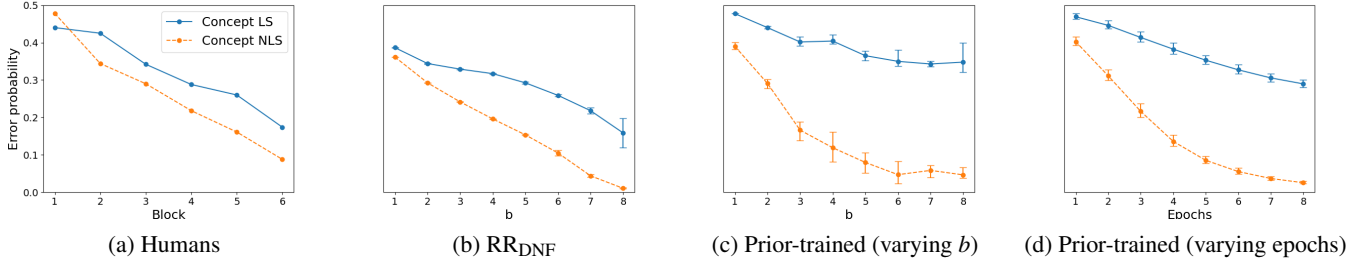| I | II | III | IV | V | VI |
|------|------|------|------|------|------|
| +000 | +000 | +000 | +000 | +000 | +000 |
| +001 | +001 | +001 | +001 | +001 | -001 |
| +010 | -010 | +010 | +010 | +010 | -010 |
| +011 | -011 | -011 | -011 | -011 | +011 |
| -100 | -100 | -100 | +100 | -100 | -100 |
| -101 | -101 | +101 | -101 | -101 | +101 |
| -110 | +110 | -110 | -110 | -110 | +110 |
| -111 | +111 | -111 | -111 | +111 | -111 |

Figure 5: Linearly and non-linearly separable concepts. Linearly separable Concept LS was more difficult to learn than Concept NLS, which is not linearly separable. Error probability in: (a) Humans, (b) Rational Rules, (c) Prior-trained modified model varying outlier parameter $b$, (d) Prior-trained baseline model varying number of epochs per episode during inference
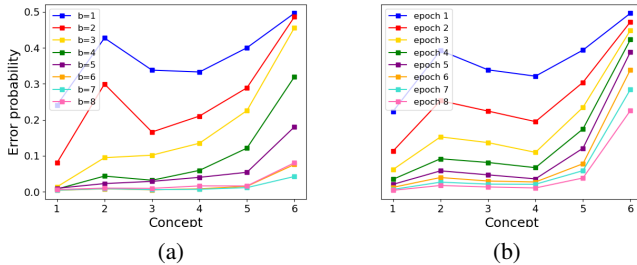


Figure 6: Error probability of the prior-trained model in learning concepts I-VI in Shepard et al. (1961). (a) Prior-trained modified model varying outlier parameter $b$. (b) Prior-trained baseline model varying number of epochs per episode.

## Shepard, Hovland and Jenkins (1961)

Shepard et al. (1961) compared difficulty in learning the six concepts shown in Table 3. The concepts have three Boolean features and are divided into four positive and four negative examples. The human performance in learning these concepts indicates the following ranking of difficulty in learning six concepts shown in Table 3: $I < II < III = IV = V < VI$. The Rational Rules error rates are $0\%, 17\%, 24\%, 24\%, 25\%, 48\%$, for $b = 3$. The ranking is consistent for all values of $b$ other than $b = 1$, where there is an inversion: $II > III$. Our prior-trained model exhibits the same trend, including the inversion $II > III$ for $b = 1$ and $b = 2$, which we show in Figure 6 for the modified model meta-trained varying $b$ and for the baseline model meta-trained with $b = 1$ and varying the number of epochs per episode at test time.

## Medin, Altom, Edelson, and Freko (1982)

This experiment follows the same setup as our experiment based on Medin and Schaffer (1978), meta-training our modified model with $b = 1$ and $b = 7$. In Table 4, we test our model on the category structure of Medin et al. (1982). The $RR_{DNF}$ model explains most of the variance in human judgments in the final stage of learning: $R^2 = 0.95$ when $b = 7$ and correlation with human judgments after one training block is $R^2 = 0.69$ when $b = 1$. Our model has $R^2 = 0.56$ when $b = 7$

Table 4: The category structure of Medin et al. (1982), with initial and final block mean human responses of McKinley and Nosofsky (1993), and the predictions of the Rational Rules model and our model at $b = 1$ and $b = 7$.

| Object | Feature values | Human initial block | $RR_{DNF}$ $b = 1$ | MAML $b = 1$ | Human final block | $RR_{DNF}$ $b = 7$ | MAML $b = 7$ |
|---|---|---|---|---|---|---|---|
| A1 | 1111 | 0.64 | 0.84 | 0.84 | 0.96 | 1 | 0.98 |
| A2 | 0111 | 0.64 | 0.54 | 0.67 | 0.93 | 1 | 0.97 |
| A3 | 1100 | 0.66 | 0.84 | 0.83 | 1 | 1 | 0.98 |
| A4 | 1000 | 0.55 | 0.54 | 0.66 | 0.96 | 0.99 | 0.96 |
| B1 | 1010 | 0.57 | 0.46 | 0.32 | 0.02 | 0 | 0.03 |
| B2 | 0010 | 0.43 | 0.16 | 0.15 | 0 | 0 | 0.02 |
| B3 | 0101 | 0.46 | 0.46 | 0.31 | 0.05 | 0.01 | 0.03 |
| B4 | 0001 | 0.34 | 0.16 | 0.15 | 0 | 0 | 0.02 |
| T1 | 0000 | 0.46 | 0.2 | 0.22 | 0.66 | 0.56 | 0.14 |
| T2 | 0011 | 0.41 | 0.2 | 0.26 | 0.64 | 0.55 | 0.32 |
| T3 | 0100 | 0.52 | 0.5 | 0.45 | 0.64 | 0.57 | 0.3 |
| T4 | 1011 | 0.5 | 0.5 | 0.48 | 0.66 | 0.56 | 0.38 |
| T5 | 1110 | 0.73 | 0.8 | 0.72 | 0.36 | 0.45 | 0.66 |
| T6 | 1101 | 0.59 | 0.8 | 0.74 | 0.36 | 0.44 | 0.79 |
| T7 | 0110 | 0.39 | 0.5 | 0.49 | 0.27 | 0.44 | 0.53 |
| T8 | 1001 | 0.46 | 0.5 | 0.51 | 0.3 | 0.43 | 0.63 |

and $R^2 = 0.66$ when $b = 1$; see the Discussion for a potential explanation for why the $b = 7$ result does a poorer job of replicating $RR_{DNF}$ than our other results do.

## Discussion

Humans can learn logical concepts rapidly, but it is natural to expect that this task will be difficult for artificial neural networks. First, these concepts are defined using discrete, symbolic, compositional rules, but neural networks are not a natural fit for symbolic domains (Fodor & Pylyshyn, 1988). Second, in the human experiments that we considered, such concepts were learned from small numbers of examples, yet standard neural networks usually require large quantities of training data. Indeed, we found that standard neural networks did a poor job at acquiring the concepts we considered. Nonetheless, we also found that inductive bias distillation made it possible for neural networks to perform well at learning logical

concepts: after we distilled a structured prior probability distribution into a neural network, it was able to learn logical concepts from few examples and in ways that aligned closely with human results across several experiments. These results therefore show that it is possible to develop connectionist instantiations of probabilistic models.

## Potential additional biases

The prior that we distilled into a neural network encodes a bias for simplicity: rules that can be expressed with a short logical description have a higher prior probability than rules that require a long description. Such a bias is well-supported by human experiments (Feldman, 2000; Neisser & Weene, 1962; Goodman et al., 2008), but it is far from the only bias involved in human concept learning. For instance, people also display a shape bias (Landau, Smith, & Jones, 1988), a whole-object bias (Markman, 1994), a basic-level bias (Rosch, Mervis, Gray, Johnson, & Boyes-Braem, 1976), and a mutual exclusivity bias (Markman & Wachtel, 1988). One direction for future work would be to incorporate these additional biases into our model, which could be done by augmenting the probabilistic model from which we sampled concepts such that the distribution of concepts that it produces reflects these biases. For instance, a shape bias could be instantiated by introducing an asymmetry between features so that concepts based on the shape feature would be sampled more often than concepts based on other features.

## Modeling the timecourse of learning

When our neural network model learns a concept, it does so incrementally, updating its parameters after each example it encounters. This incrementality means that the neural network can model the timecourse of learning more naturally than the Rational Rules model, which considers the entire set of training examples at once, such that modeling multiple points during the course of learning requires creating multiple separate instances of the model.

As one example, in some of the human experiments (Medin & Schwanenflugel, 1981; Shepard et al., 1961), participants were trained on the category using a blocked learning paradigm: each example in the training set was presented once per block, and blocks were repeated multiple times. To replicate this experiment with the Rational Rules model, a different instance of the model must be created for each step in the learning process. Specifically, as shown in Figure 5b, successive steps are modeled by increasing the $b$ hyperparameter, though the same result could also be achieved by fitting each instance to differing numbers of copies of training examples: the Rational Rules model with outlier parameter $b$ presented with $N$ identical blocks of examples is equivalent to the model presented with only one block, but with parameter $b' = b \cdot N$.

Like Rational Rules, the neural network can accommodate differences in $b$ (Figure 5c), but unlike Rational Rules it also accommodates an approach in which the same copy of the model is incrementally trained on multiple repeats of the training set (Figure 5d)—an approach that better captures incremental concept learning in humans.

## Challenges of learning a prior from data

One shortcoming of our approach arises from the fact that the neural network model learns its inductive bias from data. Therefore, the neural network is likely to have approximated this bias very well for the scenarios that were well-represented in the data, but the approximation may be poor in rarer scenarios. In our case, the frequency with which the network encountered a given concept was equal to that concept's prior probability, so the degree to which our model approximates the Rational Rules model is likely to decrease when a concept with a low prior probability is involved, since our model has seen few examples of how Rational Rules would generalize for such concepts.

We had one result where the neural network indeed provided a poor fit to Rational Rules, namely the $b = 7$ setting in Table 4. The low performance in this case is likely due to the fact that the rule to be learned, namely $(f_3 = 1 \wedge f_4 = 1) \vee (f_3 = 0 \wedge f_4 = 0)$, has a low prior probability (because its description needs to include a large number of feature specifications) and thus is outside the space in which the neural network's parameters are well-estimated. In contrast, the $b = 1$ setting in that table involves the same data but with a higher outlier probability; increasing the outlier probability leads Rational Rules to ignore some of the training examples as outliers in order to learn a simpler (i.e., higher-prior) rule, and having a high prior for the rule to be learned facilitates a stronger performance from the neural network model.

Enabling neural networks to generalize to low-probability scenarios is a challenging problem. Meta-learning is an active area of research in machine learning, and technological advances in this area may provide ways to overcome this problem, but for now it remains an important limitation for data-driven approaches to imparting inductive biases.

## Conclusion

While Bayesian models excel in certain aspects of concept learning, their traditional reliance on explicitly symbolic representations limits their implementational and algorithmic plausibility. On the other hand, neural networks do not require such representations but struggle to learn from small numbers of examples. Our approach bridges this gap by distilling probabilistic models into neural networks, leveraging the strengths of both paradigms. This integration enables neural networks to learn structured, compositional concepts from limited numbers of examples, providing new avenues for developing theories of concept learning that merge insights from Bayesian models and neural networks.

# References

Antoniou, A., Edwards, H., & Storkey, A. (2019). How to train your MAML. In *International Conference on Learning Representations*.

Binz, M., Dasgupta, I., Jagadish, A. K., Botvinick, M., Wang, J. X., & Schulz, E. (2023). Meta-learned models of cognition. *Behavioral and Brain Sciences*, 1–38.

Bloom, P. (2002). *How children learn the meanings of words*. MIT press.

Bruner, J. S., Goodnow, J. J., & George, A. (1956). *A study of thinking*. John Wiley & Sons.

Chen, X., Wang, X., Changpinyo, S., Piergiovanni, A., Padlewski, P., Salz, D., . . . Soricut, R. (2023). PaLI: A jointly-scaled multilingual language-image model. In *The Eleventh International Conference on Learning Representations*.

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*.

Feldman, J. (2000). Minimization of Boolean complexity in human concept learning. *Nature*, *407*(6804), 630–633.

Finn, C., Abbeel, P., & Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the International Conference on Machine Learning*.

Fodor, J. A., & Pylyshyn, Z. W. (1988). Connectionism and cognitive architecture: A critical analysis. *Cognition*, *28*(1-2), 3–71.

Goodman, N. D., Tenenbaum, J. B., Feldman, J., & Griffiths, T. L. (2008). A rational analysis of rule-based concept learning. *Cognitive Science*, *32 1*, 108-54.

Grant, E., Finn, C., Levine, S., Darrell, T., & Griffiths, T. (2018). Recasting gradient-based meta-learning as hierarchical Bayes. In *International Conference on Learning Representations*.

Griffiths, T. L., Callaway, F., Chang, M. B., Grant, E., Krueger, P. M., & Lieder, F. (2019). Doing more with less: meta-reasoning and meta-learning in humans and machines. *Current Opinion in Behavioral Sciences*, *29*, 24–30.

Hoiem, D., Gupta, T., Li, Z., & Shlapentokh-Rothman, M. (2021). Learning curves for analysis of deep networks. In *International Conference on Machine Learning* (pp. 4287–4296).

Kumar, S., Dasgupta, I., Cohen, J., Daw, N., & Griffiths, T. (2021). Meta-learning of structured task distributions in humans and machines. In *International Conference on Learning Representations*.

Lake, B. M. (2019). Compositional generalization through meta sequence-to-sequence learning. *Advances in Neural Information Processing Systems*, *32*.

Lake, B. M., & Baroni, M. (2023). Human-like systematic generalization through a meta-learning neural network. *Nature*, *623*(7985), 115–121.

Lake, B. M., & Piantadosi, S. T. (2020). People infer recursive visual concepts from just a few examples. *Computational Brain & Behavior*, *3*, 54–65.

Lake, B. M., Salakhutdinov, R., & Tenenbaum, J. B. (2015). Human-level concept learning through probabilistic program induction. *Science*, *350*(6266), 1332–1338.

Landau, B., Smith, L. B., & Jones, S. S. (1988). The importance of shape in early lexical learning. *Cognitive Development*, *3*(3), 299–321.

LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, *521*(7553), 436–444.

Markman, E. M. (1994). Constraints on word meaning in early language acquisition. *Lingua*, *92*, 199–227.

Markman, E. M., & Wachtel, G. F. (1988). Children's use of mutual exclusivity to constrain the meanings of words. *Cognitive Psychology*, *20*(2), 121–157.

McClelland, J. L., Botvinick, M. M., Noelle, D. C., Plaut, D. C., Rogers, T. T., Seidenberg, M. S., & Smith, L. B. (2010). Letting structure emerge: connectionist and dynamical systems approaches to cognition. *Trends in Cognitive Sciences*, *14*(8), 348–356.

McCoy, R. T., Grant, E., Smolensky, P., Griffiths, T. L., & Linzen, T. (2020). Universal linguistic inductive biases via meta-learning. *Proceedings of the 42nd Annual Conference of the Cognitive Science Society*, 737–743.

McCoy, R. T., & Griffiths, T. L. (2023). Modeling rapid language learning by distilling Bayesian priors into artificial neural networks. *arXiv preprint arXiv:2305.14701*.

McKinley, S. C., & Nosofsky, R. M. (1993). Attention learning in models of classification. *Unpublished manuscript*.

Medin, D., Altom, M. W., Edelson, S. M., & Freko, D. (1982). Correlated symptoms and simulated medical classification. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *8*(1), 37-50.

Medin, D., & Schaffer, M. (1978, May). Context theory of classification learning. *Psychological Review*, *85*(3), 207–238.

Medin, D., & Schwanenflugel, P. (1981, May). Linear separability in classification learning. *Journal of Experimental Psychology: Human Learning and Memory,*, *7*, 355–368.

Müller, S., Hollmann, N., Arango, S. P., Grabocka, J., & Hutter, F. (2022). Transformers can do Bayesian inference. In *International Conference on Learning Representations*.

Neisser, U., & Weene, P. (1962). Hierarchies in concept attainment. *Journal of Experimental Psychology*, *64*(6), 640.

Nosofsky, R. M., Palmeri, T. J., & Mckinley, S. C. (1994). Rule-plus-exception model of classification learning. *Psychological Review*, *101 1*, 53-79.

Piantadosi, S. T., Tenenbaum, J. B., & Goodman, N. D. (2016). The logical primitives of thought: Empirical foundations for compositional cognitive models. *Psychological Review, 123(4), 392–424*.

Rosch, E., Mervis, C. B., Gray, W. D., Johnson, D. M., & Boyes-Braem, P. (1976). Basic objects in natural categories. *Cognitive Psychology*, *8*(3), 382–439.

Schmidhuber, J. (1987). *Evolutionary principles in self-referential learning.* Diploma thesis, Technische Universitat Munchen, Germany.

Shepard, R. N., Hovland, C. I., & Jenkins, H. M. (1961). Learning and memorization of classifications. *Psychological monographs: General and applied*, *75*(13), 1.

Snell, J., Swersky, K., & Zemel, R. (2017). Prototypical networks for few-shot learning. *Advances in Neural Information Processing Systems*, *30*.

Thrun, S., & Pratt, L. (Eds.). (1998). *Learning to learn*. USA: Kluwer Academic Publishers.

Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al. (2016). Matching networks for one shot learning. *Advances in Neural Information Processing Systems*, *29*.

Wang, J. X. (2021). Meta-learning in natural and artificial intelligence. *Current Opinion in Behavioral Sciences*, *38*, 90–95.

Xu, F., & Tenenbaum, J. B. (2007). Word learning as Bayesian inference. *Psychological Review*, *114*(2), 245.

Yang, Y., & Piantadosi, S. T. (2022). One model for the learning of language. *Proceedings of the National Academy of Sciences*, *119*(5), e2021865119.

Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., & Wu, Y. (2022). CoCa: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research*.