

# Interpreting Expert Annotation Differences in Animal Behavior

Megan Tjandrasuwita  
Caltech

Jennifer J. Sun  
Caltech

Ann Kennedy  
Northwestern University

Swarat Chaudhuri  
UT Austin

Yisong Yue  
Caltech

## Abstract

*Hand-annotated data can vary due to factors such as subjective differences, intra-rater variability, and differing annotator expertise. We study annotations from different experts who labelled the same behavior classes on a set of animal behavior videos, and observe a variation in annotation styles. We propose a new method using program synthesis to help interpret annotation differences for behavior analysis. Our model selects relevant trajectory features and learns a temporal filter as part of a program, which corresponds to estimated importance an annotator places on that feature at each timestamp. Our experiments on a dataset from behavioral neuroscience demonstrate that compared to baseline approaches, our method is more accurate at capturing annotator labels and learns interpretable temporal filters. We believe that our method can lead to greater reproducibility of behavior annotations used in scientific studies. We plan to release our code.*

## 1. Introduction

Supervised algorithms for animal behavior quantification have become a powerful tool for characterizing the structure of behavior and its regulation by genes and the brain [9, 22, 19, 7]. However, different individuals perceive and describe the world in different ways, and this can create significant inter-annotator and inter-lab differences in the behavioral annotations used to construct such supervised classifiers. In image recognition, variability across individuals have been shown to produce different object categorizations [5] or labels for the same image data [13]. Similarly, annotator variability has been observed in animal behavior studies, even among experts studying the same behaviors [14, 22]. To improve reproducibility and annotator consensus in behavioral experiments, we propose a novel method for automatically generating interpretations of human behavior annotations.

Existing behavior classification models are typically

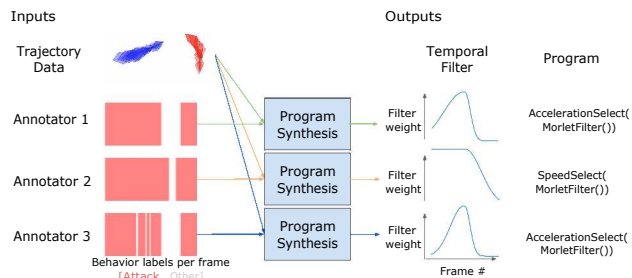


Figure 1. **Overview.** Given trajectory data and behavior labels, we use program synthesis to learn a programmatic description with temporal filters. These programs can be used to compare differences across annotators.

black-box models trained to reproduce human annotations. While these models can achieve high accuracy in the hands of individual labs, it is difficult to interpret differences between models or training sets produced by different individuals [22, 19]. Previous studies have proposed methods for post-hoc interpretation of trained models [17, 20], but the large number of dimensions and parameters in modern machine learning models can make it difficult to understand how annotators use specific features to annotate behavior.

To overcome these limitations, we use program synthesis to generate programmatic descriptions from behavior annotations, which can be interpreted without the need for post-hoc analysis. Program synthesis learns symbolic models from domain-specific languages [24, 25, 23, 3]. We introduce a domain-specific language for behavior classification, which includes learnable temporal filters and feature selections to identify behaviorally relevant features of animal movement. We incorporate our setup into an existing program synthesis method [23] to jointly search through the combinatorially large space of program architectures and optimize parameters. Our approach produces a program with temporal filters for modeling expert annotations, which domain experts qualitatively found to be interpretable for behavior analysis. Our contributions are:

- We combine program synthesis with temporal filtering to generate explanations of behavior annotations.
- We demonstrate our approach on an animal behavior dataset annotated by nine expert annotators.
- We integrate our interpretable programs with an existing tool from domain experts (Bento [22]).

## 2. Related Work

**Behavior Modeling.** Automated behavior quantification has enabled scalable analysis of behavioral data in neuroscience and ethology [9, 1]. These methods often begin with tracking of anatomically defined keypoints from recorded videos of behaving animals [22, 18]. Domain-specific features are then computed from trajectory data and used to train behavior classifiers in the form of neural networks or large random forests [22, 19], which are not easily interpretable. Instead, in our approach, we will search through these domain-specific features using program synthesis to produce programs.

**Interpretable Models.** Existing interpretability techniques in machine learning generally follow one of two approaches: creating post-hoc explanations of black-box models [17, 20, 11] or learning inherently interpretable models [16, 23, 8]. We focus on the second approach, and apply techniques in program synthesis [23] to learn a program to model behavior annotations. We compare our method against models with different levels of interpretability, from shallow decision trees to 1-D Convolutional Networks.

We note that there exists a discussion on when machine learning models are interpretable [21, 10, 15, 12]. In our work, we focus on our target users, who are domain experts in neuroscience. We work with domain experts to design a DSL which is qualitatively interpretable for them.

## 3. Approach

We consider program learning in the context of sequence classification. We train a program that predicts behavior annotations at each frame from trajectory data, and use this program as a description of an annotator’s annotation style.

### 3.1. Problem Formulation

We adopt a problem formulation similar to NEAR [23]. A program is written in a domain-specific language (DSL) and is defined as  $(\alpha, \theta)$ , where  $\alpha$  is a discrete program architecture and  $\theta$  is a vector of real-valued parameters. We denote the semantics of an architecture by  $\mathcal{S}(\alpha, \theta)$ , which is a function parameterized by  $\theta$  and applied to input  $x$ .

Our goal is to find a program that is both accurate (low prediction error) and interpretable (low structural cost), which we formulate as solving the following optimization problem:

$$(\alpha^*, \theta^*) = \arg \min_{(\alpha, \theta)} (\mathcal{S}(\alpha) + \zeta(\alpha, \theta)). \quad (1)$$

Here,  $\zeta(\alpha, \theta) = E_{(x, y)} [\mathbb{I}(\mathcal{S}(\alpha, \theta) \neq y)]$  is the standard notion of prediction error. Since interpretability is a motivating factor, we incentivize short programs by penalizing structural complexity  $\mathcal{S}(\alpha)$ , defined as follows. We let each rule  $r$  in our DSL carry a non-negative real cost  $\mathcal{S}(r)$ . The structural cost of an architecture  $\alpha$  is  $\mathcal{S}(\alpha) = \sum_{r \in R(\alpha)} \mathcal{S}(r)$ , where  $R(\alpha)$  is a multiset of rules used in  $\alpha$ .

**Program Synthesis.** We search over program architectures in a top-down manner. The search is analogous to building a graph, where the nodes consists of both partial and complete architectures that are type-consistent with the DSL. The complete architectures are required to be goal nodes. The edges each represent a single-step application of the DSL rules, and are formed between either two partial architectures or a partial and a complete architecture.

In our approach, we use the program synthesis algorithm NEAR [23], which learns differentiable programs using an admissible neural heuristic. We note that any program synthesis approach could work within our framework.

### 3.2. Learnable Temporal Filters

We develop a DSL from which program synthesis methods can find interpretable programs, based on the Morlet wavelet [4, 6]. To learn temporal information, our DSL includes a Morlet Filter operation that maps a sequence of vectors to a single vector by taking a weighted sum of the input sequence. The Morlet Filter, denoted by  $\psi$ , first does a one-to-one mapping between frames  $1, \dots, n$  in the sliding window to values  $x_1, \dots, x_n$ , where  $x_i \in [-\pi, \pi]$   $\forall i = 1, \dots, n$ .  $\psi$  is then evaluated at each  $x_i$  and is defined as:

$$\psi(x; s, w) = e^{-0.5 \left( \frac{x}{s/w} \right)^2} \cos(wx), \quad (2)$$

where  $x \in [-\pi, \pi]$ .

The Morlet Filter is parameterized by  $s, w$ , where  $w$  determines the width of the filter and  $s$  controls the wavelet frequency. In our experiments, we use a generalization of the symmetric Morlet Filter by allowing the form of the Morlet Filter to differ between the frames preceding and following the predicted frame. Specifically, the left (preceding) Morlet Filter is parameterized by  $s_1, w_1$  whereas the right (following) is parameterized by  $s_2, w_2$ , resulting in the asymmetric Morlet Filter that we include in our DSL.

Our DSL also includes affine transformations of the following form, where  $W$  is a matrix of weights,  $x$  is a feature vector, and  $b$  is a learned bias:

$$T(x) = W^T x[i_1, \dots, i_n] + b. \quad (3)$$

Given a full feature vector  $x$ , the transformation selects a subset of features at indices  $i_1, \dots, i_n$  and applies a simple

linear layer to the feature subset. For the purpose of interpretability, we limit  $n$  to be 1 or 2, i.e. the transformations either select a single feature, or the two same features for the resident and intruder mice.

Within our DSL, the Morlet Filter operation is differentiable with respect to parameters  $s$ ,  $w$ , allowing the shape of the filter to be discovered through gradient optimization. Similarly, the weights and bias  $W$ ,  $b$  of each affine transformation  $T$  are amenable to gradient descent.

**Disjunctions.** Our DSL allows *disjunctions* of two or more Morlet Filter operations. The output of a complete Morlet Filter program (the Morlet Filter applied to a sequence of feature vectors, followed by an affine transformation) is a logit. A disjunction combines the predictions of each filter by summing up the outputted logits. In order to reduce variability in the programs found by the disjunction, we perform separate runs of NEAR to find each filter in the disjunction. Once a filter in the disjunction is found, its weights are frozen when discovering the structure and optimizing the parameters of the subsequent filters. This encourages each subsequent filter to explain variance in the dataset that has not been captured by the previous filters.

## 4. Experiments

### 4.1. Dataset

We use a subset of the MARS [22] dataset for studying annotator variability, which consists of ten 10-minute videos at 30Hz of socially interacting mice from a standard resident-intruder assay. These videos are independently annotated for three behaviors of interest by each of nine domain experts. As input, we use a subset of domain-specific features from the MARS dataset: 1) for both mice, we compute head-body angle, body axis ratio, speed, acceleration, tangential velocity, social angle, 2) across mice, we compute area ellipse ratio, whether resident is facing intruder, and minimum distance of resident nose to intruder body. Here, we consider two binary classification tasks: interact vs. no-interact, and aggression vs. no-aggression. Interaction is defined as frames on which one mouse is sniffing, attacking, or mounting the other; aggression is defined as periods of high-intensity biting, chasing, or grappling.

### 4.2. Evaluation Procedure

We compare the performance of our discovered programs with the following baselines: 1) Decision Trees, a popular choice for both performance and interpretability; 2) 1D Convolutional Neural Networks, a black-box model that is well-suited for processing temporal signals.

**Decision Trees (DT).** Decision trees are constructed by finding yes/no questions that split the data into the most homogeneous groups. We implement DTs using XGBoost [2], a popular framework for training tree ensembles, and

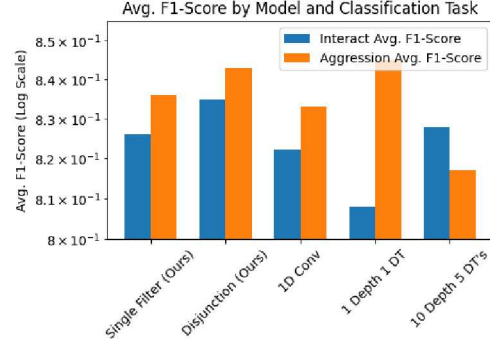


Figure 2. **Performance of Models Trained on 100% Training Data.** Bars reflect mean F1 score of each model when trained and tested separately on each of the nine annotators in the MARS dataset.

test tree classifiers of varying complexity. As input to the decision trees, we pass handcrafted temporal features, produced by convolving our 15 behavior features or their first or second derivatives with a Gaussian filter with standard deviations of 8, 30, or 120 frames. This produced 135 total features: 15 original features \* 3 derivative orders (0, 1, and 2) \* 3 filter widths.

**1D Convolutional Networks (1D Conv).** In a similar manner to a Morlet Filter, a 1D convolutional neural network produces a weighted sum of a given sequence of vectors- however unlike the Morlet Filter, weights are not constrained to have any specific temporal structure. The 1D Conv Net learns a set of weights to convolve with each input feature over time, and the logits from all features are summed for the output predictions.

**Evaluation Details.** We defined a window of  $\pm 5$  seconds centered about the frame for which behavior was to be predicted, and extracted features of animal poses within this window. We then downsampled data from 30Hz to 6Hz, producing vectors of length 61 for each of the 15 features.

We evaluated all models using the F1-score, defined as the harmonic mean of Precision and Recall. We selected 6 videos for training (106k frames), 2 for validation (40k frames), and 2 for test (39k frames). To compare data efficiency, we sub-sampled the training data by randomly sampling trajectories of 1000 frames to achieve desired fractions of the training set size. The sampling also retained a similar class distribution as the full training set. For every data fraction (1%, 10%, 50%), we create three different random samples and train all models three times for each sample. The results are reported on the average across these nine repeats, and across the nine annotators.

### 4.3. Results

**Accuracy.** Synthesized programs with a disjunction of two filters achieve the highest F1 score for detection of in-

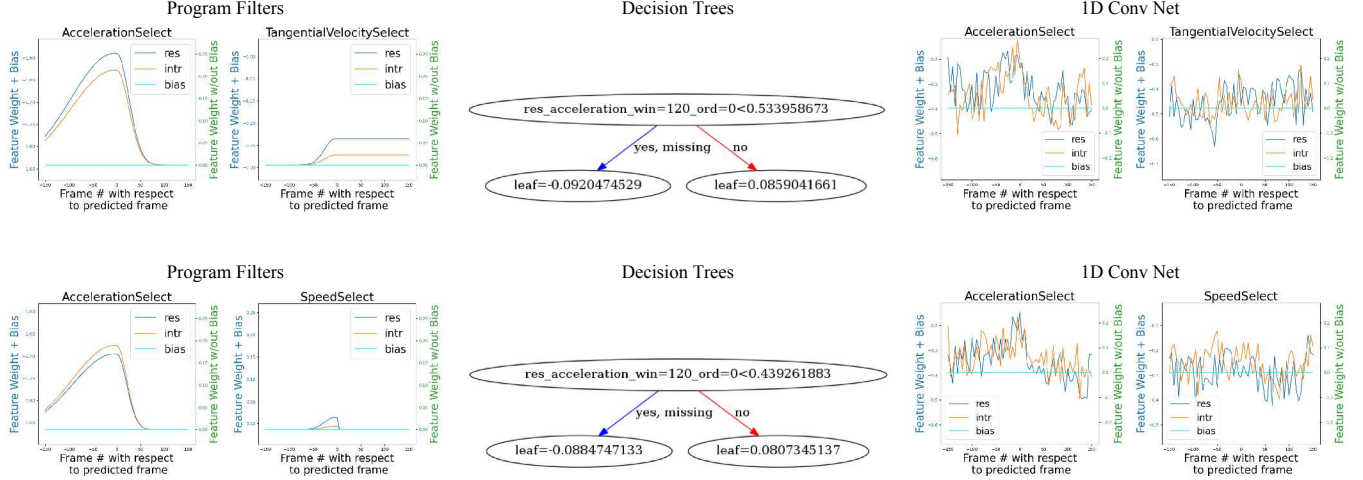


Figure 3. **Comparing models for two annotators.** Each row represents the visualized model trained on aggression vs. nonaggression annotations for one annotator. *Left:* The program filters are from the learned disjunctions, and shows the weight applied at each timestamp for normalized trajectory features from program synthesis. *Center:* Depth 1 decision tree with branches. *Right:* Neural network weights on a subset of input features, matching each annotator’s disjunction features.

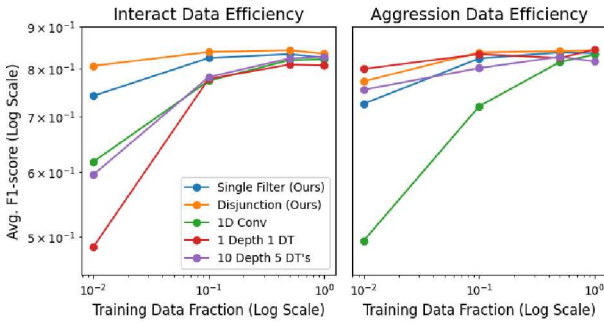


Figure 4. **Data Efficiency on Behavior Sequence Classification.** F1-score averaged across annotators vs. training data fraction on the aggression vs. non-aggression task (left) and interact vs. other task (right).

teraction, and are comparable to the Decision Tree (DT) for detection of aggression (Figure 2). Programs with a single filter had slightly lower F1 compared to the disjunction. For the DTs, the single depth 1 DT is much simpler than 10 depth 5 DTs. Single DT performs better on aggression, which implies that thresholding on one feature is able to classify aggression accurately and the deeper DT is more prone to overfitting. On the other hand, a more complex DT is needed to perform better on interaction.

In terms of data efficiency, disjunctions also remain the highest performing model on interact vs. other (Figure 4). On aggression, disjunctions are comparable to the single depth 1 DT. Because of increased model complexity, the 1D Conv Net is generally less data efficient compared to our model. We verified that the variance in performance of both disjunctions and Morlet Filters are either less than or

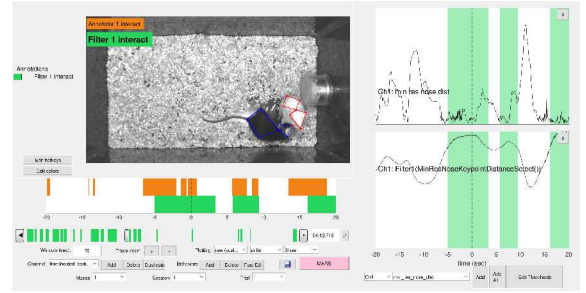


Figure 5. **Visualization of our learned filters in Bento** [22]

comparable to variance found in the baseline models.

**Interpretability.** We next visualized our models and baselines (Figure 3). All three models include some aspect of temporal filtering of the data, however we argue that visualization and interpretation of this filtering is clearest for the disjunctions. The Conv Net filters appear as noisy versions of the disjunction filters, but without the disjunction filters as reference it is difficult for a domain expert to discern their structure. Filtering in the decision tree is implicit (in the names of the features used), and interpreting the numerical thresholds and leaf values is challenging. In contrast, the smoothness of the disjunction filters makes them easy to read, and their asymmetry around the predicted frame allows them to produce a variety of temporal structures. For domain experts to visualize our model more easily, we also added support for visualization of our trained models and their output within Bento [22] (Figure 5).

## 5. Conclusion

We propose a method, based on program synthesis, for learning programmatic descriptions of behavior anno-



tations. We show that our method is accurate compared to baseline methods and that the programs we learn are qualitatively interpretable to domain experts. Automated behavior quantification systems for animal studies are often trained and evaluated on human-provided labels. As a result, human variability will affect their performance. Programmatic explanations for annotated behavior can help us interpret annotation differences towards improving reproducibility of behavioral studies.

## References

- [1] David J Anderson and Pietro Perona. Toward a science of computational ethology. *Neuron*, 84(1):18–31, 2014. [2](#)
- [2] Tianqi Chen, Tong He, Michael Benesty, Vadim Khotilovich, Yuan Tang, Hyunsu Cho, et al. Xgboost: extreme gradient boosting. *R package version 0.4-2*, 1(4), 2015. [3](#)
- [3] Kevin Ellis, Daniel Ritchie, Armando Solar-Lezama, and Joshua B Tenenbaum. Learning to infer graphics programs from hand-drawn images. *arXiv preprint arXiv:1707.09627*, 2017. [1](#)
- [4] Dennis Gabor. Theory of communication. *Journal of the Institution of Electrical Engineers-Part III: Radio and Communication Engineering*, 93(26):429–441, 1946. [2](#)
- [5] Ryan Gomes, Peter Welinder, Andreas Krause, and Pietro Perona. Crowdclustering. *Neural Information Processing Systems*, 2012. [1](#)
- [6] A Grossmann, Richard Kronland-Martinet, and J Morlet. Reading and understanding continuous wavelet transforms. In *Wavelets*, pages 2–20. Springer, 1990. [2](#)
- [7] Robert Evan Johnson, Scott Linderman, Thomas Panier, Caroline Lei Wee, Erin Song, Kristian Joseph Herrera, Andrew Miller, and Florian Engert. Probabilistic models of larval zebrafish behavior reveal structure on many scales. *Current Biology*, 30(1):70–82, 2020. [1](#)
- [8] Jongbin Jung, Connor Concannon, Ravi Shroff, Sharad Goel, and Daniel G. Goldstein. Simple rules for complex decisions, 2017. [2](#)
- [9] Mayank Kabra, Alice A Robie, Marta Rivera-Alba, Steven Branson, and Kristin Branson. Jaaba: interactive machine learning for automatic annotation of animal behavior. *Nature methods*, 10(1):64, 2013. [1, 2](#)
- [10] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. Interpreting interpretability: Understanding data scientists’ use of interpretability tools for machine learning. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2020. [2](#)
- [11] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, 2018. [2](#)
- [12] Isaac Lage, Emily Chen, Jeffrey He, Menaka Narayanan, Been Kim, Samuel J Gershman, and Finale Doshi-Velez. Human evaluation of models built for interpretability. In *AAAI Conference on Human Computation and Crowdsourcing*, 2019. [2](#)
- [13] Thomas A Lampert, Andre’ Stumpf, and Pierre Gancarski. An empirical study into annotator agreement, ground truth estimation, and algorithm evaluation. *IEEE Transactions on Image Processing*, 25(6):2557–2572, 2016. [1](#)
- [14] Xubo Leng, Margot Wohl, Kenichi Ishii, Pavan Nayak, and Kenta Asahina. Quantitative comparison of drosophila behavior annotations by human observers and a machine learning algorithm. *bioRxiv*, 2020. [1](#)
- [15] Zachary C Lipton. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018. [2](#)
- [16] Yin Lou, Rich Caruana, and Johannes Gehrke. Intelligent models for classification and regression. In *ACM SIGKDD international conference on Knowledge discovery and data mining (KDD)*, 2012. [2](#)
- [17] Scott M Lundberg and Su-In Lee. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017. [1, 2](#)
- [18] Alexander Mathis, Pranav Mamidanna, Kevin M Cury, Taiga Abe, Venkatesh N Murthy, Mackenzie Weygandt Mathis, and Matthias Bethge. Deeplabcut: markerless pose estimation of user-defined body parts with deep learning. *Nature neuroscience*, 21(9):1281–1289, 2018. [2](#)
- [19] Simon RO Nilsson, Nastacia L Goodwin, Jia J Choong, Sophia Hwang, Hayden R Wright, Zane Norville, Xiaoyu Tong, Dayu Lin, Brandon S Bentzley, Neir Eshel, et al. Simple behavioral analysis (simba): an open source toolkit for computer classification of complex social behaviors in experimental animals. *BioRxiv*, 2020. [1, 2](#)
- [20] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ”why should i trust you?” explaining the predictions of any classifier. In *ACM SIGKDD international conference on knowledge discovery and data mining (KDD)*, pages 1135–1144, 2016. [1, 2](#)
- [21] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019. [2](#)
- [22] Cristina Segalin, Jalani Williams, Tomomi Karigo, May Hui, Moriel Zelikowsky, Jennifer J. Sun, Pietro Perona, David J. Anderson, and Ann Kennedy. The mouse action recognition system (mars): a software pipeline for automated analysis of social behaviors in mice. *bioRxiv* <https://doi.org/10.1101/2020.07.26.222299>, 2020. [1, 2, 3, 4](#)
- [23] Ameesh Shah, Eric Zhan, Jennifer J. Sun, Abhinav Verma, Yisong Yue, and Swarat Chaudhuri. Learning differentiable programs with admissible neural heuristics. In *NeurIPS*, 2020. [1, 2](#)
- [24] Lazar Valkov, Dipak Chaudhari, Akash Srivastava, Charles Sutton, and Swarat Chaudhuri. Houdini: Lifelong learning as program synthesis. *arXiv preprint arXiv:1804.00218*, 2018. [1](#)

- [25] Halley Young, Osbert Bastani, and Mayur Naik. Learning neurosymbolic generative models via program synthesis. In *International Conference on Machine Learning*, pages 7144–7153. PMLR, 2019. [1](#)